

UNIVERSITATEA "POLITEHNICA" TIMIȘOARA

Facultatea de Automatică și Calculatoare

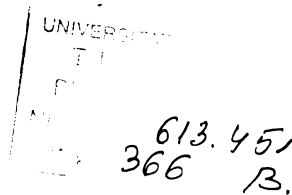
Contribuții la proiectarea sistemelor de operare distribuite

Teză de doctorat

Autor: Ing. Carmen Mușatescu

Conducător științific: Prof. univ. dr. ing. Ioan Jurcă

BIBLIOTECA CENTRALĂ
UNIVERSITATEA "POLITEHNICA"
TIMIȘOARA



Timișoara
1997

G2-52

681.3.066

PREFAȚĂ

Lucrarea de față se înscrie în tendințele actuale de migrare de la prelucrarea centralizată la prelucrarea distribuită, migrare care presupune modificarea modului de proiectare a aplicațiilor, modificarea fluxului prelucrării și controlului datelor precum și adaptarea structurilor organizatorice la noile structuri informaționale, mai ieftine și mai eficiente. Prelucrarea distribuită își propune să pună în valoare, într-un mod cât mai eficient, avantajele descentralizării, astfel încât să răspundă unor noi cerințe ale utilizatorilor, printre care se enumeră: (1) instrumente software pentru comunicare și partajare a informațiilor; (2) luarea deciziilor într-o manieră cooperativă; (3) disponibilitatea permanentă a serviciilor; (3) posibilitatea modificării progresive a echipamentelor hardware și a modulelor software pentru a putea urmări evoluția, foarte rapidă în ultimul domeniu deceniu, a tehnologiei.

O mare parte din aria prelucrării distribuite este acoperită de rețelele locale de date. Ele au cunoscut și continuă să cunoască o dezvoltare extrem de puternică, materializată atât în numărul de rețele locale existente cât și în numărul de aplicații efectiv implementate. Se apreciază că acesta este, probabil, domeniul cel mai dinamic al sistemelor de prelucrare distribuită.

Obiectivul lucrării de față se înscrie pe linia proiectării unor subsisteme software care să poată fi utile în proiectarea de aplicații distribuite în rețele locale. Implementarea acestor suporturi a fost realizată în sistemele de operare Unix, Microsoft Windows - utilizându-se funcții Win 32, implementările sunt valabile atât în Windows NT cât și în Windows 95 - și în Novell NetWare.

Lucrarea reprezintă rezultatul activității desfășurate de autor pe parcursul a 7 ani, în cadrul Catedrei de Calculatoare a Facultății de Automatică, Calculatoare și Electronică, Universitatea din Craiova. În cadrul lucrării, am abordat atât aspectele teoretice și problematica sistemelor de operare distribuite cât și aspectele practice legate de concepția, proiectarea și implementarea unor subsisteme ale sistemelor de operare distribuite, care să asiste utilizatorul în dezvoltarea de aplicații distribuite.

Pe toată durata elaborării tezei, am beneficiat de îndrumarea atentă, exigentă și competentă a conducătorului științific, prof. dr. ing. Ioan Jurcă, care, cu generozitate și înțelegere mi-a conturat domeniul pe care l-am abordat în activitatea de doctorat. Pentru sprijinul profesional și moral, pentru amabilitatea deosebită pe care mi-a acordat-o și nu în ultimul rând pentru bibliografia pe care mi-a pus-o la dispoziție, pe tot parcursul elaborării lucrării, țin să-mi exprim întreaga recunoștință, considerație și stimă, alături de cele mai respectuoase mulțumiri. Acest lucru are o rezonanță mai profundă întrucât apreciez că prof. dr. ing. Jurcă și-a adus o contribuție esențială în formarea personalității mele în calitate de cadru didactic la Catedra de Calculatoare, Facultatea de Automatică, Calculatoare și Electronică, Universitatea din Craiova.

Aduc, de asemenea, mulțumiri întregii Catedre de Calculatoare a Facultății de Automatică și Calculatoare, Universitatea "Politehnica" Timișoara pentru modul în care m-a primit în mijlocul său pe parcursul activității de doctorat, pentru încurajările susținute pe care mi le-a adresat.

Datarez mulțumiri și Colectivului Catedrei de Calculatoare, Facultatea de Automatică, Calculatoare și Electronică, Universitatea din Craiova pentru răbdarea de care a dat dovadă în situațiile în care am "monopolizat" uneori rețelele locale aflate în dotare.

Autorul,

CUPRINS

INTRODUCERE.....	intro -1
PARTEA I: SISTEME DE OPERARE DISTRIBUITE	
CAPITOLUL 1. SISTEME DE OPERARE DISTRIBUITE. SUPORTUL OFERIT PRELUCRĂRILOR DISTRIBUITE, PROBLEME SPECIFICE, MODELE DE PROIECTARE	1
1.1. Introducere.....	2
1.2. Prelucrarea distribuită a datelor.....	2
1.2.1. Obiectivele sistemelor distribuite.....	2
1.2.2. Modele actuale în prelucrarea distribuită.....	3
1.2.3. Interacțiunea prelucrare distribuită - sistem de operare	4
1.2.3.1. Proiectarea unui suport pentru conectarea aplicațiilor în mediu distribuit ...	5
1.2.3.2. OSF DCE - Strategie și arhitectură.....	8
1.3. Sisteme de operare pentru rețele de calculatoare și sisteme de operare distribuite.....	10
1.3.1. Specificul sistemelor de operare pentru rețele de calculatoare	11
1.3.1.1. Limbajul de comandă.....	11
1.3.1.2. Procesele agent.....	12
1.3.2. Specificul sistemelor de operare distribuită.....	13
1.3.2.1. Gestionarea proceselor	13
1.3.2.2. Gestiunea fișierelor	18
1.4. Modele în proiectarea sistemelor de operare distribuite	24
1.4.1. Modelul proces	24
1.4.2. Modelul obiect în proiectarea sistemelor de operare distribuite.....	25
 CAPITOLUL 2. COMUNICAȚII ÎN SISTEME DISTRIBUITE. PROTOCOALE.	 27
2.1. Conceptul de arhitectura deschisă. Stratificarea ierarhică a protocoalelor.....	27
2.1.1. Modelul arhitectural ISO/OSI.....	28
2.1.2. Nivele dependente de rețea. Prezentare generală.....	31
2.1.3 Nivelul transport, principalul furnizor de servicii pentru sistemele de operare.....	38
2.1.4. Nivele suport pentru aplicații.....	47
2.1.5. Soluții de implementare a modelului OSI.....	49
2.1.6. Importanța modelului OSI pentru sisteme distribuite.....	50
2.2. Exemple de protocoale.....	51
2.2.1. Familia de protocoale TCP/IP.....	51
A. Protocoale IP. Noua generație de protocol IP, versiunea 6.....	52
B. Protocolul TCP.....	55
C. Comunicarea prin socluri în rețele UNIX și Microsoft Windows (Win32).....	57
2.2.2. Familia de protocoale Novell IPX/SPX	60
2.3. Aspecte specifice ale suportului de comunicație pentru sistemele de operare distribuite..	61
2.4. Protocoale cu caracter universal și protocoale cu caracter restrâns.....	63

CAPITOLUL 3. PARADIGME ALE INTERACȚIUNII DINTRE PROCESE ÎN SISTEME DISTRIBUITE	64
3.1. Modelul client-server.....	65
3.1.1. Prezentare generală a modelului client/server.....	65
3.1.2. Concurența la nivelul serverelor și clienților.....	67
3.1.3. Protocoale pentru modelul client-server.....	73
3.1.4. Sisteme client-server.....	75
3.2. Comunicarea prin mesaje.....	77
3.2.1. Probleme specifice ale comunicării prin mesaje.....	77
3.2.2. Primitive cu și fără blocare.....	78
3.2.3. Primitive cu și fără bufferare.....	79
3.2.4. Primitive sigure și nesigure.....	80
3.2.5. Forme structurate ale comunicării prin mesaje.....	81
3.2.6. Adresare directă și indirectă.....	81
3.2.7. Excepții.....	83
3.2.8. Exemplu. Comunicația prin porturi în Mach.....	83
3.3. Comunicarea între procese prin conducte.....	85
3.3.1. Prezentarea generală a comunicației prin conducte.....	85
3.3.2. Comunicarea prin conducte în Microsoft WIN32.....	85
3.4. Modelul apelului de procedură la distanță.....	86
3.4.1. Prezentarea generală a modelului.....	87
3.4.2. Transmiterea parametrilor și rezultatelor.....	88
3.4.3. Legarea clienților la servere.....	89
3.4.4. Probleme de concurență.....	91
3.4.5. Semantica RPC în prezența defectelor.....	91
3.4.6. Protocoale pentru RPC.....	92
3.4.7. Extensii ale mecanismului de bază RPC.....	93
3.4.8. Aspecte critice în proiectarea comunicării prin RPC.....	94
3.4.9. Exemple.....	96
A. Modelul RPC în DCE	96
B. Microsoft RPC.....	98
3.5. Comunicații de grup.....	101
3.5.1. Probleme specifice de comunicații de grup.....	101
3.5.2. Primitive și protocoale pentru comunicații de grup.....	101
3.5.3. Exemple: MCL, ISIS, Amoeba.....	102
 CAPITOLUL 4. SINCRONIZAREA ÎN SISTEME DE OPERARE DISTRIBUITE	 105
4. 1. Ordonarea evenimentelor în sisteme distribuite	105
4.2. Soluții centralizate pentru sincronizare.....	106
4.2.1. Utilizarea unui ceas fizic global.....	106
4.2.2. Utilizarea unui proces central.....	106
4.2.3. Contori de evenimente și secvențiatori.....	107
4.3. Soluții distribuite pentru sincronizare.....	109
4.3.1. Sincronizarea ceasurilor.....	109
4.3.2. Algoritmi distribuți pentru excludere reciprocă.....	112
4.3.2.1. Algoritmul lui Lamport.....	112
4.3.2.2. Algoritmul lui Ricard și Agrawala.....	114
4.3.2.3. Algoritm pentru excludere reciprocă în accesul la K resurse identice.....	116
4.3.2.4. Alți algoritmi.....	118
4.3.2.5. Considerații finale.....	118

4.3.2.6. Implementarea algoritmulor distribuției pentru excludere reciprocă în rețele Unix, Microsoft Windows (Win32) și Novell NetWare.....	119
4.3.3. Algoritmi care utilizează un token circulant pentru excludere reciprocă.....	120
4.3.3.1. Sisteme structurate în inel.....	120
4.3.3.2. Algoritmul pentru sisteme nestructurate în inel.....	122
4.3.3.3. Secvențiator circulant.....	123
4.3.3.4. Implementarea algoritmulor care utilizează un token circulant în rețele Unix, Microsoft Windows (Win32) și Novell NetWare.....	124
4.3.4. Algoritmi broadcast.....	125
4.3.4.1. Semafoare distribuite.....	125
4.3.4.2. Contori de evenimente distribuiți.....	127
4.3.4.3. Secvențiatori distribuți.....	129
4.3.4.3. Implementarea algoritmulor broadcast în rețele Unix și Microsoft Windows (Win32).....	130
4.4. Algoritmi de alegere.....	131
4.4.1. Algoritmul Bully.....	131
4.4.2. Algoritm în inel.....	132
CAPITOLUL 5. PROIECTAREA ȘI IMPLEMENTAREA UNOR MECANISME SUPORT PENTRU DEZVOLTAREA DE APLICAȚII DISTRIBUITE ÎN REȚELE UNIX, MICROSOFT WINDOWS (WIN32) ȘI NOVELL NETWARE	133
5.1. Proiectarea unui toolkit pentru controlul accesului concurrent la o bază de date în rețele locale Unix, Windows (Win32) sau Novell NetWare.....	133
5.1.1. Soluții bazate pe semafoare.....	133
A. Prima soluție bazată pe semafoare.....	134
B. A doua soluție cu semafoare.....	135
C. A treia soluție cu semafoare.....	136
5.1.2. Soluție bazată pe un proces coordonator.....	137
5.1.3. Implementări în rețele Unix, Microsoft Windows (Win32) și Novell NetWare.....	139
5.1.3.1. Implementarea soluțiilor bazate pe semafoare în rețele Novell Netware.....	139
5.1.3.2. Implementarea soluțiilor bazate pe semafoare în rețele Unix și Windows (Win32).....	139
5.1.3.3. Implementarea soluției bazată pe un proces coordonator în rețele Novell Netware.....	140
5.1.4. Compararea soluțiilor propuse și concluzii.....	140
5.2. Proiectarea și implementarea unui toolkit pentru comunicații de grup	141
5.2.1. Cerințe de proiectare.....	141
5.2.2. Varianta semidistribuită.....	145
5.2.2.1. Arhitectura generală a sistemului TGC.....	145
5.2.2.2. Protocolul de comunicație.....	146
5.2.2.3. Descrierea serverelor locale și a serverului central.....	147
5.2.3. Alte variante	159
5.2.3.1. Variante complet distribuite	159
A. Varianta în care se folosește transmisia broadcast	159
B. Varianta în care nu se folosește transmisie broadcast.....	160
5.2.3.2. Varianta în inel.....	161
5.2.4. Implementarea variantelor complet distribuite.....	161
5.2.4.1. Implementarea variantei cu transmisie broadcast.....	161
5.2.4.2. Implementarea variantei cu transmisie în inel.....	166
5.2.5. Comparatie între variantele propuse și concluzii.....	171

**PARTEA a II a: CRITERII DE PROIECTARE A SISTEMELOR DE PROGRAMARE
DISTRIBUITĂ BAZATĂ PE OBIECTE (SPDBO). PROIECTAREA UNUI SUPPORT
SOFTWARE BAZAT PE MODELUL TRANZACȚIILOR ATOMICE PENTRU SPDBO
ÎN REȚELE UNIX ȘI WINDOWS**

CAPITOLUL 6. SISTEME DE PROGRAMARE DISTRIBUITĂ BAZATE PE OBIECTE (SPDBO). CRITERII DE PROIECTARE	161
6.1. Concepte fundamentale.....	161
6.2. Structura obiectelor.....	162
6.3. Gestionarea obiectelor.....	164
6.3.1. Gestionarea acțiunilor.....	164
6.3.2. Sincronizarea acțiunilor.....	164
6.3.3. Siguranța obiectelor.....	165
6.4. Interacțiunea obiectelor.....	166
6.4.1. Localizarea obiectelor.....	166
6.4.2. Apelul de operații.....	166
6.4.3. Detectarea eșecurilor în apelarea operațiilor.....	167
6.5. Reprezentarea obiectelor în memorie.....	168
CAPITOLUL 7. CONTROLUL CONCURENȚEI ȘI RECUPERAREA OBIECTELOR	169
7.1. Introducere.....	169
7.2. Modelul tranzacțional.....	170
7.2.1. Consistența obiectelor. Noțiunea de tranzacție.....	170
7.2.2. Modelul de lucru pentru sistemul distribuit de obiecte.....	171
7.2.3. Cerințele modelului sistemului distribuit de obiecte.....	173
7.3. Realizarea atomicității tranzacțiilor	173
7.3.1. Protocoale pentru comiterea tranzacțiilor.....	173
7.3.2. Protocolul CCR în OSI.....	175
7.3.3. Implementarea tranzacțiilor atomice.....	176
7.3.4. Recuperarea obiectelor.....	177
7.3.4.1. Protocolul Write-Ahead Logging.....	178
7.3.4.2. Tehnicile Value-Logging și Operation Logging.....	178
7.3.4.3. Tehnicile Redo-Only Logging și Undo-Only Logging.....	179
7.3.4.4. Structura sistemului de recuperare.....	179
7.4. Controlul concurenței.....	180
7.4.1. Cerințele mecanismului de control al concurenței.....	180
7.4.2. Serializarea tranzacțiilor.....	180
7.4.3. Mecanisme tradiționale pentru controlul concurenței.....	184
7.4.4. Mecanisme bazate pe blocarea cu lock-uri.....	184
7.4.4.1. Mecanismul blocării în două faze (2PL).....	185
7.4.4.2. Algoritmi de blocare ierarhici.....	186
7.4.4.3. Grafuri aciclice de lock-uri.....	188
7.4.5. Mecanisme bazate pe marca de timp.....	188
7.4.5.1. Concepte de bază.....	188
7.4.5.2. Algoritmul fundamental pentru controlul concurenței	189
7.4.5.3. Optimizări ale algoritmului de bază.....	190
7.4.5.4. Algoritmul conservativ de ordonare a mărcilor de timp.....	191
7.4.6. Controlul optimist al concurenței.....	192

7.4.6.1. Concepte de baza.....	192
7.4.6.2. Algoritmi pentru controlul optimist al concurenței.....	192
7.4.7 Extinderea mecanismelor tradiționale pentru controlul concurenței.....	194
7.4.7.1 Cerințe ale controlului concurenței în aplicațiile avansate.....	194
7.4.7.2 Extinderea mecanismelor bazate pe serializabilitate și renunțarea la serializabilitate.....	195
7.4.8. Mecanisme pentru controlul concurenței în aplicații avansate.....	195
7.4.8.1 Mecanismul blocării altruiste.....	195
7.4.8.2 Mecanismul validării instanței.....	197
7.4.9. Mecanisme pentru controlul concurenței bazate pe informații semantice.	
Sincronizarea accesului la tipuri de date abstracte.....	98
7.4.9.1. Relația de dependență.....	198
7.4.9.2. Specificarea tipurilor abstracte.....	200
7.4.9.3. Metode de sincronizare.....	205
CAPITOLUL 8. CONTROLUL INTERBLOCĂRII OBIECTELOR	208
8.1. Metode de control al interblocării.....	208
8.2. Prevenirea interblocărilor.....	208
8.3. Detecția interblocărilor în sisteme distribuite.....	209
8.3.1 Modelul de lucru.....	209
8.3.2 Clase de algoritmi pentru detecția interblocării.....	210
8.3.2.1. Clasificarea algoritmilor.....	210
8.3.2.2 Modele de interblocare.....	210
8.3.2.3 Clase de algoritmi distribuiți pentru detecția interblocării.....	211
8.3.3 Algoritmi centralizați pentru detecția interblocării.....	211
8.3.4. Algoritmul ierarhic pentru detecția interblocării.....	212
8.3.5. Algoritmi distribuiți pentru detecția interblocării.....	212
8.3.5.1. Algoritm edge-chasing pentru detecția interblocării în modele AND.....	212
8.3.5.2 Algoritm cu difuzie pentru detecția interblocării în modele OR.....	213
CAPITOLUL 9. PROIECTAREA UNUI SUPORT SOFTWARE PENTRU PROGRAMARE DISTRIBUITĂ BAZATĂ PE OBIECTE SI PE TRANZACȚII ATOMICE	214
9.1. Cerințe de proiectare.....	214
9.2. Modelul sistemelor DOSTP.....	217
9.3. Arhitectura generală a sistemului DOSTP.....	218
9.3.1. Procesele utilizator ale sistemului DOSTP.....	219
9.3.2. Managerul tranzacțiilor.....	219
9.3.3. Serverele de comunicații.....	219
9.3.4. Serverele de nume.....	220
9.3.5. Managerul central.....	220
9.3.6. Managerii de obiecte.....	220
9.3.7. Sistemul de recuperare.....	220
9.4. Proiectarea elementelor componente ale sistemului DOSTP.....	221
9.4.1. Structuri de date ale sistemului.....	221
9.4.2. Integrarea proceselor de aplicație în sistemul DOSTP.....	224
9.4.3. Proiectarea managerilor de tranzacții.....	227
9.4.4. Proiectarea serverelor de comunicații.....	234
9.4.5. Proiectarea serverelor de nume.....	239
9.4.6. Proiectarea managerului central.....	240
9.4.7. Proiectarea managerilor de obiecte.....	241

CAPITOLUL 10. IMPLEMENTAREA SISTEMULUI DOSTP ÎN REȚELE UNIX ȘI MICROSOFT WINDOWS (WIN32)	244
10.1. Particularități de implementare.....	244
10.1.1. Particularități de implementare în rețele Unix.....	244
10.1.2. Particularități de implementare în rețele Microsoft Windows (Win 32).....	246
10.2. Exemple de desfășurare a comunicațiilor în sistem.....	247
10.2.1. Conectarea nodurilor și a managerilor de obiecte la sistemul tranzacțional.....	247
10.2.2. Inițierea unei tranzacții.....	248
10.2.3. Exemple de implicare a obiectelor în tranzacții.....	248
10.2.4. Exemple de desfășurare a protocolului two-phase commit.....	250
10.2.5. Exemple de tratare a situațiilor de recuperare.....	251
10.3. Exemplu de implementare al unui manager de obiecte.....	253
10.3.1. Descrierea managerului de obiecte.....	253
10.3.1.1. Organizarea managerului de obiecte.....	254
10.3.1.2. Concurența la nivelul managerului de obiecte.....	255
10.3.1.2. Recuperarea obiectelor.....	257
10.3.2. Conectarea obiectelor în sistemul DOSTP.....	258
10.3.2.1. Descrierea clasei CObjectManager.....	258
10.3.2.2. Descrierea clasei CDispatcher.....	258
10.3.2.3. Integrarea automată a obiectelor în sistemul tranzacțional.....	258
10.3.3. Particularități ale implementării în sistemele Unix și Windows.....	260
10.3.4. Exemplu de aplicație care apelează managerii de obiecte de tipul proiectat.....	261
10.4. Concluzii și direcții de dezvoltare ulterioară.....	262
 CAPITOLUL 11. CONCLUZII GENERALE ȘI CONTRIBUȚII ORIGINALE	 263
11.1. Concluzii generale.....	263
11.2. Contribuții originale, valoare aplicativă și direcții de dezvoltare ulterioară.....	269
 ANEXE	
 ANEXA A. IMPLEMENTAREA UNUI TOOLKIT PENTRU COMUNICAȚII DE GRUP. MANAGERUL DE GRUPURI PENTRU VARIANTELE DISTRIBUITE	 anexaA-1
 ANEXA B. SINCRONIZARE ÎN SISTEME DISTRIBUITE	
B1. Exemple de programe client care, pentru realizarea excluderii reciproce în rețele locale, folosesc un server care implementează unul din algoritmi prezentați în capitolul 4, la 4.3.2 și 4.3.3.....	anexaB-1
B2. Implementarea algoritmului Lamport pentru obținerea excluderii reciproce - varianta Unix.....	anexaB-4
B3. Implementarea algoritmului Ricard și Agrawala pentru obținerea excluderii reciproce varianta Novell NetWare.....	anexaB-7
B4. Implementarea algoritmului pentru obținerea excluderii reciproce în accesul la k resurse -varianta Windows.....	anexaB-9
B5. Implementarea algoritmului bazat pe un token circulant pentru obținerea excluderii reciproce- varianta Novell NetWare.....	anexaB-11
B6. Implementarea algoritmului Chandy pentru obținerea excluderii reciproce -varianta Unix.....	anexaB-12
B7. Implementarea semafoarelor distribuite -varianta Windows.....	anexaB-15
B8. Implementarea contorilor de eveniment distribuiți -varianta UNIX.....	anexaB-18

B9. Implementarea secvențiatorilor distribuți -varianta UNIX.....	anexaB-20
B10. Implementarea excluderii reciproce folosind secvențiatori și contori de eveniment distribuți- varianta Unix.....	anexaB-24
ANEXA C. PROIECTAREA UNUI TOOLKIT PENTRU CONTROLUL ACCESULUI LA O BAZĂ DE DATE ÎN REȚELE NOVELL, UNIX ȘI WINDOWS	
C1. Implementarea primei soluții bazate pe semafoare în rețele Novell.....	anexaC-1
C2. Implementarea celei de-a doua soluții bazate pe semafoare în rețele Novell.....	anexaC-3
C3. Implementarea celei de-a treia soluții bazate pe semafoare în rețele Novell.....	anexaC-5
C4. Implementarea soluției bazate pe un proces coordonator în rețele Novell.....	anexaC-6
C5. Implementarea celei de-a treia soluții bazate pe semafoare în rețele UNIX, și Windows.....	anexaC-13
ANEXA D. IMPLEMENTAREA SISTEMULUI DOSTP ÎN REȚELE UNIX ȘI WINDOWS	
D1. Definierea clasei CObject.....	anexa D-1
D2. Definierea clasei CVector.....	anexa D-1
D3. Definierea clasei CObjectManager.....	anexa D-2
D4. Definierea clasei CDispatcher.....	anexa D-3
D5. Definierea macroinstrucțiunilor pentru conectarea automată a obiectelor în sistemul DOSTP.....	anexa D-4
ANEXA F. COMPARAȚIE ÎNTRE PERFORMANȚELE VARIANTELOR DISTRIBUITE ȘI SEMIDISTRIBUITE ALE SISTEMULUI TGC	anexa F-1
ANEXA E. COMPARAȚIE ÎNTRE PERFORMANȚELE TRANZACȚIILOR DISTRIBUITE ȘI CELE ALE TRANZACȚIILOR LOCALE ÎN SISTEMUL DOSTP	anexa E-1
BIBLIOGRAFIE.....	biblio-1

INTRODUCERE

Se poate considera că aplicațiile distribuite își datorează apariția și dezvoltarea lor actuală, progreselor notabile realizate în tehnologiile microelectronicii, calculatoarelor și comunicațiilor, progrese care au condus rapid la un raport performanță/cost avantajos.

Problemele unui sistem de operare distribuit sunt parțial similare cu cele ale unui sistem de operare clasic, centralizat. Diferențele de bază se datorează însă unor aspecte și probleme noi apărute, diferite de cele specifice sistemelor centralizate; acestea constau în acuitatea problemelor de *sincronizare*, care, în contextul unui sistem distribuit nu se pot rezolva fără *comunicare*, întârzierilor aleatorii în comunicații ceea ce generează, în consecință, natura probabilistică a informațiilor de stare, materializate prin viziuni diferite de la proces la proces, și nu în ultimul rând eterogenității elementelor de prelucrare care compun sistemul.

Se poate aprecia că modificările în structura și funcțiile sistemelor de operare distribuite au fost rezultatul presiunii din mai multe direcții:

(A) pe de o parte, modificările de arhitectură ale procesoarelor au avut implicații majore în proiectarea sistemelor de operare. În acest sens se pot distinge, cel puțin trei aspecte majore ale modificărilor: (1) creșterea vitezei de lucru a procesoarelor, care acutizează problemele de acces la memorie și timpul irosit pentru intrarea/ieșirea din nucleul sistemului de operare; (2) prelucrarea multiprocesor; (3) spații de adresa (virtuale) foarte mari. Un impact major l-a avut aspectul (1); acesta a generat două tendințe: folosirea masivă a *thread*-urilor și deplasarea unor funcții ale sistemelor de operare de la nivelul nucleului la nivelul utilizator; ultima tendință furnizează sistemului de operare și o flexibilitate crescută. Aspectul (3) implica posibilitatea folosirii unui singur spațiu global de adrese virtuale, partajat de toate aplicațiile; aceasta implică, evident, probleme referitoare la fragmentarea spațiului de adrese și protecția sa dar simplifică problemele de adresare.

(B) pe de altă parte, la cealaltă extremă, modificările calitative de la nivelul programelor de aplicație impun noi cerințe sistemelor de operare. Printre acestea pot fi enumerate: (1) *distribuția masivă*, care ridică probleme legate de localizarea resurselor, securitate, autentificare, criptare și eterogenitate și *comunicații de grup* (2) *tranzacțiile atomice*, care sunt realizate din ce în ce mai mult la nivelul sistemelor de operare distribuite, fie implementate la nivelul nucleului, fie în servere separate și accesibile prin intermediul unor biblioteci puse la dispoziția utilizatorilor (3) *aplicațiile multi-media* care introduc o caracteristică timp-real și deci impun reconsiderarea aspectelor specifice prelucrării timp-real.

Lucrarea de față are în vedere două scopuri:

(1) un studiu teoretic asupra problematicei specifice sistemelor de operare distribuite și interacțiunii sistem de operare distribuit -prelucrare distribuită. Am considerat necesar acest lucru, deoarece se apreciază că acest domeniu se află în prima fază a maturității sale, aceea de identificare a paradigmatelor și problemelor sale fundamentale fără ca acestea să fie complet soluționate, nici sub aspect teoretic nici sub aspect practic. Problemele clasice ale sistemelor de operare capătă o altă natură în contextul elementelor multiple de prelucrare și al distribuirii datelor și controlului și, în consecință, necesită metode noi de soluționare teoretică și practică.

(2) proiectarea și implementarea unor mecanisme suport pentru dezvoltarea de aplicații distribuite în rețele locale UNIX, Microsoft Windows (Win 32) și Novell NetWare. Aceste mecanisme sunt puse la dispoziția utilizatorilor sub forma unor biblioteci de programe care furnizează API-urile necesare pentru: (a) controlul accesului concurrent la o bază de date (b) comunicații de grup (c) utilizarea unui sistem distribuit de obiecte prin intermediul tranzacțiilor atomice.

Lucrarea este structurată în două mari părți: prima parte cuprinde aspectele teoretice referitoare la problemele specifice unui sistem de operare distribuit și proiectarea mecanismelor suport (a) și (b) de mai sus, iar partea a doua descrie suportul teoretic necesar proiectării unui mecanism de tipul (c) și apoi prezentarea soluțiilor de proiectare și implementare a acestuia în rețele locale.

Cuprinsul celor unsprezece capitole, în care este structurată lucrarea, este următorul:

Partea I: SISTEME DISTRIBUITE

Capitolul 1: Sisteme de operare distribuite. Suportul oferit prelucrărilor distribuite: probleme specifice, modele de proiectare.

În primul capitol se urmărește stabilirea unei definiții a sistemelor de operare distribuite, se precizează obiectivele generale urmărite de sistemele cu prelucrare distribuită și criteriile pe baza cărora se definesc sistemele distribuite și se reliefează principalele avantaje ale sistemelor distribuite în raport cu sistemele centralizate. Sunt prezentate de asemenea, modelele actuale din prelucrarea distribuită, modelul ierarhic, modelul federativ, modelul partajării resurselor, modelul omogen și modelul fluxului de date, termenul de prelucrare distribuită fiind asimilat în accepțiunea de *execuție a unui program, sub forma unor procese concurente, în mai multe noduri*. Este prezentat de asemenea, comparativ, specificul sistemelor de operare pentru rețelele de calculatoare și cel al sistemelor de operare distribuite. Sunt analizate funcțiile de gestiune a proceselor și de gestiune a fișierelor ale sistemelor de operare distribuite și sunt prezentate cele două modele de proiectare ale unui sistem distribuit: *modelul proces* și *modelul obiect*, diferențele majore între cele două modele de proiectare constând în modul în care sunt implementate conceptul de entitate funcțională și sincronizarea. O atenție deosebită este acordată interacțiunii dintre sistemul de operare distribuit - prelucrare distribuită. În acest sens, sunt prezentate două exemple de suporturi software puse la dispoziția utilizatorilor pentru dezvoltarea de aplicații distribuite. Un prim exemplu îl constituie mediul OSF DCE; sunt arătate arhitectura generală DCE, serviciile fundamentale DCE și interacțiunea dintre servicii. Cel de-al doilea exemplu este sistemul TDL, un suport software proiectat pentru conectarea programelor în rețele UNIX, cunoscându-se modul de comunicare al acestora în cadrul unei aplicații.

Capitolul 2. Comunicații în sisteme distribuite. Protocoale.

Obiectivul capitolului 2 îl constituie problemele comunicației în rețele de calculatoare, ale implementării software-ului de comunicație precum și aspecte specifice ale protocoalelor de comunicație în sisteme distribuite. Este prezentat modelul arhitectural ISO/OSI pentru interconectarea sistemelor deschise, insistându-se pe importanța pe care o prezintă nivelul transport pentru sistemele de operare, fiind furnizorul principal de servicii. Se arată două soluții de implementare a modelului OSI și importanța modelului OSI pentru sistemele distribuite. Ca exemple de protocoale utilizate în rețelele locale de calculatoare și folosite și în suporturile software proiectate și prezentate de autor pe parcursul lucrării, sunt amintite pe scurt protocolul IP, insistându-se pe propunerea IP versiunea 6, TCP, IPX/SPX realizându-se legătura cu modul de utilizare al acestora în sistemele de operare UNIX, Windows și Novell NetWare. Comparativ cu protocoalele cu caracter universal folosite în rețelele de calculatoare sunt reliefate și aspectele specifice ale suportului de comunicație pentru sistemele distribuite.

Capitolul 3. Paradigme ale interacțiunii dintre procese în sisteme distribuite.

Capitolul 3 este dedicat prezentării paradigmei interacțiunii dintre procese în sisteme distribuite. Este prezentat *modelul client-server* insistându-se pe proiectarea serverelor; se arată situațiile în care trebuie folosit un server concurent multiproces, un server concurent cu *select*, sau un server iterativ, făcându-se o comparație între performanțele acestora. Modelul client-server poate fi implementat utilizând ambele clase de protocoale prezentate în capitolul 2: protocoale cu caracter universal și protocoale cu caracter restrâns. Pentru obținerea de performanțe superioare este necesar însă, un protocol specific, deci cu caracter restrâns; se exemplifică schimbul de mesaje necesar într-un astfel de protocol client-server.

Un alt aspect descris în acest capitol este cel al comunicării prin mesaje între procese; sunt prezentate problemele specifice ale comunicării prin mesaje, primitive cu și fără blocare, primitive cu și fără blocare, sigure și nesigure, forme structurate ale comunicării prin mesaje și un exemplu de implementare a comunicației prin porturi - comunicația din sistemul Mach.

Ca o altă paradigmă a interacțiunii dintre procese în sisteme distribuite este prezentată comunicarea la distanță între procese prin conducte, exemplificările fiind făcute în Microsoft Win32 în sistemul de operare Windows NT și Windows 95, în ultimul nefiind însă implementată decât partea de client; totodată se propune o soluție pentru utilizarea conductelor în servere multifir în Win32.

Modelul apelului de procedură la distanță este dezbătut pe larg, avându-se în vedere aspecte precum transmiterea parametrilor și rezultatelor, legarea clientilor la servere, probleme de concurență, semantica RPC în prezența defectelor, protocoale pentru RPC. Sunt prezentate și extensii ale mecanismului de bază RPC și analizate aspectele critice în proiectarea comunicației prin RPC. Exemplificările sunt făcute pentru SUN RPC, DCE RPC și Microsoft RPC.

Ca un ultim model de comunicație în sistemele distribuite sunt prezentate comunicațiile de grup, cu exemplificări pentru MCL, ISIS și Amoeba.

Capitolul 4. Sincronizarea în sisteme de operare distribuite.

Capitolul 4 analizează una dintre cele mai importante probleme în sistemele distribuite: aceea a sincronizării proceselor. Termenul de *sincronizare* este folosit în sisteme distribuite pentru a referi trei probleme distincte, dar în ultimă instanță înrudite: (1) sincronizarea între emițătorul și receptorul unui mesaj; (2) specificarea și controlul acțiunilor unor procese care cooperează în cadrul unei aplicații distribuite; (3) serializarea accesului concurrent la obiectele partajate între mai multe procese, problemă cunoscută și în bazele de date sub numele de *controlul concurenței*. În capitolul 4 se tratează aspecte privitoare numai la problemele de tipul (2); problemele de tipul (1) au fost dezbătute în capitolele 2 și 3 iar problemele de tipul (3) vor fi abordate în capitolul 7.

Mecanismele de sincronizare pentru sistemele distribuite vor trebui să determine, printr-o metoda oarecare, o ordonare totală sau parțială a unui set de evenimente. Capitolul 4 prezintă mai întâi aspecte ale ordonării evenimentelor într-un sistem distribuit și apoi, pe baza acestora, două clase mari de mecanisme de sincronizare: *centralizate* și *distribuite*. Unele dintre aceste mecanisme se bazează pe conceptul de *excludere reciprocă* care presupune o ordonare a unor secvențe de cod executate de procesele concurente; aceste mecanisme pot utiliza diferite tehnici precum: un *proces central coordonator* sau un *token circulant*. Alte mecanisme se bazează pe înregistrarea evenimentelor semnificative în cursul unor prelucrări asincrone; acestea utilizează obiecte precum *contori de evenimente* și *secvențiatori* sau *semafoare distribuite*. Sistemizarea algoritmilor distribuți pentru sincronizare s-a făcut în capitolul 4 astfel: (a) în clasa algoritmilor distribuți pentru excludere reciprocă s-au prezentat algoritmul lui Lamport, algoritmul lui Ricard și Agrawala și algoritmul Raymond, pentru excludere reciprocă în accesul la K resurse identice (b) în clasa algoritmilor care utilizează un token circulant a fost prezentat un algoritm pentru sisteme în care nodurile sunt ordonate într-un inel logic și algoritmul Chandy (pentru sisteme în care nodurile nu sunt ordonate într-un inel logic) precum și o variantă de implementare a conceptului de secvențiator circulant; (c) în clasa algoritmilor care utilizează transmisii broadcast s-au proiectat și implementat semafoare distribuite, contori de evenimente distribuți și secvențiatori distribuți.

Se remarcă faptul că, pentru toți algoritmi prezentați, au fost realizate implementări în rețele locale Unix, Microsoft Windows și Novell NetWare. Unele dintre aceste implementări (din lipsă de spațiu, nu toate implementările) sunt prezentate în anexe. Testele s-au făcut în rețele mixte Unix + Windows, respectiv rețele Novell.

În final, în capitolul 4 se prezintă doi algoritmi de alegere a coordonatorului: algoritmul Bully și algoritmul în inel.

Capitolul 5. Proiectarea și implementarea unor mecanisme suport pentru dezvoltarea de aplicații în rețele Unix, Microsoft Windows (Win32) și Novell NetWare.

În finalul primei părți a lucrării, în capitolul 5 se prezintă proiectarea și implementarea a două mecanisme suport pentru dezvoltarea de aplicații în rețele UNIX, Microsoft Windows (Win32) și Novell NetWare: primul mecanism are ca scop controlul accesului concurrent la o bază de date iar cel de-al doilea constă într-un toolkit pentru comunicații de grup. Proiectarea și implementarea acestor mecanisme folosește conceptele teoretice prezentate în capitolele anterioare: protocoale de comunicații prezentate în capitolul 2 (TCP, UDP, IPX), comunicare prin mesaje, concurența la nivelul serverelor, comunicații de grup (capitolul 3), algoritmi pentru obținerea excluderii reciproce, semafoare distribuite (capitolul 4).

Pentru controlul accesului concurrent la o bază de date s-au propus patru soluții: trei soluții bazate pe semafoare și o soluție bazată pe un proces coordonator. Toate soluțiile au fost implementate în rețele mixte Unix + Windows (Win32) și Novell NetWare. Pentru soluțiile bazate pe semafoare în rețele Unix și Windows s-au folosit semafoarele distribuite, proiectate și implementate în capitolul 4, iar pentru rețele Novell NetWare s-a folosit faptul că sistemul de operare NetWare pune la dispoziția programelor de aplicație conceptul de semafor. Pentru proiectarea unui toolkit pentru comunicații de grup s-au propus patru variante: o variantă semidistribuită, două variante complet distribuite în care în una se folosește transmisia broadcast și o variantă în care nodurile sunt organizate într-un inel logic. Se prezintă în detaliu sintaxa API-urilor puse la dispoziția programelor utilizator, în cadrul cerințelor impuse proiectării și se face o comparație a celor patru variante.

Partea a II-a: CRITERII DE PROIECTARE A SISTEMELOR DE PROGRAMARE DISTRIBUITĂ BAZATĂ PE OBIECTE (SPDBO). PROIECTAREA UNUI SUPORT SOFTWARE BAZAT PE MODELUL TRANZACȚIILOR ATOMICE PENTRU SPDBO ÎN REȚELE UNIX ȘI WINDOWS.

Capitolul 6. Sisteme de programare distribuită orientată pe obiecte (SPDBO). Criterii de proiectare

Partea a II-a a lucrării cuprinde suportul teoretic necesar pentru proiectarea sistemelor de programare distribuită bazată pe obiecte precum și proiectarea și implementarea în rețele locale a unui suport software pentru lucrul cu tranzații atomice în aplicații distribuite.

Partea a II-a debutează cu capitolul 6 în care se prezintă problemele majore care trebuie avute în vedere pentru proiectarea unui sistem de programare distribuită. Sunt reliefate aspecte ale componenței și structurii obiectelor, ale gestionării obiectelor incluzând probleme de sincronizare ale operațiilor asupra obiectelor și probleme de siguranță a obiectelor. Un subcapitol special este dedicat interacțiunii obiectelor, cuprinzând probleme de localizare a obiectelor și al modului de apel al operațiilor. Un alt subcapitol este destinat modului de reprezentare a obiectelor în memorie; un accent mai mare este pus pe reprezentarea obiectelor în memoria secundară, deoarece acest lucru este legat de recuperabilitatea obiectelor, o cerință esențială într-un sistem de programare distribuită bazat pe obiecte.

Capitolul 7. Controlul concurenței și recuperarea obiectelor

Se apreciază că proiectarea aplicațiilor distribuite care să ofere siguranță în funcționare poate fi obținută prin două metode generale, conceptual diferite: (1) prima metodă utilizează conceptul de **tranzacție atomică**; (2) a doua metodă utilizează copii multiple ale obiectelor necesare aplicațiilor.

Deoarece în această lucrare se are în vedere prima metodă, în capitolul 7 se prezintă întâi modelul de lucru bazat pe tranzații, apoi soluții pentru realizarea atomicității tranzațiilor în sistemele distribuite; în acest sens este prezentat protocolul *two-phase-commit* și modul lui de implementare în OSI (Commitment Concurrency and Recovery).

Se prezintă apoi tehnici pentru recuperarea obiectelor, precum *Value-Logging* sau *Operation Logging, Redo-Only Logging* sau *Undo-Only Logging* și structura generală a subsistemului de recuperare.

Un subcapitol important din capitolul 7 este dedicat controlului concurenței. După prezentarea cerințelor mecanismelor de control al concurenței și al conceptului de serializare al tranzațiilor se analizează mecanismele tradiționale pentru controlul concurenței: mecanismele bazate pe blocarea cu *lock*-uri, mecanismele bazate pe marca de timp, și variante de control optimist al concurenței. Se discută apoi extinderea mecanismelor tradiționale pentru controlul concurenței - mecanismul *blocării altruiste* și *mechanismul validării instanțelor* - și se analizează mecanismele pentru controlul concurenței bazate pe informații semantice.

Capitolul 8. Controlul interblocării obiectelor

Capitolul 8 tratează controlul interblocării în sistemele distribuite unde, datorită dispersiei informațiilor pe mai multe noduri, interblocările sunt mai greu de prevenit, evitat și detectat. Se prezintă modelul general utilizat pentru tratarea interblocării, metodele de prevenire ale interblocărilor, care se bazează pe tranzații atomice și pe ordonarea tranzațiilor prin asocierea de mărci de timp și metodele de detecție a interblocării. Algoritmii de detecție se clasifică în trei mari grupuri: algoritmi centralizați, algoritmi ierarhici și algoritmi distribuți; ultimii, la rândul lor se subîmpart în următoarele clase de algoritmi: algoritmi care transmit caile grafului dependențelor tranzației, algoritmi care urmăresc arcele, algoritmi cu difuzie și algoritmi care stabilesc starea globală.

Pentru fiecare grup de algoritmi se prezintă cite un exemplu; pentru algoritmii complet distribuți se prezintă două exemple: un algoritm *edge-chasing* pentru detecția interblocării în modelul AND și un algoritm cu difuzie pentru detecția interblocării în modelul OR.

Capitolul 9. Proiectarea unui suport software pentru programare distribuită bazată pe obiecte și pe tranzații atomice (DOSTP)

Capitolul 9 cuprinde descrierea elementelor componente, a funcționării și a modului de proiectare al unui suport software pentru un sistem de programare distribuită. Sunt prezentate în primul subcapitol cerințele impuse: suportul are drept scop principal furnizarea unor primitive pentru lucrul cu tranzații atomice asupra unor obiecte dispersate în rețea, în condițiile asigurării recuperabilității obiectelor. Se arată (în contextul conceptelor teoretice prezentate în capitolele 6,7 și 8) că sistemul proiectat (DOSTP) utilizează ca model de arhitectură *modelul integrat* și varianta *dinamică a modelului obiect activ*. Schema de sincronizare a operațiilor la nivelul unui obiect

este o schemă optimistă, iar pentru recuperarea obiectelor se folosește protocolul *write-ahead logging* și tehnicile *Value-logging* și *Redo-only-logging*. Pentru interzicerea restartării tranzațiilor la infinit, s-a optat pentru implementarea unui mecanism *Wound-die*. După prezentarea arhitecturii generale și a modelului pentru sistemul DOSTP sunt descrise apoi componentele sistemului DOSTP: serverele de comunicații, managerii de tranzații, managerii de obiecte, serverele de nume, managerul central, procesele utilizator care realizează tranzațiile. Sistemul de recuperare al obiectelor nu este implementat în DOSTP ca proces distinct; el se compune din module plasate în managerii de obiecte și în managerii de tranzații și din fișierele lor *log* aferente.

Capitolul 10. Implementarea sistemului DOSTP în rețele Unix și Microsoft Windows (Win32)

Implementarea sistemului DOSTP s-a făcut, atât în varianta Unix cât și în varianta Windows urmărind îndeaproape proiectarea prezentată în capitolul 9. Există totuși diferențe între cele două variante, generate în special de faptul că în varianta Unix elementele componente sistemului au fost realizate ca procese distincte iar în Windows au fost realizate ca *thread*-uri. Capitolul 10 are ca centru de greutate proiectarea în manieră obiectuală a managerilor de obiecte. Scopul principal în realizarea managerilor de obiecte este ascunderea detaliilor de conectare cu sistemul tranzațional. Managerii de obiecte sunt realizați astfel încât ei comunică numai cu serverul de comunicație local, fără a avea cunoștință dacă serviciile le-au fost solicitate local sau de la distanță. În capitolul 10 sunt descrise, pentru un caz particular de obiect, mecanismele de control al concurenței și de recuperare ale obiectelor de acest tip, și, pe de altă parte, se arată cum trebuie să fie structurată o aplicație care să folosească într-o tranzație obiecte de acest tip. Un alt aspect important, prezentat în detaliu în capitolul 10, este circulația mesajelor în sistem, în cele șase situații distincte care pot să apară în funcționarea sistemului DOSTP: (1) conectarea unui nod la sistemul DOSTP; (2) integrarea sau eliminarea unui obiect din sistem; (3) inițierea unei noi tranzații; (4) implicarea unui obiect local și a unui obiect de la distanță într-o tranzație; (5) realizarea protocolului *two-phase commit*; (6) recuperare după o defecțiune a unui nod și a reintrării sale în sistem.

Capitolul 11. Concluzii și contribuții originale.

Capitolul 11 al lucrării prezintă contribuțiile originale pe care le aduce autorul în cadrul tezei și unele elemente referitoare la valoarea aplicativă și a direcțiilor de dezvoltare viitoare a cercetărilor efectuate.

Lucrarea conține patru anexe, în care sunt descrise aspecte de detaliu din capitolele 4,5, 9 și 10.

Capitolul 1.

SISTEME DE OPERARE DISTRIBUITE. SUPPORTUL OFERIT PRELUCRĂRILOR DISTRIBUITE, PROBLEME SPECIFICE, MODELE DE PROIECTARE

În acest capitol sunt prezentate aspecte legate de specificul sistemelor de operare distribuite în raport cu cele pentru rețelele de calculatoare, aspecte legate de interacțiunea sistem de operare distribuit - aplicație distribuită, probleme relative la gestiunea proceselor și fișierelor într-un sistem distribuit (problemele legate de comunicare și sincronizare într-un sistem distribuit sunt prezentate mai pe larg în capitolele următoare) iar în final se prezintă modelele uzuale de proiectare a unui sistem distribuit: *modelul proces și modelul obiect*.

1.1. Introducere

Un *sistem de operare* poate fi definit ca o colecție organizată de programe care îndeplinesc două mari funcții:

1. gestionează resursele hardware ale sistemelor de calcul, implementând algoritmi prin care se încearcă o optimizare a performanțelor; resursele gestionate sunt: procesoarele, memoria, dispozitivele de intrare/ieșire, fișierele, liniile de comunicație.
2. crează o interfață între utilizator și sistemul de calcul, extinzând setul de operații disponibile utilizatorului, minimizând efortul uman de programare și simplificând modul de lucru cu calculatorul prin automatizarea operațiilor uzuale în exploatarea echipamentelor.

Aceste funcții presupun alocarea și partajarea resurselor hardware între utilizatori, prevenirea interferenței între utilizatori, facilitarea executării operațiilor de intrare/ieșire, recuperarea erorilor, contabilizarea utilizării resurselor sistemului de calcul, facilitarea execuției de operații concurente în sistem, organizarea datelor într-o manieră optimă în vederea garantării securității lor și accelerării accesului la ele, gestiunea liniilor de comunicație în rețele de calculatoare, implementarea unui limbaj de comandă prin care sistemul interacționează cu utilizatorii.

Evoluția de peste 40 de ani a sistemelor de calcul s-a reflectat și în sistemele de operare în așa fel încât se poate spune că sistemele de operare au urmărit îndeaproape dezvoltarea echipamentelor de calcul. Rețelele locale de date devenit, în urma evoluției echipamentelor de calcul, cea mai avansată transpunere în practică, pe sisteme operaționale, a conceptelor prelucrării complet distribuite. În utilizarea rețelelor proliferază aplicații precum transferul fișierelor, poșta electronică, utilizări ale bazelor de date. Lucrul în rețea impune un nou model în proiectarea sistemelor de operare: *sistemul client-servant*. Clienții, distribuiți oriunde într-o rețea de calculatoare, utilizează servicii puse la dispoziție de servanți (componente software/hardware ale rețelei); servanții sunt în general dedicați unui anumit tip de serviciu specific, precum tipărire, accesul la fișiere, la baze de date, la informațiile de adresare în sistem. Sistemele de operare pentru sistemele distribuite au ca obiectiv major asigurarea unei interfețe unitare și omogene pentru utilizator, prin intermediul căreia acesta să poată avea acces la totalitatea resurselor sistemului, independent de gradul de eterogenitate al acestuia. Astfel, sistemele de operare încorporează totalitatea mecanismelor necesare pentru execuția distribuită a unei aplicații scrise într-un limbaj de programare de nivel înalt.

Se propune următoarea definiție pentru un *sistem de operare distribuit* ([Mus93d][DPV88]):

Definiție:

Se numește *sistem de operare distribuit*, un sistem de operare care este executat pe mai multe unități de prelucrare, independente, și care trebuie să asigure următoarele:

- (1) controlul alocării resurselor mediului distribuit, pentru a permite utilizarea lor într-un mod cât mai eficient; (2) perceperea mediului distribuit de către utilizatori ca pe un procesor virtual și nu ca pe o colecție de unități de prelucrare distincte, conectate printr-un sistem de comunicație; acest procesor virtual va servi ca suport utilizatorilor pentru dezvoltarea programelor în limbaje de nivel înalt;
- (3) comunicații sigure între unitățile de prelucrare din cadrul mediului distribuit;
- (4) ascunderea dispersării resurselor;
- (5) furnizarea unor mecanisme de protecție a resurselor împotriva acceselor neautorizate.

1.2. Prelucrarea distribuită a datelor

Anii 1990, cel puțin până în momentul de față, au marcat începutul perioadei calculului distribuit. Prelucrările se execută în paralel fie pe sisteme multiprocesor, fie în rețele de calculatoare. Rețelele de calculatoare pot fi configurate dinamic, pot fi interconectate. Interconectivitatea este facilitată prin introducerea conceptului de *sistem deschis*.

Un *sistem distribuit* este format dintr-un ansamblu de noduri legate între ele prin linii de comunicație; acestea reprezintă unica modalitate de comunicație între noduri. Fiecare nod este compus din unul sau mai multe procesoare cu acces la o memorie comună.

Sistemele distribuite sunt realizate ca un ansamblu unitar hardware-software care prezintă atât o *distribuire fizică cit și logică*. Enslow [DPV88] sugerează un spațiu tridimensional de distribuire pentru această clasă de sisteme: distribuirea la nivelul hardware-ului, distribuirea controlului și distribuirea bazelor de date. Conform definiției adoptată de el, un sistem cu prelucrare distribuită a datelor trebuie să posede cinci proprietăți:

1. *Dispersarea geografică* a resurselor fizice și logice aflate în legătură permanentă.
2. *Descentralizarea controlului și a informației de stare*; este o proprietate care conduce la deosebiri importante între sistemele de operare ale sistemelor centralizate și cele distribuite. Abordările de până acum în ceea ce privește sistemele de operare pentru rețelele locale de calculatoare au urmat două mari direcții: preluarea unui sistem de operare multiprogramat, tradițional și deci cunoscut, cărui i se adaugă capacități de comunicație în rețea sau dezvoltarea, de la zero, a unui nou sistem de operare.
3. *Multitudinea de resurse de prelucrare fizice și logice* utilizabile concomitent și care pot fi alocate dinamic unor sarcini specifice. Aceste resurse pot fi procesoare de diverse tipuri și fabricații conectate între ele, în diverse configurații și arhitecturi: de asemenea se pot utiliza (static sau dinamic) echipamente periferice diverse, partajate sau alocate corespunzător. Corespunzător multitudinii de resurse fizice, există în sistemele distribuite o varietate de resurse logice (software de comunicație, sisteme de gestiune a fișierelor sau bazelor de date, software suport pentru dezvoltarea de programe) care conlucrează la rezolvarea sarcinilor specifice ale aplicațiilor în context distribuit.
4. *Autonomia cooperativă a procesoarelor și a resurselor logice*; care se exprimă prin posibilitatea unei resurse de a rezolva în mod independent o cerere de serviciu locală și de a rezolva prin cooperare cu celelalte resurse o cerere de serviciu globală. Autonomia procesoarelor este o caracteristică fundamentală care decide, în ultimă instanță apartenența unui sistem la clasa sistemelor distribuite. Componentele de prelucrare trebuie să fie autonome. Ele cooperează cu alte componente, dar nu le controlează. Eventualele situații de concurență se rezolvă prin algoritmi specializați care se execută de către fiecare componentă de prelucrare. Necesitatea sincronizării activităților în sistemul distribuit, pentru evitarea conflictelor, constituie o altă diferență importantă față de sistemele clasice cu control centralizat.
5. *Transparența la nivelul utilizatorilor* care cer servicii prin nume. Fără a se preocupa de resursele care participă la realizarea serviciului cerut. Software-ul sistemelor distribuite identifică fizic o resursă necesară unui anumit serviciu (acțiune) indiferent de locația geografică a acesteia.

1.2.1. Obiectivele sistemelor distribuite

Apariția și dezvoltarea sistemelor distribuite a răspuns unor necesități ivite în domeniul prelucrării automate a datelor. Avântul din ultimii ani a unor astfel de sisteme a fost posibil datorită progreselor rapide înregistrate în domeniul echipamentelor de calcul și a celor de telecomunicație. Avantajele principale ale sistemelor distribuite în raport cu sistemele centralizate sunt următoarele:

1. *Autonomia locală combinată cu posibilitatea de comunicare*: O caracteristică generală în exploatarea sistemelor de calcul s-a constatat a fi [BW86] tendința utilizatorilor de a ezita în încredințarea informațiilor proprii sistemelor mari, centralizate, care deseori, pot compromite securitatea datelor; pe de altă parte, aplicații de largă utilizare, ca de exemplu cele financiare, este de dorit să se deslășoare pe grupuri de utilizatori. O altă problemă dificil de soluționat în folosirea unui sistem centralizat s-a dovedit a fi împartirea costurilor legate de acesta. Autonomia locală a sistemelor distribuite elimină aceste impedimente și, pe de altă parte, cu ajutorul facilităților de comunicație se realizează transmiterea de informații la distanță. Un alt factor subiectiv, care a impulsionat descentralizarea prelucrării, îl constituie necesitatea punerii în valoare a unor caracteristici umane specifice activității creatoare: capacitatea de acțiune independentă, sentimentul competiției în dialog direct cu calculatorul. Acestea sunt puse mai bine în valoare în grupuri relativ mici decât în cele foarte mari, impuse de caracterul tehnico-organizatoric complex al marilor sisteme centralizate.
2. *Utilizarea în comun a unor resurse*: Aceste resurse pot fi procesoare cu anumite caracteristici specifice, echipamente periferice, fișiere și baze de date. Un proces poate migra din procesorul local pe un alt nod dacă nodul local nu posedă resursele cerute, dacă prezintă defectuni în funcționare, sau dacă pe un nod situat la distanță se obține un timp de răspuns mai bun. Gestiunea ansamblului de resurse utilizate în comun se face în scopul optimizării folosirii lor, optimizării timpului de răspuns, controlului încărcării rețelei.
3. *Șiguranta în funcționare*: Existența unui număr mare de procesoare face posibilă continuarea activității unui sistem distribuit și în cazul în care un procesor își închează funcționarea corectă. Această facilitate este prezentă și la nivelul software: în acest scop, fișierele și bazele de date pot fi replicate, pachetele de programe software pot fi accesibile în mai multe noduri, calcule executate parțial pot fi continuate pe alt nod.

4. *Performante superioare.* Conexiuni multiple între diversele procesoare permit execuția concurentă a unor porțiuni din programe (procese) ceea ce, evident, va permite scurtarea duratei de execuție totale a programului. Localizarea datelor și procesoarelor disponibile se face după principiul *minimei apropieri* față de nodul local.

5. *Extensibilitate și adaptabilitate a configurațiilor* (atit hardware cit și software). Un sistem distribuit poate fi mult mai ușor modificat decit un sistem centralizat, prin adăugarea sau îndepărtarea unor componente, datorită modularității sale, atit în ceea ce privește structura hardware cit și cea software. Modificările vizează mărirea performanțelor sistemului.

În aceste condiții, este motivată tendința actuală ca sistemele distribuite să înlocuiască, treptat sistemele centralizate. Prelucrarea distribuită își propune să pună în valoare în modul cel mai eficient, avantajele descentralizării prelucrării, astfel încit să răspundă necesităților utilizatorului și aplicațiilor sale, fără pierderea controlului asupra standardelor, calității datelor, procedurilor de colectare, manipulare și prelucrare a unor volume mari de date, avantaje specifice sistemelor centralizate. În același timp, prelucrarea distribuită reprezintă un mijloc de a satisface dorința utilizatorului de apropiere față de sistemul de calcul, evitind dezavantajele unei descentralizări excesive.

În formularea criteriilor pe baza cărora se definesc actual sistemele cu prelucrare distribuită s-au avut în vedere de la început o serie de obiective pe care acestea trebuie să le atingă. Acestea sunt:

- (1) siguranță în funcționare sporită;
- (2) disponibilitate mărită în cazul defectării unor componente;
- (3) grad înalt de modularitate;
- (4) performanțe de sistem mai ridicate;
- (5) timp de răspuns bun în cazul suprîncărcărilor temporare;
- (6) utilizarea în comun a resurselor.

În prezent, pentru ca toate aceste facilități să poată fi disponibile, proiectanții de sisteme distribuite sunt confrunțați cu probleme vizînd următoarele direcții: (1) comunicația dintre procesoare și sincronizarea acestora; (2) gestiunea resurselor din sistem; (3) distribuirea resurselor; (4) distribuirea încărcării sistemului, care include rezolvarea problemei planificării globale a proceselor pe nodurile sistemului; (5) gestiunea proceselor; (6) gestiunea operațiilor de intrare-ieșire; (7) gestiunea memoriei; (8) controlul sistemului distribuit astfel încit să se conserve, caracteristica sa cea mai importantă: fiabilitate ridicată și disponibilitate în cazul defectării unor componente; (9) securitatea informațiilor; (10) mascarea eterogenității sistemului (aparută datorită eterogenității procesoarelor din nodurile rețelei și diferitelor convenții de reprezentare a datelor); (11) existența unor limbaje de programare care să includă posibilitatea de distribuire a proceselor și realizarea unei autonomii cooperative transparente la nivelul utilizatorului.

1.2.2. Modele actuale în prelucrarea distribuită

Termenul de *prelucrare distribuită* semnifică execuția unui program sub forma unor procese concurente, în mai multe noduri.

În prezent, în prelucrările distribuite sunt folosite următoarele modele: *modelul ierarhic*, *modelul federativ*, *modelul partajării resurselor*, *modelul omogen* și *modelul fluxului de date*:

1. *Modelul ierarhic* este specific sistemelor distribuite în timp real. El presupune un nivel inferior în cadrul căruia se colectează datele de intrare și se controlează dispozitivele fizice care achiziționează datele; nivelul imediat următor realizează gestiunea datelor și analizează informațiile de stare în conformitate cu anumite cerințe, putînd iniția chiar o reconfigurare a nivelului inferior în cazul detectării unor valori diferite de cele normale; nivelul superior este destinat controlului interactiv de către un operator uman.

2. *Modelul federativ* se compune dintr-un nivel centralizat de macroprocese și mai multe nivele de miniprocese distribuite; poate fi considerat ca un tip special de relație *master-slave*. Acest model încearcă să execute în nodul local cit mai multe procese: la depășirea capacității nodului, procesele aditionale sunt executate la distanță, în alte noduri iar la terminarea prelucrării, rezultatele sunt redirectate către nodul inițiator al prelucrării. Caracteristica principală a acestui model este coexistența mai multor noduri centrale fiecare maximizînd la nivelul său prelucrările de tip *master-slave*. Modelul poate plasa procese la distanță și la cerere. O extensie a modelului este cea în care controlul execuției se face centralizat, într-un nod, în timp ce procesele *slave* sunt distribuite astfel încit, prin execuția lor concurentă, să se obțină o prelucrare rapidă.

3. *Modelul partajării resurselor* pleacă de la ipoteza că toate resursele (procesoare, periferice) sunt date, în mod singular, ca disponibile pentru utilizatori care, după alocare, le utilizează conform cerințelor lor. Adicional, se introduce ca resurse partajabile anumite funcții software generice, prevăzute cu interfețe standard. În acest model, utilizatorul execută prelucrările asupra resurselor printr-o secvență de apel a acestor funcții. Modelul necesită un sistem de operare global care să controleze partajarea diverselor elemente distribuite din sistem și conduce la noțiunile de obiect, gestiunea obiectelor și programare la nivel de obiecte.

4. *Modelul omogen* se aplică unui ansamblu de calculatoare aproximativ de aceeași caracteristici și resurse. În acest model se urmărește o distribuție cit mai mare a proceselor generate de un program. Distribuția proceselor se face pe baza unui algoritm care minimizează (sau maximizează) diferiți parametri precum: *gradul de paralelism*, *încărcarea egală a procesoarelor*, *timpul de răspuns*, *încărcare maximă pe un cluster*.

5. *Modelul fluxului de date* presupune prelucrarea datelor în etape (începînd de la dispozitivele de intrare; datele sunt procesate de dispozitive specifice pînă la obținerea rezultatului). Acest tip de prelucrare distribuită necesită limbaje specializate și chiar procesoare specializate.

613.451
366 B

1.2.3. Interacțiunea prelucrare distribuită - sistem de operare

O tendință generală actuală în folosirea rețelelor atât locale cit și generalizate este de a se cupla stații de lucru sau calculatoare de diverse fabricații, eventual cu sisteme de operare diferite. Utilizatorii unui astfel de mediu eterogen sunt confrunțați cu probleme referitoare la alocarea resurselor la diversele procese (în special procesoare de diferite caracteristici). Ei trebuie să creeze procese, să aloce procesoare la procese, să specifice modul de comunicație între procese. Aceste activități constituie *programarea de la nivelul comutării proceselor* (PMS=Processor-Memory Switch Level Programming) față de programarea *tradițională* plasată la nivelul programării procesorului (ISP=Instruction Set Processor Programming). Într-un mediu eterogen, procesoarele nu sunt singura resursă critică; în afară de procesoare specializate, mai trebuie alocate și alte resurse precum linii de comunicație, buferi de date (fig 1.2.3.1).

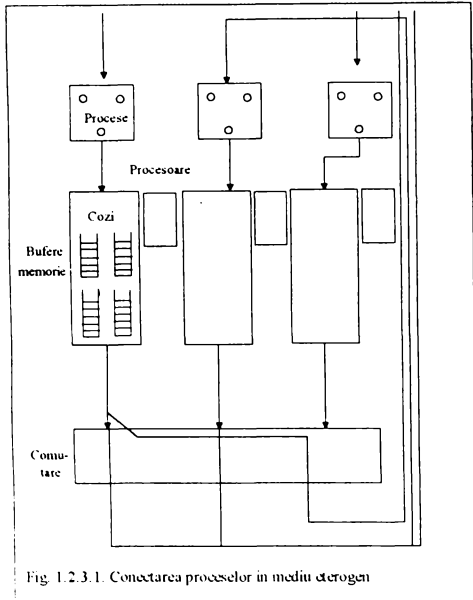


Fig 1.2.3.1. Conectarea proceselor in mediu eterogen

În mod curent, utilizatorii unui mediu eterogen dezvoltă programe în manieră tradițională: se scriu procese separate ca programe în diferite limbaje (Pascal,C) implementate pe anumite procesoare și apoi printr-un program separat se încarcă aceste programe pentru execuție. Descoari, aceste programe sunt scrise astfel încât înglobează în ele specificații privind disponerea lor pe procesoare sau chiar privind structura rețelei. Acest mod de lucru face dificilă sau chiar imposibilă utilizarea aplicației într-un mediu diferit, complică o dezvoltare ulterioară a aplicației care ar schimba structura programelor și suportul pentru siguranță la defecte (prin restartarea unei noi copii a programului pe un procesor diferit). O manieră nouă de abordare a acestor probleme este de a separa activitatea de programare a programelor de activitatea de cooperare a lor în codul aplicației. Vor exista astfel două nivele de programare în cadrul unei aplicații: *un nivel de programare al programelor componente ale unei aplicații (ISP)* și *un nivel de programare a aplicației care utilizează programele dezvoltate anterior (PMS)*. [BDW89]

Dezvoltarea suportului software pentru nivelul PMS este, calitativ, mult diferită față de nivelul IPS: necesită alte limbaje, metodologii, chiar medii de programare și, pe de altă parte, interacțiunea cu sistemul de operare distribuit.

Pot fi enumerate următoarele realizări notabile în această direcție:

1. *Limbajul Durra* este proiectat pentru a constitui suportul software al nivelului PMS [BDW89]. Limbajul oferă posibilitatea dezvoltării aplicațiilor prin rafinări succesive după modelul metodei Spiral [BDW89]. Utilizatorii modelului Spiral identifică selectiv componentele de risc ale aplicației, stabilesc cerințele lor, și apoi realizează fazele de proiectare, codificare și testare. Pe parcursul deslășurării acestui proces pot fi identificate și alte componente de risc cărora li se poate da o prioritate mai mare. Limbajul asigură construcții specifice pentru simularea unei aplicații, plecând de la o sumară decompoziție în taskuri specificate prin interfața lor cu celelalte taskuri și prin caracteristici ale execuției. După realizarea decompoziției, aplicația poate fi simulată utilizând un *master task* care va fi înlocuit mai târziu prin implementări specifice. Scopul acestei simulări inițiale este de a detecta acele taskuri a caror comportare afectează în cel mai mare grad performanțele întregului sistem. Utilizatorii pot experimenta diverse alternative de specificații comportamentale ale taskurilor până se obțin rezultatele impuse. În acest moment, se poate trece la înlocuirea descrierilor acestor taskuri generale cu descrieri mai detaliate cuprinzând alte taskuri interne și cozi, utilizând facilitățile de descriere ale taskurilor puse la dispoziție de Durra. Procesul de simulare continuă, luându-se în considerare noile taskuri și continuând decompoziția până la obținerea unei structuri interne satisfăcătoare și a unei comportări corecte a aplicației, pe ansamblu, fiind posibile și reluarea de tip *backtracking*. Codificarea taskurilor se realizează numai după ce s-au definitivat specificațiile lor de descriere, acestea s-au dovedit a fi acceptabile și nu mai este necesară continuarea decompoziției.

2. Lee și Goldwasser sunt autorii limbajului *DICON*, destinat controlului execuției concurente a programelor secventiale scrise în limbajele C și PROLOG. [LG85]

3. Magee și Kramer sunt autorii limbajului *CONIC* destinat reconfigurării dinamice a unui sistem în timp real. *CONIC* restricționează însă programarea taskurilor la o extensie a limbajului Pascal care include primitive pentru transferul mesajelor [MK83].

4. Belzile propune limbajul *RNET* pentru construcția aplicațiilor distribuite în timp-real. Un program *RNET* constă într-o specificație a configurației și un cod procedural care sunt compilate și legate cu un nucleu pentru execuție. Sunt permise facilități legate de execuția în timp-real precum specificarea unor timpi limită și a unor întirzieri folosite pentru monitorizarea și planificarea proceselor [Bel89].

1.2.3.1. Proiectarea unui suport pentru conectarea aplicațiilor în mediu distribuit

O aplicație distribuită în cadrul sistemului proiectat [Mus96][BDW89] TDL, constă într-un set de taskuri care comunică prin cozi de mesaje (fig. 1.2.3.2.1.). Descrierile taskurilor și declarațiile de tipuri specifică modul în care taskurile sunt lansate în execuție și executate ca procese concurente, tipul datelor prin care comunică procesele, precum și cozile de date necesare proceselor producătoare și consumatoare. Rezultatul unei compilări a unei descrieri de aplicație în limbajul TDL este un set de directive de alocare resurse și de planificare.

Există trei faze distincte în dezvoltarea unei aplicații distribuite utilizând limbajul: crearea unei biblioteci de taskuri, crearea unei aplicații utilizând biblioteca de taskuri și execuția aplicației (fig 1.2.3.2.2.)

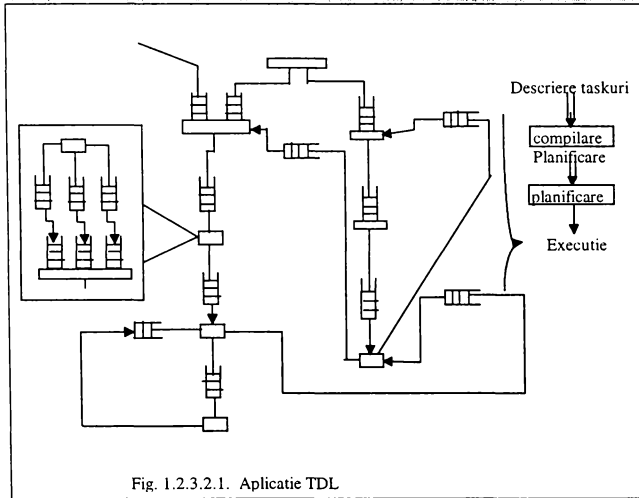


Fig. 1.2.3.2.1. Aplicatie TDL

În prima fază, folosind limbaje de programare adecvate, se scriu diversele taskuri care vor fi executate concurrent în mediul eterogen. Pentru fiecare task implementat se scrie și o descriere a taskului în limbajul TDL; pentru un anume task pot exista mai multe implementări, depinzând de limbajul de programare folosit, tipul procesorului, caracteristici de performanță și alte proprietăți. Fiecare implementare a unui task este însoțită de o descriere a taskului care este compilată și catalogată în biblioteca.

Descrierea unui task include specificatii referitoare la performanțele implementării, tipuri de date pe care taskul le produce sau consumă, porturile utilizate pentru comunicația cu alte taskuri și alte atribute statistice ale implementării.

A doua fază constă în crearea descrierii aplicației. Din punct de vedere sintactic, descrierea aplicației este o descriere a unui task și poate fi catalogată în biblioteca la fel ca un task obișnuit. Aceasta permite scrierea unor ierarhii de descrieri de taskuri. După analiza descrierii aplicației, analizorul generează un set de comenzi de alocare a resurselor și de planificare a execuției, comenzi care vor fi interpretate de planificator. Ultima fază constă în execuția aplicației. Planificatorul încarcă implementările taskurilor în nodurile rețelei și execută comenzile.

Descrierile de taskuri, de fapt entitățile sintactice utilizate în construcția aplicației, contin următoarele tipuri de informații: (1)porturile de interfață cu alte taskuri ; (2). atributele taskului. In figura 1.2.3.2.3. se prezintă structura generală a descrierii unui task.

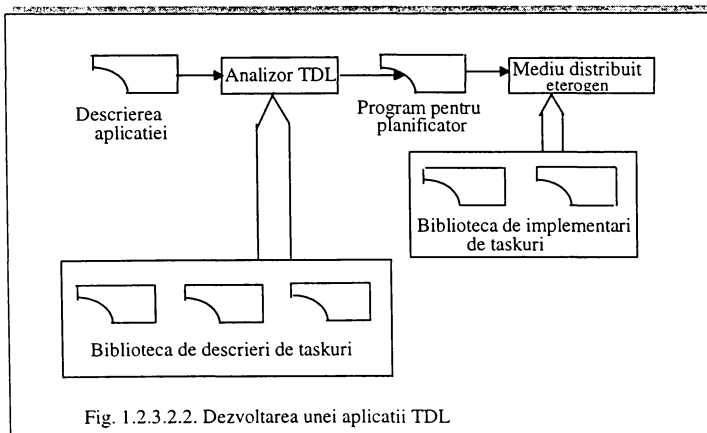


Fig. 1.2.3.2.2. Dezvoltarea unei aplicatii TDL

```

task nume_task
  ports          # pentru comunicatia dintre un proces [i o coada
                 declaratii de porturi
  attributes     # pentru a specifica proprietati ale taskului
                 atribute valoare
end nume_task
    
```

Fig.1.2.3.2.3. Structura unei descrieri de task

Partea de interfață cu alte taskuri conține descrierea porturilor de interfață specificând tipul datelor care se transferă în/din port și direcția, astfel:

```
ports
  in1:in string;
  out1,out2:out int;
```

Partea de atribute specifică diverse informații precum: autor, număr de versiune, limbaj de programare, nume fișier, tip procesor sub forma:

```
attribute_
  author="xxxx";
  implementation = "nume_program";
  queue-size=25;
```

Folosirea părții de atribute pentru a selecta un task din biblioteca de taskuri este ilustrată în fig. 1.2.3.2.4..

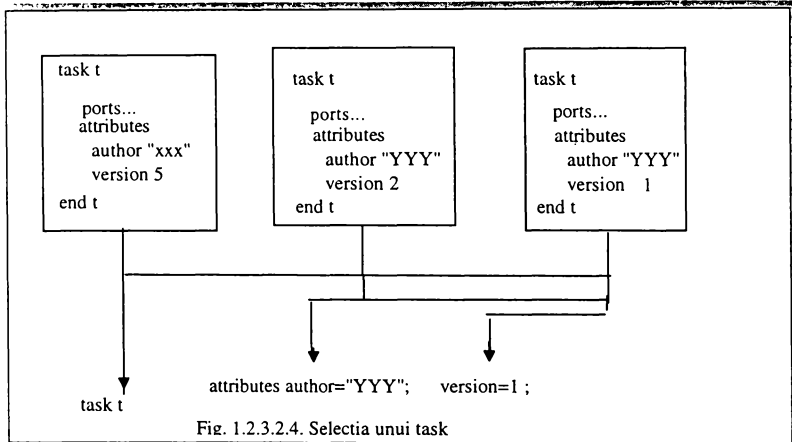


Fig. 1.2.3.2.4. Selectia unui task

Descrierea unei aplicații se face precizând taskurile componente și cozile prin care ele comunica. Există două tipuri de declarații într-o aplicație: declarațiile din secțiunea de procese, introdusă de cuvântul cheie *process* și declarațiile din secțiunea de cozi, introdusă de cuvântul cheie *queues*:

1. O declarație de forma:

```
nume-proces: task nume_task;_
```

crează un exemplar al taskului selectat, ca un proces nou.

2. O declarație de forma:

```
nume-coada[dim_coadă]: nume-port-1>transformari-date>nume-port-2
```

crează o coadă în care datele se acumulează de la un port de ieșire al unui proces (*nume-port-1*) către un port de intrare al unui alt proces (*nume-port-2*) și în care datele sunt transferate în conformitate cu *transformari-date*.

Există trei componente active în mediul de execuție : *taskurile de aplicații*, *serverul și planificatorul* . În figura 1.2.3.2.5.se prezintă relația dintre aceste componente.

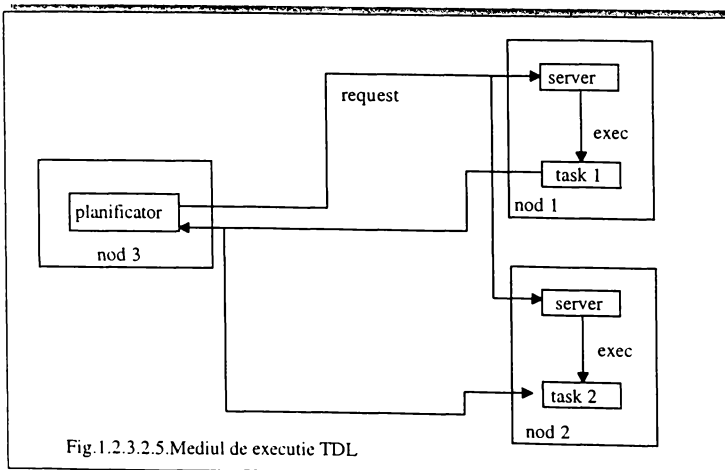


Fig.1.2.3.2.5.Mediul de execuție TDL

După analiza declarațiilor de tip, a descrierii taskurilor și a descrierii aplicației, se poate trece la execuția aplicației prin lansarea în fiecare nod a unui proces server și într-un nod separat al planificatorului. Planificatorul, care este apclat primind drept argument numele fișierului rezultat în urma analizei, are rolul de a starta taskurile, de a stabili liniile de comunicație și de a monitoriza execuția aplicației. În plus, implementează taskuri predefinite precum: *broadcast* și *merge*. În urma fazei de analiză se obțin instrucțiuni pentru planificator care descriu modul în care se execută procesele concurente, porturile și cozile utilizate pentru comunicație între procese tipurile de date ce se transferă între procese și posibilele reconfigurări ale aplicației. Serverul are rolul de a lansa taskul în nodul respectiv, la indicația planificatorului.

În orice nod al rețelei în care se execută un task, se activează și un exemplar de server.

În momentul lansării unui task, planificatorul îi furnizează, via server, informații precum: adresa nodului pe care se execută planificatorul, socket-ul Unix ce va fi folosit pentru comunicație task-planificator, și un număr întreg, identificator al taskului. Taskurile de aplicație utilizează interfața pentru a comunica între ele.

Taskurile de aplicație, în mod tipic, realizează următoarele operații: (1). Stabilesc comunicația cu planificatorul; (2). Cer identificatorii porturilor; (3). Trimit și recepționează date; (4). Termină comunicația cu planificatorul.

Implementarea actuală (prima versiune) a fost realizată folosind protocolul TCP/IP în rețele Unix; se preconizează ca versiunea a 2-a să se extindă la rețele mixte Unix și MS Windows (Win 32).

Se prezintă, ca exemplu, în figurile 1.2.3.2.6. și 1.2.3.2.7. un program pentru următoarea aplicație simplă: taskul A, un task producător generează un sir de date de tip *string* către un task de tip predefinit *broadcasting* care trimite mai departe aceste siruri către alte două taskuri consumatoare identice.

```
task taskA;
  attributes
    processor="xantipa";
    implementation="string_producer";
  ports
    out1:out string;
end taskA;

task taskB;
  attributes
    processor="alpha";
    implementation="string_consumer";
  ports
    in1:in string;
end taskB;
```

Fig.1.2.3.2.6..Descrierea taskurilor și a tipurilor de date

```
task main
  structure
    process
      p1:task taskA;
      p2:task taskB;
      p3:task taskB;
      pb:task broadcast
        ports
          in1:string;
          out1,out2:out string;
        end broadcast;
    queues
      qb1:p1.out1> >pb.in1;
      qb2:pb.out2> >p2.in1;
      qb3:pb.out2> >p3.in1;
  end main
```

Fig.1.2.3.2.7..Descrierea aplicației

1.2.3.2. OSF DCE - Strategie și arhitectura

Mediul de programare distribuită DCE (Distributed Computer Environment) [Tan96][BPK91][BGH93], creat la inițiativa organizației OSF (Open Software Foundation) este o colecție de componente software integrate în sistemul de operare UNIX; scopul DCE este de a furniza un suport pentru dezvoltarea și execuția aplicațiilor distribuite în medii distribuite omogene. OSF DCE este inclus în mediul de prelucrare *deschisă* propus de OSF, alături de sistemul de operare Unix OSF/1, interfața grafică OSF/Motif și OSF DME (Distributed Management Environment). Din punctul de vedere al sistemelor distribuite, DEC și DME sunt de cea mai mare importanță. În timp ce DCE constituie baza dezvoltării aplicațiilor distribuite și oferă servicii direct utilizatorului, DME abordează problemele gestionării rețelei și sistemului, alături de suportul pentru protocolul SNMP. DCE tinde să devină astăzi un standard industrial pentru prelucrarea distribuită.

Arhitectura generală a mediului DCE este reprezentată în fig. 1.2.3.3.1. Toate componentele se bazează pe

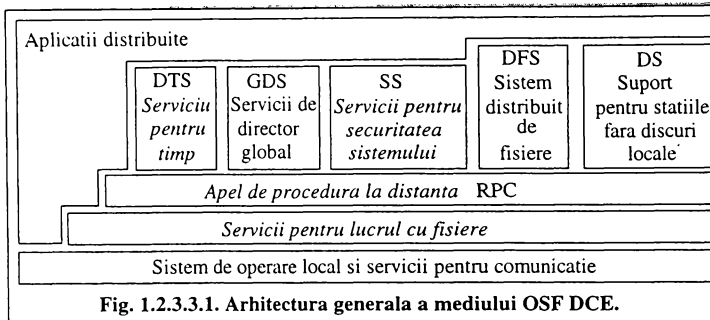


Fig. 1.2.3.3.1. Arhitectura generală a mediului OSF DCE.

serviciile sistemului de operare Unix și pe serviciile de comunicație TCP/IP. Aplicațiile distribuite dispun de serviciile fundamentale DCE via interfața de programare a limbajului C; celelalte servicii DCE sunt furnizate implicit prin intermediul serviciilor fundamentale sau prin intermediul serviciilor modificate ale sistemului de operare. Serviciile fundamentale DCE sunt:

(1) *Serviciile pentru lucrul cu fire*, care reprezintă o implementare portabilă a proceselor usoare (fire) în conformitate cu standardul POSIX 1003.4a. Firele fac posibilă prelucrarea concurentă într-un spațiu de adrese partajat (vezi subcapitolul 1.3.2.1.) și sunt utilizate în special de RPC pentru implementarea serverelor multifir și pentru implementarea apelurilor asincrone de proceduri, fără blocare;

(2) *DCE RPC* este suportul pentru comunicațiile între sisteme eterogene. O cerere pentru serviciu a unui client se transmite printr-un apel de procedură obișnuit, confirmat prin primirea rezultatelor de la server. Conversia de parametri și transmiterea efectivă de mesaje se realizează prin intermediul componentelor denumite *RPC stub*, aflate atât în mediul clientului cât și al serverului (DCE RPC este descris în detaliu în subcapitolul 3.4.9.A); acestea se generează automat plecând de la o scurtă descriere a procedurilor (asemănătoare prototipurilor funcțiilor din limbajul C). DCE oferă în acest scop limbajul IDL (Interface Definition Language). În plus, DCE mai oferă și posibilitatea mai multor tipuri de semantice a apelurilor, securitatea apelurilor de proceduri, cu autentificare și autorizare bazate pe serviciul de securitate DCE, posibilitatea utilizării unor nume globale de servere, pe baza standardului de servicii directorate X.500;

(3) *CDS (Cell Directory Service)* asigură gestiunea numelor în cadrul unui domeniu. Aceasta constă în maparea numelor în adrese și actualizarea informațiilor după nume; în plus, CDS stă la baza mapării adreselor de servere în cererile clienților RPC. CDS utilizează pentru creșterea eficienței replicarea și cache-ingul. O interfață de programare avansată este oferită prin X/Open Directory Service Interface;

(4) *Serviciul de securitate* implementează autentificarea, autorizarea și criptarea. Aceste mecanisme sunt puternic integrate în DCE RPC; de exemplu, clienții și serverele se pot autentifica reciproc, serverele pot verifica dinamic liste de control ale accesului pentru autorizarea clienților, iar mesajele RPC pot fi criptate la cerere;

(5) *DTS (Distributed Time Service)* implementează sincronizarea ceasurilor în mediul distribuit (o problemă comună tuturor mediilor distribuite, vezi subcapitolul 4.3.1.); se garantează că ceasurile locale ale nodurilor participante sunt sincronizate în interiorul unui interval și, eventual, cu o sursă externă exactă. Acest serviciu este impotrânt pentru implementarea algoritmilor distribuiți, care se bazează pe marca de timp; este, de asemenea, folosit și de alte componente DCE.

Alte servicii DCE sunt: (1) *GDS (Global Directory Service)* care extinde CDS-ul, furnizând facilitățile de nume în afara domeniilor; se bazează pe standardul X.500. Ca urmare, face posibilă cooperarea nu numai cu servere DCE, dar și cu alte servere X.500. Este prevăzută și alternativa Internet Domain Name Service pentru nume globale; (2) *DFS (Distributed File System)* implementează gestiunea fișierelor într-o manieră distribuită. Fișierele pot fi localizate pe diverse servere și pot fi replicate; programele client pot lucra cu fișierele independent de localizarea lor, ca și cu sistemul local de fișiere Unix. Accesul la fișiere este optimizat prin mecanismul de cache la nivel de fișier în nodurile clienților. Interoperabilitatea cu Network File System este asigurată printr-o interfață NFS/DFS; în plus, printr-o componentă specială (Diskless Support, DS) se asigură funcțiile de boot, swap și lucru cu fișiere și pentru stațiile fără disc local.

DCE furnizează structurarea sistemelor distribuite pe celule, cu scopul de a menține în limite rezonabile dimensiunea domeniilor administrative. O celulă cuprinde nodurile atașate unei rețele locale dar, în mod obișnuit, se definește în funcție de consideratiile de organizare și mai puțin în funcție de structura rețelei fizice; ca urmare, o celulă apare ca un set de noduri gestionate de o autoritate comună. Multe din serviciile DCE sunt optimizate special pentru accesul intra-celular; traversarea graniței celulei, către servere din afara celulei impune autentificare specială și

proceduri de autorizare diferite de cele pentru accesul intra-celula. Sistemul distribuit de fișiere oferă localizarea total transparentă pentru accesul intra-celula; la traversarea celulelor se cere însă specificarea explicită a numelor de celule.

În final se prezintă în fig. 1.2.3.3.2. un exemplu de aplicație distribuită în mediul OSF DCE. Comunicarea între procesele aplicației se prezintă în fig. 1.2.3.3.3.; aplicația constă în controlul producției și vânzării unei firme. Fiecare celulă conține un set de servere (DTS, DFS, CDS, SS); comunicația intra-celula și inter-celule se realizează via RPC.

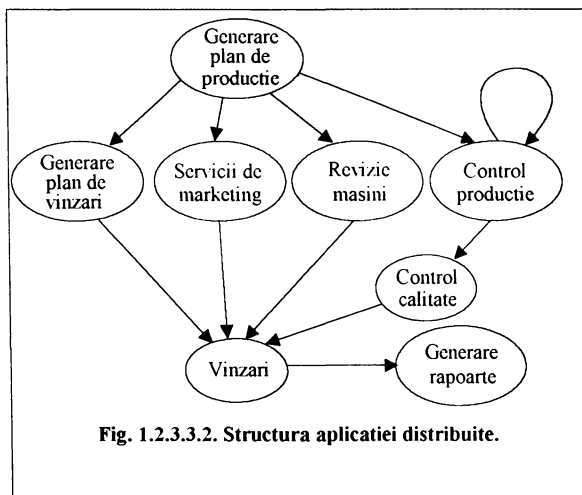


Fig. 1.2.3.3.2. Structura aplicației distribuite.

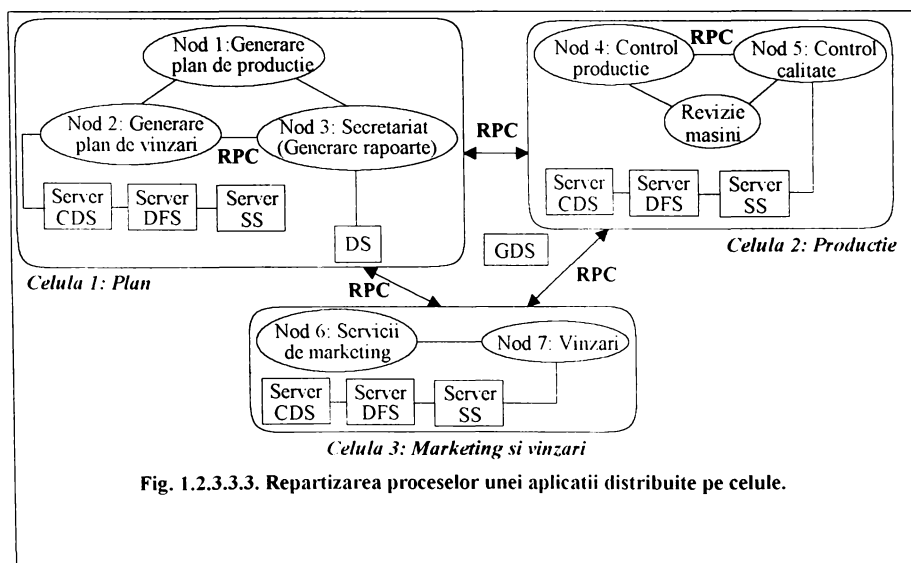


Fig. 1.2.3.3.3. Repartizarea proceselor unei aplicații distribuite pe celule.

1.3. Sisteme de operare pentru rețele de calculatoare și sisteme de operare distribuite

Într-o rețea de calculatoare o resursă ce trebuie gestionată este însăși rețeaua (comunicațiile în rețea). În arhitectura stratificată a rețelelor propusă de modelul de referință OSI, fiecare nivel dintr-un nod al rețelei va avea de implementat protocolul sau specific. Aceste protocoale apar astfel ca o prelungire a sistemului de operare în fiecare nod al rețelei. Pe de altă parte, unul din scopurile proiectării rețelelor este acela de a oferi utilizatorilor și servicii situate la distanță, localizate în rețea și nu în nodul de proveniență a cererii. Rămâne în sarcina sistemului de operare să identifice astfel de cereri, să le localizeze și să detecteze accesul neautorizat la serviciile oferite restrictiv în rețea (implementând astfel un mecanism de validare a cererilor). Deoarece într-o rețea de calculatoare pot fi legate, în general, calculatoare de diverse tipuri, provenind de la firme producătoare diferite, este posibil ca fiecare nod să fie dotat cu propriul sau sistem de operare. Aceste sisteme de operare locale vor gestiona resursele locale, ale nodului respectiv, fără a avea cunoștințe despre rețeaua în întregime. Pentru a beneficia de serviciile rețelei există în această situație două soluții: o prima soluție este de a adăuga fiecărui sistem de operare local un proces nou care va deservei cererile pentru acces la serviciile din rețea; altă soluție este de a proiecta un alt sistem de operare, care va înlocui sistemele de operare din toate nodurile. Procesul cu rol de interfață adăugat sistemului de operare în fiecare nod al rețelei se numește *proces agent*. Totalitatea proceselor agent din rețea formează *sistemul de operare al rețelei de calculatoare (NOS-network operating system)*. Folosind prima soluție, se extind astfel sistemele de operare existente; folosind a doua soluție, se obține un sistem de operare unic, care optimizează gestiunea tuturor resurselor din rețea (*DOS-distributed operating system*). Cele două soluții, NOS și DOS au drept scop furnizarea de servicii aplicațiilor, servicii legate de gestiunea proceselor, gestiunea comunicațiilor, gestiunea perifericelor, gestiunea memoriei și gestiunea fișierelor. Diferența majoră dintre cele două soluții este în maniera de abordare a gestiunii acestor resurse. Modelul NOS tratează resursele fiecărui nod ca fiind locale aceluiași nod și accesibile prin cereri explicite la intervenția proceselor agent. Modelul DOS tratează toate resursele distribuite în rețea într-o manieră globală, controlul și gestiunea tuturor acestor resurse realizându-se folosind strategii unice în toată rețeaua. De exemplu, dacă un utilizator inițiază execuția unui proces într-un calculator situat la distanță, în modelul NOS el trebuie să se conecteze la rețea, printr-o cerere explicită și să comunice nodului respectiv cererea sa. În sistemul de operare al calculatorului situat la distanță, procesul inițiat va fi perceput ca un proces local și va fi gestionat în consecință. În modelul DOS, inițierea unui proces este lăsată complet în sarcina sistemului de operare. Sistemul de operare examinează blocul de control al procesului pentru a determina eventualele cerințe specifice, după care planificatorul de procese, în funcție de încărcarea nodurilor ia în considerare procesul creat. Planificarea proceselor se face cu scopul de a optimiza gestiunea proceselor din întreaga rețea. Aceasta implică o transparență totală față de utilizatori care percep rețeaua ca un tot unitar și nu ca o colecție de dispozitive autonome care cooperează. Tipic, un sistem de operare distribuit este compus dintr-un nucleu, ale cărui copii sunt dispuse în toate nodurile și dintr-o ierarhie de componente software pentru gestiunea resurselor, situate deasupra acestui nucleu. Aceste componente pot fi organizate urmând modelele calculului distribuit (deoarece sistemul de operare este un exemplu de prelucrare distribuită): singular, federativ, replicat sau factorizat. În fig.1.3.1. este prezentată comparativ organizarea celor două tipuri de sisteme de operare pentru rețea: NOS și DOS [Mus93d][Cri89][Tan96].

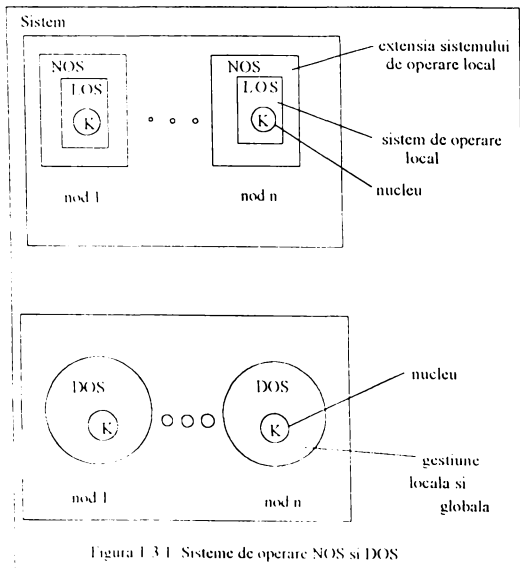


Figura 1.3.1. Sisteme de operare NOS și DOS

Diferența dintre sistemele de operare distribuite și cele pentru rețele de calculatoare poate fi reliefată sugestiv dacă se considera nivelul pe care îl ocupa comunicația între procese (IPC) în arhitectura sistemului de operare (fig. 1.3.2.): cu cât IPC este situat la nivele mai inferioare, cu atât crește transparența localizării resurselor în rețea. Astfel, în fig. 1.3.2.c, dacă IPC este situat chiar în nivelul de gestionare al proceselor se obține maximum de transparență și sunt satisfăcute astfel cerințele unui sistem de operare distribuit.

Spre deosebire de sistemele centralizate, controlul în sistemele de operare distribuite se realizează pe baza de negociere, consens și compromis între diversele componente răspindite în rețea. Într-un sistem de operare centralizat deciziile care se iau în gestiunea resurselor, se bazează pe o cunoaștere exhaustivă, în orice moment stării curente a sistemului; în contrast, sistemul de operare distribuit își bazează deciziile pe o apreciere, niciodată exhaustivă a stării rețelei, pe informații statistice și procentuale și folosește în acest scop metode și tehnici ca metoda *divide et impera*, teoria decizilor, teoria grafurilor, teoria jocurilor, algoritmi de control al concurenței și bazele de date.

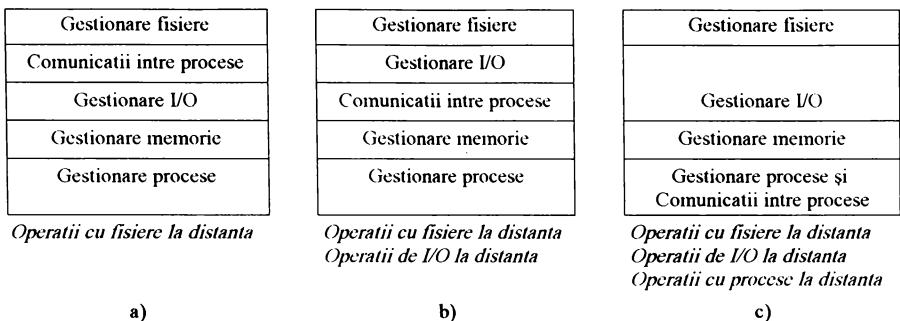


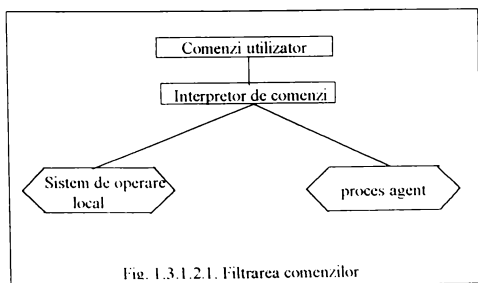
Fig.1.3.2. Nivele de transparenta obtinute în funcție de plasarea nivelului de comunicare între procese

1.3.1. Specificul sistemelor de operare pentru rețele de calculatoare

O caracteristică importantă a rețelelor de calculatoare, în special a rețelelor generalizate este posibilitatea conectării în rețea a diverse tipuri de calculatoare. Apare astfel necesară proiectarea sistemului de operare pentru rețeaua de calculatoare, sistem care să permită însă execuția sistemelor de operare locale, în nodurile rețelei. El trebuie să contină limbajul de comandă prin care utilizatorii să-și prezinte cererile în rețea și procesele agent prin care aceste cereri sunt acceptate, transmise și recunoscute în rețea.

1.3.1.1. Limbajul de comandă

Instrucțiunile limbajului de comandă sunt cereri de servicii adresate sistemului de operare de către utilizatori. Ele sunt tratate de interpretorul de comenzi care se constituie fie ca nivelul superior al sistemului de operare fie ca un proces separat, distinct față de sistemul de operare. Limbajul de comandă al sistemului de operare pentru rețea va cuprinde setul de cereri de servicii posibile în rețea prin care utilizatorii au acces la facilitățile oferite de rețea. Comenzile trebuie să fie, pe de o parte, distincte față de comenzile sistemului de operare local, dar pe de altă parte trebuie să refere accesul la distanță într-o manieră cât mai transparentă pentru utilizator. De exemplu, un transfer de fișiere între două noduri poate fi văzut ca un caz particular al unei comenzi de tip *copy*: distincția între accesul local și accesul la distanță se face prin argumentele comenzii. Dacă nivelul de transparentă este mai scăzut, într-o comandă de tipul "*copy sursă destinație*", a cărei semnificație este copiază fișierul "sursă" situat la distanță în fișierul local "destinație", se poate specifica explicit localizarea fișierului "sursă" prin nume de cale, de exemplu "nod/nume-utilizator/sursă". Dacă, dimpotrivă, nivelul de transparentă este mare, în timpul "sursă" se va specifica numai numele fișierului, fără localizare. Rămine în sarcina sistemului de operare să găsească adresa nodului care conține fișierul ceea ce implică pastrarea unor structuri de date în care să se memoreze corespondența dintre numele locale ale fișierelor și localizarea lor. Sistemul de operare al rețelei poate recunoaște și apeluri sistem adresate rețelei în timpul execuției unui program. De exemplu apelul sistem: *fork(nod1/nume-utilizator/nume-fișier.nod2)* poate fi interpretat ca având următoarea semnificație: cu fișierul /nume-utilizator/nume-fișier conținând forma executabilă a unui program din nodul "nod1", se crează și execută un proces în nodul "nod2". Responsabilitatea sistemului de operare al rețelei este să capteze astfel de cereri care implică accesul la distanță, să le interpreteze și să le transmită mai departe nivelelor inferioare de protocol.



Utilizator 1	nod UKC	login IMC.OW	cont 179600	acces mail
Utilizator 2				

Fig 1.3.1.2.2. Structuri de date utilizate pentru conectarea la distanta

1.3.1.2. Proceele agent.

Într-o rețea de calculatoare cu sisteme de operare eterogene, fiecare nod se conectează la rețea prin intermediul unui proces agent. Comenzile cu efect local sunt separate și tratate de sistemul de operare local; cele ce implică acces la distanță sunt transmise procesului agent (fig. 1.3.1.2.1.), pentru a putea le transmite mai departe către nodul care le poate îndeplini. Procesul agent menține structuri de date care îi permit să facă legătura între cele trei elemente ale comenzii: utilizator, nod la distanță, serviciu. Pentru fiecare utilizator trebuie să existe o indicație relativă la nodurile la care se poate conecta: parole, conturi folosite la conectare, tipurile de servicii permise (fig.1.3.1.2.2).

O categorie importantă de servicii oferite utilizatorilor de către sistemul de operare al rețelei se referă la manipularea fișierelor. Sistemul de operare trebuie să furnizeze utilizatorilor facilitatea de a accesa fișiere la distanță, cu verificarea autenticității cererilor, într-o manieră mai mult sau mai puțin transparentă. Dacă nivelul de transparență este mic, utilizatorul va indica și denumirea simbolică a nodului la distanță implicat în transfer; sarcina procesului agent este doar de a determina adresa fizică a acestui nod, rămânând în sarcina sistemului de operare local nodului să determine fișierul implicat în transfer prin analiza numelui de calcul al fișierului furnizat de procesul agent. Dacă nivelul de transparență în localizarea fișierelor este mare, procesul agent menține un catalog al fișierelor accesibile - via proces agent - în rețea. Utilizatorii nu trebuie să cunoască localizarea acestor fișiere în rețea. Mai mult, dacă procesul agent menține pentru fiecare utilizator un catalog cu fișierele accesibile în rețea, el poate fi realizat astfel încât să ascundă numele fișierelor și alte convenții locale față de utilizatorii situați la distanță. Procesul agent va păstra o baza de date care memorează corespondența nume virtual al fișierului-nume nod, nume fișier. (fig. 1.3.1.2.3.).

Pentru implementarea procesului agent trebuie avute în vedere două probleme: mecanismul interacțiunii cu sistemul de operare local și localizarea procesului agent. Dacă sistemului de operare local îi este caracteristic interacțiunea între procese folosind apelul de procedură atunci și comunicarea cu procesul agent se va face tot prin apel de procedură. Dacă tehnica de interacțiune între procese, folosită de sistemul de operare local se bazează pe transmiterea de mesaje atunci și interacțiunea cu procesul agent se face cu ajutorul mesajelor.

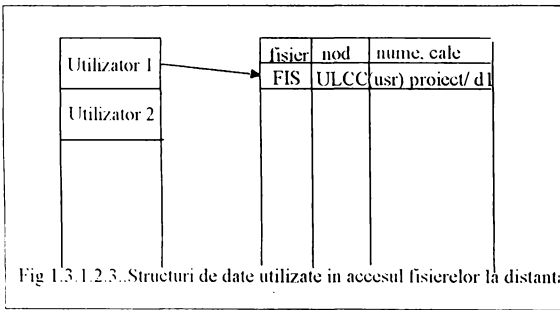


Fig 1.3.1.2.3. Structuri de date utilizate în accesul fișierelor la distanță

În ceea ce privește localizarea procesului agent există trei variante folosite: (1) procesul agent poate fi implementat ca un proces separat în fiecare nod al rețelei; (2) poate fi implementat ca o componentă a nucleului sistemului de operare local; (3) poate fi plasat într-un nod special;

În prima variantă, (fig.1.3.1.2.4.) toate comunicațiile dintre procesele utilizator și procesul agent se fac prin intermediul nucleului sistemului de operare local. Aceasta înseamnă că pentru transmiterea unui mesaj în rețea, au loc trei interacțiuni utilizator-nucleu: prima interacțiune este apelul sistem al utilizatorului către procesul agent, a doua este transmiterea mesajului la procesul agent de către nucleu, iar a treia este transmiterea de către agent a mesajului în rețea, transmitere care se face folosind protocoalele din nucleul sistemului de operare. Aceasta triplă interacțiune poate conduce la un timp de răspuns uneori prea mare. Situația poate fi și mai dezavantajoasă dacă procesul agent este tratat ca orice proces utilizator și evacuat din memoria internă. Totuși, această variantă are avantajul că nu necesită creșterea dimensiunii nucleului.

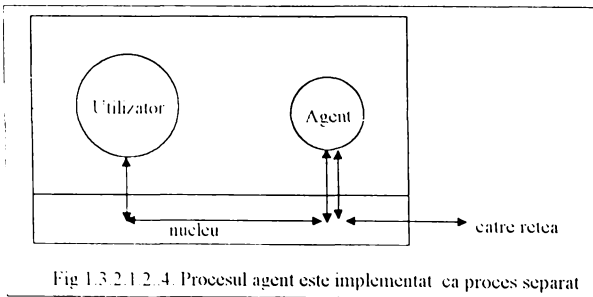


Fig 1.3.1.2.4. Procesul agent este implementat ca proces separat

Varianta a doua (fig. 1.3.1.2.5) evită traficul suplimentar de mesaje generat în prima variantă dar dimensiunea nucleului este marit considerabil; această creștere a dimensiunii poate fi uneori intolerabilă.

În varianta a treia (fig. 1.3.1.2.6) agentul este rezident într-un nod special.

El constă de fapt într-o colecție de module, cite unul pentru fiecare sistem de operare local al fiecărui nod. Din aceste noduri, accesul la rețea se face numai prin

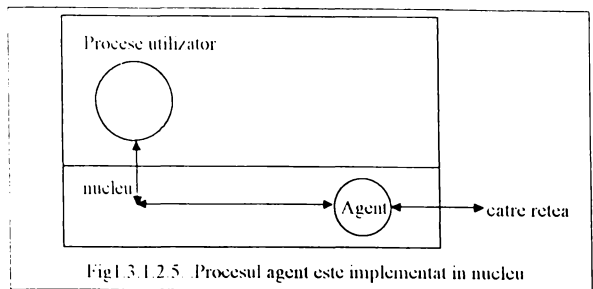


Fig 1.3.1.2.5. Procesul agent este implementat în nucleu

intermediul nodului agentului; de aceea sunt necesare rutine de protocol la nivelele inferioare ale rețelei care să asigure emisia și recepția mesajelor de la nodul agent. Metoda are dezavantajul că este vulnerabilă la defecțiunile nodului agent dar și avantajul că simplifică implementarea. Ea se aplică numai în situația în care agentul asigură un anumit tip de servicii în rețea (gestiunea fișierelor); în aceste situații pot exista mai mulți agenți, fiecare asigurând un anumit tip de servicii. Această variantă poate fi clasificată în modelul general *client-servant*.

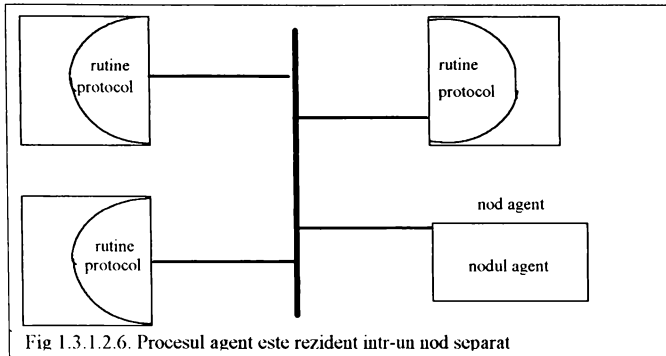


Fig 1.3.1.2.6. Procesul agent este rezident într-un nod separat

1.3.2 Specificul sistemelor de operare distribuită

În acest subcapitol se prezintă probleme legate de gestionarea procesoarelor și a fișierelor într-un sistem de operare distribuit. Problemele legate de comunicație, interacțiunea între procese și sincronizarea proceselor într-un sistem distribuit sunt prezentate în capitolele următoare, 2, 3 și 4.

1.3.2.1. Gestionarea proceselor

Într-un sistem de operare distribuit, funcția de gestionare a proceselor trebuie să furnizeze strategii și mecanisme pentru: (1) operații asupra proceselor, atât a proceselor locale, cât și a celor de la distanță: crearea, distrugerea, terminarea, suspendarea, modificarea atributelor, modificarea stării proceselor; (2) execuția concurentă a proceselor; (3) planificarea pentru execuție a proceselor; (4) lansarea proceselor în execuție pe procesoare la distanță; (5) migrarea proceselor; (6) încărcarea echilibrată a procesoarelor [Mus93d].

Fată de sistemele centralizate, în sistemele distribuite, funcția de gestionare a proceselor are un handicap în plus: ea trebuie realizată în condițiile unor informații de stare incomplete, deoarece, în astfel de sisteme, mesajele, utilizate pentru comunicarea între procese a informațiilor de stare sunt întârziate (uneori chiar pierdute); starea transmisă prin mesaje este întotdeauna starea dintr-un moment de timp anterior și nu din momentul curent. În plus, în sistemele distribuite, pentru a se îmbunătăți timpul de răspuns, încărcarea procesoarelor și partajarea resurselor, sunt necesare noi tipuri de operații cu procesele: lansarea lor în execuție pe un procesor de la distanță și migrarea proceselor.

Există două tipuri de unități de lucru pentru procesor furnizate de sistemele de operare distribuite ([Gos91] [Tan93] [Tan96]): *procese grele* și *procese usoare* sau *fire*. Diviziunea proceselor în cele două categorii se bazează pe gradul de partajare al memoriei. Astfel, *procesele grele*, cărora sistemul de operare le asociază blocuri de control sau descriptori de proces au spații de adresa proprii, atât pentru cod, cât și pentru date și stivă, dar și tabelă de descriptori fișiere, timere, semnale, semafoare proprii. *Procesele usoare* partajează datele și au numai stivă proprie (fig. 1.3.2.1.); ele partajează doar variabile locale, fișiere deschise, timere, semnale.

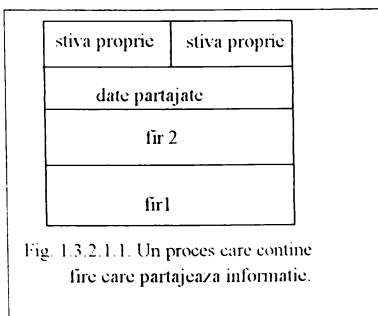


Fig. 1.3.2.1.1. Un proces care conține fire care partajează informație.

Firele au fost introduse pentru: (1) a se putea realiza execuția concurentă și în cazul apelurilor sistem cu blocare; (2) a se reduce regia de sistem datorată comutării contextului proceselor; (3) a se asigura tratarea concurentă a unor acțiuni ale utilizatorilor (de exemplu, în sistemul de gestiune al ferestrelor); (4) a se partaja informații în cadrul programelor utilizator. Deoarece firele partajează memoria, accesul la datele comune se face în regim de excludere reciproca. Pentru sincronizarea firelor (inclusiv în accesul la secțiunile critice) se folosesc în general semafoare binare (numite *mutex-uri*) și *variabile condiție* (asemănătoare variabilelor condiție folosite pentru sincronizarea cu monitorare); fiecare variabilă condiție, cînd este creată, se asociază unui *mutex*. Diferența dintre *mutex-uri* și variabilele condiție este aceea că *mutex-urile* se folosesc pentru blocarea pe termen scurt (pentru blocarea accesului la o secțiune critică), în timp ce variabilele condiție se utilizează pentru blocarea

pe termen lung (pentru a se aștepta disponibilitatea unor resurse). Evitarea interblocării firelor se face prin utilizarea combinată a mutex-urilor și a variabilelor condiție (fig. 1.3.2.1.2.) [Tan93]. [Tan96]). Variabila condiție se utilizează pentru a achiziționa resurse, deoarece la efectuarea unei operații *wait* asupra ei se eliberează automat și atomic mutex-ul asociat (1.3.2.1.2.a). Când mai târziu, firul care deține resursa, o eliberează, execută o operație *wake-up* (*variabila_condiție*), deblochează firul (sau toate firele) care așteaptă la condiție și mutex-ul este alocat din nou.

```

lock (mutex)
    *verifica structurile de date
    while (*resursa_ocupata)
        wait (variabila_conditie)
    - *marcheaza_resursa_ocupata
unlock (mutex)
(a)

lock (mutex)
    marcheaza_resursa_libera
unlock (mutex)
wakeupt (variabila_conditie)
(b)
    
```

Fig. 1.3.2.1.2. Utilizarea combinată a mutex-urilor și a variabilelor condiție

Se apreciază ([Gos91]) că puține sisteme de operare distribuite permit lucrul cu fire și deci partajarea memoriei între unități de lucru ale aceluiași proces. Astfel, sistemul Charlotte, creat la Universitatea Wisconsin-Madison, de către o echipă de cercetători condusă de Finkel ([Fin89]) este un exemplu de sistem care nu permite partajarea memoriei prin intermediul firelor; în schimb, sistemul V, creat la Universitatea Stanford, de către o echipă de cercetători condusă de Chertoff ([Che88]) reunește toate procesele cu același spațiu de adrese sub numele de *echipă de procese*; astfel, serverele pot utiliza o zonă de memorie comună, care implementează de obicei buffere. Sistemul Mach, creat la Universitatea Carnegie-Mellon de către o echipă de cercetători condusă de Rashid, permite lucrul cu task-uri și fire în cadrul task-urilor; în Mach este tipic un număr mic de taskuri și multe fire în cadrul unui task. Mach permite partajarea memoriei între task-uri, atât sub forma *read/write* cit și sub forma *copy-on-write*. Taskul este unitatea de alocare pentru resurse, atât în ceea ce privește spațiul de adrese paginat, cit și în ceea ce privește protecția accesului la resurse; drept rezultat, fiecare task își partajează memoria cu taskul părinte în conformitate cu un atribut de mostenire.

Sistemele de operare distribuite trebuie să permită efectuarea de operații asupra proceselor la distanță, într-un mod transparent pentru utilizator. Suportul pentru astfel de operații trebuie să includă primitive și facilități pentru crearea, distrugerea, suspendarea și continuarea unui proces, pentru alocarea/dealocarea de porturi pentru comunicație, ca și pentru lansarea și migrarea unui proces pe un alt procesor. Un model pentru crearea proceselor la distanță este oferit de sistemul V [Che88]. Datorită facilităților de comunicație dintre procese, un proces poate fi creat, prin intermediul apelului sistem, *CreataProces()* și la distanță; acesta construiește și transmite un mesaj nucleului nodului destinație care creează procesul, asociindu-i un identificator (pid) unic. Noul proces creat rămâne în starea *awaiting_reply* pînă la transmisia de către părinte a datelor inițiale care permit începerea execuției.

În ceea ce privește *crearea proceselor* există mai multe implementari posibile: (1) execuția părintelui se deslășoară concurrent cu cea a proceselor fii (sincronizarea execuției între aceștia se poate face cu construcții de tipul *fork/join*) sau procesul părinte așteaptă pînă se termină fii săi; (2)(a) procesele fii partajează numai un subset din datele părintelui sau toate variabilele în comun. Iată două exemple: (1) sistemul distribuit LOCUS, o extensie a sistemului Unix în mediu distribuit; în LOCUS, apelurile Unix *fork()* și *exec()* au fost extinse astfel încît: *fork()* poate crea un proces la distanță, cu același cod ca și apelantul, dar cu stivă și date separate iar *exec()* a fost extins pentru a permite unui proces să migreze pe alt nod. În plus au fost introduse două apeluri noi: *run()*, al cărui efect logic este combinația *fork()+exec()* și care permite inițierea execuției unui program fie local, fie la distanță, și *migrate()* care permite unui proces să-și modifice nodul în care se execută; (2) sistemul Sprite, creează un nou proces, cu apelul *Procfork()* care poate avea două comportamente, fie analog apelului *fork()* din Unix, fie poate crea noului proces să-și partajeze datele cu procesul părinte (fig.1.3.2.3).

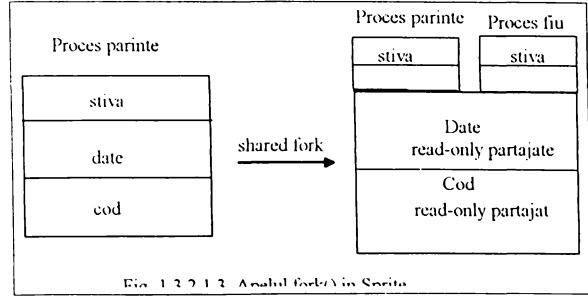


Fig. 1.3.2.3. Apelul fork() în Sprite

Un alt aspect legat de gestiunea proceselor este *crearea și distrugerea porturilor*, în sistemele în care comunicația între procese se face prin porturi. Un port poate fi creat la crearea proceselor sau la cererea explicită a proceselor. La crearea unui proces, procesul părinte poate asocia fiului un port de date și un port de servicii; acestea vor fi utilizate pentru comunicație

cu sistemul. O problemă care se pune însă este cum procesul începe comunicație cu alte procese: o soluție este ca procesul fiu să comunice numai cu sistemul și să ceară acestuia crearea de porturi când trebuie să comunice cu alte procese; altă soluție este ca parintele să-i ofere fiului un set de porturi pe care acesta să le poată folosi imediat. Un port poate exista fie pînă la distrugerea procesului proprietar, fie pînă la distrugerea procesului. În mod obisnuit, numai creatorul portului poate primi mesaje în acel port; drepturile pentru emisie pot fi transmise de către creator la orice proces cu care se dorește comunicarea.

În ceea ce privește **execuția la distanță a proceselor** există următoarele aspecte: (1) nivelele la care trebuie furnizate facilități pentru execuție la distanță: (a) la nivelul programelor, în care instrumentul principal este apelul de proceduri la distanță și, mai rar, particularități de partajare a memoriei; (b) la nivelul limbajului de comandă. (2) handicapul pe care îl prezintă execuția la distanță în mediile eterogene: sunt necesare traduceri datorită modificării arhitecturii hardware la trecerea pe un alt nod, dificultatea de a asigura un spațiu de nume globale, dificultatea de a descrie serviciile de executat, datorită diferențelor sintactice și semantice în liniile de comandă ale sistemelor diferite; (3) cerințele impuse unei execuții la distanță: (a) existența unui mecanism care să identifice nodurile slab încărcate; (b) posibilitatea ca utilizatorul de la un nod să poată stopa execuția proceselor altui utilizator, care determină degradarea performanțelor proceselor proprii (preemptibilitate); (c) execuția la distanță să poată fi realizată cu rezultate comparative cu execuția locală; (4) selecția unui nod (server) pentru execuția la distanță a unui proces se face fie pe baza propagării informațiilor despre un server disponibil sau a cererilor pentru servere a clienților, fie pe baza informațiilor furnizate de sistemul de gestiune al resurselor. În primul caz, se utilizează o mică bază de date, locală fiecărui nod care conține date relative la serverele în care nodul este interesat și date despre clienții cărora le poate oferi resurse. Numai două operații sunt disponibile asupra acestei baze de date: adăugare și ștergere; ambele sunt realizate în urma recepției unor mesaje broadcast transmise de un server care dorește să se facă disponibil. Pentru a selecta un server disponibil, un client folosește baza de date locală, eventual după actualizarea ei în urma transmiterii broadcast a cererii.

Migrarea proceselor într-un sistem distribuit este facilitatea care asigură relocarea dinamică a proceselor de pe un procesor pe alt procesor. Operația poate fi inițiată fără stirea procesului respectiv, și a proceselor cu care acesta interacționează. Migrarea proceselor este justificată prin următoarele cerințe: (1) repartizarea eficientă a proceselor pe procesoare; (2) creșterea gradului de paralelism ale operațiilor; (3) execuția proceselor cu cerințe speciale de resurse pe procesoarele care satisfac cel mai bine aceste cerințe; (4) execuția proceselor care lucrează cu fișiere de dimensiuni mari pe nodurile care contin acele fișiere, pentru a reduce astfel traficul de mesaje în sistem; (5) creșterea siguranței de funcționare a sistemului. Dintre problemele numeroase pe care le ridică migrarea proceselor, se enumera: (1) determinarea contextului unui proces și detașarea procesului de acesta, transferul procesului și conectarea sa într-un nou context (cod, date, stiva, registre, cozi de mesaje, semafoare); (2) tratarea mesajelor transmise procesului după migrarea sa; (3) asigurarea transparenței migrării proceselor; acest lucru este uneori dificil de obținut deoarece starea procesului poate fi distribuită în mai multe tabele ale sistemului de operare, și deci dificil de extras și apoi de creat în nodul destinație; (4) asigurarea unei strategii de migrare care definește *cînd* are loc migrarea unui proces și *unde* are loc, în scopul obținerii unor performanțe superioare; (5) determinarea momentului în care procesul care migrează poate să-și înceapă execuția pe nodul destinație; această problemă este generată de necesitatea copierii întregului spațiu de adresă al unui proces. O excepție de la această regulă o constituie sistemul Spice, creat la Universitatea Carnegie-Mellon (Zay87), care oferă posibilitatea ca un proces să-și înceapă execuția la destinație aproape imediat, utilizînd un transfer *copy-on-reference* (fig.1.3.21.4); (6) migrarea proceselor în medii eterogene necesită în afară de salvarea și restaurarea stării contextului proceselor și interpretarea consistentă a datelor. În asemenea situații se pun probleme de traducere a datelor; o soluție a acestor probleme este folosirea unei reprezentări standard pentru transportul datelor.

În ceea ce privește tratarea mesajelor primite de un proces care urmează să migreze există trei soluții: (1) **soluția redirectării mesajelor**. În momentul în care un proces este suspendat pentru migrare, se memorează toate mesajele primite și

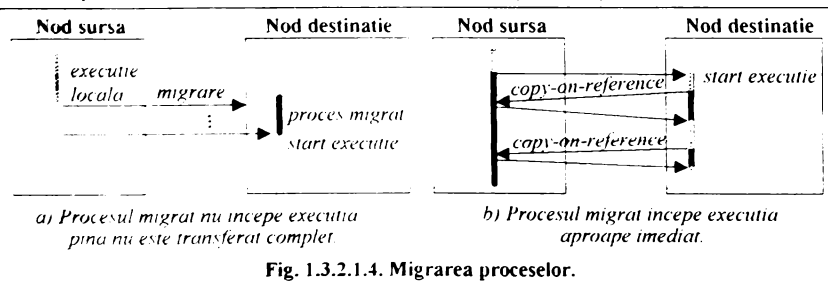


Fig. 1.3.21.4. Migrarea proceselor.

se trimite la noua destinație a procesului. după migrare; nu se înfăptuiesc procesele emițătoare despre noua adresă a procesului migrat, astfel că nodul sursa va trebui să retransmită toate mesajele la destinație (nodul în care procesul a migrat); (2) **soluția prevenirii pierderii mesajelor**. În acest caz se înfăptuiesc procesele care comunică cu procesul ce urmează a migra despre această intenție. Pe timpul suspendării procesului pentru migrare, mesajele se buferează; în urma terminării migrării procesului se transmit mesaje de tip *slăsit migrare* și următoarele se transmit direct la noua adresă; (3) **soluția recuperării mesajelor pierdute**. În acest caz, nu se salvează mesajele transmise în timpul migrării (ele sunt astfel, pierdute); totuși se memorează numele și adresele proceselor emițătoare, urmînd ca după migrare procesul respectiv să restabilească legătura cu aceste procese în vederea recuperării mesajelor.

Mecanismul de implementare al migrării poate fi împartit în două faze: *faza de colectare de informații statistice și faza de transfer a procesului*. Colectarea de informații statistice, în general relative la încărcarea procesorului (exprimată în numărul de procese gata de execuție), a liniilor de comunicație (exprimate în numărul de mesaje în așteptare) se poate face fie periodic, fie la apariția unor evenimente precum terminarea proceselor sau recepția mesajelor. Faza de transfer a procesului prezintă particularități legate de sistem. Se prezintă pentru exemplificare, etapele necesare pentru transferul unui proces în sistemul DEMOS/MP, creat la Universitatea din California, Berkeley. În DEMOS/MP transferul unui proces, care are loc după un proces de negociere, în urma căruia se stabilește nodul destinatar, se desfășoară astfel: (1) se marchează procesul în migrare, starea procesului se reține (cu excepția stării *ready*, când procesul este scos din lista *ready*) pentru că procesul își va menține starea în urma migrării; mesajele în tranzit în timpul migrării vor fi plasate în coada de mesaje a procesorului; (2) deoarece în DEMOS/MP, nucleul nodului destinatar este cel care realizează transferul, se trimite de către nodul sursă un mesaj care conține informațiile necesare pentru migrare: dimensiunea și localizarea contextului procesului, informații relative la swapping; (3) se alocă un nou context al procesului la nodul destinatar; (4) se transferă contextul; (5) se transferă codul procesului; (6) se transferă toate mesajele în tranzit din timpul migrării; totodată nodul sursă modifică adresa procesului astfel încât să reflecte noua localizare a procesului; (7) se înlătură identificatorul procesului din toate tabelele nodului sursă; (8) se restartează procesul în nodul destinație.

Sistemul Emerald exemplifică un alt aspect al migrării: acela al unității de migrare. În Emerald ([Jul88]) unitatea de distribuție și mobilitate este *obiectul*. Unele obiecte conțin procese, altele conțin numai date. Principalul scop al sistemului Emerald a fost de a experimenta mobilitatea în programarea distribuită, subiect al mobilității fiind atât procesele cât și datele.

Un alt aspect important în gestionarea procesoarelor este *alocarea proceselor la procesoare*. Există două soluții: în prima soluție, programatorul asociază procesele la procesoare în mod explicit; în acest caz sistemul de operare este simplu, dar programatorul trebuie să fie implicat în toate deciziile de alocare. În cea de a doua soluție, distribuirea configurației programului este controlată de sistem: sistemul repartizează procesele pe procesoare și poate provoca migrarea acestor procese. Soluția prezintă avantajul unei transparente complete relativ la utilizator, dar și dezavantajul unui sistem de operare mai sofisticat. O strategie de echilibrare a încărcării procesoarelor are două aspecte: *strategia de partajare* a procesoarelor care asigură încărcarea tuturor procesoarelor, ocolind situația în care unele procesoare sunt libere și altele supraîncărcate și *strategia de echilibrare a încărcării* care asigură egalarea sarcinii de lucru pentru toate procesoarele. Ambele strategii utilizează migrarea proceselor și pot fi divizate în trei: (1) *strategia de transfer*, care determină când trebuie executat un proces local sau la distanță; (2) *strategia de selecție*, care determină care proces trebuie transferat; (3) *strategia de localizare*, care determină pe care nod trebuie selectat procesul care trebuie transferat.

Sistemele care folosesc partajarea procesoarelor pot utiliza o gamă largă de algoritmi pentru planificării încărcării începând cu variante complet centralizate și terminând cu variante complet distribuite.

În *varianta centralizată*, sistemul colectează informații despre localizarea nodurilor fără sarcină de lucru, procesele care așteaptă pentru execuția la distanță și de cât timp, parametrii proceselor care se execută și unde se execută și, în funcție de aceste informații, decide ce proces va primi un anumit procesor. Soluția centralizată are însă două dezavantaje: primul este legat de vulnerabilitatea la defecțiunile hardware iar al doilea de acuratelyea informațiilor colectate, greu de menținut datorită întârzierilor introduse de comunicații.

În ceea ce privește *variantele distribuite*, se enumără patru algoritmi cunoscuți și relativ ușor de implementat: (1) algoritmul care organizează procesoarele într-o ierarhie logică ([Wit80], [Gos91], [Tan93], [Tan96]); (2) algoritmul care organizează procesoarele într-un inel logic [Gos91]; (3) algoritmul programelor de tip *worm*; (4) algoritmul *planificării Condor*; dintre aceștia, vor fi prezentați primii doi.

Algoritmul (1) utilizează o împărțire a nodurilor într-o ierarhie logică (fig. 1.3.2.1.5.) în care unele noduri sînt servanți, accesibili global, iar altele gestionează sarcină de lucru a servanților; conexiunile din figură nu reprezintă liniile de comunicație fizice, ci conexiunile logice în scopul gestionării nodurilor. Un nod manager gestionează încărcarea servanților din subarborile său. Alocarea procesoarelor, pentru un job care necesită

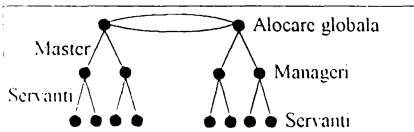


Fig. 1.3.2.1.5. O ierarhie de procesoare.

m procese se realizează astfel (se presupune, pentru simplificare, că procesoarele sunt monoprogramate): (1) fiecare manager cunoaște aproximativ numărul de servanți disponibili din subarborile său; numărul exact nu este cunoscut, deoarece informațiile se transmit periodic. Pentru un job cu *m* procese vor fi necesare *m* procesoare; (2) dacă numărul de procesoare este suficient, adică $S > m$, managerul rezervă un număr de *m* servanți și creează contextul pentru *m* procese; (3) dacă numărul de

procesoare nu este suficient, managerul transmite cererea, către predecesorul său din arborele de organizare logică a procesoarelor. Acestea are informații pentru un subarbor mai mare (numărul de nivele nu este limitat); dacă nici predecesorul nu poate satisface cererea, aceasta este transmisă și mai departe (la predecesorul predecesorului) s.a.m.d. Există două probleme pe care le pune acest algoritm: (1) prima se referă la situația în care un manager are o defecțiune. În acest caz, dacă în arbore există un predecesor al acestui manager, acesta comunică cu servanții și selectează unul dintre aceștia pentru a îndeplini în continuare funcția de manager. Dacă nu există un predecesor al managerului, atunci se folosește faptul că rădăcina structurii ierarhice nu este unică; nodul duplicat a fost introdus tocmai pentru protecția împotriva defecțiilor. Nodurile de pe nivelul rădăcină comunică periodic, transmitându-și informații relative la subarborii lor; prin unare, nodul duplicat poate selecta un alt manager pentru nivelul rădăcinii; (2) a doua problemă se referă la faptul că, dacă mai mulți manageri alocă procesoare, atunci poate apărea o interblocare; aceasta implică un mecanism de sincronizare între manageri.

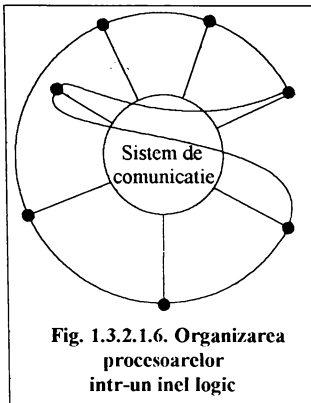


Fig. 1.3.2.1.6. Organizarea procesoarelor într-un inel logic

Cel de al doilea algoritim organizează procesoarele într-un inel logic (figura 1.3.2.1.6.); fiecare nod poate fi și manager și servent. Pentru a asigura alocarea procesoarelor se transmite un token prin inelul logic. Dacă un job care necesită m procese apare în oricare nod, el trebuie să aștepte sosirea token-ului în nodul respectiv. Alocarea procesoarelor unui job ce necesită m procese se desfășoară astfel (se presupune, la fel ca în algoritmul precedent, că nodurile sunt monoprogramate): (1) deoarece nodurile nu pastrează informații despre disponibilitatea altor noduri din sistemul distribuit, nodul cu rol de manager (cel care trebuie să aloce procesoare pentru job) transmite un mesaj broadcast, pentru a se informa despre disponibilitatea procesoarelor. Dacă un procesor este disponibil, trimite un mesaj de confirmare și managerul îl alocă. Dacă numărul de procesoare disponibile este mai mic decât numărul de procese, atunci managerul eliberează procesoarele alocate deja. Pe de altă parte, dacă un nod care și-a transmis disponibilitatea nu primește în decursul unui interval de timp fixat o cerere de creare al unui context de proces, procesorul este eliberat. (2) se transmite token-ul următorului proces din inel; (3) dacă numărul de procesoare disponibile (servanți) a fost găsit suficient, atunci managerul execută toate operațiile necesare pentru crearea contextului proceselor în nodurile servanților; (4) procesoarele servanți sunt eliberate după terminarea proceselor respective.

O altă variantă a algoritmului în inel este aceea în care token-ul transportă informații despre disponibilitatea procesoarelor în inel. La primirea token-ului, un nod manager, contorizează numărul de procesoare disponibile, d , și dacă $d \leq k$ atunci înregistrează localizarea lor și actualizează informațiile din token, după care transmite token-ul mai departe. Acest algoritim (în ambele variante) ridică două probleme: pierderea token-ului și distrugerea inelului în urma defectării unui nod; pentru acest tip de probleme, soluțiile sunt prezentate în subcapitolul 4.4.

Strategia de echilibrare a încărcării procesoarelor spre deosebire de strategia de partajare a procesoarelor trebuie să asigure egalarea sarcinii de lucru a procesoarelor. Strategia de partajare a procesoarelor utilizează migrarea proceselor numai cînd un procesor devine liber; în schimb, strategia de echilibrare poate utiliza migrarea chiar dacă un nod nu este liber. O clasificare a algoritmilor de echilibrare a încărcării este prezentată în fig. 1.3.2.1.7. ([Gos91]).

Algoritmii statici au drept scop determinarea procesorului optim pentru un proces în condițiile în care procesul nu

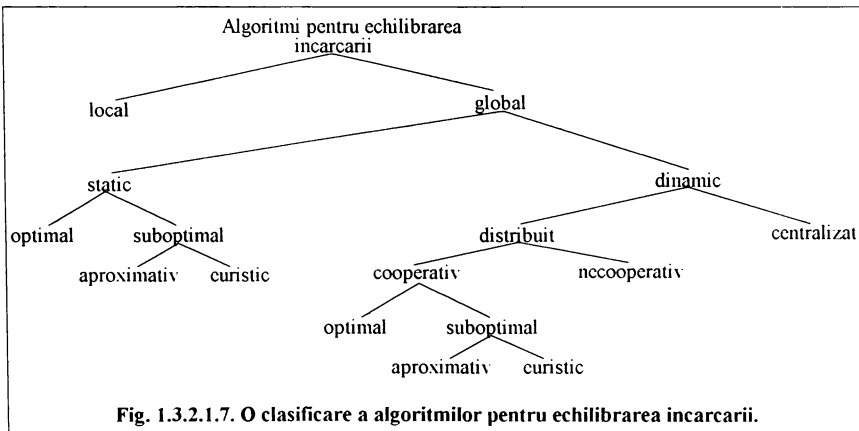


Fig. 1.3.2.1.7. O clasificare a algoritmilor pentru echilibrarea incarcarii.

va mai putea fi mutat de pe acel nod; deciziile sunt luate pe criterii probabilistice, fără a se lua în considerare starea curentă a sistemului distribuit. Echilibrarea încărcării în mod static se reduce astfel la o problemă de matematică: determinarea stării sistemului și a cerințelor de resurse în condițiile în care trebuie optimizați anumiți indici de performanță, pe baza unor funcții luate drept criteriu (minimizarea timpului total de execuție al unui proces, maximizarea utilizării resurselor, maximizarea productivității procesorului -a numărului de procese executate în unitatea de timp). Algoritmii statici au fost dezvoltati pe baza teoriei cozilor, a metodei de programare *branch and bound* și a teoriei grafurilor. Soluțiile optime pot fi foarte costisitoare în ceea ce privește timpul și de aceea se implementează în general soluții suboptimale: curistice sau aproximative. Soluțiile aproximative nu explorează întreg spațiul soluțiilor, scopul fiind determinarea unei soluții satisfăcătoare. Echilibrarea încărcării în mod dinamic are drept scop egalarea sarcinii de lucru a procesoarelor pe măsură ce procesele sunt create, suspendate și reluate, deciziile fiind luate utilizînd informații relative la starea curentă a sistemului. Găsirea unei soluții dinamice este mai complicat de realizat: este necesară colectarea de informații de stare, migrarea proceselor; totuși, o soluție dinamică este mai performantă decît una statică, deoarece ține cont de valoarea reală a indicilor pentru încărcarea sistemului.

1.3.2.2. Gestiunea fișierelor

Diferențele majore între sistemele de gestiune ale fișierelor pentru sisteme centralizate și cele pentru rețelele de calculatoare decurg din problemele de fiabilitate puse de ultimele deoarece trebuie prevenită apariția inconsistențelor datorate unei avarii accidentale a unui nod (client sau server), a liniei de comunicații sau datorate accesului concurrent la fișierele din rețea [Mus95d][Mus96f].

În contextul gestiunii distribuite a fișierelor se definesc următorii termeni: un *serviciu* este o entitate software care se execută pe unul sau mai multe noduri și care furnizează o anumită funcție, un *server* este un proces care se execută pe un singur nod și implementează un serviciu iar un *client* este un proces care poate invoca un serviciu utilizând un set de operații care formează *interfața clientului*.

În rețelele de calculatoare, pentru gestiunea fișierelor există, în general, două soluții:

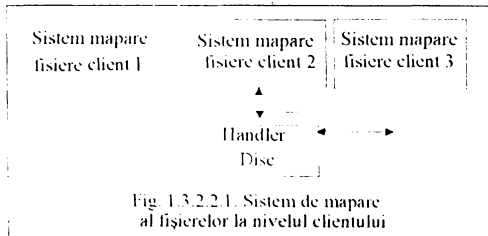
1. **Soluția semidistribuită** prin care se distribuie toate cererile de acces la fișierele din rețea către unul sau mai multe servere de fișiere (un calculator dedicat).

2. **Soluția complet distribuită** în care fiecare nod mentine sistemul sau local de fișiere și comunică cu orice alt nod al rețelei; localizarea fișierelor poate sau nu poate fi transparentă pentru utilizatori. În primul caz sistemul de gestiune mentine tabele și folosește algoritmi pentru a găsi un fișier oriunde în rețea. Un astfel de sistem de fișiere distribuit permite efectuarea de operații la distanță asupra fișierelor, în aceeași manieră ca și asupra celor locale. Un exemplu este sistemul distribuit de fișiere proiectat de Sun Microsystems, NFS (Network File System) și implementat pe microcalculatoare până la supercomputere (în rețele eterogene) în care fiecare nod poate fi server și/sau client; un sistem de fișiere al unui server apare ca parte componentă a sistemului de fișiere al unui client la distanță; scrierea și citirea unui bloc se face prin transmiterea sa prin rețea: deschiderea fișierelor se face la distanță. Alte exemple sunt sistemul de gestiune fișiere al sistemului de operare distribuit LOCUS (compatibil Unix) dezvoltat la UCLA ([215],[1],[8]) în care localizarea fișierelor, transparentă pentru utilizatori, se face printr-un protocol bine stabilit iar pentru protecție se mențin mai multe copii pentru un fișier și Apollo Domain dezvoltat de Apollo Computer Inc.

Un *sistem distribuit de fișiere (DFS)* este o implementare distribuită a modelului clasic time-sharing de partajare a fișierelor; sistemul de fișiere al Unix-ului este desori luat ca model.

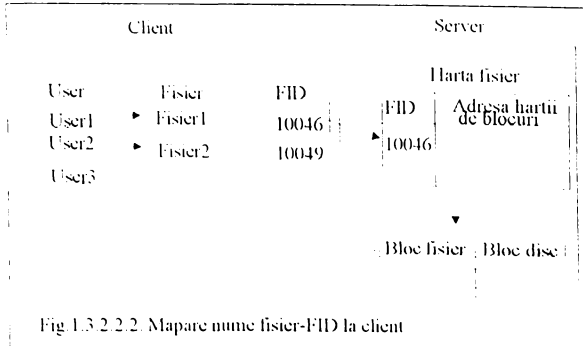
O clasificare a serverelor de fișiere este prezentată în următorul tabel ([Gos91]).

	Server simplu	Server universal	Suport pentru baze de date
Unitate de acces la date	<i>fișier</i>	<i>secvența de octeți</i>	<i>înregistrare</i>
Unitate de blocare	<i>fișier</i>	<i>fișier</i>	<i>înregistrare</i>
Domeniul unei actualizări atomice	<i>fișier</i>	<i>mai multe fișiere</i> <i>mai multe servere</i>	<i>mai multe fișiere</i> <i>mai multe servere</i>



Cea mai simplă organizare a unui server pentru fișiere este ca un remote disc al fiecărui client. În acest model, fiecare client are alocat un disc virtual, care este de fapt o zonă din discul serverului și pe care îl adresează la fel ca pe un disc local. Rețeaua de comunicații are rolul de a simula un controler de disc la distanță, serverul furnizând numai execuția comenzilor *read, write* la nivel de bloc la fel ca și un driver de disc. Întregul sistem de gestiune al fișierelor la distanță se află în nodul clientului, la fel ca și pentru discurile locale.

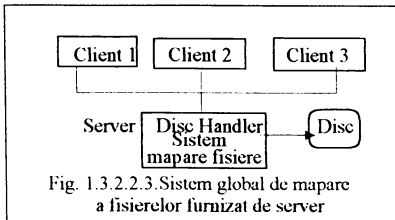
Următoarea treaptă de complexitate în organizarea serverelor pentru fișiere este ca acestea să gestioneze fișiere, dar cataloagele care inventariază fișierele unui client să fie menținute de către client (fig. 1.3.2.2.1.).



Serverul gestionează alocarea și dealocarea spațiului pe disc (dinamica) și furnizează clientului o interfață pentru comunicație, care include operații (creare, stergere fișier, citire, scriere, poziționare într-un fișier) implementate într-un protocol de comunicație. La nivelul serverului se mențin identificatori globali ai fișierelor, accesibil eventualelor tuturor clienților din rețea. La crearea fișierului se generează și se returnează clientului un identificator de fișier, FID, (în general număr de 8-64 biti); clientul poate transmite FID-ul altor clienți pentru a avea acces la acei fișier și înserează FID-ul în catalogul fișierelor

propriu (figura 1.3.2.2.2).

Pentru a evita situația în care rămân fișiere inaccesibile în urma esuării înregistrării FID-ului lor de către clienți, serverul trebuie să furnizeze o comandă prin care un client să poată cere o listă completă a fișierelor sale. Se pot încorpora în FID, alături de un număr aleator și informații relative la localizarea fișierelor făcând imposibilă deplasarea fișierului. Unele rețele, într-o astfel de organizare a serverului oferă clienților și serviciile de catalog (*directory service*) [Tan93] (fig. 1.3.2.2.3); aceste servicii sunt cumulate într-un proces separat, rezident în nodul serverului, care apare ca un client pentru server și ca server pentru client. Izolarea serviciilor de catalog face posibilă menținerea mai multor scheme de evidență și denumire a fișierelor utilizând însă aceleași servicii de server.

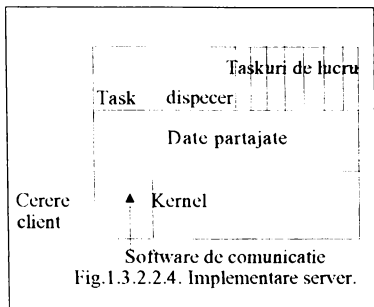


Această structurare este echivalentă cu următoarea treaptă de organizare a unui server -aceea în care serverul furnizează un sistem de gestiune fișiere complet, posibil al unui sistem de operare cunoscut și poate trata alături de comenzile pentru manipularea fișierelor și comenzi pentru creare și stergere cataloage, creare de legături.

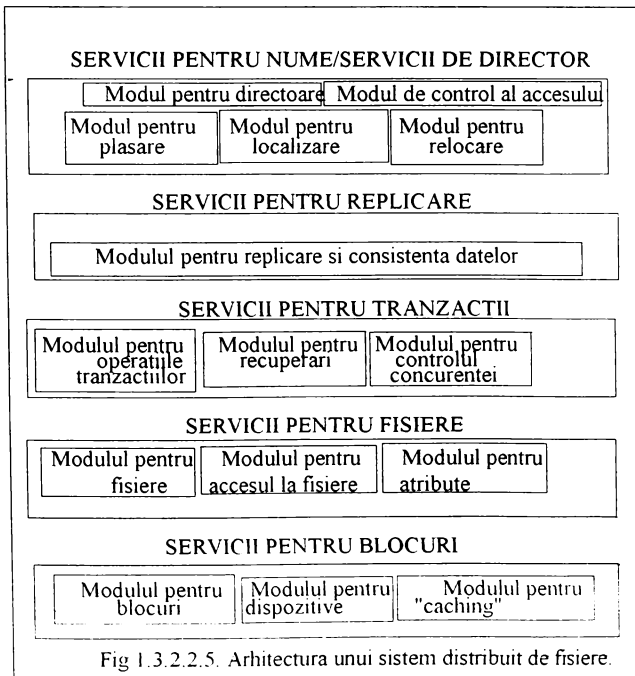
În ceea ce privește organizarea serverului, cea mai simplă organizare este de a-l considera ca un singur proces având următoarea structură:

```
do true -->
  GetMessage(buffer);
  DoWork(buffer,result);
  SendReply;
od
```

Dezavantajul este că prelucrarea unei cereri necesită în general mai multe accesuri la disc, timp în care serverul nu mai poate accepta noi cereri.



O variantă mai complexă (fig 1.3.2.2.4) este aceea în care serverul este descompus în mai multe *task-uri de lucru* concurente care partajează structuri de date comune precum tabele ale sistemului de fișiere și buffere, un *kernel* care acceptă mesajele de la clienți și un *dispatcher* care administrează încărcarea *task-urilor de lucru* (fiind în același spațiu de adrese se transmite taskurilor o referință la mesaj).



O analiză a sistemelor distribuite pentru gestiunea fișierelor din [Gos91],[Nut92],[Mul89],[Tan93] reliefează elemente comune care permit propunerea arhitecturii generale prezentate în figura 1.3.2.2.5. pentru un astfel de sistem.

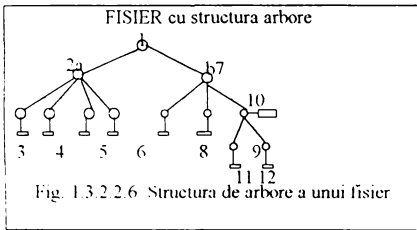
Serviciile furnizate de sistemele distribuite pentru gestiunea fișierelor sunt asigurate de următoarele module:

1. Modulul pentru dispozitive care asigură execuția operațiilor de I/O și gestiunea bufferelor.
2. Modulul pentru blocuri care alocă blocurile de disc.
3. Modulul pentru implementarea mecanismului de "caching" care păstrează rezultatele obținute recent prin operații cu discul sau operații la distanță.
4. Modulul pentru fișiere care asociază identificatori fișierelor.
5. Modulul pentru accesul la fișiere care asigură citirea și scrierea datelor în fișiere.
6. Modulul pentru atribute care citește și scrie atributele asociate fișierelor.
7. Modulul pentru operațiile tranzațiilor care furnizează facilitatea lucrului cu tranzații atomice.
8. Modulul pentru recuperare care furnizează mecanismul de recuperare după caderea unor componente.
9. Modulul pentru controlul concurenței care asigură consistența datelor în fișierele partajate.
10. Modulul pentru replicare și consistența datelor care asigură replicarea fișierelor și pastrarea consistenței copiilor.
11. Modulul pentru directoare care asigură corespondența între numele date de utilizatori fișierelor și identificatorii fișierelor.
12. Modulul de control al accesului care asigură verificarea drepturilor de acces pentru operația cerută.
13. Modulul pentru plasare care selectează serverul ce va primi o nouă copie a fișierului.
14. Modulul pentru localizare care găsește serverul ce gestionează un anumit fișier.
15. Modulul pentru relocare care îmbunătățește performanțele sistemului prin relocarea fișierelor.

Separarea serviciilor pentru blocurile discurilor de cele pentru fișiere este o caracteristică a majorității sistemelor distribuite de fișiere; această separare asigură, pe de o parte posibilitatea combinării a diferite metode de stocare și, pe de o parte, utilizarea unei stocari duale (*stable storage*).

Din punctul de vedere al utilizatorilor, serverele pentru fișiere se caracterizează prin: *structura fișierelor, atributele lor și operațiile permise* asupra fișierelor [Tan88][Mus96f][Mus96d].

Există trei *modele de structuri de fișiere* utilizate de serverele pentru fișiere: primul model este *fișierul brut*, a cărui structură internă nu este cunoscută de server; singurele operații posibile sunt scrierea și citirea întregului fișier; următorul model este fișierul constituit dintr-o *secvență ordonată de articole*, eventual nu toate de aceeași dimensiune. Serverul poate furniza operații precum adăugare, înlocuire sau ștergere de articole. Un fișier Unix poate fi privit ca având articole de un octet, accesibile prin poziție. Modelul cel mai complex este *fișierul ierarhic*, organizat pe o structură de arbore. Aceasta poate fi o structură de tip *secvențial indexat* (VISAM - Virtual Indexed Sequential Access Method) în care fiecare articol se identifică printr-o cheie, sau o structură arborescentă mai generală în care fiecare nod al unui arbore poate conține un articol, amindouă sau nici unul. Un exemplu este prezentat în fig.1.3.2.2.6.: un arbore cu 12



Atribut	Tip	Setat la creare	Setat de user	Setat de server
Nume fișier	String	x	x	
Operații permise	Harta biti	x		
Control acces	Lista			
Numar de cont	INT	x		
Data creerii	DATA	x		
Data ultimei modificari	DATA	x		
Data ultimei citiri	DATA	x		
Data ultimei modificari a atributelor	DATA			
Proprietar	UserID	x		
Identitatea ultimului modificador	UserID	x		
Identitatea ultimului cititor	UserID	x		x
Identitatea ultimului modificador al atributelor	UserID	x		
Tipul continutului	Obiect	x		
Cheie de criptare	String	x		
Dimensiune	Long	x		
Dimensiune maxima evaluata	Long	x	x	
Calificatori legali	String	x	x	
Utilizare proprie	String	x	x	

Fig. 1.3.2.2.7. Atribute asociate fișierelor

noduri, două având etichete (a și b) și nouă noduri având articole (toate nodurile exceptând 1, 2 și 7). Prin liniarizarea arborelui, nodurile se ordonează (în exemplu prin parcurgere în Depth First), nu numai pentru o cautare mai rapidă a nodurilor necetichetate în scopul ștergerii sau înlocuirii lor, dar și pentru adresarea convenabilă a subarborilor, în scopul efectuării de operații la nivel de subarbori.

Atributele asociate fișierelor diferă mult de la un server la altul. Unele atribute, stabilite la crearea fișierelor rămân neschimbate, altele pot fi modificate explicit prin operații sau actualizate automat de server.

Standardul OSI [10] prevede pentru fișiere atributele prezentate în fig. 1.3.2.2.7. precum și operațiile prezentate în figura 1.3.2.2.8.

Operatie	Aplicabila fisier	Aplicabila articol	Bitmap
Creare fisier	x		
Stergere fisier	x		x
Selectie fisier	x		
Deselectie fisier	x		
Inchidere fisier	x		
Citire attribute	x		x
Modificare attribute	x		x
Localizare		x	
Citire		x	x
Inserare		x	x
Inlocuire		x	x
Adaugare		x	x

Fig 1. 3.2.8. Operatii asociate fisierelor.

Pentru comunicarea între clienți și serverele de fișiere se folosesc protocoale care trebuie să satisfacă două cerințe specifice rețelelor: (1) controlul fluxului de mesaje dintre server și client astfel încât nici unul să nu fie saturat (2) controlul erorilor datorate transmisiei prin liniile de comunicație, sau defectării serverului sau clientului. Astfel, unele servere impun clienților lor să emită cel puțin o cerere într-un interval de timp maxim fixat. În absența unor astfel de mesaje de tip *keep alive*, serverul abortează tranzacția asociată clientului. Trebuie luate în considerare două aspecte: clienți care nu răspund și cereri duplicate (răspunsul serverului nu a fost recepționat de client care lansează o nouă cerere).

Se utilizează în general următoarele protocoale: (1) *protocolul în doi pași*, bazat pe modelul *send - request / receive - response* (LOCUS, DEMOS); (2) *protocolul în trei pași*, bazat pe modelul *send - request / receive - response / send - acknowledgement* (RRA) (SWALLOW); (3) *protocolul într-un singur pas, no-wait send, (single shot)* care impune ca toate cererile să fie *idempotente* (execuția repetată are același efect ca și o singură execuție); (4) *apel de procedură la distanță* (NFS). Parametrii fiecărui apel de procedură executate în cadrul protocolului NFS, bazat pe RPC, conține toate informațiile necesare execuției unei cereri, sau, altfel spus, serverul nu reține nici o informație din cadrul cererii precedente; astfel de servere se numesc *servere care nu păstrează informații de stare* (stateless servers) și prezintă avantajul unei recuperari simple după o avarie a serverului (se repetă cererea clientului). Pe de altă parte, însă, un server care nu păstrează informații de stare are și următoarele dezavantaje: (1) este un model inadecvat pentru tratarea acceselor concurente a mai multor clienți; (2) cererile sunt mai lungi, având încorporați mulți parametri.

Mecanismul utilizat pentru a retrage resursele alocate de server unui client care nu mai răspunde sunt: (1) *timeout* la nivelul unei tranzacții (SWALLOW); (2) *timeout* la nivelul cererilor; (3) transmiterea de mesaje *keep-alive*.

Un alt aspect important al gestiunii distribuite a fișierelor este *transparența*, ceea ce implică faptul că clienții pot accesa fișiere la distanță utilizând aceleași operații aplicabile fișierelor locale. Aceasta înseamnă că interfața clienților într-un sistem distribuit este uniformă, nu distinge fișierele locale de cele la distanță; este sarcina sistemului să localizeze fișierele și să realizeze transportul datelor. Legată de transparență este problema numelor fișierelor. Regula de formare a numelor fișierelor stabilește o corespondență între obiectele logice cu care operează utilizatorii și obiectele fizice (blocurile fizice ale fișierului); ea trebuie să asigure, pe de o parte *transparența localizării* (numele unui fișier nu trebuie să conțină indicații relative la locația sa fizică) și pe de altă parte *independența față de localizare* (numele unui fișier nu trebuie să se modifice când fișierul este mutat).

Independența față de localizare, adeseori numită și *migrarea fișierelor* sau *mobilitatea fișierelor* este o proprietate mai puțin întâlnită la sistemele distribuite, care oferă în general o mapare statică, și transparența localizării (Locus, NFS, Sprite). Independența față de localizare este oferită de Andrews și câteva sisteme distribuite experimentale (Eden, de exemplu).

Există trei tehnici pentru denumirea directorilor și fișierelor:

(1) cea mai simplă tehnică este prin combinarea numelui nodului și fișierelor, care garantează astfel unicitatea în tot sistemul: *host: nume_local*. Evident că tehnica nu oferă nici transparența localizării și nici independența față de

localizare; totuși, pot fi folosite aceleași operații atât pentru fișierele la distanță, cât și pentru cele de la distanță. Structura DFS-ului constă într-o colecție de componente izolate, care împreună formează un sistem de fișiere. (2) a doua tehnică, popularizată de SUN NFS se bazează pe operația de *montare* a directoarelor de la distanță într-un spațiu de nume local; după această atașare, fișierele pot fi denumite într-o manieră transparentă din punct de vedere al

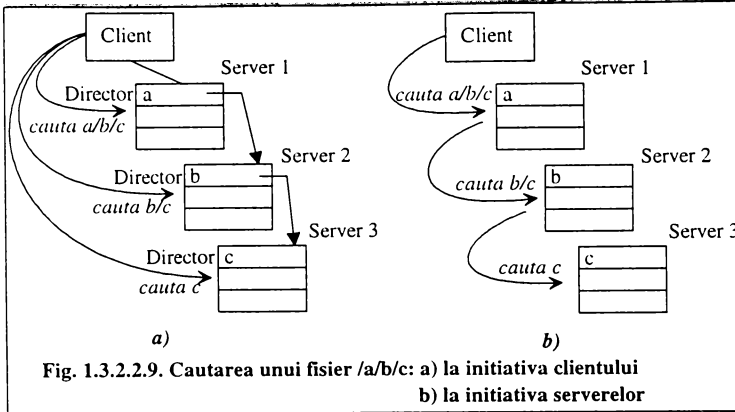


Fig. 1.3.2.2.9. Cautarea unui fișier /a/b/c: a) la inițiativa clientului
b) la inițiativa serverelor

localizării. O proprietate a acestei tehnici, uneori considerată un dezavantaj, este aceea că spațiile de nume paratașate în sistem nu sunt identice pentru toate nodurile; (3) a treia tehnică oferă o structură globală de nume în sistem, care cuprinde toate nodurile; întregul spațiu de nume este vizibil tuturor clienților. În ceea ce privește implementarea acestei tehnici, se specifică următoarele aspecte: (a) *translatarea căilor*, în fig. 1.3.2.2.9.a și 1.3.2.2.9.b se

prezintă două modalități de determinare a localizării unui fișier /a/b/c; în fig. 1.3.2.2.9.a cautarea se realizează la inițiativa clientului (varianta adoptată în Andrew și Locus), iar în fig. 1.3.2.2.9.b căutarea se realizează la inițiativa serverului (variantă adoptată în NFS și Sprite); (b) *identificatorii fișierelor*, care reprezintă numele fișierelor la nivelul de sistem , uzual întregi de lungime fixă, se pot obține prin concatenarea următoarelor câmpuri: adresa serverului care a creat fișierul, o marcă de timp de la ceasul serverului sau un întreg care reprezintă valoarea unui contor al serverului incrementat la crearea fiecărui fișier. Asemenea identificatori prezintă avantajul că identifică în mod unic un singur obiect fișier în întreg fișierul, oferă independență de localizare (adresa serverului este necesară numai pentru asigurarea unicității, nu și pentru localizare), pot fi transmiși între procese și de la nod la nod, fără a fi nevoie de transformarea lor la fiecare pas. (c) folosirea unei proceduri de *cache-ing* pentru informațiile de localizare la nivelul clienților. Exemple sunt sistemul Andrews, sistemul Apollo Domain și sistemul Sprite care utilizează un mecanism similar *cache-ing*-ului, așa numitele *tabele de prefix*, iar când serverul cautat nu se află în aceste tabele, transmite mesaje broadcast.

Un alt aspect important într-un DFS este *partajarea fișierelor*. În acest sens, există următoarele implementări; (1) utilizarea *modelului semantic Unix*; în acest caz toate cererile de citire și scriere se adresează direct serverului, care le procesează secvențial. Unele DFS-uri, precum Locus, Sprite, încearcă să emuleze această semantică. Soluția este aplicabilă în practică în special în cazul unui singur server și în lipsa *cache-ing*-ului la nivelul clienților; (2) o alternativă la semantică Unix este întârzierea propagării modificărilor în fișier, până în momentul închiderii fișierelor (*semantică de sesiune*). Problema pe care o ridică acest model este situația în care doi sau mai mulți clienți modifică independent o copie *cache* a fișierului. O soluție este ca, la închiderea fișierului, acesta să fie înapoiat serverului. Acest model este însă adecvat pentru aplicațiile la care serializarea acceselor este esențială (baze de date); aceste aplicații trebuie să-și coordoneze accesul explicit. (3) O implementare diferită de cele de mai sus este cea a *fișierelor partajabile imuabile*; acestea, după creare, nu mai pot fi alterate, iar numele lor nu mai pot fi reutilizate. În acest caz, fiecare fișier este reprezentat printr-o istorie de fișiere imuabile; se crează câte o nouă versiune, la fiecare încercare de modificare; (4) o extindere a modelului *semantică de sesiune* este *modelul tranzacțiilor atomice*: efectul sesiunilor este echivalent cu execuția lor serială, într-o anumită ordine. Blocarea fișierului pe parcursul unei sesiuni este o soluție pentru realizarea acestui model.

Un alt aspect al unui DFS este implementarea *accesului la distanță*; în acest sens există următoarele soluții:

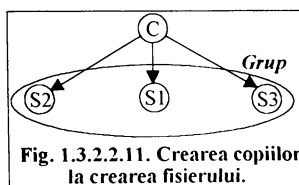
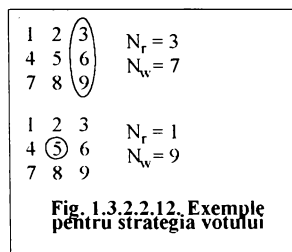
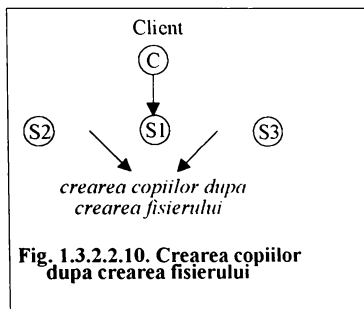
(1) rezolvarea cererilor clienților numai la nivelul serverului (servicii la distanță): serverul recepționează și tratează cererile și transmite rezultatele înapoi la client. În această implementare există o corespondență directă între cererile clienților și traficul mesajelor;

(2) utilizarea unui mecanism de *cache-ing* (la nivelul serviciului și/sau nivelul clienților). Există următoarele decizii importante relative la proiectarea unui mecanism de *cache-ing*: (a) *granularitatea datelor* menținute în *cache*: întregul fișier sau blocuri din fișier. De exemplu, sistemul Andrews, după o versiune în care menținea în *cache* întreg fișierul, a fost rescris într-o versiune în care în *cache* se mențin câte 64Kb din fișier. Prin creșterea unității de memorare în *cache*, crește evident probabilitatea ca următorul acces să fie realizat local, dar crește și probabilitatea apariției problemelor de consistență. Pe de altă parte, deoarece unitatea de transfer în rețele este relativ mică (pachetele Ethernet au aproximativ 1,5Kb) o unitate de memorare în *cache* mare necesită operații de împachetare la emisie și despachetare la recepție. În mod tipic *cache-ing*ul la nivel de bloc utilizează tehnica *read-ahead* pentru citirea fișierelor. (b) *localizarea cache*-ului: în memorie sau pe un disc local. Localizarea pe disc are avantajul supraviețuirii datelor în urma unei defecțiuni; pe de altă parte, însă localizarea în memorie permite lucrul cu stăii fără discuri, iar datele pot fi accesate mai rapid. (c) *propagarea modificărilor din cache*. Există mai multe soluții: cea mai simplă este scrierea imediată a datelor pe server, la orice modificare (strategia *write-through*) avantajul fiind siguranța. Această strategie este însă echivalentă cu realizarea serviciilor de scriere la distanță; numai citirea beneficiază de fapt de *cache-ing*. Altă soluție este scrierea cu întârziere a datelor pe server (strategia *delayed-write*)

delayed-write) care are două variante: scrierea pe server se face periodic (de exemplu în Sprite, la fiecare 30s), fie la închiderea fișierului - *write-on-close*; (d) *validarea datelor din cache*. Utilizarea *cache*-ingului la nivelul clienților introduce problema consistenței datelor. Există două soluții de bază pentru verificarea validității datelor din *cache*: verificarea la inițiativa clientului și verificarea la inițiativa serverului. În cazul primei soluții, frecvența verificării este crucială: ea poate fi făcută la fiecare acces la un fișier sau numai la deschiderea fișierului sau poate fi făcută periodic. În cazul celei de-a doua soluții, serverele înregistrează ce parte a fișierelor se află în *cache*-ul fiecărui client. La detectarea unui conflict (când cel puțin un client accesează fișierul pentru scriere), dacă de exemplu se utilizează semantica de sesiune, la închiderea oricărui fișier modificat, serverul anunță ceilalți clienți să-și descarce *cache*-ul ca invalid. Dacă se implementează însă o semantică mai restrictivă (de exemplu semantica Unix) atunci serverul trebuie să înregistreze deschiderea oricărui fișier și modul de lucru cu acesta (read sau write); serverul trebuie să interzică deschiderea fișierelor într-un mod conflictual, pur și simplu dezactivând *cache*-ul pentru acel fișier și furnizând execuția serviciilor la distanță (ca de exemplu în Sprite).

Multe implementări au fost realizate prin combinarea celor două soluții de mai sus (*soluții hibride*). Astfel, în LOCUS și NFS implementarea se bazează pe soluția 1 dar este completată, pentru mărirea performanțelor cu un mecanism de *cache*-ing. Pe de altă parte, implementarea accesului la distanță în Sprite se bazează pe *cache*-ing, dar, așa cum s-a aratat, în anumite situații se poate adopta și soluția serviciului la distanță.

Pentru creșterea performanțelor unui DFS, se folosește desigur *replicarea fișierelor*. Sunt trei mari avantaje oferite de replicare: (1) creșterea siguranței de funcționare, având mai multe copii ale unui fișier; (2) disponibilitatea unui fișier și în cazul avariei unui server; (3) divizarea cererilor de acces pe mai multe servere, evitându-se "sufocarea" unui singur server. Principalul aspect legat de replicare este *transparența*. În acest sens există trei implementări: (1) controlul copiilor se face de către client; (2) la crearea fișierului, se realizează o singură copie, urmînd ca celelalte să fie realizate ulterior (verificîndu-se dacă nu cumva fișierul s-a modificat între timp (fig. 1.3.2.2.10)); (3) se creează toate copiile la crearea fișierului utilizîndu-se comunicațiile de grup (fig. 1.3.2.2.11).



În ceea ce privește modificarea copiilor fișierelor există următoarele soluții: (1) *strategia master-slave* sau a *copiei primare* în care, pentru fiecare fișier replicat se desemnează un *server master* care conține singura copie ce poate fi actualizată. Strategia este adecvată în aplicații în care fișierele nu se actualizează frecvent și în care actualizările pot fi realizate numai într-un punct central. Un exemplu de aplicare al acestei strategii este serviciul Sun Yellow Pages; (2) *strategia votului* în care ideea de bază este de a obliga clienții să obțină permisiunea de la mai multe servere pentru a scrie sau citi un fișier replicat ([Tan88][Giff79]). Astfel, Gifford propune ca pentru ambele operații să existe câte un *quorum* (*read quorum* și *write quorum*) N_r , respectiv N_w pe care clienții trebuie să-l obțină pentru a executa operațiile respective. După obținerea acestor permisiuni, fișierul este modificat și i se asociază o nouă veriune. În fig. 1.3.2.2.12 se prezintă două exemple.

1.4. Modele în proiectarea sistemelor de operare distribuite

Există două modele de proiectare a unui sistem de operare distribuit: *modelul proces* și *modelul obiect*. Diferențele majore între cele două modele de proiectare constau în modul în care sunt implementate conceptul de entitate funcțională și sincronizarea.

1.4.1. Modelul proces

Modelul proces propune sistemul de operare distribuit ca fiind un ansamblu de procese, utilizate pentru controlul resurselor din rețea (fig. 1.4.1). Fiecărei resurse din rețea i se asociază un proces pentru a o gestiona. Un proces este văzut ca o zonă de cod și de date, și se află într-o anumite stare.

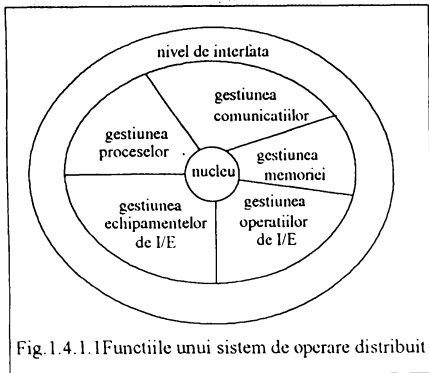


Fig. 1.4.1.1 Funcțiile unui sistem de operare distribuit

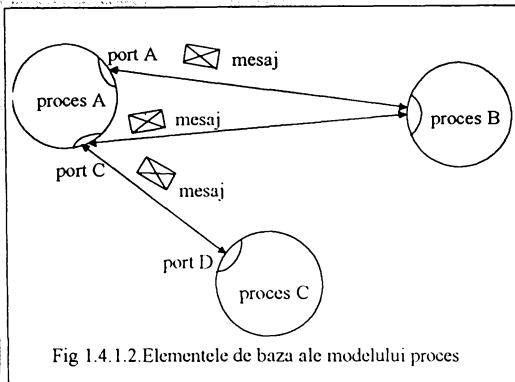


Fig. 1.4.1.2. Elementele de baza ale modelului proces

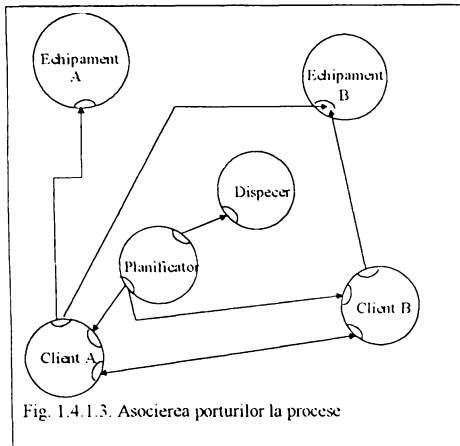


Fig. 1.4.1.3. Asocierea porturilor la procese

Procesele sistemului de operare gestionează resursele și interacționează; ele sunt utilizate de procesele aplicațiilor rezultând astfel un sistem operațional. Procesele sistem sunt denumite *procese publice* sau *servere* (în contrast cu procesele lor utilizator denumite *clienți*). Procesele server sunt de două tipuri: *utilizabile serial* sau *paralel*, în funcție de tipul serviciului implementat. Procesele client și server comunică prin mesaje pentru a-și putea transmite informații de control, de stare și de date. Toate acțiunile din sistem sunt îndeplinite via mesaje. Perechea (*proces, mesaj*) este conceptul de bază în acest model, entitatea care îl distinge de modelul obiect; procesul devine mecanismul care implementează diverse politici de gestiune și control iar mesajele asigură cooperarea și sincronizarea interprocese (fig. 1.4.1.2).

Principalele mecanisme folosite în modelul proces pentru interconectarea logica/fizică a proceselor sunt *porturile* și *canalele*. Conceptul de *port* definește o entitate logică. Realizarea sa fizică poate consta într-un buffer adresabil, încărcat prin hardware la detectia unui identificator particular. Porturile pot fi de natură *statică* sau *dinamică*. Porturile a căror alocare se

face static sunt adecvate cind o pereche de procese comunică prin porturile sursa-destinație pe tot parcursul execuției lor; este cazul situațiilor cind procesele sunt create asemenea corutilor sau cind procesele distribuite necesită un nivel ridicat de sincronizare (vezi subcap. 3.2.6). În contrast, porturile dinamice sunt create, asociate, utilizate și conectate de o manieră dinamică în timpul execuției proceselor. Un exemplu tipic este utilizarea dinamică a porturilor alocate unui proces server care controlează un echipament specializat. Procesul este accesibil pentru comunicare oricărui alt proces din sistem care necesită serviciile sale și care posedă drepturi de utilizare a echipamentului. Portul alocat procesului server este ținut permanent activ: în momentul în care un proces client dorește să utilizeze echipamentul, își creează un port (sau utilizează un port deja creat) pe care îl conectează la portul procesului server, stabilește parametrii de control și utilizează apoi echipamentul prin transfer de mesaje via cele două porturi. La încheierea utilizării, echipamentul periferic este eliberat și conexiunea dintre cele două porturi este ruptă. Prin conectarea proceselor de o manieră corectă se pot implementa diverse sincronizări între procese precum și o multitudine de schimburi de date care pot furniza, în final, importante servicii de cooperare în rețea. De exemplu, pentru a-și alocă echipamente de intrare-ieșire și de unități de prelucrare, două procese client A și B trebuie să creeze porturile și conexiunile din fig. 1.4.1.3.

Transferul de mesaje între procese asigură nu numai schimbul de informații ci și sincronizarea între procese. Mesajele pot fi transmise de la un proces A la un proces B printr-un mecanism de cutie poștală sau utilizând primitive pentru emisie-recepție de mesaje. Mesajele pot fi transmise cu blocare, cind emițătorul sau receptorul așteaptă confirmarea efectuării operației de emisie, respectiv recepție, ceea ce asigură și sincronizarea între procese, sau fără blocare. Pentru a face mecanismul flexibil, se prevede, în general și o variantă de transmisie fără blocare, alături de operațiile implementate cu blocare (vezi subcap. 3.2.). Mecanismele fără blocare, ca de exemplu, cutiile poștale nu asigură sincronizarea implicită. Mesajele sunt transmise dintr-o cutie poștală în cealaltă, receptorului i se semnalează că s-a primit un mesaj, după care, acesta extrage mesajul și își continuă prelucrarea. Acest mecanism se utilizează în cazul proceselor slab cuplate, care achiziționează date unul de la celălalt în mod ocazional. Mesajele reprezintă astfel mecanismul prin care procesele își comunică unele altora informații de control, de date, de stare; schema care implementează transmisia trebuie să asigure siguranța în funcționare. Ea trebuie să semnaleze proceselor cind s-a transmis un mesaj, cind s-a recepționat, și cind s-a pierdut un mesaj. Aceasta implică posibilitatea ca nivelurile inferioare ale rețelei să poată detecta și corecta un mesaj incorect, eventual prin retransmitere.

În general, la proiectarea unui sistem de operare distribuit se are în vedere o anumită politică în conformitate cu care se tratează sincronizarea proceselor. Această politică încadrează mecanismul de sincronizare într-una din următoarele trei clase: *sincronizare slabă, sincronizare puternică sau o combinație a acestora*. O sincronizare puternică se referă la situația unui proces care nu își poate continua prelucrarea pînă cind un alt anumit proces nu atinge punctul de sincronizare. Este cazul corutinelor. O altă formă de sincronizare puternică o constituie folosirea semafoarelor. Semafoarele furnizează proceselor garanția că se pot "întîlni" într-un anumit punct dacă acest lucru este necesar. Pentru asigurarea unei sincronizări slabe, o schemă folosind cutii poștale este suficientă. În cadrul acesteia, se rezervă și asociază o cutie poștală fiecărui proces sau unui grup de procese. Cînd un proces intenționează să comunice cu un altul, trimite mesajul în cutia poștală asociată destinatarului; concomitent se setează și un indicator care semnalează că procesul destinat să primească un mesaj. Cînd este posibil, procesul destinat să scanează cutia poștală și extrage mesajul respectiv.

1.4.2. Modelul obiect în proiectarea sistemelor de operare distribuite

Un sistem de operare proiectat după modelul *obiect* are o funcționalitate similară cu modelul proces, dar conceptele utilizate pentru construcția sa sunt diferite. Modelul obiect prevede o structură de bază denumită obiect compusă din două parti: o specificitate externă și un corp al obiectului. *Corpul obiectului* este cunoscut și accesibil numai obiectului însuși și reprezintă sfera sa de control. *Specificitatea externă* conține structurile sale accesibile din exterior. Conceptul de obiect semnifică atât entități hardware cit și software și este o generalizare a tipurilor de date abstracte: un obiect se caracterizează, în ceea ce privește reprezentarea sa, prin structurile sale de date (fișiere, directoare, liste, cutii poștale, resurse hardware, module executabile) și prin setul invariant de operații ce se pot executa asupra acestor structuri de date. Obiectele pot fi și incubate unele în altele.

În modelul obiect este posibil să coexiste diferite tipuri de comunicații pentru comunicațiile interobiect și întraobiect. Controlul și cooperarea între obiecte se realizează prin achiziția de capacități. Astfel, accesul la o resursă în sistem se face prin achiziția unei capacități de la managerul resursei respective: la terminarea utilizării, capacitatea este eliberată și poate fi pusă la dispoziția altor cereri.

Interacțiunea în întreg sistemul se bazează în întregime pe operații executate asupra obiectelor. Tipic, obiectele implementează un număr mic de operații, dar prin combinarea în diverse structuri a obiectelor se pot obține funcții complexe. Se obține astfel o flexibilitate și modularitate sporită într-un astfel de sistem.

Această implementare conduce la structurarea sistemului ca o colecție de resurse și interfețe la resurse. Sistemul apare ca o colecție de obiecte mai mult sau mai puțin dependente unele de altele. Se spune că sistemul de operare creat cu un astfel de model capătă o structură orizontală [Jon78]. Sistemul apare însă și ca o ierarhie de abstracțiuni. Aici intervine conceptul de menținere din limbajele de programare orientate pe obiecte și, mai general, conceptul de modularitate, care conduce și la structurarea pe verticală a sistemului. Cele mai complexe obiecte apar utilizatorului ca o interfață disponibilă în programare, pe baza căreia acesta poate să-și construiască aplicații. Avantajul acestei implementări este că această interfață este, la rîndul ei prevăzută cu un alt set de interfețe (posibil ierarhic) pentru accesul la obiecte. Utilizatorul poate continua o astfel de dezvoltare ierarhică și scopul extensibilității modelului este atins. Modelul obiect prezintă avantajul unui protocol uniform pentru accesul la obiecte, care interacționează, în general prin transferul de mesaje.

Luînd în considerare implementarea paralelismului în modelul orientat pe obiecte, o clasificare [DG89][Mus93d] a sistemelor de operare distribuite este reprezentată în fig. 1.4.3.1.

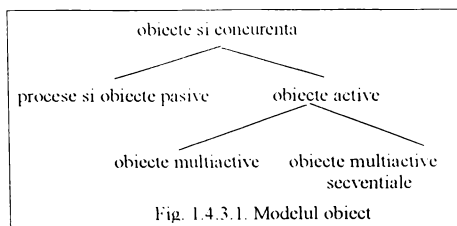


Fig. 1.4.3.1. Modelul obiect

Modelul *proces și obiecte pasive* semnifică de fapt prezența simultană a două modele: proces și obiect. Entitățile active sunt reprezentate ca procese, iar cele pasive ca obiecte. Avantajul reprezentării entităților pasive ca obiecte rezidă în structurarea datelor. Exemple sunt sistemele Emerald [CC91] [DG89]. Totuși, acest model este considerat a fi la granița dintre modelele proces și obiect [DG89] datorită următoarelor motive: (1). Modelul trebuie descris folosind și entitatea de proces, pierzîndu-si

astfel caracteristică de a fi structurat pe baza unei singure entități (obiectul); (2). Există două forme de interacțiune între componentele sistemului: transferul de mesaje și cooperarea între procese. Ultima formă se reprezintă în mecanisme tradiționale partajării unei memorii comune (semafoare, monitoare). Cooperarea între procese se specifică independent de prezența obiectelor.

Modelul *obiecte active* pastrează cele două caracteristici esențiale ale modelului obiect: încapsularea și abstractizarea. Un obiect poate fi și o entitate activă; el capătă proprietățile unui proces sau unui grup de procese. Aceste obiecte active se execută în paralel. Interacțiunea între obiectele active se descrie uniform printr-un singur mecanism: transfer de mesaje. La nivelul unui obiect, recepția mesajelor va fi implementată printr-un *controller* (uneori denumit *guardian*). Sincronizarea este integrată în mecanismul de transfer. Modelul *obiecte active* poate fi în continuare divizat în două submodele: *obiecte multiactive și obiecte active secvențiale*.

În submodelul *obiecte multiactive* un obiect activ apare ca un set de activități, executate concurent și corespunzând mesajelor recepționate de către obiect. O activitate este creată și inițiată dinamic la recepția unui mesaj și distrusă după încheierea tratării mesajului. Toate activitățile concurente ce se execută în cadrul unui obiect partajează o zonă de date comune, interne obiectului. Gestiunea acestor activități (proces) este complet disociată de gestiunea obiectelor. Submodelul este caracteristic sistemelor de operare distribuite Argus, Eden și Mach [BF85], [DG89]. O consecință a acestei implementări este necesitatea sincronizării activităților interne, care partajează același context de execuție. Există, de aceea, în acest model, două nivele de cooperare: între obiecte, pe baza transferului de mesaje și între activitățile interne, pe baza partajării unei memorii comune. Astfel, se folosesc semafoare în Argus și monitoare în Eden. Soluția implementată în submodelul *obiecte multiactive* oferă avantajul paralelismului la nivelul obiectelor precum și în interiorul obiectelor dar nu prezintă avantajul unui mecanism uniform de proiectare. Prezența a două tipuri de cooperare în sistem crește complexitatea și face dificilă dezvoltarea taskurilor de aplicație.

Submodelul *obiecte active secvențiale* prevede un obiect activ ca un singur proces secvențial, care tratează mesajele recepționate. După recepție, mesajele sunt înălțuite în interiorul obiectului și tratate secvențial, în conformitate cu o anumită strategie (în ordinea datei de sosire-deci FIFO, în ordinea datei de emisie sau folosind marca de timp propriu).

Capitolul 2.

COMUNICAȚII ÎN SISTEME DISTRIBUITE. PROTOCOALE.

Acest capitol prezintă probleme specifice ale comunicației în rețele de calculatoare, ale implementării software-ului de comunicație precum și aspecte specifice ale protocoalelor de comunicație în sisteme distribuite.

2.1. Conceptul de arhitectura deschisă. Stratificarea ierarhică a protocoalelor.

Deși modul de conectare fizică al calculatoarelor care comunică într-o rețea la mediul de transmisie ce constituie suportul rețelei poate diferi considerabil de la un tip de rețea la altul, se poate aprecia ca schimbul de informație între calculatoare în rețele poate fi reprezentat la modul general ca în figura 2.2.1

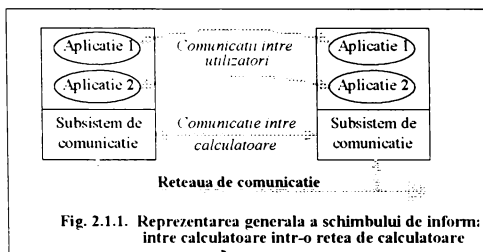


Fig. 2.1.1. Reprezentarea generală a schimbului de inform: între calculatoare într-o rețea de calculatoare

La baza întregului schimb de informație poate fi o rețea locală LAN, un număr de astfel de rețele interconectate sau o rețea WAN. Totuși, indiferent de tipul rețelei de comunicație a datelor, fiecărui calculator trebuie să i se atașeze un software de comunicație, denumit *sistem de comunicație*, pentru a controla stabilirea unor canale de comunicație între calculatoare și pentru a controla fluxul de mesaje în cadrul acestor canale. Implementarea și structurarea sistemului de comunicație este dependentă de tipul rețelei de comunicație. Pe de altă parte, acest software

trebuie integrat într-un sistem de operare cu API-uri specifice și apelat din limbaje de programare (cu caracteristici diferite în ceea ce privește reprezentarea datelor); ca atare sistemul de comunicație este dependent și de sistemul de operare folosit.

Proiectarea și implementarea sistemului de comunicație ridică probleme mult mai mari decât proiectarea sistemului de I/O al calculatoarelor într-un sistem de operare, datorate în special fiabilității mai reduse a mediului de comunicație și partajarea lui de către toate echipamentele legate în rețea. Dacă detaliile hardware pot fi separate într-un modul software driver, analog sistemelor de I/O, și parțial, în controller-ul de rețea (fig 2.1.2), trebuie însă rezolvate și alte probleme specifice precum *asigurarea unui serviciu de livrare a mesajelor sigur și dirijarea mesajelor*.

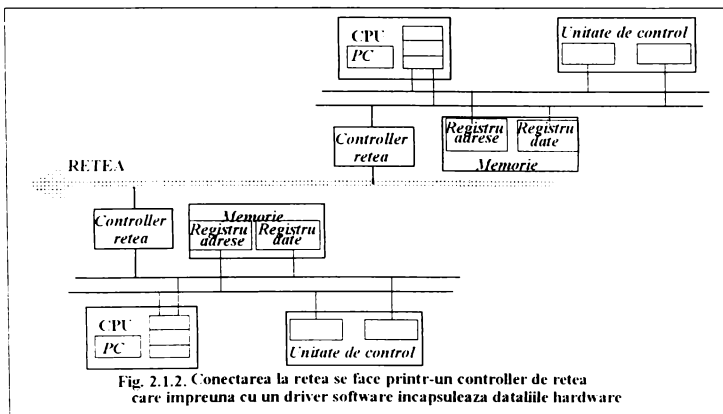


Fig. 2.1.2. Conectarea la rețea se face printr-un controller de rețea care împreună cu un driver software încapsulează detaliile hardware

Complexitatea sistemului de comunicație poate fi redusă prin separarea lui în două mari părți, fiecare la rândul ei fiind structurată pe nivele: o parte înglobează funcțiile care permit utilizarea mediului fizic de

comunicare și cealaltă parte cuprinde funcțiile legate de schimbul de mesaje între procese (fig. 2.1.3), oferind astfel o comunicare virtuală.

Numarul de nivele, numele fiecărui nivel și funcțiile aferente fiecărui nivel difera de la un tip de rețea

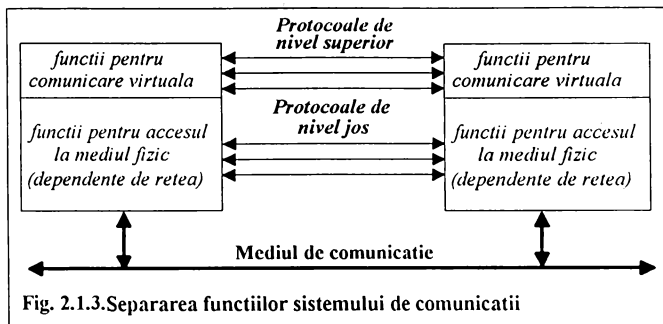


Fig. 2.1.3. Separarea funcțiilor sistemului de comunicații

la altul. Interfața dintre două nivele adiacente definește *operațiile primitive* și *serviciile* oferite de nivelul inferior celui superior. Comunicația între nivelele de pe aceleași ranguri este guvernată de un set de reguli, denumite *protocoale*; acestea controlează modul în care se desfășoară comunicația. Un *protocol* poate fi definit ca un set formal de reguli și convenții care guvernează formatul și sincronizarea mesajelor în cadrul unui schimb de mesaje între două sau mai multe procese. Proiectarea unui protocol necesită construirea scenariului de bază în care se definesc mesajele și secvențele de bază din cadrul unui schimb și se specifică procedurile pentru recuperarea erorilor; acestea introduc noi tipuri de mesaje și măresc lista de secvențe de mesaje acceptate într-un schimb.

Protocoalele trebuie să respecte *cerințe cantitative* și *calitative* relative la interconectarea proceselor. Cerințele calitative se referă la mecanismele de sincronizare, detecția și recuperarea erorilor, evitarea interblocărilor, evitarea depășirii capacității buffere-lor, asigurarea secvențierii mesajelor, detectarea mesajelor duplicate, mecanisme pentru asigurarea securității datelor, mecanisme pentru prioritate, mecanisme de contabilizare, garantarea transmiterii mesajelor, transformări de cod/format și compatibilitate cu caracteristicile sistemului de operare și ale mașinii fizice. Cerințele cantitative se referă la întârzieri de transmisie, productivitatea (cantitatea de date transferate în unitatea de timp), probabilitatea unui eșec în tratarea erorilor. Totalitatea nivelelor și protocoalelor asociate lor definesc *arhitectura rețelei*.

2.1.1. Modelul arhitectural ISO/OSI

Datorită diversificării echipamentelor de calcul, dar și a mediilor de comunicație (rețele locale, rețele de comunicare prin radio sau satelit, rețele digitale cu servicii integrate) și cerinței de interconectare a rețelelor și echipamentelor de proveniență diferită, în scopul alegerii unei variante optime pentru o aplicație, a devenit necesară elaborarea unor standarde internaționale care să guverneze schimbul de informație între calculatoare, rețele, terminale și chiar procese, creîndu-se astfel posibilitatea interconectării unor sisteme *deschise*. Ca atare, Organizația Internațională de Standarde (ISO) a propus un set de reguli în cadrul unui model de referință pentru interconectarea sistemelor deschise (ISO/OSI RM). Un sistem *deschis* este în accepțiunea ISO (SI) un sistem cu o arhitectură standardizată public, care poate să fie conectat și deci să comunice cu orice alt sistem deschis care utilizează același model de arhitectură.

Scopurile propunerii modelului ISO OSI au fost următoarele [Sta89]: (1)-posibilitatea interconectării sistemelor realizate de diversi producători; (2)- realizarea unui cadru conceptual coordonării și definirii de standarde pentru interconectarea sistemelor; (3)-identificarea domeniilor în care sunt necesare dezvoltări și îmbunătățiri ale standardelor; (4)-furnizarea unei referințe comune pentru menținerea consecvenței tuturor standardelor.

Setul de protocoale propus de ISO/OSI este stratificat ierarhic pe șapte nivele, fiecare nivel furnizându-i servicii nivelului superior, îmbunătățind într-un fel calitatea serviciilor oferite de nivelele inferioare (2.1.1.1.a, b). Separarea pe nivele a protocoalelor și serviciilor s-a realizat avîndu-se în vedere obiective precum: (1) un număr cit mai redus de nivele cu puține interacțiuni între ele, astfel încît modificarea funcțiilor unui nivel să nu afecteze celelalte module; (2) comasarea funcțiilor similare în același nivel și în concordanță cu necesitățile de standardizare sau standardele deja în folosință.

Terminologia OSI este ilustrată în fig.2.1.1.2. Într-un sistem există una sau mai multe entități în fiecare nivel. O *entitate* de la nivelul N (entitate (N)) implementează funcțiile nivelului (N) și protocolul de comunicație cu entitățile de la nivelul N din alte sisteme. Exemple de entități pot fi *procesele* (într-un sistem multitasking) și *procedurile*. Pot fi mai multe entități într-un nivel dacă acest lucru este eficient pentru un anumit sistem, eventual diferite, fiecare corespunzînd unui protocol standard de la acel nivel. Fiecare entitate

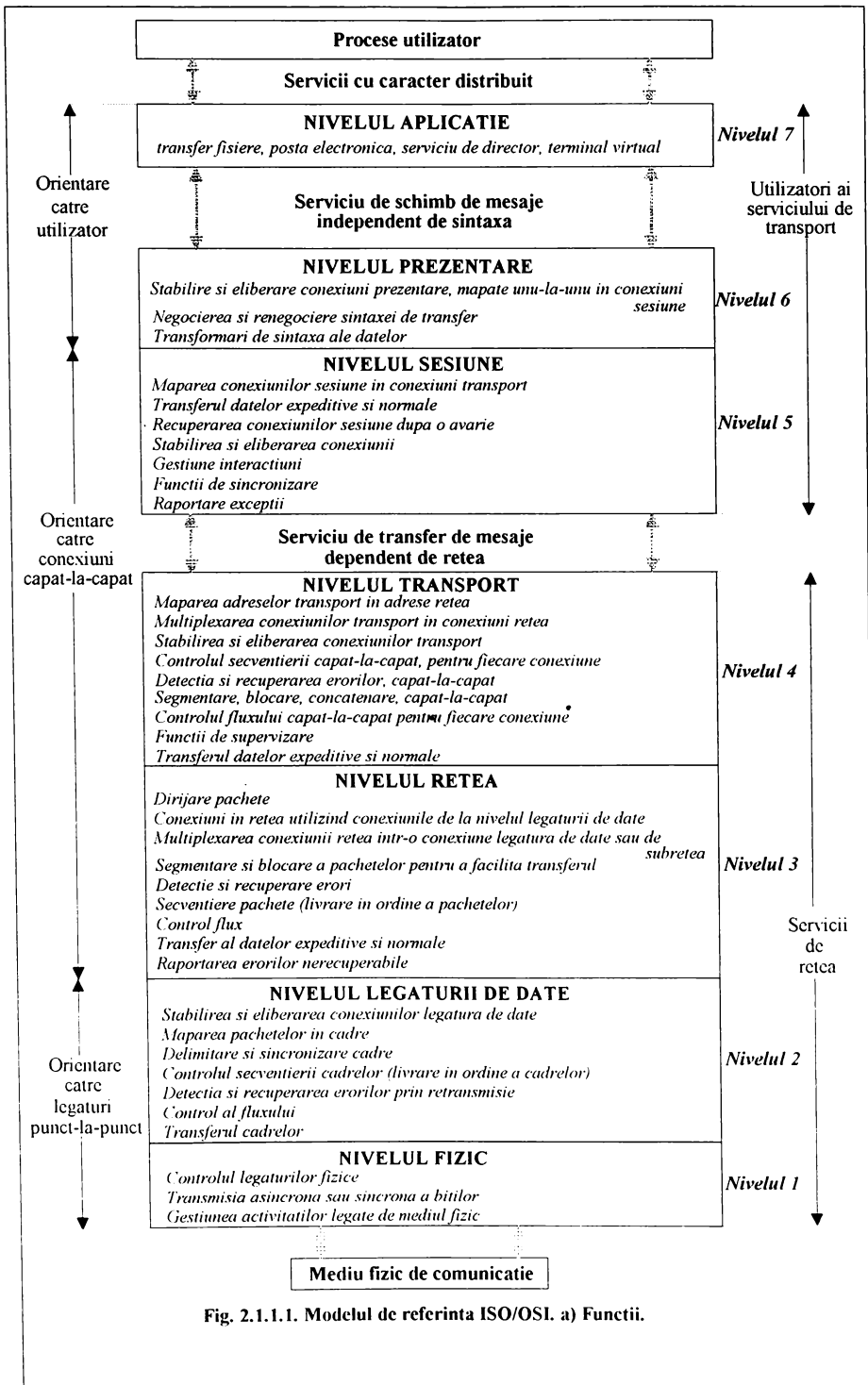


Fig. 2.1.1.1. Modelul de referinta ISO/OSI. a) Functii.

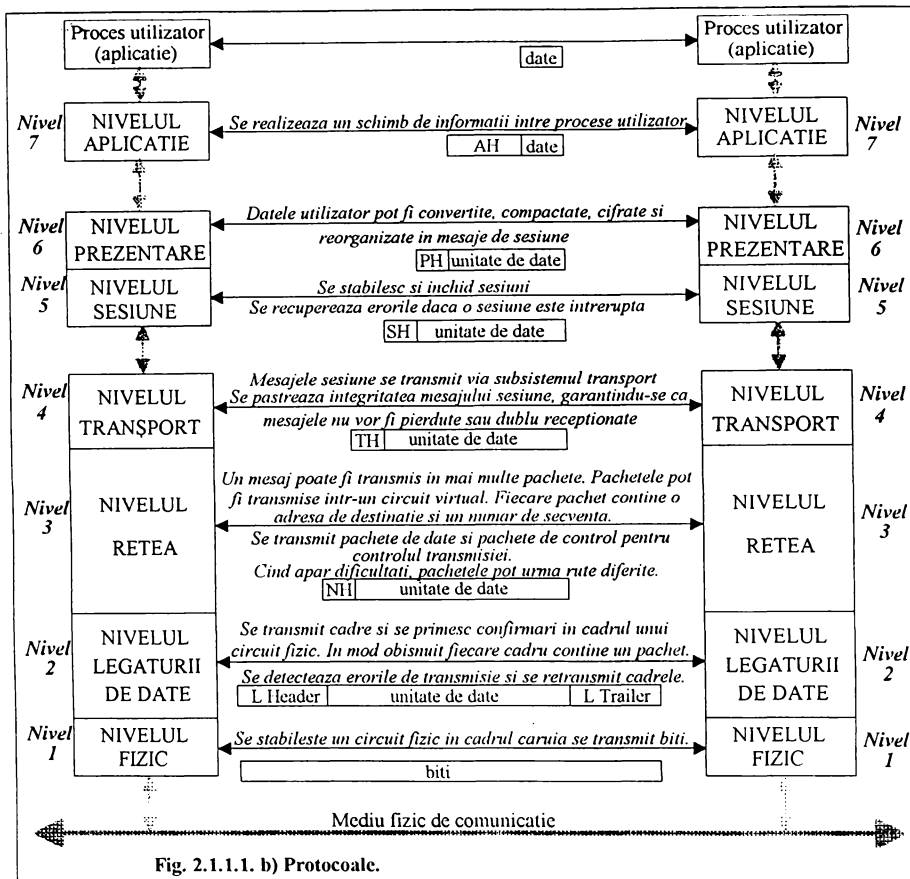


Fig. 2.1.1.1. b) Protocoale.

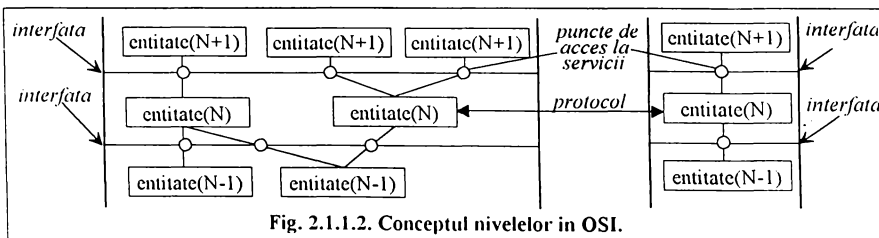


Fig. 2.1.1.2. Conceptul nivelor in OSI.

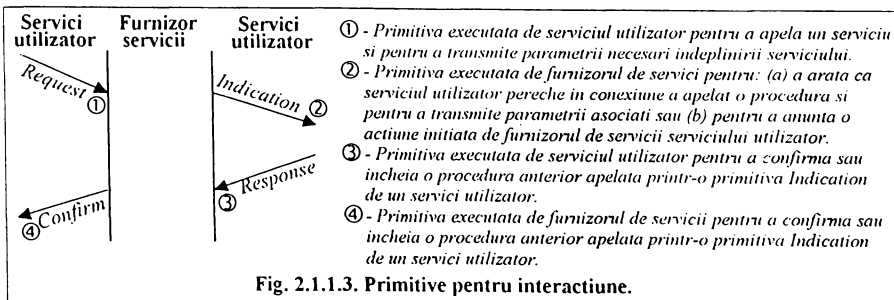


Fig. 2.1.1.3. Primitive pentru interactiune.

comunică cu nivelul superior printr-o interfață; aceasta este realizată printr-unul sau mai multe puncte de intrare la serviciu (Service Access Points=SAP). Entitățile de la un nivel N furnizează servicii entităților de la nivelul (N+1) printr-un set de primitive. O *primitivă* specifică funcția de executat și se utilizează pentru a transmite informații de control și date. Forma de realizare a primitivei este dependentă de implementare (un exemplu este un *apel de procedură*). Există 4 tipuri de primitive utilizate pentru a defini interacțiunea între nivele adiacente ale arhitecturii (fig.2.1.1.3).

Datele se transmit între entități pereche sub forma *unităților de date ale protocolului (PDU)*. O entitate (N) poate transmite o PDU la nivelul N-1, care o percepe ca pe o *unitate de date de serviciu (SDU)* ce trebuie sau nu mapată în mai multe PDU de la nivelul N-1. Dacă maparea este necesară, atunci entitatea pereche (N-1) trebuie să asambleze PDU-urile primite, înainte de a le livra entității (N). În fig.2.1.1.1 b se ilustrează principiile OSI relativ la protocoale și unitățile de date ale protocoalelor.

Arhitectura generală propusă pe șapte nivele de protocoale este un *model de referință* care descrie un ansamblu abstract de interdependențe, un cadru general în care pot fi definite și dezvoltate standarde, fără a impune restricții în ceea ce privește implementarea, topologia și tehnologia de interconectare.

O implementare a ierarhiei de protocoale ce pot fi considerate, de fapt, o mașină abstractă ce realizează operațiile de I/O cu rețeaua, este reprezentată în figura 2.1.1.4.

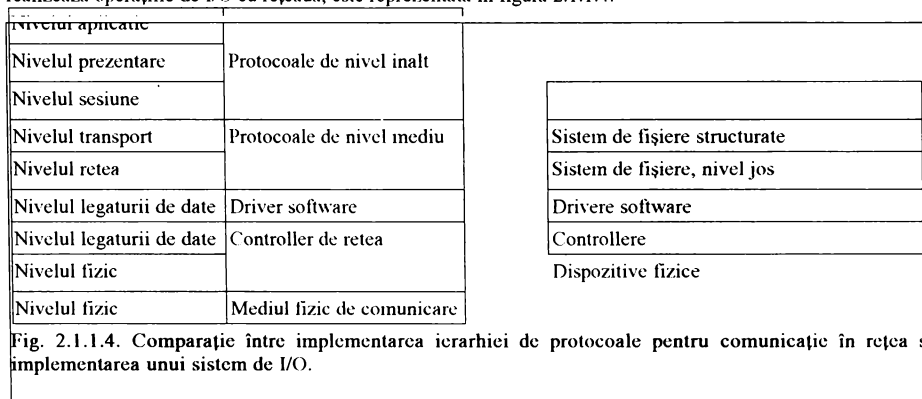


Fig. 2.1.1.4. Comparație între implementarea ierarhiei de protocoale pentru comunicație în rețea și implementarea unui sistem de I/O.

Nivelul legăturii de date este parțial implementat în controller-ul de rețea și parțial într-un driver software. Nivelul fizic este implementat în controller-ul de rețea și în mediul fizic de comunicare. Protocoalele de nivel mediu (nivelul de rețea și nivelul de transport) furnizează aplicațiilor independența de hardware, iar cele de nivel înalt extind serviciile puse la dispoziția programelor.

Făcînd o comparație între organizarea ierarhică tradițională a sistemului de I/O al calculatoarelor (drive-uri de disc, sisteme de fișiere, etc.) reprezentată în figura 2.1.1.4 și organizarea ierarhică a sistemului de comunicație se poate considera că nivelul nivel transport, uneori chiar nivelul rețea, corespund nivelului de fișiere, deși există și implementări în care, parțial sau total, în special nivelul rețea, este implementat în driver. Se poate spune că, analog sistemului de fișiere care furnizează o abstractizare a operațiilor de I/O, protocoalele la nivel mediu furnizează o abstractizare a operațiilor de I/O cu rețeaua de comunicație.

2.1.2. Nivele dependente de rețea. Prezentare generală

Funcțiile nivelurilor pot fi descrise pe scurt astfel:

A. **Nivelul fizic** definește modalitatea în care succesiunea de biți este transmisă pe mediul fizic de comunicație, respectînd ordinea biților și fără a garanta o transmisie sigură. Standardele de la acest nivel au în vedere aspecte electrice, funcționale, procedurale și mecanice ale transmisiei pe mediul fizic. Aspectele principale sunt legate de modul de codificare al biților, alegerea mediului fizic pentru transmisie (cablu torsadat, cablu coaxial, fibre optice, frecvențe radio), modul de stabilire și desființare a conexiunilor fizice, modul de transmisie (duplex, semiduplex, simplex), nivele electrice utilizate pentru transmisie, definirea pinilor și asignarea lor la conectori, numărul de biți pe secundă care pot fi transmiși. Standardele definite la acest nivel trebuie să asigure conectarea echipamentelor atât în rețele locale cît și în rețele municipale și teritoriale.

B. **Nivelul legăturii de date** are ca funcție principală transformarea canalului de comunicație furnizat de nivelul fizic, într-un canal de comunicare sigur, între două noduri adiacente. În acest scop, datele care urmează să fie transmise nivelului fizic, primite de la nivelul superior nivelului legăturii de date, sunt divizate în *cadre*, a caror structură generală este prezentată în figura 2.1.2.1. În general un cadru este compus din: (1) indicator, o succesiune specială de biți care delimitează începutul și sfîrșitul unui cadru (și care nu poate apare în interiorul datelor utile); (2) informații de control, care pot include tipul cadrului (se disting cadrele de

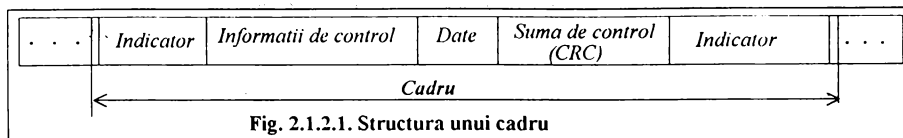


Fig. 2.1.2.1. Structura unui cadru

control de cadrele purtătoare de date utile), numărul de secvență (cadrele se numerează pentru ca receptorul să poată identifica eventualele cadre duplicate transmise ca urmare a unei erori de time-out), lungimea, adresa sursei și adresa destinației; (3) datele utile; (4) un câmp CRC care conține suma de control ciclic. Cadrele sunt transmise individual. La recepția lor, se verifică suma de control ciclic și, în cazul unei erori, se poate cere retransmisia cadrului.

O altă funcție a nivelului legăturii de date este *controlul fluxului datelor*, adică a ratei de transmisie a pachetelor dintre două noduri. Controlul fluxului este necesar deoarece: (1) recepția datelor se face în mod obișnuit într-un driver care lucrează cu ajutorul întreruperilor; dacă întreruperile sunt dezactivate în momentul primirii datelor, acestea pot fi pierdute; (2) driverul primește datele adresate oricărui proces local; ca atare păstrează datele în bufferele sale locale pînă cînd procesul destinat cerere sistemului de operare local să-i livreze un mesaj. Cînd nu mai există spațiu în buffere, driverul nu mai poate recepționa date fără a suprascrie alte date, nelivrate unui proces destinat. Pentru controlul fluxului se utilizează frecvent două soluții [Nut92]:

(1) *mecanismul confirmării cadrelor (protocolul stop-and-wait)*, în cadrul căruia transmiterea unui cadru se face numai după primirea unui mesaj de confirmare a cadrului precedent; în termenii comunicației interprocese în sistemele de operare, este analog unui apel send() sincron (fig. 2.1.2.2.). Sincronizarea între receptor și emițător se realizează printr-un cadru special de tip confirmare (ACK), care nu conține date utile. Pentru a preveni așteptarea la nesfârșit a unei confirmări "pierdute" sau a confirmării unui cadru "pierdut" (în ambele cazuri se presupune că nu a ajuns la destinație cadrul) se folosește un mecanism de time-out și se numerează cadrele.

```

Emitatorul transmite_cadrul();
Emitatorul seteaza un timer la valoarea TIMEOUT_pentru_a_urmari_transmisia();
Emitatorul asteapta o intrerupere_de_receptie_cadru_sau_de_time-out();
if(Emitatorul_a_primit_ACK) {
    Emitatorul reseteaza timerul();
    continue;
}
if(Timerul_a_generat_TIMEOUT)
    Retransmite_cadrul_care_a_produc_time-out();
    
```

(a) Secventa de cod executata pentru transmisie

```

Receptorul accepta un cadru neduplicat();
Receptorul transmite un cadru ACK();
    
```

Fig. 2.1.2.2. Controlul fluxului prin protocolul stop and wait

```

(1) if( Es-Ei<dimensiunea_maxima_a_ferestrei ) {
    Emitatorul transmite_cadrul(Es);
    Emitatorul seteaza un-timer la-valoarea TIMEOUT_pentru_a_urmari_transmisia_acestui_cadru();
    Incrementeaza_Es;
    }
(2) if( Emitatorul_a_primit_ACK ) {
    Emitatorul reseteaza timerul_pentru_cadrele_de_dinaintea_celui_pentru_care_a_primit_ACK();
    Incrementeaza_Ei;
    }
(3) if( un-timer_a_generat_TIMEOUT ) {
    Reseteaza_Es_Ei;
    Retransmite_toate_cadrele_de_dupa_cel_care-a_generat_time-out()
    }
    
```

(a) Secventa de cod executata la transmisie:

- 1) la primirea unui cadru pentru transmisie;
- 2) la generarea unei intreruperi de receptie de la driver;
- 3) la generarea unei intreruperi de time-out de un timer.

```

Receptorul accepta-un_cadru_neduplicat();
if( Ri = numarul_de_secventa_al_cadrului ) {
    Incrementeaza_Ri_mod_dimensiunea_maxima_a_ferestrei + 1();
    Receptorul transmite_confirmarea-ACK_pentru_acest_cadru();
} else Rejecteaza_cadrul();
    
```

(b) Secventa de cod executata la receptie

Fig. 2.1.2.3. Controlul fluxului prin protocolul cu fereastra glisanta.

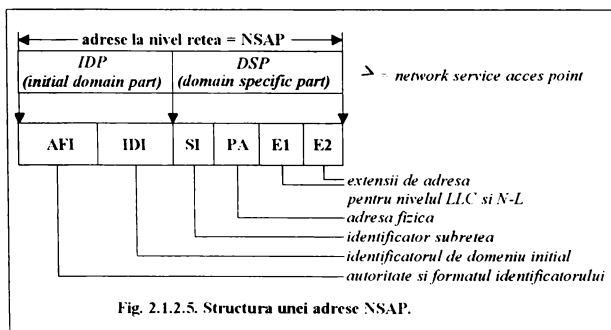
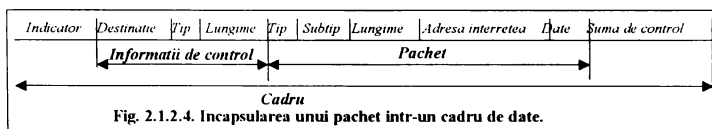
(2) *protocolul cu fereastra glisantă (sliding window protocol)* este o generalizare a celui de mai sus, în sensul că admite ca emițătorul să aibă N cadre în tranzit înainte de primirea unei confirmări. În acest caz, dacă se numerează cadrele, emițătorul va lucra cu doi pointeri, Ei și Es, care poartă la cadrul cu numărul

cel mai mic care a fost transmis, dar pentru care nu s-a primit confirmare, respectiv la primul cadru netransmis; astfel, toate cadrele cu numerele între E_i și E_s , care alcătuiesc fereastra curentă a emițătorului, sunt în tranzit. Protocolul permite "pierderea" confirmărilor unui cadru, deoarece confirmarea cadrului i provoacă automat confirmarea cadrelor $i-1, i-2, \dots, E_i$. La emisie, se va folosi pentru fiecare cadru neconfirmat câte un timer. Emitatorul va putea avansa cu maximum N cadre fără a primi o confirmare; când se atinge limita ferestrei, se întrerupe transmisia cadrelor pînă la primirea unei confirmări sau a generării unui semnal de time-out.

Există și o altă variantă a protocolului, în care, în cazul apariției unei erori time-out nu este necesar să se retransmită toate mesajele din bufferul emițătorului, ci numai mesajul nerecepționat corect (sau neconfirmat); aceasta presupune că și receptorul să lucreze cu o fereastră de recepție. Trebuie însă să fie mai mică decît numărul maxim de secvență al cadrelor (N); în general are valoarea $(N+1)/2$.

C. *Nivelul rețea* are drept scop final crearea posibilității de distribuție a informației în rețele logice, în care adresele să fie generalizate, fără a depinde de spațiile de adrese specifice diferitelor tipuri de protocoale de la nivelul legăturii de date. Pentru a realiza acest scop, la nivelul rețea trebuie să se realizeze interconectarea rețelelor cu arhitecturi diferite și dirijarea pachetelor de date între nodul sursa și cel destinație, indiferent dacă acestea sunt sau nu adiacente (se vor folosi noduri intermediare).

Adresele din pachete sunt interpretate de protocoalele de la nivelul rețea; acestea utilizează un spațiu de adrese diferit de cel utilizat de cadrele nivelului legăturii de date; un cadru trebuie să specifice ca adresă destinație numai adresa unui nod adiacent, pe cînd, la primirea unui pachet, nivelul rețea poate să transmită mai departe pachetul, încapsulîndu-l într-un cadru (fig. 2.1.2.4), permițînd crearea unui efect logic de conectare a tuturor nodurilor mai multor rețele interconectate. Rezulta că nivelul rețea trebuie să furnizeze un mecanism uniform de adresare care să aibă în vedere diversitatea rețelelor de comunicație și faptul că acestea pot fi distribuite geografic pe țări, continente și aflate în proprietatea sau gestiunea unor autorități diferite.



Pentru adresare la nivelul rețea, ISO a propus o structură ierarhizată a adreselor, reprezentată în figura 2.1.2.5. Fiecare adresă este compusă din două părți: *partea de domeniu inițial* ("Initial domain part" IDP) și *partea specifică domeniului* ("domain specific part"- DSP). Pentru a reflecta numeroasele alternative de interconectare a rețelelor de diferite tipuri, fiecare cu propria sa schema de adresare, cîmpul IDP este rîndul lui împărțit în două părți: *autoritatea și identificatorul formatului* (AFI) și *identificatorul formatului inițial* (IDI). Cîmpul AFI specifică atât modul de interconectare al rețelelor, dar și formatul cîmpului IDI asociat; de exemplu, dacă interconectarea se face prin rețele internaționale cu comutare de pachete, IDI identifică adresa porții locale cu care subrețeaua este conectată și va cuprinde o secvență de maximum 14 cifre alocate în conformitate cu recomandarea CCITT X.121, iar cîmpul AFI va conține valoarea zecimală 36 dacă adresa IDI este reprezentată în forma BCD sau valoarea zecimală 37, dacă IDI este reprezentată în forma binară. În tabelul din figura 2.1.2.6 se arată și alte valori pentru AFI. Pentru toate aceste valori, adresa IDI este compusă dintr-un număr de alte subadrese ierarhice care reprezintă codul țării și codul național al rețelei. Cîmpul DSP cuprinde, analog, un număr de subadrese ierarhice: identificatorul subrețelei (SI), care identifică subrețeaua particulară (I.AN), la care este conectată stația destinație, adresa fizică (PA) care precizează adresa fizică a subrețelei identificată de SI.

Formatul IDI	Valoare AFI		Lungime IDI (cifre)	Lungime DSP (cifre zecimale)
	Zecimal	Binar		
Rețele internaționale cu comutare de pachete X.121	36	37	14	24
Rețele telex internaționale F.69	40	41	8	30
Rețele internaționale telefonice comutate E.164	42	43	12	26
Rețele internaționale digitale cu servicii integrate	44	45	15	23
	48	49	Nula	38

Fig. 2.1.2.6. Valori pentru cîmpurile unei adrese NSAP

Dirijarea pachetelor este una dintre funcțiile esențiale ale nivelului rețea care, indiferent de algoritmul folosit pentru dirijare, se poate rezuma la următoarele operații: (1) nodul emițător trebuie să utilizeze o tabela de dirijare (locală), cu ajutorul căreia să poată selecta un prim nod intermediar în calea către destinație; (2) nodul emițător trimite pachetul încapsulat într-un cadru nodului intermediar; (3) nodul intermediar determină dacă pachetul trebuie dirijat în continuare către destinație, utilizînd propria sa tabela de rutare sau a ajuns la destinație. Algoritmii de dirijare a pachetelor, care trebuie să satisfacă cerințe ([Tan88] [And94]) de simplitate, robustețe, stabilitate, optimalitate, nepărtinire (algoritmul nu trebuie să favorizeze o categorie de utilizatori în defavoarea celorlalți, cu excepția cazului cînd se atribuie prioritățile), pot fi clasificați, după reacția la modificările traficului sau topologiile rețelei, în algoritmi *adaptivi și neadaptivi*. Din clasa algoritmilor neadaptivi, care realizează o dirijare statică se enumeră *dirijarea fixă* (în variantele dirijarea cea mai scurtă și dirijarea multicai), *dirijarea în flux, dirijarea incidentală și dirijarea ierarhică*, iar din clasa algoritmilor adaptivi: *dirijarea centralizată, dirijarea izolată* (algoritmi "hotpotato" și "backward learning"), *dirijarea delta și dirijarea distribuită*. Pentru rețelele de largă întindere geografică este importantă dirijarea ierarhică.

Indiferent de algoritmul de dirijare utilizat, trebuie evitată degradarea performanțelor rețelei, ca urmare a supraîncărcării ei. În acest scop, o altă funcție a nivelului rețea este *evitarea congestionării rețelei*, ca poate apare datorită dispariției unei legături în rețea, a tranziției la noi tabele de dirijare sau a sufocării unor noduri utilizate pentru a calcula căi de dirijare. Există două categorii de soluții pentru a evita congestionarea rețelei: (1) evitarea supraîncărcării fiecărui nod și (2) menținerea traficului în jurul unei nivel fixat pe liniile de legătură dintre două procese utilizator (pe calea unui circuit virtual), pe toate liniile de legătură din jurul unui nod sau pe liniile de legătură dintre noduri și comutatoarele de pachete. Din algoritmi de categoria (1), de control a supraîncărcării fiecărui nod, se enumeră ([For88][Tan88]): prealocarea zonelor tampon în cadrul unui circuit virtual, ignorarea pachetelor pentru care nu există spațiu pentru memorare și utilizarea unor pachete de permisiune (în număr limitat în rețea) pentru a putea transmite pachete de date. Dintre algoritmi de categoria (2), de control al fluxului de mesaje, se specifică doi algoritmi: transmiterea de către comutatoarele de pachete a unor pachete de șoc nodurilor sursă, cu scopul de a le impune micșorarea ratei de transfer a pachetelor pe liniile din jurul comutatoarelor de pachete (care își monitorizează activitatea liniilor de ieșire) și impunerea transmisiei unui pachet numai după primirea confirmării recepției celui transmis anterior.

O altă situație care trebuie luată în considerare la nivelul rețea este *blocarea definitivă a nodurilor*: aceasta poate apare în diverse situații, cele mai posibile fiind: (1) cînd două comutatoare de pachete trebuie să-și transmită reciproc cite un pachet, dar ambele au toate bufferele ocupate sau (2) cînd un comutator de pachete, care folosește la nivelul legătură de date un protocol de fereastră glisantă, varianta selectivă. (în care se memorează în buffere mesajele recepționate corect, în limita ferestrei, chiar dacă s-a pierdut un mesaj anterior) iar pentru alocarea bufferelor sale locale la circuitele virtuale folosește o strategie FIFO globală, are toate bufferele locale ocupate cu pachete provenind, pentru un anumit circuit virtual, de la același mesaj, cu numerele de secvență mai mari decît ale altor pachete așteptate pentru care nu mai există loc în buffere. Dintre soluțiile utilizate pentru evitarea blocării definitive se enumeră: (1) utilizarea în fiecare nod a M + 1 buffere, M fiind lungimea maximă a căilor rețelei; (2) asocierea unei informații de vechime fiecărui pachet; (3) transmiterea pachetelor obținute prin fragmentarea unui mesaj numai după alocarea unui număr suficient de buffere în care să se poată memora pachetele mesajului multipachet.

Unul dintre cele mai discutate aspecte ale proiectării nivelului rețea a fost *tipul serviciilor care trebuie furnizate de acest nivel nivelului sau superior, cel de transport: servicii orientate sau neorientate pe conexiune*. Problema poate fi extinsă și în cazul rețelelor locale: ce tip de servicii trebuie să ofere nivelului transport nivel legăturii de date (eventual, prin intermediul nivelului de rețea, dacă acesta există). Pe de altă parte, deoarece arhitectura modelului OSI pentru subrețeaua de comunicații prevede ca ultim nivel nivel rețea, rezultă că serviciul oferit de acesta este serviciul oferit de subrețea în general. Ca atare, ISO a hotărît ca modelul OSI să ofere la nivelul rețea ambele tipuri de servicii: (1) servicii orientate pe conexiune; (2) servicii neorientate pe conexiune.

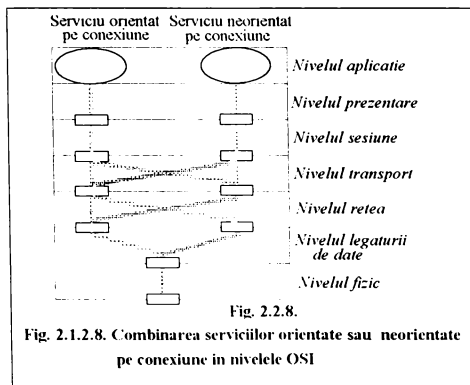
Tipul (1) se poate caracteriza astfel [Tan88][Gos91]: (1) Între 2 entități care comunică, la cererea sursei, se stabilește mai întâi un canal logic sau o conexiune. Această conexiune primește un identificator care rămâne valabil până la terminarea datelor de transmis și desființarea explicită a conexiunii, de către oricare dintre cele două entități. (2) La stabilirea unei conexiuni, cele două entități de la nivelul transport și nivelul rețea care furnizează serviciul stabilesc *parametrii* (adresele sursei și destinației, și indicatori pentru utilizarea sau nu a confirmărilor și a datelor expeditive) și *calitatea serviciilor* în cursul unui proces de negociere. (3) Comunicația este permisă în ambele sensuri, păstrându-se ordinea pachetelor, garantându-se transmisia și controlul fluxului datelor astfel încât comunicația este posibilă și între calculatoare de performanțe diferite. Pentru tipul (1) de servicii sunt prevăzute primitive pentru stabilirea conexiunii pentru trimiterea și recepționarea datelor și pentru desființarea conexiunii.

În contrast, *serviciile neorientate pe conexiune* transmit individual pachete care, nefiind integrate într-o conexiune, trebuie să încorporeze, întotdeauna, adrese destinație; nu se garantează sosirea în secvență a pachetelor. Pentru acest tip de servicii sunt prevăzute primitive numai pentru transmiterea și recepționarea datelor.

Aspect	Stabilire conexiune	Adresa destinație	Secvențiere pachete	Control erori	Control flux	Este posibilă negocierea de opțiuni	Se utilizează identificatori de conexiune
Serviciu orientat pe conexiune	Necesar	Numai la stabilirea conexiunii	Garantată	Realizat la nivelul rețea	Realizat la nivelul rețea	Da	Da
Serviciu neorientat pe conexiune	Imposibil	La orice pachet	Negarantată	Realizat la nivelul transport	Nu este furnizat de nivelul rețea	Nu	Nu

Fig. 2.1.2.7. Principalele diferențe între serviciile orientate pe conexiune și cele neorientate pe conexiune

În figura 2.1.2.8. se prezintă repartizarea celor două tipuri de servicii (orientate și neorientate pe conexiune) pe toate nivelele OSI. Se observă că este posibil să se dispună de servicii orientate pe conexiune la nivelul rețea sau transport, chiar când nivelele inferioare nu sunt orientate pe conexiune. De exemplu, un serviciu orientat pe conexiune la nivelul transport poate fi construit într-o rețea locală pe baza unui protocol neorientat pe conexiune tratând controlul erorilor, controlul fluxului, stabilirea conexiunilor la nivelul rețea sau transport.



Ambele tipuri de servicii au avantaje și dezavantaje. Astfel, serviciile neorientate pe conexiune sunt mai rapide (nu este necesară stabilirea conexiunilor), au protocoale simple, dar și dezavantajele că nu asigură controlul fluxului și livrarea în secvență a pachetelor, erorile (pachete pierdute sau duplicate) trebuie tratate la nivelul superior, iar adresa destinație încarcă toate pachetele. Pe de altă parte, deși dezavantajele de mai sus sunt înfăurate la serviciile orientate pe conexiune, acestea sunt puțin eficiente, datorită

regiei de sistem sporite.

În cazul rețelilor cu comutarea pachetelor, organizarea subrețelei de comunicație poate fi făcută folosind *tehnica rețelilor virtuale* sau a *datagramelor*. *Tehnica circuitelor virtuale* se utilizează, în general, în legătură cu serviciul orientat pe conexiune, deoarece presupune dirijarea tuturor pachetelor dintre un nod sursă și un nod destinație pe aceeași cale, fixată inițial printr-un pachet de stabilire a circuitului; de aceea, pe de o parte, în antetul pachetului se precizează numărul circuitului și nu adresa completă a destinatarului și, pe de altă parte, fiecare nod de comunicație păstrează un tabel cu circuitele virtuale care îl pareurg.

Tehnica datagramelor se utilizează în legătură cu serviciul neorientat pe conexiune, deoarece presupune dirijarea pachetelor individual, independent de calea urmată de pachetele anterioare între o aceeași sursă și o aceeași destinație; în acest caz, fiecare pachet trebuie să conțină în header adresa completă a destinației, irosindu-se mai mult spațiu decât numărul circuitului virtual din cealaltă tehnică.

Deși nu se conformează modelului OSI, este posibil ca unele aplicații să fie construite direct deasupra nivelului rețea; în acest caz, ele percep rețeaua ca un mediu de comunicație nesigur, și de aceea trebuie să trateze erorile (de pierdere sau duplicare a pachetelor) și să controleze fluxul de date și primirea în secvență a pachetelor. Există categorii de aplicații (de exemplu, aplicațiile legate de voce), pentru care este preferat serviciul neorientat pe conexiune, fiind mai rapid, chiar dacă nu garantează 100% transmisia fără erori a pachetelor; în acest caz, pierderea unui mesaj nu este catastrofă (destinația se poate dispensa de datele conținute într-un pachet pierdut sau poate să le determine printr-un mecanism de interpolare).

O altă funcție a nivelului rețea, în concepția ISO, este *interconectarea rețelelor* (LAN-LAN, LAN-WAN, WAN-WAN, LAN-WAN-LAN); aceasta impune următoarele cerințe: (1) suport pentru realizarea legăturilor fizice între rețele conectate; este necesar, în acest caz, la minimum, nivelul fizic și nivelul legăturii de date; (2) suport pentru dirijarea pachetelor între procese aflate în noduri aparținând unor rețele de tipuri diferite; (3) suport pentru un serviciu de monitorizare a utilizării rețelelor interconectate și a dispozitivelor de interconectare (porți) și pentru memorarea unor informații de stare; (4) satisfacerea cerințelor de mai sus în condițiile în care nu sunt permise modificări ale arhitecturii modelului OSI sau ale rețelelor componente. Aceasta înseamnă că interconectarea rețelelor trebuie să rezolve problema diferențelor multiple dintre rețelele conectate, precum: scheme de adresare diferite, diferențe între dimensiunea pachetelor (divizarea unui pachet mai mare provenind dintr-o rețea în pachete mai mici pentru a putea fi trimise în altă rețea este numită *segmentare*), mecanisme de acces la rețea diferite, valori diferite pentru timeout, tehnici de dirijare diferite; parametrii pentru raportarea stării și performanțelor diferiți, utilizarea circuitelor virtuale într-o rețea și a datagramelor în alta. Complexitatea soluțiilor de interconectare a rețelelor, în general (nu neaparat OSI) este prezentată în figura 2.1.2.9.

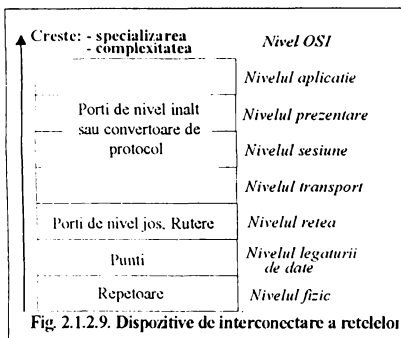


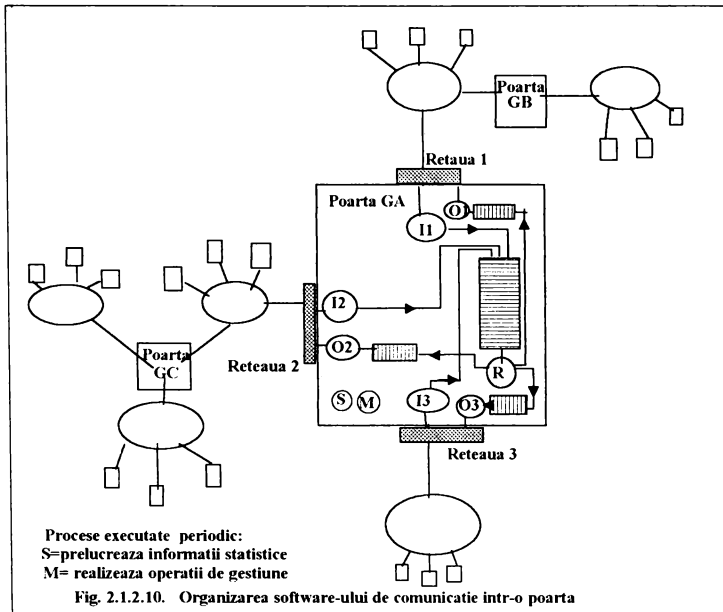
Fig. 2.1.2.9. Dispozitive de interconectare a rețelelor

În concepția OSI ([Sta89][Tan88]), pentru a putea realiza o funcție de interconectare a rețelelor, acolo unde este nevoie, nivelul rețea este împărțit în 3 subnivele: (a) subnivel de acces la rețeaua de comunicații, cu rolul de a realiza funcțiile nivelului rețea în conformitate cu protocolul subrețelei conectate; (b) nivelul de îmbogățire a rețelei, cu rolul de a uniformiza deosebirile dintre subrețelele conectate; (c) nivelul de interconectare a rețelelor, cu rolul de a realiza funcțiile nivelului rețea în subrețeaua de comunicație (în special dirijarea pachetelor - cu cele două variante, datagrame sau circuit virtual - preferându-se, de obicei un algoritmul de dirijare ierarhică).

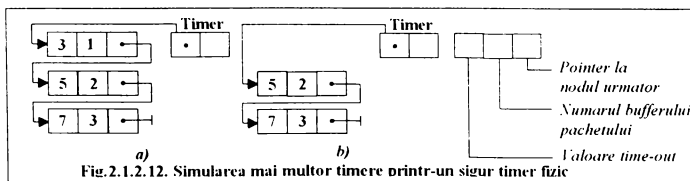
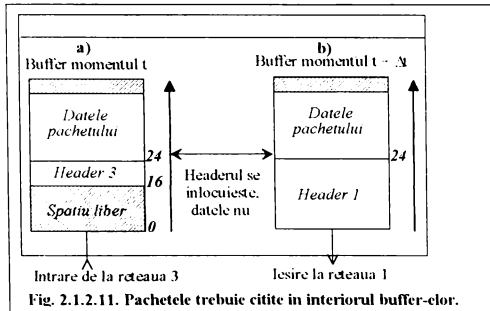
Proiectarea *software-ului de comunicație al punctelor*, care conectează rețelele locale, ca și al porturilor, trebuie să primească o atenție specială datorită restricțiilor de timp

care se pun. Deși rapidă, varianta în care software-ul este organizat pe *întreruperi (interrupt driven)* - fiecare recepție a unui pachet provoacă declansarea unei întreruperi -, fiind nestructurată, este mai puțin preferabilă variantei în care software-ul este organizat pe *procese* sau și mai bine pe *fire*, ținând seama că sistemele de operare actuale permit lucrul cu fire. În figura 2.1.2.10. se prezintă organizarea software-ului unei porți GA care conectează trei rețele. În acest caz, software-ul de comunicație este organizat în 9 procese. Procesele de intrare 11, 12, 13, având prioritatea cea mai mare, pentru a nu se pierde pachete, sunt executate la recepția unui pachet: ele au rolul de a verifica suma de control a pachetului și de a-l depune într-o coadă de mesaje asociată procesului R, procesul care realizează dirijarea pachetelor. Procesul R este executat atunci când nu este activ nici un proces de prioritate mare (11, 12, 13), sau de prioritate medie (01, 02, 03) și când coada de mesaje asociată nu este vidă. În acest caz, rolul său este de a prelua primul pachet din coada, de a determina calea pe care trebuie dirijat, de a-l depune în coada corespunzătoare unuia dintre procesele de ieșire (01, 02, 03). În momentul activării, un proces de ieșire preia din coada sa un pachet și inițiază transmisia pachetului, după care se autoblochează; este reactivat la declansarea întreruperii de sfârșit de transmisie, moment în care încearcă transmisia următorului de pachet din coadă, s.a.m.d. Procesul R are și rolul de a filtra pachetele: astfel, pachetele de control nu sunt transmise mai departe. Procesele S și M, de prelucrare statistică și de gestiune sunt executate periodic, cu cea mai mică prioritate.

Software-ul de comunicație expus mai sus, bazat pe conceptul de proces este structurat, dar cu o eficiență relativ scăzută. Pentru a-i crește eficiența, pot fi aduse următoarele îmbunătățiri: (1). Cozile de pachete (structuri de tip FIFO), prin care comunică procesele trebuie implementate, dacă sistemul de operare permite, în zone de memorie partajate de toate procesele care trebuie să aibă acces la ele: se alege în acest fel copierea mesajelor dintr-un spațiu de adrese al unui proces în spațiul de adrese al altui proces. (2). În locul proceselor, se pot folosi fire; acestea prezintă avantajul că nu se iesește timp procesor pentru comutarea proceselor și partajează structurile de date. (3). Se evită complicațiile care pot apărea datorită diferenței dintre dimensiunea pachetelor (și header-elor) în rețele diferite. O soluție este de a se recepționa pachetele în



buffere, începînd nu de la adresa de start, ci de la o oarecare distanță, astfel încît bufferul să poată conține și cel mai mare header, în locul celui recepționat. Figura 2.1.2.11 prezintă un exemplu în care se presupune că cele trei rețele au dimensiunea header-elor de 24 octeți, respectiv 16 și 8; în acest caz, pachetele din rețeaua 3 vor fi recepționate în buffer cu un deplasament de 16 octeți față de adresa sa de start (fig. 2.1.2.10.a), astfel încît, dacă pachetul va trebui transmis în rețeaua 1, cu un header de 24 octeți, va exista loc în buffer pentru construcția unui header mai mare (fig. 2.1.2.10.b).



listei la setarea și resetarea unui timer. (5) În unele situații, de exemplu, în transmisia între rețele 802.3 și 802.5 este necesară inversarea ordinii biților într-un octet. Cea mai simplă soluție este cea hardware, prin conectarea bitului de intrare 0 la bitul de ieșire 7, bitul 1 la bitul 6, s.a.m.d., dar dacă acest lucru nu este posibil se poate folosi o tabelă cu 256 intrări, indexată după valorile unui octet (a i-a intrare reprezintă octetul de valoare i avînd inversată ordinea biților).

2.1.3. Nivelul transport, principalul furnizor de servicii pentru sistemele de operare

Se apreciază în [Tan88][Sta89][Mul89][Nut92] că nivelul transport, nivelul superior al protocoalelor de nivel mediu este centrul de greutate în ierarhia de protocoale. Este considerat interfața dintre nivelele dependente de aplicație și nivelele dependente de rețea, cerințele de deasupra acestui nivel focalizându-se pe cerințe specifice utilizatorilor și aplicațiilor. Scopul său este de a furniza nivelului superior un serviciu de transfer sigur al datelor (fără date eronate, duplicate sau pierdute și respectând ordinea în care au fost emise), indiferent de tipul rețelei; el maschează astfel detaliile de implementare ale operațiilor de comunicație în rețea punând la dispoziția nivelului sesiune o interfață relativ simplă care înglobează toate facilitățile de transfer de mesaje. Fara nivelul transport nu ar avea sens toată ierarhia de nivele și protocoale deoarece scopul final al unei arhitecturi stratificate este de a furniza utilizatorului (și aplicațiilor) servicii de transfer de mesaje independente de tipul rețelei utilizate, sub forma unui set standard de primitive.

Un aspect important legat de acest lucru este faptul că, de foarte multe ori, utilizatorii preferă să-și scrie aplicațiile folosind direct primitive ale nivelului transport, nefiindu-le necesare serviciile nivelelor de deasupra nivelului transport. Un exemplu sugestiv este cazul sistemului de operare Unix care nu implementează, pentru lucrul în rețele Unix, nivelele de deasupra nivelului de transport. Aplicațiile Unix folosesc operațiile de I/O cu rețeaua, analog operațiilor de lucru cu fișiere, scriind și citind dintr-o rețea ca într-o conductă *pipe* dar folosind protocolul nivelului de transport; acest mod de lucru a fost ulterior adoptat și de Microsoft în Windows NT și Windows 95; mai mult chiar, se poate lucra cu structuri de tip *remote pipe* folosind protocolul de transport.

Un alt aspect important este faptul că nivelul transport este implementat de obicei în interiorul sistemelor de operare, spre deosebire de nivelele inferioare lui care sunt implementate sub forma de drivere software.

Un ultim aspect care trebuie reliefat este faptul că studierea nivelului de transport OSI este utilă și pentru situația în care trebuie proiectate aplicații care folosesc direct nivelul rețea (este cazul rețelelor locale) deoarece pot fi găsite soluții pentru problemele care trebuie rezolvate în acest caz la nivelul aplicației (controlul fluxului datelor în special).

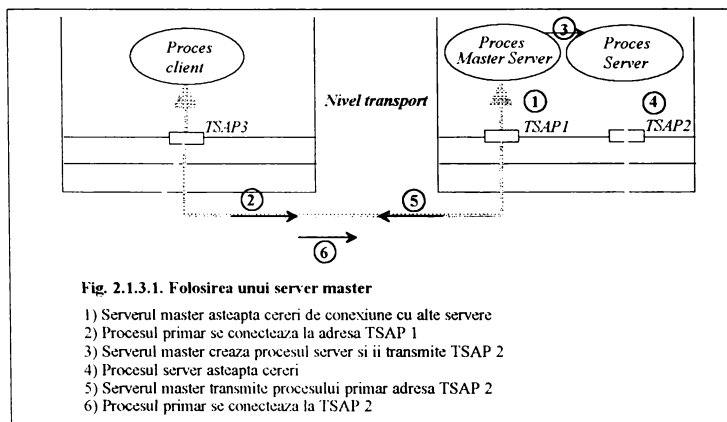
Funcțiile nivelului de transport sunt legate de adresare, realizarea transferului datelor, segmentarea mesajelor primite la nivelul sesiune, respectarea ordinii de emisie a mesajelor și la recepție, controlul erorilor la recepționarea mesajelor, controlul fluxului mesajelor capăt la capăt, stabilirea și desființarea conexiunilor și multiplexarea conexiunilor. Ca și la nivelul rețea, sunt posibile două tipuri de servicii: *orientate și neorientate pe conexiune*, dar sunt caracteristice serviciile orientate pe conexiune. În ambele tipuri de servicii, utilizatorul nivelului transport trebuie să identifice utilizatorul nivelului transport pereche printr-o adresă de nivel transport (*adresa TSAP = transport service acces point*)

Un utilizator al nivelului transport înseamnă un proces în cadrul unui sistem de operare (se consideră că nodurile rețelei sunt gestionate de sisteme de operare multitask), și deci comunicația asigurată de nivelul transport este o comunicație între două procese: un proces primar care apelează un proces secundar dintr-un alt nod. De aceea, la nivelul transport se specifică o adresă ierarhică formată din adresa de rețea, adresa de nod și un port de comunicație pentru proces (rețea, nod, port). Modul în care procesul se conectează la port este specific sistemului de operare local (un exemplu este sistemul Unix, în care portul de comunicație corespunzător unui proces este stabilit printr-un *socket*). Adresa de rețea și adresa de nod vor fi transmise mai departe nivelului rețea, sub forma unei *adrese NSAP* care identifică unic nodul pereche iar *portul* va fi transmis ca informație în cadrul pachetului de date (la stabilirea unei conexiuni pentru serviciile orientate pe conexiune și respectiv în fiecare pachet de date pentru serviciile neorientate pe conexiune). Pentru determinarea portului de comunicație corespunzător procesului secundar (care fiind apelat pentru îndeplinirea unui serviciu se numește *server*), un proces primar (numit *client*, fiind apelant al unui serviciu) poate folosi următoarele metode [Sta89][Tan88]:

1. Procesul primar cunoaște adresa completă corespunzătoare procesului secundar; aceasta este fie (a) o caracteristică de sistem, fie (b) este o adresă făcută publică în rețea (*well-known-adresse*) ca aparținând unui proces server ale cărui servicii sunt uzuale. Pentru cazul (a) poate fi dat următorul exemplu: în fiecare nod al unei rețele există câte un proces care colectează informații statistice; periodic, un proces master, situat pe un anumit nod, centralizează aceste informații conectându-se la toate procesele de colectare din celelalte noduri. Porturile de comunicare ale acestora nu trebuie să fie publice, dar pot fi fixate și cunoscute de procesul master. În cazul (b) pot fi incluse toate serverele care oferă servicii uzuale într-o rețea, transfer de fișiere (ftp, deschiderea de sesiune de lucru de la distanță-Telnet).

2. În ambele cazuri, 1.a și 1.b procesul secundar trebuie să existe, să fie activ și să asculte la adresa sa TSAP fixată, în vederea primirii unei cereri. În cazul (b) când numărul de servere cu adrese publice în rețea devine prea mare este ineficient să se păstreze active aceste servere. În acest caz poate fi folosită o altă

metoda: procesele server sunt create și activate numai la primirea unei cereri de către un proces server master. Acest server master va fi activ în orice nod în care pot fi activate servere pentru servicii uzuale; el va aștepta cereri de conectare la aceste servere la o adresă TSAP1 făcută publică. Când va primi o cerere pentru un server, va crea procesul server, transmitându-i o adresă TSAP2 la care aceasta să se conecteze pentru a prelua cererile și va transmite această adresă clientului, după care se va deconecta de la client, așteptând alte cereri de conectare. După stabilirea în acest mod a conexiunii dintre client și server, serverul execută cererile primite de la client iar la desființarea conexiunii de către client se autodistruge. Succesiunea de operații de mai sus este prezentată și în figura 2.1.3.1.



3. O altă metodă de a furniza adrese TSAP este folosind un *server de nume (name server sau directory server)*, a cărui adresă TSAP este publică în rețea. Pentru a afla adresa TSAP corespunzătoare unui server, orice client stabilește întâi o conexiune cu serverul de nume, furnizează un nume generic global al serverului, primește adresa TSAP a acestuia și se deconectează de la serverul de nume, conectându-se în continuare la adresa TSAP determinată. Serverul de nume trebuie să mențină o tabelă cu corespondența nume server - adresă TSAP; orice server nou creat trebuie înregistrat de serverul de nume.

Din punctul de vedere al unui proiectant al nivelului transport, prezintă importanța calitatea serviciilor de rețea care îi pot fi oferite deoarece de aceasta va depinde complexitatea sa; cu cât serviciile de rețea care depind de tipul rețelelor și ale subrețelelor de comunicații sunt mai sigure, cu atât complexitatea nivelului de transport va fi mai mică. În acest sens, serviciile de rețea pot fi clasificate astfel: (A). Rețele cu rata erorilor nesemnaltate acceptabilă și cu rata erorilor semnalate, de asemenea acceptabilă (se considera că de acest tip de rețea se apropie rețelele locale LAN) [Tan88][Sta89]; (B). Rețele cu rata erorilor nesemnaltate acceptabilă dar cu rata erorilor semnalate inacceptabilă (în care se încadrează multe rețele WAN, de exemplu rețele publice X.25); (C). Rețele cu rata erorilor nesemnaltate și semnalate inacceptabilă (rețele WAN care oferă servicii neorientate pe conexiune, interconectările de rețele). În acest context, o eroare este definită ca un pachet pierdut sau duplicat, o eroare semnalată este o eroare detectată și nerecuperată de nivelul rețea și, ca urmare semnalată nivelului transport iar o eroare nesemnaltată este o eroare nedetectată de nivelul rețea și, ca urmare nerecuperată și neraportată. Chiar și în cazul rețelelor sigure, de tipul A, trebuie luate în considerare încă două aspecte care aduc complicații nivelului transport: limitarea dimensiunii unui pachet și secvențierea pachetelor (recepționarea pachetelor trebuie făcută în aceeași ordine în care au fost emise).

A. Gestiunea transferului datelor

Gestiunea transferului datelor la nivelul transport este evident influențată de tipul de rețea utilizat (A, B sau C) deoarece, pe de o parte, mecanismele de control al fluxului și de confirmare a pachetelor în rețele de tipul A sunt minimale pe când în rețelele de tipul C ele sunt complexe, iar, pe de altă parte, în rețelele de tipul C trebuie luate în considerare și erorile posibile în cursul transferului proceselor și deci construiți algoritmi pentru retransmisia pachetelor și verificarea corectitudinii lor.

S-ar putea spune că, în intenția de a corecta imperfecțiunile nivelului rețea, nivelul transport va avea în ceea ce privește transferul datelor, funcții similare nivelului legaturii de date: detectarea erorilor, controlul fluxului; complexitatea soluțiilor adoptate la nivelul transport este însă mai mare, deoarece, dacă la nivelul legaturii de date două noduri comunică direct printr-un canal fizic, la nivelul transport două noduri comunică prin intermediul unei întregi subrețele, în care pachetele pot fi mai ușor pierdute, întârziate și deci duplicate. Mecanismele pentru controlul fluxului și pentru rezervare de buffere sunt necesare în ambele nivele: datorită

Însă numărul mult mai mare de conexiuni existente la nivelul transport decât numărul de linii de la nivelul legăturii de date este inefficient ca la nivelul transport să se aloce individual, pentru fiecare conexiune, câte un set de buffere, ca la nivelul legăturii de date. Strategia de rezervare a bufferelor este integrată în *mecanismul de control al fluxului*; acesta este mai complex decât controlul fluxului de la nivelul legăturii de date deoarece pe de o parte (a) implică interacțiunea a trei nivele, nivelul utilizator al nivelului transport, nivelul transport și nivelul rețea, și pe de altă parte, (b) timpul de transfer al datelor între două entități transport din noduri diferite este variabil, depinzând de încărcarea rețelei. Figura 2.1.3.2. ilustrează aspectul (a) de mai sus: la

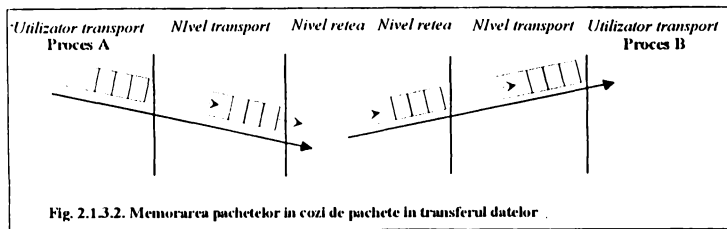
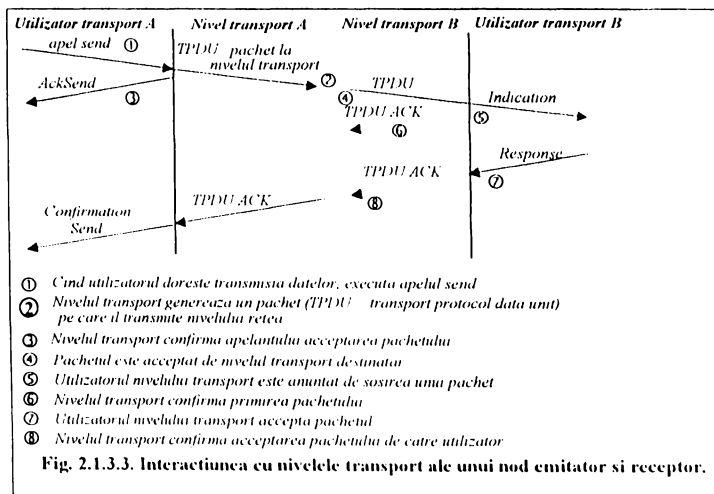


Fig. 2.1.3.2. Memorarea pachetelor în cozi de pachete în transferul datelor.

transmiterea unui pachet de date de către procesul A (procesul A este utilizatorul nivelului de transport) către procesul B sunt implicate patru cozi de pachete de date. Astfel,

procesul A generează pachetul și îl depune în coada de pachete pentru a fi trimis către B; procesul A poate aștepta la operația send până: (1) are permisiunea nivelului transport din propriul nod (controlul fluxului la nivel de interfață) și (2) are permisiunea nodului B (controlul fluxului capăt la capăt). După ce nivelul de transport a preluat pachetul, îl va memora într-o coadă până când are permisiunea de la nivelul rețea de a-l trimite, după care pachetul este trimis la nodul destinație de către nivelul rețea. La nodul destinație, nivelul rețea trebuie să păstreze în coadă pachetul până când are de la B permisiunea de a-l livra; analog B va memora pachetul până va avea permisiunea destinației (procesul B) de a-i livra pachetul. Figura 2.1.3.3. ilustrează interacțiunile posibile, ca efect al întâzierilor de mai sus. Se observă că există trei cazuri posibile de a confirma o operație *send*: (1) la acceptarea datelor pentru transfer de către nivelul de transport al emițătorului (3 în fig.2.1.3.3.); (2) la acceptarea datelor de către nivelul transport al receptorului (6 în fig.2.1.3.3.); (3) la acceptarea datelor de către utilizatorul datelor din nodul destinat (8 în fig. 2.1.3.3.).

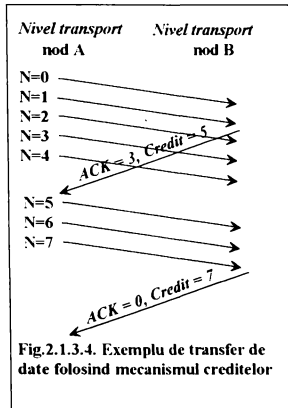


- ① Când utilizatorul dorește transmitia datelor, executa apelul send
- ② Nivelul transport generează un pachet (TPDU - transport protocol data unit) pe care îl transmite nivelului rețea
- ③ Nivelul transport confirma apelandului acceptarea pachetului
- ④ Pachetul este acceptat de nivelul transport destinat
- ⑤ Utilizatorul nivelului transport este anunțat de sosirea unui pachet
- ⑥ Nivelul transport confirma primirea pachetului
- ⑦ Utilizatorul nivelului transport accepta pachetul
- ⑧ Nivelul transport confirma acceptarea pachetului de către utilizator

Fig. 2.1.3.3. Interacțiunea cu nivelele transport ale unui nod emițător și receptor.

O concluzie a discuției de mai sus este faptul că există două motive pentru care nivelul transport al unui nod destinat poate cere emițătorului micșorarea ratei de transmisie a datelor: atunci când nivelul transport al destinației nu poate face față fluxului de date de la emițător sau când utilizatorul nivelului transport al destinației nu poate face față fluxului de date. Ambele cauze se manifestă în ultima instanță prin epuizarea capacității de memorare de la nivelul transport. În această situație, pentru controlul fluxului datelor emițătorului se pot aplica următoarele metode: (1) receptorul nu execută nimic special; (2) receptorul refuză să primească alte pachete; (3) se utilizează un algoritim cu fereastra glisantă; (4) se utilizează o schema cu credite. Abordarea (1) este puțin eficientă; ea nu are ca efect micșorarea ratei de emisie a pachetelor de către emițător, ci dimpotrivă, acesta nemaiprimit confirmările va trebui să retransmită, pe lângă noile date și

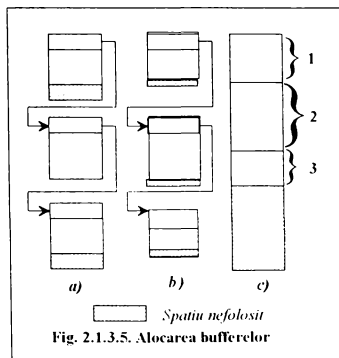
datele neconfirmate încă. Metoda (2) transferă problema controlului fluxului la nivelul rețea: nivelul transport al receptorului refuză să mai primească date de la nivelul rețea și acesta va trebui să rezolve controlul fluxului la nivelul său, nivelul rețea al emițătorului refuzând să mai primească date de la emițător. Metoda (3) este cunoscută de la nivelul legăturii de date. Ea lucrează bine într-o rețea de tipul A, dar într-o rețea de tipul B sau C poate genera pachete duplicate nedorite. Metoda (4) este, în schimb, convenabilă în orice tip de rețea. Spre deosebire de (3) ea decuplează confirmarea pachetelor de problema controlului fluxului. Pe scurt, metoda constă tot într-un algoritim cu fereastra glisantă, dar, în cursul transferului, se transmite emițătorului numărul de buffere ale receptorului ramase disponibile (evident în limita dimensiunii ferestrei); pentru a nu se genera noi pachete pentru credite, acestea se trimit în corpul confirmărilor. Prin urmare, pentru a confirma N pachete și a acorda credite, un nivel transport receptor va trimite o confirmare de forma (ACK N+1, Credit M) ceea ce înseamnă confirmarea pachetelor 1...N și acordarea de credite pentru pachetele N+1, ..., N+M (numai



acestea pot fi transmise). Pentru a crește sau descrește creditul la X, fără ca alte pachete să fi fost recepționate, receptorul va trimite o confirmare de forma (ACK N+1, Credit X); pentru a confirma un nou pachet primit, receptorul va trimite (ACK N+2, Credit M-1). În fig. 2.1.3.4. se prezintă un exemplu de transfer de date folosind acest algoritim. Chiar dacă algoritmul este utilizat într-o rețea de tipul C, el funcționează corect și eficient: dacă se pierde un pachet de confirmare, emițătorul, în urma unei cereri de timeout va retransmite pachetul și va primi confirmarea (și creditul). Poate apare însă o interblocare când, după recepționarea unei confirmări de tipul (ACK N, Credit 0), practic închiderea ferestrei, receptorul trimite emițătorului (ACK N, Credit M) dar confirmarea este pierdută; emițătorul așteaptă un nou credit de la receptor, care însă l-a trimis. Ca atare, pentru a înlătura interblocarea, trebuie setat un timer la trimiterea unui pachet de confirmare; când timpul a expirat, se repetă confirmarea (dacă se duplică cea precedentă nu este nici o eroare).

Pentru ambele metode de control al fluxului, la nivelul receptorului și emițătorului trebuie utilizată o metodă pentru *rezervarea bufferelor*, având în vedere că nivelul transport va avea de tratat mai multe conexiuni. Dacă se ține cont de faptul că acestea se stabilesc și desființează dinamic, atunci rezulta că este eficient ca *alocarea bufferelor* să se facă dinamic, la stabilirea conexiunilor. Bufferele pot fi alocate separat, pentru fiecare conexiune sau, per global, pentru toate conexiunile dintre două noduri. Este eficient ca pentru conexiunile cu trafic mare (de exemplu transfer de fișiere) să se utilizeze buffere dedicate, iar pentru conexiunile cu trafic variabil (de exemplu cel produs de un terminal) să se rezerve bufferele pe măsură ce sunt necesare. Un alt aspect legat de buffere este *dimensiunea lor* (fig.2.1.3.5.). Pot fi utilizate: (a) buffere de dimensiune fixa cu dezavantajul irosirii de spațiu dar cu o

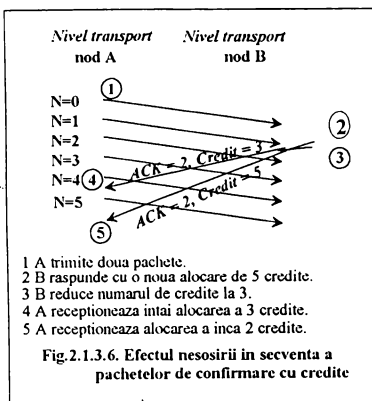
gestionare simplă (b) buffere de dimensiune variabila, sau, pentru o mai bună utilizare a memoriei se poate alocă (c) un singur buffer, mare, circular pentru fiecare conexiune (metodă eficientă în cazul conexiunilor cu trafic mare).



Alocarea dinamică a bufferelor determină decuplarea mecanismului confirmărilor de cel al controlului fluxului ceea ce implică utilizarea algoritmului cu credite. Pentru controlul fluxului datelor poate fi eficientă și ajustarea ferestrei [Tan89]. Dimensiunea ferestrei depinde de capacitatea de transmisie a subrețelei: astfel, dacă subrețeaua poate transmite P pachete/sec și durata ciclului unui pachet în rețea (incluzând transmisia, memorarea în cozi, prelucrarea la nodul destinație și returnarea confirmării) este C atunci fereastra emițătorului ar trebui să fie P*C, pentru ca timpul ciclului unui pachet să fie folosit complet. Fereastra trebuie ajustată periodic deoarece capacitatea rețelei depinde de traficul pachetelor. În acest scop se pot măsura parametrii P și C: pentru C se alege o valoare medie, iar P se

evaluează prin împărțirea numărului de confirmări primite într-o anumită perioadă de timp la aceea perioadă de timp (în timpul măsurătorii, emițătorul va trimite pachete la rata maximă de transmisie pentru a fi sigur că rata de primire a confirmărilor este limitată de capacitatea rețelei și nu de rata de transmisie mai lentă).

Alte două aspecte ale transferului de date caracteristice tuturor tipurilor de servicii de rețea A. B. C sunt: (1) pastrarea ordinii de emisie a pachetelor și la recepție și (2) fragmentarea pe pachete a unui mesaj primit de la nivelul superior. Aspectul (1) implică nu numai numerotarea în secvență a pachetelor din cadrul unei conexiuni, dar și a pachetelor de confirmare, dacă se folosesc credite. Fig. 2.1.3.6. arată efectul

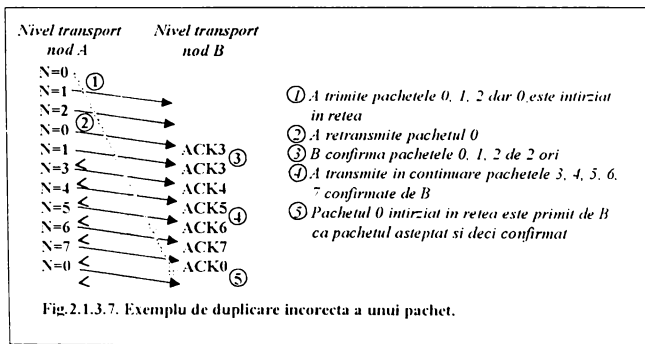


dezastruos al recepționării pachetului (ACK 2, Credit=3) înaintea pachetului (ACK 2, Credit=5).

Aspectul (2) trebuie luat în considerare pentru două cazuri distincte, având în vedere că aplicațiilor le sunt caracteristice două tipuri de transfer de date: (a) orientate pe flux, când utilizatorul nivelului transport acceptă datele în flux continuu, fără ca eventualele puncte de oprire să aibă semnificație și (b) orientate pe blocuri, când receptorul acceptă datele la nivel de bloc (un bloc este un mesaj) așa cum au fost transmise de emițător. În cazul b), dacă un bloc (TSDU în terminologia OSI) depășește dimensiunea unui pachet, el trebuie fragmentat, și reasamblat la nodul recepție înainte de livrarea la utilizator. Pentru a se putea controla ordinea sosirii pachetelor unui bloc la destinație nu este nevoie de numerotarea pachetelor din interiorul blocurilor (aceasta ar conduce la o dubla numerotare deoarece există numărul de secvență al pachetelor dintr-o conexiune) ci este suficient ca ultimul pachet al unui bloc să fie marcat cu un indicator

special de sfârșit de bloc, EOT; toate pachetele următoare pînă la primul pachet cu indicatorul EOT inclusiv, vor fi considerate ca făcînd parte din următorul bloc s.a.m.d..

Verificarea corectitudinii datelor (obligatorie în cazul rețelelor din clasa C pentru a compensa eventualele erori nesemnificate ale pachetelor, în decursul parcurgerii subrețelei) se poate face prin adăugarea unei sume de control în headerul fiecărui pachet transmis. În cazul rețelelor de clasa C, implicații mari o are strategia de retransmisie a pachetelor neconfirmate. Aspectul critic în ceea ce privește strategia de retransmisie este stabilirea valorii timeout-ului; o valoare fixată nu răspunde cerințelor de adaptare la încărcarea rețelei; de aceea este indicată o valoare medie bazată pe măsurarea întâzierilor confirmărilor. Chiar și o schemă adaptivă pentru stabilirea timeout-ului găsește tot o valoare aproximativă deoarece destinația nu trimite confirmările imediat după primirea pachetelor (se folosesc algoritmi cu fereastră glisantă), încărcarea rețelei este variabilă iar măsurătorile pot fi imprecise pentru că dacă un pachet este retransmis nu se poate distinge între confirmarea pachetului inițial și a celui duplicat. Un alt aspect critic legat de strategia de retransmitere este că secvența numerelor pachetelor dintr-o fereastră nu trebuie să se repete de cît după cel puțin timpul maxim de viață al unui pachet în rețea. Fig. 2.1.3.7. ilustrează un

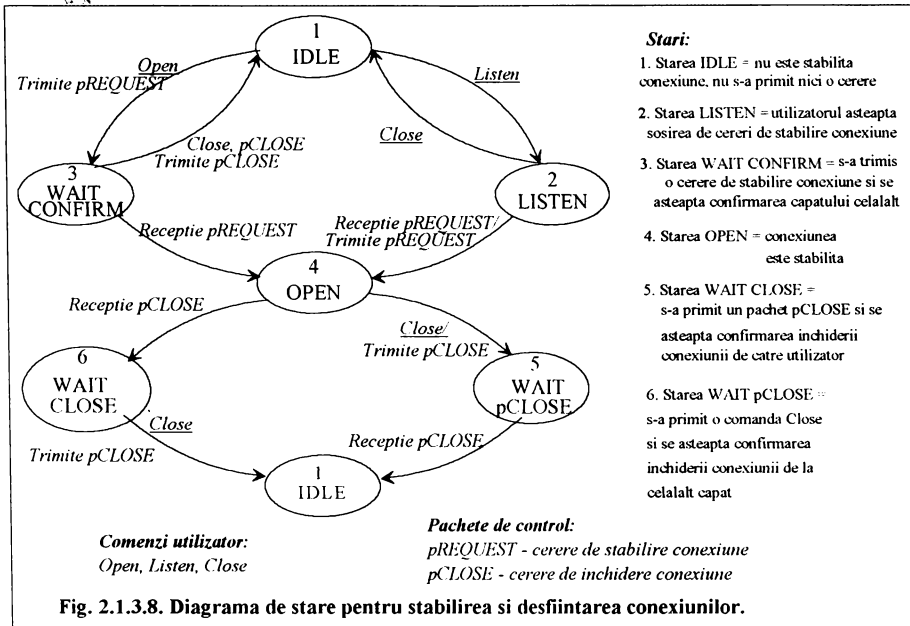


contraexemplu pentru aceasta cerință. Se presupune, pentru simplitate că se folosește pentru controlul fluxului un algoritmul cu fereastră glisantă de dimensiune 3 (deci un spațiu de 8 numere). Se observă că dacă pachetul 0 apare la destinație după epuizarea unei secvențe complete, dereglează întregul flux al datelor.

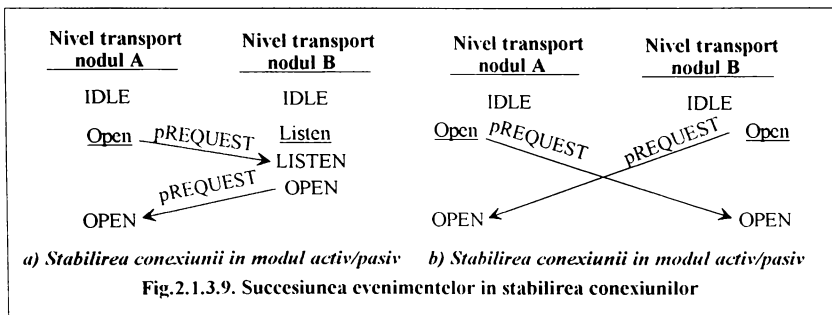
B. Stabilirea și desființarea conexiunilor

În scopul implementării serviciilor orientate pe conexiuni sunt necesare, pentru toate tipurile de rețele, A, B, C proceduri pentru stabilirea, respectiv desființarea unei conexiuni. Stabilirea unei conexiuni servește, în principal, pentru 3 scopuri: (1) permite fiecărui capăt să se asigure că celălalt există; (2) permite negocierea unor parametri opționali (dimensiunea pachetelor, dimensiunea ferestrei, calitatea serviciului) (3) declanșează alocarea resurselor necesare conexiunilor (buffer, intrări în tabela de conexiuni).

Deoarece la nivelul transport pot exista mai multe conexiuni deschise pe aceeași mașină, fiecărei conexiuni trebuie să i se asocieze un număr care se depune în fiecare pachet, în scopul recunoașterii apartenenței pachetelor de către receptor. Conexiunile se stabilesc prin acord reciproc, în urma unor comenzi de la utilizatorul nivelului transport și a unor pachete de control, așa cum se prezintă în diagrama din fig. 2.1.3.8. Se observă că utilizatorul semnalează că așteaptă o cerere de conexiune printr-o comandă Listen, care declanșează tranziția IDLE -> LISTEN (un exemplu ar fi un server de fișiere); utilizatorul poate închide



conexiunea trimițând o comandă *Close*. Dacă în starea LISTEN se primește o cerere de stabilire conexiune pREQUEST (pachet de control *request for connection*) care specifică ca destinație utilizatorul aflat în așteptare, atunci se stabilește o conexiune. În acest caz, nivelul transport trimite o confirmare pREQUEST nivelului transport pereche și transferă conexiunea în starea OPEN. Un utilizator poate deschide o conexiune executând o comandă *Open*, care declanșează trimiterea unui pachet pREQUEST nivelului transport pereche. Recepția unui pachet pREQUEST de la nodul pereche stabilește o conexiune. Conexiunea poate fi desființată prematur dacă utilizatorul local execută o comandă *Close* sau nivelul transport pereche refuză conexiunea trimițând un pachet pCLOSE. Figura 2.1.3.9. arată succesiunea evenimentelor în stabilirea conexiunilor. Se observă că stabilirea conexiunii se face corect chiar și în cazul cînd ambele capete încearcă să o stabilească aproximativ în același timp.



În cazul în care se recepționează un pachet pREQUEST cînd utilizatorul nu a cerut așteptarea cererilor de stabilire a conexiunii există trei soluții: (1) nivelul transport poate rejecta cererea trimițând un pachet pCLOSE (2) cererea este memorată într-o coadă de cereri pînă utilizatorul execută o comandă *Listen* (3) nivelul transport poate întrerupe sau semnaliza într-un mod oarecare utilizatorului că există o cerere de stabilire conexiune în așteptare. În ultimul caz, nu mai este strict necesară o comandă *Listen*; ea poate fi înlocuită printr-o comandă *Accept*, care este văzută ca un semnal de la utilizator către nivelul transport că accepta cererea de conexiune (aceasta este soluția adoptată de ISO).

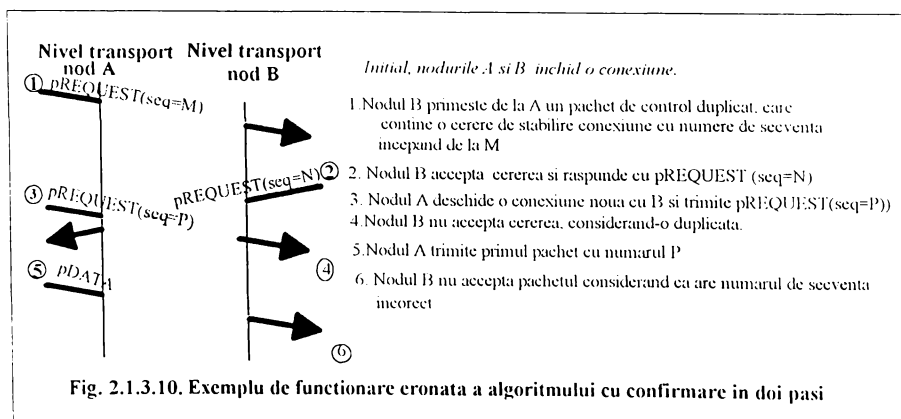
Desființarea conexiunilor se face similar; desființarea poate fi cerută de orice capăt, de ambele capete sau de nivelul transport însuși la ambele capete. Conexiunea se desființează prin acord reciproc; strategia permite o desființare bruscă sau lentă. În ultimul caz, o conexiune în starea 5 (se așteaptă pCLOSE) trebuie să continue să accepte date pînă la recepția unui pachet pCLOSE.

Diagrama prezentată în fig.2.1.3.8. este valabilă pentru rețele de tipul A. Pentru celelalte tipuri de rețele trebuie luate în considerare aspecte legate de (1) nesosirea în secvență a pachetelor (2) generarea de către nivelul rețea a unei primitive N-RESET pentru a semnala o situație catastrofală în funcționare (rețea supraincercată, probleme hardware, software) (3) erori în transmisie și apariția de pachete duplicate.

Situațiile de tipul (1) pot provoca recepția unui (sau mai multor) pachete de date înainte de recepția acceptării conexiunii trimisă de capatul pereche; în acest caz este necesar să se depună aceste pachete într-o coadă de pachete pînă la recepția acceptării conexiunii. O situație asemănătoare poate apare și la desființarea conexiunilor: după ce un nivel transport a trimis ultimul pachet de date, trimite pachetul de control pCLOSE, care poate ajunge la destinație înaintea ultimului pachet de date. O soluție pentru a evita astfel de probleme este ca pachetul pCLOSE să conțină numărul de secvență de după ultimul pachet transmis; astfel, nivelul transport destinație va aștepta eventualele pachete întârziate și după ce a recepționat pCLOSE, fără închiderea conexiunii.

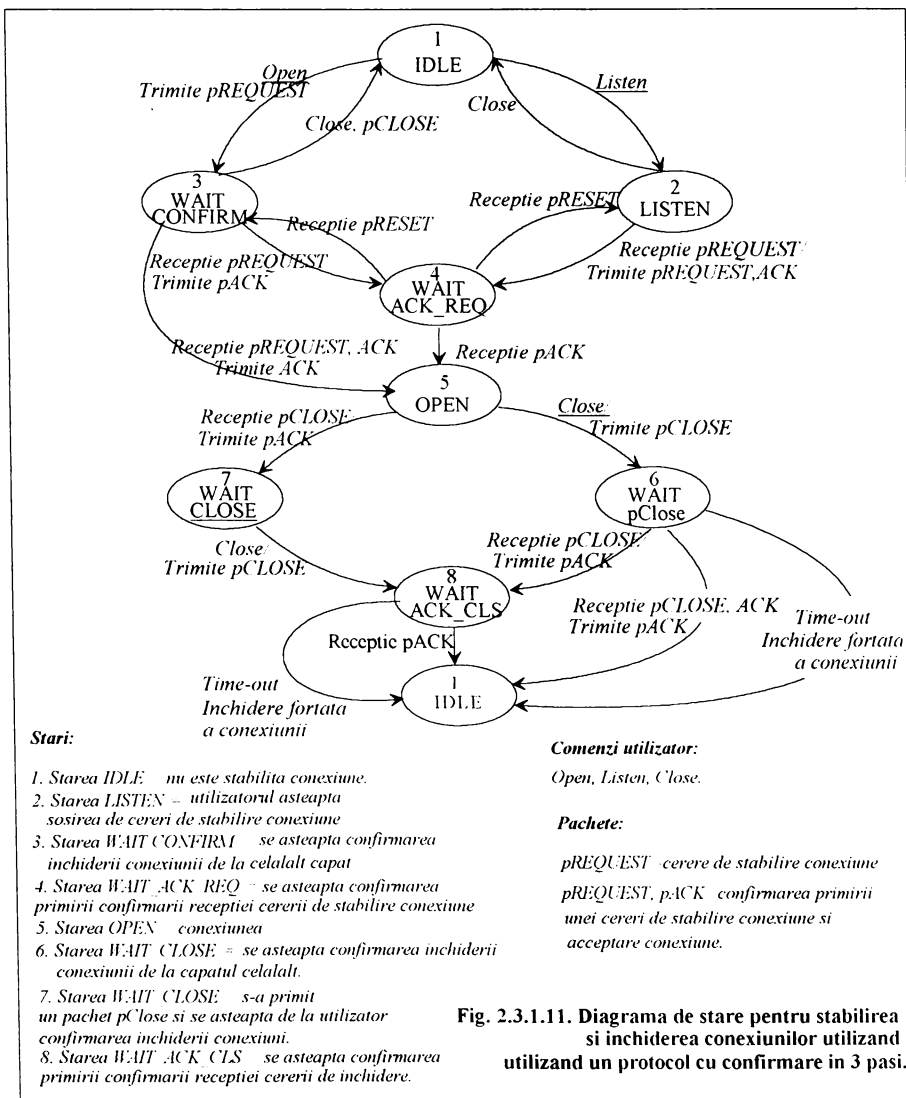
Se consideră că situațiile de tipul (2) pot apare foarte rar în rețele de tip A; în cazul apariției lor se închid conexiunile fără a se mai încerca recuperarea lor. Pentru casele B, C se încearcă o resincronizare după un N-RESET în vederea determinării de către fiecare capăt a pachetelor trimise, și recepționate (ceea ce înseamnă că trebuie reținute copii ale pachetelor trimise, dar neconfirmate, pentru a putea fi retransmise după un N-RESET) și stabilirea unei noi conexiuni la nivelul rețea.

Situațiile de tipul (3) pot genera probleme la stabilirea conexiunilor. Iată un exemplu: dacă după închiderea unei conexiuni, mai există pachete duplicate și dacă se deschide imediat o altă conexiune între aceleași entități. Se poate întâmpla ca pachetele duplicate, întârziate aparținînd primei conexiuni să fie acceptate ca valabile și în cadrul celei de-a doua conexiuni; analog și confirmările întârziate pot fi acceptate în alte conexiuni și crea probleme. Există mai multe soluții pentru aceste situații: (a) extinderea numerelor de secvență ale pachetelor și după închiderea conexiunilor. Aceasta înseamnă că nivelul transport trebuie să cunoască numerele ultimelor pachete ale conexiunilor închise, astfel încît la stabilirea unei noi conexiuni să aleagă un alt număr de secvență pentru primul pachet pe care îl transmite capătului celalalt în pREQUEST. Acest procedeu este simetric; fiecare capăt trebuie să-și declare la stabilirea conexiunii numărul de secvență al primului pachet; ele nu trebuie să fie egale, pentru cele două capete; (b) utilizarea unui identificator de conexiune separat pentru fiecare tranzație, ales de capătul ce inițiază conexiunea (de exemplu, la stabilirea unei conexiuni se incrementează identificatorul) care se depune în pachete; la închiderea fiecărei conexiuni, nivelul transport actualizează o tabela a conexiunilor terminate (fiecare intrare conține identitatea celui alt capăt și identificatorul conexiunii). Ambele procedee de mai sus se comportă bine numai dacă nu apare o avarie a nodului, cînd se pierde atît numerele de secvență ale ultimelor pachete ale conexiunilor închise cît și identificatorii conexiunilor închise. În acest caz, singura soluție este limitarea duratei de viață a pachetelor în rețea și deci, după avaria unui nod se așteaptă timpul maxim de viață al pachetelor după care se pot deschide noi conexiuni fără probleme. Pentru limitarea duratei de viață a pachetelor se folosesc 3 metode: (a) limitarea circulației pachetelor în subrețea (se previne astfel circulația în bucla închisă a pachetelor) (b) folosirea unui contor al nodurilor intermediare parcurse de un pachet, care se depune în fiecare pachet; acesta se inițializează



cu o valoare maximă permisă și se decrementează, în fiecare nod intermediar. Când contorul ajunge la zero, protocolul legaturii de date refuză să primească pachetul. (c) folosirea unei marci de timp pentru fiecare pachet. La crearea unui pachet se inițializează marca cu valoarea timpului curent, urmînd ca pachetul să fie ignorat cînd devine mai vechi decît o anumită perioadă de timp stabilită. Metoda este însă mai dificil de implementat deoarece necesită ca ceasurile nodurilor să fie sincronizate.

Pachetele duplicate pot avea însă și alte efecte dezastruoase la stabilirea conexiunilor. Figura 2.1.3.10. arată o secvență de pachete de control în urma căreia algoritmul de stabilire a conexiunilor prezentat, și numit *algoritmul cu confirmare în 2 pași (two-way handshake)* nu funcționează corect dacă ignoră duplicatele cererilor de stabilire conexiuni după stabilirea conexiunilor. Eroarea poate fi corectată dacă fiecare capăt confirmă cererea de stabilire și numărul de secvență al celuilalt. Algoritmul este cunoscut sub numele de *algoritm cu confirmare în 3 pași*; el conduce la modificarea diagramei de stari la cea prezentată în figura 2.3.1.11.



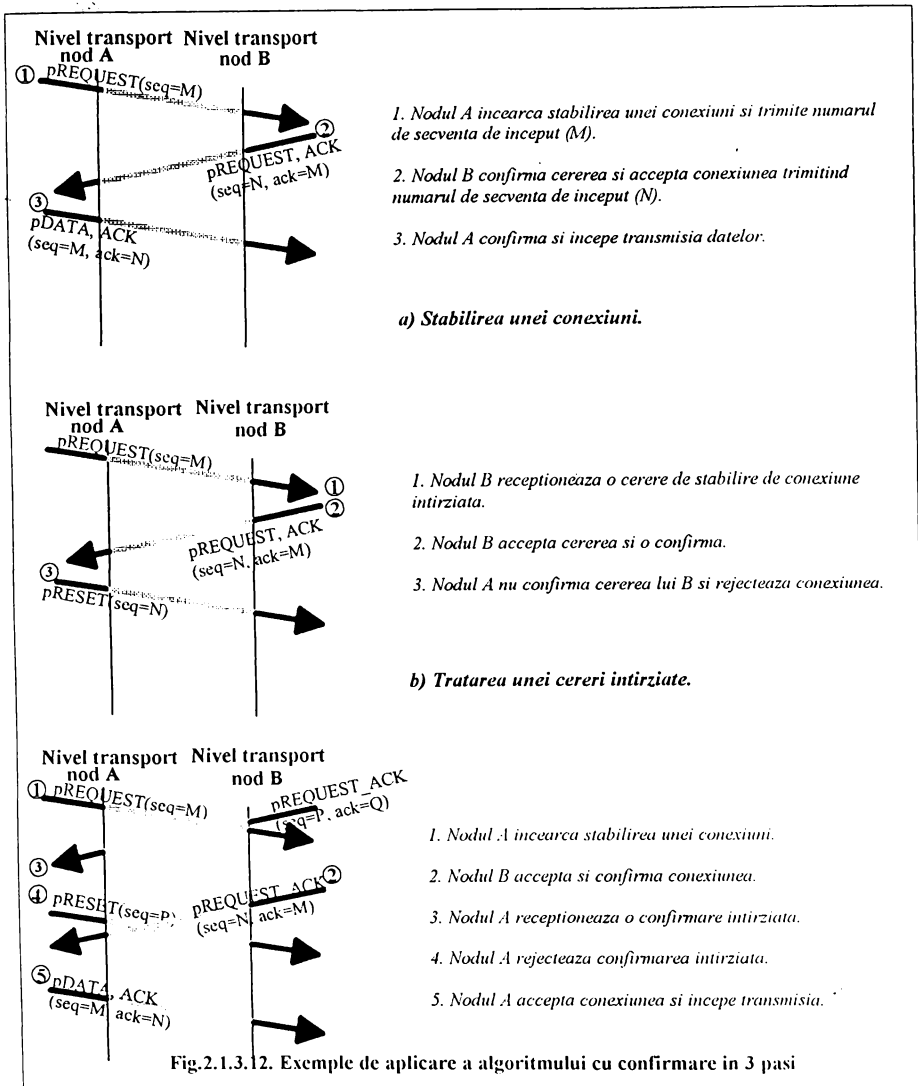


Fig.2.1.3.12. Exemple de aplicare a algoritmului cu confirmare in 3 pasi

Față de diagrama din fig. 2.1.3.8, apare o stare în plus pînă la stabilirea conexiunii, starea `WAIT_ACK_REQ` și un tip de pachet de control nou, pachetul `PRESET` pentru a reseta capătul celălalt cînd s-a detectat recepția unei cereri de stabilire `pREQUEST` duplicate.

Figura 2.1.3.12 prezintă aplicarea algoritmului de confirmare în trei pași în trei cazuri: se observă că în cazurile b) și c) de apariție a unui pachet duplicat a unei cereri de stabilire conexiune, respectiv a unui pachet duplicat a confirmării unei cereri, algoritmul lucrează corect.

Pentru simplificarea diagramei de stare din fig. 2.1.3.11, nu au fost reprezentate tranzițiile în care se transmite `PRESET`. Pentru acestea, regula de bază este: se transmite `PRESET` dacă starea conexiunii nu este `OPEN` și dacă se primește o confirmare `ACK` invalidă (care refera o cerere ce nu s-a trimis).

Pentru închiderea unei conexiuni se poate folosi o soluție asemănătoare: fiecare capăt trebuie să confirme explicit fiecare pachet `pCLOSE` al celuilalt. Numerele de secvență se folosesc astfel: (a) pachetul `pCLOSE` conține numărul de secvență plus unu al ultimului pachet de date trimis; (b) pachetul `ACK` conține numărul de secvență primit în pachetul `pCLOSE`. Pentru a elimina interblocarea care poate apare în cazul

pierderii de pachete, trebuie folosite mecanisme de timeout (fig.2.1.3.11); după un număr oarecare de retransmisii, conexiunea este închisă forțat.

În cazul avariei unui nod (și repunerii în funcțiune apoi a lui) toate informațiile despre conexiuni vor fi distruse, iar conexiunile afectate devin "deschise pe jumătate". Capătul activ al conexiunii poate închide conexiunea cu ajutorul unui timer care masoară timpul total de așteptare a unei confirmări după un număr maxim de retransmisii; dacă se depășește acest timp se închide conexiunea și se semnalează utilizatorului nivelului transport închiderea forțată. În cazul refacerii imediate a nodului cazut, închiderea conexiunii se poate face și altfel: nivelul transport trebuie să returneze RESET X la orice pachet; capătul activ al conexiunii, va verifica validitatea numărului de secvență X (pentru a depista dacă RESET este răspunsul pentru un pachet întârziat); în cazul unui RESET activ, nivelul transport va declanșa închiderea forțată a conexiunii. Resincronizarea între cele două noduri nu poate fi făcută însă la nivelul transport, deoarece nu se poate determina câte pachete au fost transmise corect.

2.1.4. Nivele suport pentru aplicații

Nivelele superioare ale modelului ISO/OSI - nivelul *sesiune*, *prezentare* și *de aplicații* - realizate pe

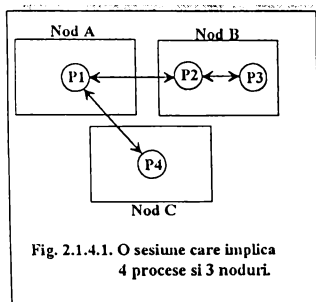


Fig. 2.1.4.1. O sesiune care implica 4 procese și 3 noduri.

baza serviciilor de comunicație independente de rețea furnizate de nivelul transport au sarcina de a furniza unui proces utilizator facilitățile necesare schimbului de informații în cadrul unui mediu de prelucrare distribuită a informației. Ele sunt referite colectiv ca *nivele suport pentru aplicații* sau *nivele specifice aplicațiilor*.

Nivelul sesiune extinde comunicația dintre două procese, prezentată la nivelul transport; astfel, la acest nivel, comunicația între procese apare în contextul unei aplicații, divizată în subtaskuri, fiecare task putând fi mapat în unul sau mai multe procese. Nivelul sesiune gestionează comunicația între o pereche de astfel de procese (fig.2.1.4.1) oferind suport pentru stabilirea și terminarea sesiunilor de transmisii, între procese și pentru organizarea și sincronizarea dialogului dintre ele. Serviciile acestui nivel, definite de ISO 7498, sunt descrise în fig. 2.1.4.2.

<u>Servicii</u>	<u>Descriere</u>
(1) Stabilirea unei conexiuni sesiune	Folosita pentru conectarea entitatilor prezentare
(2) Eliberarea unei conexiuni sesiune	Eliberarea se poate face după transmiterea în ambele direcții a unitatilor de date (SSDU) în așteptare
(3) Serviciu pentru date în carantina	Permite entității prezentare de a cere carantina (bufferare) pentru una sau mai multe unitati de date (SSDU)
(4) Transferul datelor normale și a datelor expeditive	Transfer de date
(5) Gestiunea interaciunilor	Dialogul între 2 procese poate fi în ambele direcții simultan, în ambele direcții alternind sau într-o singură direcție
(6) Sincronizare	Permite entitatilor prezentare definirea de puncte de sincronizare și resincronizare în raport cu aceste puncte
(7) Raportarea excepțiilor	Anunta aparitia unei erori entității prezentare

Fig. 2.1.4.2. Serviciile nivelului sesiune

Se apreciază [Tan88] că nivelul sesiune este implementat mai puțin în diverse arhitecturi de rețea, fiind o invenție OSI.

Nivelul prezentare are drept focar de interes sintaxa datelor schimbate între entitățile aplicație; scopul său este de a rezolva diferențele în ceea ce privește formatul și reprezentarea datelor. Acest lucru este necesar deoarece la nivelul aplicație, prezintă importanță *semantica datelor*, iar de la nivelul sesiune până la nivelul fizic datele sunt văzute ca șiruri de octeți. Este sarcina nivelului prezentare de a furniza reprezentarea datelor de la nivelul aplicație într-o formă binară, cu aceeași semnificație, la ambele capete, de aceea centrul său de interes este *sintaxa datelor* de la nivelul aplicație. Deoarece este practic imposibil de a separa complet semantica datelor de sintaxa lor, ISO a propus următoarea soluție: (fig.2.1.4.3.): (1) la nivelul aplicație informația este reprezentată într-o *sintaxă abstractă*, care specifică datele la nivel de tipuri și valori de date,

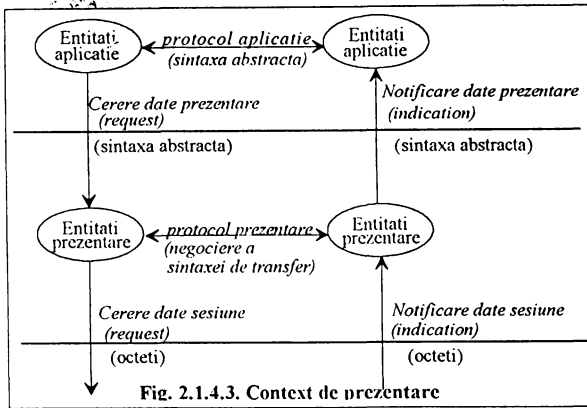


Fig. 2.1.4.3. Context de prezentare

independent de orice reprezentare specifică (procoloale la nivel de aplicație vor descrie unitățile lor de date în termenii sintaxei abstracte) (2) nivelul prezentare comunică cu nivelul aplicație în termenii sintaxei abstracte (3) nivelul prezentare translatează datele reprezentate în sintaxa abstracta într-o sintaxa de transfer, sub o forma binara, utilizată pentru interacțiunea cu nivelul sesiune. Sintaxa de transfer definește reprezentarea datelor schimbate între entitățile prezentare; translatarea sintaxa abstracta -> sintaxa de transfer se face cu ajutorul unor reguli de codificare, care specifică reprezentarea valorii fiecărui tip de

date. Fiecare entitate prezentare recunoaște sintaxa abstractă a utilizatorului său și una sau mai multe sintaxe de transfer disponibile pentru codificare (de exemplu, pentru o aceeași sintaxa abstractă pot exista mai multe variante ale unei sintaxe de transfer: o variantă care realizează codificarea datelor, alta care realizează criptarea, alta care realizează și criptare și codificare s.a.m.d); alegerea unei sintaxe de transfer se realizează prin negociere între două entități prezentare. Combinația dintre o sintaxa abstractă și o sintaxa de transfer este denumită context de prezentare. ISO reunește definițiile pentru o sintaxă abstractă (ASN.1) în standardul DIS 8824, iar regulile de codificare în standardul DIS 8825 (Basic Encoding Rules).

În arhitectura OSI funcțiile nivelului prezentare pot fi rezumate astfel: (1) negociere pentru acceptarea unei sintaxe de transfer convenabila pentru reprezentarea datelor la ambele capete (2) conversia datelor la acest format la ambele capete (3) maparea cererilor de servicii ale aplicațiilor în funcții, precum controlul dialogului și al sincronizării în primitivele corespondente ale serviciului sesiune. Deoarece înainte de a fi transmise nivelului sesiune, datele pot fi criptate și comprimate, criptarea și comprimarea datelor sunt alte aspecte legate de nivelul prezentare.

Nivelul aplicație ca nivel de graniță între mediul OSI și procesele utilizator de aplicație, furnizează servicii acestor procese (aflate deci în mediul exterior arhitecturii OSI) în contextul unei prelucrări distribuite. Nivelul aplicație este compus din cateva module separate, denumite elemente de servicii. Unele dintre aceste

elemente sunt comune multor aplicații și de aceea referite ca elemente de servicii comune aplicațiilor (CASE= Common Application Service Elements); ele includ funcții precum gestiunea conexiunilor aplicație (modulul ACSE=Association Control Service Element) și coordonarea acțiunilor entităților aplicație (modulul CCR= Commitment Concurrency and Recovery). Alte elemente de servicii au fost prevazute special pentru a furniza servicii specifice și de aceea sunt referite ca elemente de servicii ale unor aplicații specifice (Specific Application Service Elements); ele includ servicii precum: (1) accesul dintr-un proces utilizator la un sistem de fișiere aflat la distanță - modulul FTAM (File Transfer Access and Management) - (2) comunicarea a două procese utilizator care gestionează dialoguri cu terminale locale într-o manieră standardizată, independentă de caracteristicile și proveniența terminalului - modulul VT (Virtual Terminal) - (3) furnizarea unei facilități generale de poșta electronică - modulul MHS (Message Handling Service) - (4) obținerea adresei de rețea a utilizatorilor pornind de la numele lor simbolice - modulul DS (Directory Services). Structura generală a nivelului aplicație este reprezentată în figura 2.1.4.4.

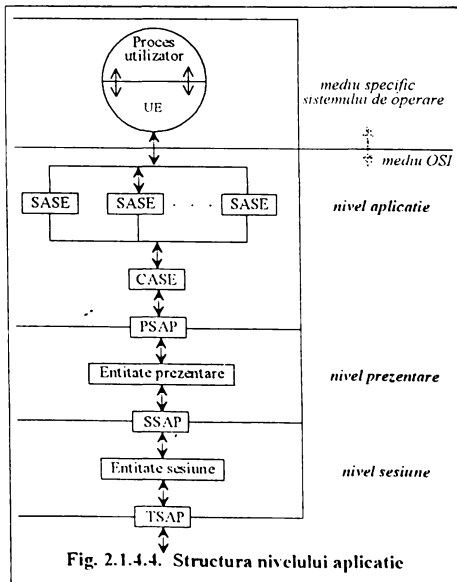


Fig. 2.1.4.4. Structura nivelului aplicatie

Interfața dintre serviciile furnizate de nivelul aplicație și procesele utilizator, denumită în general *elementul utilizator (UE= User Element)* sau *agentul utilizator (UA= User Agent)* sau *interfața utilizator (UI= User Interface)* se implementează printr-un set de funcții sau proceduri de bibliotecă, linkeditate cu procesele utilizator. Primitivele interfeței dintre procesele utilizator și UE nu trebuie să fie aceleași primitive (standard) furnizate de un element serviciu al entității aplicație. Neexistând aceasta restricție, pot fi folosite primitive specifice sistemului de operare local; este în sarcina UE-ului să realizeze corepondența necesară. Acesta a fost

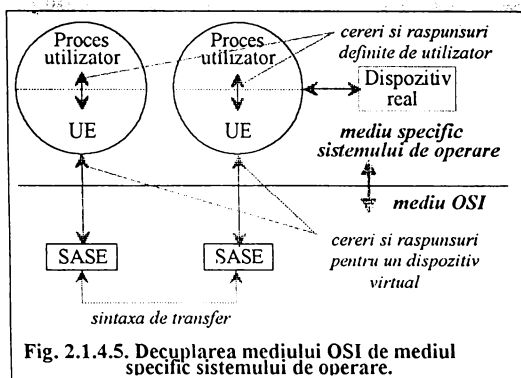


Fig. 2.1.4.5. Decuplarea mediului OSI de mediul specific sistemului de operare.

de altfel unul din scopurile ISO pentru nivelul aplicație: decuplarea mediului OSI de mediul specific sistemului de operare. Decuplarea se realizează prin următoarea manieră generală, comună tuturor elementelor de servicii ale unor aplicații specifice (fig.2.1.4.5): în primul rind, pentru fiecare element se definește un dispozitiv virtual, împreună cu un set de primitive standard pentru utilizarea lui. Toate cererile și răspunsurile de la nivelul mediului OSI se fac folosind primitivele standard. Pentru maparea primitivelor relative la dispozitivul real (cereri și răspunsuri), de la nivelul proceselor utilizator în primitive standard se folosește UE.

2.1.5. Soluții de implementare a modelului OSI

Deși modelul OSI prevede împărțirea funcțiilor pe șapte nivele, numărul de nivele într-o implementare efectivă în sistemul de operare ca și diviziunea funcțiilor pe nivele poate diferi într-o anumită măsură. Se prezintă în acest subcapitol două implementări posibile ale modelului OSI, conceptual diferite.

Prima implementare se prezintă în figura 2.1.5.1 (se face abstracție de gestiunea memoriei). Un proces poate comunica fie cu un proces local, fie cu unul de la distanță. Interfața 1 transmite un apel fie modulului de comunicații interprocese, fie modulului corespunzător nivelului sesiune, în funcție de localizarea procesului destinație: local sau la distanță; interfața 2 furnizează servicii de I/O.

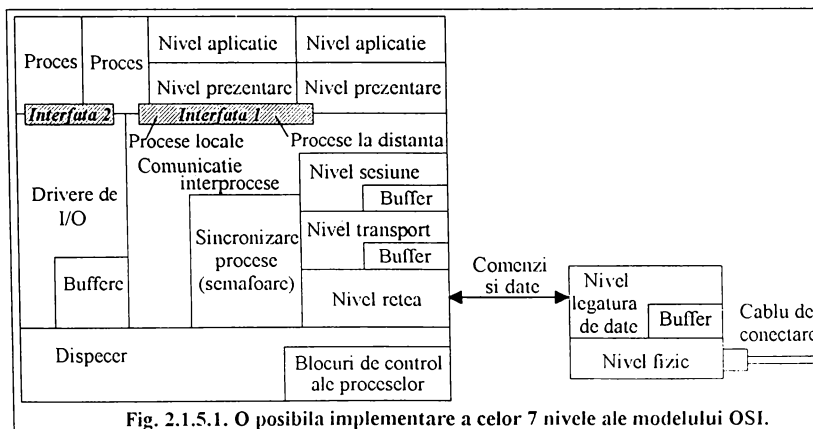


Fig. 2.1.5.1. O posibilă implementare a celor 7 nivele ale modelului OSI.

A doua implementare utilizează un procesor separat pentru funcțiile de comunicație, considerând că timpul necesar pentru acestea este considerabil și, ca urmare, reduce performanțele sistemului în raport cu aplicațiile. Arhitectura generală a unui astfel de sistem este prezentată în figura 2.1.5.2; subsistemul de comunicație este implementat pe o placă separată a calculatorului gazda. Aceasta conține în afară de procesorul rezervat pentru comunicație, o memorie locală, necesară pentru protocoalele software (informații de stare), o memorie partajată, necesară pentru bufferele de mesaje (care trebuie transmise de la un nivel la

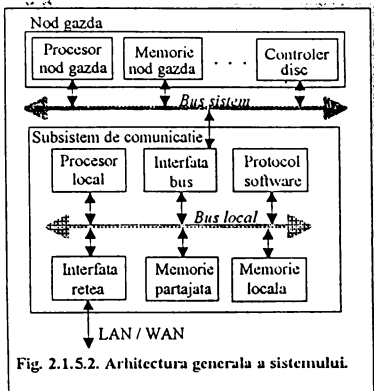


Fig. 2.1.5.2. Arhitectura generala a sistemului.

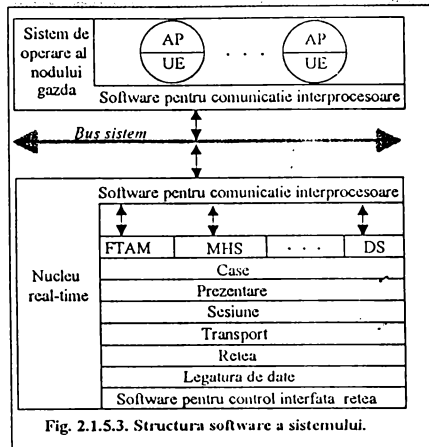


Fig. 2.1.5.3. Structura software a sistemului.

altul) și circuite de interfațare cu rețeaua de comunicații. La memoria partajată trebuie să aibă acces interfața cu rețeaua, procesorul local și procesorul nodului gazdă; în acest mod toate bufferele pot fi accesate direct de fiecare, astfel încât între nivele nu se mai transmit mesaje, ci se transmit direct pointeri la mesaje. Structura software-ului acestei arhitecturi este prezentată în figura 2.1.5.3; subsistemul de comunicație este implementat sub forma de procese separate: câte un proces pentru fiecare entitate de protocol plus un proces pentru gestiunea lor și funcții de timer. Procesele pot comunica între ele prin liste FIFO sau cutii poștale.

Comunicarea interprocese este implementată de un nucleu local real-time, responsabil și pentru planificarea la execuție a proceselor și tratarea întreruperilor ca și pentru gestionarea transferului de date la/de la interfața cu rețeaua. Comunicarea între subsistemele de comunicație și nodul gazdă se face printr-un modul software de comunicație interprocesor, pentru care există câte o copie în fiecare procesor. Pentru a se asigura sincronizarea transferului de informație între cele două procesoare, se folosesc întreruperi astfel: cînd procesorul nodului gazdă dorește, de exemplu, să transmită informații, scrie mai întâi informațiile în bufferul asociat entității aplicație respective și generează apoi o întrerupere pentru procesorul subsistemului de comunicație. Rutina de tratare a întreruperii determină entitatea SASE implicată, căreia îi transmite în coada de intrare pointerul la mesajul transmis.

2.1.6. Importanța modelului OSI pentru sisteme distribuite.

Ansamblul de protocoale OSI a fost creat în scopul interconectării sistemelor deschise; ele permit comunicarea între noduri fără a impune o fază de negociere a protocoalelor folosite de fiecare nod (acest lucru a fost scopul pe care și l-a propus ISO).

Intenția ISO/OSI nu a fost de a furniza standarde pentru sistemele distribuite. Se apreciază ([Gos91][Mul89]) că standardele internaționale conțin prea multe prevederi pentru a oferi o comunicație eficientă pentru sistemele de operare distribuite (vezi secțiunea 2.2); pe de altă parte acestea pot să nu fie necesare datorită următoarelor motive: (1) sistemele de operare distribuite sunt în general proiectate de o echipă de specialiști, lucrînd în strînsă colaborare (2) pe fiecare mașină se execută o copie a nucleului aceluiași sistem de operare (3) un proces client care trimite o cerere cunoaște exact formatul și structura cererilor acceptate de server, ca și al răspunsurilor de la server.

Totuși, se apreciază ([Mul89]) că se impune o considerabilă presiune, atît din partea unor organizații europene cît și din motive politice, asupra proiectanților de sisteme distribuite de a utiliza protocoale ISO/OSI. Pe de altă parte, efortul de standardizare pe plan internațional este continuat pentru a produce standarde pentru prelucrări distribuite deschise (Open Distributed Processing); efortul a fost susținut de ISO (International Standards Organization), ECMA (European Computer Manufacturers Association) și ANSA (Advanced Networked Systems Architecture Project).

2.2. Exemple de protocoale

2.2.1. Familia de protocoale TCP/IP

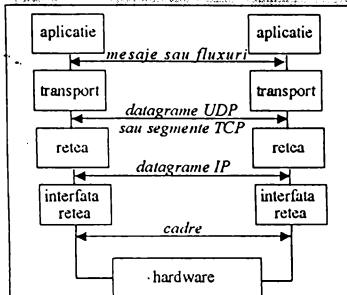


Fig. 2.2.1.2. Nivelele protocoalelor TCP/IP.

Deoarece familia de protocoale TCP/IP, este, la ora actuală utilizată pe scară largă, începând de la calculatoare personale și terminând cu supercomputere atât în LAN-uri cât și în WAN-uri, dar și în sisteme de operare distribuite experimentale ca suport pentru primitivile de comunicare interprocese, se prezintă în această secțiune elementele caracteristice, insistându-se pe protocolul IP versiunea 6 considerat ca o nouă generație de protocol interrețele.

Familia de protocoale TCP/IP și raportarea acestora la arhitectura ISO/OSI, este prezentată în figura 2.2.1.1, iar în figura 2.2.1.2 se prezintă structurarea protocoalelor pe nivele.

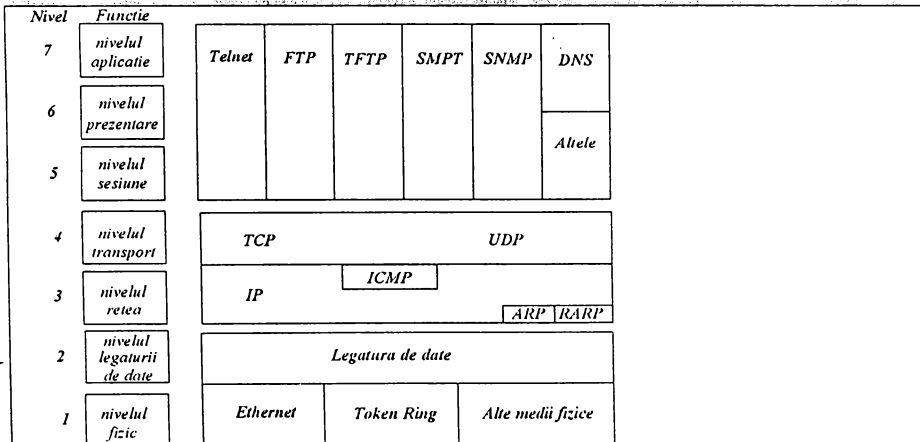


Fig. 2.2.1.1. Raportul OSI - TCP/IP.

ARP - Address Resolution Protocol - Mapeaza adresele IP în adrese fizice.

RARP - Reverse Address Resolution Protocol - Mapeaza adresele fizice în adrese IP.

IP - Internet Protocol - furnizeaza serviciul de livrare al pachetelor pentru TCP, UDP, ICMP.

ICMP- Internet Control Message Protocol - controleaza transmitia mesajelor de eroare și control întrenoduri s porti
Mesajele (datagrame IP) sunt generate de software-ul de rețea, nu de procesele utilizator.

UDP - User Datagram Protocol - protocol neorientat pe conexiune, care nu garanteaza livrarea sigură a datagramelor.

TCP - Transmission Control Protocol - protocol orientat pe conexiune, garanteaza o transmisie sigură full-duplex, sub forma de flux de octeti.

TELNET -Protocol pentru emularea unui terminal

FTP - File Transfer Protocol - protocol pentru transferul fișierelor între două noduri; ofera un set bogat de facilități pentru autentificare, conversie date, parcurgere directoare.

TFTP - Trivial File Transfer Protocol - protocol pentru transferul fișierelor între două noduri; nu ofera toate facilitățile protocolului FTP.

SMTP - Simple Mail Transfer Protocol - protocol pentru posta electronica; furnizeaza specificații pentru interacțiuni și formatul mesajelor de control.

SNMP - Simple Network Management Protocol - protocol pentru gestiunea rețelei; permite schimbul de informații între noduri relativ la configurație și stare. Informația disponibilă este definită de un set de obiecte gestionate (statistici relativ la trafic, starea rețelelor conectate, contor al erorilor întâlnite; conținutul tabelelor de rutare IP). Standardul MIB (Management Information Base) definește numele și semnificativitatea variabilelor pe care serverele SNMP le mențin.

DNS - Domain Name System - Un server DNS translateaza numele simbolice în adresa IP.

A. Protocoale IP. Noua generație de protocol IP, versiunea 6.

Protocolul IP versiunea 4, aflat în exploatare, în momentul curent, folosește pachete de maxim 64Ko; header-ul unui pachet, format dintr-o parte fixă de 20 octeți și o parte opțională de lungime variabilă, este prezentat în fig.2.2.1.3.

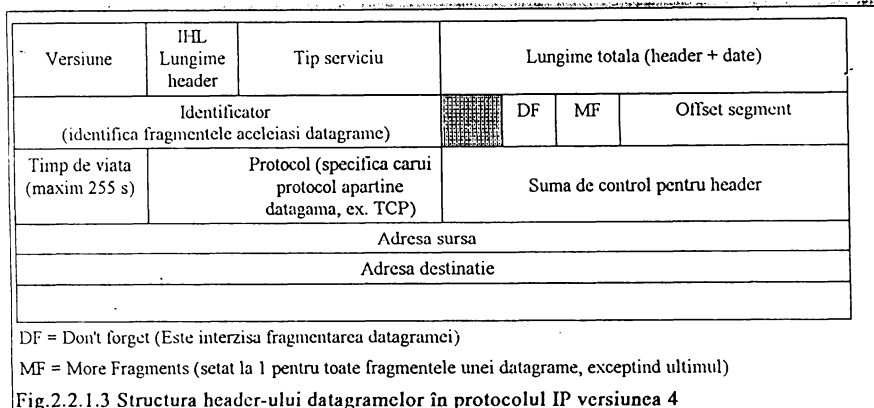
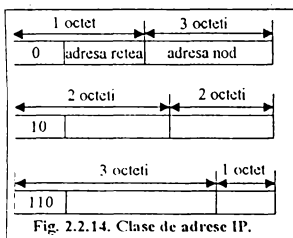
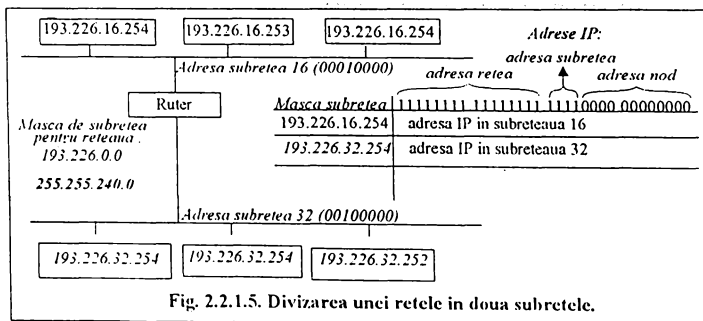


Fig.2.2.1.3 Structura header-ului datagramelor în protocolul IP versiunea 4



Exista trei clase de adrese IP (fig.2.2.1.4) și următoarele adrese rezervate: *adresele de rețea* (cele care au adresa nodului 0, de exemplu 129.0.0.0), *adresele broadcast* (cele care au biții din adresa nodului setați la 1) și *adresele loopback* (127.0.0.0 și 127.0.0.1).

Folosind rutere, rețelele pot fi divizate în subrețele (fie din motive de distanță, fie din motive de reducere a încărcării din cerințe de izolare a unei rețele sau de creșterea securității); în acest caz, porțiunea de adresă nod a adresei IP este divizată în două părți: <adresa IP> = <Adresa rețea> <Adresa subrețea> <Adresa nod>. O mască de subrețea indică modul de divizare al adresei de nod (fig.2.2.1.5).



Se prezintă, în continuare, în acest subcapitol, pe scurt, noul protocol Internet (Internet Protocol=IP) așa cum este documentat în RFC 1752 *The Recommendation for the IP Next Generation Protocol*. [Hin96]. Recomandarea a fost aprobată de IESG (Internet Engineering Steering Group) în 17 noiembrie 1994, și propusă ca standard RFC în decembrie 1995 de către IETF (Internet Engineering Task Force). Denumirea oficială aleasă a fost Internet Protocol version 6 (IPv6, versiunea protocolului IP folosita curent este 4); noua versiune, desemnată ca o revoluție dar nu ca o radicală schimbare față de IPv4, este construită pe arhitectura IPv4 și are scopul de a rezolva cerințele legate de extinderea impresionantă actuală a Internet-ului, problemele de rutare, ca și rezolvarea unor noi cerințe luate în considerare, precum securitatea, autoconfigurarea, servicii în timp real.

1. Cerințe impuse unui nou protocol Internet și realizate de IPv6

Există o serie de cerințe pe care trebuie să le respecte un nou protocol Internet, rezolvate în totalitate de IPv6, motiv pentru care acesta a fost desemnat ca standard al noii generații de protocoale Internet (de fapt, IPv6 este rezultatul evoluției mai multor propuneri de protocoale la IETF, între anii 1992-1994, (*The Internet Protocol-PIP, The Simple Internet Protocol-SIP, TP/IX, TCP and UDP with Bigger Addresses-TUBA, IP Address Encapsulation-IPAE, Simple Internet Protocol Plus-SIPP-, Common Architecture for the Internet-CATNIP.*[Hin96]). Iată câteva dintre aceste cerințe: (1)*compatibilitate cu IPv4*; IPv6 a fost proiectat astfel încât să poată interacționa cu IPv4. Mecanisme specifice protocolului IPv4 (adrese stil IPv4, reguli pentru checksum) au fost înglobate și în IPv6 în scopul realizării tranziției de la IPv4 la IPv6; (2)-*facilități extinse în ceea ce privește posibilitățile de adresare și de rutare*. IPv6 mărește dimensiunea adreselor IP, (de la 32 biți la IPv4) la 128 de biți, în scopul furnizării mai multor nivele ierarhice de adresare, a unui număr foarte mare de noduri adresabile și a unei mai simple configurări de adresă. IPv6 introduce un nou tip de adrese, *anycast address* pentru a identifica orice nod dintr-un set de noduri la care poate fi trimis un pachet cu o astfel de adresă de destinație, tehnică utilă pentru controlul traseului traficului. Pe de altă parte, IPv6 adaugă un câmp, *scope*, adreselor multicast îmbogățind astfel gama de adresare și permite folosirea unui *mechanism de adrese locale* care face posibilă instalarea *plug and play*. Structura adreselor adoptată de IPv6 este realizată astfel încât să poată exprima și adresele altor protocoale precum IPX și NSAP, facilitând migrarea de la aceste protocoale la IPv6.; (3)-*adaugarea de noi funcționalități precum servicii real-time*, autentificare, asigurarea securității, integrității și confidențialității datelor.

Toate cerințele de mai sus au fost impuse de amploarea actuală a Internetului (80 000 de rețele în ianuarie 1996, număr care se dublează la aproximativ 12 luni), gama foarte largă de computere pe care le cuprinde (de la PC-uri la supercomputere, majoritatea incluse în rețele locale) ca și de domenii de interes pe care le acoperă (începând de la institutii de învățământ -de la școli primare pînă la universități- și companii comerciale și terminând cu domenii nou intrate în sfera Internetului sau în curs de a intra precum rețele locale bazate pe frecvențe radio, pentru care se pun probleme specifice, datorită celorlalte ordine de mărime mai mici ale vitezei de transmisie impuse de disponibilitatea frecvențelor, alocarea spectrului, ratei de erori crescute sau canalele de televiziune, video la cerere, care având în vedere diminuarea diferențelor dintre computere și televiziune tind să devină o imensă piață pentru Internet ce va necesita astfel un protocol dotat cu mari posibilități de rutare și adresare ca și de autoconfigurare [Hin96]).

2. Structura pachetelor Header IPv6

IPv6 este proiectat ca fiind succesorul lui IPv4; protocolul constă în două părți: *header-ul de baza*, prezentat în figura 2.2.1.6. și *header-urile extensiilor IPv6*. Multe dintre extensii nu sunt examinate sau procesate de ruterele de pe traseul pachetelor; aceasta se realizează numai la sosirea la destinație (spre deosebire de IPv4, unde fiecare ruter era obligat să analizeze toate opțiunile). O altă îmbunătățire față de IPv4 este aceea că dimensiunea header-elor extensiilor IPv6 este variabilă, iar dimensiunea lor totală într-un pachet nu este limitată la 40 de octeți. Pentru a menține alinierea la multiplu de 8 octeți, header-urile extensiilor au întotdeauna o dimensiune multiplu de 8 octeți. Header-urile extensiilor definite pînă la ora actuală se referă la următoarele opțiuni IPv6: rutare, fragmentare și reasamblare, autentificare, securitate date, opțiunea *hop-by-hop*, opțiuni pentru destinație (care urmează să fie analizate la nodul destinație).

Versiune	Prioritate	Eticheta traseu	
Lungime		Header-ul următor	Limita hop
Adresa sursa (16 octeți)			
Adresa destinație (16 octeți)			

Fig.2.2.1.6. Formatul unui header IPv6

Versiune=cîmp de 4 biți în care se înscrie numărul versiunii (6).

Prioritate =cîmp de 4 biți care conține valoarea priorității.

Eticheta traseu =cîmp de 24 de biți; semnificația sa este legată de facilitatea *calității serviciului*.

Lungime =cîmp de 16 biți care conține lungimea restului pachetului (care urmează după header) exprimată în octeți, ca un întreg fără semn.

Header-ul următor =cîmp de 8 biți care identifică printr-un selector tipul header-ului care urmează imediat după header-ul IPv6.

Limita hop =cîmp de 8 biți care conține un întreg fără semn ce se decrementează la fiecare nod care transmite mai departe acest pachet. Pachetul este descărcat dacă limita atinge valoarea 0.

Adresa sursă=cîmp de 128 de biți care conține adresa nodului care a emis inițial pachetul.

Adresa destinație =cîmp de 128 de biți care conține adresa nodului destinație al pachetului (dacă este prezent un header pentru opțiunea de rutare, acest cîmp poate să nu conțină adresa nodului final).

3. Adrese IPv6

IPv6 asociază adresele (exprimate pe 128 biți) interfețelor (sau seturilor de interfețe) și nu nodurilor. O aceeași interfață poate avea asignată mai multe adrese IPv6, de orice tip. Sunt definite următoarele trei tipuri de adrese:

A. adrese unicast: este cazul în care o adresă corespunde unei singure interfețe, astfel încât un pachet trimis la o adresă unicast este recepționat numai de interfața corespunzătoare acelei adrese. Există mai multe formate de adrese unicast asignabile în IPv6. Acestea sunt: adrese unicast *global_provider_based* (tip 1), *geographic_based* (2), *NSAP* (3), *IPX* (4), *size_local_use* (5), *link_local_use* (6) și adrese unicast pentru noduri IPv4 (7).

Adresele tip 1 se folosesc pentru comunicațiile globale. Formatul lor este prezentat în figura 2.2.1.7.:

010	REGISTRY ID	PROVIDER ID	SUBSCRIBER ID	SUBNET ID	INTERFACE ID
3 biți	n biți	m biți	o biți	p biți	125 - n - m - o - p biți

Fig. 2.2.1.7.. Formatul unei adrese unicast tip 1.

Adresele unicast tip 5 și 6 au semnificație numai pentru rutare locală (la nivel de SUBNET sau SUBSCRIBER). Formatul lor este prezentat în figurile 2.2.1.8. și 2.2.1.9.:

111111010	0	SUBNET ID	INTERFACE ID
10 biți	n biți	m biți	118-n-m biți

Fig. 2.2.1.8. Formatul unei adrese unicast tip 5.

111111010	0	INTERFACE ID
10 biți	m biți	118 - n biți

Fig. 2.2.1.9. Formatul unei adrese unicast tip 6.

Pentru ambele tipuri de adrese (5 și 6), INTERFACE ID este un identificator unic în domeniul în care este utilizat; în majoritatea cazurilor se utilizează o adresă pe 48 de biți IEEE-802 MAC; cimpul SUBNET ID identifică o subrețea specifică, astfel încât combinația SUBNET ID + INTERFACE ID furnează o adresă ce va fi utilizată local într-o rețea particulară tip Internet, fără a fi nevoie de vreo alocare de adrese anterioară. Aceasta utilizare locală permite organizațiilor neconectate încă în Internet să opereze fără a avea nevoie de un prefix de adresă din spațiul global de adrese Internet; dacă, ulterior, aceste organizații se conectează la rețeaua globală Internet, pot utiliza în continuare cimpurile SUBNET ID + INTERFACE în combinație cu un prefix global, REGISTRY ID + PROVIDE ID + SUBSCRIBER ID, pentru a genera o adresă globală.

Compatibilitatea IPv6 cu IPv4 în ceea ce privește adresele se realizează prin prevederea formatului de tip 7 (adrese IPv6 compatibile cu IPv4) prezentat în figura 2.2.1.10., dar și printr-un alt format de adresă IPv6, folosit numai pentru a reprezenta adresele nodurilor IPv4 (care nu suportă IPv6) ca adrese IPv6; acest tip de adresă este numit "*IPv4-mapped IPv6 address*" și reprezentat în figura 2.2.1.11.

0000...0000	0000	Adrese IPv 4
80 biți	16 biți	32 biți

Fig. 2.2.1.10. Adrese IPv 6 compatibile IPv 4.

0000...0000	FFFF	Adresa IP. 4
80 biți	16 biți	32 biți

Fig. 2.2.1.11. Adrese IPv 6 mapate IPv 4.

B. Adrese anycast; este cazul în care o adresă se asignează unui set de interfețe (în mod tipic aparținând unor noduri diferite). Un pachet trimis la o adresă anycast va fi recepționat de una dintre interfețele desemnate, și anume de cea mai "apropiată", măsurarea distanței fiind făcută de protocoalele de rutare.

Adresele anycast se alocă din spațiul de adrese unicast, utilizând oricare din formatele descrise mai sus. Astfel adresele anycast sunt sintactic indistingibile de cele unicast; în schimb, atunci când o adresă unicast se atribuie mai multor noduri (transformându-se în adresă anycast), nodurile respective trebuie configurate explicit pentru a "cunoaște" că sunt noduri anycast.

Adresele anycast, atunci când sunt utilizate într-o secvență de rutare, permit unui nod să-și aleagă traseul pe care să-l urmeze pachetul (*source selected policies*); în acest caz, o adresă anycast corespunde unui furnizor de servicii Internet (PROVIDER ID).

C. Adrese multicast; este cazul în care o adresă se asignează unui grup de interfețe. Orice interfață poate aparține la oricâte grupuri multicast. Formatul unei adrese multicast este descris în figura 2.2.1.12.:

11111111	000T	SCOP	GROUP ID
8 biți	4 biți	4 biți	112 biți

Fig. 2.2.1.12. Formatul unei adrese multicast.

Semnificația cimpurilor este următoarea: - T indică o asignare permanentă (T=0) sau temporară (T=1). - SCOP selectează printr-o valoare pe 4 biți gama de vizibilitate a adresei multicast. Sunt recunoscute următoarele codificări: 0 = Rezervat, 1 = Vizibilitate la nivel de nod, 2 = Vizibilitate la nivel de link, 3,4 Neasignat, 5 = Vizibilitate la nivel de site, 6,7 = Neasignat, 8 -D=vizibilitate la nivel de organizație, E = Vizibilitate la nivel global, F = Rezervat.

În IPv6 nu există adrese broadcast; funcția lor a fost preluată de multicast.

B. Protocolul TCP

Protocolul TCP acceptă mesaje de lungime arbitrară pe care le transmite în segmente de maxim 64Ko; fiecare octet transmis are asociat un număr de secvență de 32 biți. Structura unui segment TCP este prezentată în fig.2.2.1.13. Elementele caracteristice ale protocolului sunt prezentate în fig.2.2.1.14 printr-o comparație cu protocolul de clasă 4 de la nivelul transport OSI (TP4 OSI).

32 biți							
Port sursa				Port destinație			
Numar de secventa (pentru fiecare octet)							
Numar de confirmare							
Lungime header TCP	URG	ACK	EOM	RST	SYN	FIN	Fereastra
Suma de control				Pointer de urgenta			
Optiuni							
Date							

Fig. 2.2.1.13. Structura unui segment TCP

Semnificatia indicatorilor:

URG= setat la 1 dacă se utilizează Pointer de Urgență; acesta indică deplasamentul față de octetul cu numărul de secvență curent la care se găsesc datele urgente (expeditiv este termenul OSI) datele urgente sunt cele care trebuie transmise înaintea celor bufferate, de exemplu CTRL/C, CTRL/D, CTRL/S

ACK= setat la 1 indică folosirea câmpului Numar de confirmare

EOM= indică Sfârșit de mesaj

RST= folosit la resetarea unei conexiuni (de exemplu în cazul recepționării unui SYN duplicat sau a căderii unui nod)

SYN= folosit pentru stabilirea conexiunilor

FIN= folosit la încheierea unei conexiuni

Semnificatia porturilor:

Adresele porturilor sursa, destinație plus adresele IP corespunzătoare formează o conexiune. Asocierea unui proces la o pereche (port, adresa IP) se realizează cu ajutorul unui *soclu* (socket). O conexiune cuprinde 2 socluri.

Fig. 2.2.1.14. Elemente caracteristice ale protocolului TCP

Element caracteristic	TCP	OSI TP4	Semnificatie
Numar de tipuri de segmente		1	9 În TCP se folosește același header la toate segmentele. În OSI se folosesc 9 tipuri de header
Conexiuni permise între aceleași două adrese TSAP		1	2 În TCP o conexiune se identifica printr-o pereche de TSAP, astfel încât este permisă o singură conexiune full duplex; în OSI sunt permise 2 conexiuni full-duplex independente între aceleași două TSAP
Format adrese	adresa pe 32 biți	nedefinit	deosebire în formatul adreselor
Calitatea serviciilor	opțiuni specifice	negociere	În TCP nu există negocierea calității serviciilor din OSI
Date în cererea de stabilire conexiune	nepermise	permise	În OSI se permite deschiderea conexiunii, de exemplu, în funcție de valoarea unei parole, precizată în câmpul de date din CR; în TCP nu se permite transfer de date la stabilirea conexiunii
Flux	la nivel de octet	la nivel de mesaj	TCP furnizează un flux de octeți; OSI un flux de mesaje
Date importante	date urgente	date expeditiv	TCP folosește câmpul Pointer de urgenta, pentru a indica, un număr de octeți, urgent de transmis, din segmentul curent. OSI are două fluxuri independente de mesaje, multiplexate, date obisnuite și expeditiv.

Control explicit al fluxului	Întotdeauna	Uneori	TCP utilizează un control explicit al fluxului permitând receptorului să însoțească o confirmare cu o fereastră care specifică domeniul de numere de secvență în care transmitatorul poate furniza date fără a permite o nouă confirmare
Confirmări cumulate	Da	Nu	TP4 poate utiliza o schemă cu credite
Închiderea conexiunii	Fără pierdere de date	Abrupt	Recuperarea datelor după o închidere abruptă a unei conexiuni se face în nivelul sesiune în OSI; în TCP se folosește un protocol cu confirmare în 3 pași pentru închiderea unei conexiuni

În TCP stabilirea și conexiunilor se face printr-un algoritmul cu confirmare în 3 pași conform diagramei de stări din fig. 2.2.1.15. Stabilirea unei conexiuni presupune două funcții de deschidere conexiune: una pasivă de așteptare a unei cereri de stabilire conexiune și una activă de cerere de stabilire conexiune. Diagrama este asemănătoare cu cea prezentată în secțiunea 2.1.3.

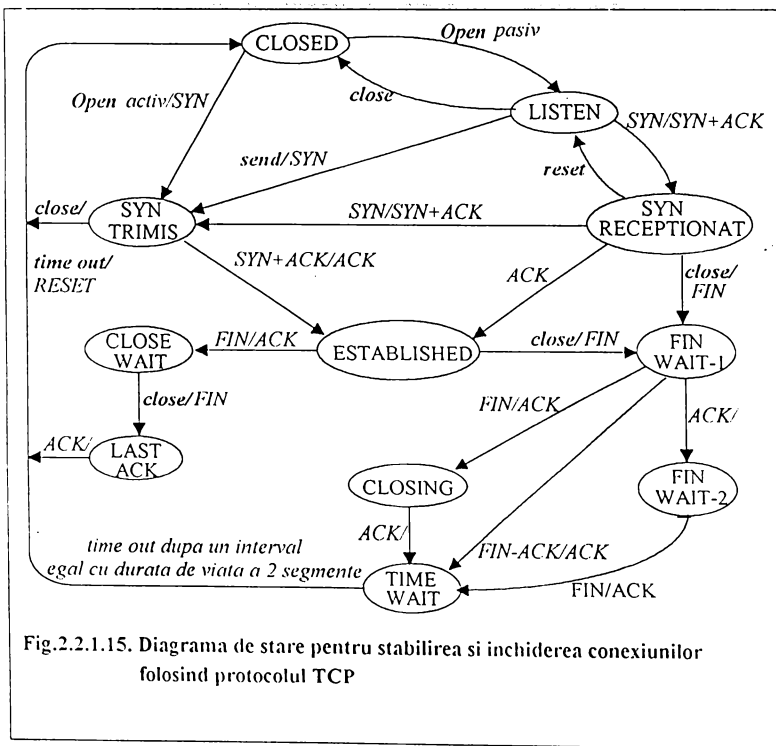
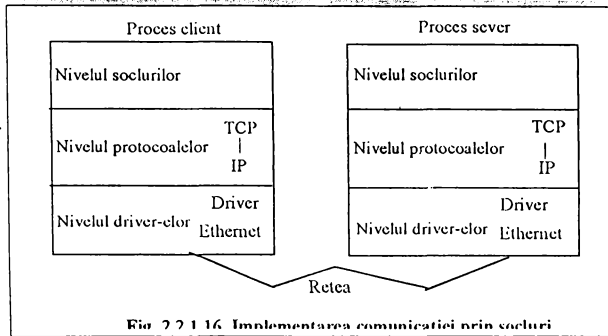


Fig.2.2.1.15. Diagrama de stare pentru stabilirea și închiderea conexiunilor folosind protocolul TCP

C. Comunicarea prin socluri în rețele UNIX și Microsoft Windows (Win32)

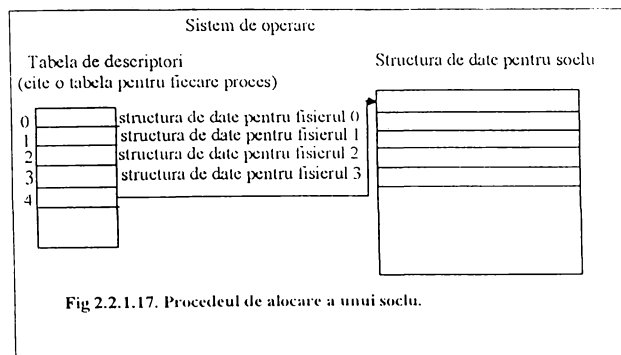
Socurile sunt un mecanism de comunicație interprocese în rețele UNIX, introdus de 4.3 BSD ca o extensie sau generalizare a comunicației prin conducte ("pipes"). Conductele facilitează comunicarea între procese de pe același sistem; soclurile facilitează comunicarea între procese de pe sisteme diferite sau de pe același sistem. O dată cu introducerea comunicării prin socluri și în System V Release 4, soclurile au devenit cea mai importantă interfață de programare (API) pentru rețelele Unix. Deoarece foarte multe companii de calculatoare precum Sun Microsystem, Digital Equipment Corporation și Tektronix au adoptat această interfață de programare soclurile au devenit *de facto* un mecanism standard de comunicație interprocese în rețele Unix. Structura subsistemului de comunicație prin socluri, într-un sistem de operare Unix, este prezentată în figura 2.2.1.16 [Bac89].



Subsistemul conține trei părți: nivelul soclurilor, nivelul protocoalelor și nivelul driverelor. Nivelul soclurilor furnizează interfața dintre apelurile sistem și nivelele inferioare. Nivelul protocoalelor conține modulele utilizate pentru comunicație (TCP și IP), iar nivelul driverelor conține driverele care controlează plăcile de conectare în rețea.

Comunicația între procese prin socluri utilizează modelul client-server (vezi 3.1.1.): un proces server (*listener*) ascultă pe un soclu (un capăt al unei linii de comunicație între două procese) - iar un proces *client* comunică cu serverul, care poate fi situat pe un alt nod, la celălalt capăt al liniei de comunicație, pe un alt soclu. Nucleul sistemului de operare memorează date despre linia de comunicație deschisă între cele două procese [Mus96j].

Se poate face o analogie între apelurile sistem pentru lucrul cu fișiere în UNIX și structurile de date pe care le generează nucleul pentru gestiunea fișierelor cu apelurile sistem pentru comunicația între procese prin socluri și structurile de date pe care le menține nucleul relativ la liniile de comunicație. Astfel, analog descriptorilor de fișiere care se obțin de la sistemul de operare ca urmare a unor apeluri *open()* sau *creat()* există descriptorii de socluri care se obțin ca urmare a unui apel sistem *socket()*. Nucleul Unix alocă descriptorii de socluri în aceeași tabelă de descriptori în care alocă descriptorii de fișiere [Bac89]. Deci o aplicație nu poate avea un descriptor de fișier și un descriptor de soclu de aceeași valoare. Când o aplicație

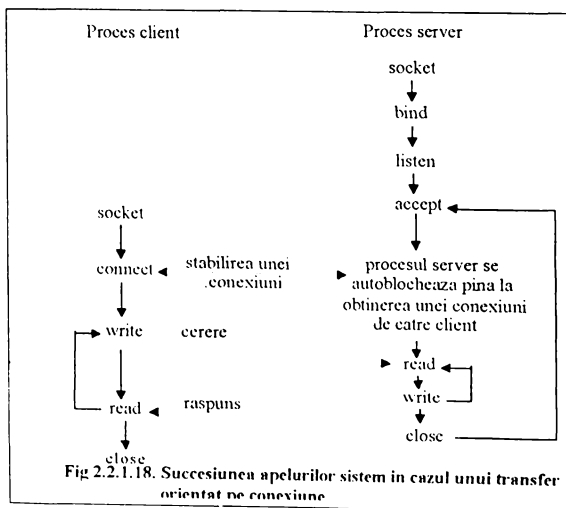


apelează `socket()`, nucleul aloacă o structură de date pentru noua linie de comunicație și completează într-o intrare liberă din tabela de descriptori un pointer la această structură de date. Figura 2.2.1.17 exemplifică acest procedeu. În urma apelului `socket()`, structura de date creată este completată doar parțial, aplicația fiind obligată să apeleze în continuare și alte funcții pentru a completa în întregime această structură de date înainte de a utiliza socketul [Mus96][Mus91b].

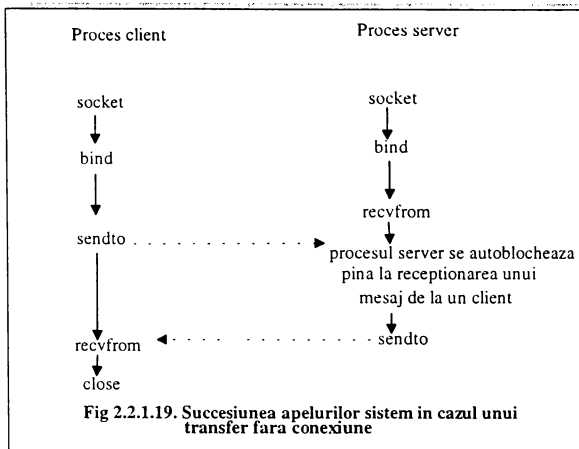
Socketurile care au caracteristici comune în ceea ce privește comunicația, precum convenții de nume, formatul adreselor utilizate de protocoalele de comunicație se grupează în *domenii* (familii). Există astfel, pentru comunicația în rețele, *Internet domain* și *Xerox NS domain*, iar pentru comunicația între procese de pe același sistem Unix, *Unix system domain*. Fiecare socket are un tip specificat care determină modul de transfer al datelor, în mod obișnuit, printr-un *circuit virtual* (transfer orientat pe conexiune, STREAM SOCKET în terminologia Berkeley), sau prin *datagrame* (transfer neorientat pe conexiune). Un circuit virtual permite transferul garantat (bidirecțional) al datelor, fără duplicarea și cu secvențierea pachetelor. Datagramele nu garantează nici transmisia în secvență a pachetelor nici nduplicarea lor, dar au avantajul că nu necesită stabilirea conexiunii. Mecanismul de comunicație prin socketuri asociază implicit anumite protocoale de comunicație pentru fiecare combinație *domeniu -tip socket* validă. De exemplu, protocolul TCP furnizează un transfer orientat pe conexiune (circuit virtual), iar UDP furnizează un transfer bazat pe datagrame.

Cuplul de informații {*domeniu, tip socket, protocol*}, relative la un socket se furnizează sistemului ca parametri ai apelului `socket()` care returnează un descriptor de socket. Socketul astfel creat nu are însă asociat nici un proces (nu i se asociază la crearea nici adresa sistemului local nici a celui de la distanță). Pentru a specifica adresa proprie la un capăt al liniei de comunicație între două procese, procesul server trebuie să execute apelul sistem `bind()`. Când se creează un socket nu se specifică dacă este activ (pentru a fi utilizat de un client), sau pasiv (pentru a fi utilizat de un server). Serverele care folosesc o transmisie orientată pe conexiune trebuie să execute (după `bind()`) apelul sistem `listen()` pentru a pune socketurile lor într-o stare pasivă, gata să primească cereri de stabilire a unei conexiuni. Majoritatea serverelor constau într-o buclă infinită în care preiau o cerere, o tratează și returnează un răspuns. Chiar dacă tratarea cererii durează foarte puțin, câteva milisecunde, se poate ca o altă cerere să sosească exact în acest interval. Pentru a se evita pierderea cererilor sosite la nivelul serverului, acesta trebuie să transmită apelului `listen()` un argument prin care se cere sistemului de operare să lanșeze cererile primite pe un socket. Astfel, argumentele apelului `listen()` specifică ce socket se setează ca pasiv și dimensiunea maximă a cozii de cereri pentru acel socket. În scopul recepționării cererilor de conectare de clienți, un server trebuie să execute în continuare un apel `accept()`. `Accept()` creează câte un socket nou pentru fiecare cerere de conexiune și returnează apelantului descriptorul noului socket creat. Serverul utilizează noul socket exclusiv pentru comunicația cu clientul respectiv, în timp ce pe vechiul socket acceptă noi cereri de conectare.

Referitor la client, acesta trebuie să execute un apel `connect()` pentru a stabili o conexiune cu un server. Clientul trebuie să furnizeze sistemului informații despre identitatea serverului (adresa IP și numărul portului, de exemplu). După stabilirea conexiunii, schimbul de informații dintre client și server se poate face cu apelurile `read()` / `write()`, iar în final, distrugerea conexiunii se face cu apelul `close()`. Figura 2.2.1.18 ilustrează secvența de apeluri sistem executată de un server și un client în cazul comunicării prin TCP.



Pentru o transmisie prin UDP secvența de apeluri este prezentată în figura 2.2.1.19.



În acest caz, clientul nu trebuie să se conecteze la server; el trebuie doar să trimită către server un mesaj sub forma unei datagramme cu ajutorul apelului sistem *sendto()* care primește ca parametru adresa sistemului destinație. Analog, serverul nu trebuie să accepte o conexiune ci doar să recepționeze mesaje datagramme de la un client, cu ajutorul apelului sistem *recvfrom()*, care blochează serverul până la primirea unui mesaj.

Specificațiile Windows Sockets definesc o interfață de programare pentru rețea în Microsoft Windows bazată pe modelul socketurilor Berkely, standardul de facto TCP/IP. Există însă câteva API-uri Windows Sockets care nu respectă întru-totul convențiile Berkely, datorită dificultăților de implementare în mediul Windows ; aceste deosebiri se referă la faptul că două funcții de lucru cu socketurile au fost înlocuite complet (*open()* și *ioctl()* cu *closesocket()* și *ioctlsocket()*), funcția *select()* a fost modificată în sensul că fiecare set de socketuri este totuși reprezentat de tipul *fd_set* dar acesta este implementat ca un tablou de tipul SOCKET (tip nou introdus în Windows pentru descriptorii socketurilor), codurile de eroare nu sunt accesibile via variabila *errno* ci via funcția *WSAGetLastError()*. Pe de altă parte au fost introduse o serie de API-uri noi, care, cu excepția *WSAStartup()* și *WSACleanup()* nu sunt obligatorii iar unele sunt disponibile și într-o formă compatibilă cu funcțiile Berkeley, dar sub forma *WSAAsyncXXX()* ele sunt recomandate în special datorită necesității respectării paradigmelor de programare Windows.

Aplicațiile descrise în anexele B4, B7 utilizează specificațiile Windows Sockets.

2.2.2. Familia de protocoale Novell IPX/SPX

Cele două protocoale IPX/SPX au fost introduse de Novell în produsele NetWare. IPX este protocolul de bază Novell NetWare și este recomandat de Novell ca protocolul de comunicație pentru aplicațiile client-server. Protocolul IPX este o implementare a protocolului Internetwork Datagram Packet (IDP) proiectat de Xeros Corporation. Este un protocol neorientat pe conexiune și suportă numai mesaje de tip datagrame. Pe de altă parte, protocolul SPX este orientat pe conexiune; livrarea pachetelor necesită stabilirea unei conexiuni, emisia/recepția mesajelor printr-un dialog adecvat și închiderea sesiunii. Astfel, se poate spune că interfața de programare în IPX corespunde nivelului rețea iar interfața de programare în SPX corespunde nivelului sesiune[Mus96j].

```

unsigned short  checksum;    /*high-low 1's complement*/
unsigned short  packetLen;   /*high-low*/
unsigned char   transportControl;
unsigned char   packetType;
unsigned long   destNetwork; /*high-low*/
unsigned char   destNode[6];
unsigned short  destSocket;  /*high-low*/
unsigned long   sourceNetwork; /*high-low*/
unsigned char   sourceNode[6]; /*high-low*/
unsigned short  sourceSocket; /*high-low*/
    
```

Fig. 2.2.2.1. Structura header-ului pachetelor IPX

(cîmpul destNode din structura IPX) este un număr pe 6 octeți care conține adresa fizică a nodului destinație și care poate fi setată la 0x FF FF FF FF FF FF FF pentru a trimite un pachet broadcast. Soclul destinație (cîmpul destSocket din structura IPX) este un cîmp pe 2 octeți conținând adresa soclului procesului destinație.

```

unsigned char   connectionControl; /*bit flags*/
unsigned char   dataStreamType;
unsigned short  sourceConnectionID; /*high-low*/
unsigned short  destConnectionID; /*high-low*/
unsigned short  sequenceNumber; /*high-low*/
unsigned short  acknowledgeNumber; /*high-low*/
unsigned short  allocationNumber /*high-low*/
    
```

Headerul unui pachet SPX are în plus cele 7 cîmpuri de mai sus.

Fig. 2.2.2.2. Structura header-ilor pachetelor SPX

Un număr de retransmisii neconfirmate se presupune că aplicația destinație nu mai "ascultă" sosirea pachetelor și întrerupe conexiunea.

Înainte de a utiliza fie funcțiile IPX, fie funcțiile SPX, aplicațiile trebuie să pregătească și să completeze o structură de date numită EventControlBlock (ECB). Numai după aceea se pot apela funcțiile IPX sau SPX cu pointeri către ECB, care conține adresa pachetului IPX sau SPX.

Un sistem client-server tipic constă într-un proces server care deschide un soclu (Novell IPX API: *IPXOpenSocket*) și așteaptă să recepționeze un pachet de la orice client (Novell IPX API: *IPXListenForPacket*). Serverul poate de asemenea să "anunțe" prezența sa (Novell Bindery API: *AdvertiseService*) astfel încît atunci cînd un client trebuie să se conecteze la server, el poate să găsească serverul (Novell Bindery API: *ScanBindery*) și să determine adresa sa (Novell IPX API: *IPXGetLocalTarget*). În final, clientul poate trimite o cerere de serviciu la server (Novell IPX API: *IPXSendPacket*) [Mus95c][Mus96j][Nov90a][Nov90b].

Aplicațiile din anexele B3, B, C1, C2, C3, C4 utilizează protocolul de comunicație IPX; exemplul din anexa C4 constă într-o aplicație client-server într-o rețea Novell.

Structura unui pachet IPX este identică cu structura unui pachet Xeros Internetwork Standard (XNS) care constă într-un header de 30 octeți urmat de 0 pînă la 576 octeți în zona de date (fig 2.2.2.1.).

Un mesaj poate fi trimis către orice nod din rețea prin plasarea lui în zona de date a unui pachet. Fiecare pachet trebuie să specifice adresa de destinație, sub forma unui cuplu {rețea destinație, adresă nod, soclu}. Rețeaua destinație (cîmpul destNet din structura IPX) este un număr pe 4 octeți care specifică numărul rețelei locale și este 0 pentru un nod pe aceeași rețea locală ca și nodul emițător. Adresa nodului destinație

Un pachet SPX conține în plus în header 12 octeți pentru a urmări secvența pachetelor și numărul confirmărilor (fig 10) și poate conține pînă la 534 octeți de date. Aplicațiile care transmit pachete SPX trebuie să stabilească conexiuni cu aplicațiile destinație; pachetele transmise pentru care nu s-a primit confirmare sunt retransmise după expirarea unui interval de timeout. După un

2.3. Aspecte specifice ale suportului de comunicație pentru sistemele de operare distribuite

S-a arătat în secțiunea 2.1.3 că, în ceea ce privește sistemul de operare, importanța cea mai mare o prezintă nivelul transport, care poate oferi funcții apelabile direct din procese. Serviciile și mecanismele necesare pentru a implementa un protocol de transport pot fi clasificate, concluzionând secțiunea 2.1.3, astfel:

(1) *Mecanisme pentru gestiunea conexiunilor* care implică aspecte capăt-la-capăt în ceea ce privește: (a) alocarea, sincronizarea și dealocarea de resurse; (b) negocierea unor resurse și a unor moduri de desfășurare a operațiilor. Gestiunea unei comunicații *sigure* poate fi obținută prin următoarele mecanisme: (a) *utilizarea unor algoritmi cu confirmare în 2 sau 3 pași* pentru deschiderea și închiderea conexiunilor, algoritmi apreciați însă drept costisitori ca timp, deoarece necesită schimburi de mesaje (b) *folosirea timerelor*.

(2) *Mecanisme pentru controlul erorilor* care implică aspecte legate de pierderea, duplicarea sau deteriorarea mesajelor; pot fi utilizate următoarele mecanisme: (a) *sume de control* pentru detectarea mesajelor eronate; pachetele eronate sunt ignorate iar recuperarea lor se face prin utilizarea de confirmări și retransmisii; (b) *confirmări explicite* pentru detectarea mesajelor pierdute; o problemă este determinarea time-out-ului; (c) *confirmări implicite* pentru detectarea mesajelor pierdute; mecanismul se bazează pe presupunerea că un răspuns al unui server este confirmat de următoarea cerere sau neconfirmat de time-out; (d) *numerotarea mesajelor și a pachetelor în mesaje*, pentru a detecta mesajele lipsă sau cele care nu sunt în secvență.

(3) *Controlul fluxului* care implică un acord aprioric sau o negociere capăt la capăt asupra mecanismului folosit; pot fi utilizate: (a) ferestre glisante explicite; (b) ferestre glisante implicite (de exemplu cu o singură cerere/răspuns în așteptare la un moment dat); (c) presupuneri privind rata de transmisie și de recepție.

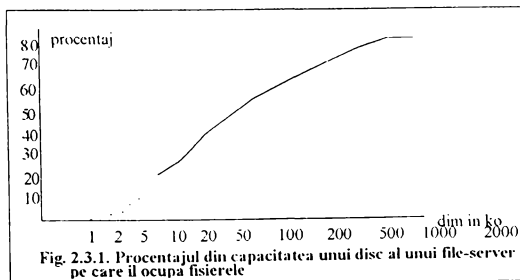
(4) *Stabilirea dimensiunii* mesajelor care diferă de la un protocol la altul și implică segmentarea mesajelor. Se apreciază ([Gos91][Mul89]) ca pierderile de timp legate de controlul erorilor și fluxului pachetelor în mesaje și de segmentarea/reasamblarea trebuie luate în considerare.

Ca aspecte specifice în vederea implementării acestor servicii și mecanisme pentru sistemele de operare distribuite pot fi enumerate următoarele:

(1) Primul aspect de care trebuie să se țină seama este că sistemele distribuite sunt implementate în general în LAN-uri, astfel încât cea mai mare parte a traficului de mesaje rămâne într-o singură rețea locală; chiar dacă sistemul cuprinde mai multe LAN-uri interconectate, proiectanții sistemului încearcă, în general, să facă în așa fel încât traficul dominant să fie cel din rețeaua locală.

(2) Un alt aspect important pornește de la aprecierea dimensiunii datelor care se vor transfera în sistemul distribuit; de exemplu, pentru minimizarea traficului în rețea, bazele de date moderne folosesc tehnici pentru alocarea optimă a datelor și alegerea judicioasă a locațiilor care vor rezolva cererile lor. Plecând de la ideea (firească) că dimensiunea datelor transferate într-un sistem distribuit este legată de dimensiunea fișierelor existente în sistem (în ultima instanță ceea ce se transferă sunt fișiere de date, parțial sau total) și realizându-se o statistică a dimensiunii fișierelor și capacității de memorare ocupată de ele în [Mul89] se ajunge la următoarea concluzie: dacă fișierele unui file-server sunt citite cu aceeași frecvență, 95% din timpul de transfer se consuma pentru citirea unor fișiere cu dimensiunea mai mare decât a unui pachet (presupunând că un pachet are o dimensiune sub 2Ko). Acest lucru se poate constata din fig.2.3.1 unde numai 5% din capacitatea de memorare este ocupată de fișiere sub 2Ko; fișierele mai mari de 24Ko care reprezintă 10% din totalul fișierelor ocupa mai mult de 60% din capacitatea de memorare, iar fișierele mai mari de 1Mo, deși puțin (0,1% din totalul fișierelor) ocupa totuși 17% din capacitate.

Concluzia este că proiectarea unui protocol de transport trebuie făcută astfel încât să dea rezultate atât pentru transferuri voluminoase cit și pentru transferuri mici. Ca urmare, timpul necesar pentru pregătirea unui transfer (rezervarea resurselor pentru a obține o productivitate mai bună) trebuie să fie mic pentru a nu contribui la scăderea eficienței transferurilor mici și pe de altă parte eficiența transmisilor voluminoase trebuie să fie mare deoarece aceste transmisii ocupa o fracțiune mare din timpul total de transfer. Rezultă că sunt necesare două tipuri de protocoale: unul pentru transferuri mici, cu *întârzieri minime* și altul pentru transferuri voluminoase cu *productivitate maximă*. Este exact ceea ce s-a standardizat: două clase de protocoale, una orientată pe conexiuni (pentru transferuri voluminoase) și alta neorientată pe conexiuni (pentru transferuri mici). Exista însă și



sisteme distribuite [Mul89][Gos91] pentru care s-a proiectat un singur protocol de transport eficient atât pentru mesaje mici cât și pentru mari.

(3) Legat de aspectul (2) este și următorul: o mare parte din comunicații se datorează transmiterii unor cereri de execuție operații la distanță (de exemplu se cere unui server de la distanță crearea unui fișier sau a unei ferestre) și se așteaptă rezultatul, care ca și cererea este un mesaj de dimensiune foarte mică; deoarece, pentru continuarea execuției, procesul client este blocat în așteptarea rezultatului, este clar că trebuie minimizată întârzierea în preluarea rezultatului.

(4) Un alt aspect este cel legat de tratarea erorilor. Erorile de tip date corupte, pierdute sau afară din secvență pot fi recuperate prin retransmisii; erorile grave precum căderea unui nod, a unei linii de comunicații pot fi recuperate mult mai greu. Primul tip de erori este tratat în OSI la nivelul transport; erorile din al cel de-al doilea tip sunt recuperate la nivelul sesiune. În sistemele distribuite, al caror suport sunt rețele locale, clasificate de cele mai multe ori în tipul A de rețele (sigure), frecvența erorilor din ambele tipuri de mai sus este cam de același ordin de mărime și de valoare, relativ mică; de aceea nu se tratează la nivele diferite. La nivelul cel mai înalt, cel de aplicație, se va testa întotdeauna execuția cu succes a unei operații la distanță.

În vederea proiectării unor protocole de comunicații rapide, se prezintă următoarele tehnici:

(1) *Integrare protocol-planificator.* S-a reliefat mai sus, la (2) și (3) importanța minimizării timpului de obținere a rezultatului unei cereri; în [Mul89] se arată că în sistemul de operare Amoeba, într-o stație SUN conectată într-o rețea Ethernet, din aproximativ 1.4ms, timpul de obținere a răspunsului la cerere nula, 400μs reprezintă transmitia mesajelor pe liniile de comunicație plus depunerea lor în memorie și generarea întreruperilor, alte 500μs timpul de execuție al protocolelor și alte 500μs timpul necesar comutării procesorului la procesele client și server. Măsurătorile au fost făcute în situația în care protocolul este încapsulat într-un proces de rețea separat; în aceasta situație procesorul trebuie comutat la fiecare recepție și emisie. O implementare eficientă impune includerea codului protocolului în nucleul sistemului de operare și parțial chiar în rutinele de întrerupere dacă este posibil, reducând astfel la minim comutarea contextului.

(2) *Controlul fluxului între un emițător rapid și un receptor lent* poate fi făcut prin introducerea unei mici întârzieri înainte de transmiterea unui pachet; aceasta întârziere poate fi determinată chiar în timpul transmisiilor. Metoda a fost implementată în protocolul VMTP (Versatile Message Transport Protocol), care se bazează pe așa numitele pachete în rafală, *packet blasts*. În VMTP o datagrama mare este descompusă în grupuri de cel mult 32 pachete; fiecare pachet dintr-un grup este însoțit de un număr de secvență de 32 biți, dintre care numai unul are valoarea 1. La recepție, numerele de secvență sunt verificate printr-o operație OR; dacă toate pachetele au fost recepționate, toți biții rezultatului trebuie să fie 1; dacă sunt biți 0, pachetele corespunzătoare vor fi retransmise (*retransmisie selectivă*). Presupunem că, având un emițător mai rapid decât receptorul, la transmitia pachetelor dintr-o datagrama, fiecare al 2-lea pachet se pierde. Mecanismul de retransmisie va retransmisia printr-o altă rafală (blast), astfel încât ar urma să se piardă pachetele multiplu de 4 s.a.m.d. Pentru a ocoli acest lucru, VMTP lucrează cu o întârziere între pachete, *interpacket gap*, inițial această este setată la 0 de către emițător dar când se constată că s-au pierdut pachetele pare se mărește gap-ul până când nu se mai pierd pachete (s-a constatat că valoarea întârzierii este în jurul a câteva sute de microsecunde).

(3) *Utilizarea datagramelor mari și a unei facilități de scater-gather.* Permițând datagrame mari, se poate reduce semnificativ overhead-ul de transmisie. În plus operațiile de copiere ale header-ului și a datelor unui pachet într-un/dintr-un buffer pot fi eliminate dacă se permite ca acestea să fie în zone diferite de memorie (scater-gather); mai mult, se poate ține seama și de structura cererilor: de exemplu, Amoeba permite ca parametrii unei cereri pentru un server de fișier (cod operație, nume fișier, poziția în fișier) să fie plasate într-un buffer separat de bufferul de date. Evitarea operațiilor de copiere a mesajelor din/in zona de date a protocolului și a celei a receptorului/emițătorului poate fi făcută și cu ajutorul mecanismului de gestiune a memoriei, prin partajarea paginilor. Apar complicații legate de faptul că mesajele trebuie aliniate la granițe de pagini iar o pagină poate conține numai mesaje complete; aceasta înseamnă că compilatoarele trebuie să rezerve bufferele de mesaje respectând aceste condiții.

(4) *Utilizarea timer-elor de analiză.* Având în vedere următoarele aspecte ale utilizării timer-elor: (a) startarea și oprirea unui timer necesită timp. S-a aratat că mai multe timere pot fi implementate folosind un singur timer fizic și o lista sortată după valorile de time-out, astfel încât startarea și oprirea unui timer înseamnă operații de inserare și ștergere dintr-o lista, evident cu dezactivarea întreruperii de timer în timpul operației; (b) timerele expiră arareori, numai când mesajele sunt pierdute, la avaria unui nod sau a unei linii de comunicații; (c) valorile la care sunt setate timerele sunt utilizate pentru recuperarea erorilor: ele nu trebuie să fie extraordinar de exacte - se poate concluziona că timerele pot fi înlocuite cu un fir separat al nucleului, activat periodic, al cărui rol este de a analiza structurile de date ale protocolului (la care trebuie să poată avea acces); el compară starea curentă a protocolului cu cea anterioară și dacă nu se constată un progres, inițiază procedura de serviciu (de exemplu retransmisia unui pachet). Acest fir se execută cu o prioritate mică în background.

2.4. Protocoale cu caracter universal și protocoale cu caracter restrâns

Protocoalele de transport proces-către-proces continuă să constituie un domeniu important de cercetare în cadrul sistemelor distribuite avînd în vedere faptul că protocoalele de transport stau la baza implementării primitivelor de comunicație interprocese în sistemele de operare distribuite. Există două direcții în cadrul acestor cercetari: (1) Proiectarea și implementarea unor protocoale specifice unui sistem de operare distribuit; (2) Reexaminarea proiectării și implementării protocoalelor cu caracter universal.

Pentru ambele tipuri de protocoale există argumente pro și contra. Astfel, *protocoalele cu caracter universal* oferă avantaje precum portabilitate, independența aplicațiilor, dar și dezavantaje care provin din ignorarea aspectelor specifice ce apar în sistemele distribuite, prezentate în secțiunea anterioară; aceste dezavantaje constau într-o complexitate excesivă (protocoalele de transport la nivelul ISO/OSI sunt un exemplu în acest sens) cu repercursiuni asupra performanțelor, orientarea preponderentă pe variante de orientate pe conexiune, lipsa unui suport pentru implementarea unor facilități speciale în rețelele locale precum partajarea resurselor multicasting.

Protocoalele cu caracter restrans au ca principal avantaj performanțe mai bune, (întîrzieri mici, printr-o execuție eficientă) dar și dezavantajul aplicabilității numai într-un anumit tip de rețele și/sau arhitectură de aplicații. Aceste protocoale cu caracter restrîns se obțin în urma luării în considerare a următoarelor simplificări:

(1) minimizarea modului orientat pe conexiune; în acest sens, există următoarele soluții: (a) utilizarea unui protocol neorientat pe conexiune ori de cîte ori este cazul; (b) utilizarea așa numitelor protocoale cu conexiuni "lightweight", în cadrul carora o conexiune este deschisă de un server la recepția primei cereri a unui client, și închisă la recepția unei indicații de sfîrșit de transfer; (c) stabilirea unei conexiuni numai pe timpul necesar satisfacerii unei cereri

(2) minimizarea numărului de pachete schimbate la stabilirea conexiunilor și cofirmarea lor pe baza unor presupuneri asupra caracteristicilor de eroare ale rețelei

(3) minimizarea circulației confirmărilor prin considerarea recepției răspunsului la o cerere sau a recepției unei cereri imediat după o altă cerere drept confirmări (presupunînd că nu există pachete duplicate sau în afara secvenței).

În ceea ce privește *optimizarea protocoalelor cu caracter general* se enumeră următoarele soluții:

(1) Utilizarea unei alte arhitecturi pe nivele. În [Gos91], în urma analizei mai multor sisteme de operare distribuite contemporane se propune arhitectura din fig.2.4.1.

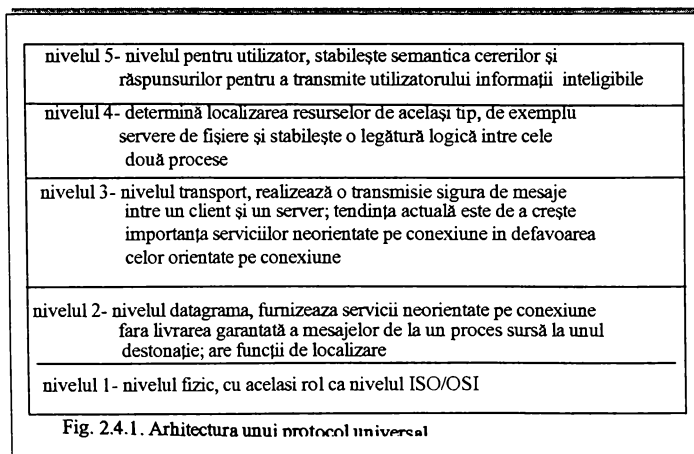


Fig. 2.4.1. Arhitectura unui protocol universal

(2) Implementarea protocoalelor parțial în sistemul de operare, parțial în drivere și hardware; procesele utilizator pot să interacționeze cu nivelul transport, celelalte 3 nivele mai inferioare sunt însă neaccesibile.

(3) Suport pentru transferul direct al datelor direct din spațiul de memorie utilizator în rețea și invers.

Capitolul 3.

PARADIGME ALE INTERACȚIUNII DINTRE PROCESE IN SISTEME DISTRIBUITE

În mod uzual, software-ul de sistem și de aplicație au o structură modulară, ierarhizată, pe nivele (fig. 3.1); software-ul de sistem ocupă nivelele inferioare, iar cel de aplicație nivelele superioare.

Într-un sistem distribuit, o aplicație poate fi structurată în mai multe procese, repartizate pe mai multe procesoare, dispuse pe mai multe noduri într-o rețea. Ca atare, se pot identifica două clase de interfețe între procese [Mul89]: interfețe între procese de pe nivele diferite și interfețe între procese din același nivel.

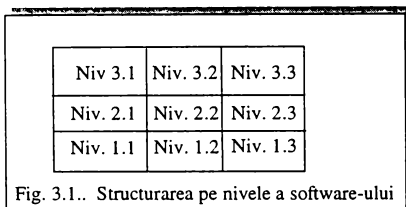


Fig. 3.1.. Structurarea pe nivele a software-ului

În sistemele tradiționale, nedistribuite, interfețele între nivele includ interfața cu sistemul de operare (incluzând sistemul de gestiune fișiere, ca un subset), interfețele cu diferite biblioteci de sistem și interfețele cu un sistem de gestiune ferestre. Mai mult chiar, în unele sisteme, interfața cu sistemul de gestiune ferestre apare ca o interfață între procese care comunică printr-o rețea (de exemplu, sistemul X-window, MIT). Pe de altă parte, fiecare nivel din structura generală din fig. 3.1 se poate spune că oferă servicii nivelului superior. Nivelul definește o anumită clasă de obiecte și operații asupra acestor obiecte. Astfel, sistemul de operare furnizează ca obiecte: procese, fișiere, directoare etc; sistemul X-window furnizează obiecte precum ferestre, hărți de biți etc; o bibliotecă matematică furnizează obiecte precum numere complexe, matrici etc. În mod tradițional, interfața dintre nivele a fost implementată procedural (interfețele cu bibliotecile sunt în general bazate pe apeluri de proceduri; interfețele cu sistemele de operare sunt dotate cu proceduri *stub*) sau sub formă de comunicare prin mesaje (Microsoft Windows). Într-un sistem distribuit, procesele din diverse noduri trebuie să comunice prin intermediul mesajelor transmise prin mediul fizic de comunicație. Interfața dintre procesele situate în noduri diferite poate avea însă același format ca și interfața dintre procesele locale. Performanțele subsistemului de comunicație între procese depinde, pe de o parte, de primitivele de comunicație între procese (interfața de comunicație) și pe de altă parte, de protocolul de comunicație, care furnizează aceste primitive. Cele mai cunoscute clase de primitive utilizate pentru comunicația între procese în sistemele de operare distribuite sunt: (1) *comunicarea prin mesaje*, care se asociază modelului client/server și la care fluxul de mesaje poate fi unidirecțional sau bidirecțional; (2) *apel de procedură la distanță (RPC)*, care oferă un mecanism puternic tipizat pentru transferul bidirecțional al informației; (3) *tranzacții*, care oferă un mecanism complex pentru executarea și sincronizarea de operații asupra unor obiecte partajate, asigurând recuperarea informației în cazul apariției de erori hardware și software; (4) *comunicare prin conducte (pipes)* ca o generalizare a modelului conductelor locale introdus de sistemele Unix.

Aceste clase de primitive permit implementarea a patru tipuri de procese într-un program distribuit: *filtre, clienți, serveri și perechi*. Un *filtru* este un proces care transformă datele; el primește datele în flux pe canalele sale de intrare, execută prelucrări asupra acestor valori și trimite apoi rezultatele la canalele sale de ieșire. Exemple binecunoscute sunt o mulțime de comenzi ale sistemului Unix (comenzile care prelucrează fișierele text).

Un *client* este un proces care comută un proces server dintr-o stare pasivă într-o stare activă, trimițându-i cereri. Astfel, un client inițiază o activitate în momentele alese de el; în majoritatea cazurilor el este întârziat până la îndeplinirea serviciului. Un *server* este un proces infinit care, în mod uzual deserveste mai mulți clienți.

Procesele *pereche* sunt o colecție de procese identice care cooperează pentru a furniza un serviciu sau pentru a calcula un rezultat. De exemplu, două procese pereche pot gestiona fiecare o copie a unui fișier și interacționa pentru a păstra cele două copii consistente. Un alt exemplu de procese pereche este dat de procesele care interacționează pentru a rezolva o problemă de calcul paralel, fiecare rezolvând o parte a acestei probleme.

3.1. Modelul client-server

Multe din sistemele de operare distribuite folosesc modelul client/server în scopul repartizării funcțiilor de control în procesele din diverse noduri (fig. 3.1.1).

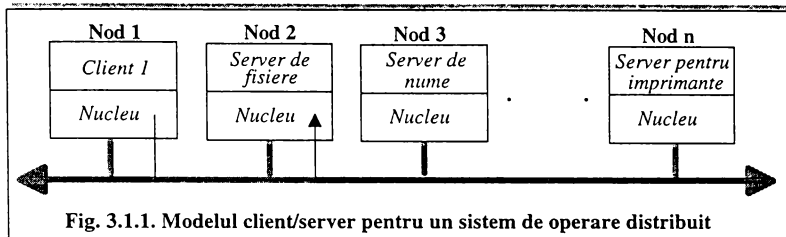


Fig. 3.1.1. Modelul client/server pentru un sistem de operare distribuit

Se poate afirma că există o tendință actuală a sistemelor de operare de a plasa cât mai multe funcții de la nivelul nucleului sistemului de operare la nivelul proceselor server, reducându-se astfel nucleul la o formă minimală.

Următoarele sisteme de operare distribuite experimentale au fost realizate pe baza modelului client/server[Gos91]: Eden, V Kernel, Amoeba; caracteristica lor principală, comună, este aceea că resursele sunt controlate centralizat, prin servere.

Se apreciază ([Gos91]) că există trei probleme majore ale modelului client/server: (1) dacă un nod pe care se află unul sau mai multe servere este indisponibil, întregul sistem devine indisponibil; o astfel de soluție nu oferă un grad prea mare de siguranță, și, pentru unele funcții de control critice într-un sistem (precum cele ale unui server de nume sau de fisiere) poate deveni inadecvată; (2) dacă numărul clienților este foarte mare pentru un server, acesta poate intra într-un regim de "sufocare"; cu cât numărul clienților crește, cu atât efectul este exacerbat; (3) pentru îmbunătățirea performanțelor modelului client/server se poate folosi, local, un mecanism de caching pentru anumite funcții ale unui server.

3.1.1. Prezentare generală a modelului client/server

Modelul client-server este cea mai cunoscută și folosită paradigmă în programarea distribuită. Conform acestei paradigme, un proces *server* (în sensul de servent) îndeplinește un anumit serviciu pentru un proces *client*. Clientul cere serviciul prin emiterea unui mesaj către server. Serverul, la rândul lui, execută ciclic următoarele acțiuni: recepționează o cerere de serviciu de la un client, îndeplinește serviciul și emite un mesaj de răspuns spre client. Acest model, în care informația circulă în ambele sensuri între procesele implicate este, evident, asimetric [Mus96j].

Procesul server este un proces *persistent*; el există întodeauna în sistem, așteptând pasiv ca un proces client să-i utilizeze serviciile. Paradigma client-server divizează deci procesele care își transferă informații în două categorii: client și server, în funcție de acțiunea pe care procesul o execută - fie așteaptă o comunicație (server), fie inițiază o comunicație (client). De fapt, motivația fundamentală a paradigmei client-server provine din cunoscuta problema de "rendez-vous". Modelul client-server rezolvă această problemă impunând ca dintre oricare două procese care vor să comunice, unul să-și înceapă execuția și să aștepte (un timp nedefinit) până la contactarea sa de către alt proces. Soluția este importantă mai ales în rețelele în care comunicația se face prin TCP/IP deoarece TCP/IP nu furnizează nici un mecanism care să creeze automat procese la primirea unui mesaj. Serverul este inițiat ca un proces autonom în rețea. Structura de bază a unui proces server este

```
server(){
    initializari;
    while (TRUE){
        asteapta_o_cerere(...);
        trateaza_cererea(...);
    }
}
```

Figura 3.1.1. Structura unui proces server (I)

reprezentată în figura 3.1.1.

Această structură de bază sugerează o interfață prin datagrame deoarece serverul poate accepta o cerere de la orice client. Este posibil, în această implementare a procesului server, ca deservirea unei cereri să se facă și în cadrul unei conexiuni și să necesite un timp destul de mare, deoarece, pe de o parte este necesar să se

stabilească un circuit virtual între client și server și, pe de altă parte, un proces client poate monopoliza serverul un timp nedefinit. De aceea nodul care implementează serverul ar trebui să memoreze cererile către server care sosesc în timp ce acesta tratează o cerere (într-un context de multiprogramare).

Modelul poate fi însă extins astfel încât să folosească multiprogramarea. Se inițiază întâi un proces special, *listener*, a cărui singură sarcină este de a accepta cereri de servicii și de a le transmite altor procese care deserveșc cererile. Structura de bază a procesului server obținut în această implementare este prezentată în figura 3.1.2.

```
listener(){
  socket client, server;
  virtual_circuit circuit_virtual;
  initializari;
  while (TRUE){
    client <-- asteapta_o_cerere(socket_impus);
    server <-- creare_socket();
    circuit_virtual <-- accept(client, server);
    creare_proces(server, ..., circuit_virtual, ...);
  }
}
server(...,circuit_virtual, ...){
  while (*serviciul este cerut){
    read(circuit_virtual,...);
    indeplinește_serviciul(...);
    write(circuit_virtual,...);
  }
  close(circuit_virtual);
}
```

Figura 3.1.2. Structura unui proces server (II)

Fiecare cerere a unui client generează crearea unui nou proces server conectat la procesul client printr-un circuit virtual. Clientul poate monopoliza acest server un timp nedefinit, fără a provoca blocarea altor clienți care pot emite în același timp cereri către procesul listener; acesta creează câte un proces server pentru fiecare cerere.

În acest al doilea model, procesele server sunt de fapt proceduri identice care se execută asupra unor structuri de date partajate (tabloul de cereri recepționate). Comutarea contextului între procesele server poate deveni un factor dominant și perturbator în deservirea cererilor. De aceea modelul poate fi perfecționat prin folosirea *firelor* ("threads"). Firele nu au alocate resursele necesare unui proces; ele pot opera însă asupra resurselor alocate procesului care îl încapsulează.

```
client()
{
  message msg;
  socket client, server;
  virtual_circuit circuit_virtual;
  ....
  send_datagram(NAME_SERVER, "nume server");
  msg = receive(NAME_SERVER);
  server = msg.adr_server;
  client = creare_socket();
  circuit_virtual = open(client, server);
  while ()
  {
    write(circuit_virtual, ...);
    read(circuit_virtual, ...);
    analizeaza_raspusul();
  }
  close(circuit_virtual);
  ...
}
```

Figura 3.1.3. Structura unui proces client

În figura 3.1.2 clientul își transmite cererea la un "socket impus", în sensul de o adresă cunoscută în rețea. Astfel, se impune proceselor client să "cunoască" o adresă fixată în rețea, pentru a-și putea adresa cererile de servicii la acea adresă. Dacă adresa este cunoscută atunci ea se poate specifica direct în programul sursă al clientului, dar, în acest caz, o dată compilat programul client, serverul nu mai poate fi mutat pe alt nod sau poate fi returnată de un server de nume.

Dacă procesul client intenționează să trimită o cerere de servicii unui server de nume pentru a proiecta numele serverului dorit într-o adresă atunci serverul de nume trebuie să memoreze un catalog de servere și adresele lor. Acest scenariu pentru luarea contactului cu un server și pentru adresarea cererilor este descris în figura 3.1.3.

În diverse aplicații software, serverele trebuie să poată avea acces la fișiere, porturi, la care în mod normal sistemul de operare interzice accesul unui client obisnuit. În aceste cazuri serverele se execută cu privilegii speciale și trebuie luate măsuri de protecție ca aceste privilegii speciale să nu fie transmise clienților care apelează aceste servere. De exemplu, un server care se execută ca un program privilegiat trebuie să conțină o funcție care să verifice dacă un fișier poate fi accesat de un client sau nu.

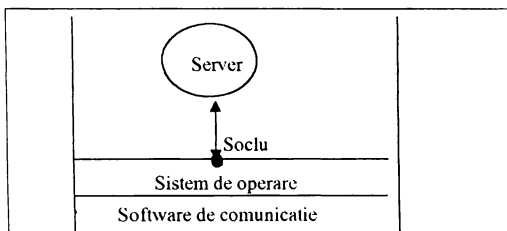
În general, serverul trebuie să asigure verificări pentru: (1). *Autentificarea clienților* - verificarea identității unui client. (2). *Autorizarea clienților* - determină dacă unui client îi este permis accesul la serviciul pe care serverul îl furnizează. (3). *Securitatea datelor* - garantarea păstrării integrității datelor (ele pot fi compromise de clienți în mod intenționat sau neintenționat) (4). *Păstrarea secretului datelor* - garantarea interzicerii accesului unor clienți la datele particulare ale unui client. (5). *Protecția datelor* - garantarea că aplicațiile nu pot abuza de resursele sistemului

3.1.2. Concurența la nivelul serverelor și clienților. Exemple.

S-au prezentat în secțiunea precedentă două modele ale procesului server. Termenul de *server iterativ* desemnează implementarea unui server care procesează la un moment dat o singură cerere, iar termenul de *server concurrent* implementarea unui server care procesează mai multe cereri la un moment dat. Se va arăta în acest subcapitol că în cazul multor servere este suficientă realizarea unei concurențe aparente; astfel, dacă un server execută puține calcule în raport cu operațiile de I/O pe care trebuie să le execute, acesta este posibil să se implementeze ca un singur proces care utilizează operații de I/O asincrone, pentru a permite comunicații simultane pe mai multe canale. Din perspectiva unui client, serverul apare comunicând concurrent cu mai mulți clienți. În acest sens, termenul de *server concurrent* se referă numai la faptul că serverul gestionează concurrent cererile nu și la implementarea serverului care utilizează procese concurrente.

Se poate defini ca *timpul de procesare* a unei cereri de către un server, timpul necesar unui server pentru a trata o cerere izolată, iar ca *timpul de raspuns* întârzierea produsă de la lansarea unei cereri de către client pînă la primirea răspunsului.

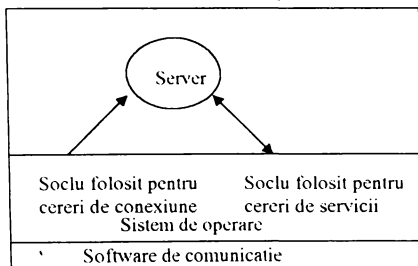
Serverele iterative (fig 3.1.2.1.) procesează câte o cerere, pe rînd. Dacă o nouă cerere apare în timp ce serverul este ocupat cu tratarea alteia, sistemul memorează nouă cerere; la terminarea tratării cererii curente, se extrage nouă cerere de tratat s.a.m.d.. Dacă lungimea medie a cozii de cereri este N , atunci timpul de răspuns pentru o cerere poate fi apreciat în medie ca fiind mai mare de $N/2 + 1$ ori decît timpul de procesare a unei cereri. Deoarece timpul de răspuns crește proporțional cu N , multe implementări restricționează valoarea lui N la o valoare mai mică (de exemplu 5) și sugerează ca programatorii să folosească servere concurrente în cazul în care o coadă de cereri, de o astfel de dimensiune mică nu este suficientă. Un server iterativ trebuie analizat și din alt punct de vedere: acela al suprasaturării sale. Astfel, un server iterativ desemnat a deservi C clienți, fiecare trimițînd K cereri pe secundă va trebui să aibă un timp de procesare pentru o cerere mai mică decît $1/(C*K)$ secunde. Dacă serverul nu poate trata cererile cu această rată cerută, atunci există pericolul depășirii cozii sale de cereri și deci a pierderii unor cereri. Pe de alta parte, un server iterativ este ușor de proiectat, depanat și verificat; varianta iterativă poate fi aleasă cînd execuția iterativă furnizează un răspuns suficient de rapid pentru încărcarea estimată. În mod obișnuit, serverele iterative se asociază unui protocol neorientat pe conexiune, deși pot fi scrise variante de servere iterative și bazate pe protocole orientate pe conexiune.



(a) Server iterativ, protocol neorientat pe conexiune

Algoritm:

- (1) Se creează un soclu și se leagă (bind) acest soclu la adresa stabilită pentru serviciul oferit
- (2) În mod iterativ, se citește cererea următoare (recvfrom), se tratează cererea și se trimite răspunsul clientului



(b) Server iterativ, protocol orientat pe conexiune

Algoritm:

- (1) Se creează un soclu și se leagă (bind) acest soclu la adresa stabilită pentru serviciul oferit
- (2) Se așteaptă o nouă cerere de conexiune pe soclu creat și se obține un alt soclu
- (3) În mod iterativ, se citește o cerere, se tratează cererea și se trimite răspunsul clientului
- (4) Când se termină de tratat cererile clientului, se închide conexiunea (socul obținut la doi) și se trece la pasul 2

Fig. 3.1.2.1. Servere iterative

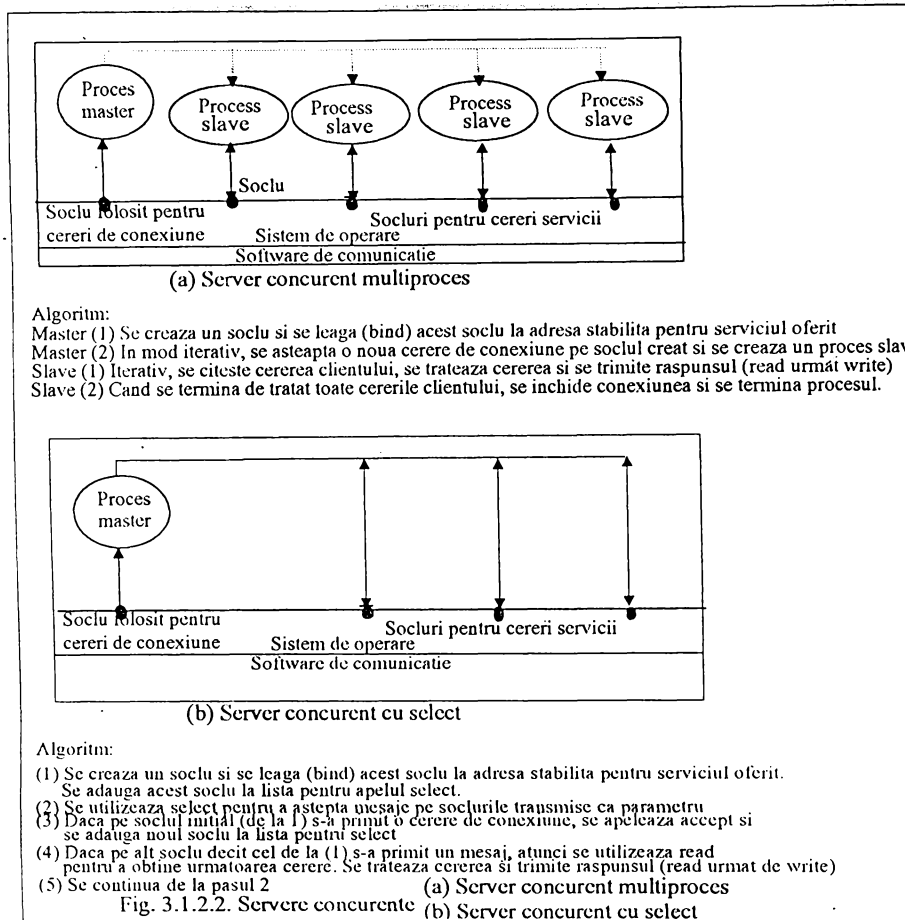
- (a) Server iterativ, protocol neorientat pe conexiune
- (b) Server iterativ, protocol orientat pe conexiune

Algoritm:

Există două soluții, conceptual diferite pentru a realiza tratarea concurență a cererilor într-un server. O primă soluție, cea mai cunoscută este de a utiliza mai multe procese sau fire în cadrul serverului. În cadrul unui server neorientat pe conexiuni va exista deci un singur *proces server master* care începe inițial execuția; procesul *master* deschide un soclu pe

tip de server, *serverele concurrente*, este acela al timpului de răspuns; serverele concurrente pot deservi mai mulți clienți oferind un timp de răspuns rapid. Concurența îmbunătățește timpul de răspuns dacă: (1) generarea unui răspuns implică operații de I/O; prin suprapunerea operațiilor de I/O corespunzătoare unei cereri, cu execuția de către procesor a altei cereri, se poate câștiga timp; (2) timpul de procesare al cererilor variază mult, în raport de conținutul cererilor; tehnica *timeslicing* permite tratarea cererilor scurte în același timp cu cele lungi fără a fi necesară așteptarea terminării cererilor lungi. (3) serverul se execută pe un calculator cu mai multe procesoare.

Există două soluții, conceptual diferite pentru a realiza tratarea concurență a cererilor într-un server. O primă soluție, cea mai cunoscută este de a utiliza mai multe procese sau fire în cadrul serverului. În cadrul unui server neorientat pe conexiuni va exista deci un singur *proces server master* care începe inițial execuția; procesul *master* deschide un soclu pe

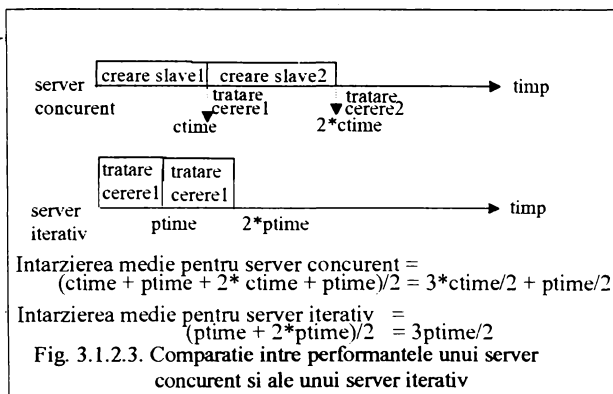


un port cunoscut în sistem, așteaptă cerere și crează un proces *slave* pentru a trata fiecare cerere. Procesul master nu comunică direct niciodată cu un client; aceasta va fi responsabilitatea proceselor slave. Fiecare proces slave comunică cu un client; el își încetează execuția când termină de tratat cererea. Deoarece crearea de noi procese este o operație costisitoare în timp, serverele neorientate pe conexiune folosesc mai rar implementări concurente. În schimb, serverele orientate pe conexiuni, în cazul cărora o conexiune dintre server și un client cuprinde mai multe cereri, implementează concurența la nivelul conexiunilor și nu la nivelul fiecărei cereri. Ca și în cazul unui server neorientat pe conexiune, există un proces master și procese slave. Procesul master nu comunică direct cu un client niciodată; la orice cerere de deschidere a unei noi conexiuni, procesul master crează un proces slave care să gestioneze conexiunea. În timp ce procesele slave vor interacționa cu clienții, procesul master va aștepta alte conexiuni (figura 3.1.2.2.). Pentru multe servicii, este suficient un singur program, care să conțină codul corespunzător atât procesului master cât și proceselor slave. În cazurile în care scrierea proceselor slave reclamă un program independent, pentru o mai bună proiectare, deparare și umarare, în Unix trebuie utilizat un apel *execve()* după fiecare cerere de proces (*fork()*).

A doua soluție pentru a realiza tratarea aparent concurrentă a unor cereri dar într-un singur proces, se poate realiza dacă încărcarea totală a serverului (numărul de cereri) nu depășește capacitatea sa de prelucrare. În acest caz, serverul se constituie într-un singur proces care utilizează apelul *select* pentru a gestiona asincron canalele de comunicații. Fie, de exemplu, un singur proces server care are deschise mai multe conexiuni TCP pentru mai mulți clienți. La sosirea unui mesaj pe o conexiune, procesul se deblochează, tratează cererea și trimite un răspuns, după care se autoblochează, așteptând din nou sosirea de mesaje pe una dintre conexiuni. Atât timp cât procesul este suficient de rapid pentru a satisface cererile prezentate serverului, această versiune de server este cel puțin atât de rapidă cât și versiunea serverului multiproces. Ideea de bază este ca într-un singur proces să se înglobeze funcțiile atât ale procesului master cât și a proceselor slave prin folosirea

apelului *select*. Serverul creează câte un soclu pentru fiecare conexiune pe care trebuie să o gestioneze și apoi execută *select()* pentru a aștepta sosirea de mesaje pe unul din socluri. Apelul *select()* primește ca parametrii setul de descriptori ai soclurilor pe care se așteaptă mesaje și transmite o mască de biți care specifică care descriptori din set sunt gata. Pentru a distinge între operații master și slave, serverul utilizează descriptorii; dacă devine gata un descriptor corespunzător unui soclu master atunci serverul execută aceeași operație pe care ar realiza-o un proces master: apelează *accept()* pentru acel soclu pentru a obține o nouă conexiune. Dacă devine gata un descriptor corespunzător unui soclu slave, atunci serverul efectuează aceeași operație pe care ar realiza-o un proces slave: apelează *read* pentru a obține mesajul, tratează cererea și trimite răspunsul (3.1.2.2.b.).

Alegerea între un server iterativ și unul concurrent este dificilă deoarece cererile utilizatorilor, vitezele de prelucrare, posibilitățile de comunicare se pot modifica rapid. Legat de concurență, mai trebuie abordat un aspect important: acela al *nivelului de concurență*, definit ca numărul total de procese pe care le-a generat la un moment dat. Nivelul de concurență al unui server variază pe parcursul existenței lui, dar ceea ce interesează este nivelul de concurență al unui server maxim atins. În mod obișnuit, un server concurrent, bazat pe un protocol orientat pe conexiune creează câte un proces pentru fiecare conexiune cerută de un client. Evident, un server nu poate crea însă un număr arbitrar de conexiuni, deoarece, pe de o parte, fiecare implementare TCP limitează numărul de conexiuni active posibile și, pe de altă parte, fiecare sistem de operare limitează numărul de procese care pot fi create în sistem. Când serverul atinge una dintre aceste limite, sistemul va refuza serverului crearea de noi procese. Prin urmare, serverele care creează câte un proces nou la fiecare cerere, fie de conexiune, fie de serviciu, (care urmează astfel o strategie *demand-driven* când nivelul de concurență crește la *cerere*), dacă nu limitează explicit nivelul de concurență, prezintă riscuri în funcționarea într-un sistem de operare cu o încărcare mai mare. Strategia *demand-driven* în ceea ce privește concurența cererilor, desi prezintă avantajul economisirii resurselor (bufferele, blocurile de control pentru procese nu se alocă decât atunci când sunt necesare) poate să nu producă însă întotdeauna rezultate optime. Astfel, dacă timpul de procesare al unei cereri, *ptime*, este mai mic decât timpul necesar pentru crearea unui proces, *ctime*, versiunea iterativă a unui server poate fi mai performantă decât versiunea concurrentă în ceea ce privește întârzierea medie de deservire a cererilor (figura 3.1.2.3.). Întârzierea medie care se obține pentru procesarea a două cereri de către serverul iterativ este $3 * ptime / 2$, mai mică decât $3 * ctime / 2 + ptime / 2$, întârzierea pentru procesarea a două cereri de către serverul concurrent. Implementarea concurrentă a serverului va putea deservi $1 / ptime$ cereri pe unitatea de timp, în timp ce implementarea iterativă va putea deservi $1 / ctime$ cereri pe unitatea de timp; evident că dacă rata de sosire a cererilor *r* este $1 / ctime < r < 1 / ptime$ numai o implementare iterativă a serverului funcționează corect; cea concurrentă pierde mult timp pentru crearea proceselor și rejectează cereri când cozile folosite de protocoalele software ajung la capacitatea maximă. Concluzia este că, în practică, nu trebuie utilizate variante concurente decât atunci când timpul de procesare a unei cereri depășește timpul de creare a unui proces.



Pe de altă parte, trebuie cautate soluții alternative la strategia *demand-driven* în ceea ce privește concurența. O soluție care limitează nivelul maxim de concurență al serverelor, controlează întârzierile în deservirea cererilor și oferă performanțe chiar și când timpul de creare al proceselor este relativ mare, constă în prelocarea proceselor concurente, pentru a evita pierderile de timp la crearea lor. Astfel, la inițializare, serverul master creează toate cele P procese slave care își încep execuția așteptând sosirea unei cereri. În Unix, un server master neorientat

pe conexiune va crea și soclul asociat portului rezervat serviciului respectiv; acest soclu va fi moștenit de procesele slave care execută *recvfrom()* pentru a obține o datagramă și adresa clientului. După crearea soclului și a proceselor slave, procesul master se poate termina sau se poate transforma el însuși în ultimul proces slave, evitând astfel crearea încă unui proces. Analog se procedează și cu server master orientat pe conexiune: se creează soclul pe portul rezervat și se creează apoi cele P procese slave. Acestea execută *accept()* pe soclul moștenit; o cerere de conexiune va debloca un singur slave care va prelua noul soclu, și va trata cererile pe aceasta conexiune. Când clientul termină cererile, procesul slave închide conexiunea și reia apelul *accept()* pentru a prelua o nouă cerere de conexiune.

Deși prealocarea proceselor poate îmbunătăți eficiența, ea nu rezolvă însă toate problemele. Crearea proceselor noi crește overhead-ul datorat software-ului de comunicație. O a doua soluție, care poate crește

eficiența serverelor este opusă celei precedente, ca constând în amânarea alocării proceselor. Ideea de bază este următoarea: în loc de a decide a priori o implementare iterativă sau concurrentă, se permite serverului să masoare timpul de procesare a unei cereri și să aleagă dinamic între o tratare iterativă sau concurrentă.

Pentru a implementa dinamic alocarea întârziată a proceselor slave, algoritmul este următorul: serverul master recepționează o cerere, setează un timer (*alarm()* în Unix) și începe tratarea cererii, în mod iterativ. Dacă serverul termină tratarea cererii înainte de expirarea timer-ului, atunci anulează timer-ul; dacă, însă, timer-ul expira înainte ca serverul să termine tratarea cererii, serverul crează un proces slave care tratează în continuare cererea (apelul *fork()* în Unix permite, pe lângă moștenirea soclurilor și continuarea codului parintelui exact din punctul în care a fost întrerupt când timer-ul a expirat).

Cele două soluții prezentate mai sus pot fi combinate. Un server își poate începe execuția fără a avea procese prealocate și poate utiliza alocarea întârziată a proceselor. Se așteaptă recepția unei cereri, și se crează un proces slave numai dacă tratarea necesită un timp mai mare decât cel setat la un timer. După ce procesul slave a fost creat și s-a terminat prelucrarea unei cereri, procesul slave nu se termină imediat ci așteaptă recepționarea unei noi cereri. Problema soluției combinate nu constă în decizia când se crează procese slave ci când un proces slave trebuie sau nu terminat. O variantă este ca numărul de procese slave de la un moment dat să nu depășească o limită fixată; acest lucru poate fi realizat dacă procesele slave partajază o variabilă (într-o zona de memorie partajată) care indică nivelul de concurență curent și care se testează după tratarea unei cereri pentru a se determina dacă procesul se termină sau persistă. O altă variantă pentru a controla concurența se bazează pe măsurarea perioadei de inactivitate a unui server. Algoritmul este următorul: procesul slave startează un timer înainte de a trece la așteptarea unei noi cereri; dacă timerul expira înainte de recepționarea unei cereri, procesul slave se termină.

Ca o concluzie a celor prezentate pînă acum se pot face următoarele observații:

1. *Cînd se folosește o varianta iterativă și cînd o variantă concurrentă a unui server?* Serverele iterative sunt mai ușor de proiectat, implementat și depanat, în schimb serverele concurente furnizează răspunsuri mai rapide la cereri. Se utilizează o variantă iterativă dacă timpul de procesare a unei cereri este scurt și soluția iterativă oferă timp de răspuns suficient de rapid aplicației respective.

2. *Cînd se folosește un protocol orientat pe conexiune și cînd un protocol neorientat pe conexiune?*

Se utilizează un protocol neorientat pe conexiune dacă protocolul de la nivelul aplicației tratează pierderile de mesaje, sosirea pachetelor în secvența (în general, aplicațiile nu realizează, însă, acest lucru) sau dacă procesele client sunt localizate în aceeași rețea locală cu serverul, rețea care nu ridică probleme de fiabilitate.

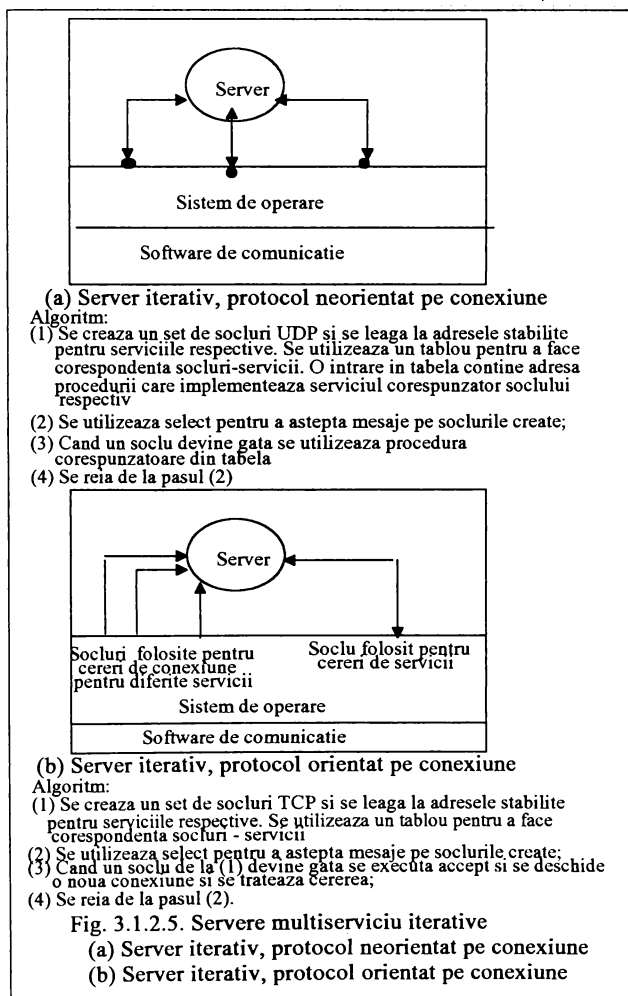
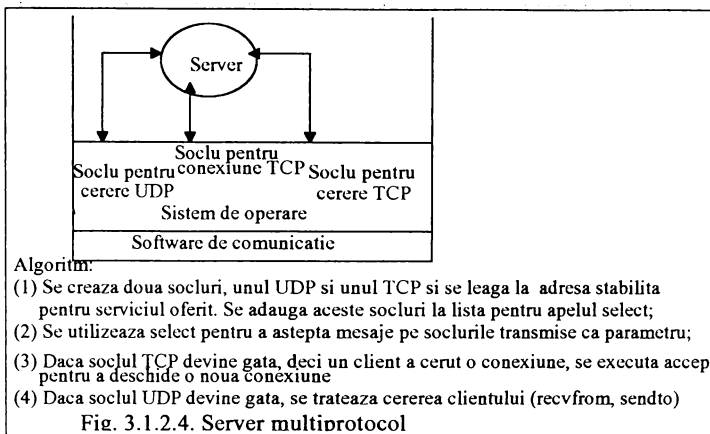
Dacă clienții și serverul se afla într-o rețea metropolitană se va utiliza întotdeauna un protocol orientat pe conexiune.

3. *Cînd se folosește o concurență aparentă în server și cînd o concurență reală, multiproces?*

Un singur proces server poate gestiona mai multe conexiuni utilizînd operații de I/O asincrone; un server multiproces folosește sistemul de operare pentru a obține automat tratarea concurrentă a cererilor. Un server implementat ca un singur proces (concurență aparentă) se utilizează dacă serverul trebuie să partajeze sau transmită date între conexiuni. Un server multiproces se utilizează dacă fiecare proces slave poate lucra izolat, obținîndu-se astfel un nivel mare de concurență (în special pe sisteme multiprocesor).

Folosirea operațiilor de I/O asincrone pentru a gestiona simultan mai multe conexiuni este o soluție și pentru realizarea de servere multiprotocol sau/și multiserviciu. Serverele multiprotocol sunt adecvate pentru situațiile în care un serviciu trebuie realizat atît pentru protocol TCP cît și pentru UDP; dacă același serviciu s-ar implementa în două servere separate (unul UDP, altul TCP) se irosesc resurse ale sistemului în mod inutil (intrări în tabela de procese), mai ales dacă se are în vedere că standardul TCP/IP definește zeci de servicii. Un server multiprotocol constă într-un singur proces care utilizează operații de I/O pentru a gestiona comunicatii, fie pe UDP, fie pe TCP (fig. 3.1.2.4.) și care pentru tratarea cererilor prevede o procedură unică, indiferent dacă cererea a sosit via UDP sau via TCP. Evident că serverul din fig. 3.1.2.4. poate fi extins pentru a permite concurența aparentă între cereri care se recepționează pe mai multe conexiuni TCP sau pe mai multe socluri UDP.

Serverele multiserviciu sunt adecvate pentru situațiile în care înglobarea mai multor servicii într-un program reduce codul total obținut față de situația în care aceste servicii ar fi înglobate în servere separate. Serverele multiserviciu pot fi și ele iterative (fig. 3.1.2.5.), concurente (3.1.2.6.) sau chiar multiprotocol sau unele servicii pot fi tratate iterativ și alte concurent, în funcție de natura serviciului. Un server multiserviciu multiprotocol este denumit *suverserver*.



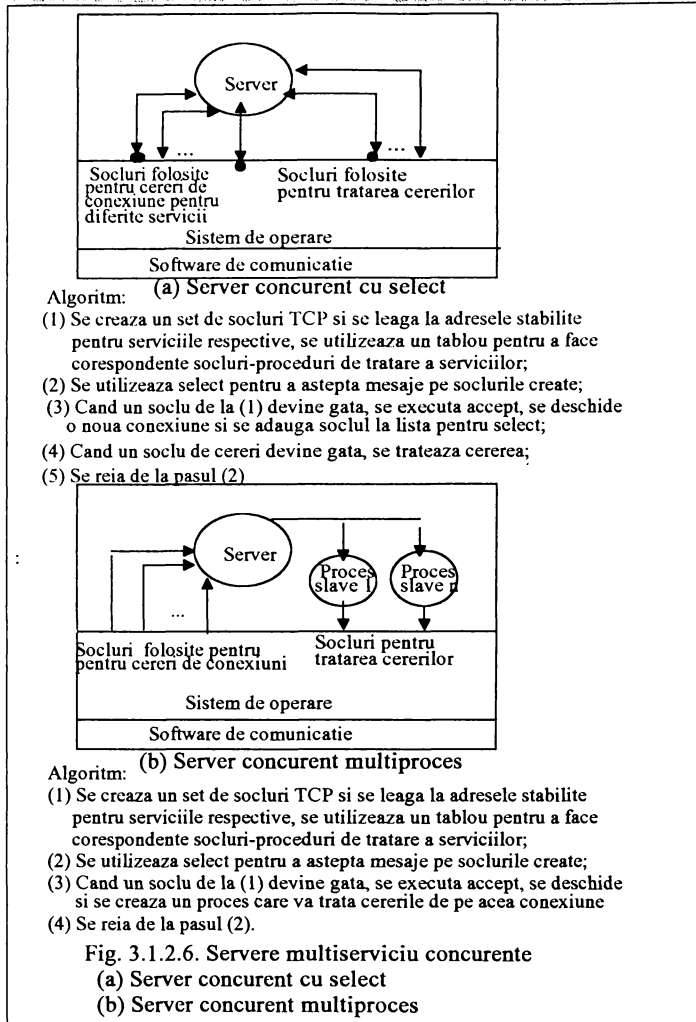
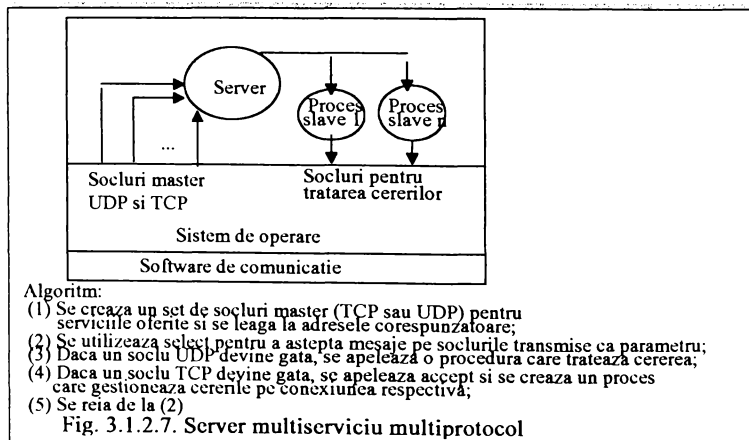
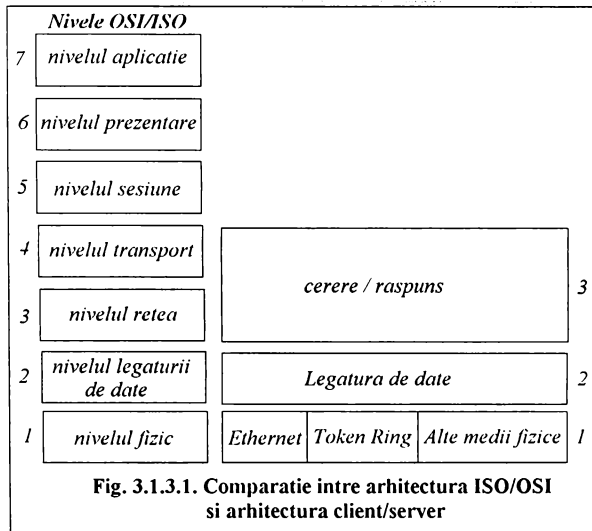


Fig. 3.1.2.6. Servere multiserviciu concurente
 (a) Server concurrent cu select
 (b) Server concurrent multiproces



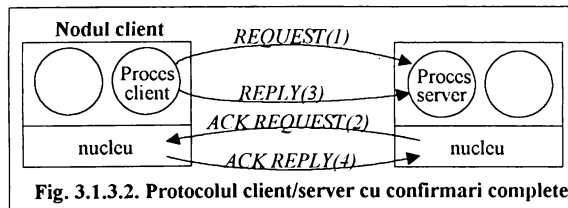
3.1.3. Protocele pentru modelul client-server

Modelul client-server poate fi implementat utilizând ambele clase de protocele prezentate în subcapitolul 2.4: protocele cu caracter universal și protocele cu caracter restrans. O caracteristică a primei clase de protocele este aceea că prezintă rețeaua de comunicație ca pe un dispozitiv de I/O; în consecință, comunicația la nivelul proceselor în rețele se face utilizând *modelul fișier*. Aceasta înseamnă deschiderea unei conexiuni (open), utilizarea de primitive de citire sau scriere (read, write) și, în final, închiderea conexiunii (close). Modelul fișier este propus atât în protocelele TCP/IP, cât și în ISO/OSI. S-au aratat în 2.3 și 2.4 dezavantajele utilizării unui protocol cu caracter general. În sensul obținerii unei performanțe superioare, pentru modelul client/server, se poate adopta un protocol cu caracter restrans, a cărui interfață cu procesele utilizator constă numai într-un mesaj REQUEST transmis de la procesul client la server, urmat de obținerea răspunsului de la server, printr-un mesaj REPLY; acest model de comunicație se mai numește și *modelul operației la distanță (remote operation model)*. Structura pe nivele a acestui model este foarte simplă (fig. 3.1.3.1).



Pentru detectarea și recuperarea erorilor (folosind, evident, mecanismul confirmărilor) pot fi folosite unele simplificări prezentate în secțiunea 2.3. Astfel, varianta completă de confirmare a ambelor mesaje REQUEST, REPLY (fig. 3.1.3.2) poate fi simplificată în modul următor: dacă REPLY este primit de client într-un interval de timp fixat, el este considerat și confirmare; dacă acest interval este depășit, se retransmite cererea; (3) La primirea unui REQUEST, serverul startează un timer; dacă serverul trimite REPLY fără a depăși timpul fixat de timer, se consideră acest REPLY ca și confirmare.

Dacă timerul generează timeout înainte de trimiterea REPLY-ului, atunci se va trimite o confirmare pentru REQUEST și apoi REPLY. Se vor utiliza astfel confirmări numai pentru mesajele de REPLY (fig. 3.1.3.3); acest lucru este util pentru ca serverele să execute numai o singură dată o cerere repetată. Dacă nu se confirmă de către client un REPLY și se retransmite cererea, serverul nu va mai executa cererea ci, constatând că a primit o cerere deja tratată, va trimite imediat REPLY.



Un alt aspect legat de micșorarea numărului de confirmări este confirmarea la nivel de mesaje, și nu de pachete (considerând un mesaj segmentat în pachete) și transmiterea mesajelor în flux de pachete (mecanismul *packet blasts*). Un astfel de mecanism este utilizat în protocolul VMTP propus de V kernel.

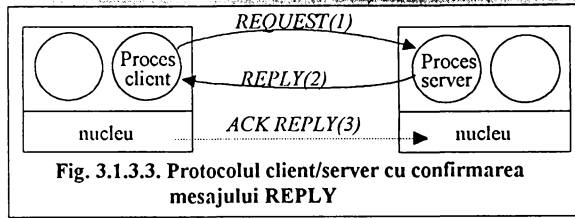


Fig. 3.1.3.3. Protocolul client/server cu confirmarea mesajului REPLY

În [Tan93][Tan96] se apreciază că, pentru controlul erorilor la minimum, sunt necesare șapte tipuri de mesaje (fig. 3.3.4). La nivelul proceselor sunt vizibile, evident, numai REQUEST și REPLY.

Tip mesaj	De la	Catre	Semnificatie
REQUEST (REQ)	Client	Server	Clientul cere executia unei cereri serverului
REPLY (REP)	Server	Client	Raspuns pentru client de la server
Acknowledg. (ACK)	Oricare	Celalalt	Confirmarea unui mesaj anterior
Are you alive (AYA)	Client	Server	Test pt. un server care nu raspunde la comenzi
Am am alive (IAA)	Server	Client	Raspuns la un test de prezenta
Try again (TA)	Server	Client	Serverul nu poate prelua o cerere (fie nu are loc sa o memoreze, fie nu o poate prelua nefiind intr-o stare in care poate prelua cerer
Address unknown (AU)	Server	Client	Nu exista un proces server asociat acestei adre

Fig. 3.1.3.4. Tipuri de mesaje

În fig. 3.1.3.5 se arată și cateva protocoale client/server care utilizează aceste mesaje.

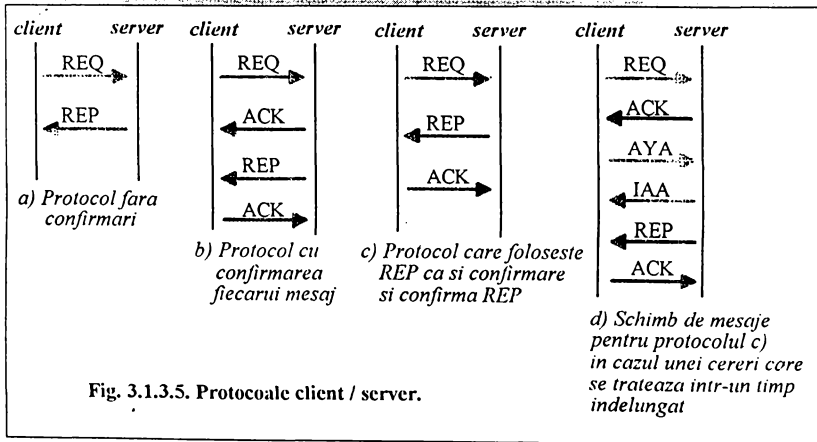


Fig. 3.1.3.5. Protocoale client / server.

3.1.4. Sisteme client-server

Se poate considera că, în aplicațiile distribuite, unul sau mai mulți clienți, și unul sau mai multe servere, împreună cu suportul pentru comunicarea interprocese formează un sistem compus (Client Server System sau, pe scurt, CSS) care permite efectuarea de calcule și prelucrări distribuite, analiza și vizualizarea datelor. Într-un astfel de sistem, un client este un proces care interacționează cu utilizatorul și are următoarele caracteristici [Mus96j]:

1. *Prezintă o interfață cu utilizatorul* (User Interface = UI), singurul mijloc prin care utilizatorul își poate prezenta cererile de consultare și regăsire a datelor, dar și de vizualizare a rezultatelor cererilor sale. În mod tipic, clienții interacționează cu utilizatorii printr-o interfață grafică (Graphical User Interface = GUI). Deoarece într-un CSS pot coexista mai mulți clienți, pot exista astfel și mai multe variante de GUI. Unele CSS-uri pot fi prevăzute și cu o interfață pentru controlul resurselor sistemului.

Tipic, interfața grafică a sistemului de operare definește interfața GUI dintre client și utilizator. Astfel, într-un mediu UNIX, sistemul XWindows a devenit, "de facto", interfața grafică standard.

2. *Formează, pentru una sau mai multe comenzi ale utilizatorului una sau mai multe cereri către server.* Serverul și clientul pot utiliza un limbaj comun pentru prezentarea cererilor. Nu toate comenzile utilizatorului se regăsesc, într-o corespondență de unu la unu cu cererile adresate serverului.

Un client poate utiliza un mecanism de "caching" și diverse tehnici de optimizare pentru a reduce fluxul de cereri către server; de asemenea poate realiza verificările legate de autentificare și autorizare (controlul accesului la date).

3. *Comunică cu serverul printr-un mecanism de comunicație între procese.* Un client ideal ascunde complet utilizatorului mecanismul de comunicație interprocese.

4. *Execută o analiza a datelor trimise ca răspuns de către server în urma unei cereri și apoi le prezintă utilizatorului.* Complexitatea acestei analize variază de la un CSS la altul.

Într-un CSS, un server este un proces sau un set de procese, toate aflate pe același nod care furnizează un anumit serviciu unuia sau mai multor clienți. Serverul are următoarele caracteristici:

a) Un server furnizează un serviciu unui client. Natura și complexitatea acestui serviciu este definit de scopul fiecărui CSS. Astfel, un server poate efectua prelucrări minore ale datelor (un server de fișier sau un server de imprimante), sau prelucrări complexe ale datelor (un server pentru baze de date sau pentru procesarea imaginilor).

b) Un server răspunde la cererile și comenzile pe care le primește de la clienții sai. Astfel, un server nu va iniția niciodată un dialog cu unul dintre clienții sai. El acționează astfel ca un depozit de date (server de fișiere), de cunoștințe (server de baze de date) sau furnizor de servicii (server de imprimante).

Cu toate acestea, un server poate adresa la rândul său cereri de servicii altui server, deci devine, la rândul său, un client pentru un alt server. De exemplu, un server de fișier poate cere obținerea datei și timpului curent (pe care le asociază unui fișier) unui server de dată și timp. Situația serverului care devine astfel client

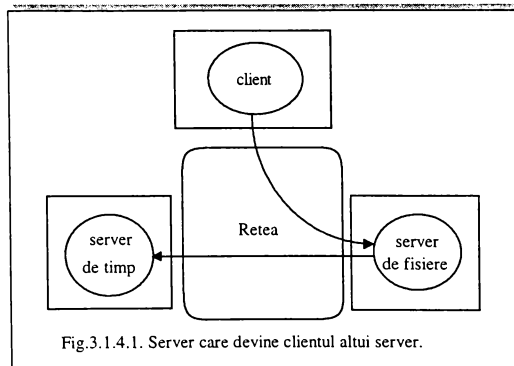


Fig.3.1.4.1. Server care devine clientul altui server.

este ilustrată în figura 3.1.4.1

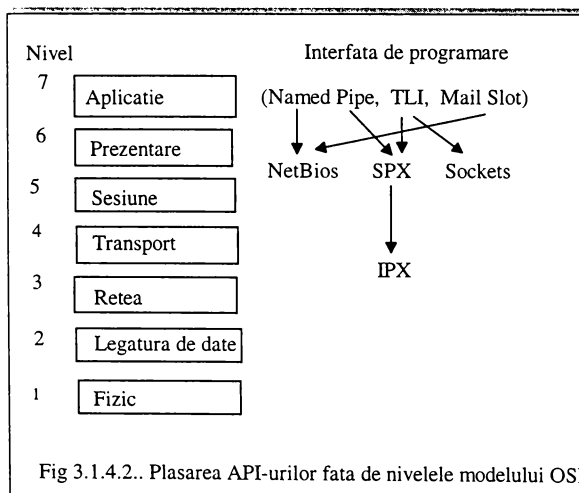
c) Un server ideal ascunde întreaga componentă a unui sistem CSS clienților și utilizatorilor sai. Un astfel de server poate comunica cu clienții săi indiferent de tehnologia de comunicație (hardware și software). Astfel, un client DOS ar trebui să poată comunica cu un server UNIX, indiferent deci de sistemul de operare și de tehnologia de conectare a clienților la server.

d) Un server poate memora informații relative la interacțiunile sale cu clienții (*stateful server*) sau poate să nu păstreze nici o astfel de informație de stare (*stateless server*). Servere *stateful* sunt motivate de următoarele observații:

- prin păstrarea unor informații în server se reduce dimensiunea mesajelor pe care clientul și serverul le interschimbă și serverul răspunde mai rapid cererilor;
- informațiile de stare memorate permit serverelor să-și "amintească" cererea precedentă unei cereri curente și să construiască un răspuns incremental.

În schimb, serverele *stateless* oferă o fiabilitate mai mare decât serverele *stateful*, pentru că, în cazul ultimelor, informațiile de stare devin incorecte când se pierd sau se duc în dublu sau când nodul clientului este oprit și repornit.

Ca o concluzie a celor de mai sus se poate spune că arhitectura client-server divide o aplicație distribuită în procese separate, localizate pe noduri diferite, conectate într-o rețea formând astfel ceea ce se cheamă "loosely coupled system". Un proiectant al unei aplicații divide aplicația sa în subtaskuri care urmează să fie executate fie de procese client, fie de procese server în condițiile unor restricții impuse de sistemul CSS însuși și cu o funcționalitate dependentă de sistemul de operare. Cu cât API-urile puse la dispoziția sistemului de operare sunt mai avansate, cu atât codul aplicației va fi mai mic. Pentru proiectantul care dezvoltă o aplicație într-o rețea de calculatoare, folosind modelul client-server, prezintă deci o mare importanță interfața de programare (API-urile) pusă la dispoziție de sistemul de operare. În figura 3.1.4.2. se prezintă plasarea în nivelele modelului OSI a unor astfel de interfețe de programare cunoscute .

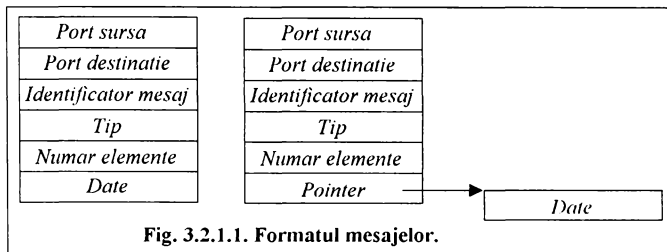


3.2. Comunicarea prin mesaje

Comunicarea prin mesaje este asociată cu modelul client/server și constituie o extensie a comunicării locale interprocese prin mesaje; este o formă de comunicare în care procesele utilizator generează explicit mesajele transmise prin subsistemul de comunicație și utilizează explicit mecanismele de livrare și recepție ale mesajelor.

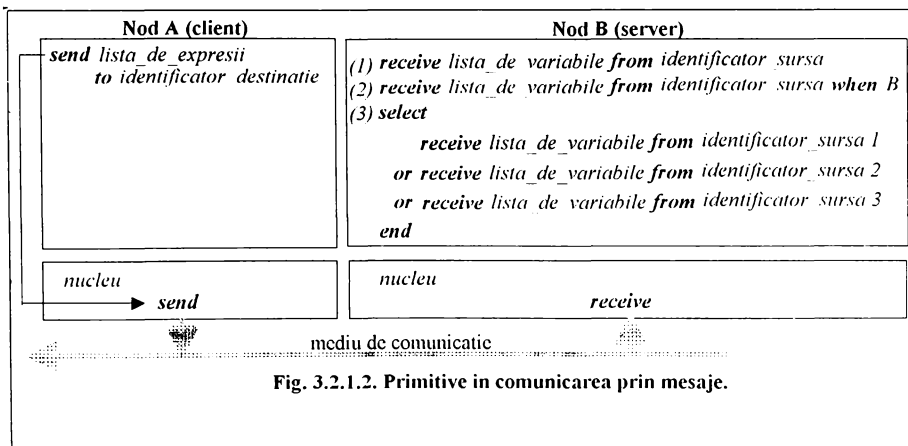
3.2.1. Probleme specifice ale comunicării prin mesaje

Un *mesaj* poate fi definit ca o structură de date compusă dintr-un header de lungime fixă și o zonă de date de lungime variabilă (sau fixă, uneori), care poate fi gestionată de procese, inclusiv livrată către o destinație sau recepționată de la o sursă (fig. 3.2.1.1).



Se enumeră următoarele probleme specifice ale comunicării prin mesaje;

- (1) mesajele pot fi complet *nestructurate* sau *structurate*; mesajele structurate sunt de preferat, pentru că: (a) ușurează tratarea unor părți din mesaj (port sursă, port destinație) de către nucleul sistemului de operare și (b) în sistemele eterogene, numai datelor cărora li se atașează tipul pot fi transmise într-o manieră transparentă;
- (2) există două aspecte importante legate de primitivile care implementează comunicarea prin mesaje: stabilirea setului de primitivă și semnificația acestora. În fig. 3.2.1.2 se prezintă primitivile de bază în comunicarea prin mesaje. Se observă că: (a) clientul specifică prin primitivă adresa destinației (serverul) și mesajul și (b) serverul specifică de la cine primește mesajul și bufferul în care se memorează mesajul recepționat; (c) nu este necesară stabilirea unei conexiuni.
- (3) sunt permise în unele implementări și primitivă selective (fig. 3.2.1.2).
- (4) semnificația primitivelor separă primitivile în primitivă *cu și fără blocare*, *sigure și nesigure*, *cu și fără bufferare*.

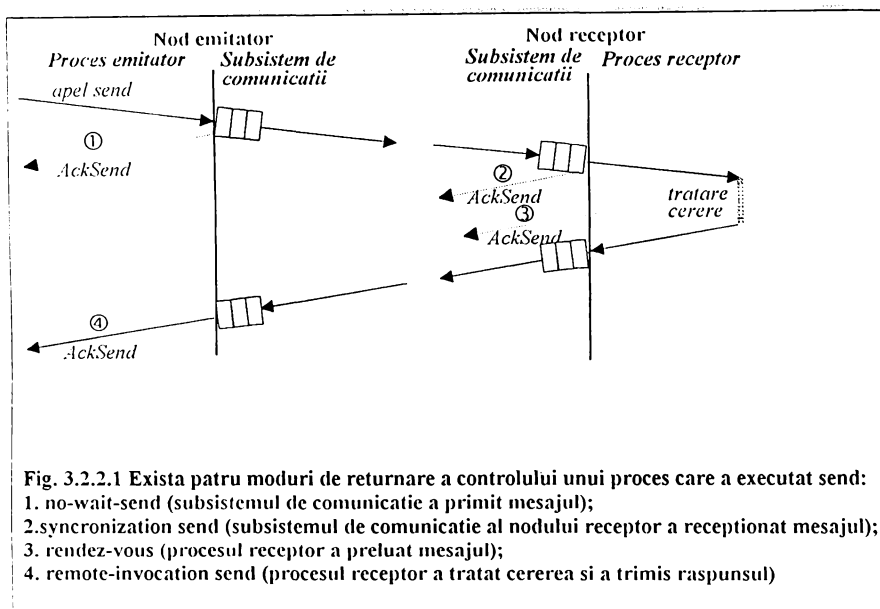


3.2.2. Primitive cu și fără blocare

O primitiva de comunicare interprocese este *fără blocare (asincrona)* dacă execuția sa nu determină întârzierea apelantului; în caz contrar, primitiva este *cu blocare*. Scopul introducerii primitivelor fără blocare este exploatarea paralelismului obținut datorită acestor primitive (procesul își continuă execuția în paralel cu transmiterea mesajelor).

Semantica *cu blocare, fără blocare* pentru operația *send* se realizează observând că programul apelant poate obține controlul după execuția operației *send* în următoarele 3 moduri (fig. 3.2.2.1):

- (1) semantica *no-wait-send* sau *asincron-send* - procesul emitor așteaptă doar preluarea mesajului de către subsistemul de comunicație;
- (2) semantica *synchronization send* - procesul emitor așteaptă pînă cînd nodul destinație primește mesajul;
- (3) semantica *rendez-vous* - procesul emitor așteaptă pînă cînd procesul receptor primește mesajul;
- (4) semantica *remote-invocation send* sau *remote procedure call* - procesul așteaptă pînă cînd primește un răspuns de la procesul receptor, în urma tratării cererii.



Modul (2) este mai puțin folosit; pentru o rețea sigură, el este echivalent cu (1). Modul (3) este folosit în sistemele centralizate (dacă procesul receptor a primit mesajul, atunci nu există motive pentru care să nu îl poată prelucra - decât în cazul caderii nodului, dar atunci va dispărea și emitorul -), dar în sistemele distribuite nu se poate garanta că, după primirea mesajului, receptorul va trata cererea și va returna răspunsul. În practică, se folosesc cele două extreme (1) și (4) și uneori (3). Modulurile (2)-(4) se consideră cu blocare, iar (1) fără blocare. Unii autori ([Tan96]) clasifică chiar și modul (1) în categoria primitivelor cu blocare, considerând fără blocare o primitiva *send* care primește controlul imediat; în acest caz, nu pot fi utilizate două primitive *send* consecutive, fără ca apelantul să se asigure că subsistemul de comunicație are bufferul liber pentru a prelua cererea. De aceea, pentru acest mod (complet fără așteptare la execuția unui *send*) subsistemul de comunicație anunță apelantul printr-o întrerupere că bufferul său este disponibil pentru un nou *send*.

În ceea ce privește operațiile de *recepție fără blocare* se obișnuiește să se furnizeze aplicațiilor următoarele variante:

- (1) o primitiva *receive* care livrează subsistemului de comunicație adresa la care se depune mesajul recepționat; pentru a se determina dacă s-a recepționat sau nu un mesaj se folosește o altă primitiva sau se generează o întrerupere la primirea mesajului;
- (2) o primitiva *receive* care preia un mesaj, dacă există disponibil, și întoarce un cod de eroare dacă nu există un mesaj (eventual, după un mic timp de așteptare); aceasta se completează fie cu o primitiva *wait*, care permite așteptarea sosirii unui mesaj, fie cu o primitiva *test*, care permite testarea stării operației de recepție;
- (3) o primitiva care livrează subsistemului de comunicație adresa unei proceduri ce se va executa în momentul sosirii unui mesaj.

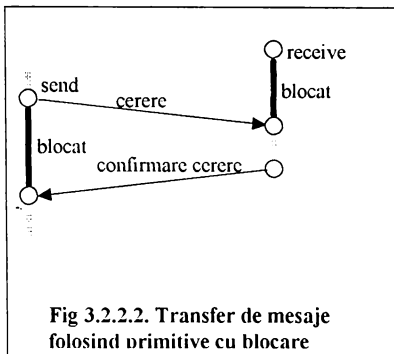


Fig. 3.2.2.2 arată un transfer de mesaje folosind primitive cu blocare.

Ambele clase de primitive, cu și fără blocare, au avantaje și dezavantaje. Comparativ cu primitivele cu blocare, cele fără blocare oferă posibilitatea paralelismului în execuție, o maxima flexibilitate și pot fi folosite și în aplicațiile în timp real, dar au și dezavantajele că conduc la programe nestructurate, greu de urmărit, dependente de timp și necesita bufferare atât la nivelul receptorului, cât și la nivelul emitorului. Unele sisteme oferă pentru emisie varianta fără blocare, iar pentru recepție varianta cu blocare. În sistemele care oferă variantele cu blocare, paralelismul poate fi obținut prin folosirea firelor: fiecare cerere receptionată se tratează printr-un fir separat. Firele pot fi create dinamic, la primirea cererilor sau se poate menține un ansamblu de fire, create static, la lansarea în execuție (dar în acest caz, pot fi tratate

simultan un număr fixat de cereri - egal cu numărul firelor); structurile de date comune firelor trebuie accesate în regim de excludere reciprocă (deci trebuie protejate prin semafoare sau alte mecanisme).

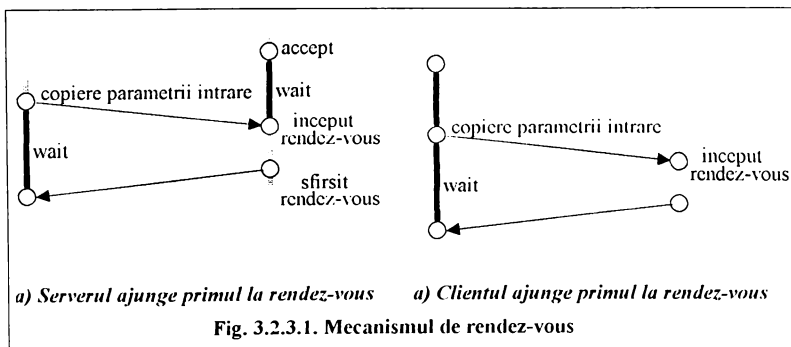
Programele din anexa A3 (client și server) exemplifică utilizarea primitivelor fără blocare. Programul client folosește o variantă de emisie care se poate încadra în tipul (1) de primitive *send* fără blocare (*no wait send*), iar programul receptor folosește o variantă de recepție analoga variantei (3) de mai sus. Recepția unui mesaj este semnalată printr-o întrerupere; la generarea întreruperii, se lansează în execuție funcția declansată ca Event-Synchronous-Routine (ESR), care are rolul de a prelua și memora mesajul într-un buffer liber. Serverul preia succesiv câte un mesaj din buffere, eliberează buferul respectiv, tratează cererea și trimite răspunsul înapoi la client; în tot acest timp, dacă sosește o altă cerere, ea este memorată pentru o tratare ulterioară.

3.2.3 Primitive cu și fără bufferare

S-a aratat în secțiunea precedentă că, în cazul primitivelor fără blocare este necesara menținerea unor buffere, atât la nivelul clienților, cât și la nivelul serverelor, pentru a permite memorarea mai multor cereri încă netrimise, respectiv încă nereceptionate. Utilizarea bufferelor ridică probleme specifice: alocare, eliberare, acces în regim de excludere reciproca. O modalitate des întâlnită în cadrul sistemelor de operare distribuite de a rezolva problema bufferării mesajelor este de a defini o structura de date de tip *cutie postală sau port*, la care se trimite toate mesajele aparținând unui destinatar; procesul destinatar crează cutia postală printr-o primitiva pusă la dispoziție de nucleu. Problema spațiului limitat nu este însă eliminată nici prin acest procedeu (cel mult apare mai rar situația bufferelor pline); ca atare, emitorul trebuie blocat la o operație *send*, pînă la primirea unui mesaj de confirmare că mesajul a fost preluat.

Exemplul din anexa 1 arată o modalitate de bufferare la nivelul serverului. Dacă numărul de clienți este foarte mare, s-ar putea să se depășească capacitatea de memorare a bufferelor; în acest caz, clienții repeta cererea, după o mică întârziere, pentru a nu sufoca serverul.

Intr-un sistem care nu folosește bufferarea, procesele se sincronizează în momentul transferului. Aceasta strategie este o extensie a problemei rendez-vous, cunoscuta din sistemele centralizate. Semantica pentru rendez-vous între procese aflate pe noduri diferite (la distanță) este aceeași ca și în procesele locale (fig.



3.2.3.1). Strategia este adevărată pentru modelul client/server și presupune circulația informației în ambele sensuri.

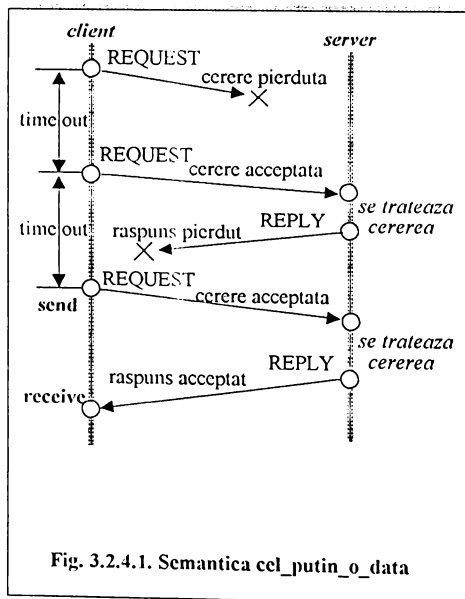
Pentru ca un client să declanșeze rendez-vous cu un server, sunt necesare următoarele:

- (1) clientul trebuie să cunoască definiția interfeței serverului;
- (2) clientul trebuie să cunoască localizarea serverului: identificatorul procesului și adresa nodului. Identificatorul procesului server poate fi determinat, având în vedere următoarele considerente:
 - (a) în fiecare nod se execută un server de nume, al cărui identificator este cunoscut; orice server este înregistrat de serverul de nume local, printr-un nume și informațiile de localizare;
 - (b) serverele de nume pot comunica între ele, tot prin rendez-vous; în acest mod, orice client poate determina localizarea unui server de pe orice nod.

3.2.4. Primitive sigure și nesigure

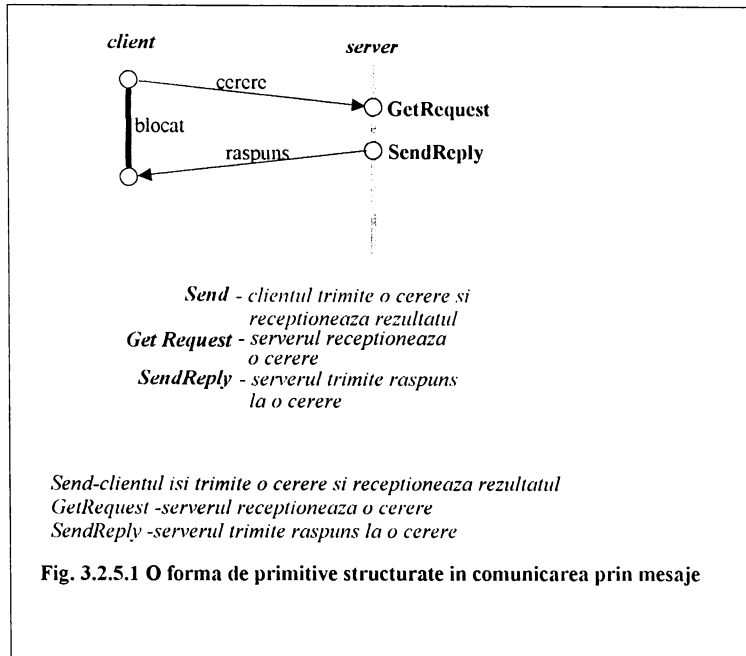
În cadrul unei comunicații interprocese sigure, la execuția primitivelor *send*, subsistemul de comunicație rezolvă problema mesajelor pierdute pe baza confirmărilor și retransmisilor, iar la execuția primitivelor *receive* se transmite confirmări. Aceasta înseamnă că în urma unei operații *send* sigure, terminate cu succes, procesul emitor este sigur că mesajul a fost recepționat (fiind confirmat). Comunicatiile folosind primitive nesigure, care nu garantează livrarea mesajelor, pot fi făcute însă sigure prin translatarea intrului mecanism de confirmare și retransmisii la nivelul proceselor. Unii autori [Mul89] chiar sugerează plasarea mecanismelor de recuperare a erorilor la nivel înalt. Pentru ambele soluții de recuperare a erorilor, există două tipuri de semantica posibilă, definite; în funcție de acțiunea pe care o întreprinde un server, la primirea unei cereri duplicate: (1) semantica *cel-puțin-o-dată*, prin care se asigură că o cerere a unui client este tratată cel puțin o dată (Fig. 3.2.4.1); (2) semantica *exact-o-dată*, prin care se asigură execuția numai o dată a codului de tratare a unei cereri; aceasta este absolut necesară în cazurile în care tratarea aceleiași cereri de mai multe ori conduce la erori (cererile a caror tratare se poate repeta fără prejudicii se numesc *idempotente*; un exemplu de cerere neidempotentă este o operație *append* într-un fișier).

Pentru a se implementa primitive de tipul (2), serverele trebuie să mențină o tabelă cu identificatorii cererilor tratate și răspunsurile trimise clienților pentru aceste cereri; la primirea unei cereri duplicate, serverul nu mai execută codul de tratare a cererii respective, ci doar transmite rezultatul.



3.2.5 Forme structurate ale comunicării prin mesaje

Pentru a se obține performanțe superioare în comunicația prin mesaje, aplicată modelului client-server, ținând cont că un mesaj REQUEST al unui client este urmat de un mesaj REPLY, se pot combina primitive *send* și *receive* de la nivelul clientului într-o singură primitivă *send*. Acest lucru simplifică bufferingul mesajelor, deoarece cererea poate rămâne în bufferul clientului, iar răspunsul serverului poate fi memorat în acest buffer și protocolul, deoarece controlul crorilor și al fluxului interpretează răspunsul la o cerere ca și o confirmare. Setul de primitive folosit este reprezentat în fig. 3.2.5.1.



3.2.6. Adresare directa și indirecta

Comunicația între procese prin mesaje dispune de două tehnici de adresare: *adresarea directă*, în cadrul căreia emițătorul specifică explicit procesul destinatar și *adresarea indirectă* în cadrul căreia emițătorul specifică drept destinație un nume global, denumit *cutie poștală* sau uneori *port*.

Adresarea directă în cele două variante, *simetrică* și *asimetrică* (fig. 3.2.6.1.a,b) are următoarele caracteristici: (1) între două procese care doresc să comunice, și care trebuie să-și cunoască reciproc identitatea, se stabilește automat o cale de comunicare; (2) între două procese care comunică, există numai o cale de comunicație; (3) comunicația este bidirecțională. *Adresarea directă* prezintă avantaje precum o implementare simplă și posibilitatea controlului momentelor de transmisie a mesajelor. În schimb, modificarea numelui unui proces impune modificarea în alte procese a tuturor referințelor la acesta; în plus, adresarea directă nu permite trimiterea unei cereri la mai multe servere identice. O variantă a adresării directe adoptată în două sisteme de operare distribuite experimentale ([Gos91]), DEMOS/MP și Charlotte este *comunicația prin legături*. O *legătură* este definită ca un canal de comunicație unidirecțional sau bidirecțional la un proces utilizator, proces de sistem sau nucleu. Fiecărui proces i se atașează o tabelă de legături, care identifică în fiecare intrare o linie de comunicație; procesul poate crea, distruge legături și transfera un capăt al unei legături la alt proces (fig. 3.2.6.1.). O legătură poate fi privită ca o conexiune, într-un protocol orientat pe conexiuni, în care se poate lucra cu sau fără blocare, cu sau fără buffering de mesaje.

În cadrul adresării *indirecte*, o cutie poștală poate apărea ca destinație în orice primitivă *send*, sau ca sursa în orice primitivă *receive*. Deși o astfel de comunicație este adecvată pentru modelul client-server, implementarea ei fără un suport specializat în rețea (un mesaj trebuie transmis de subsistemul de comunicație către mai multe noduri) poate fi neperformantă. Ca atare, în sistemele de operare distribuite se implementează,

Send (P_process, mesaj);# se trimite un mesaj procesului P.
 Receive (Q_process, buffer);# se receptioneaza un mesaj de la procesul Q
 a) Adresare directa simetrica

Send (P_process, mesaj);# se trimite un mesaj procesului P.
 Receive(id, buffer);# se receptioneaza un mesaj de la orice proces; id precia numele procesului emitator
 b) Adresare directa nesimetrica

Send (link, mesaj);# se trimite un mesaj pe legatura link.
 Receive(link, buffer);# se receptioneaza un mesaj din legatura link.
 Send (link, buffer, link_end);# se transfera un capat al unei legaturi pe legatura link si se depune in buffer
 c) Adresare directa cu legaturi

Send (P_port, mesaj);# se transmite un mesaj in portul P.
 Receive (P_port, mesaj);# se receptioneaza un mesaj din portul P.
 Send (P_port, port);# se transporta dreptul de proprietate la un port
 d) Adresare indirecta

Fig. 3.2.6.1. Primitive pentru adresarea directa si indirecta.

in general, un tip simplificat particular de cutie poștală numit *port*. Un *port* este un obiect al nucleului, protejat, catre care procesele pot să transmită, respectiv extrage mesaje. Comunicarea prin porturi are următoarele caracteristici: (1) asociază o cale de comunicare la mai mult de două procese; (2) între două procese se pot stabili mai multe căi de comunicare, fiecare corespunzând unui port; (3) comunicația poate fi fie unidirecționala, fie bidirecționala. Accesul pentru extragerea mesajelor din porturi a proceselor poate fi condiționat prin drepturi de recepție (de exemplu, prin obținerea unei capabilitati). Porturile trebuie declarate; declararea unui port înseamnă asocierea la port a unei cozi de mesaje, gestionată în mod normal printr-un algoritim FIFO (ordinea FIFO poate fi modificată prin mesaje urgente, care se tratează în mod special). Porturile pot avea ca proprietar fie un proces (1), fie sistemul de operare (2). Cazul (1) implică următoarele: (a) creatorul portului este și proprietarul; el poate extrage mesaje din port, spre deosebire de celelalte procese care pot doar trimite mesaje la acest port; (b) portul este atașat procesului proprietar, astfel încât distrugerea acestui proces înseamnă și distrugerea portului; oricărui alte procese, care transmit în continuare mesaje la acest port, va trebui să li se semnaleze o eroare (eventual printr-un mecanism de tratare excepții); (c) există 2 modalități de stabilire a proprietarului și utilizatorilor unui port: fie se permite unui procesor, care devine astfel proprietar, declararea unei variabile port, fie se declară un port partajat și apoi se stabilește proprietarul; (d) proprietarul unui port poate, printr-o primitivă, transmite altui proces drepturile de proprietate (fig. 3.4.6.1.d). Cazul (2) presupune că sistemul de operare pune la dispoziție primitive pentru crearea unui port, (procesul apelant este și proprietarul portului) distrugerea unui port și trimiterea și recepția mesajelor. În cazul în care un proces încearcă să trimită mesaje într-un port în care nu mai există loc liber, există trei soluții: (1) procesul este blocat pînă la plasarea mesajului; (2) procesului i se semnalează o eroare; (3) mesajul este acceptat de nucleu și trimis de acesta atunci cînd mesajul poate fi plasat.

Porturile pot fi *statice* (porturi prin care o pereche de procese comunică pe tot parcursul duratei de viață) sau *dinamice* (porturi create, asociate și utilizate pentru o perioadă de timp). Porturile statice sunt adecvate pentru o interacțiune între procese de tipul corutinelor (fig. 3.2.6.2.) sau proceselor distribuite care necesita partajarea unor informații împreună cu un anumit nivel de sincronizare.

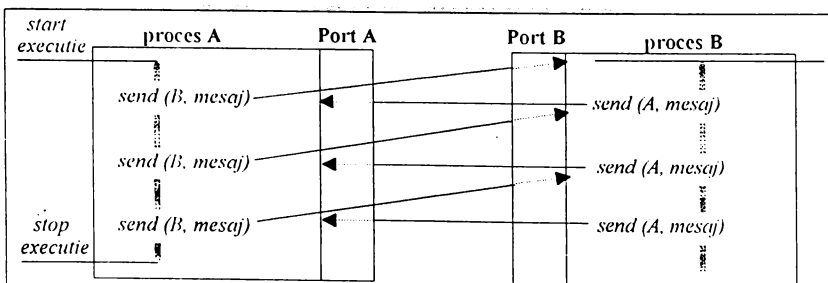


Fig. 3.2.6.2. Comunicarea prin porturi folosita pentru a obtine o interactiune de tip corutine.

3.2.7. Excepții

În aceasta secțiune se iau în considerare probleme legate de terminarea procesului emițător sau receptor sau de pierderea mesajelor în rețea. Dacă procesul receptor se termină fără ca emițătorul să cunoască acest lucru, emițătorul va continua să transmită mesaje. Pentru primitivile cu bufferare, execuția continuă, eventual până la depășirea capacității bufferelor. Dacă emițătorul trebuie să știe dacă mesajul a fost preluat de receptor, atunci trebuie folosit un protocol cu confirmări. Când se folosește o primitivă *send* fără bufferare, pentru ca emițătorul să nu se blocheze la infinit, sistemul de operare trebuie fie să termine emițătorul, fie să semnaleze că receptorul s-a terminat. Analog trebuie procedat în cazul când un server este blocat în așteptarea unui mesaj de la un client care s-a terminat. Un exemplu în acest sens este scrierea de către un proces într-un soclu, la capatul căruia procesul server este terminat, în sistemele UNIX. În acest caz, sistemul de operare generează pentru procesul client semnalul SIGPIPE a cărui tratare implicita este terminarea procesului.

În ceea ce privește pierderea mesajelor, în secțiunea 3.1.3. s-a arătat că prin mecanismul timeout din protocolul software acestea pot fi depistate. Problema care se ridică aici este la care nivel se face recuperarea erorilor. Există 3 soluții: (1) sistemul de operare este responsabil pentru detecție și pentru transmisie; (2) procesul emițător este responsabil pentru detecție și retransmisie; (3) sistemul de operare este responsabil pentru detecție, el semnalează emițătorului că mesajul a fost pierdut și retransmisia este lăsată pe seama procesului emițător (dacă aceasta o dorește). În general pentru alegerea unei soluții trebuie luate în considerare următoarele observații: (1) acțiunea depinde de și de specificul aplicației; (2) erorile necritice pot fi ignorate; cele care duc însă la distrugerea unei componente a sistemului trebuie tratate; (3) în cazurile în care o eroare dintr-un server are efecte globale, trebuie aleasă o comunicare prin tranzacții atomice (vezi capitolul 7).

3.2.8 Exemplu. Comunicație prin porturi în Mach

Arhitectura sistemului distribuit Mach, creat la Carnegie-Mellon de o echipă de cercetători conduși de Richard Rashid ([Ras88] [Ras81]) este prezentată în fig. 3.2.8.1.

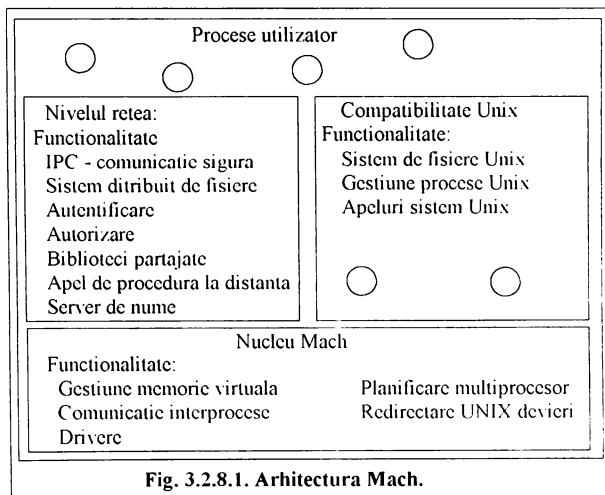


Fig. 3.2.8.1. Arhitectura Mach.

Comunicația interprocese în Mach se bazează în întregime pe *porturi și mesaje*. Un *port* este definit ca un canal de comunicație, implementat printr-o coadă de mesaje protejată de nucleu, de lungime finită. În Mach, porturile sunt referințele pentru obiecte: porturile reprezintă obiectele; Mach este proiectat utilizând modelul obiect (vezi subcapitolul 1.4.2). Un port poate avea mai mulți emițători, dar un singur receptor. Drepturile de acces la un port constau în acces pentru citire, scriere, și dreptul de proprietar. Când într-un task se creează un nou port, taskul primește toate cele trei drepturi de acces la port; un fir al taskului nu poate referi decât un port la care taskul respectiv are dreptul de acces. Un task poate transmite drepturile de acces la un port altui task, utilizând mesaje care conțin capacități pentru port. Dacă un task transmite dreptul de recepție la un port, el pierde acest drept (dreptul de emisie la un port nu se pierde); în perioada de transmisie a dreptului de recepție, toate mesajele transmise la acel port sunt reținute de nucleu, care le transmite apoi noului task ce a primit dreptul de recepție. Un *mesaj* în Mach este obiect de date, utilizat pentru comunicație între fire. Mesajele nu sunt limitate ca dimensiune, dar au header de lungime fixă care conține: *msg size*, dimensiunea mesajului de transmis sau dimensiunea maximă a mesajului de recepționat, *msg local port*.

msg_remote_port, numele porturilor locale sau la distanță la/de la care mesajul trebuie transmis/ recepționat, *msg_id* un identificator al tipului mesajului; în partea de lungime variabilă, mesajele pot conține și pointeri sau capabilități pentru porturi. Primitivele de lucru cu porturile sunt următoarele: (1) *msg_send(mesaj, optiune, timeout)* pentru transmiterea unui mesaj la portul specificat în headerul mesajului; argumentul *timeout* specifică timpul maxim de așteptare pentru cazul în care portul destinație este plin cu mesaje. Receptorul poate utiliza *msg_local_port* pentru a trimite un răspuns emițătorului (cu excepția cazului cînd *msg_local_port* este PORT_NULL); (2) *msg_receive(mesaj, optiune, timeout)* pentru recepția unui mesaj de la portul specificat în headerul mesajului, *msg_local_port* sau în grupul de porturi implicite; *timeout* este intervalul maxim de așteptare a sosirii unui mesaj; (3) *msg_rpc(mesaj, optiune, rcv_dim, send_timeout, receive_timeout)* pentru trimiterea unui mesaj, urmat de recepția unui răspuns, utilizînd același buffer; *msg_remote_port* specifică portul la care mesajul trebuie transmis, iar *msg_local_port* specifică portul la care mesajul este recepționat; *send_timeout* și *receive_timeout* sunt timpurile de așteptare la emisie, respectiv la recepție. Pentru toate primitivele (1), (2), (3) *optiune* specifică condițiile de eroare în care primitivele trebuie să fie terminate.

Primitiva *msg_rpc()* oferă facilitati de apel de procedură la distanță prin intermediul porturilor.

Pentru manipularea drepturilor de acces la porturi și recepția mesajelor de control, Mach mai prevede următoarele primitive: (1) *port_allocate(task, port)* crearea unui port, obținându-se toate cele trei tipuri de drepturi; identificatorul portului se returnează în port; (2) *port_deallocate(task, port)* distruge drepturile de acces la port; dacă taskul era proprietar, portul este distrus și se anunta toate taskurile cu drept de emisie la acest port; dacă taskul avea dreptul de recepție, el pierde acest drept, care se transmite proprietarului; (3) *port_enable(task, port)* adauga acest port la grupul de porturi implicite pentru *msg_receive*; (4) *port_disable(task, port)* elimină acest port din grupul de porturi implicite; (5) *port_messages(task, porturi, contor_porturi)* returnează pentru un task specificat, un tablou de porturi active cu mesaje în așteptare; (6) *port_status(task, port, num_msg, backlog, receptor, proprietar)* returnează informații de stare relative la port: numărul de mesaje în așteptare (*num_msg*), numărul de mesaje care mai pot fi primite fără a se bloca emițătorul (*backlog*), indicatorii ai drepturilor de acces, pentru recepție (*receptor*), respectiv dreptul de proprietar (*proprietar*); (7) *port_backlog(task, port, backlog)* returnează numărul de mesaje care mai pot fi trimise fără a bloca emițătorul la acest port.

În Mach, mesajele pot fi trimise sincron sau asincron; în ultimul caz, recepționarea unui mesaj se notifică printr-un semnal Unix; se pot folosi însă și fire pentru tratarea operațiilor de recepție asincrone. Pentru transmiterea mesajelor la distanță se folosesc serverele de rețea astfel încît se asigură transparența schimbului de mesaje între procese (fig.3.2.8.2). Nucleul Mach nu recunoaște o transmisie la distanță; din punctul de vedere al nucleului, toate comunicațiile interprocese se fac local. Serverele de rețea, în afară rolului extinderii mecanismului de comunicație la procesele din rețea, asigură și siguranța transmisiei și conversia datelor la transmisie.

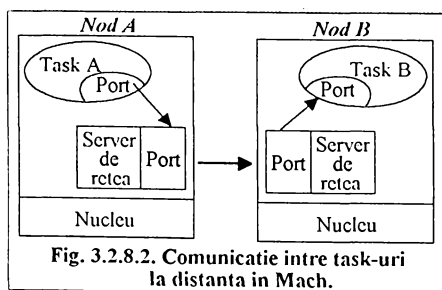


Fig. 3.2.8.2. Comunicatie între task-uri la distanță în Mach.

3.3. Comunicarea între procese prin conducte

Comunicarea între procese prin conducte (*pipes*) este un mecanism care permite transferul datelor între procese utilizând strategia FIFO, concomitent cu sincronizarea proceselor implicate în transfer. Ea poate fi aplicată și modelului client server prezentat anterior deși nu îi este caracteristică.

3.3.1. Prezentarea generală a comunicației prin conducte

În mod tradițional, conductele au fost implementate utilizând sistemul de gestiune al fișierelor. Proprietățile principale ale conductelor sunt acelea ca: (1) permit schimbul de informație între procese care nu-și cunosc identitatea; (2) permit transmiterea datelor în avalansa. (3) Iesirea unuia sau mai multor procese poate fi conectată la intrarea unuia sau mai multor procese.

În Unix o conductă apare ca un flux de octeți; nu există mecanism de transmisie out-of-band iar bufferul ei este de capacitate limitată. Când un proces încearcă să citească date dintr-o conductă vidă, se blochează până când conducta este completată cu volumul de date necesar. Un proces care încearcă să citească date dintr-o conductă la capatul căreia nu mai există nici un proces care scrie primește indicația end-of-file; în schimb dacă un proces încearcă să scrie într-o conductă la capatul căreia nu există cititori se generează semnalul SIGPIPE. Un proces care așteaptă pentru citire dintr-o conductă, sau altul care așteaptă pentru scriere (când bufferul conductei este plin) poate fi întrerupt. Există două tipuri de conducte: *conducte fără nume* și *cu nume*; diferența între cele două este dată de modul în care sunt create și accesate inițial. Astfel, în sistemul UNIX conductele fără nume se creează cu apelul sistem *pipe* în timp ce conductele cu nume se creează cu apelul *mknod*; conductele fără nume nu pot fi accesate decât din procesele descendente ale procesului care le-a creat. În ambele cazuri, accesul la conducte se face cu apelurile *read*, *write*. O conductă cu nume este luată în evidență de sistemul de gestiune al fișierelor; unul sau mai multe procese pot să deschidă conducta specificând modul *read*, în timp ce altele pot specifica la deschidere modul *write*. După deschidere, o conductă cu nume se comportă exact ca și o conductă fără nume.

Pentru comunicația la distanță prin pipe sistemul de operare distribuit LOCUS a extins ambele tipuri de conducte.

3.3.2 Comunicarea prin conducte în Microsoft WIN32

Microsoft WIN32 furnizează pentru comunicația la distanță *conducte cu nume*, pentru a căror implementare utilizează servicii orientate pe conexiune; ca atare trebuie să existe un procedeu pentru ca aplicațiile server să aștepte cererile de conectare. Acest procedeu constă în: (1) *crearea unei conducte* cu nume, executând apelul *creatNamedPipe()*; (2) așteptarea cererilor de conectare, după execuția apelului *ConnectNamedPipe()*. Conductele au câteva caracteristici specifice care stabilesc cu apelul *creatNamedPipe()*: (1) *direcția*. O conductă poate fi *bidirecțională sau unidirecțională* client-la-server sau server-la-client; (2) *modul de scriere*. În conducte se pot scrie *mesaje* (modul *mesaj*) sau *octeți* (modul *byte*) și se pot citi mesaje și octeți. (3) *modul citire*. Pentru o conductă din care se citește în modul byte, operațiile de scriere și citire nu trebuie să manipuleze simultan același număr de octeți; în schimb, pentru o conductă mod mesaj, fiecare scriere în conductă definește un mesaj, astfel că citirile trebuie efectuate la nivel de mesaje. Încercarea de a citi un număr mai mare de date decât un mesaj, are ca efect returnarea numai a mesajului; în schimb încercarea de a citi parțial un mesaj returnează eroare; (4) *modul de blocare*. Există două moduri de lucru, cu blocare sau fără; ele se referă la blocarea apelantului unei funcții *ConnectNamedPipe()*. Modul fără blocare este rar folosit, fiind prevăzut pentru compatibilitate cu conductele din OS/2. În schimb, ca și fișierele standard, conductele pot fi deschise pentru operații de I/O asincrone (overlapped).

Pentru contactarea unui server printr-o conductă, un client trebuie să execute apelul *creatFile()*, care deschide conducta cu nume, (numele fiind scris sub o formă specială: \\<nume_nod>\PIPE\<nume_conducta>) implicit în modul cu blocare și citire la nivel de octet. Pentru a modifica aceste setări implicite, se poate utiliza funcția *SetNamedPipeHandleState()*. După deschiderea conductei, clientul și serverul pot utiliza pentru schimbul de date funcțiile *ReadFile()* și *WriteFile()* (diferențele față de utilizarea unui fișier și o conductă fiind în acest caz transparente) sau funcția *TransactNamedPipe()*, care permite scrierea unui mesaj într-o conductă și așteptarea unui răspuns (utilă pentru modelul client server).

Serverele pot fi proiectate să funcționeze *multifir*. În acest sens, o soluție este să se creeze mai multe instanțe ale unei conducte, iar fiecărei instanțe să i se asocieze un fir în cadrul căruia se așteaptă cereri de conectare și se deservească cereri de servicii. Schema acestui mod de lucru este prezentată în figura 3.3.2.1.

```

const short NO_PIPES=3;
HANDLE hPipeInstance[NO_PIPES], hPipeThread[NO_PIPES];
for (int i=0; i<NO_PIPES; ++i) {
    hPipeInstance[i] = CreatNamedPipe (...);
    if (hPipeInstance[i] != INVALID_HANDLE_VALUE) {
        hPipeThread[i] = CreatThread (NULL, 0, NamedPipeThread, hPipeInstance[i], ...);
    }
}
DWORD WINAPI NamedPipeThread (HANDLE hPipeInstance) {
    DWORD dwError;
    BYTE byBuffer[4096];
    DWORD dwBytes;
    while (ConnectNamedPipe (hPipeInstance, NULL) || GetLastError() == ERROR_PIPE_CONNECTED)
        while (ReadFile (hPipeInstance, byBuffer, sizeof(byBuffer), &dwBytes, NULL)) {
            //trateaza cererea
        }
    DisconnectNamedPipe (hPipeInstance);
}
ExitThread (dwError = GetLastError());
return dwError;
}

```

Fig. 3.3.2.1. Utilizarea conductelor in servere multifir

3.4. Modelul apelului de procedură la distanță

Modelul de comunicație bazat pe *apel de procedură la distanță (RPC=Remote Procedure Call)* permite unui proces să apeleze și să execute o procedură de la distanță ca și pe o procedură locală. Pe timpul execuției procedurii de la distanță, apelantul este blocat, pînă cînd se receptionează rezultatul. Prin RPC se permite astfel extinderea apelului de proceduri, folosit în limbajele de programare la mediile distribuite. Modelul RPC este atractiv pentru că oferă un mecanism de comunicație în cadrul programelor distribuite, care se poate caracteriza astfel: (1) este simplu și familiar (datorita similarității cu apelurile de proceduri locale); (2) este general (apelul de procedură este folosit în toate aplicațiile); (3) poate fi implementat eficient.

Modelul RPC simplifică construcția programelor distribuite, permițînd ignorarea de către acestea a detaliilor de comunicație și a erorilor de transmisie.

Modelul RPC a fost utilizat în multe sisteme distribuite începînd cu Argus ([LS83]), HRPC (University of Washington, [BC85]), Eden [ABL85], Spice (Carnegie Melon University, [JRT85]) și terminînd cu variante pentru sisteme Unix, precum Xerox Courier RPC, creat la Xerox Parc pentru a fi utilizat în mediile MESA și Cedar [BN84], SUN RPC (SUN Microsystems [Sun85]), V Kernel creat de Cheriton la Stanford University, Apollo RPC (Hewlett Packard) și mai recente OSF DCE (Open Software Fondation Distributed Computing Environment) și Microsoft Windows NT și Windows '95. Printre acestea, Argus prezintă o caracteristică specială: aceea că integrează modelul RPC cu mecanismul tranzacțiilor atomice. Unele sisteme din cele de mai sus (Argus, Cedar) au fost inițial proiectate pentru a lucra cu un singur limbaj, în timp ce altele au fost destinate pentru lucrul cu diverse limbaje și chiar în medii eterogene.

Deși simplitatea modelului RPC îl recomandă ca o paradigmă principală pentru comunicația în aplicații distribuite, performanțele oferite de acesta pot să nu fie totuși adecvate pentru unele aplicații. Aceasta deoarece, sistemele RPC, în majoritate, sunt proiectate să furnizeze un răspuns cît mai rapid unei cereri, astfel încît apelurile individuale sunt deservite relativ repede, dar productivitatea este mică. Pentru aplicațiile în care se cere o productivitate mare de deservire a cererilor, rămănînd pe planul doi rapiditatea răspunsului, modelul RPC s-ar putea să nu ofere performanțele dorite.

3.4.1. Prezentarea generală a modelului

Ideal este ca un apel de procedură la distanță să aibă aceeași semantică și sintaxa ca și un apel de procedură locală. În practică însă, acest ideal este dificil de realizat datorită următoarelor cauze: (1) erorile de transmisie sunt greu de mascat complet; (2) procesele au spații de adrese diferite; (3) procesele au durate de viață diferite.

Modelul RPC trebuie să fie dotat cu un protocol de transport, care să furnizeze transmisia parametrilor procedurii ca și a rezultatului; uneori chiar termenul de RPC este folosit pentru a desemna acest protocol de transport.

Comunicația prin mesaje și extensiile sale (formele structurate de mesaje) furnizează servicii rareori situate deasupra nivelului de transport al modelului ISO/OSI. Aceste funcții sunt în general încapsulate în programele de aplicații. Spre deosebire de comunicația prin mesaje, protocolul RPC este responsabil și de maparea caracteristicilor de limbaj ale unui apel de procedură în comunicațiile oferite de nivelul transport al rețelei.

Deși comunicația prin mesaje și apelul de procedură la distanță sunt analoge (în secțiunea 3.2.2. fig. 3.2.2.1. se arată patru moduri de interacțiune a proceselor de la distanță, după modul de returnare a controlului apelantului, apelul de procedură fiind clasificat ca al 4-lea mod), ele diferă prin următoarele: (1) în timp ce la comunicația prin mesaje, procesele care interacționează pot fi vazute ca parteneri, la RPC relația dintre ele este de tip master-slave; (2) în comunicația prin mesaje, procesele completează explicit cîmpurile unui mesaj înainte de transmisie, ca și despachetarea rezultatelor dintr-un mesaj.

Andrews și Schneider [AS83], identifică două moduri de a specifica un serviciu al unui server într-un mecanism RPC (fig. 3.4.1.1.):

(1) serviciul se declară ca o procedură dintr-un limbaj secvențial.

```
remote procedură serviciu ( in parametrii_cu_valoare:
                           out parametrii_rezultat)
    corp
end
a)declarare ca o procedură
accept serviciu ( in parametrii_cu_valoare:
                  out parametrii_rezultat) -> corp
b)declarare ca o instrucțiune
```

Fig. 3.4.1.1. Două moduri de specificare a serviciilor în modelul RPC

În acest caz procedura se implementează într-un proces server, acesta așteaptă sosirea unui mesaj conținând valorile argumentelor de la un proces client, asignează aceste valori parametrilor procedurii, execută procedura, și returnează un mesaj de răspuns conținând valorile parametrilor rezultat. Procesul client va fi blocat pe parcursul execuției procedurii chiar dacă nu există parametrii rezultat.

(2) serviciul se declară ca o instrucțiune, care poate fi plasată oriunde în procesul server. Interacțiunea între procesul server și client este de tipul rendez-vous.

În continuarea acestui capitol se va considera modelul RPC numai pentru primul mod.

Implementarea modelului RPC este prezentată sub o formă generală în fig. 3.4.1.2. Se observă că sunt necesare cinci elemente pentru a se obține efectul execuției unei proceduri locale într-un mediu distribuit: procesul client, codul stub pentru client, mecanismul de transport, codul stub pentru server și procesul server. Mecanismul de transport implementează transmisia sigură de mesaje; acesta poate utiliza protocoale de la nivelul transport, însă tendința implementărilor actuale este de a utiliza un protocol (neorientat pe conexiune) bazat pe datagrame scrise special pentru RPC, care să asigure o transmisie sigură a mesajelor. Modulele stub au ca scop împachetarea și despachetarea parametrilor în mesaje.

Modelul RPC poate fi interpretat și ca o extindere a primitivelor sigure, cu blocare, structurate *Send*, *GetRequest*, *SendResponse*. Spre deosebire de aceste primitive, datorită simplității și totodată generalității sale, modelul RPC poate fi însă extins la mai mulți clienți și la mai multe servere, ca și la medii eterogene.

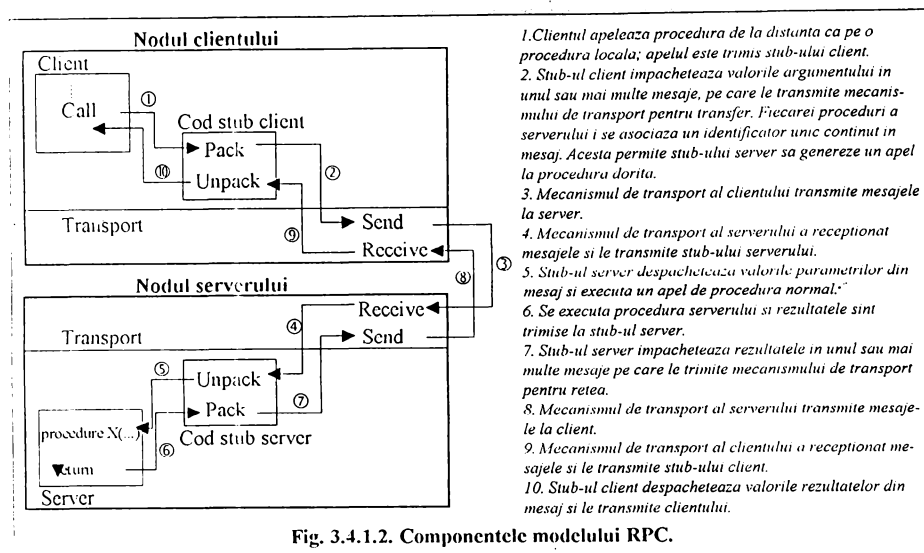


Fig. 3.4.1.2. Componentele modelului RPC.

3.4.2. Transmiserea parametrilor și rezultatelor

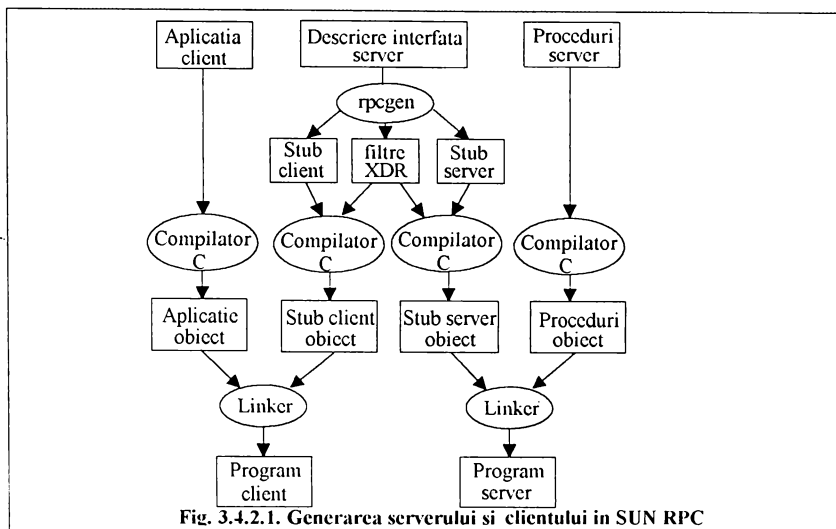
O problema cheie a modelului RPC este transmiterea parametrilor și reprezentarea parametrilor și rezultatelor în mesaje.

Parametrii pot fi transmiși prin *valoare* sau *referință*; dacă transmiterea prin valoare nu pune nici o problema (valorile parametrilor se copiază direct în mesaj), în schimb transmiterea parametrilor prin referință și a pointerilor ridică probleme. O soluție este de a înlocui transmiterea parametrilor prin referință prin transmitere prin copiere; dacă parametrul este doar intrare pentru server, atunci nu mai este necesară copierea lui la loc în client iar dacă parametrul este doar ieșire nu este necesară transmiterea valorii lui la server.

Ca atare, se asociază fiecărei declarații a unei proceduri o specificație formală numită *interfață* (prin care se precizează tipul parametrilor, dacă sunt de intrare sau de ieșire (sau ambele) și dimensiunea lor (maximă); această specificație este folosită de un compilator special care generează codul stub, stub-ul se linkidează împreună cu clientul, respectiv serverul, pentru a avea acces la spațiul lor de adrese. În fig. 3.4.2.1 se arată ca exemplu generarea programelor client și server în SUN RPC.

Definirea interfeței procedurilor serverului are mai multe scopuri: (1) primul este cel arătat mai sus, de a permite serverului (care exportă procedurile) și clientului (care importă) să se pună de acord asupra parametrilor de intrare și de ieșire. Impachetarea și despachetarea parametrilor în/din mesaje ca și dispecerizarea apelurilor la proceduri în stub-ul server se va genera pe baza definiției interfeței; de aceea interfața conține caracteristici ale procedurilor furnizate de server care trebuie să fie vizibile clientului (numele procedurilor și tipurile parametrilor lor). Unele sisteme permit procedurilor apelate la distanță să semnaleze excepții; aceste tipuri de excepții trebuie documentate și ele uneori în descrierea interfeței. Mecanismele RPC includ, în general, și un limbaj de interfață, IDL (Interface Description Language) pentru a permite programatorilor să definească interfața. Un astfel de limbaj de interfață furnizează posibilitatea de a defini tipuri scalare (char, boolean, int, real), facilități de definire a tipurilor structurate (tablouri, stringuri, înregistrări), facilități pentru utilizarea pointerilor. Definirea interfeței într-un astfel de limbaj este procesată de un compilator special. (2) al doilea scop este că declararea interfeței servește ca bază pentru verificarea tipurilor parametrilor la compilare.; (3) reprezentarea datelor la cele două capete poate fi diferită, declarația interfeței poate fi folosită ca bază pentru generarea conversiei de cod necesară.

Problema conversiei între reprezentări diferite ale datelor la cele două capete poate fi rezolvată în două moduri: (a) prin definirea unei reprezentări standard pentru fiecare tip de date (din mesaje) și conversia mesajelor, în fiecare capăt, la reprezentarea standard respectivă. Aceasta soluție oferă portabilitate între diverse sisteme și permite programelor de aplicații ignorarea unor detalii de reprezentare a datelor precum ordinea octetilor (*big* sau *little endian*). (b) Clientul trimite mesajele utilizând propria sa reprezentare a datelor (eventual cu cativa indicatori pentru identificarea reprezentării utilizate) iar serverul realizează conversiile necesare. Aceasta soluție are dezavantajul că serverul trebuie să poată accepta și converti orice reprezentare particulară (problemele cresc când se adaugă noi calculatoare sau limbaje în sistem), dar și avantajul că nu este necesară nici o conversie în cazul comunicației între două sisteme cu aceeași arhitectură.



Ambele soluții sunt utilizate. Astfel, SUN RPC și Xerox Courier folosesc prima soluție; SUN RPC definește reprezentarea standard XDR, care impune o ordine big-indian și o dimensiune minimă pentru orice tip de 32 biți (aceasta înseamnă că pentru a transmite un întreg de 16 biți, se folosesc 32 biți în ordinea big-endian) iar Xerox Courier RPC definește o reprezentare tot cu ordinea big-endian dar cu dimensiunea minimă de 16 biți. Apollo RPC și DCE RPC folosesc soluția a doua. Apollo RPC permite mai multe formate (big sau little endian, formatele IEEE, VAX, IBM și Cray pentru numere reale); emițătorul transmite în formatul sau propriu (dacă este unul din formatele suportate) și receptorul realizează conversia (soluția a doua); DCE RPC permite, analog, mai multe formate și face conversia tot la receptor (metoda se numește *receiver makes right*). Ambele soluții utilizează ceea ce se numește *implicit typing*, adică se transmite valoarea unei variabile (într-un format standard) și nu tipul acelei variabile, în contrast cu tehnica ISO OSI și ASN.1, care transmite în mesaje tipul fiecărui câmp de date (codificat pe 8 biți) și apoi valoarea aceluși câmp (*explicit typing*).

3.4.3. Legarea clienților la servere

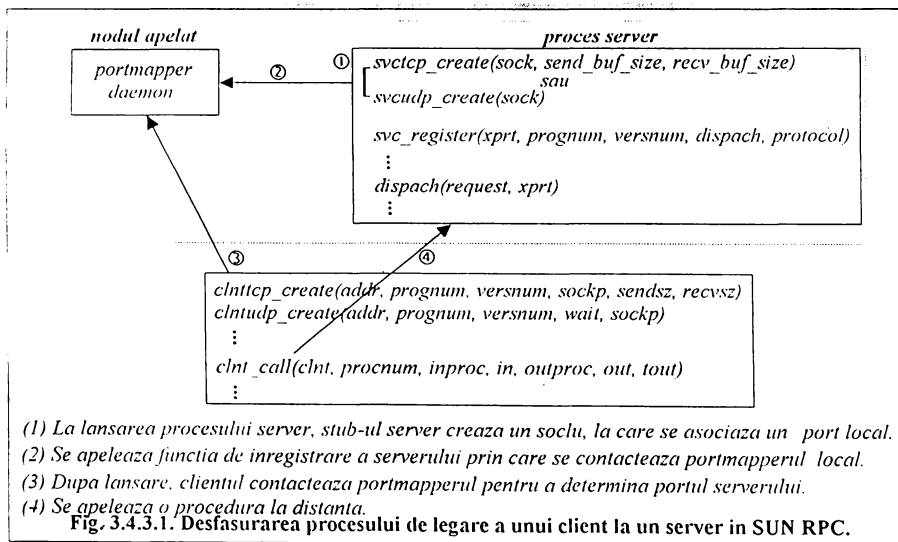
Legarea clienților la servere se referă la contactarea aceluși sistem de la distanță care deține procedura de executat. Există două componente ale procesului de legare: (1) localizarea calculatorului pentru serviciul dorit; (2) localizarea procesului server de pe acel calculator.

O primă tehnică impune ca fiecare client să cunoască nodul care trebuie contactat; un superserver de pe acest nod, a cărui adresă este publică, înregistrează adresele (porturile) tuturor serverelor disponibile de pe nod. O altă tehnică este transmiterea unui mesaj broadcast (sau multicast) în vederea obținerii localizării de la superservere. O tehnică mai generală este utilizarea unei baze de date centralizate, furnizată de o autoritate centrală în sistem (un server de nume denumit *binder*). La începutul execuției sale, orice server va exporta interfața sa; procedura este denumită *înregistrarea serverului în baza de date*. Pentru înregistrare, serverul livrează numele său, un număr de versiune, un identificator unic (în general de 32 biți) și un *handle* pentru a fi localizat; handle-ul este o caracteristică a sistemului și poate fi o adresă IP, o adresă Ethernet, un identificator de proces. Pot fi livrate și alte informații, ca de exemplu legate de autorizarea clienților. Funcția de înregistrare a serverului este pusă la dispoziția programelor de o bibliotecă de funcții run-time; de asemenea, pot fi furnizate și funcții pentru deînregistrarea unui server, sau obținerea handle-ului și identificatorului unic al unui server ale cărui nume și versiune se cunosc.

Legarea clientului la server se face de către stub-ul client, când clientul cere pentru prima oară execuția unei proceduri la distanță; stubul client trimite un mesaj la binder cerind importul interfeței serverului X, versiunea Y (X, Y sunt cunoscute, stub-ul client fiind compilat din definiția interfeței serverului, un model al procesului de compilare și legare este prezentat în figura 3.4.2.1). Serverul de nume returnează clientului handle-ul și identificatorul unic al serverului; în continuare clientul utilizează handle-ul ca adresă de destinație a mesajelor, în timp ce identificatorul unic este inclus în fiecare mesaj, și folosit pentru a distinge un server dintre mai multe executate pe aceeași mașină. Aceasta tehnică, a exportului și importului interfeței serverului oferă un grad mare de flexibilitate: serverul de nume poate asocia clienții la servere în așa fel încât să echilibreze încărcarea în sistem, poate contribui la mecanismul de autentificare al clienților, sau poate testa

periodic serverele înregistrate, pentru a deinregistra automat serverele care nu mai raspund la apeluri. Tehnica are însă și dezavantajul că necesită o circulație de mesaje în plus pentru exportul și importul interfețelor; dacă se menține un singur server de nume pentru toti clienții, într-un număr mare, serverul de nume poate fi sufocat de cereri, astfel încât sunt necesare mai multe servere de nume, a caror sincronizare pentru actualizarea informațiilor crește și mai mult regia de sistem.

Se prezintă în continuare, ca exemplu, modul de legare al clienților la servere adoptat în SUN RPC[Ste90] (RPC este o componentă în ONC= Open Network Computing). Mecanismul SUN RPC este asigurat de 3 elemente: (1) *rpcgen*, compilatorul care, primind descrierea interfeței procedurilor generează codurile stub (server și client); (2) *XDR (eXternal Date Representation)* standardul de codificare a informației pentru compatibilitate între sisteme diferite; (3) o bibliotecă de funcții run-time. Un serviciu, implementat de un server, este reprezentat de doi identificatori (numere întregi fără semn pe 32 biți): numărul de program și versiunea lui. Serviciile se înregistrează într-un superserver denumit *portmapper* cu următoarele informații: numărul de program, numărul versiunii, protocolul UDP sau TCP, numărul portului, numele serviciului și numele rutinei de dispecerizare. *Portmapper*-ul este lansat ca un proces daemon la startarea sistemelor; printr-o comandă Unix, folosind diverse opțiuni, se poate afișa lista serviciilor înregistrate, a serviciilor ce folosesc porturi UDP sau TCP, se pot deinregistra servicii. Procedurilor unui serviciu li se asignează numere. Numărul 0 se asignează unei proceduri speciale, care nu primește argumente și nu calculează rezultate doar returnează controlul apelantului, generate automat de compilatorul *rpcgen*, și folosită pentru a verifica dacă serviciul este în execuție sau nu. Apelul unei proceduri la distanță se face după crearea unui handle client (UDP sau TCP) care precizează adresa serverului (portul poate fi determinat de *portmapper*), numele de program, versiunea și soclul clientului; după obținerea handle-ului se pot apela proceduri, indicându-se handle-ul clientului, numărul procedurii, argumentele și rutinele *xdr* (filtre *xdr*) care serializează argumentele, respectiv deserializează rezultatele. Procesul de legare, descris în figura 3.4.3.1, decurge astfel: (1) la lansarea procesului server, stub-ul server creează un soclu, la care se asociază un port local; (2) se apelează funcția de înregistrare a serviciului (*svcregister*) prin care se contactează *portmapper*-ul local pentru înregistrarea serverului; (3) se lansează clientul, care mai întâi, contactează *portmapper*-ul de pe nodul indicat (*clnt_create()*) pentru a determina portul serverului; (4) la apelul unei proceduri de la distanță, stub-ul client serializează argumentele și le transmite stub-ului server, folosind portul determinat la (3).



Procesul *portmapper* de pe nodul serverului este contactat de clienți pentru a se localiza un program cu un anumit număr și cu o anumită versiune. Clientul trebuie să specifice în plus doar nodul serverului și protocolul de transport. Procesul *portmapper* furnizează și suport pentru tratarea cererilor tip broadcast către un anumit program. O cerere broadcast poate primi mai multe răspunsuri, de la fiecare nod care conține programul respectiv; răspunsurile conțin și numerele porturilor asociate.

SUN RPC furnizează și modalități de autentificare a clienților: formă implicită (fără autentificare), autentificare UNIX și autentificare DFS. Autentificarea UNIX a clienților provoacă transmiterea în fiecare apel RPC a următoarelor câmpuri: o marca de timp, numele nodului, identificatorul de utilizator al clientului, identificatorul de grup efectiv al clientului și lista identificatorilor tuturor grupurilor la care aparține clientul. Este sarcina serverului să examineze aceste câmpuri pentru a determina dacă poate răspunde unui client sau nu.

3.4.4. Probleme de concurență

Mecanismul RPC este util în programarea distribuită, în medii eterogene chiar, dar nu incurajează calculul paralel, deoarece atunci când un client apelează o procedură de la distanță, el se blochează pe parcursul execuției procedurii. Aceasta poate cauza performanțe slabe dacă clienții și serverele nu utilizează mecanismul firelor; de exemplu un server poate apela, la rândul lui, pentru implementarea unui serviciu, alte proceduri din alte servere și poate deservi alți clienți în timp ce așteaptă răspunsurile de la aceste servere. Folosirea firelor este o soluție adecvată în multe dintre problemele care apar datorită blocării într-un apel de procedură la distanță. Un server poate executa concurrent mai multe apeluri ale clienților, fiecare apel într-un fir. Analog, un client poate crea mai multe fire, pentru a fi executate concurrent; dacă unul dintre acestea se blochează într-un apel de procedură la distanță, celelalte își pot continua execuția.

3.4.5 Semantica RPC în prezența defectelor

Spre deosebire de procedurile locale, pentru care se știe întotdeauna că, la returnarea controlului, după apelul lor, au fost executate exact o dată, în cazul procedurilor executate la distanță, dacă nu se returnează controlul după un anumit interval de timp este greu de precizat de câte ori au fost executate. Există 5 situații care trebuie considerate în acest sens:

- (1) Clientul nu poate localiza serverul (un exemplu este un client care încearcă să contacteze o versiune învechită a unui server, ce nu mai există); eroarea poate fi raportată prin generarea unei excepții (cu dezavantajul că nu toate limbajele dispun de mecanismul excepțiilor).
- (2) Mesajul de apel al procedurii se pierde. În acest caz, dacă nu se primește răspunsul (sau o confirmare) după un interval de timeout, se poate retransmite mesajul. Dacă și după un număr maxim de încercări nu se primește răspuns, se poate trage concluzia că serverul a cazut (și se poate adopta aceeași soluție ca la punctul 1).
- (3) Mesajul de răspuns la procedura se pierde. În acest caz fie că se încearcă o asemenea structură a cererilor încât să fie idempotente, fie se asignează fiecărei cereri un număr de secvență iar serverul reține rezultatele ultimelor cereri astfel încât la detectarea unei cereri duplicate, nu se mai execută procedura ci se transmite direct rezultatul. Un exemplu: în SUN RPC, dacă se folosește ca protocol de transport protocolul UDP, se retransmite automat o cerere dacă este necesar; timeout-ul se poate specifica de client (în funcțiile `clntudp_create()`, `clnt_call()`). După un număr maxim de încercări nereușite se returnează o eroare procesului client. Dacă se folosește ca protocol de transport protocolul TCP, se returnează o eroare clientului dacă conexiunea este terminată de server. Fiecărei cereri de execuție a unei proceduri a unui client i se asociază un identificator pe 32 biți care este trimis de client serverului; serverul returnează acest identificator clientului, astfel că răspunsul la o cerere nu se livrează decât după ce clientul verifică identificatorul primit. În plus, serverele UDP pot menține o zonă cache a răspunsurilor trimise, astfel încât cererile duplicate nu generează noi execuții.
- (4) Nodul serverului suferă o avarie gravă. Deoarece în această situație există posibilitatea ca (1) serverul să fi suferit avaria înainte de a executa cererea sau (2) după execuție, dar înainte de a transmite răspunsul, se disting 3 semantici posibile: (1) *semantica cel-puțin-o-dată* care garantează că procedura a fost executată cel puțin o dată, dar posibil de mai multe ori; este situația ideală pentru cereri idempotente; pentru cereri neidempotente există mai multe variante, dintre care cea mai folosită este varianta *ultima-din-mai-multe*; (2) *semantica cel-mult-o-dată* care garantează că cererea cel mult o dată dar posibil nici o dată; (3) *semantica exact-o-dată* este cerința ideală, dar este dificil de obținut.
- (5) Nodul clientului suferă o avarie gravă. În această situație, dacă clientul a cerut execuția unei proceduri la distanță, și între timp a suferit o avarie, serverul poate rămâne în execuție un timp nedefinit; astfel de procese active, fără ca rezultatele execuției lor să fie așteptate de vreun proces se numesc *orfanii*. Procesele orfane trebuie distruse deoarece, în afară de faptul că irosesc timp procesor, pot bloca accesul la resursele pe care le detin (semafoare, fișiere). Există patru tehnici de distrugere a orfanilor [Gos91][Tan89][Tan96]: (1) *tehnica exterminării*; aceasta impune ca, la orice mesaj RPC, stub-ul client să înregistreze într-un fișier jurnal (care trebuie să supraviețuiască avariei) informații despre cerere. După reinițializarea nodului clientului, jurnalul este verificat și orfanii distruși explicit. În afară de faptul că necesită resurse speciale (fișier jurnal), tehnica poate să nu funcționeze în unele cazuri particulare, precum partiționarea unei rețele (datorită unei părți care nu funcționează, de exemplu) când se poate localiza un server orfan, dar nu se poate distruge sau când un orfan a lansat în execuție alte apeluri RPC, aparând astfel orfanii nepoți, care nu mai pot fi detectați; (2) *tehnica expirării* este analogă, ca principiu, tehnicii time-out; la începerea execuției unei proceduri, serverul are alocat pentru execuție un timp maxim T. Dacă se depășește acest timp, stub-ul server cere stub-ului client o nouă cuantă T; în cazul în care clientul nu mai există, serverul se distruge. Tot ceea ce trebuie să facă un client care se reinițializează, este să aștepte un interval de timp T; astfel este sigur că orfanii sunt distruși; (3) *tehnica reîncercării* divizează timpul în perioade, numite epoci; când un client este reinițializat, el trimite (prin broadcast) către toate nodurile un mesaj de start al unei noi epoci. La primirea unui astfel de mesaj, toți orfanii

sunt distruși. Evident, și aceasta tehnică nu detectează toți orfanii dacă rețeaua este partitionată în urma unui defect hardware; totuși când vor putea returna un răspuns vor raporta o epoca diferită de cea curentă; (4) tehnica *reincercării lente* este o variantă mai puțin restrictivă a tehnicii reincercării. Astfel, când un nod primește un mesaj de început al unei noi epoci verifică dacă are servere în execuție și, în caz afirmativ, clienții lor. Serverele se distrug numai dacă clienții lor nu pot fi localizați.

3.4.6. Protocoale pentru RPC

Semantica unui apel de procedură la distanță nu trebuie să fie legată de implementarea nici unui protocol; rutinele RPC care se apelează la nivelul stub-urilor trebuie să fie disponibile indiferent de clasa protocolului de transport utilizat. Astfel, o primă clasă de protocoale pe care se pot construi apeluri ale procedurilor la distanță sunt protocoalele de nivel transport care asigură un transfer de date în flux, sigur, de exemplu TCP. Protocolul TCP este utilizat în SUN RPC, DCE RPC, Microsoft RPC. Protocoalele bazate pe conexiuni, cum este TCP, prezintă avantajul că implementarea RPC nu trebuie să trateze defectele, nu trebuie să utilizeze confirmări; protocolul orientat pe conexiune se ocupă de aceasta. Când aplicația lucrează într-o rețea de mare întindere geografică (deci nu o rețea locală) aceasta este soluția cea mai bună. Cealaltă clasă de protocoale care pot fi utilizate de implementările RPC cuprinde protocoalele nesigure, la nivel de datagramă, de exemplu UDP; și protocolul UDP poate fi utilizat în SUN RPC, DCE RPC, Microsoft RPC. Utilizarea unui protocol cunoscut cum este UDP (de fapt IP, care sta la baza lui) oferă avantajul unui protocol deja elaborat, disponibil pe mai multe arhitecturi, folosit aproape în totalitate de sistemele Unix, care asigură comunicația directă între multe rețele deja existente dar și dezavantajul că protocolul IP, nefiind proiectat ca un protocol end-user nu oferă o eficiență ridicată. Alternativa este de a folosi un protocol specializat, care să fie proiectat având ca scop principal utilizarea lui în rețele locale. Birell și Nelson ([BN84]) au propus un protocol, adoptat apoi de sistemul Sprite ([OCD88]), având drept scop implementarea protocolului în comunicații prin RPC. În cadrul acestui protocol, apelurile de proceduri sunt împartite în două grupuri: *simple* și *complicate*. Cele simple se disting de cele complicate prin faptul că argumentele lor, ca și rezultatele încap într-un singur pachet și prin frecvența lor ridicată. Un pachet trimis printr-un apel de procedură simplă trebuie să cuprindă un identificator al apelului, numărul procedurii apelate și argumentele; un pachet de răspuns cuprinde același identificator din apel și rezultatele. Pentru apeluri simple, pachetele cu rezultate sunt folosite drept confirmări ale pachetelor cu cereri. Există trei situații când se retransmite o cerere simplă, în urma detectării unui time-out (când nu s-a primit răspuns după un interval maxim): (1) cererea s-a pierdut; (2) nu s-a pierdut nici cererea, nici răspunsul, dar execuția durează mai mult; în această situație se transmite o confirmare explicită de către nivelul de transport; (3) s-a pierdut răspunsul; procedura nu se mai execută încă o dată, se trimite doar răspunsul. Apelurile complicate se utilizează în 2 cazuri: (1) execuția procedurii durează mult sau apelurile sunt rare. În acest caz există două strategii pentru retransmisii și confirmări: (a) emițătorul unui pachet este responsabil pentru retransmisia unui pachet până la primirea confirmării; după primirea confirmării, așteaptă rezultatul, eventual trimițând câte un pachet de probă, pentru a vedea că receptorul mai există; (b) receptorul transmite automat o confirmare dacă nu poate genera un pachet în intervalul de timp maxim pentru o

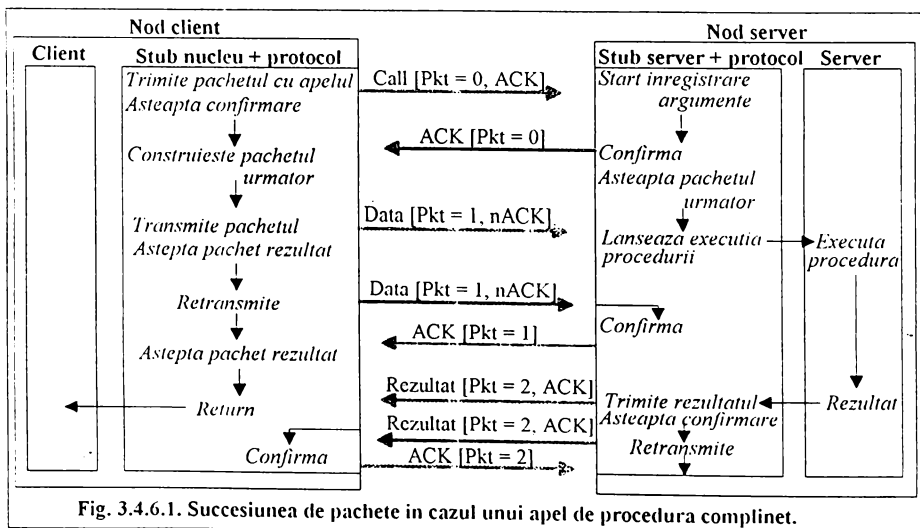


Fig. 3.4.6.1. Succesiunea de pachete in cazul unui apel de procedura completat.

etranmisie; (2) al doilea caz cînd se considera un apel complicat este cel în care argumentele sau rezultatele unei proceduri nu încap într-un singur pachet. În acest caz se confirmă numai ultimul pachet din seria argumentelor sau răspunsurilor iar pachetele dintr-o serie se numerează; în acest mod, în cazul pachetelor pierdute se evita retransmisia întregii serii de pachete.

În figura 3.4.6.1 se arată succesiunea de pachete interschimbate în cadrul unei proceduri la distanță considerată complicată deoarece execuția ei durează un interval de timp mai mare decît cel stabilit ca maxim pentru a iniția o retransmisie iar argumentele și rezultatele nu încap într-un singur pachet.

3.4.7. Extensii ale mecanismului de bază RPC

Mecanismul RPC de bază a fost extins în Sprite Network Operating Systems [OCD88][Gos91] pentru a permite comunicații între mai mulți clienți și mai multe servere; în Sprite, mecanismul RPC este implementat la nivelul de nucleu către nucleu într-un mod asemanator cu ceea ce s-a descris mai înainte. Este posibil însă, ca un nucleu client să ofere servicii la mai mult de un proces care a apelat o procedură la distanță și că un server să deservească mai mulți clienți în același timp. Ambele extensii pot fi implementate pe baza

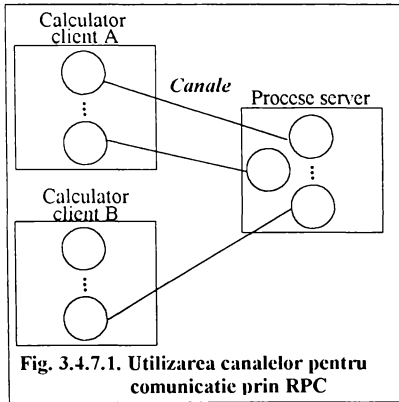


Fig. 3.4.7.1. Utilizarea canalelor pentru comunicatie prin RPC

notiunii de *canale* care conectează procesele client la procesele server; aceste canale, rezervate în număr fix pentru fiecare calculator, sunt multiplexate între clienți (fig.3.4.7.1). Un calculator server are mai multe proecese server; procesele server păstrează informații despre canalele utilizate de procesele client și nu despre procesele client. Un proces client poate utiliza un canal atunci cînd el este liber. Nucleul cu rolul de server este cel care asignează un proces server la un canal al unui client, pentru a răspunde la o serie de apeluri de proceduri; selecția procesului server care trebuie să răspundă la o cerere a unui client se face pe baza identificatorului procesului server care a fost utilizat pentru a trată ultima cerere.

O altă extensie care a fost adusă mecanismului RPC în sistemele de operare distribuite experimentale este posibilitatea de a difuza o cerere RPC la mai multe servere, pentru a obține un grad sporit de siguranta în sistemele în care nodurile se pot defecta independent. Astfel, în sistemul de operare IIPC University of Rochester [BF85] s-a introdus conceptul de apeluri de *proceduri replicate*, alături de acela de *trupa*. O *trupa* este un set de module replicate; -trupele se pot executa pe mai multe calculatoare. Cînd o trupa, reprezentînd un modul client, cere un serviciu, fiecare membru al trupeii server execută procedura necesară și fiecare membru al trupeii client primește un răspuns. Un program construit pe baza apelurilor de proceduri replicate va continua să lucreze în prezența defectelor atîta timp cît există cel puțin o replică (copie) al fiecărui modul. Un ringmaster se folosește pentru a lega (bind) nume la trupe.

Fig. 3.4.7.2 arată două trupe și un ringmaster. Fiecare trupa este implementată ca un obiect IIPC cu o

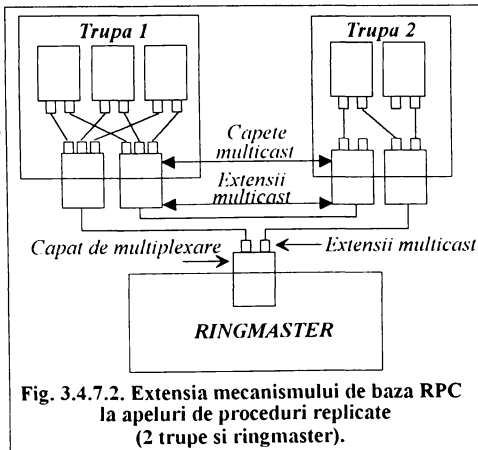


Fig. 3.4.7.2. Extensia mecanismului de baza RPC la apeluri de proceduri replicate (2 trupe si ringmaster).

interfața multicast. În interiorul trupeii, membrii ei sunt interconectați la capatul interfeței multicast. Cele două trupe sunt conectate în exterior prin intermediul extensiei interfeței multicast; fiecare trupa are o interfață similară și pentru comunicația cu ringmasterul. Acesta are o interfață multiplexată pentru a comunica cu un număr variabil de trupe; fiecare componentă a interfeței este o extensie multicast. Toate trupele au interfețe externe analoge, indiferent de cati membri replicati au. Complexitatea membrilor replicati ai unei trupe se izolează la interiorul trupeii (guvernată de un shell), în timp ce complexitatea datorată unui număr variabil de trupe și interacțiunilor lor se izolează la nivelul imediat superior. Astfel, în IIPC se permite trupelor să fie construite independent de utilizarea lor și să fie utilizate independent de constructia lor.

3.4.8 Aspecte critice în proiectarea comunicării prin RPC.

Tannenbaum a identificat ([Tan93][Tan96]) în desfășurarea unui apel de procedură la distanță etapele prezentate în fig. 3.4.8.1. Aceste etape formează *calea critică*. Unele aspecte ale acestor etape, considerate critice, pot avea o influență considerabilă asupra timpului total de execuție al unei proceduri la distanță.

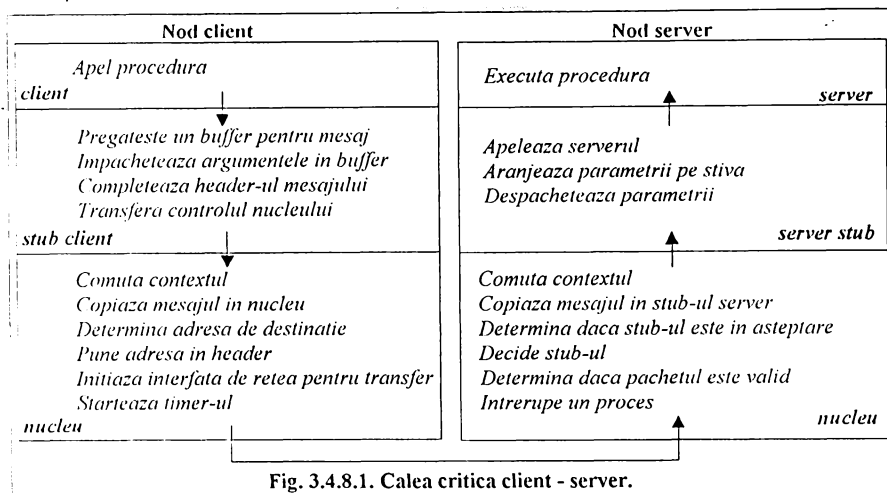


Fig. 3.4.8.1. Calea critică client - server.

Astfel, problema copierii conținutului mesajelor de transmis are o importanță majoră deoarece în funcție de structurarea software-ului și de facilitățile hardware pot fi necesare de la 2 pînă la 8 faze de copiere a mesajelor; în acest sens o importanță mare o prezintă următoarele aspecte: (1) posibilitatea ca interfața de rețea să acceseze prin DMA zona de memorie din stub-ul client care conține mesajul de transmis; (2) prezența facilității *scatter/gather* la interfața de rețea (facilitatea *scatter/gather* a fost definită în secțiunea 2.3.); aceasta este utilă deoarece sistemul de operare poate construi numai header-ul pachetelor, lasînd mesajul propriu zis în stub-ul client; (3) dacă se utilizează memorie virtuală, se poate include pagina de memorie care conține un mesaj recepționat în spațiul de memorie al stub-ului server. Alte aspecte care pot aduce o îmbunătățire a performanțelor mecanismului RPC este folosirea algoritmilor de analiză (*sweep timers*, vezi secțiunea 2.3.) și reducerea regiilor de sistem datorată gestiunii proceselor în nodul serverului.

Pentru a se exemplifica direcțiile în care se încearcă optimizarea mecanismelor RPC, se prezintă în continuare modalitățile de optimizare ale mecanismului RPC Peregrine (realizat la Rice University, Houston, Texas, [JW93]), ale cărui performanțe sunt remarcabile (timpul de trecere pentru o procedură nulă este de 573us față de 1100us la Amoeba și 1097us la Cedar) și apoi se va prezenta o comparație între câteva implementări RPC alături de performanțele lor (publicate).

Peregrine RPC a fost proiectat în vederea obținerii unor performanțe superioare în ceea ce privește comunicația interprocese, menținînd facilități legate de comunicația prin proceduri precum argumente și rezultate de tipuri arbitrare. Implementarea (și măsurătorile) a fost făcută pe baza unor stații Sun 3/60 conectate într-o rețea Ethernet, prin controlere de rețea AMD Am/7990 LANCE. Protocolul RPC se situează imediat deasupra protocolului IP și este asemănător protocolului Cedar RPC, cu excepția faptului că utilizează un protocol blast pentru mesajele multipachet. În Peregrine, este responsabilitatea clientului să convertească argumentele la reprezentarea serverului și rezultatele la reprezentarea sa proprie, eliberînd astfel serverul de această sarcină, care poate conduce în unele situații la degradarea performanțelor serverului.

În Peregrine RPC se realizează creșterea performanțelor prin urmărirea obiectivelor de mai jos: (1) Argumentele și rezultatele se transmit direct din spațiul de adrese utilizator (client sau server) evitînd astfel copierile intermediare. Se utilizează facilitatea gather DMA a controlerului de rețea (controler Ethernet); primind o listă de segmente, fiecare segment fiind specificat prin adresa de început și lungimea, controlerul poate transmite un singur pachet ale cărui date constau în combinarea segmentelor de mai sus. (2) Nu se realizează conversia datelor (argumente și rezultate) cînd clientul și serverul utilizează aceeași reprezentare (spre deosebire, de exemplu de SUN RPC care utilizează întotdeauna reprezentarea standard). (3) Pachetele se trimit utilizînd headere deja alocate și completate evitînd astfel aceste operații la fiecare apel. (4) Nu se salvează starea firului între apeluri în server; deci nu este nevoie de salvarea și restaurarea registrelor între apeluri în cadrul serverului și nici de salvarea stivei firului. (5) Argumentele sunt mapate în spațiul de adrese al serverului, deci nu sunt copiate. Se utilizează facilitatea de a remapa paginile într-un spațiu de adresa numai prin modificarea intrărilor corespondente din tabela de pagini, procedură suficientă pentru Sun 3/60 (alte

Partea întâi: Sisteme de operare distribuite

sisteme mai necesită și modificarea intrărilor corespunzătoare din TLB= translation lookaside buffer). (6) Argumentele apelurilor multipachet sunt transmise încât nu este necesară copierea lor în calca critică; ele sunt fie copiate în paralel cu transmisia în rețea, fie sunt mapate direct în spațiul serverului.

În figura 3.4.8.2 se prezintă o comparație între modalitățile utilizate pentru optimizarea performanțelor de câteva sisteme RPC și sisteme bazate pe comunicația prin mesaje dar care construiesc apeluri sistem asemănătoare apelurilor de proceduri. În cea de-a doua categorie pot fi considerate comunicațiile din Amoeba și din V Kernel.

Pentru compararea performanțelor sistemelor RPC se folosesc: timpul de trecere al unui apel de procedură nulă și productivitatea obținută în urma unei serii de apeluri de proceduri cu argumente mari.

Sistem	Caracteristici ale implementării	Dacă se folosește facilitatea gather DMA și header template	Cum se face conversia datelor	Cum se transmit argumentele în spațiul de adrese al serverului	Cum se realizează execuția procedurilor la nivelul serverului
<i>V System</i>	Nucleul se bazează pe comunicații prin mesaje	În unele implementări; se folosește un header template		Prin copiere; pentru argumente multipachet copierea se face în paralel cu transmisia	
<i>Sprite</i>	RPC folosit pentru comunicații nucleu către nucleu	Într-un mod limitat; se folosește un header template	la nivelul receptorului	Prin copiere; pentru argumente multipachet copierea se face în paralel cu transmisia	Un set de fire daemon ale nucleului tratează cererile; stiva și registrele firelor se salvează și restaurează
<i>X-Kernel</i>	Implementarea se compune dintr-un număr de protocoale de nivel jos	Fiecare protocol folosește un header template	Datele se transmit în formatul emițătorului, cu un bit indicând arhitectura lui	Argumentele mici prin copiere; pentru cele mari se evită copierea	Procedura se execută de un fir server selectat dintr-un set de fire prealocate; firele persistă în urma execuției; registrele și stiva se salvează și restaurează
<i>Firefly</i>	Implementarea s-a făcut pe un sistem multiprocesor	Se folosește transmisia prin DMA direct din bufferul utilizatorului; se folosește un header template construit în faza de legare	Se negociază reprezentarea datelor în faza de legare	Se evită copierea	Un set de fire tratează cererile
<i>Amoeba</i>	Nucleul se bazează pe comunicații prin mesaje; RPC este constituit pe baza unor primitive sincrone prin stub-uri care se pot genera și automat	În nucleu; Amoeba reconstruiește complet header-ul pachetelor fără a utiliza un template la fiecare transmisie	Conversiile sunt lăsate în seama stub-ului	Prin copiere	Se folosesc stub-uri
<i>Sun RPC</i>	Implementarea s-a făcut pe baza unor funcții de bibliotecă scrise deasupra nucleului Unix	Nu este cazul	Se folosește transmisia într-o reprezentare standard a datelor	Prin copiere	Se folosesc stub-uri

Fig. 3.4.8.2. Comparație între câteva sisteme bazate pe RPC și pe comunicații prin mesaje.

3.4.9. Exemple

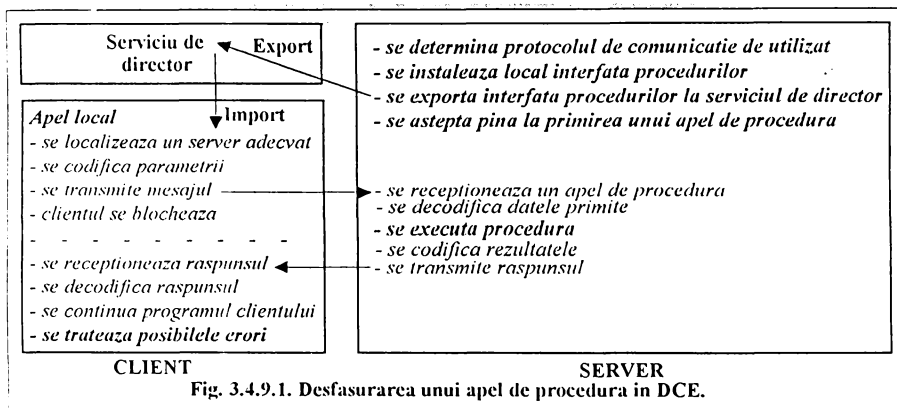
A. Modelul RPC în DCE

Transmiterea parametrilor și reprezentarea datelor

În DCE, RPC este un mecanism de bază; implementarea DCE RPC folosește un limbaj specific pentru definiția serverelor, Interface Definition Language (IDL) care este de fapt un supraset al părții declarative din limbajul C (corespunzând porțiunii de fișier header al unui program C). Funcțiile run-time RPC sunt oferite tot în limbajul C (similar altor componente DCE). IDL oferă pentru specificarea tipurilor parametrilor procedurilor facilitățile limbajului C. Diferențele între reprezentarea datelor la client și server se rezolvă prin tehnica *receiver makes right*, aceasta implică transmiterea datelor în reprezentarea emițătorului și adaptarea lor la reprezentarea receptorului la nodul destinat.

Legarea clienților la servere

Desfasurarea unui apel de procedură la distanță în DCE este reprezentată în figura 3.4.9.1. La initializarea unui server se determina ce protocol se utilizează (TCP/IP sau UDP/IP), se înregistrează interfața procedurilor oferite la un manager local de RPC-uri, se exportă informații despre server la CDS, și, în final, se așteaptă sosirea unor cereri. Pentru a apela o procedură la distanță, clientul folosește o sintaxa normală, ca pentru un apel local. Totuși, datorită stub-ului client generat din IDL, se execută o rutină internă; aceasta contactează serviciul de director CDS pentru a localiza un server adecvat (informațiile de intrare sunt numele logic al serverului și interfața dorită, iar informația de ieșire este adresa serverului). Apoi se codifică parametrii de intrare și se transmit serverului.



Concurența

Existind fire, este posibil să se implementeze servere multifir. Aceasta implică doar setarea unui parametru adecvat la initializarea serverului; apoi un ansamblu de servere concurente este alocat (static). Programatorul este cel care trebuie să țină seama, însă, de corecta sincronizare a firelor, în cazul modificării datelor partajate. În ceea ce privește clientul, firele trebuie startate explicit, pentru a realiza apeluri asincrone la mai multe servere.

Semantica RPC

Aplicațiile pot alege între diferite tipuri de semantica. Astfel, semantica implicată, cel-mult-o-dată asigură execuția unui apel chiar dacă mesajele sunt pierdute temporar (printr-o tehnică de retransmisie combinată cu detectia mesajelor duplicate). Deși caderea nodurilor nu este tolerată, pierderea mesajelor poate fi mascată în acest mod. Alte semantici care pot fi selectate furnizează garanții mai puține în prezența defectelor, dar realizează o creștere a eficienței transmisiei.

Securitatea

Siguranța comunicațiilor RPC se asigură prin serviciul de securitate. Întâi, clientul și serverul execută un protocol de autentificare în cooperare cu serverul pentru securitate. În această fază, ambii își validează reciproc identitatea pe baza unui mecanism de criptare cu cheie privată. În faza a doua, când trebuie executată o procedură, serverul verifică autorizarea unui client pe baza unei liste de control al accesului locală. În plus, datele pot fi opțional criptate, pentru a asigura o protecție completă și în cursul comunicațiilor.

Dezvoltarea aplicațiilor distribuite cu DCE RPC

Dezvoltarea aplicațiilor distribuite cu DCE RPC necesită trei etape [BGH93]: definirea interfeței IDL, implementarea serverului și implementarea clientului. În continuare, prezentarea acestor etape se va face pe baza unui mic exemplu care constă într-un client ce adresează cereri unui server de gestiune a produselor într-o întreprindere.

Prima etapă constă în definirea unei interfețe IDL care trebuie să specifice toate procedurile serverului. Interfața constă într-un header cu un număr de interfață unic (generat automat) și cu informații despre versiune și într-un corp care cuprinde definiția funcțiilor ca în limbajul C, cu specificarea tipurilor parametrilor (de intrare și/sau de ieșire). Pentru exemplul luat, definiția interfeței este cea din figura 3.4.9.2.

A doua etapă este implementarea serverului (fig.3.4.9.3). Faza de inițializare a serverului este implementată conform celor expuse anterior.

Implementarea clientului este independentă de aspectele de distribuție (fig.3.4.9.4).

```

|
|  uuid(876c2a20-100a-145c-1678-040083d67921)
|  version (1.0)
|
|  interface DateProductie { // interfata pentru date productie
|  import "globaldef.idl" // import al definitiilor generale
|  const long ProdMax=30 // numar maxim de produse
|  typedef [string] char *string // tipul string
|  typedef struct {
|      string NumeProdus // nume produs
|      string DetaliiProdus // informatii despre produs
|      Plan PlanDeProductie // tip definit in globaldef.idl
|  } DescriereProdus // tip de date pentru descriere produs
|  long CerereProdus // procedura ce poate fi apelata de la distanta
|  [in] string NumeProdus[ProdMax] // parametru intrare
|  [in] DescriereProdus *pd[ProdMax] // parametru de iesire
|  [out] long *stare; // parametru de iesire
|  }
    
```

Fig.3.4.9.2 Definirea interfeței

```

#include "dateproductie.h" // generat de compilatorul IDL
#define DirServer "../ServerProductie" // numele directorului serverului
#define MaxApelConc 5 // numar maxim de apeluri concurente
main() {
    unsigned stare; // stare returnata
    rpc_binding_vector_t *bVec; // vector de handle-rc
    /*** Alte initializari ***/
    /*** Precia vectorul handle ***/
    rpc_server_inq_bindings(&bVec,&stare);

    /*** Se inregistreaza local interfata ***/

    rpc_cp_register (DateProductie_v1_0_s_ifspec, bVec,NULL,NULL,&stare);

    /*** Se exporta interfata la CDS ***/
    rpc_ns_binding_export (rpc_c_ns_syntax_default, DirServer,
        DateProductie_v1_0_s_ifspec, bVec,
        NULL, &stare);

    /*** Accepta cereri concurente ***/
    rpc_server_listen (MaxApelConc, &stare);
}
    
```

Fig.3.4.93 Serverul

```
#include "dateproductie.h"
main() {
string produse[ProdMax];□           // nume produse
DescriereProdus *dp[ProdMax];      // descrieri cerute
long stare,rq;
PrciaDate (produse);                // functie specifica aplicatiei
rq=CerereProdus (produse,dp,&stare);// RPC
// ... Se verifica starea si se trateaza erorile
```

Fig. 3.4.9.4 Clientul

B. Microsoft RPC

Microsoft RPC se bazează pe DCE RPC; difera totusi de aceasta prin urmatoarele:

- (1) În biblioteca de funcții runtime, denumirile sunt modificate în concordanta cu conventiile Microsoft; de exemplu funcția `rpc_string_binding compose()` devine `RPCStringBindingCompose()`;
- (2) funcțiile raportează informații de stare prin valoarea lor de retur; în OSF se folosește o variabila de ieșire;
- (3) Tipurile de date primese nume, utilizandu-se `typedef` pentru tipurile structurate;
- (4) Unele funcții primese argumente în plus, specifice sistemului NT, precum descriptori de securitate;
- (5) Sunt adăugate funcții și macro-uri noi pentru tratarea excepțiilor.

Microsoft RPC cuprinde componentele tradiționale mecanismelor RPC: (1) compilatorul IDL (versiunea Microsoft se numește MIDL). Compilatorul primește ca intrare două fișiere (.IDL și optional .ACF) pentru a genera codurile stub pentru server și client; acestea sunt apoi compilate și linkeditate cu aplicațiile client și server (fig. 3.4.9.5); (2) bibliotecile de funcții runtime, utilizate de aplicațiile server pentru a-și anunța serviciile și de aplicațiile client, pentru a localiza și contacta serverele. Bibliotecile Microsoft necesare, per total, sunt: `rper4.lib` (furnizează funcțiile runtime), `rpens4.lib` (furnizează funcțiile pentru serviciile de nume), `libent.lib` (furnizează funcțiile C standard) și `Kernel32.lib` (furnizează funcțiile pentru fire).

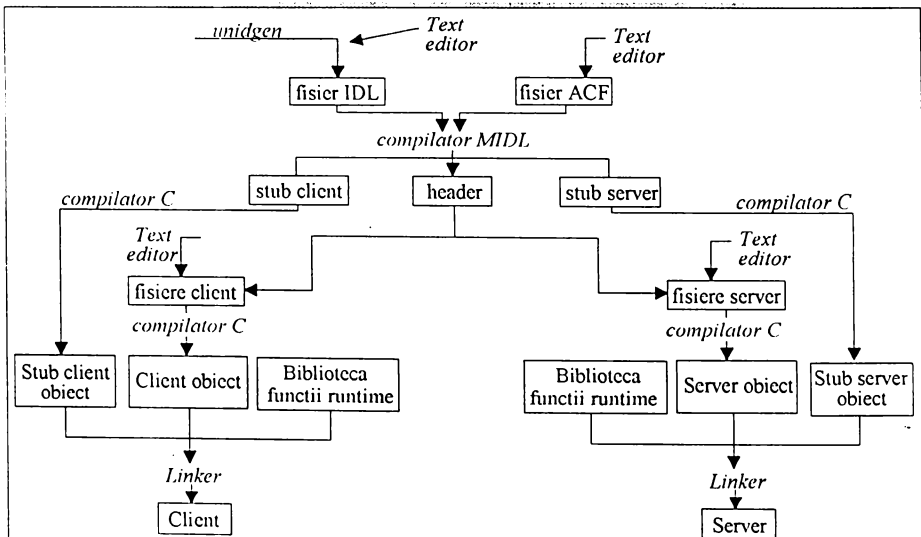
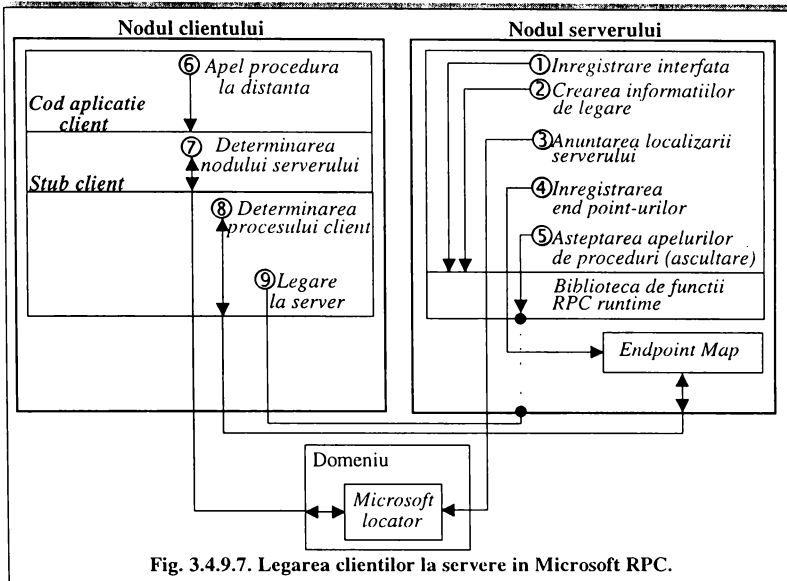
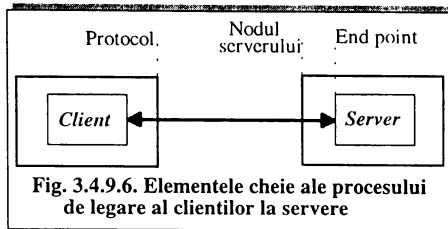


Fig. 3.4.9.5. Dezvoltarea aplicațiilor client-server pe baza mecanismului RPC in Microsoft RPC.

Legarea clienților de servere

Fig. 3.4.9.6 arată informațiile cheie în procesul de legare al clienților la servere; acesta include următoarele operații (fig. 3.4.9.7):



A. La initializarea unui server:

- (1) Înregistrarea interfeței serverului, în scopul realizării legăturii dintre handle-ul interfeței, produs de compilatorul MIDL și identificatorul sau unic UUID (UUID este o prescurtare de la Universal Unique Identifier) de către rutinele RPC runtime; mai târziu clienții vor putea determina dacă sunt compatibili cu serverul.
- (2) Crearea unor informații pentru legare; pentru aceasta trebuie stabilit unul sau mai multe protocoale de comunicație și eventual end-point-ul la care va asculta serverul sosirea apelurilor de proceduri. Pentru fiecare protocol posibil, end-point-ul are un format specific. Astfel, pentru *ncacn_np* (Named Pipes), end-point-ul este un nume de canal în formatul `\PIPE\<nume>canal`, pentru *ncacn_cp_tcp* (conexiune PCP/IP) este portul TCP (string) iar pentru *ncacn_spx* (conexiune SPX) este un soclu cu valoare între 1 și 65536 (string). Deoarece Windows 95 nu are implementată partea de server pentru Named pipes, protocolul *ncacn_np* nu poate fi folosit în rețea în Windows 95.
- (3) Plasarea informațiilor pentru legare într-o baza de date gestionată printr-un serviciu denumit Name Service, pentru ca clienții să aibă acces la ele (Microsoft denumește serviciul de nume *locator*). Biblioteca de funcții runtime conține un set de funcții generale (Name Service Independent Routines), care comunică cu locatorul pentru a plasa informații în baza de date; datorită nivelului de abstractizare al acestor funcții (cel mult independente de un serviciu de nume particular), Microsoft RPC poate fi configurat astfel încât să utilizeze, de exemplu, când comunică cu sisteme DCE, serviciul DCE Call Directory Service. Microsoft RPC nu impune aplicațiilor distribuite folosirea serviciului de nume, deși pentru aplicațiile complexe, această variantă este recomandată. Există și alte două modalități de a lega clienții la servere: (a) legarea directă, clientul cunoaște localizarea serverului, în codul clientului și serverului; (b) utilizarea unei metode specifice aplicației prin care serverul își anunță identitatea clienților.

(4) Plasarea informațiilor despre end-point-uri într-o bază de date specială a nodului serverului, numită *harta end-point-ului* (End-Point Map); un server special, End-Point Mapper este responsabil pentru tratarea unei cereri RPC cu un handler de legare parțial determinat. End-Point Mapper-ul încearcă să localizeze un server local, care să ofere o interfață identificată prin același UUID; dacă reușește acest lucru, atunci completează structura de date asociată handle-ului de legare. Din acest moment, handle-ul fiind construit, următoarele cereri nu mai contactează End Mapper-ul.

B. La apelul unei proceduri de către un client (Etapile 6-9 corespund operațiilor executate de client pentru legarea la server)

(6) Clientul execută un apel de procedură la distanță;

(7) Prin intermediul locator-ului, stub-ul client determină nodul serverului;

(8) Prin intermediul End-point Mapper-ului, clientul determină identitatea serverului.

Dezvoltarea aplicațiilor cu Microsoft RPC

Dezvoltarea aplicațiilor cu Microsoft RPC se face urmărind aceleași etape prezentate în secțiunea anterioară. Aceste etape pot fi ilustrate cu un exemplu simplu. În fig. 3.4.9.8 se arată configurația fișierului. IDL al exemplului prezentat în fig. 3.4.9.9 prezintă codul funcției main a serverului. Codul clientului este independent de aspectele de distribuție.

```

/* fisierul exemplu.idl */
[
  uuid(47883460-bafb-11cf-a9do-0000c012d29d), /* atribut uuid */
  version(1.0),
  pointer_default (ref)                       /* tratarea implicita pentru pointeri este p
]
interface exemplu                            /* interfata pentru exemplu */
{
  const unsigned short size=20,              /* dimensiunea maxima a vectorilor */
  typedef long array[size],                  /* tip de date */
  void sum (
    [in] array a1, /* pointer intrare */
    [in] array a2, /* parametru intrare */
    [out] array a /* parametru iesire */
  );
}

```

Fig. 3.4.9.8 Continutul fișierului exemplu.idl

```

/* fisierul server.c */
#include <stdio.h>
#include "exemplu.h" /* generat de compilatorul MIDL */
main() {
  unsigned long stare; /* cod de eroare */
  rpc_binding_vector_t *bvec /* set de handle pentru legare */
  unsigned char * name; /* nume pentru serviciu */
  /* inregistrarea interfetei */
  stare=RpcServerRegisterIf (exemplu_v1_o_s_ifspec, NULL, NULL);
  /*** cod pentru verificare stare ***/
  /* crearea informatiilor de legare */
  stare=RpcServerUseAllProtseqs (RPC_C_PROTSEQ_MAX_REQS_DEFAULT, NULL);
  /*** cod pentru verificare stare ***/
  /* obtinerea vectorului de handle-re; se incheie astfel etapa 2 din fig. 3.4.9.7 */
  stare=RpcServerInqBindings (&bvec);
  /*** cod pentru verificare stare ***/
  /* se face cunoscut in baza de date nodul serverului */
  name= (unsigned char *)getenv("EXAMPLE SERVER ENTRY");
  STARE=RpcNsBindingExport(RPC_C_NS_SYNTAX_DEFAULT, name, exemplu_v1_o_s_ifspec, bvec
  NULL);
  /* cod pentru verificare stare */
}

```

Fig. 3.4.9.9. Codul funcției main() a serverului

3.5. Comunicații de grup

O caracteristică a modelului RPC prezentat în secțiunea 3.4 este aceea că implică în comunicare două părți: clientul și serverul. Există însă aplicații în care este necesară implicarea mai multor procese, pentru care modelul RPC nu mai este adecvat. Ca exemplu de astfel de aplicații pot fi enumerate: serverele replicate de fișiere care asigură un serviciu de gestiune fișiere tolerant la defecte, sistemul de marketing electronic, protocoale two-phase, prelucrări distribuite, aplicații tip conferință. În această secțiune vor fi prezentate mecanisme de comunicații de grup, în cadrul cărora un mesaj poate fi trimis mai multor receptori prin aceeași operație.

3.5.1. Probleme specifice de comunicații de grup

Un *grup multicast* este un set de procese care reprezintă destinația aceluiași mesaj ([GS91]). Mesajele pot fi generate de către una sau mai multe surse, iar procesele destinație pot fi rezidente pe unul sau mai multe noduri, nu neapărat distincte. *Comunicațiile de grup (sau multicast)* se referă la comunicații de tipul unul-la-mai-multi, în contrast cu comunicația punct-la-punct prezentată până acum. Grupurile sunt entități dinamice; pot fi create și distruse grupuri, iar componența grupurilor poate fi modificată. Un proces poate fi membru la mai multe grupuri în același timp. În cadrul unui grup, procesele pot fi egale sau poate fi creată o ierarhie între procese; de exemplu, un proces poate îndeplini rolul de coordonator, în timp ce celelalte pot îndeplini aceeași sarcină de lucru. Primul model prezintă o simetrie; dacă unul dintre procese dispăre, grupul devine mai mic, dar continuă să lucreze. În cel de-al doilea, dispariția coordonatorului pune în pericol întregul grup. În schimb, în primul model, luarea unei decizii este o operație complicată, care implică un proces de votare.

Pentru *gestiunea grupurilor* (crearea, stergerea grupurilor, modificarea componenței grupurilor) există două soluții:

(1) o soluție *centralizată*, în cadrul căreia un server de grupuri primește și tratează toate cererile care se referă la componența grupurilor; serverul de grupuri trebuie să mențină o bază de date cu grupurile create și componența lor. Soluția, deși simplă și ușor de implementat, prezintă dezavantajul că, în cazul dispariției serverului de grupuri, gestiunea grupurilor este compromisă, iar grupurile trebuie reconstruite.

(2) o soluție *distribuită*, în cadrul căreia un proces transmite un mesaj tuturor proceselor din toate grupurile la intrarea sa într-un grup, ca și la parasirea unui grup. Grupurile pot fi *inchise*, caz în care numai membrii unui grup pot transmite mesaje grupurilor sau *deschise*, caz în care orice proces din sistem poate transmite mesaje la orice grup. Grupurile închise sunt utile în aplicații de prelucrare paralelă (de exemplu, procesele care implementează un joc de sah pot construi un grup închis: au același scop, deși nu interacționează cu exteriorul), iar grupurile deschise sunt caracteristice serverelor replicate, când un proces care nu este membru al grupului serverului (un client) trebuie să poată să le trimită un mesaj.

O problemă specifică comunicațiilor de grup este *adresarea grupurilor*; dacă sistemul hardware (rețeaua) suportă adresă multicast, adresa fiecărui grup poate fi asociată unei adrese multicast. Dacă nu există această posibilitate, atunci se pot transmite mesajele prin broadcast, rămânând ca fiecare nucleu al unui nod să ignore mesajul, dacă nici un proces din acel nod nu face parte din grupul destinație. Dacă nu este disponibilă nici transmisia broadcast, atunci nucleul emițătorului va trimite mesajul, individual, la toate procese care alcatuiesc grupul destinație. O altă tehnică de adresare a membrilor grupului destinație difera total de toate cele mai sus, prin faptul că nu oferă transparență. Dacă în toate tehnicile de mai sus, procesele utilizator care apelează primitivele de transmisie, nu trebuie să cunoască componența grupului destinație, în această tehnică, apelantul trebuie să specifice o listă cu toate adresele componentelor grupului destinație (de exemplu, adresa IP); în acest caz, un parametru al primitivei send va fi un pointer la această listă. O variantă a celor două tehnici prezentate mai sus este ca fiecare mesaj să conțină un predicat (expresie booleană), care să poată fi evaluat la destinație; dacă valoarea obținută este TRUE, mesajul este acceptat, altfel ignorat.

3.5.2. Primitive și protocoale pentru comunicații de grup

Primitivele pentru comunicațiile de grup pot fi comune cu cele pentru comunicațiile punct la punct. Totuși, dacă mecanismul de comunicație punct la punct este RPC, acest lucru nu este posibil. Ca și pentru comunicațiile punct la punct, primitivele pentru comunicațiile de grup pot fi cu și fără blocare, cu și fără bufferare, sigure și nesigure.

În general, comunicațiile de grup trebuie să furnizeze primitive pentru: (1) crearea și distrugerea grupurilor; (2) inserarea unui nou membru într-un grup, eliminarea unui membru dintr-un grup; (3) atribuirea de pseudonume grupurilor, stergerea de pseudonume, căutarea unui grup după pseudonume; (4) transmiterea și recepția de mesaje către/de la procese aparținând unui grup multicast. Trebuie menționat că alăturarea sau parasirea la/unui grup trebuie făcută sincron cu transmisia de mesaje. Aceasta înseamnă că după eliminarea

unui membru al unui grup, nu mai pot fi primite mesaje de la acest proces; analog, la alăturarea unui proces unui grup, el poate primi imediat mesaje.

O caracteristică a protocoalelor pentru comunicațiile de grup este *atomicitatea*: mesajele sunt trimise unui grup, numai dacă ele sunt recepționate corect de către toți membrii grupului; situații de genul în care numai unii membrii au recepționat un mesaj, iar ceilalți nu l-au recepționat (de exemplu datorită unei erori *overrun*) nu sunt permise. Următoarele două metode pot fi folosite pentru implementarea atomicității:

(1) la primirea mesajului, fiecare membru al grupului destinație va trimite o confirmare emițătorului; dacă, după un interval de timeout, emițătorul nu a primit toate confirmările, retransmite mesajul acolo unde este cazul. Metoda nu păstrează însă atomicitatea în cazul unui incident al nodului emițător;

(2) o metoda care păstrează atomicitatea în orice situație a fost propusă de Joseph și Birman [BJ87]: emițătorul trimite mesajul, și, după intervalul de timeout, îl retransmite acolo unde este necesar. Când un proces primește mesajul, verifică dacă este un mesaj nou: în caz afirmativ, îl transmite în continuare tuturor membrilor grupului (cu retransmisii, dacă este cazul), altfel ignora mesajul. Cu aceasta metoda, indiferent câte noduri suferă incidente sau câte mesaje se pierd, procesele care rămân vor primi mesajul.

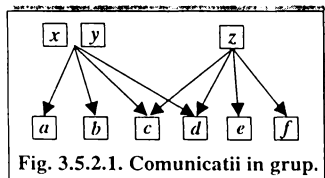
Pentru multe aplicații, protocolul multicast trebuie să furnizeze unele garanții relativ la ordinea de livrare a mesajelor la destinație. Aceste garanții, prezentate în ordine crescândă a cerințelor sunt [GS91][Tan93][Tan96]:

(1) Ordonarea mesajelor în cadrul unei singure surse: dacă mesajele m_1 și m_2 sunt trimise de aceeași sursă la același grup multicast, atunci toate procesele grupului destinație trebuie să primească mesajele în aceeași ordine.

(2) Ordonarea mesajelor în cazul mai multor surse: dacă mesajele m_1 și m_2 sunt adresate către același grup multicast, atunci toate procesele grupului destinație trebuie să primească mesajele în aceeași ordine (chiar dacă mesajele provin de la surse diferite).

(3) Ordonarea mesajelor în cazul unor grupuri nedisjuncte: analog cu (2), chiar dacă mesajele provin de la surse diferite și chiar dacă sunt adresate unor grupuri distincte, dar care au procese comune.

Cerința (1) (rezolvată uneori de protocolul de comunicație) poate fi îndeplinită dacă se ordonează (eventual numerotează) mesajele emițătorului și nu se permite recepția mesajului decât în ordine crescătoare. Cerințele (2) și (3) sunt mai greu de îndeplinit. O soluție este de a atribui o marcă de timp fiecărui mesaj și de a livra mesajele în ordinea mărcilor de timp. Fie, de exemplu grupurile (fig. fig. 3.5.2.1): $P = \{a, b, c, d\}$ și $Q = \{c, d, e, f\}$. Presupunem că nodul X trimite un mesaj m_1 cu marca de timp T_1 către grupul P. Când nodul C, de exemplu, recepționează mesajul m_1 , nu îl poate livra imediat procesului destinație, deoarece nu are certitudinea că mesajul are cea mai mică marcă de timp față de alte



mesaje nelivrate încă. Pentru îndeplinirea cerinței (2), nodul C trebuie să compare marca de timp a mesajului cu marcele tuturor mesajelor trimise către P (în cazul de mai sus, ale mesajelor de la Y). Pentru îndeplinirea cerinței (3), nodul C trebuie, în plus, să compare marca de timp și cu cele ale mesajelor trimise către Q (deoarece nodul C face parte din două grupuri, P și Q). Dacă potențialele surse de mesaje pentru C nu se cunosc, comparațiile se fac relativ la mesajele tuturor nodurilor din sistem.

Birman și Joseph [BJ87] au propus altă soluție similară protocolului two-phase commit (vezi capitolul 7). Fiecare nod menține cîte o coadă cu priorități pentru fiecare proces. Emițătorul trimite mesajul către destinațiile multicast, care îi atribuie un număr de prioritate propriu, cel mai mare pentru acel proces; mesajul este marcat *nelivrabil* și memorat în coadă. Fiecare receptor returnează emițătorului numărul de prioritate; emițătorul alege cel mai mare număr de prioritate și îl trimite înapoi la receptorii, care înlocuiesc vechiul număr de prioritate cu acesta și marchează mesajul *livrabil*. Fiecare receptor va ordona în acest moment coada sa; cînd un mesaj din varful cozii devine *livrabil* el este transmis. Aceasta soluție prezintă avantajul că asigură cerințele (1), (2) și (3) fără a necesita cunoașterea potențialelor surse de mesaje.

O soluție centralizată este următoarea [CM84]: toate sursele transmit către un nod central, care atribuie numere de secvență mesajelor și le transmite destinației. Nodul central poate fi identificat printr-un token care circula prin sistem (din punct de vedere al transmisiei mesajelor, faptul că nodul central își modifică în timp localizarea nu are importanță).

3.5.3. Exemple: MCL, ISIS, V-system, Amoeba

Se prezintă în cadrul acestei secțiuni trei exemple de comunicații de grup, care diferă complet, atît din punct de vedere al concepției cît și ca implementare. MCL și ISIS se materializează sub forma unor toolkit-uri pentru comunicații de grup, pe cînd Amoeba implementează comunicațiile de grup la nivelul nucleului.

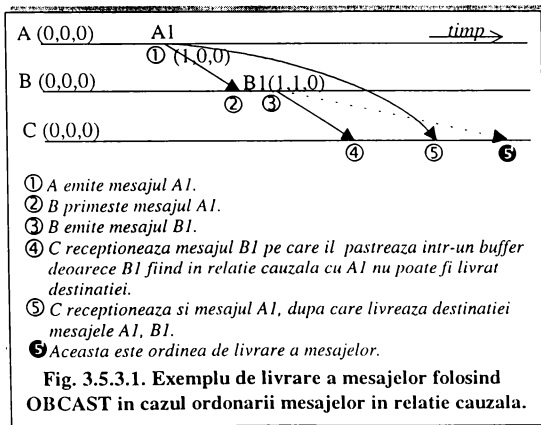
A. MCL

Un toolkit pentru comunicații de grup dezvoltat pentru sisteme Unix este MCL; el cuprinde un set de 8 primitive de comunicații. Sistemul MCL este implementat prin intermediul a două procese daemon care se execută în fiecare nod: un proces daemon se ocupa cu transmitia mesajelor la alte noduri, iar celalalt numit daemon de distribuție cu distribuția mesajelor recepționate la procesele destinație. În MCL, pentru a transmite mesaje multicast se folosește transmitia broadcast.

B. ISIS

Un alt toolkit, cunoscut și mult referit în literatura de comunicații de grup, este sistemul ISIS (Universitatea Lornell, Ithaca, New York [BJ87][BJ90]); acesta constă într-un set de programe care se pot execută în sistemul de operare Unix și în alte câteva sisteme de operare cunoscute. În ISIS sunt definite următoarele trei primitive pentru multicast [BJ90]: (1) Primitiva *GBCAST* (*group broadcast*), cea mai restrictivă, impune respectarea totală a condițiilor de ordonare a mesajelor prezentate la 3.5.2. Primitiva se folosește pentru gestiunea membrilor grupurilor, inclusiv în cazul unui membru care dispare în urma unei avarii grave a nodului sau. Argumentele primitivei sunt mesajul și identificatorul grupului de procese, care este tradus în setul de adrese ale membrilor grupului, pe baza unei liste locale cu componența grupurilor, ce se menține în fiecare nod. Deoarece execuția primitivelor *GBCAST* păstrează ordinea mesajelor, rezultă că, deși listele grupurilor nu sunt actualizate simultan, toți membrii grupurilor vor percepe aceeași secvență de modificare a grupurilor. (2) Primitiva *BCAST* este mai puțin restrictivă decât cea precedentă, fiind folosită pentru transmitia mesajelor obisnuite. Ea se utilizează în situațiile în care nu este necesară o ordonare totală a mesajelor multicast, în sensul că dacă două mesaje sunt primite într-o ordine la un nod, ele vor fi recepționate în aceeași ordine la toate destinațiile, dar nu neaparat în ordinea în care au fost emise. De exemplu, dacă un grup de procese se utilizează pentru a gestiona o structură tip coada, replicată, atunci *BCAST* se poate utiliza pentru a transmite operațiile asupra cozii; deoarece ordinea operațiilor va fi aceeași la toate nodurile, copiile cozii vor ramane consistente. Primitiva are sintaxa *BCAST*(msg, label, dest); msg și dest reprezintă mesajul, respectiv destinația iar label este un sir de caractere care specifică comportarea primitivei: numai primitivle *BCAST* care au aceeași valoare pentru label păstrăza aceeași ordine de livrare a mesajelor la noduri. Pentru implementarea ambelor primitive *GBCAST* și *BCAST*, ISIS folosește un protocol asemanator protocolului two-phase-commit, descris în secțiunea 3.5.2, completat pentru *GBCAST* cu un mecanism care elimina, mesajele provenind de la un proces membru disparut, înainte de a se livra un *GBCAST* care să anunte caderea procesului. (3) Primitiva *OBCAST* (ordered broadcast) este cea mai puțin restrictivă. Avind sintaxa *OBCAST*(msg, dest) primitivele *OBCAST* asigură transmitia atomică a mesajelor și ordonarea lor într-un nod numai în cazul unei relații cauzale între două primitive *OBCAST*. Relația cauzală este referita aici în sensul

Lamport (vezi secțiunea 4.1): două evenimente sunt în relație cauzală dacă comportarea celui de-al doilea poate fi influențată de primul; astfel, dacă nodul A trimite un mesaj la B și apoi B analizând mesajul primit trimite un alt mesaj la C, cele două operații de emisie sunt în relație cauzală; în caz contrar, de exemplu dacă A emite un mesaj la C și, aproximativ în același timp, C emite un mesaj la D, evenimentele sunt concurente. Protocolul pentru implementarea primitivei *OBCAST* se poate descrie astfel: fiecărui proces membru al unui grup i se asociază un vector având atatea componente cți membri sunt în acel grup; componenta i a vectorului este numărul asociat ultimului mesaj primit de membrul i (initial vectorul este setat la 0). Cind un proces trimite un mesaj, se incrementează componenta corespunzătoare



lui și se trimite și vectorul în cadrul mesajului. Fie V_i aceasta componentă și L_i componenta corespunzătoare a vectorului unui receptor. Algoritmul decide cind se livrează mesajul procesului sau se bufferează mesajul pentru a fi livrat mai tarziu. Figura 3.5.3.1 conține un exemplu de ordonare a mesajelor în cazul a doua operații de emisie de mesaje pentru un grup format din 3 procese. Se observă din figură că, dacă presupunem că sursa unui mesaj este j , condițiile de livrare a mesajului sunt: (1) $V_j \leq L_i + 1$; aceasta impune deci ca mesajul primit sa fie într-adevar urmatorul în secvența de mesaje primite de la sursa j , altfel spus, sa nu fi fost un mesaj pierdut; (2) $V_i < L_i$, pentru orice $i < j$; aceasta impune că emițătorul sa nu fi primit un mesaj pe care nodul receptor nu l-a primit.

În figura 3.5.3.5 se prezintă un exemplu de interacțiune între un client și un grup de procese în care se folosesc toate tipurile de primitive.

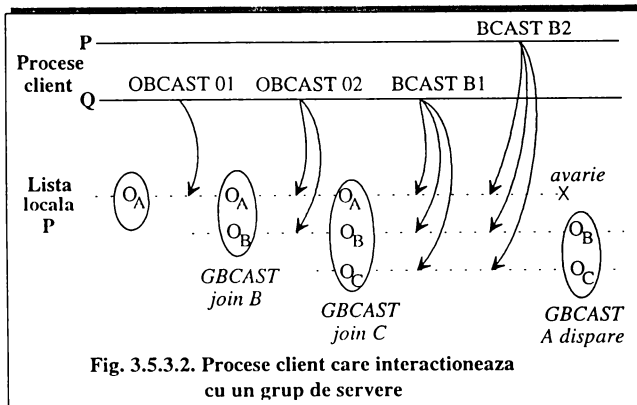


Fig. 3.5.3.2. Procese client care interacționează cu un grup de servere

C.Amoeba

În sistemul de operare Amoeba, alături de forme de bază de comunicație RPC, sunt prevăzute și comunicații de grup [Tan93][Tan96]. În Amoeba grupurile sunt închise; singura formă de acces a unui client la serviciul furnizat de un grup este ca clientul să execute un apel de procedură la unul din serverele grupului; dacă este necesar, membrii grupului pot comunica, în continuare între ei, folosind primitive pentru comunicații în grup, pentru a stabili cine execută serviciul cerut. Primitivele pentru comunicații în grup în Amoeba sunt prezentate în fig.3.5.3.3. Se asigură o transmisie sigură, atomică a mesajelor către membrii unui grup, chiar și în situații de mesaje pierdute sau depășirea capacității bufferelor. Amoeba implementează pentru comunicațiile în grup un protocol care utilizează metoda centralizată prezentată la sfârșitul secțiunii 3.5.2.

Primitiva	Semnificație
<i>Create Group</i>	Crează un grup nou și setează anumiți parametri, ca de exemplu, numărul maxim de membri disparuți ai unui grup pentru care grupul continuă să lucreze
<i>Join Group</i>	Permite introducerea unor membri noi într-un grup
<i>Leave Group</i>	Permite unui membru părăsirea unui grup
<i>Send To Group</i>	Transmite singur un mesaj membrilor unui grup
<i>Receive From Group</i>	Permite recepționarea unui mesaj de la un grup; blocare, dacă nu este disponibil nici un mesaj
<i>Reset Group</i>	Se utilizează la revenirea după o avarie gravă și specifică numărul minim de membri necesari pentru a forma un nou grup; dacă nucleul reușește să contacteze numărul respectiv de membri grupul se reconstruiește

Fig. 3.5.3.3. Primitive pentru comunicația de grup în Amoeba.

Capitolul 4.

SINCRONIZAREA ÎN SISTEME DE OPERARE DISTRIBUITE

Termenul *sincronizare* este folosit în sisteme distribuite pentru a referi trei probleme distincte, dar înrudite: (1) sincronizare între emițătorul și receptorul unui mesaj; (2) specificarea și controlul acțiunilor unor procese care cooperează în cadrul unei aplicații distribuite; (3) serializarea accesului concurent la obiectele partajate între mai multe procese; această problemă este cunoscută și în bazele de date sub denumirea de *controlul concurenței*. În capitolul de față, se vor trata aspecte privitoare numai la problemele de tipul (2); problema (1) a fost tratată în capitolul precedent, prezentându-se în secțiunea 3.2.2. patru moduri posibile de sincronizare a unui receptor cu un emițător, (în cadrul operației *send*), iar aspecte referitoare la controlul concurenței vor fi prezentate în contextul unui sistem de obiecte distribuite în cap. 7.

O concluzie a acestui capitol va fi aceea că sincronizarea în sisteme distribuite implică întotdeauna și comunicare; *sincronizarea nu se poate realiza fără comunicare*.

O prima definiție pentru sincronizare, având în vedere cele de mai sus, este următoarea:

Definiție:

Sincronizarea este o ordonare temporală a unui set de evenimente produse în timp de două sau mai multe procese concurente.

Datorită distribuirii resurselor, întârzierilor de transmisie, lipsa unei memorii comune și deci lipsa unor informații globale, algoritmi și metodele pentru sincronizare din sistemele centralizate nu pot fi utilizate în sistemele distribuite.

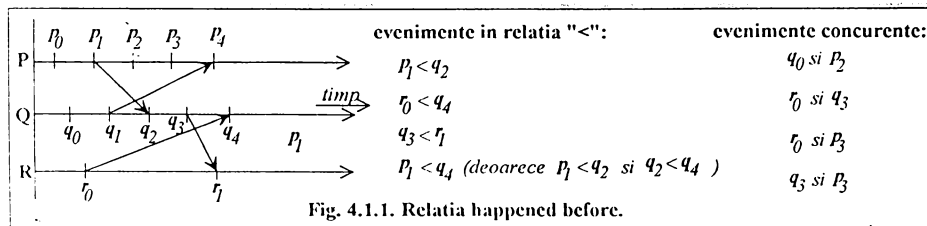
Mecanismele de sincronizare pentru sistemele distribuite vor trebui să determine, printr-o metodă oarecare, o ordonare totală sau parțială a unui set de evenimente. Capitolul de față va prezenta mai întâi aspecte ale ordonării evenimentelor într-un sistem distribuit și apoi, pe baza acestora, două clase mari de mecanisme de sincronizare: *centralizate* și *distribuite*. Unele dintre aceste mecanisme se bazează pe conceptul de *excludere reciprocă* care presupune o ordonare a unor secvențe de cod executate de procesele concurente; aceste mecanisme pot utiliza diferite tehnici precum: *un proces central coordonator* sau un *token circulant*. Alte mecanisme se bazează pe înregistrarea evenimentelor semnificative în cursul unor prelucrări asincrone; acestea utilizează obiecte precum *contori de evenimente* și *secvențiatori*, sau *semnale distribuite*.

4. 1. Ordonarea evenimentelor în sisteme distribuite

În sistemele centralizate este întotdeauna posibil să se determine ordinea în care au loc două evenimente deoarece există o memorie și un ceas comun tuturor proceselor. În sistemele distribuite, neexistând nici memorie comună și, în general, nici un ceas global, comun, este uneori imposibil să se distingă că un eveniment a avut loc înaintea altuia. O soluție de a ordona evenimentele, în sisteme distribuite, fără a utiliza ceasuri fizice, se bazează pe relația *happend-before* (*s-a întâmplat înainte*). Aceasta relație este o relație de ordine parțială pentru evenimentele sistemului distribuit, dedusă din legea cauzalității ([BBK91][Tan96][Gos91]).

Deoarece execuția unui proces se face secvențial, toate evenimentele aceluși proces vor fi total ordonate. De asemenea, conform legii cauzalității, un mesaj poate fi recepționat numai după ce a fost emis. Prin urmare, relația *happend-before*, notată cu $<$, poate fi definită peste un set de evenimente astfel (presupunând că emisia și recepția unui mesaj constituie un eveniment): (1) Dacă A și B sunt evenimente în același proces și A a fost executat înaintea lui B, atunci $A < B$; (2) Dacă A este un eveniment de emisie al unui mesaj într-un proces și B este un eveniment de recepție al aceluși mesaj de către un alt proces, atunci $A < B$; (3) Dacă $A < B$ și $B < C$ atunci $A < C$;

Deoarece un eveniment nu se poate precede pe el însuși, rezultă că relația " $<$ " este o relație nereflexivă, de *ordine parțială*. Dacă două evenimente A și B nu sunt în relație $<$ (adică A nu l-a precedat pe B și nici B nu l-a precedat pe A) atunci cele două evenimente sunt concurente: nici unul din aceste evenimente nu îl poate afecta, cauzal, pe celălalt. Relația *happend-before* este reprezentată printr-un exemplu, utilizând o diagramă timp-spațiu în fig. 4.1.1. Evenimentele p_1, q_1, r_0 corespund unor emisii de mesaje iar p_4, q_4, r_1 unor



recepții de mesaje; se observă că două evenimente A și B sunt *concurente* dacă și numai dacă nu există o cale de la A la B sau de la B la A.

4.2. Soluții centralizate pentru sincronizare

Soluțiile centralizate pentru sincronizare desemnează o entitate unică a sistemului distribuit ca entitate coordonatoare și, implicit, nodul pe care se afla aceasta, ca nod coordonator. O astfel de entitate poate fi un ceas fizic global sau un proces coordonator unic în sistem.

Definiție:

O entitate de sincronizare se spune că este *coordonatoare* sau *centrală* dacă: (1) posedă un nume unic, cunoscut de toate procesele care o folosesc pentru sincronizare; (2) în orice moment ea poate fi accesată de oricare proces.

Un algoritim centralizat este caracterizat prin două proprietăți: (1) deciziile se iau numai la nivelul nodului coordonator; (2) toate informațiile necesare deciziilor sunt concentrate în nodul coordonator.

Dezavantajele algoritmilor centralizați constau în vulnerabilitatea lor ca urmare a defectării nodului coordonator și în circulația intensă de mesaje în jurul nodului coordonator.

Unele sisteme, proiectate cu toleranță la defecte, furnizează, printr-o tehnică specială de recuperare, posibilitatea înlocuirii entității coordonatoare (există algoritmi specifici de alegere a unui coordonator pentru aceste cazuri, vezi subcapitolul 4.4). Între două avarii grave, care au ca urmare pierderea coordonatorului și alegerea altuia, definiția de mai sus rămâne valabilă.

4.2.1. Utilizarea unui ceas fizic global

Soluția utilizării unui *ceas fizic global* presupune utilizarea unui ceas fizic unic în sistemul distribuit, analog sistemelor centralizate. Pentru a ordona evenimentele, procesele construiesc mărcile de timp utilizând ceasul fizic global. Simplă, această soluție are numeroase dezavantaje ([Gos91]):

- (1) valoarea corectă a mărcilor de timp depinde în întregime de posibilitatea recepționării corecte, în toate momentele, a valorii ceasului fizic;
- (2) erorile de transmisie pot altera ordinea evenimentelor;
- (3) este necesară estimarea apriori sau posteriori a întârzierilor de transmisie;
- (4) gradul de perfecțiune al ceasului influențează constrîngerile care se pot respecta în sistem.

4.2.2. Utilizarea unui proces central

Soluția utilizează un proces unic în sistemul distribuit, coordonator, care recepționează cererile de sincronizare ale tuturor proceselor din sistem. Sincronizarea care se poate obține astfel este de tipul excludere reciprocă și necesită trei mesaje: două pentru intrarea într-o secțiune critică și unul pentru ieșire, astfel: (1) procesul care cere intrarea într-o secțiune critică trimite o cerere coordonatorului; (2) coordonatorul verifică dacă sunt alte procese în secțiunea critică. În caz negativ, trimite un mesaj de răspuns, de acceptare a intrării în secțiunea critică; în caz afirmativ, cererea este introdusă într-o coadă de așteptare; (3) la primirea mesajului de răspuns, procesul intră în secțiunea critică; (4) la ieșirea din secțiunea critică, procesul trimite coordonatorului un mesaj de terminare secțiune critică; (5) la primirea mesajului, coordonatorul preia din coada de așteptare următoarea cerere de intrare într-o secțiune critică și trimite un mesaj de răspuns aceluia proces care a lansat cererea.

Utilizînd o planificare FIFO, nu poate avea loc aminarea nedefinită a proceselor la intrarea în secțiunea critică. Soluția are dezavantajul că necesită trei mesaje pentru o secțiune critică și că, în cazul dispariției coordonatorului, trebuie aplicat un algoritim de alegere al unui nou coordonator; după desemnare, acesta trebuie să interogheze toate procesele pentru a reconstitui coada de așteptare.

În ceea ce privește utilizarea soluției procesului central pot fi date multe exemple. În bazele de date distribuite, de exemplu, controlul concurenței se poate baza pe un *proces unic*, care primește toate cererile proceselor ce inițiază acțiuni; sincronizarea acțiunilor este realizată de procesul unic. Un alt exemplu este

algoritmul centralizat pentru *detectia interblocărilor* în sistemele distribuite; acesta utilizează un proces central care construiește un graf global al conflictelor (vezi subcapitolul 7.5.3.3.).

4.2.3. Contori de evenimente și secvențiatori

Coordonarea efectivă a proceselor poate fi atinsă dacă se menține o relație de ordine parțială pentru toate evenimentele proceselor, mai precis o relație de ordine peste toate evenimentele înrudite. Acest lucru se realizează cu ajutorul *contorilor de evenimente*.

Definiție:

Se numește *contor de evenimente* o variabilă de tip întreg, inițial cu valoarea 0 și care ia valori în ordine strict crescătoare.

Un contor de evenimente se manipulează prin următoarele operații: (1) *advance*. Operația *advance(E)* se utilizează pentru a anunța apariția unui eveniment legat de E, și provoacă incrementarea cu 1 a lui E. (2) *read*: Operația *read(E)* se folosește de către un proces pentru a returna valoarea curentă a lui E; (3) *await*. Operația *await(E, v)* provoacă suspendarea procesului apelant atâta timp cât $E < v$.

Un contor de evenimente este analog cu un *ceas global* care înaintea la fiecare apariție a unui eveniment semnificativ, prin execuția explicită a operației *advance*, care apare astfel ca un eveniment de semnalizare. Un proces se sincronizează cu *ceasul global* executând operația *await*, el rămâne astfel blocat pînă cînd *ceasul global* atinge un timp predefinit, sau, altfel spus, pînă cînd se execută v apeluri *advance*. Cînd sincronizarea este mai complexă, se poate folosi și operația *read* care returnează numărul de apeluri *advance* executate. Aceste primitive permit rezolvarea multor probleme de sincronizare, clasice în sistemele de operare: problema producator-consumator ([Gos91], [Nut92]), problema readers-writers ([Gos91]).

Reed și Kanodia au arătat pentru prima oară ([Gos91] [RK79], [Lan83]) cum contorii de evenimente pot fi extinși la sisteme distribuite. Ei au considerat că execuția primitivelor *advance*, *read* și *wait* constituie evenimente și au extins relația de ordine parțială $<$ a lui Lamport (prezentată în 4.1) astfel ([Gos91]):

- (1) Dacă A și B sunt evenimente în același proces, și A a fost executat înaintea lui B, atunci $A < B$;
- (2) Dacă A este un eveniment de emisie al unui mesaj dintr-un proces și B este evenimentul de recepție al aceluiași mesaj, atunci $A < B$;
- (3) Dacă $A < B$ și $C < D$ atunci $A < D$ sau $C < B$;
- (4) Pentru toate evenimentele A, $A < A$ nu este adevărată.

Ca și în modelul lui Lamport, două procese A și B sunt concurente dacă relațiile $A < B$ și $B < A$ nu sunt adevărate. Relația impusă de cele patru condiții de mai sus a fost notată \prec și definită astfel: a "precede" pe b ($a \Rightarrow b$) dacă, și numai dacă, nu este adevărată relația $b < a$. Pe baza ei, Reed și Kanodia au specificat dependențele de timp dintre primitive astfel. Există posibilitatea ca aceste operații, aplicate asupra aceluiași contor să fie concurente. Dependențele de timp sunt:

- (5) Dacă W_E este execuția unei operații *await(E, t)* atunci trebuie să existe cel puțin t membri ai mulțurii $\{A_E, A_E$ este execuția unei operații *advance(E)* și $A_E \Rightarrow W_E\}$;
- (6) Dacă R_E este execuția unei operații *read, v = read(E)*, atunci există cel puțin v membri ai mulțurii $\{A, A_E$ este execuția operației *advance(E)* și $A_E < R_E\}$ și există cel puțin v membri ai mulțurii $\{A, A_E$ este execuția operației *advance(A)* și $A_E \Rightarrow R_E\}$.

Prin *contor de evenimente singular* se desemnează un contor de evenimente care lucrează corect dacă nu există operații *advance* concurente. Într-un sistem distribuit, un contor de evenimente poate fi construit pe baza mai multor contoare de evenimente singulare, fiecare situat pe câte un nod. Numai procesele locale corespunzătoare contoarelor singulare le pot modifica. Aceasta înseamnă că o operație *advance* asupra contorului global va avansa contorul local singular. O operație *read* asupra contorului global se realizează prin citirea contoarelor de evenimente singulare, în orice ordine și însumând valorile rezultate. Procesele care execută *advance* sunt denumite *semnalizatori*, iar cei care execută *read* *observatori*.

Contorii de evenimente stabilesc o relație de ordine parțială peste evenimentele legate ca semnificație de contor, dar nu pot fi utilizați pentru evenimentele neînrudite. În astfel de cazuri se utilizează *secvențiatorii*. Un secvențiator poate fi utilizat pentru a stabili o relație de ordine totală peste evenimentele legate de destinație.

Definiție:

Se numește *secvențiator* o variabilă de tip întreg cu valoare inițială 0, care ia numai valori strict crescătoare.

O singură primitivă este permisă asupra unui secvențiator: primitiva *ticket* (operația *ticket()*) returnează valoarea curentă a secvențiatorului și incrementează cu 1 valoarea lui. Operația *ticket* este atomică.

În fig. 4.2.3.1. se prezintă ca exemplu *problema accesului în citire și actualizare a unei baze de date*, rezolvată cu contori de evenimente și secvențiatori. Procesele care accesează baza de date printr-un nod numai de observator în ceea ce privește contorii de evenimente asociați bazei de date; aceasta înseamnă că nu pot executa operații *advance* asupra acestora și operația *ticket*.

```

Exclusive Acces()
{
/* cand E1 > E2 baza de date este */
/* accesata in acces exclusiv */
int i;
advance (E1);
i <- ticket (T); /* secventiator folosit pentru
excluderea reciproca a proceselor
care cer acces exclusiv */
await (E2, i);
/* asteptare pentru acces exclusiv */
/* Acces exclusiv in baza de date */
advance (E2);
}

Shared Acces()
{
int i;
retry:
S1 <- read (E1);
await (E2, j);
/* asteptare pentru acces partajat */
/* Acces partajat in baza de date */
S2 <- read (E1);
if ( S1!=S2 )
goto retry; /* date inconsistente*/
}
    
```

Fig. 4.2.3.1. Utilizarea indicatorilor de eveniment si a secventatorilor in problema accesului exclusiv si partajat la o baza de date

Soluția Read și Kanodia de implementare a contorilor de evenimente se bazează pe presupunerea că în sistemul distribuit se acceptă erori în comunicație, dar nu și avarii grave ale nodurilor. Bullis și Franta [BF80] au propus o soluție valabilă în rețelele care admit transmisiile broadcast și respectă următoarele cerințe: (1) comunicațiile nesigure, mesajele pot fi pierdute, duplicate sau modificate; (2) nodurile

pot suferi avarii grave, iar informațiile deținute de acestea pierdute complet; (3) avaria unui nod poate fi identificată de celelalte noduri, în situația în care un proces de pe acel nod nu mai răspunde la mesaje. Scopul implementării Bullis și Franta este, pe de o parte, de a asigura funcționarea sistemului chiar și în cazul dispariției unui nod, și pe de altă parte, de a permite revenirea unui nod în sistem după înlăturarea unei avarii. Bullis și Franta au propus ca fiecare nod să conțină o copie a fiecărui contor de eveniment. Aceasta implică transmiterea unui mesaj de broadcast la fiecare operație *advance* a unui contor de evenimente, mesaj care să fie primit de toate nodurile, pentru a asigura consistența contorilor.

Primitiva *await* a fost modificată la sintaxa: *await (E, v, timeout)* pentru a permite deblocarea unui proces care așteaptă nedefinit. Semantica noii primitive *await* este:

*așteapta pana E>=v sau timeout
dacă E>=v atunci returnează 0
al/îl returnează -1*

```

Exclusive Access ()
{
/* cand E1>=E2 baza de date este */
/* accesata in mod exclusiv */
int i;
advance (E1);
t<--ticket (T);
c<--lread (E2);
for ( i=c; i<t; i++ ) {
j<--await (E2,i,timeout);
if ( j= -1 ) return ('eroare');
/* o eroare semnifica un proces in
acces exclusiv ramas in
sectiunea critica
*/
}
/*
Acces exclusiv in baza de date
citire si scriere in baza
*/
advance (E2);
}

Shared Acces ()
{
retry:
S1<--lread (E1);
eroare<--await (E2,S1,timeout);
/* o eroare semnifica un proces in
acces exclusiv ramas in sectiunea critica
*/
if ( eroare != -1 ) {
/*
Acces partajat in baza de date
citire din baza
*/
}
S2<--max (rread (E1,1),...rread (I:1,N));
if ( ( S2 != -1 ) && ( eroare != -1 ) ) {
if ( S1= S2 ) return
else
/* baza de date ar putea
fi utilizata in timpul citirii
*/
goto retry;
}
else
/* daca S2 = -1 toate nodurile
cu acces exclusiv nu mai raspund
*/
/* baza de date este inconsistenta */
return 'eroare'
}
    
```

Fig. 4.2.3.2 Algoritmul Bullis si Franta

Primitiva *read* a fost înlocuită cu alte două primitive: *lread* (local read) și *rread* (remote read). Primitiva *lread* (*E*) returnează valoarea copiei locale a lui *E*; primitiva *rread* (*E*, *n*) returnează valoarea copiei contorului de pe nodul *n* sau -1 în cazul când nodul nu răspunde. În fig. 4.2.3.2. se prezintă același exemplu, al accesului exclusiv și partajat la o baza de date, utilizând implementarea Bullis și Franta a contorilor.

4.3. Soluții distribuite pentru sincronizare

Un *algoritm complet distribuit* este caracterizat de următoarele proprietăți: (1) toate nodurile dețin și mențin aceeași cantitate de informație; (2) toate nodurile iau decizii numai pe baza informațiilor locale; (3) toate nodurile au responsabilități egale în raport cu decizia finală; (4) toate nodurile depun același efort în realizarea deciziei finale; (5) avaria unui nod nu conduce în general la compromiterea întregului sistem.

În practică, toate aceste cinci proprietăți sunt respectate în totalitate de puțini algoritmi distribuiți. Unii algoritmi, de exemplu, pot cere ca fiecare nod să aibă aceeași responsabilitate în decizia finală, dar efortul depus de fiecare nod pentru obținerea excluderii reciproce să nu fie același; alți algoritmi pot conduce la compromiterea sistemului, chiar când un singur nod suferă o avarie gravă.

Indiferent de modul și de gradul de distribuire, există două presupuneri fundamentale, comune tuturor algoritmilor distribuiți:

Ipoțeză 1. Fiecare nod are numai o vedere parțială a întregului sistem, pe care trebuie să-și bazeze deciziile;

Ipoțeză 2. Nu există un ceas fizic global.

Ultima presupunere este o constrângere majoră, avînd în vedere că ordonarea temporală a evenimentelor este un concept fundamental din punct de vedere al sincronizării lor. A spune că într-un sistem distribuit un eveniment *A* a avut loc înainte de alt eveniment *B*, nu poate fi o afirmație sigură din două motive: (1) poate exista o întârziere de transmisie între momentul în care are loc un eveniment și momentul în care este observat; (2) dacă fiecare nod conține un ceas fizic, diferența dintre ele conduce la variații la citirea timpului.

Pentru a construi baza de dezvoltare a algoritmilor distribuiți, trebuie stabilită o metodă de ordonare a evenimentelor într-un sistem distribuit.

4.3.1. Sincronizarea ceasurilor

1. Ceasuri fizice

S-a aratat în secțiunea 4.2.1. că utilizarea unui ceas fizic unic, global într-un sistem distribuit prezintă o serie de dezavantaje, dintre care cel mai important este necesitatea aprecierii întârzierilor de transmisie. Dacă se utilizează cîte un ceas fizic în fiecare nod, atunci ele trebuie sincronizate astfel încît diferența relativă a gricaror două ceasuri fizice să se păstreze sub o limită constantă fixată. Deoarece două ceasuri nu vor funcționa niciodată cu aceeași rată, diferența dintre ele crescînd mereu, Lamport ([Lam78][BBK91][Gos91]) a prezentat o soluție a problemei sincronizării ceasurilor, prin introducerea unui algoritm, pe baza a două condiții și a două reguli.

Prima condiție descrie funcționarea corectă a unui ceas fizic, astfel:

Condiția 1: Există o constantă $k \ll 1$, astfel încît pentru toate ceasurile i : $|dC_i/dt - 1| < k$, unde $C_i(t)$ este o funcție continuă, diferențiabilă, iar dC_i/dt reprezintă rata cu care funcționează ceasul i la momentul t .

În plus, toate ceasurile trebuie sincronizate astfel încît $C_i(t) \approx C_j(t)$, pentru orice i, j, t , astfel:

Condiția 2: Există o constantă suficient de mică astfel încît pentru tot i, j , $|C_i(t) - C_j(t)| \leq \epsilon$.

Deoarece întotdeauna valorile a două ceasuri tind să se îndepărteze, este necesar un algoritm care să asigure respectarea condiției 2. Pentru a dezvolta un astfel de algoritm, se presupune că procesul P_i trimite un mesaj m la timpul t ; mesajul este primit de P_j la momentul t' , deci întârzierea de transmisie este $\delta_m = t' - t$. Procesul destinație cunoaște numai întârzierea minimă $\mu_m \geq 0$ ($\mu_m < \delta_m$). Diferența dintre aceste două întârzieri se numește *diferență nepredictibilă* $\rho_m = \delta_m - \mu_m$. Pot fi enumerate acum regulile pentru ceasurile fizice:

Regula 1: Pentru fiecare i , dacă P_i nu recepționează un mesaj la momentul t , atunci C_i este diferențiabilă la t și $dC_i(t)/dt > 0$.

Regula 2:

(a) Dacă P_j trimite un mesaj m la momentul t , atunci m conține o marea de timp $TS_m = C_j(t)$.

(b) După primirea mesajului m la momentul t' , procesul P_j stabilește $C_j(t')$ egal cu $\max(C_j(t'-0),$

$TS_m + \mu_m)$.

Un proces trebuie deci să cunoască numai ceasul său fizic și marea de timp a mesajelor recepționate. Trebuie arătat că algoritmul prezentat prin cele două reguli de mai sus respectă condiția 2. Presupunem că sistemul este modelat printr-un graf de procese puternic conectat; arcele vor reprezenta liniile de comunicație pe care pot fi trimise direct mesaje între procese. *Diametrul* acestui graf este cel mai mic număr d , astfel încît,

pentru orice pereche de procese P_i, P_j există o cale de la P_i la P_j având cel mult d arce. Vom spune că, pe un arc, se trimite un mesaj de la P_i la P_j la fiecare τ secunde, dacă pentru orice t , P_i trimite cel puțin un mesaj la P_j între momentele $t, t+\tau$. La fiecare τ secunde se trimite un mesaj de sincronizare pe fiecare arc, care conține o marcă de timp TS . După primirea mesajului, dacă este nevoie, un proces va avansa ceasul său local la valoarea conținută în mesaj. Se presupune că întârzierile de transmisie sunt limitate de o limită inferioară u și una superioară $u+z$, cunoscute. Fie K marja de eroare a ceasului, $K < 10E-6$ și ϵ diferența permisă pentru două ceasuri. Dacă $\epsilon/(1-K) < u$ și $\epsilon < \tau$ atunci este posibil să se aprecieze valoarea aproximativă pentru ϵ ca fiind $d(2K\tau+z)$.

În final trebuie precizat că, în funcție de cerințele legate de diferența dintre ceasuri și de validitatea presupunerilor legate de limitele întârzierilor, se poate decide că, fie se poate asuma riscul pierderii unor mesaje de sincronizare din când în când, datorită numărului excesiv de mare de mesaje care trebuie transmise, obținându-se astfel o sincronizare numită *probabilistică*, fie nu se poate asuma acest risc; parametrul cheie este raportul z/u .

Un algoritmul de sincronizare a ceasurilor fizice ([Tan96]) se bazează pe divizarea timpului în intervale fixe de resincronizare și pe transmiterea de mesaje broadcast ce conțin timpul curent de către fiecare nod la începutul fiecărui interval $[T_0+iR, T_0+(i+1)R]$ (T_0 este un moment în trecut stabilit de comun acord, iar R un parametru de sistem). După ce un nod a transmis mesajul broadcast, startează un timer local pentru a colecta toate celelalte mesaje în timpul unui interval S . După ce s-au primit toate mesajele se calculează timpul curent prin media valorilor primite în mesaje, eventual eliminând cele mai mari respectiv, mici m valori.

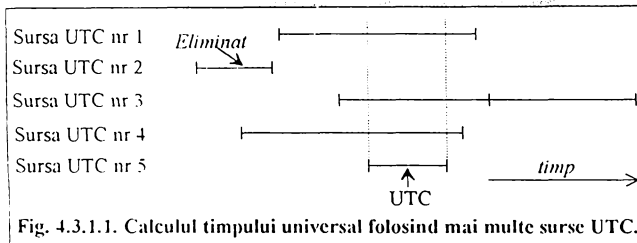


Fig. 4.3.1.1. Calculul timpului universal folosind mai multe surse UTC.

Pentru sisteme distribuite în care se cere o sincronizare cât mai perfectă cu timpul universal, UTC (Universal Coordinated Time) pot fi folosite mai multe surse UTC.

Pentru sincronizarea acestora se prezintă ca exemplu algoritmul utilizat în DCE: periodic, la începutul unui interval de timp fixat, (de exemplu 1 min) fiecare procesor dotat cu o sursă UTC transmite un mesaj broadcast către toate nodurile. Când un procesor primește mesaje de la toate nodurile cu surse UTC, elimină întâi specificațiile disjuncte (figura 4.3.1.1) și apoi calculează timpul curent ca mijlocul intersecției tuturor specificațiilor.

2. Ceasuri logice

În multe aplicații, este suficient ca toate calculatoarele să perceapă același timp, nefiind necesar ca acesta să coincidă cu timpul real. Pentru aceasta clasă de aplicații este importantă consistența ceasurilor și nu dacă au sau nu valoarea timpului real; de aceea se folosesc *ceasurile logice*, care nu implică prezența unui ceas fizic. Ceasurile logice permit extinderea relației *happend-before* la o relație totală de ordine pe evenimentele unui sistem distribuit ([Lam78]). Cerințele unei astfel de relații totale sunt: pentru orice pereche de evenimente A și B , dacă $A < B$ atunci marca de timp a lui A este mai mică decât marca de timp a lui B (reciproca nu este valabilă).

Un ceas logic LC_i se asociază unui proces P_i ; el este definit ca o funcție care asociază un număr întreg nenegativ $LC_i(a)$ fiecărui eveniment a din proces. Întregul sistem de ceasuri logice se reprezintă printr-o funcție LC care atribuie oricărui eveniment b numărul $C(b)$ unde $C(b) = LC_i(b)$ dacă b este un eveniment în procesul P_i . Pe baza definiției poate fi enunțată acum condiția de ceas logic a lui Lamport:

Definiție.

Fie C o funcție care asignează o valoare nenegativă fiecărui eveniment e_i ; pentru evenimentele e_i, e_j avem atunci: $e_i < e_j \Rightarrow LC(e_i) < LC(e_j)$

Rezultă că ceasul logic poate fi implementat printr-un contor care este incrementat între oricare două evenimente succesive dintr-un proces. Deoarece un contor furnizează valori crescătoare, el asignează valori unice fiecărui eveniment iar dacă un eveniment a precedă într-un proces P_i un eveniment b , atunci $LC(a) < LC(b)$. Marca de timp asociată unui eveniment este valoarea contorului (ceasul logic) pentru acel eveniment. Această schema asigură, într-adevar cerințele relației de ordine totală, dar numai pentru evenimentele ale aceluiași proces; ea trebuie completată având în vedere situația proceselor care-și comunică mesaje.

Fie $\{p_0, p_1, \dots, p_{n-1}\}$ un set de procese, fiecare cu un număr de evenimente și fie LC_i ceasul logic corespunzător procesului i , care asignează un număr întreg nenegativ evenimentelor procesului. Condițiile următoare asigură respectarea globală a condiției de ceas logic a lui Lamport:

(1) Dacă evenimentele e_i, e_k au loc în același proces și e_j are loc înaintea lui e_k , atunci: $LC_i(e_i) < LC_i(e_k)$

(2) Dacă e_i este un eveniment de emisie a unui mesaj de către procesul P_i , iar e_j este un eveniment de recepție a acestui mesaj de către procesul P_j , atunci: $LC_j(e_j) < LC_i(e_i)$

Pe baza acestor condiții se poate defini complet modul de implementare al ceasului logic printr-un contor. Fie LC_i contorul asociat procesului P_i ; el devine un ceas logic dacă:

(1) C_i este incrementat între două evenimente succesive ale procesului

(2) Dacă evenimentul e_i este o emisie unicast sau multicast a mesajului m de către procesul P_i , atunci mesajul i va include marca de timp a evenimentului TS_m (valoarea ceasului logic). Când procesul P_j recepționează mesajul m , setează ceasul sau logic LC_j la valoarea $\max(TS_m+1, LC_j)$.

Condiția (2) de mai sus asigură executarea corecției de timp numai prin adăugarea de valori pozitive ceasurilor, niciodată prin scădere. În acest caz, ceasul unui proces P_j va fi sincronizat cu ceasul procesului P_i când se transmite un mesaj de la P_i la P_j . Figura 4.3.1.2. reprezintă succesiunea valorilor ceasurilor logice pentru trei procese. Diagrama din fig. 4.3.1.4. reprezintă progresul logic a trei procese arătate în fig. 4.3.1.3.,

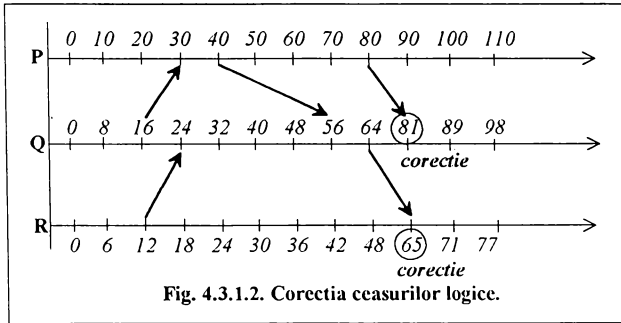


Fig. 4.3.1.2. Corectia ceasurilor logice.

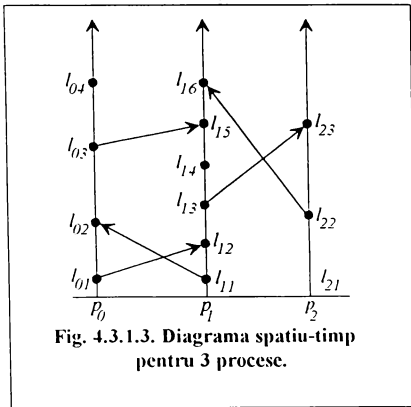


Fig. 4.3.1.3. Diagrama spatiu-timp pentru 3 procese.

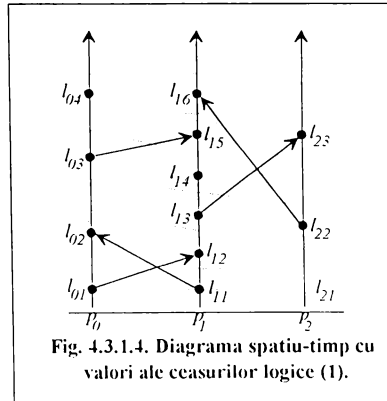


Fig. 4.3.1.4. Diagrama spatiu-timp cu valori ale ceasurilor logice (1).

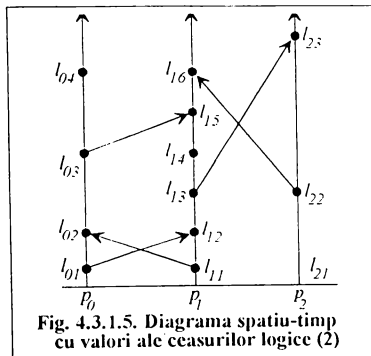


Fig. 4.3.1.5. Diagrama spatiu-timp cu valori ale ceasurilor logice (2)

sugerind o viziune consistentă a relației de ordine a evenimentelor în sistem; fig. 4.3.1.5. reprezintă de asemenea o viziune consistentă, de această dată sugerind că ceasurile locale sunt sincronizate în timp.

În final, pentru a realiza complet o relație de ordine totală, mai trebuie specificat că, în cazul a două evenimente concurente ordonarea trebuie făcută ținând cont de identitatea proceselor (procesele se ordonează după un criteriu oarecare, de exemplu în ordinea identificatorilor lor).

4.3.2. Algoritmi distribuți pentru excludere reciprocă

Deși diferă prin cerințe și obiective, algoritmi distribuți pentru excludere mutuală (Lamport [Gos91][Lam78], Ricard și Agrawala [RA81], Maekawa [MOO87]) au comune următoarele cerințe și caracteristici:

(1) Sistemul distribuit este văzut ca un sistem cu N noduri, numerotate de la 1 la N . Fiecare nod conține un proces care cere acces exclusiv la un obiect partajat. Cererea sa este comunicată și celorlalte procese.

(2) Proprietatea *pipeline* este adevărată; aceasta înseamnă că mesajele trimise de la un proces la altul sunt numerotate în ordinea transmiterii lor.

(3) Fiecare mesaj este livrat corect la destinație într-un timp finit;

(4) Rețeaua este complet conectată; aceasta înseamnă că fiecare proces poate trimite mesaje direct la oricare alt proces.

Cerințele 2 și 3 se realizează prin metode descrise în capitolul 2.

În acest model al sistemului distribuit, fiecare proces, în afară de faptul că cere un acces exclusiv la obiectul partajat, servește și ca arbitru pentru a rezolva cererile care se suprapun în timp. Pentru a arbitra cererile, orice pereche de două cereri trebuie să fie cunoscută de cel puțin unul dintre arbitrii. Să presupunem, de exemplu, că pentru a obține acces exclusiv la obiect procesul i trebuie să obțină permisiunea fiecărui membru al unui subset S_i de procese din sistem; atunci, pentru a asigura excluderea reciprocă a proceselor i și j , S_i și S_j trebuie să conțină un arbitru comun. Formal, acest lucru se exprimă astfel:

(a) *Proprietatea intersecției nenule* : $S_i \cap S_j \neq \emptyset$ pentru orice combinație i, j cu $1 \leq i, j \leq N$. Dacă se respectă această proprietate, excluderea reciprocă este garantată dacă orice proces care cere acces la obiectul partajat obține permisiunea de la fiecare membru al lui S_i . Această proprietate este condiția necesară pentru realizarea excluderii reciproce. În plus, uneori sunt valabile și următoarele proprietăți:

(b) *Regula Efortului Egal*: $|S_1| = \dots = |S_N| = K$; $|S_i|$ desemnează dimensiunea lui S_i .

(c) *Regula Responsabilității Egale*: Fiecare proces j , $1 < j <= N$ este conținut în exact același număr D de seturi S_i , $1 < i <= N$.

(d) S_i , $1 < i <= N$, conține întotdeauna procesul i .

Proprietățile (b) și (c) sunt necesare pentru un algoritm complet distribuit; proprietate (d) este inclusă pentru a reduce cu 1 numărul de mesaje de trimis și de recepționate.

În condițiile de mai sus, se poate considera că un algoritm centralizat desemnează un singur proces ca arbitru (procesul coordonator). Acest arbitru este conținut în fiecare S_i și astfel este satisfăcută proprietatea (b) cu $k=1$. Totuși, proprietatea (c) nu este satisfăcută, deoarece nodul coordonator este singurul nod conținut în fiecare S_i .

4.3.2.2. Algoritmul lui Lamport

A. Descrierea algoritmului

În plus față de cerințele prezentate mai sus pentru realizarea excluderii reciproce, problema lui Lamport necesită ca cererile de acces la un obiect să fie servite în ordinea în care au fost făcute. Algoritmul distribuit următor, cunoscut ca *algoritmul lui Lamport*, satisface toate cerințele problemei. Algoritmul impune următoarele:

(1) Mesajul în care se încapsulează o cerere de acces exclusiv de la procesul P_i conține o marcă de timp (TS_i, i) unde $TS_i \in C$, este valoarea ceasului logic pentru procesul i , din momentul trimiterii mesajului.

(2) Fiecare proces menține o coadă de cereri de acces, inițial vidă, care conține mesajele ordonate după relația \Rightarrow .

Algoritmul se definește prin cinci reguli; se presupune că fiecare regulă se realizează printr-o operație indivizibilă:

(1) Când un proces cere accesul la obiectul partajat, el depune marca în coada de cereri proprii și apoi trimite un mesaj REQUEST (TS_i, i) la fiecare proces care participă la algoritm.

(2) Când un proces P_j recepționează mesajul REQUEST (TS_i, i) returnează imediat un mesaj REPLY însoțit de marca de timp și plasează mesajul recepționat în coada de mesaje proprii.

(3) Procesul P_i are acces la obiect când una din următoarele două condiții este satisfăcută: (a) P_i are în varful cozii propria sa cerere; (b) P_i a primit un mesaj de la toate procesele cu o marcă de timp anterioară lui (TS_i, i) .

(4) Pentru a elibera obiectul, P_i înlătură mesajul REQUEST din coada sa proprie și-i trimite un mesaj RELEASE, însoțit de o marcă de timp la toate procesele.

(5) Cînd un proces P_j recepționează mesajul RELEASE de la P_i , înlătură din coadă mesajul REQUEST al acestuia.

Verificarea corectitudinii algoritmului se face pe baza observațiilor următoare:

(1) Regula (3)(b) împreună cu presupunerea că mesajele sunt recepționate în ordinea emisie, garantează că P_i recunoaște corect toate cererile care preced propria sa cerere.

(2) Deoarece relația \Rightarrow ordonează total cererile, regula (3)(a) garantează că se permite numai unui singur proces să aibă acces la un moment dat la un obiect.

Algoritmul satisface condițiile (b) și (c) pentru un algoritim distribuit (secțiunea 4.3.2.): fiecare proces menține toate informațiile și simulează un proces planificator central.

Numărul total de mesaje necesare pentru a obține și elibera accesul exclusiv la obiect este $3*(N-1)$: (N-1) mesaje REQUEST, (N-1) mesaje REPLY și (N-1) mesaje RELEASE.

B. Implementare.

O descriere generală a implementării algoritmului lui Lamport se poate face utilizînd notația SR [And91].

Excluderea reciprocă este asigurată, de fapt, de un set de procese Server $[1:n]$, cîte unul asociat fiecărui nod i . Cînd un proces $Server[i]$, în urma unei cereri de intrare în secțiune critică primită de la clientul său asociat, decide că se poate intra în secțiune critică, deblochează clientul i , care intră astfel în secțiunea critică. La ieșirea din secțiunea critică, $Client[i]$ semnalizează procesului $Server[i]$.

```

type kind = enum (REQUEST, REPLY, RELEASE) # tipuri de mesaje
chan remoteChan[1:n](kind, timeStamp: int, sender: int) #canal pentru comunicatie in retea
chan localChan[1:n](kind) # canal pentru transmiterea cererilor de intrare si de iesire din sectiunea
# critica (comunicatie locala intre clienti si server)
chan enterChan[1:n]() # canal pentru deblocarea unui client in asteptare la intrare in sectiune

```

Client [$i: 1..n$]:

```

do true →
    send localChan[i] (REQUEST) # protocol de intrare in sectiunea critica
    receive enterChan[i] ( )
    #sectiune critica...
    send localChan[i] (RELEASE) # protocol de iesire din sectiunea critica
od

```

Server [$i: 1..n$]:

```

var reqQueue: queue of (int,int) # coada de cereri ordonata dupa marcile de timp
var localTimeStamp: int # marca de timp ce va fi asignata urmatoarei cereri locale
var highTimeStamp: int := 0 # marca de timp cea mai mare dintre o marca de timp primita
# intr-un mesaj si marca de timp locala
var replyCounter: int # contorizeaza numarul de confirmari care trebuie asteptate
var inAccess: bool := false # indica accesul intr-o sectiune critica
var kind: kind, sender: int, ts: int # folosite la citirea mesajelor

```

```

do true -> # bucla infinita a serverului
    if not (empty (localChan[i])) →
        receive localChan[i](kind)
        if kind = REQUEST →
            localTimeStamp := highTimeStamp + 1
            replyCounter := n - 1
            fa j := 1 to n →
                send remoteChan[j](REQUEST, localTimeStamp, i)
            af
        [] kind = RELEASE →
            fa j := 1 to n st j ≠ i →
                send remoteChan[j](RELEASE, localTimeStamp, i)
            af
fi

```

```

fi
if not (empty(remoteChan[i])) →
    receive remoteChan[i](knd, ts, sender)
    if knd = REQUEST →
        insereaza cererea in coada queue
        highTimeStamp = max(localTimeStamp, ts)
        if sender ≠ i → send remoteChan[sender](REPLY, localTimeStamp, i) fi
    [] if knd = REPLY →
        replyCounter := replyCounter - 1
        if replyCounter = 0 and (cererea proprie este prima din coada) →
            inAccess := true;
            send enterChan[i]()
        fi
    [] knd = RELEASE →
        sterge din coada queue cererea de la sender
        if replyCounter = 0 and (prima cerere din coada este o cerere proprie) →
            inAccess := true;
            send enterChan[i]()
        fi
    fi
fi
od

```

4.3.2.2. Algoritmul lui Ricard și Agrawala

A. Descrierea algoritmului

Algoritmul lui Ricard și Agrawala rezolvă aceeași problemă de excludere reciprocă, dar este mai eficient decât soluția lui Lamport, deoarece necesită cel mult $2*(N-1)$ mesaje. Ca și în algoritmul lui Lamport, cozile de cereri sunt total ordonate de relația \Rightarrow . Algoritmul este descris prin următoarele reguli;

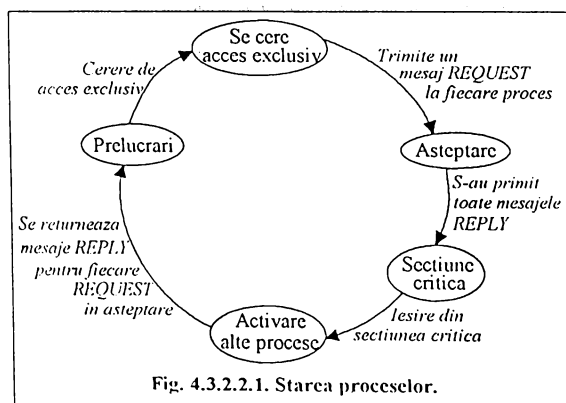
(1) Când un proces P_i cere accesul la obiectul partajat, trimite un mesaj REQUEST (TS_i, i) la următoarele procese.

(2) Când un proces P_j recepționează mesajul REQUEST execută următoarele operații:

(a) Dacă P_j este într-o secțiune critică, reține mesajul REQUEST al procesului i .

(b) Dacă P_j nu cere și el acces la obiect, returnează un mesaj REPLY (însoțit de marca de

timp).



(c) Dacă procesul cere și el acces la obiect și marca de timp a cererii sale, (TS_j, j) precede pe (TS_i, i), atunci reține mesajul REQUEST al procesului P_i ; în caz contrar, returnează un mesaj REPLY.

(3) Când procesul P_i recepționează mesajele REPLY de la toate nodurile, are acces exclusiv la obiect.

(4) Când procesul P_i eliberează obiectul, transmite un mesaj REPLY pentru fiecare mesaj REQUEST în așteptare.

Diagrama de tranziție a stărilor pentru fiecare proces este reprezentată în fig. 4.3.2.2.1.

Analiza acestui algoritm relevă următoarele proprietăți:

(1) Toate cererile vor fi total ordonate deoarece: (a) Unei cereri noi i se asignează un număr de secvență mai mare decât toate numerele de secvență ale oricăror cereri REQUEST cunoscute nodului; (b) Arbitrajul necesar în cazul unor cereri REQUEST trimise de noduri care nu au cunoștință ficcare de cererea celuiilalt se face pe baza numerelor nodurilor.

(2) Controlul este complet distribuit.

(3) Interblocări nu pot apare deoarece nu pot apare cicluri cu noduri care așteaptă la intrare în secțiunea critică.

(4) Aminarea nedefinită nu poate apare deoarece ordinea de deservire a cererilor este FIFO.

(5) Algoritmul este independent de caracteristicile comunicației în sistem.

(6) Algoritmul necesită cel mult $2*(N-1)$ mesaje pentru a deservi o cerere.

Algoritmul are totuși și două dezavantaje:

(1) Fiecare proces trebuie să cunoască identitatea tuturor celorlalte procese din sistem. Când un proces nou se alătură unui grup de procese care participă la algoritmul de excludere reciprocă, sunt necesare următoarele acțiuni: (a) Procesul trebuie să primească numele tuturor proceselor din grup; (b) Numele noului proces trebuie făcut cunoscut tuturor proceselor din grup; Ambele acțiuni (a) și (b) se fac avînd în vedere că pot exista mesaje REPLY și REQUEST care circulă în sistem în momentul cînd noul proces se alătură grupului. Pot fi utilizate comunicații de grup sau fiecare proces menține o listă cu toți membrii grupului. Algoritmul funcționează bine pentru grupuri mici, care nu își schimbă niciodată componența.

(2) Dacă unul dintre procese dispare, atunci întreaga schemă este compromisă. Această dificultate poate fi rezolvată prin monitorizare continuă a stării tuturor proceselor din sistem. Dacă un proces dispare, atunci restul proceselor sunt anunțate, pentru a nu mai trimite mesaje REQUEST aceluia nod. În momentul readucerii lui în sistem, trebuie inițializată procedura de alăturare a unui proces grupului.

B. Implementare

```

type kind = enum (REQUEST, REPLY, RELEASE)      # tipuri de mesaje
chan remoteChan[1:n](kind, timeStamp: int, sender: int) #canal pentru comunicatie in retea
chan localChan[1:n](kind)                       # canal pentru transmiterea cererilor de intrare si de iesire din sectiunea
                                                # critica (comunicatie locala intre clienti si server)
chan enterChan[1:n]()                           # canal pentru deblocarea unui client in asteptare la intrare in sectiune
    
```

Client $[i: 1..n]::$

```

do true →
    send localChan[i] (REQUEST)    # protocol de intrare
    receive enterChan[i] ()        #sectiune critica...
    send localChan[i] (RELEASE)    # protocol de iesire
od
    
```

Server $[i: 1..n]::$

```

var localTimeStamp: int           # marca de timp ce va fi asignata urmatoarei cereri locale
var highTimeStamp: int := 0       # marca de timp cea mai mare dintre o marca de timp primita
                                  # intr-un mesaj si marca de timp locala
var replyCounter: int            # contorizeaza numarul de confirmari care trebuie asteptate
var inAccess: bool := false      # indica accesul intr-o sectiune critica
var reqAccess: bool := false     # indica o cerere de intrare in sectiune critica
var later: bool := false         # indica necesitatea aminarii unui mesaj
var laterTab[1:n]: bool := ([n] false) # indica nodurile la care trebuie trimise confirmari la
                                  # iesirea din sectiunea critica
var kind: kind, sender: int, ts: int # folosite la citirea mesajelor

do true -> # bucla infinita a severului
    if not (empty (localChan[i])) →
        receive localChan[i](kind)
        if kind = REQUEST →
            localTimeStamp := highTimeStamp + 1
            reqAccess := true; replyCounter := n - 1
            fa j := 1 to n stj ≠ i →
    
```

```

        send remoteChan[j](REQUEST, localTimeStamp, i)
    af
[] kno = RELEASE →
    inAccess := false
    fa j := 1 to n st laterTab[j] = true →
        send remoteChan[j](REPLY, localTimeStamp, i)
        laterTab[j] := false
    af
fi
if not (empty(remoteChan[i])) →
    receive remoteChan[i](kno, ts, sender)
    if kno = REPLY →
        replyCounter := replyCounter - 1
        if replyCounter = 0 →
            inAccess := true; reqAccess := false
            send enterChan[i]()
        fi
[] kno = REQUEST →
    highTimeStamp := max(localTimeStamp, ts)
    later := (inAccess = true) or (reqAccess = true and ((ts > localTimeStamp)
        or (ts = localTimeStamp) and (sender > i)))
    if later = true → laterTab[sender] := true
    [] later ≠ true → send remoteChan[sender](REPLY, localTimeStamp, i)
    fi
fi
od

```

4.3.2.3. Algoritm pentru excludere reciprocă în accesul la K resurse identice

A. Descrierea algoritmului

Algoritmii prezentați până acum permit numai unui singur proces să fie în secțiunea critică. Raymond [Ray87] a elaborat un algoritm, de fapt o extensie a algoritmului Ricard și Agrawala, care permite intrarea a maxim K noduri în secțiunea critică. Aceasta problemă apare când există K obiecte identice care trebuie partajate între N procese.

Algoritmul se bazează pe observația că se permite unui proces să fie în secțiunea critică dacă nu mai mult de K-1 procese (din N-1) sunt, la acel moment, în secțiunea critică. Un proces care trimite N-1 mesaje pentru a anunța dorința de a intra într-o secțiune critică, poate intra când primește N-K mesaje REPLY, deoarece aceasta înseamnă că cel mult N-1-(N-K)=K-1 procese pot fi în secțiunea critică.

O diferență majoră între algoritmul lui Raymond și algoritmul Ricard și Agrawala constă în momentele în care pot apare mesaje REPLY, în cel de-al doilea, mesajele REPLY pot fi recepționate numai când procesul așteaptă la intrarea în secțiunea critică. În algoritmul Raymond, mesajele REPLY pot fi recepționate de un nod când este în așteptare la intrarea în secțiunea critică, când este în secțiunea critică sau după ieșirea din secțiunea critică. Pentru a distinge mesajele REPLY, trebuie menținută de către fiecare proces o structură de date care să memoreze câte mesaje REPLY mai sunt de recepționat de la fiecare proces.

În algoritmul Ricard și Agrawala, tabloul *laterTab[]* era de tip logic, deoarece un proces nu putea întârzia decât cel mult un mesaj REPLY către alt nod. În algoritmul Raymond este posibil de a amina mai multe mesaje REPLY către un nod; de aceea în tabloul *laterTab[]* elementele au valori întregi.

B. Implementare

```

const objectNo: int := k                #numarul de obiecte disponibile
type kind = enum (REQUEST, REPLY, RELEASE) # tipuri de mesaje
chan localChan[1:n](kind)              # canal pentru transmiterea cererilor de intrare si de iesire din sectiunea
                                        # critica (comunicatie locala intre clienti si server)
chan enterChan[1:n]()                  # canal pentru deblocarea unui client in asteptare la intrare in sectiune
chan remoteChan[1:n](kind, intV: int, sender: int) #canal pentru comunicare in retea

```

Client [i: 1..n]:


```

do true →
    send localChan[i] (REQUEST) # protocol de intrare
    receive enterChan[i] ( )
    #sectiune critica...
    send localChan[i] (RELEASE) # protocol de iesire
od
Server [i: 1..n]:
var localTimeStamp: int # marca de timp ce va fi asignata urmatoarei cereri locale
var highTimeStamp: int := 0 # marca de timp cea mai mare dintre o marca de timp primita
    # intr-un mesaj si marca de timp locala
var replyCounter[1:n]: int := ([n] 0) # contorizeaza numarul de confirmari care trebuie
    # primite de la celelalte noduri
var inAccess: bool := false # indica accesul intr-o sectiune critica
var reqAccess: bool := false # indica o cerere de intrare in sectiune critica
var later: bool := false # indica necesitatea aminarii unui mesaj
var laterTab[1:n]: int := ([n] 0) # indica noduri le si numarul de confirmari care trebuie trimise la
    # iesirea din sectiunea critica
var kind: kind, sender: int, intValRec: int # folosite la citirea mesajelor

do true → # bucla infinita a severului
    if not (empty (localChan[i])) →
        receive localChan[i](knd)
        if knd = REQUEST →
            localTimeStamp := highTimeStamp + 1
            reqAccess := true;
            fa j := 1 to n st j ≠ i →
                send remoteChan[j](REQUEST, localTimeStamp, i)
                replyCounter[i] := replyCounter[i] + 1
            af
        [] knd = RELEASE →
            inAccess := false
            fa j := 1 to n st laterTab[j] ≠ 0 →
                send remoteChan[j](REPLY, laterTab[j], i)
                laterTab[j] := 0
            af
        fi
    fi
    if not (empty(remoteChan[i])) →
        receive remoteChan[i](knd, intValRec, sender)
        if knd = REPLY →
            replyCounter[sender] := replyCounter[sender] - intValRec
            inregistreaza in cnt numarul de confirmari primite: este
            numarul de elemente nule din replyCounter
            if (reqAccess = true and inAccess = false and cnt ≥ n - objectNo) →
                inAccess := true; reqAccess := false
                send enterChan[i]( )
            fi
        [] knd = REQUEST →
            highTimeStamp := max(localTimeStamp, ts)
            later := (inAccess = true) or (reqAccess = true and ((ts > localTimeStamp)
                or (ts = localTimeStamp) and (sender > i)))
            if later = true → laterTab[sender] := laterTab[sender] + 1
            [] later ≠ true → send remoteChan[sender](REPLY, 1, i)
            fi
        fi
    fi
od

```

4.3.2.4. Alți algoritmi

Alți algoritmi dezvoltati pentru rezolvarea problemei accesului exclusiv într-un sistem distribuit au avut drept obiectiv micșorarea și mai mult a numărului de mesaje. Astfel, Thomson [MOO87] a propus o *regula de consens* a majorității bazată pe ideea că, pentru a se obține acces la un obiect partajat este necesar permisiunea numai a majorității proceselor. Corect implementată, niciodată două procese nu pot acapara votul majorității deoarece intersecția a doua majorități este o mulțime nenulă de arbitri. De aceea cel puțin un arbitru va vota favorabil pentru o cerere și nefavorabil pentru cealaltă. Numărul de mesaje necesar pentru a obține permisiunea de acces la un obiect partajat este, în cel mai bun caz, $n/2$. Extensii ale acestei metode, propuse de Gifford [MOO87] și Skeen permit unor procese să contribuie cu mai mult de un vot. În aceste scheme este suficient să se obțină o majoritate a voturilor, care nu este neapărat o majoritate a arbitrilor. Printr-o corectă distribuție a voturilor, algoritmul poate varia de la o formă centralizată la o formă distribuită. De exemplu, dacă unui proces i se acordă toate voturile, în timp ce celorlalte nu li se acordă nici un vot, algoritmul este redus la un control central.

Este evident că algoritmi Lamport, Ricard și Agrawala și Thomas satisfac proprietățile (a), (b) și (c) prezentate în secțiunea 4.3.2, cu $K=D=n$. În general, algoritmi cu voturi de greutate diferită, nu satisfac proprietatea (c). Dacă greutatea voturilor nu este egal distribuită printre procese, atunci algoritmul rezultat nu este total distribuit. Un algoritim care reduce și mai mult numărul de mesaje necesare pentru obținerea accesului exclusiv este algoritmul lui Mackawa [MOO87]; acesta alege pentru K și D valorile minime care satisfac condițiile (b) și (c).

4.3.2.5. Consideratii finale

În secțiunea 4.3.2 s-a considerat că sistemul distribuit se bazează pe o rețea complet conectată ceea ce înseamnă că există o linie de comunicație directă între oricare două noduri. Dacă această cerință nu este satisfăcută, la optimalitatea unui algoritim pentru excludere mutuală contribuie și un alt factor, numărul total de noduri intermediare necesare pentru a trimite un mesaj de la sursa la destinație. Pe de altă parte, unele caracteristici speciale ale rețelei pot permite și alte optimizări ale algoritimilor.

A. Topologii speciale ale rețelei

Dacă topologia rețelei permite transmiterea de mesaje broadcast, atunci, pentru a transmite mesajele REQUEST un proces poate transmite doar un singur mesaj și nu $(n-1)$.

Dacă topologia rețelei este un inel, algoritmul poate fi modificat pentru a reduce numărul de mesaje și pentru a minimiza numărul de noduri intermediare. Următoarele modificări ale algoritmului elimină necesitatea mesajelor REPLY:

(1) Nodul care cere acces la obiect trimite un mesaj REQUEST la vecinul sau din inel. Mesajele traversează într-o direcție inelul.

(2) Nodul vecin recepționează mesajul REQUEST; dacă acesta are prioritate adecvată, este trimis la vecinul următor.

(3) Când mesajul REQUEST este returnat la nodul inițiator, se permite accesul la obiect.

Astfel, sunt permise n noduri intermediare. Algoritmul este optimal relativ la numărul de noduri intermediare și numărul de mesaje.

B. Modificari dinamice ale topologiei rețelei

Un nod care se alătură sistemului distribuit, trebuie să execute următorii pași:

(1) Să interogheze nodurile active pentru a obține o listă a tuturor nodurilor participante.

(2) Să-și asigneze un număr de nod propriu, unic.

(3) Să-și plaseze numărul de nod pe listele tuturor nodurilor participante.

(4) Să-și stabilească valoarea adecvată pentru marca de timp inițială.

Pasul (2) poate fi complicat dacă se permite o asignare dinamică a numerelor de noduri. O soluție simplă este de a asigna numere fixe și unice fiecărui nod, indiferent dacă este activ sau nu.

În pasul (4), marca de timp trebuie să fie mai mare decât marca oricărui mesaj REQUEST care ar fi putut fi primit dacă nodul ar fi fost activ continuu; aceasta implică comunicare cu celelalte noduri. Până la satisfacerea acestei condiții, noul nod nu poate cere acces la obiect și trebuie să înfirzie orice mesaj REPLY necesar pentru un REQUEST primit.

Un nod care părăsește voluntar sistemul, trebuie să-și anunțe intenția celorlalte noduri. Numai după primirea unor confirmări de la toate nodurile, îi este permis să părăsească rețeaua; aceasta înseamnă că, în timp ce așteaptă confirmările, nu poate participa la comunicații într-un mod care să afecteze excluderea reciprocă (de exemplu, nu mai poate cere acces la obiect și trebuie să returneze REPLY la orice mesaj REQUEST primit).

Pentru a detecta o defecțiune a unui nod, se poate folosi un mecanism de timeout. La trimiterea unui mesaj REQUEST, nodul emițător startează un timer. Dacă timerul detectează timeout, se trimite un mesaj de probă, timp în care se consideră că nodul respectiv nu mai este în sistem. Dacă nodul nu răspunde nici la mesajul de probă, se inițiază procedura de înlăturare a unui nod.

C. Gestiunea numerelor de secvența ale mesajelor

Dacă traficul de mesaje este mare și toate nodurile sunt participante la decizia de acces la obiect, ceasurile logice vor fi sincronizate foarte strâns. Ca rezultat, regula de prioritate a proceselor va tinde să favorizeze nodurile cu numere mici. Acest lucru poate fi evitat prin utilizarea unei creșteri aleatoare a numărului de secvență; și mai bine, se poate introduce deliberat o nouă prioritate dacă se permite nodurilor cu prioritate mare să utilizeze creșteri mici ale numerelor de secvență, iar nodurile de prioritate mică să utilizeze creșteri mari.

4.3.2.6. Implementarea algoritmilor distribuți pentru excludere reciprocă în rețele Unix, MS Windows (Win32) și Novell NetWare

În această secțiune se prezintă câteva soluții de implementare a algoritmilor Lamport, Ricard & Agrawala, și Raymond în rețele locale Novell NetWare 3.11, Unix și Microsoft Windows (utilizând funcții Win32).

A. Implementarea în rețele Novell NetWare.

Implementarea unui program demonstrativ care înglobează atât serverul cât și clientul din descrierea algoritmilor SR este prezentată în anexa B. Se folosește pentru comunicație protocolul IPX și facilitățile de a se putea primi și trata mesaje asincron față de procesul în curs de execuție, utilizând rutine de tratare a întreruperilor declarate Event Service Routine (ESR). Acestea vor recepționa mesajele REQUEST, REPLY și RELEASE.

Implementarea pune la dispoziția programelor utilizator, într-o bibliotecă, următoarele rutine, pentru a căror utilizare trebuie inclus header-ul "exc.h":

(1) `int ExclAccessInitialize(void)` pentru inițializarea blocului de control al operațiilor de recepție a mesajelor, deschiderea socketurilor pentru emisie și recepție, instalarea rutinei ESR și intrarea în ascultare de pachete pentru recepție. Funcția returnează un cod de eroare adecvat dacă inițializările nu pot fi efectuate. Eroarea IPX care provoacă eșuarea inițializărilor este indicată în variabila externă *code*.

(2) `void ExclAccessOpen(void)` pentru obținerea accesului într-o secțiune critică, conform algoritmului. Rutina apelează două funcții ale sistemului, `Request()` pentru transmiterea în rețea a cererii de intrare în secțiune și `Enter()` care testează în buclă continuă valoarea unei variabile setată de rutina ESR atunci când sunt îndeplinite toate condițiile de intrare în secțiune; pentru algoritmul Ricard & Agrawala condiția de intrare în secțiunea critică este ca toate nodurile să fi trimis mesaje REPLY (*replyMesNo=0*), pentru algoritmul Raymond numărul de noduri care au trimis mesaje REPLY trebuie să fie mai mare sau egal cu numărul de noduri minus numărul de obiecte disponibile iar pentru algoritmul Lamport condiția este ca în lista de cereri de intrare, ordonată în raport de mărcile de timp, cererea proprie să fie prima și să fi fost confirmată de toate celelalte noduri.

(3) `void ExclAccessClose(void)` pentru ieșirea dintr-o secțiune critică, când trebuie trimise mesaje REPLY (în cazul algoritmilor Ricard & Agrawala și Raymond) sau mesaje RELEASE (în cazul algoritmului Lamport). Aceste mesaje pot provoca deblocarea altor clienți care așteaptă la intrarea într-o secțiune critică.

(4) `void ExclAccessEnd(void)` pentru deconectarea apelantului din sistem: închiderea socketurilor și întreruperea ascultării de pachete pentru recepție.

(5) `int SetUsersAndObjects (char users[MAX_CONJ|USER_LEN], int hosts, int noObjects)` pentru inițializarea structurilor de date ce trebuie menținute pentru fiecare nod: numărul conexiunii cu serverul Novell și adresa Ethernet. Funcția primește ca parametri de intrare numele de user sub care s-a făcut conectarea în fiecare nod care participă la algoritmul, numărul de noduri și, pentru algoritmul Raymond, numărul de obiecte disponibile (pentru restul algoritmilor, acest parametru se transmite întotdeauna 1) și returnează un cod de eroare care reprezintă numărul primului user din lista *users* care nu este conectat în rețea și al cărui număr de conexiune, evident, nu poate fi determinat.

Pentru demonstrarea corectitudinii tuturor algoritmilor pentru excludere reciprocă a fost scris un singur program utilizator general care permite alegerea algoritmului (anexa B).

B. Implementare Unix și Microsoft Windows

Sistemele de operare Unix și Microsoft Windows (Windows NT, Windows 95) permit o implementare mai apropiată de descrierea SR a algoritmilor prezentată în secțiunile 4.3.2.1., 4.3.2.2. și 4.3.2.3. Protocolul utilizat, având în vedere că se lucrează în rețele locale, este UDP.

Unul dintre scopurile implementării a fost ca ea să fie valabilă, cu modificări minime, atât în sisteme Unix cât și Microsoft Windows (deci în rețele mixte). Pentru toți algoritmi de excludere reciprocă a fost realizat un singur program client demonstrativ (aplicație XWindows, respectiv Windows) a cărui funcție este de a transmite cererile de intrare și ieșire din secțiunea critică ale utilizatorului, prin intermediul unui socket UDP local (canalul *localChan* din descrierile algoritmilor) la procesul server. Serverul și clientul din descrierea algoritmilor sunt implementate ca procese separate.

Procesul server constă într-o buclă infinită în care așteaptă recepționarea de mesaje de la clientul local, pe soclul UDP local, sau de la un proces server de pe alt nod, pe soclul *remoteChan*. Pentru a detecta când s-au primit mesaje pe aceste socluri se folosește funcția *select()*. Pentru deblocarea unui client în așteptare la intrare într-o secțiune critică se folosește, în Windows un semafor iar în Unix, un canal *pipe*; acesta constituie de fapt implementarea canalului *enterChan*.

La dispoziția unui program client se pun următoarele funcții (anexa B):

- (1) `int ExclAccessInitialize(void)` pentru deschiderea semaforului sau *pipe*-ului folosit, determinarea adresei proprii și deschiderea soclului UDP local. Funcția returnează un cod de eroare în cazul că nu se pot face inițializările.
- (2) `void ExclAccessOpen(void)` pentru obținerea accesului în secțiunea critică. Funcția transmite pe soclul local mesajul "REQUEST" serverului și blochează clientul, în urma unui apel `Win32, WaitForSingleObject()` la semaforul folosit pentru așteptare la intrarea în secțiunea critică; deblocarea clientului se face de către procesul server, printr-un apel `Win32, ReleaseSemaphore()` când au fost îndeplinite toate condițiile pentru intrare.
- (3) `char* GetMyAddress(void)` pentru obținerea adresei IP proprii. Funcția returnează adresa unui șir la care a depus adresa IP a propriului nod.
- (4) `void ExclAccessClose(void)` apelată la ieșirea dintr-o secțiune critică.

Pentru procesele server nu a fost prevăzută o interfață grafică; ele se execută, în mod normal, în background. Pentru determinarea adreselor nodurilor care participă la algoritmul s-a adoptat următoarea soluție: serverele se lansează în execuție având ca parametri un identificator numeric al serverului și numărul total de noduri care participă la algoritmul. La începerca execuției, fiecare server transmite broadcast identificatorul său și adresa sa IP proprie și așteaptă recepționarea adreselor IP ale tuturor nodurilor pe care le memorează într-o tabelă, după care intră în buclă infinită de așteptare și tratare a mesajelor.

Există și alte soluții de implementare a serverelor; s-a preferat însă soluția care utilizează *select()* pe de o parte, datorită considerațiilor prezentate la 3.1.2. și, pe de altă parte, datorită universalității ei. Soluția este valabilă, aproape identic, atât în rețele Unix cât și în rețele Windows NT și Windows 95; singura deosebire este faptul că, în Win32 nu există funcția *getdtablesize()* obligatorie însă în Unix înainte de utilizarea lui *select()*. O altă soluție ar putea fi aceea în care prelucrarea mesajelor locale și de la distanță se face în două fire separate; un fir se deblochează la recepția unui mesaj pe soclul *localChan* iar celălalt la recepția unui mesaj pe soclul *remoteChan*.

4.3.3. Algoritmi care utilizează un token circulant pentru excludere reciprocă

O altă metodă de a obține excluderea reciprocă constă în circulația unui token între procesele sistemului. Un token este un mesaj de un tip special care se transmite în întregul sistem. Posesia tokenului permite procesului respectiv să intre în secțiunea critică. Deoarece există un singur token în sistem, numai un singur proces se poate afla la un moment dat în secțiunea critică.

4.3.3.1. Sisteme structurate în înel

A. Descrierea algoritmului

Presupunem că procesele din sistem sunt organizate logic într-o structură sub forma de inel; aceasta nu înseamnă că topologia rețelei trebuie să fie neapărat inel. Se poate utiliza, de exemplu o rețea bus (Ethernet) fără nici o ordine predefinită a proceselor și, se poate stabili, prin software, o ordine oarecare a lor, de exemplu în ordine crescătoare a adreselor fizice de rețea a nodurilor lor; ceea ce conține este ca fiecare proces să cunoască adresa vecinului sau. Când se construiește inelul, procesului 0 i se repartizează tokenul. Tokenul este transmis de la procesul i la procesul $i+1$. Când un proces primește tokenul fie (1) dacă vrea să intre în secțiunea critică, pastrează tokenul pînă la ieșirea din secțiunea critică când îl transmite nodului următor (vecinului din inel), fie (2) dacă nu vrea să intre în secțiunea critică, transmite imediat tokenul vecinului.

Excluderea reciprocă este asigurată, de fapt, de un set de procese $Server[1:n]$, cite unul asociat fiecărui proces $Client[1:]$; procesele $Server$ formează inelul și partajează tokenul. Când un proces $Server[i]$ posedă tokenul, deblochează procesul său asociat $Client[i]$, care intră astfel în secțiunea critică. La ieșirea din secțiunea critică, $Client[i]$ semnalizează procesului $Server[i]$ care transmite tokenul mai departe, procesului $Server[(i \bmod n)+1]$.

Deoarece inelul este unidirecțional, amănarea nedefinită a unui proces nu poate avea loc. Numărul de mesaje necesare pentru implementarea excluderii reciproce variază de la un mesaj pe intrare, în cazul unei concurențe mari (fiecare proces dorește intrarea în secțiunea critică) la infinit, în cazul în care nici un proces nu dorește intrarea în secțiunea critică.

Există două probleme legate de acest algoritmul: (1) dacă tokenul este pierdut, el trebuie regenerat. Detectarea pierderii tokenului nu este o operație simplă, deoarece intervalul de timp dintre două apariții succesive ale tokenului nu este limitat. Ca atare, când tokenul este pierdut este necesar un algoritmul de alegere

iar procesul ales să genereze tokenul. (2) dacă un proces (*server*) dispare, trebuie restabilit un nou inel logic. Există mai mulți algoritmi de alegere și de reconstrucție a unui inel logic ([Lan77] [Lan83] [GMO82] (vezi subcapitolul 4.4); în continuare va fi prezentat algoritmul lui Lelann).

Algoritmul lui Lelann permite detecția unui token pierdut și regenerarea lui. Se presupune că fiecare proces are un timer. Acest timer este resetat la recepția *tokenului de control* (CT) sau alt token; tokenurile circulă în sistem după regula FIFO. Protocolul propus de Lelann este un alt exemplu de algoritm distribuit: principala structură de date utilizată de algoritm este o listă cu numerele de priorități ale tuturor proceselor active, denumită *lista activă* (fiecare proces menține propria sa lista activă). Protocolul se descrie astfel:

(1) Când timpul setat de timer s-a epuizat, procesul corespunzător P_i generează un *token de alegere*: (ET(i)) care conține numărul procesului. Timerul este setat din nou; se spune că s-a intrat în *faza de alegere*. Totodată în lista activă proprie, $S(i)$, se inseric i .

(2) Când un candidat primește înapoi tokenul CT înainte de propriul sau token ET, anulează faza de alegere și înlătură tokenul ET când îl primește.

(3) Când un candidat P_i primește de la vecinul din stanga un token ET(j) atunci: (1) Dacă acesta este primul token de alegere primit sau trimis, P_i creează o listă activă nouă, cu numerele i și j și trimite apoi ET(i) urmat de ET(j); (2) Dacă $i < j$, atunci P_i adaugă j la lista activă și trimite mesajul mai departe la vecinul din dreapta.

(4) Când un candidat P_i primește propriul sau ET(i), îl înlătură din inel, resetează timerul și execută următorul algoritm: Dacă $i = \max S(i)$ atunci se generează un nou token CT și se setează timerul.

Eficiența mecanismului de sincronizare cu token circulant depinde mult de dimensiunea (în timp) a secțiunii critice. Este evident că, atunci când secțiunile critice implică schimburi de mesaje între producători și consumatori acest mecanism de sincronizare nu oferă performanțe bune. În schimb, acest mecanism a fost folosit pentru dezvoltarea unui algoritm pentru detectarea terminării unei prelucrări distribuite ([And91] [DS83] [Mis83]); circulația tokenului se folosește pentru a semnaliza modificările de stare.

B. Implementare

```

type kind = enum (REQUEST, RELEASE) # tipuri de mesaje
chan remoteChan[1:n]() # canal pentru comunicare în retea
chan localChan[1:n](kind) # canal pentru transmiterea cererilor de intrare și de ieșire din secțiunea
# critica (comunicatie locala între clienti și server)
chan enterChan[1:n]() # canal pentru deblocarea unui client în așteptare la intrare în secțiune
    
```

Client [i: 1..n]:

```

do true →
    send localChan[i] (REQUEST) # protocol de intrare
    receive enterChan[i] ()
    #secțiune critica...
    send localChan[i] (RELEASE) # protocol de ieșire
od
    
```

Server [i: 1..n]:

```

var doAccess: bool := false # indica accesul într-o secțiune critica
var kind: kind # folosita la citirea mesajelor locale
do true → # bucla infinita a severului
    if not (empty (localChan[i])) →
        receive localChan[i](kind)
        if kind = REQUEST → doAccess := true:
        || kind = RELEASE →
            doAccess := false
            send remoteChan[i mod n - 1]()
        fi
    fi
    if not (empty(remoteChan[i])) →
        receive remoteChan[i]()
        if doAccess = true → send enterChan[i]()
        || doAccess ≠ true → send remoteChan[i mod n - 1]()
        fi
    fi
od
    
```

4.3.3.2. Algoritmul pentru sisteme nestructurate în inel

În contrast cu sistemele structurate în inel, în cele nestructurate în inel un proces care deține tokenu poate să-l trimită oricărui alt proces din sistem. În plus, nu este nevoie ca tokenul să circule dacă nici un proces nu cere să intre în secțiunea critică.

Algoritmul următor aparține lui Chandy ([Gos91] [Cha82] [SP88]) și garantează că fiecare proces care cere intrarea într-o secțiune critică va primi tokenul într-un interval de timp finit față de momentul lansării cererii.

Se consideră că există n procese în sistem. Tokenului i se asociază un vector $T=(T_1, T_2, \dots, T_n)$, unde T indică numărul de intrări în secțiunea critică a procesului P_i . Algoritmul constă în următoarele etapi:

(1) Când un proces P_i cere să intre în secțiunea critică pentru a m_i -a oară, trimite mesajul REQUEST(P_i, m_i) către toate celelalte procese și așteaptă apoi să primească tokenul.

(2) Când procesul P_i primește tokenul, setează $T_i = m_i$ în vectorul T al tokenului și intră în secțiunea critică.

(3) După ce procesul P_i iese din secțiunea critică, examinează coada de cereri de acces primite:

- (a) Dacă coada este vidă, procesul P_i continuă execuția normală (pînă primește mesaje de la alti procese)
- (b) Dacă coada nu este vidă, procesul P_i preia prima cerere din coadă, fie aceasta (P_j, m_j) Coada este ordonată FIFO.
- (c) Dacă $m_j > T_j$ atunci P_i trimite tokenul procesului P_j ; în caz contrar, elimină cererea, ca fiind o cerere veche, deja deservită și preia din coadă următoarea cerere, pe care o prelucurează ca mai sus.
- (d) Aceasta prelucrare continuă pînă cînd coada este vidă sau tokenul este trimis altui proces.

```

type kind = enum (REQUEST, TOKEN, RELEASE)      # tipuri de mesaje
chan localChan[1:n](kind)                       # canal pentru transmiterea cererilor de intrare si de iesire din sectiunea
                                                  # critica (comunicatie locala intre clienti si server)
chan enterChan[1:n]()                           # canal pentru deblocarea unui client in asteptare la intrare in sectiune
chan remoteChan[1:n](kind, sender: int, case kind of
    REQUEST: (accesses: int)
    TOKEN:   (T[1..n]: int)
) #canal pentru comunicatie in retea
    
```

Client [i: 1..n]:

```

do true →
    send localChan[i] (REQUEST)  # protocol de intrare
    receive enterChan[i] ( )
    #sectiune critica...
    send localChan[i] (RELEASE)  # protocol de iesire
od
    
```

Server [i: 1..n]:

```

var inAccess: bool := false      # indica accesul intr-o sectiune critica
var tokenOwner: bool := false   # indica detinerea tokenului
var reqQueue[1:n]: int := (|n| 0) # tablou care inregistreza cererile proceselor
var T[1:n]: int := (|n| 0)      # token transmis in sistem
var kind: kind, sender: int     # folosite la citirea mesajelor
var accesses: int := 0          # folosit pentru a contoriza cererile pentru sectiune critica
    
```

```

do true → # bucla infinita a serverului
    if not (empty (localChan[i])) →
        receive localChan[i](kind)
        if kind = REQUEST →
            accesses := accesses + 1
            req.Access := true;
            for j := 1 to n →
                send remoteChan[j](REQUEST, i, accesses)
            af
        || kind = RELEASE →
            in.Access := false
            reqQueue[i] := 0
    
```

```

fa j := 1 to n st reqQueue[j] ≠ 0 and (tokenOwner = true) →
    if reqQueue[j] > T[j] →
        send remoteChan[j](TOKEN, i, T)
        tokenOwner := false
    [] reqQueue[j] := 0
    if
af
fi
fi
if not (empty(remoteChan[i])) →
    receive remoteChan[i](knd, sender)
    if knd = REQUEST →
        receive remoteChan[i](reqQueue[sender])
        if (tokenOwner := true and inAccess = false) →
            if (reqQueue[i] ≠ 0 →
                inAccess := true; T[i] := reqQueue[i]
                send enterChan[i]()
            [] reqQueue[i] = 0 →
                fa j := 1 to n st reqQueue[j] ≠ 0 and (tokenOwner = true) →
                    send remoteChan[j](TOKEN, i, T)
                    tokenOwner := false
                af
            fi
        [] knd = TOKEN →
            receive remoteChan[i](T)
            T[i] := reqQueue[i]
            inAccess := true; tokenOwner := true
        fi
    fi
od

```

4.3.3.3. Secvențiator circulant

A. Descrierea algoritmului

O altă metodă bazată pe circulația dreptului de acces folosește un secvențiator circulant ([Lan83] [Gos91]). Un *secvențiator circulant* este un secvențiator care circulă permanent în inelul logic al proceselor, via unui mecanism de token circulant (token de control).

Modalitatea de a obține excluderea reciprocă de procesele care partajează secvențiatorul este următoarea: (1) La primirea tokenului, un proces poate executa câteva primitive *ticket*; (2) Apoi trimite tokenul la succesorul său în inel.

Dacă se consideră p procese și n acțiuni în așteptare pînă ce procesul p primește tokenul, strategia constă în a executa exact n consecutive operații *ticket* deci a obține n tickete consecutive. În unele cazuri însă, necesitatea de a aștepta tokenul pentru a executa o acțiune poate fi considerată inacceptabilă. Pentru aceste cazuri, Lelann prezintă diferiți algoritmi care permit selecția anticipată a ticketelor. În funcție de strategia de selecție a ticketelor, procesele consumatoare ale ticketelor pot utiliza relația de ordine stabilită prin numerele ticketelor fie pentru a evita planificării inconsistente, fie pentru a detecta apariția unor planificări inconsistente și a lua decizia necesară. De exemplu, conflictele între acțiunile corespunzătoare a două sau mai multor tranzacții asupra unor copii ale aceluiași obiect (fișiere replicate, baze de date partiționate) pot fi rezolvate în faza de commit a tranzacțiilor, deci după ce acțiunilor li s-a permis execuția, pe baza priorităților date de numerele ticketelor.

B. Implementare

```

type kind = enum (REQUEST, TICKET, RELEASE) # tipuri de mesaje
chan remoteChan[1:n](seqValue: int) # canal pentru comunicare in retea
chan localChan[1:n](kind) # canal pentru transmiterea cererilor
chan enterChan[1:n]() # canal pentru deblocarea unui client in asteptare la intrare in sectiune
chan answerChan[1:n](seqValue: int) # canal pentru obtinerea unui ticket
Client [i: 1..n]:

```



```

var ticket: int = 0          #valoarea ticketului primita de la server
do true →
    send localChan[i] (REQUEST) # protocol de intrare
    receive enterChan[i] ( )
    #sectiune critica...
    send localChan[i] (TICKET) # protocol de obtinere a ticketului
    receive answer[i] (ticket)
    send localChan[i] (RELEASE) # protocol de iesire
od
    
```

Server [i: 1..n]:

```

var doAccess: bool := false # indica accesul intr-o sectiune critica
var kind: kind # folosita la citirea mesajelor locale
var ticketValue: int = 0 # valoarea curenta a ticketului

do true -> # bucla infinita a serverului
    if not (empty (localChan[i])) →
        receive localChan[i](kind)
        if kind = REQUEST → doAccess := true;
        || if kind = TICKET → send answerChan[i](ticketValue)
        || kind = RELEASE →
            doAccess := false
            send remoteChan[i mod n + 1](ticketValue)
        fi
    fi
    if not (empty(remoteChan[i])) →
        receive remoteChan[i](ticketValue)
        if doAccess = true → send enterChan[i]( )
        || doAccess ≠ true → send remoteChan[i mod n + 1]( )
        fi
    fi
od
    
```

4.3.3.4. Implementarea algoritmilor care utilizează un token circulant în rețele Unix, MS Windows (Win32) și Novell NetWare

A. Implementare Novell NetWare.

Se utilizează același program demonstrativ prezentat la 4.3.2.6. Analog, implementarea pune la dispoziția programelor client următoarele rutine (anexa B3):

(1) `int ExclAccessInitialize(void)` care are aceleași funcții ca și în 4.3.2.6., adică deschiderea socketurilor pentru emisie și recepție, inițializarea blocului de control pentru recepția mesajelor, instalarea rutinei ESR, intrarea în ascultare de pachete pentru recepție, dar, în plus, are și funcția de a genera token-ul în nodul cu identificatorul cel mai mic. Funcția returnează un cod de eroare sau zero în cazul execuției cu succes.

(2) `void ExclAccessOpen(void)` cu rolul de a obține accesul în secțiunea critică, atunci când este primit token-ul.

(3) `void ExclAccessClose(void)` apelată la ieșirea dintr-o secțiune critică, pentru transmiterea token-ului.

(4) `void ExclAccessEnd(void)` apelată pentru deconectarea apelantului din sistem: închiderea socketurilor și întreruperea ascultării de pachete pentru recepție.

(5) `int SetUsersAndObjects (char users[MAX_CON][USER_LEN], int hosts, int noObjects)` pentru inițializarea structurilor de date ce trebuie menținute pentru fiecare nod: numărul conexiunii cu serverul Novell și adresa Ethernet. Tabela numerelor de conexiune a nodurilor care participă la algoritmul este ordonată crescător. În cazul algoritmului în incl, aceasta tabelă se utilizează pentru a determina adresa nodului vecin, iar în cazul algoritmului Chandy, aceasta tabelă se utilizează pentru a trimite token-ul la nodul cu numărul de conexiune cel mai mic dintre toate nodurile care așteaptă la intrarea într-o secțiune critică. În cazul algoritmului Chandy se asociază astfel o prioritate fiecărui nod.

B. Implementare Unix și Windows.

În cazul rețelelor Unix și Windows se utilizează același program client demonstrativ prezentat la 4.3.2.6. Spre deosebire de varianta Novell, clientul și serverul se înglobează însă în procese diferite. Canalul *enterChan* se implementază în Windows sub forma unui semafor inițializat la zero; după lansarea unei cereri de intrare în secțiunea critică, un client așteaptă blocat la acest semafor. În Unix, blocarea clienților se realizează prin intermediul unui *pipe*. În cazul secvențiatorului circulant, clientul trebuie să primească valoarea lui. În Windows, comunicarea server-client s-a realizat printr-o zonă de memorie partajată iar în Unix printr-un *pipe*. Implementarea pune la dispoziția programelor client următoarele funcții (anexa B):

La dispoziția unui program client se pun următoarele funcții:

- (1) `int ExclAccessInitialize(void)` pentru inițializare: deschiderea semaforului sau *pipe*-ului pentru sincronizarea client-server, determinarea adresei proprii și deschiderea socketului UDP local.
- (2) `void ExcAccessOpen(void)` pentru obținerea accesului la secțiunea critică.
- (3) `char* GetMyAddress(void)` pentru obținerea adresei IP proprii. Funcția returnează adresa unui șir la care a depus adresa IP a propriului nod.
- (4) `void ExclAccessClose(void)` apelată la ieșirea dintr-o secțiune critică.

Tabela adreselor IP ale nodurilor se ordonează în raport de valoarea identificatorilor numerici, pentru a se putea determina adresa nodului vecin, în cazul algoritmului în inel sau pentru a se obține prioritățile alocate nodurilor în cazul algoritmului Chandy.

4.3.4. Algoritmi broadcast

În unele tipuri de rețele locale, nodurile partajează un canal de comunicație comun. În aceste cazuri, se poate considera că fiecare nod este conectat direct cu celelalte. Deseori, în aceste tipuri de rețele, suportul software oferă o primitivă specială, *broadcast* care transmite un mesaj de la un nod la toate celelalte.

Fie $P[1..n]$ n procese și $channel[1..n]$ n canale de comunicație, câte unul pentru fiecare proces. Execuția unei primitive *broadcast* $channel(mesaj)$ are ca efect plasarea câte unei copii a mesajului pe fiecare canal $channel[i]$, asociat procesului $P[i]$; acest efect este echivalent cu execuția în paralel a n primitive *send*, fiecare trimițând un mesaj pe unul dintre cele n canale. Procesul i recepționează mesajul transmis, în urma execuției primitive *receive* $channel[i]$.

Primitiva *broadcast* poate fi exploatată în proiectarea unor algoritmi de sincronizare a proceselor în sisteme distribuite, toți bazându-se pe utilizarea primitivei *broadcast*. Problema centrală a acestor algoritmi o constituie realizarea ordonării totale a evenimentelor. Aceasta va fi rezolvată folosind ceasurile logice astfel: (1) fie LC valoarea timpului logic curent în procesul A ; când A trimite un mesaj (nu neaparat *broadcast*) se înglobează în mesaj și LC -ul curent, după care se incrementează LC ($LC:=LC+1$); (2) când A recepționează un mesaj cu marca de timp TS de la un proces B , se setează timpul logic curent LC la $\max(LC, TS)$ și, apoi se incrementează LC .

Algoritmii *broadcast* sunt adecvați pentru sincronizarea unui număr mic de procese, având în vedere că toate procesele trebuie să proceseze fiecare mesaj.

4.3.4.1. Semafoare distribuite

Semafoarele se implementează, în mod normal, în sistemele centralizate, fie utilizând variabile partajate, fie, în sistemele bazate pe transmisia de mesaje, utilizându-se un proces server (cu rol de monitor). În subcapitolul de față, se prezintă o implementare a semafoarelor în sisteme distribuite la baza căreia stă utilizarea primitivei *broadcast* [And91].

Un semafor s este un tip de date abstracte, care poate fi accesat numai prin intermediul a două operații primitive: *Down*(s) (sau *P*(s)) și *Up*(s) (sau *V*(s)). Aceste operații se sincronizează astfel încât întotdeauna se menține un *invariant* pentru semafoare: valoarea dintre numărul de operații *Down* complet executate este cel mult numărul de operații *Up* complet executate plus valoarea inițială a semaforului.

Într-o implementare a semafoarelor utilizând variabile partajate, s se reprezintă printr-o variabilă de tip întreg, negativă. Execuția unei operații *Up*(s) incrementează aceasta variabila printr-o acțiune atomică; execuția unei operații *Down*(s) este întârziată pînă cînd variabila s are o valoare pozitivă, după care se decrementează, tot printr-o acțiune atomică. Într-un sistem distribuit, este însă necesară o altă tehnică pentru a reprezenta valoarea semaforului și pentru a menține invariantul pentru semafoare. În principal este necesară contorizarea operațiilor *Down*(s) și *Up*(s) și întârzierea execuției operațiilor *Down*(s); în plus, procesele care partajează un semafor trebuie să coopereze astfel încît să păstreze invarianța semaforului chiar dacă starea programului este distribuită.

Cerințele de mai sus pot fi realizate dacă se impune ca, înainte de execuția unei operații *Down* sau *Up*, fiecare proces să transmită această intenție printr-un mesaj *broadcast* și să examineze răspunsurile primite pentru a determina dacă execuția poate continua sau nu. Fiecare proces trebuie să mențină o coadă de mesaje, *mesQueue*, și un ceas logic *logicClock*. Pentru a simula execuția unei operații *Down* sau *Up* fiecare proces

trebuie să transmită broadcast un mesaj la toate nodurile, inclusiv lui însuși; mesajul conține identitatea emițătorului, tipul operației (Down sau Up) și marca de timp, actualizată conform regulilor prezentate la începutul acestui capitol. Când un proces recepționează un mesaj Down sau Up, îl memorează în coada de mesaje proprie. Aceasta coadă se ordonează crescător, în raport de mărcile de timp conținute în mesaje iar pentru mărci de timp egale se folosește identitatea emițătorului. Deoarece operația de transmisie broadcast nu este atomică, există posibilitatea ca două mesaje broadcast emise de două procese diferite să fie recepționate în ordine diferită de celelalte procese. Mesajele broadcast emise însă de același proces vor fi recepționate de celelalte procese întotdeauna în ordinea crescătoare a mărcilor lor de timp; aceasta oferă posibilitatea luării unor decizii corecte pentru sincronizare. Astfel, presupunând că procesul conține în *mesQueue* un mesaj *m* cu marca de timp *TS*, atunci după ce procesul *i* a recepționat un mesaj cu o marcă de timp mai mare decât *TS*, de la toate celelalte noduri, este sigur că nu va mai recepționa un mesaj cu o marcă de timp mai mică; se spune că mesajul *m* este **confirmat complet**. În acest caz, orice alt mesaj din coadă, situat înaintea lui *m* va fi, de asemenea confirmat complet deoarece are o marcă de timp mai mică decât *m*. Se furnează astfel în coada *mesQueue* un **prefix stabil**: nici un alt mesaj nu va mai fi inserat în această parte a cozii.

Având în vedere observațiile de mai sus, se poate descrie acum algoritmul de implementare a semafoarelor distribuite:

(1) la execuția unei operații Up sau Down, un proces transmite un mesaj broadcast care conține identitatea sa, tipul operației și valoarea ceasului logic local.

(2) la recepționarea unui mesaj Down sau Up, un proces transmite broadcast un mesaj de confirmare (ACK). Confirmările se transmit broadcast pentru a putea fi văzute de toate celelalte noduri. Un mesaj de confirmare conține marca de timp locală, dar, spre deosebire de mesajele Down și Up, nu se memorează în *mesQueue*; scopul lor este numai de a determina când un mesaj Up sau Down este complet confirmat.

(3) fiecare proces simulează execuția operațiilor Down și Up, determinate de mesajele situate în prefixul stabil al cozii proprii *mesQueue* astfel: fiecare proces menține doi contori locali *nUp* și *nDown* care contorizează numărul de mesaje Up, respectiv Down confirmate complet, respectând invarianța semaforului. Când un mesaj Up devine complet confirmat, se incrementează contorul *nUp*. Când un mesaj Down devine complet confirmat, se incrementează *nDown* și, dacă *nUp > nDown*, se deblochează apelantul. După incrementarea contoarelor, procesul poate șterge mesajul respectiv, Up sau Down, din coada proprie, *mesQueue*.

Nodurile sistemului distribuit se pot afla în faze diferite de prelucrare ale mesajelor Up și Down, deoarece mesajele pot să fie complet confirmate în ordine diferită, dar fiecare nod va prelucra mesajele complet confirmate în aceeași ordine.

Descrierea algoritmului în notația SR este următoarea:

```

type kind    enum (DOWN, UP, ACK)    # tipuri de mesaje
chan remoteChan[1:n](kind, sender: int, timeStamp: int, sem: int) #canal broadcast
chan localChan[1:n](kind, sem: int)    # canal pentru transmiterea cererilor
chan answerChan[1:n](sem\al: int)      # canal pentru obtinerea valorii contorului
    
```

Client [*i*: 1..*n*]:

```

var sem: int           # pentru a indica numarul semaforului
var value: int        # pentru a primi valoarea semaforului de la server

do true →
    ...
    send localChan[i] (UP, sem)           # incrementare semafor
    receive answerChan[i] (value)        # citire valoare semafor
    send localChan[i] (DOWN, sem)        # sincronizare
    receive answerChan[i] (value)        # citire valoare semafor

od
    
```

Server [*i*: 1..*n*]:

```

var localTimeStamp: int = 0           # marca de timp ce va fi asignata urmatoarei cereri locale
var nUp[1:n]: int = (|n| 0)           # numarul de operatii UP pentru semafoare
var nDown[1:n]: int = (|n| 0)        # numarul de operatii DOWN pentru semafoare
var kind: kind, timeStamp: int       # folosite la citirea mesajelor
var sender: int, sem: int           # folosite la citirea mesajelor
var mesQueue: queue of (sender: int, kind: kind, ts: int, sem: int)

do true →                               # bucla infinita a severului
    if not (empty (localChan[i])) →
    
```

```

receive localChan[i](knd, sem)
broadcast remoteChan (knd, i, localTimeStamp, sem )
localTimeStamp := localTimeStamp + 1
fi
if not (empty(remoteChan[i])) →
receive remoteChan[i](knd, sender, timeStamp, sem )
localTimeStamp := max(localTimeStamp, timeStamp +1)
localTimeStamp := localTimeStamp + 1
if knd = UP or knd= DOWN→
    insereaza in coada mesQueue cererea (sender, kind, timeStamp, sem)
    broadcast remoteChan[i](ACK, i, timeStamp, evCounter)
    localTimeStamp := localTimeStamp + 1
    [] knd = ACK
        inregistreaza confirmarea primita
        fa cererile UP complet confirmate →
            inregistreaza identitatea semaforului in sem si a emitorului in sender
            sterge cererea din mesQueue
            nUp[sem] := nUp[sem] +1
            if sender = i → send answerChan[i](nUp[sem] - nDown[sem])
        af
        fa (cererile DOWN complet confirmate) and (nUp[sem] > nDown[sem]) →
            inregistreaza identitatea semaforului in sem si a emitorului in sender
            sterge cererea din mesQueue
            nDown[sem] := nDown[sem] +1
            if sender = i → send answerChan[i](nUp[sem] - nDown[sem])
        af
    fi
fi
od

```

Există câte un proces *Client* și *Server* în fiecare nod. Procesele *Client* inițiază operațiile Up și Down transmișând câte un mesaj serverelor lor locale. Procesele server implementează operațiile Up și Down. Ele recepționează mesajele transmise pe *remoteChan*, gestionează cozile de mesaje, transmit broadcast cererile de operații Up sau Down și confirmările, deblochează procesele *Client* în așteptare la o operație Down și gestionează ceasul logic local care furnizează mărcile de timp pentru mesaje.

Semafoarele distribuite pot fi utilizate pentru sincronizarea proceselor în programare distribuită într-un mod asemănător cu utilizarea lor în sistemele centralizate. De exemplu, semafoarele distribuite pot fi utilizate pentru realizarea excluderii reciproce în accesul la fișiere, baze de date, pentru a coordona accesul la fișiere replicate. În anexa C se prezintă un exemplu de aplicație implementată în rețele Unix și Microsoft Windows bazată pe semafoarele distribuite.

4.3.4.2. Contori de evenimente distribuiti

O implementare a contorilor de evenimente poate fi obținută într-un sistem distribuit și folosind mesaje broadcast. Ca și implementarea semafoarelor în sisteme distribuite, prezentată la 4.3.4.1., această soluție de implementare a contorilor de evenimente este complet distribuită. Ideea de bază este de a transmite mesaje broadcast pentru fiecare operație *Advance*; operațiile *Await* și *Read* sunt procesate local. Operațiile *Advance* trebuie confirmate prin mesaje broadcast, de către fiecare nod. Ca și implementarea semafoarelor distribuite, fiecare proces menține local o coadă de mesaje incomplet confirmate, *mesQueue*, ordonată crescător în raport de mărcile de timp ale mesajelor și valoarea contorilor.

Algoritmul de implementare a contorilor de evenimente distribuiți poate fi descris astfel:

- (1) la execuția unei operații *Advance*, un proces transmite un mesaj broadcast care conține identitatea sa și valoarea ceasului logic local.
- (2) la recepționarea unui mesaj *Advance*, un proces transmite un mesaj de confirmare (ACK), care conține valoarea ceasului logic local. Scopul mesajelor de confirmare este de a se putea determina când un mesaj *Advance* este complet confirmat.
- (3) când un mesaj *Advance* (plasat în *mesQueue*) devine complet confirmat, se incrementează contorul și se elimină mesajul din coadă.

(4) Operațiile *Read* și *Await* se procesează local, în nodul respectiv. La execuția operației *Read(contor)* se returnează valoarea contorului din momentul respectiv; la execuția operației *Await(contor, waitValue)* se deblochează apelantul numai când valoarea contorului este mai mare sau egală cu argumentul *waitValue*.

Descrierea algoritmului în notație SR este următoarea:

```

type kind = enum (ADVANCE, READ, AWAIT, ACK) # tipuri de mesaje
chan remoteChan[1:n](kind, sender: int, timeStamp: int, counter: int) #canal broadcast
chan localChan[1:n](kind, counter: int, case kind of
    ADVANCE, READ: ()
    AWAIT: waitValue: int
) # canal pentru transmiterea cererilor
chan answerChan[1:n](counterValue: int) # canal pentru obtinerea valorii contorului

Client [i: 1..n]:
var counter: int # pentru a indica numarul contorului
var value: int # pentru a primi valoarea contorului de la server
var waitValue: int # pentru a indica valoarea pentru AWAIT
do true →
    ...
    send localChan[i] (ADVANCE, counter) # incrementare contor
    send localChan[i] (READ, counter) # citire valoare contor
    receive answerChan[i] (value)
    send localChan[i] (AWAIT, counter, waitValue) # sincronizare
od

Server [i: 1..n]:
var localTimeStamp: int = 0 # marca de timp ce va fi asignata urmatoarei cereri locale
var evCounterValues[1:n]: int = {n} 0 # valoarea curenta a contorilor
var evWaitValues[1:n]: int = {n} -1 # valoarea curenta a parametrilor transmisi pentru sincronizare
var kind: kind, timeStamp: int # folosite la citirea mesajelor
var sender: int, evCounter: int # folosite la citirea mesajelor
var mesQueue: queue of (sender: int, ts: int, counter: int) #coada de cereri incomplet confirmate

do true → # bucla infinita a severului
    if not (empty (localChan[i])) →
        receive localChan[i](kind, evCounter)
        if kind = ADVANCE →
            broadcast remoteChan(ADVANCE, i, localTimeStamp, evCounter)
            localTimeStamp := localTimeStamp + 1;
        if kind = READ →
            send answer[i](evCounterValues[evCounter])
        if kind = AWAIT →
            receive localChan[i](waitValue)
            evWaitValues[evCounter] := waitValue
            if → evCounterValues[evCounter] ≥ evWaitValues[evCounter]
                evWaitValues[evCounter] := -1
                send enterChan[i]()
            fi
        fi
    fi
    if not (empty(remoteChan[i])) →
        receive remoteChan[i](kind, sender, timeStamp, evCounter)
        if kind = ADVANCE →
            insereaza in coada mesQueue cererea (sender, kind, evCounter)
            localTimeStamp := max(localTimeStamp, timeStamp + 1)
            localTimeStamp := localTimeStamp + 1
            broadcast remoteChan[i]( ACK, i, timeStamp, evCounter)

```

```

    localTimeStamp = localTimeStamp + 1
[] knđ = ACK
    inregistreaza confirmarea primita
    fa cererile ADVANCE complet confirmate →
        inregistreaza identitatea contorului in evCounter
        sterge cererea din mesQueue
        evCountValues[evCounter] := evCountValues[evCounter] + 1
        if evWaitValues[evCounter] > 0 and
            ( evCountValues[evCounter] ≥ evWaitValues[evCounter] ) →
                evWaitValues[evCounter] := -1
                send enterChan[i]()
        fi
    fi
fi
od

```

Există câte un proces *Client* și *Server* în fiecare nod. Procesele *Client* comunică numai cu procesele *Server* locale prin intermediul canalului *localChan*. Procesele *Server* implementează operațiile *Advance*, *Read* și *Await*. Ele recepționează mesaje pe canalul *localChan* (de la user-ul asociat) și pe canalul *remoteChan* (de la celelalte servere) și asigură transmitia broadcast a mesajelor *Advance*, recepția mesajelor broadcast (*Advance* și de confirmare), gestiunea cozilor de mesaje și deblocarea proceselor *Client* în așteptare la o operație *Await*.

4.3.4.3. Secvențiatori distribuți.

Implementarea complet distribuită a secvențiatorilor poate fi făcută analog semafoarelor distribuite și contorilor de eveniment distribuiti, prin utilizarea mesajelor broadcast. În cazul secvențiatorilor există o singură operație (primitivă) care trebuie implementată: primitiva *Ticket()*, care returnează valoarea curentă a secvențiatorului.

Implementarea presupune și în acest caz două procese în fiecare nod: procesul *Client* care transmite cererea de execuție a unei operații *Ticket* la procesul *Server* corespunzător, prin intermediul unui canal local. Operația *Ticket* este transmisă apoi mai departe, la procesele *Server* din toate nodurile, prin transmitia unui mesaj broadcast. Incrementarea valorii secvențiatorului se face numai pentru operațiile *Ticket* complet confirmate.

- Descrierea algoritmului în notatia SR este următoarea:

```

type kind = enum (TICKET, ACK) # tipuri de mesaje
chan remoteChan[1:n](kind, sender: int, timeStamp: int) # canal broadcast pentru comunicare in retea
chan localChan[1:n](kind) # canal pentru transmiterea cererilor de obtinere a unui ticket
chan answerChan[1:n](seqValue: int) # canal pentru obtinerea valorii ticketului

Client [i: 1..n]:
var value: int
do true →
    send localChan[i] (TICKET) # protocol de obtinere a unui ticket
    receive answerChan[i] (value)
od

Server [i: 1..n]:
var localTimeStamp: int = 0 # marca de timp ce va fi asignata urmatoarei cereri locale
var ticketValue: int = 0 # valoarea curenta a ticketului
var knđ: kind, timeStamp: int, sender: int, ts: int # folosite la citirea mesajelor
var mesQueue: queue of (sender: int, ts: int) # coada de mesaje incomplet confirmate

do true ->
    if not (empty (localChan[i])) →
        receive localChan[i](knđ)
        if knđ = TICKET →

```

```

broadcast remoteChan(TICKET, i, localTimeStamp )
    localTimeStamp := localTimeStamp + 1;
fi
fi
if not (empty(remoteChan[i])) →
    receive remoteChan[i](knd, sender, timeStamp )
    if knd = TICKET →
        inscreaza in mesQueue cererea
        localTimeStamp := max(localTimeStamp, timeStamp +1)
        localTimeStamp := localTimeStamp + 1
        broadcast remoteChan[i]( ACK, i, timeStamp)
        localTimeStamp = localTimeStamp + 1
    || knd = ACK
        inregistreaza confirmarea primita
    fa cererile TICKET complet confirmate →
        sterge cererea din coada
        ticketValue := ticketValue + 1
        if sender = i → send answer[i](ticketValue) fi
    fa
fi
fi
od

```

4.3.4.3. Implementarea algoritmilor broadcast în rețele Unix și Microsoft Windows (Win32)

Implementarea canalului broadcast (*remoteChan* din descrierea SR a algoritmilor) s-a realizat prin setarea unui soclu UDP ca soclu broadcast, utilizându-se *setsockopt()*. Procesele server recepționează mesaje pe două socluri: *localChan*, soclu UDP local, pe care se primesc cereri de la clienții locali și *remoteChan*, soclu UDP pe care se primesc mesajele broadcast trimise de celelalte servere. Toate implementările (pentru semafoare distribuite, contori de evenimente distribuiți și secvențiatori distribuiți) utilizează o listă de mesaje, *mesQueue* în care se memorează cererile clienților incomplet confirmate, ordonată după valoarea mărcilor de timp ale mesajelor; în momentul în care o cerere a unei operații devine complet confirmată (constatare care se face examinând mesajele de confirmare primite pentru fiecare cerere), cererea respectivă este scoasă din lista.

În cazul algoritmului pentru semafoare distribuite s-a considerat necesar ca, în urma unei operații Up sau Down, serverul să returneze clientului și valoarea asociată semaforului respectiv. Comunicarea între server și client s-a realizat în Windows printr-o zonă de memorie partajată iar în Unix printr-un *pipe*. Sincronizarea dintre client și server s-a realizat în Windows printr-un semafor iar în Unix prin același pipe prin care se realizează și comunicarea valorii semaforului. Analog se realizează comunicarea și sincronizarea server-client și în cazul contorilor de eveniment distribuiți (prin operațiile *Read* și *Wait*) și în cazul secvențiatorilor distribuiți (operația *Ticket*).

În cazul algoritmului pentru semafoare distribuite și a contorilor de eveniment distribuiți s-a implementat un tablou de semafoare, respectiv un tablou de contori, astfel că fiecare cerere de operație va avea ca parametru și numărul semaforului, respectiv numărul contorului.

La dispoziția unui program client sunt puse următoarele funcții de bibliotecă:

A. pentru semafoarele distribuite:

1. **int SemaphoreInitialize(void)** pentru crearea și atașarea la zona de memorie partajată, deschiderea semafoarelor, deschiderea soclurilor. Funcția returnează un cod de eroare.
2. **int SemaphoreUp(int sem)** pentru execuția unei operații Up asupra semaforului *sem*. Funcția returnează valoarea semaforului sau -1 dacă *sem* nu este un semafor valid.
3. **int SemaphoreDown(int sem)** pentru execuția unei operații Down asupra semaforului *sem*. Funcția returnează valoarea semaforului din momentul deblocării clientului sau -1 dacă *sem* nu este un semafor valid.
4. **void SemaphoreEnd(void)** pentru deconectarea de la sistem.

B. pentru contorii de eveniment distribuiți:

1. **int EventCountersInitialize(void)** pentru crearea și atașarea la zona de memorie partajată deschiderea semafoarelor, deschiderea soclurilor. Funcția returnează un cod de eroare.
2. **int EventCountersRead(int counter)** pentru determinarea valorii contorului de eveniment indicat de parametru sau -1 dacă *counter* nu este valid.
3. **int EventCountersAdvance(int counter)** pentru execuția unei operații Advance(counter).

4. `int EventCoutersAwait(int couter, int val)` pentru execuția unei operații `Await(couter, val)`.
5. `void EventCoutersEnd(void)` pentru deconectarea de la sistem.

C. pentru secvențiatorii distribuți:

1. `int SequencerInitialize(void)` pentru inițializari.
2. `void SequencerTicket(void)` pentru obținerea unui ticket.
3. `void SequencerEnd(void)` pentru deconectarea de la sistem.

Ca exemple, au fost elaborate șase programe client care ilustrează, interogativ, prin intermediul unor interfețe grafice `XWindows`, respectiv `Windows`, modul de lucru cu semnalele distribuite, contorii de eveniment distribuți și secvențiatorii. Un alt program client a fost scris pentru a demonstra folosirea atât a contorilor de eveniment cât și a secvențiatorilor într-un singur program, care controlează excluderea reciprocă folosind contori și secvențiatori (anexa B).

4.4. Algoritmi de alegere

Numeroși algoritmi pentru sisteme distribuite necesită un mecanism prin care să se desemneze un proces coordonator sau un proces inițiator care să realizeze funcții necesare tuturor proceselor din sistem. Aceste funcții pot fi, de exemplu, un mecanism centralizat de sincronizare, gestiunea unui grafic global pentru detecția interblocării sau controlul unui dispozitiv de intrare/ieșire. Dacă procesul coordonator dispăre datorită unor avarii a nodului sau, sistemul trebuie să-și continue execuția prin restartarea unei noi copii a coordonatorului pe un alt procesor disponibil. Algoritmii care determină unde trebuie restartată o nouă copie a coordonatorului se numesc *algoritmi de alegere*.

Se consideră, în general, că, pentru simplificare, algoritmii de alegere presupun: (1) o corespondență unu la unu între procese și procesoare; (2) fiecare proces are asociată o prioritate; se presupune că procesorul P_i are prioritatea i ; (3) coordonatorul este întotdeauna procesul cu numărul de prioritate cel mai mare; (4) cînd procesul coordonator dispăre, algoritmul trebuie să aleagă procesul activ cu cel mai mare număr de prioritate. Acest număr trebuie apoi cunoscut de toate procesele din sistem. În plus, un algoritm de alegere trebuie să furnizeze și posibilitatea ca un proces care se alătură sistemului să identifice coordonatorul curent.

În acest subcapitol vor fi prezentați doi algoritmi de alegere: primul este aplicabil în sisteme în care un proces poate trimite mesaje la toate celelalte procese iar al doilea pentru sisteme în care procesele sunt organizate în inel. Ambii algoritmi necesită n^2 mesaje pentru alegere.

4.4.1. Algoritmul Bully

Presupunem că un proces P_i trimite o cerere la care coordonatorul nu răspunde într-un interval de timp T_i ; în această situație, P_i presupune că procesul coordonator a dispărut, și încearcă să aleagă noul coordonator.

Algoritmii următori, datorat lui Garcia-Molina ([SP88][GMO82]) propune următoarea desfășurare a procedurii de alegere:

(1) Procesul P_i trimite un mesaj `ELECTION` la fiecare proces cu un număr de prioritate mai mare ca al său;

(2) P_i așteaptă apoi răspunsurile de la aceste procese un interval de timp T_i ; dacă nu se primește nici un răspuns de la aceste procese se presupune că toate au dispărut și P_i trebuie să devină noul coordonator. Procesul P_i restartează o nouă copie a coordonatorului și trimite un mesaj `INFORM` pentru a informa toate procesele active (cu prioritatea mai mică decît a sa) cine este noul coordonator.

(3) Dacă se primește totuși un mesaj în intervalul T_i , P_i așteaptă încă un interval de timp T_i recepția unui mesaj `INFORM` (alt proces, mai prioritar s-a ales coordonator). Dacă nu se primește mesajul `INFORM` în intervalul T_i , P_i presupune că procesul mai prioritar a eșuat în alegere și restartează algoritmul.

Rezultă că, în orice moment, în timpul execuției, P_i poate recepționa următoarele două mesaje de la un proces P_j : (1) `INFORM`: P_j este noul coordonator ($j > i$); P_i memorează această informație; (2) `ELECTION`: P_j a început procedul de alegere ($j < i$); P_i trimite un răspuns la P_j și începe propria sa procedură de alegere. Procesul care termină procedura de alegere va fi cel cu numărul de prioritate cel mai mare și noul coordonator. La revenirea în sistem a unui proces, se execută din nou același algoritm. Dacă nu există procese active cu numere de prioritate mai mari, procesul nou intrat forțează toate procesele cu numere mai mici de prioritate să-l accepte ca nou coordonator, chiar dacă există un coordonator activ.

Un exemplu de execuție al acestui algoritm este prezentat în figura 4.4.1.

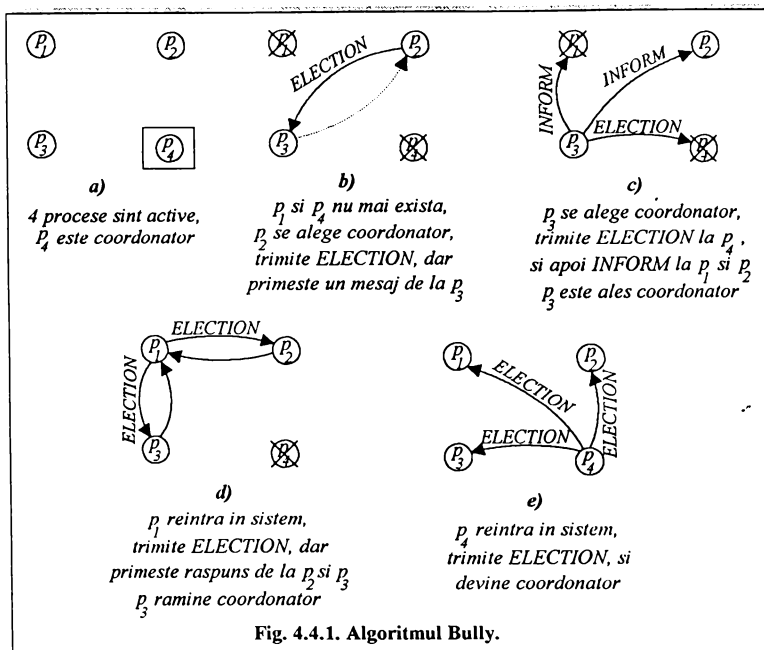


Fig. 4.4.1. Algoritm Bully.

4.4.2 Algoritm în inel

Tanenbaum descrie în [Tan96][Tan93] următorul algoritm în inel pentru alegerea coordonatorului:

- (1) Când un proces P_i detectează dispariția coordonatorului, construiește un mesaj ELECTION(i) conținând numărul procesului pe care îl trimite vecinului din dreapta (succesorul în inel). Dacă succesorul a dispărut, mesajul se trimite la succesorul succesorului s.a.m.d.
- (2) La fiecare pas, procesul emițător își înscrie numărul în mesaj.
- (3) Când procesul inițiator primește mesajul ELECTION în care și-a înscris numărul, schimbă tipul mesajului în COORDINATOR.
- (4) Mesajul COORDINATOR circulă o dată în întreg inelul pentru a informa toate procesele despre identitatea coordonatorului (cel mai mare număr din lista din mesaj) și despre identitatea proceselor din sistem. Algoritmul se încheie când mesajul COORDINATOR a circulat o dată prin tot inelul.

Capitolul 5.

PROIECTAREA ȘI IMPLEMENTAREA UNOR MECANISME SUPT PENTRU DEZVOLTAREA DE APLICAȚII DISTRIBUITE ÎN REȚELE UNIX, MICROSOFT WINDOWS (WIN32) ȘI NOVELL NETWARE

În acest capitol se prezintă proiectarea și implementarea mai multor variante pentru două toolkit-uri având următoarele două scopuri:

- (1) controlul accesului concurrent la o bază de date în rețele locale Unix, Windows sau Novell NetWare;
- (2) comunicații de grup în rețele locale Unix, Windows sau Novell NetWare.

5.1. Proiectarea unui toolkit pentru controlul accesului concurrent la o bază de date în rețele locale Unix, Windows (Win32) sau Novell NetWare

Una dintre cele mai importante probleme de sincronizare este coordonarea accesului mai multor utilizatori la un obiect comun, în mod uzual o bază de date, a cărei structură internă nu este specificată, dar care trebuie accesată fie pentru a fi citită, fie pentru a fi modificată [Nu92][Gos91]. Există deci două categorii de utilizatori: utilizatorii care modifică obiectul și deci, care pe parcursul modificării trebuie să aibă acces exclusiv și utilizatorii care citesc obiectul și care, pe parcursul citirii, pot partaja obiectul cu alți utilizatori în citire.

Există două variante ale problemei: varianta în care utilizatorii care au acces în citire sunt mai prioritari și varianta în care utilizatorii care au acces exclusiv pentru scriere sunt mai prioritari.

În acest subcapitol vor fi date soluții pentru ambele variante ale problemei, aceste soluții se vor implementa în rețele Unix, Windows și Novell.

Vor fi prezentate două categorii de soluții: soluții bazate pe semafoare și soluții bazate pe un proces central, coordonator. Semafoarele sunt disponibile în sistemul de operare Novell NetWare iar pentru rețelele Unix și Windows vor fi utilizate semafoarele distribuite a căror proiectare și implementare a fost prezentată în subcapitolul 4.3.4.

5.1.1. Soluții bazate pe semafoare

Soluțiile bazate pe semafoare pleacă de la ideea că, sincronizarea proceselor care cer acces exclusiv, pe de o parte, și pe de altă parte, sincronizarea lor cu procesele care cer acces partajat se poate face numai cu semafoare.

Starea sistemului poate fi descrisă prin următoarele elemente: (1) doi contori globali: unul care contorizează procesele ce sunt în citire la un moment dat (acces partajat) și denumit *noSharedAccess*, și al doilea contor care contorizează procesele care fie așteaptă pentru scriere, fie sunt în scriere; denumit *noExclAccess*; (2) o variabilă logică, *inExclAccess* care indică dacă o operație de scriere este în curs de desfășurare; (3) două liste: o listă a proceselor aflate în așteptare pentru citire, *sharedAccessList* și o listă a proceselor aflate în așteptare pentru scriere, *exclAccessList*. Pentru lucrul cu aceste liste se presupune că se dispune de două rutine care inserază un proces *p* în listă și respectiv extrage un proces *p* din listă (*GetList(p, lista)*, *PutList(p, lista)*) și o funcție care testează dacă o listă este vidă sau nu (*IsEmpty(lista)*).

Trebuie interzisă apariția situațiilor în care: (1) mai mult decât un proces este în acces exclusiv în scriere; (2) un proces este în scriere, în timp ce altul este în citire. Prima condiție va fi îndeplinită implicit datorită modului în care au fost alese variabilele de stare; considerând cazul în care procesele care cer acces exclusiv sunt mai prioritare, a doua condiție poate fi exprimată prin (*not inExclAccess or NoSharedAccess*). La sfârșitul unei operații de scriere sau al unui ansamblu de operații de citire, dacă ambele liste de procese nu sunt vide, se va da prioritate proceselor în așteptare la scriere.

În general, procesele vor avea structura următoare:

```

ReadProcess[i:1..n]::
do true →
    OpenSharedAccess(i)
    #citirea bazei de date
    CloseSharedAccess(i)
od
WriteProcess[i:1..n]::
do true →
    OpenExclAccess(i)
    #scriere in baza de date
    CloseExclAccess(i)
od

```

Se vor prezenta, în continuare, trei soluții de codificare a rutinelor *OpenSharedAccess*, *CloseSharedAccess*, *OpenExclAccess* și *CloseExcAccess* folosind semafoare.

A. Prima soluție bazată pe semafoare

Pentru a controla accesul la variabilele de stare va fi introdus un semafor care asigură excluderea reciprocă a celor n procese în accesul la aceste variabile, *SharedDataProtect*. Fiecrui proces k i se asociază un semafor *WaitingAccess[k]* pentru a controla accesul la baza de date. Ținând cont și de semnificațiile variabilelor de stare, descrierea rutinelor *OpenSharedAccess*, *CloseSharedAccess*, *OpenExclAccess* și *CloseExcAccess* este următoarea:

```

var NoSharedAccess: int :=0;      # contorizeaza numarul de procese in acces partajat
var NoExclAccess: int :=0;       # contorizeaza numarul de procese in asteptare pentru acces exclusiv si in
                                # acces exclusiv
var InExclAccess: bool :=false   # indica daca un proces este in cursul unui acces exclusiv
var ExclAccessList, SharedAccessList: queue of processes :=∅, ∅ # lista de procese
var SharedDataProtect: semaphore := 1      # controleaza accesul la variabilele de stare
var WaitingAccess[1..n]: semaphore :=([n] 0) # coordoneaza blocarea si deblocarea proceselor

```

```

OpenSharedAccess[i:1..n]::
    Down(SharedDataProtect) # se protejeaza actualizarea variabilelor de stare
    if NoExclAccess = 0 →
        NoSharedAccess:= NoSharedAccess + 1 # inca un proces in acces partajat
        Up(WaitingAccess[i]) # se pregateste validarea accesului
    [] NoExclAccess ≠ 0 → PutList(SharedAccessList,i) # se intirzie citirea din baza de date
    fi
    Up(SharedDataProtect) # se deblocheaza variabilele de stare
    Down(WaitingAccess[i]) # se incearca intrarea in acces partajat
    return

```

```

CloseSharedAccess[i:1..n]
    Down(SharedDataProtect) # se protejeaza actualizarea variabilelor de stare
    NoSharedAccess:= NoSharedAccess - 1 # s-a terminat un proces care a citit baza de date
    if NoSharedAccess=0 and not empty(ExclAccessList) → #nici un proces nu citește baza de date
        GetList(ExclAccessList,k) # se extrage din lista un proces care asteapta pentru scriere
        InExclAccess:= true # se marcheaza intrarea in acces exclusiv
        Up(WaitingAccess[k]) # se deblocheaza procesul k pentru a intra in acces exclusiv
    fi
    Up(SharedDataProtect) # se deblocheaza variabilele de stare
    return

```

```

OpenExclAccess[i:1..n]
    Down(SharedDataProtect) # se protejeaza actualizarea variabilelor de stare
    NoExclAccess:= NoExclAccess + 1 # inca un proces care asteapta pentru scriere
    if NoSharedAccess=0 and not InExclAccess → # nu exista nici un proces in citire sau scriere
        InExclAccess:= true # se marcheaza accesul exclusiv

```

Capitolul 5: Proiectarea și implementarea unor mecanisme suport pentru dezvoltarea de aplicații distribuite

```
    Up(WaitingAccess[i])    # se pregatește validarea accesului
[] NoSharedAccess ≠ 0 →
    PutList(Excl.AccessList, i) # se întârzie scrierea
fi
Up(SharedDataProtect)     # se deblochează variabilele de stare
Down(WaitingAccess[i])    # se încearcă intrarea în acces partajat
return

CloseExclAccess[i:1..n]:
    Down(SharedDataProtect) # se protejează actualizarea variabilelor de stare
    NoExclAccess := NoExclAccess - 1 # s-a terminat un proces care a scris în baza de date
    InExclAccess:=false     # se marchează ieșirea din acces exclusiv
    if not empty(Excl.AccessList) → # dacă există procese care așteaptă pentru scriere
        GetList(Excl.AccessList,k) # se extrage din listă un proces care așteaptă pentru scriere
        InExclAccess:= true       # se marchează intrarea în acces exclusiv
        Up(WaitingAccess[k])     # se deblochează procesul k pentru a intra în acces exclusiv
    [] empty(Excl.AccessList) →
        do not empty(SharedAccessList) → #cât timp există procese în așteptare pentru citire
            GetList(SharedAccessList,k) # se extrage un proces din listă
            NoSharedAccess:= NoSharedAccess + 1 # se contorizează
            Up(WaitingAccess[k])       # și se deblochează pentru a intra în citire
        od
    fi
    Up(SharedDataProtect)     # se deblochează variabilele de stare
    return
```

Se observă că un proces poate trece de `Down(WaitingAccess[k])` dacă fie baza de date a fost liberă, fie a fost activat de un alt proces în `CloseExclAccess()` sau `CloseSharedAccess()`. Toate procesele în așteptare pentru scriere vor fi activate, pe rând, la sfârșitul fazelor de scriere, ele având prioritate față de procesele în așteptare pentru citire. Procesele în așteptare pentru citire vor fi activate la sfârșitul unei faze de scriere numai dacă lista `Excl.AccessList` este vidă. Listele conțin identificatorii proceselor în așteptare: rutinele `GetList()` și `PutList()` pot fi realizate pentru a implementa o anumită strategie de planificare a proceselor.

B. A doua soluție cu semafoare

Soluția prezentată mai sus poate fi ușor adaptată pentru diverse strategii de planificare a proceselor. În schimbul obținerii acestui avantaj, trebuie menținute, însă, explicit, liste de procese, `ExclAccessList` și `SharedAccessList`. Cu costul micșorării flexibilității se poate renunța la gestionarea explicită a listelor de procese dacă se fac următoarele observații: (1) fiecărui set de secțiuni critice îi corespunde implicit o listă de procese în așteptare (2) un set de secțiuni critice pot fi blocate, în mod selectiv, astfel încât mai mult decât un proces dintr-o anumită clasă poate avea acces simultan la baza de date sau un singur proces din altă clasă să aibă acces exclusiv. Mecanismul de blocare selectivă utilizează următoarele două operații:

`DownWithCondition(exclusive, waiting: semaphore; var counter: int):`

```
    Down(exclusive)
    if counter = 0 → Down(waiting) fi
    counter:= counter + 1
    Up(exclusive)
    return
```

`UpWithCondition(exclusive, waiting semaphore; var counter: int):`

```
    Down(exclusive)
    counter:= counter + 1
    if counter = 0 → Up(waiting) fi
    Up(exclusive)
    return
```

Operația `DownWithCondition(exclusive, waiting, counter)` va bloca o operație `Down(waiting)`; totuși apelul următor `DownWithCondition()` cu aceiași parametri nu va produce blocare. Operația `UpWithCondition(exclusive, waiting, counter)` este operația corespunzătoare de deblocare.

Utilizând aceste operații, codificarea rutinelor de intrare și ieșire din secțiunile critice este:

```
var NoSharedAccess: int :=0      # contorizeaza numarul de procese in acces partajat
var NoExclAccess: int :=0       # contorizeaza numarul de procese in asteptare pentru acces
                                # exclusiv si in acces exclusiv
var SharedDataProtect: semaphore :=1 # folosit pentru protectia datelor in operatiile
                                # DownWithCondition() si UpWithCondition()
var Exclusive: semaphore :=1 # folosit pentru accesul exclusiv la datele partajate de procesele care cer citire
var ExclDataProtect: semaphore :=1 # folosit pentru protectia datelor in operatiile
                                # DownWithCondition si UpWithCondition
var StopExclAccess: semaphore :=1 # folosit pentru blocarea proceselor care cer scriere
var StopSharedAccess: semaphore :=1 # folosit pentru blocarea proceselor care cer citire
```

```
OpenSharedAccess[i:1..n]:
    Down(Exclusive)
    Down(StopSharedAccess)
    DownWithCondition(SharedDataProtect, StopExclAccess, NoSharedAccess)
    Up(StopSharedAccess)
    Up(Exclusive)
    return
```

```
CloseSharedAccess[i:1..n]
    UpWithCondition(SharedDataProtect, StopExclAccess, NoSharedAccess)
    return
```

```
OpenExclAccess[i:1..n]
    DownWithCondition(ExclDataProtect, StopSharedAccess, NoExclAccess)
    Down(StopExclAccess)
    return
```

```
CloseExclAccess[i:1..n]:
    Up(StopExclAccess)
    UpWithCondition(ExclDataProtect, StopSharedAccess, NoExclAccess)
    return
```

Operația *DownWithCondition()* din *OpenSharedAccess()* va genera blocarea unui proces care execută o operație *Down(StopExclAccess)* în *OpenExclAccess()* dar nu și alte procese care execută *DownWithCondition()* în *OpenSharedAccess()*; ca atare, un proces care are acces în citire la baza de date va bloca alte procese care cer scriere, dar nu și alte procese care cer citire. Operația *Down(StopExclAccess)* va bloca alte procese care cer citire

C. A treia soluție cu semafoare

Se prezintă acum o variantă în care procesele care cer acces partajat sunt mai prioritare. Singurul caz în care un proces care cere acces partajat trebuie să aștepte este acela în care baza de date este ocupată de un proces în scriere. Un proces care cere scriere nu poate conține accesul exclusiv decât atunci când nici un proces care cere citire nu este în execuție sau așteptare. Codificarea algoritmului este următoarea:

```
var NoSharedAccess : int := 0      # contorizeaza numarul de procese in acces partajat
var Exclusive : semaphore :=1     # folosit pentru protectia variabilelor comune
var PrioritySharedAccess : semaphore :=1 # folosit pentru a obtine prioritatea proceselor care cer citire
var WaitingExclAccess : semaphore :=1 # folosit pentru a bloca procesele care cer citire
```

```
OpenSharedAccess [i:1..n]:
    Down(Exclusive)                # se protejeaza actualizarea variabilelor de stare
    NoSharedAccess := NoSharedAccess +1 # inca un proces care citeste din baza de date
    if NoSharedAccess = 1 → Down(PrioritySharedAccess) fi # se blocheaza accesul proceselor
                                # care cer scriere
    Up(Exclusive)                  # se deblocheaza variabilele de stare
    return
```

```
CloseSharedAccess [i:1..n]:
```

```

Down(Exclusive)
NoSharedAccess := NoSharedAccess - 1    # s-a terminat un proces care e citit din baza de date
if NoSharedAccess = 0 → Up(PrioritySharedAccess) fi    # se deblocheaza accesul exclusiv
Up(Exclusive)    # se deblocheaza variabilele de stare
return
    
```

OpenExclAccess [i:1..n]:

```

Down(WaitingExclAccess)    # se blocheaza alte procese care cer scriere
Down(PrioritySharedAccess)    # se blocheaza alte procese care cer citire
return
    
```

CloseExclAccess [i:1..n]:

```

Up(PrioritySharedAccess)    # se deblocheaza intii procesele care asteapta pentru citire
Up(WaitingExclAccess)    # se deblocheaza apoi procesele care asteapta pentru scriere
return
    
```

Se observă că: (1) semaforul *Exclusive* se folosește pentru a proteja incrementarea și testarea variabilei *NoSharedAccess* (2) un proces care citește baza de date blochează alte procese care cer scriere prin *Down (PrioritySharedAccess)* dar nu și alte procese care cer citire (3) un proces care a obținut accesul exclusiv blochează alte procese care cer acces pentru scriere la *Down (WaitingExclAccess)* și alte procese care cer acces pentru citire la *Down (PrioritySharedAccess)* (4) Prioritatea proceselor care cer citire se obține datorită observațiilor 2 și 3.

5.1.2. Soluție bazată pe un proces coordonator

Toate soluțiile de sincronizare prezentate anterior au în comun faptul ca operațiile pentru obținerea accesului în secțiunile critice sunt realizate prin apelul unor rutine executate în cadrul proceselor apelate: astfel, dacă un proces este întrerupt în timp ce încearcă să obțină accesul într-o secțiune critică, baza de date devine blocată pe un interval nedefinit. Acest dezavantaj poate fi înlăturat dacă se separă toate operațiile pentru obținerea accesului și se înglobează într-un proces separat față de apelant, asociat bazei de date și denumit *manager*. O operație pentru obținerea accesului se apelează printr-un mesaj transmis procesului *manager*, care returnează un mesaj de răspuns atunci când apelantul poate intra în secțiune critică.

Codificarea procesului *manager* poate fi realizat prin analogie cu prima soluție cu semafoare. În plus față de variabilele de stare, se mai menține o variabilă *CliNo* care conține identificatorul clientului apelant și o variabilă *Request* care identifică tipul operației cerute. În aceste condiții, procesul *manager* poate fi codificat astfel:

```

type Kind = enum (r,R,w,W)    # tipuri de operatii
chan managerChan(kind, cliNo: int)    # canal pentru transmiterea mnesajelor la manager
chan clientChan[1:n] ()    # canale ale clientilor pentru primirea raspunsurilor de la manager
    
```

Manager:

```

var NoSharedAccess: int = 0    # contorizeaza numarul de procese in acces partajat
var NoExclAccess: int = 0    # contorizeaza numarul de procese in asteptare pentru
    # acces exclusiv si in acces exclusiv
var InExclAccess: bool = false    # indica daca un proces este in cursul unui acces exclusiv
var ExclAccessList, SharedAccessList: queue of processes := {}, {} # liste de procese
var CliNo: int    # folosit ca identificator de proces
var Request: kind    # folosit pentru a identifica tipul operatiei cerute
    
```

OpenSharedAccess()

```

if NoExclAccess = 0 →    #indica daca un proces este in cursul unui acces exclusiv
    NoSharedAccess := NoSharedAccess + 1
    send clientChan[CliNo]() # deblocare apelant
[] NoExclAccess ≠ 0 → PutList(SharedAccessList, CliNo)
fi
return
    
```

CloseSharedAccess()

Capitolul 5: Proiectarea și implementarea unor mecanisme suport pentru dezvoltarea de aplicații distribuite

```

NoSharedAccess := NoSharedAccess - 1      # s-a terminat un proces care a citit din baza de date
send ClientChan[CliNo]()                  # se deblocheaza clientul apelant
if NoSharedAccess = 0 and not empty (ExclAccessList) →
    GetList(ExclAccessList, CliNo)
    InExclAccess := true
    send ClientChan[CliNo]()              # se deblocheaza un client care a cerut scriere
fi
return
OpenExclAccess()
NoExclAccess := NoExclAccess + 1          # un alt proces in asteptare pentru scriere
if NoSharedAccess = 0 and not InExclAccess →
    InExclAccess := true
    send ClientChan[CliNo]()              # se deblocheaza clientul
|| not ( NoSharedAccess = 0 and not InExclAccess) →      # se amina accesul exclusiv
    PutList(ExclAccessList, CliNo)
fi
return
CloseExclAccess()
NoExclAccess := NoExclAccess - 1          # s-a terminat un proces care a scris in baza de date
InExclAccess := false
send Client Chan[CliNo]()                 # se deblocheaza clientul apelant
if not empty (ExclAccessList) →
    GetList (ExclAccessList, CliNo)
    InExclAccess := true
    send ClientChan [CliNo]()            # se deblocheaza un client care a cerut scriere
|| empty (ExclAccessList) →
    do not empty (SharedAccessList) →
        GetList (SharedAccessList, CliNo)
        NoSharedAccess := NoSharedAccess + 1
        send ClientChan[CliNo]()
    od
fi
return
Manager()
do true →
    receive managerChan(Request, CliNo)    # se primeste un mesaj care contine o cerere de
                                            # executie a unei operatii
    case Request of
        r: OpenSharedAccess()
        R: CloseSharedAccess()
        w: OpenExclAccess()
        W: CloseExclAccess()
    esac
od

```

Procesele client pot fi codificate astfel:

```

Client [i:1..n]:
do true →
    send managerChan (r, i)                # protocol pentru intrare in sectiunea critica pentru
    receive clientChan[i]()                # citire. la fel pentru scriere
    # citire din baza de date              # citire
    send managerChan(R, i)                 # protocol pentru iesirea din sectiunea critica
    receive clientChan[i]()
od

```

5.1.3. Implementări în rețele Unix, Microsoft Windows (Win32) și Novell Netware

Soluțiile bazate pe semafoare și prezentate la 5.1.1 au fost implementate în rețele Novell Netware, având în vedere că acest sistem de operare furnizează mecanismul semafoarelor și în rețele Unix și Microsoft Windows, având în vedere că în acest sistem de operare a fost implementat mecanismul semafoarelor distribuite (vezi 4.3.4.1 și 4.3.4.4.)

5.1.3.1. Implementarea soluțiilor bazate pe semafoare în rețele Novell Netware

Implementarea soluțiilor bazate pe semafoare în rețele Novell Netware este prezentată în anexa C. Ceea ce este caracteristic tuturor acestor trei implementări este faptul că variabilele comune tuturor proceselor (variabilele de stare) au fost implementate utilizând semafoare. Prin urmare, pentru incrementarea contoarelor *NoSharedAccess*, *NoExclAccess* vor fi executate funcțiile Netware *SignalSemaphore()*, pentru incrementare, respectiv *WaitOnSemaphore()* pentru decrementare, iar pentru examinarea valorilor lor, funcția *ExamineSemaphore()*. Implementarea soluției a doua și a treia cu semafoare este mai directă, deoarece nu necesită liste de procese. Implementarea primei soluții necesită realizarea a două liste de procese; s-a optat ca un proces, asociat unui singur nod, să fie identificat prin numărul de conexiune asociat nodului respectiv. Deoarece și listele de procese sunt variabile comune, este necesar ca ele să fie realizate utilizând tot semafoare; astfel, pentru implementarea listelor se folosește un tablou de semafoare, a cărui dimensiune maximă este egală cu cel mai mare număr de conexiune posibil. Faptul că un proces este în așteptare pentru scriere va fi semnalat prin două operații *SignalSemaphore()* asupra semaforului din tabloul de semafoare de pe poziția corespunzătoare numărului de conexiune asociat procesului; analog, printr-o singură operație *SignalSemaphore()* va fi semnalat un proces care așteaptă la citire.

Toate trei soluțiile oferă unui program utilizator următoarele funcții de bibliotecă:

- (1) *int ExclOrSharedAccessInitialize()* apelată pentru deschiderea semafoarelor și determinarea numărului de conexiune propriu. Funcția returnează un cod de eroare sau 0 pentru execuție fără erori.
- (2) *int OpenSharedAccess(void)* apelată pentru obținerea accesului în mod partajat.
- (3) *int OpenExclAccess(void)* apelată pentru obținerea accesului în mod exclusiv.
- (4) *int CloseSharedAccess(void)* apelată pentru terminarea accesului în mod partajat.
- (5) *int CloseExclAccess(void)* apelată pentru terminarea accesului în mod exclusiv.
- (6) *int ExclOrSharedAccessEnd(void)* apelată pentru deconectarea de la sistem.

Toate funcțiile returnează un cod de eroare sau zero pentru execuție cu succes.

5.1.3.2 Implementarea soluțiilor bazate pe semafoare în rețele Unix și MSWindows (Win32)

Implementarea în rețele Unix și Windows a soluției a treia bazată pe semafoare, în care procesele care cer acces partajat sunt mai prioritare, este prezentată în anexa C.

Se folosește serverul pentru semafoare distribuite a cărui proiectare și implementare a fost prezentată la 4.3.4.1. și 4.3.4.4.

Pe baza funcțiilor *Semaphore_XX()* puse la dispoziția programului client au fost concepute, conform algoritmului, funcțiile:

- (1) *void OpenSharedAccess(void)* apelată pentru obținerea accesului în mod partajat.
- (2) *void OpenExclAccess(void)* apelată pentru obținerea accesului în mod exclusiv.
- (3) *void CloseSharedAccess(void)* apelată pentru terminarea accesului în mod partajat.
- (4) *void CloseExclAccess(void)* apelată pentru terminarea accesului în mod exclusiv.

Esențial pentru implementarea variabilelor comune prin semafoare distribuite este faptul că funcțiile *Semaphore_Up()* și *Semaphore_Down()* returnează valoarea semaforului din momentul începerii execuției funcției.

Pentru a folosi funcțiile de mai sus, un program client (cu sau fără interfață grafică) trebuie linkeditat cu aceste funcții și, de asemenea trebuie executat în *background* serverul pentru semafoare distribuite. Semafoarele care simulează variabilele comune sunt inițializate la începutul programului prin apelul explicit al funcțiilor *Semaphore_Down* și *Semaphore_Up*.

În varianta prezentată în anexă se lucrează cu un singur client într-un nod, atașat unui server pentru semafoare distribuite.

5.1.3.3. Implementarea soluției bazată pe un proces coordonator în rețele Novell Netware

Protocolul folosit pentru comunicarea clienți-proces coordonator (server) este protocolul IPX. Identificarea atât a clienților cât și a procesului coordonator se face pe baza numerelor de conexiuni; în orice mesaj emis către server, clienții înglobează și numărul propriu de conexiune, pentru a putea fi identificați pentru o deblocare ulterioară.

Serverul constricte o coadă de cereri *queue[`MAX_NO_PACKETS`]*; sosirea unui mesaj de la un client declanșează execuția unei rutine declarate ca *EventServiceRoutine*, care depune în coada de cereri mesajul sosit. Se ocolește astfel situația ca, datorită unei supraîncărcări a serverului, acesta să piardă mesaje. Listele de procese sunt implementate ca tablouri de conexiuni ale proceselor client. Funcția *GetPriProcess()* extrage un proces client din listă conform priorității proceselor client, alese egale cu conexiunile lor.

Procesele client li se pun la dispoziție aceleași funcții ca și la 5.1.3.1:

- (1) int *ExcOrSharedAccessInitialize()* apelată pentru deschiderea soclurilor, conectarea rutinei ESR, determinarea adresei Ethernet a serverului cunoscându-se user-ul la care acesta a fost lansat, determinarea conexiunii proprii.
- (2) int *OpenSharedAccess(void)* pentru obținerea accesului în mod partajat, prin schimb de mesaje cu serverul.
- (3) int *OpenExclAccess(void)* apelată pentru obținerea accesului în mod exclusiv.
- (4) int *CloseSharedAccess(void)* apelată pentru terminarea accesului în mod partajat.
- (5) int *CloseExclAccess(void)* apelată pentru terminarea accesului în mod exclusiv.
- (6) int *ExcOrSharedAccessEnd(void)* apelată pentru deconectarea de la sistem.

Pentru demonstrarea tuturor soluțiilor implementate în rețele Novell Netware a fost scris un singur program client care se linkidează cu setul de funcții pus la dispoziție de soluția alcasă pentru verificare.

5.1.4. Compararea soluțiilor propuse și concluzii

În urma analizei celor patru soluții propuse și a implementării lor în rețele locale se pot enunța următoarele concluzii:

- (1) Toate soluțiile prezentate necesită implementarea unor variabile de stare; acestea pot fi implementate fie separat, în cadrul unui proces coordonator, fie pot fi simulate ca "variabile comune" prin asimilarea lor cu valorile unor semafoare distribuite.
- (2) Prin faptul că mențin o listă de procese, prima soluție cu semafoare și soluția procesului central pot planifica procesele pentru execuție după diverși algoritmi.
- (3) Pe de altă parte, a doua și a treia soluție cu semafoare, negestionând liste, sunt mai rapide.
- (4) Implementarea soluției bazată pe un proces coordonator se face având în vedere raportul dintre timpul necesar tratării unei cereri și rata de sosire a cererilor, în conformitate cu cele arătate la 3.1.2., pentru a elimina posibilitatea apariției situației de "sufocare" a serverului.
- (5) Toate soluțiile prezintă posibilitatea apariției amânării continue a unui proces din cele două categorii de procese (în funcție de cui se acordă prioritate: proceselor care cer acces exclusiv sau celor care cer acces partajat). Amânarea la infinit a proceselor poate fi înlăturată dacă, la terminarea unui proces care cere acces exclusiv (cazul primelor două soluții) se au în vedere ambele liste de procese în așteptare.
- (6) Deoarece mesajele care se transmit în sistem, sunt de lungime foarte mică, având semnificație numai pentru sincronizare, s-a optat ca mesajele să fie codificate ASCII astfel încât comunicarea este valabilă în rețele eterogene (fără a se pune problema de "big" sau "little "endian"). Singura cerință este cea a sistemului de operare (Unix sau Windows).
- (7) Având în vedere că sistemul Unix este multiuser, o cerință în plus ar fi eliminarea restricției unui singur client pentru un singur server pe un nod. Restricția provine de fapt de la implementarea semafoarelor distribuite. Ea poate fi eliminată ușor dacă se folosește pentru fiecare client, pentru sincronizare, un canal *pipe* rezervat special, pentru Unix; aceste canale pot fi alocate static, la inițializarea serverului, dar limitând astfel numărul de clienți de pe un nod, sau dinamic, dar în acest caz însă mai trebuie introdusă și o cerere specială, a fiecărui client, de "conectare" la sistem.
- (8) Identificarea clienților s-a făcut pe baza numerelor de conexiune în Novell și pe baza unor identificatori transmiși prin linia de comandă în Unix și Windows. Dacă sistemul se extinde pentru a permite mai multe procese client într-un nod se poate folosi aceeași metodă, pentru a identifica canalele *pipe* asociate clienților pe baza acestor identificatori sau se pot folosi identificatorii asociați proceselor de către sistemul de operare ("pid"-urile Unix ale proceselor).
- (9) Toate soluțiile prezentate se bazează pe respectarea protocolului de obținere a accesului de către procesele client.

5.2. Proiectarea și implementarea unui toolkit pentru comunicații de grup

Acest capitol înglobează aspectele teoretice prezentate în capitolele 2,3 și 4 în cadrul proiectării unui toolkit pentru comunicații de grup în rețele Unix și Windows, denumit TGC.

5.2.1. Cerințe de proiectare

Înainte de a prezenta arhitectura sistemului TGC se enumeră câteva cerințe de proiectare impuse:

- (1) Sistemul va folosi modelul grupurilor *deschise*; aceasta înseamnă că, pentru a putea transmite mesaje la toate procesele componente ale unui grup, un proces conectat la TCG nu trebuie să fie membru al aceluși grup.
- (2) Un proces trebuie să poată aparține simultan mai multor grupuri.
- (3) Fiecare grup multicast poate avea asociate mai multe pseudonime (*alias*-uri).
- (4) Sistemul trebuie să se execute în mod *utilizator*;
- (5) Protocolul de comunicație trebuie să fie independent de mașina hardware;
- (6) Sistemul trebuie să genereze un trafic de mesaje minim.
- (7) Sistemul trebuie să furnizeze o comunicație sigură.

TGC trebuie să furnizeze primitive relative la:

- (1) crearea, alăturarea la și părăsirea unui grup;
- (2) atribuirea de pseudonime grupurilor, dezafectarea unor pseudonime, căutarea unui grup;
- (3) transmisia și recepționarea de mesaje la/de la procese aparținând unui grup multicast;
- (4) conectarea și deconectarea de la sistem.

Înainte de apelul oricărei operații, orice client trebuie să se conecteze la sistem; conectarea la sistem înseamnă obținerea unui identificator, care este apoi livrat în corpul fiecărei cereri, pentru identificarea clienților.

Se prezintă în continuare o descriere în detaliu a primitivelor furnizate de TGC:

- (1) `int CreateGroup(idGroup *pgroup, char *name)`

Scop:

Asigură crearea unui nou grup multicast, al cărui nume se indică în apel; procesul apelant devine automat un membru al grupului. Identificatorul noului grup se stabilește de către un server central și este unic; el este returnat apelantului.

Parametri:

name - reprezintă numele stabilit pentru noul grup.

pgroup - pointer la identificatorul noului grup creat

Valoare returnată:

Funcția returnează valoarea *SUCCESS* în cazul în care crearea a reușit, sau un cod de eroare în caz contrar. Codurile de eroare care pot fi returnate sunt:

[EBADF]	Format eronat al parametrului de intrare (<i>nume</i> indică un șir nul)
[EGROUP_ALREADY_CREATED]	Grupul există deja
[ETOO_MANY_GROUPS]	Spațiu insuficient pentru înregistrarea unui nou grup
[ETOO_MANY_HOSTS_IN_GROUP]	Grupul este deschis pe maximum de noduri posibil
[ECLIENT_NOT_CONNECTED]	Procesul nu s-a conectat la sistem

- (2) `int JoinGroup(idGroup group)`

Scop:

Permite unui proces atașarea la un grup multicast, *group*. Orice proces trebuie să execute această primitivă pentru a putea transmite sau recepționa mesaje în contul acestui grup.

Parametri:

group - identificatorul grupului la care se dorește atașarea

Valoare returnată:

Funcția returnează valoarea *SUCCESS* în cazul în care alăturarea a reușit, sau un cod de eroare în caz contrar. Codurile de eroare care pot fi returnate sunt:

[ETOO_MANY_CONC_GROUPS]	Procesul este atașat la numărul maxim de grupuri
[EFULL_GROUP]	Spațiu insuficient pentru înregistrarea unui nou proces în acest grup
[ECLIENT_ALREADY_CONNECTED]	Procesul există deja în grup
[EWRONG_DEST_GROUP]	Grupul la care procesul dorește să se atașeze nu este creat

[ECLIENT_NOT_CONNECTED] Procesul nu s-a conectat la sistem

(3) **int LeaveGroup(idGroup group)**

Scop:

Permite părăsirea unui grup grup de către un proces. Primitiva se execută de către un client în momentul în care nu mai dorește apartenența la un grup. După părăsirea unui grup, procesul respectiv nu mai poate primi mesaje adresate aceluși grup, chiar dacă ele sunt memorate deja în buffer-ul de mesaje al aceluși proces.

Parametri:

group - identificatorul grupului la care procesul aparține

Valoare returnată:

Funcția returnează valoarea *SUCCESS* în cazul în care operația a reușit, sau un cod de eroare în caz contrar. Codurile de eroare care pot fi obținute sunt:

[WRONG_DEST_GROUP] Grupul nu există
[ECLIENT_NOT_IN_DEST_GROUP] Procesul nu aparține grupului
[ECLIENT_NOT_CONNECTED] Procesul nu s-a conectat la sistem

(4) **int CreateAliasGroup(idGroup group, char *alias)**

Scop:

Permite crearea unui pseudonim pentru grupul group; nu sunt permise aceleași pseudonime pentru două grupuri. Procesul apelant nu trebuie să fie membru al grupului.

Parametri:

group - identificatorul grupului pentru care se dorește crearea pseudonimului

alias - pseudonimele

Valoare returnată:

Funcția returnează valoarea *SUCCESS* în cazul în care operația a reușit, sau un cod de eroare în caz contrar. Codurile de eroare care pot fi obținute sunt:

[EBADF] Format eronat al parametrului de intrare (alias indică un șir nul)
[WRONG_DEST_GROUP] Grupul nu există
[ETOO_MANY_ALIASES] Spațiu insuficient pentru înregistrarea unui nou pseudonim
[WRONG_ALIAS] Există un grup care are acest pseudonim
[ECLIENT_NOT_CONNECTED] Procesul nu s-a conectat la sistem

(5) **int DestroyAliasGroup(idGroup group, char *alias)**

Scop:

Permite distrugerea unui pseudonim pentru grupul group. Procesul apelant nu trebuie să fie membru al grupului.

Parametri:

group - identificatorul grupului pentru care se dorește distrugerea pseudonimului

alias - pseudonimele

Valoare returnată:

Funcția returnează valoarea *SUCCESS* în cazul în care operația a reușit, sau un cod de eroare în caz contrar. Codurile de eroare care pot fi obținute sunt:

[EBADF] Format eronat al parametrului de intrare (alias indică un șir nul)
[WRONG_DEST_GROUP] Grupul nu există
[WRONG_ALIAS] Pseudonimele nu există pentru acest grup
[ECLIENT_NOT_CONNECTED] Procesul nu s-a conectat la sistem

(6) **int GetIdGroup(idGroup *pgroup, char* alias)**

Scop:

Permite căutarea identificatorului unui grup cînd se cunoaște un pseudonim.

Parametri:

alias - un pseudonim al grupului pentru care se dorește obținerea identificatorului

pgroup - la adresa pgroup se returnează, în caz de succes, identificatorul grupului.

Valoare returnată:

Funcția returnează valoarea *SUCCESS* în cazul în care operația a reușit, sau un cod de eroare în caz contrar. Codurile de eroare care pot fi obținute sunt:

[EBADF] Format eronat al parametrului de intrare (alias indică un șir nul)
[WRONG_ALIAS] Pseudonimele nu există

Capitolul 5: Proiectarea și implementarea unor mecanisme suport pentru dezvoltarea de aplicații distribuite

[*ECLIENT_NOT_CONNECTED*] Procesul nu s-a conectat la sistem

(7) `int SendToGroup(idGrup dest_gr_id, idGrup source_gr_id,
char *mes, int counter, pf filter)`

Scop:

Permite transmiterea unui mesaj tuturor proceselor dintr-un grup multicast indicat în parametrul *dest_group*. Deoarece procesul care transmite mesajul poate fi conectat simultan la mai multe grupuri, el trebuie să indice sistemului, în parametrul *source_gr_id*, grupul sursă în contul căruia se transmite mesajul. Mesajul constă într-un șir de octeți de lungime *contor*. În mesajul transmis de sistemul TCG se transmite și numele utilizatorului, indicat la conectarea la sistem prin funcția *InitTCG()*.

Parametri:

dest_gr_id - identificatorul grupului la care se dorește transmisia mesajului

source_gr_id - identificatorul grupului la care aparține procesul care transmite mesajul, grup în contul căruia se transmite mesajul

counter - dimensiunea mesajului

mes - mesajul de transmis

pfilter - pointer la o funcție utilizator apelată de sistemul TGC înainte de emisia fiecărui mesaj, pentru codificarea lui; valoarea NULL indică absența ei.

Valoare returnată:

Funcția returnează valoarea *SUCCESS* în cazul în care operația a reușit, sau un cod de eroare în caz contrar. Codurile de eroare care pot fi obținute sunt:

[<i>EWRONG_SOURCE_GROUP</i>]	Grupul sursa nu există
[<i>EWRONG_DEST_GROUP</i>]	Grupul destinație nu există
[<i>ECLIENT_NOT_IN_SOURCE_GROUP</i>]	Procesul nu aparține grupului sursă
[<i>EMES_TOO_LONG</i>]	Dimensiunea mesajului depășește limita maximă
[<i>ECLIENT_NOT_CONNECTED</i>]	Procesul nu s-a conectat la sistem

(8) `int ReceiveFromGroup(idGrup sgr_id, idGrup dgr_id, int
*pcounter, char *mes, int timeLimit, pf pfilter)`

Scop:

- Permite unui proces să recepționeze *pcounter* mesaje de la procesele unui grup. Un proces poate recepționa un mesaj numai dacă este membru al grupului la care s-a transmis mesajul receptiv. *Pfilter* este o funcție furnizată de procesul apelant, apelată de TGC pentru fiecare mesaj recepționat. Rolul ei este de a decodifica mesajele. TGC memorează toate mesajele primite pentru un client pînă cînd acesta execută *ReceiveFromGroup()*.

Parametri:

sgr_id - identificatorul grupului de la care se dorește recepția de mesaje. Poate fi folosit și identificatorul generic *ANY*, cu sensul că se recepționează mesaje de la orice grup

dgr_id - identificatorul grupului către care se transmite mesajul și la care procesul apelant trebuie să fie atașat

mes - adresa la care se vor depune mesajele recepționate

pcounter - adresa unui întreg, la care, inițial, procesul a depus numărul de mesaje ce se dorește recepționate și la care sistemul returnează numărul de mesaje recepționate.

pfilter - pointer la o funcție utilizator apelată de sistemul TGC după recepția fiecărui mesaj, pentru decodificarea lui; valoarea NULL indică absența ei.

timeLimit - Poate avea următoarele valori:

- *INFINITE*, cînd apelantul este pus în așteptare pînă la primirea a *counter* mesaje,
- o valoare întreagă *n* care semnifică punerea în așteptare a apelantului pentru maximum *n* secunde, dacă în bufler nu există *pcounter* mesaje; după *n* secunde se livrează toate cele *pcounter* mesaje sau cite s-au găsit în bufler
- *NOW* care semnifică returnarea imediată a controlului cu livrarea a maximum *pcounter* mesaje.

La returnarea controlului, la *pcounter* se returnează numărul de mesaje recepționate

Valoare returnată:

Funcția returnează valorile:

-*SUCCESS*, în cazul în care s-a recepționat cel puțin un mesaj și nu mai sunt alte mesaje în așteptare în buffer,
-*MORE_AIES*, în cazul în care s-a recepționat cel puțin un mesaj și mai sunt și alte mesaje în așteptare în buffer,
-*NO_WAIT_MESSAGES* în cazul în care nu s-a recepționat nici un mesaj.
-un cod de eroare în caz contrar. Codurile de eroare care pot fi obținute sunt:
[*EWRONG_SOURCE_GROUP*] Grupul sursa nu există
[*EWRONG_DEST_GROUP*] Grupul destinație nu există
[*ECLIENOT_IN_DEST_GROUP*] Procesul nu aparține grupului destinație
[*ENO_BUF*] Nu există spațiu suficient pentru înregistrarea procesului ca blocat, în așteptare de mesaje.
[*ECLIENOT_CONNECTED*] Procesul nu s-a conectat la sistem

(8) **int GetNoOfMessages(idGroup source_gr_id, idGroup dest_gr_id, int *pcounter)**

Scop:

Permite unui proces să determine la *pcounter* numărul de mesaje, aflate în buffer, în așteptare pentru recepție.

Parametri:

source_gr_id - identificatorul grupului de la care se dorește recepția de mesaje. Poate fi folosit și identificatorul generic *ANY*, cu sensul că se recepționează mesaje de la orice grup

dest_gr_id - identificatorul grupului către care s-a transmis mesajul și la care procesul apelant trebuie să fie atașat

pcounter - adresa unui întreg la care sistemul returnează numărul de mesaje recepționate.

Valoare returnată:

Funcția returnează valoarea *SUCCESS* sau un cod de eroare. Codurile de eroare care pot fi obținute sunt:

[*EWRONG_SOURCE_GROUP*] Grupul sursă nu există
[*EWRONG_DEST_GROUP*] Grupul destinație nu există
[*ECLIENOT_IN_DEST_GROUP*] Procesul nu aparține grupului sursă
[*ECLIENOT_CONNECTED*] Procesul nu s-a conectat la sistem

(9) **int InitTGC(char *user, char *autoProto, char *autoString)**

Scop:

Permite conectarea unui proces la sistemul TGC.

Parametri:

user - numele utilizatorului

autoProto - specifică protocolul de autorizare utilizat de serverul local pentru a determina dacă se permite unui proces să aibă acces la TGC sau nu

autoString - șir (parola) specific autorizării

Valoare returnată:

Funcția returnează valoarea *SUCCESS* în cazul în care conectarea a reușit, sau un cod de eroare în caz contrar. Codurile de eroare care pot fi obținute sunt:

[*ENO_AUTHORIZATION*] Conectarea nu este autorizată în urma verificării parolei
[*ETOO_MANY_CLIENTS*] Sunt conectate deja maximum de procese la serverul local
[*ECLIENOT_CONNECTED*] Procesul nu s-a conectat la sistem

(10) **int CloseTGC(void)**

Permite deconectarea unui proces de la serverul local.

Valoare returnată:

Funcția returnează valoarea *SUCCESS* în cazul în care operația a reușit, sau un cod de eroare în caz contrar:

[*ECLIENOT_CONNECTED*] Procesul nu s-a conectat la sistem

Toate funcțiile descrise mai sus vor fi puse la dispoziția aplicațiilor utilizator într-o bibliotecă de funcții.

5.2.2. Varianta semidistribuită

Se propune ca arhitectură a sistemului TGC o arhitectură semidistribuită.

5.2.2.1. Arhitectura generală a sistemului TGC

Întreaga proiectare se bazează pe interacțiunea client-server dintre procese. În fiecare nod TGC există câte un server local care comunică cu un server global, situat pe un nod fixat în rețea (nodul central). Clienții comunică cu serverele locale de pe nodurile lor, utilizând funcții puse la dispoziție de TGC prin intermediul unei biblioteci.

Arhitectura sistemului este prezentată în figurile 5.2.2.1.1. și 5.2.2.1.2.

Această structură distribuită oferă următoarele avantaje:

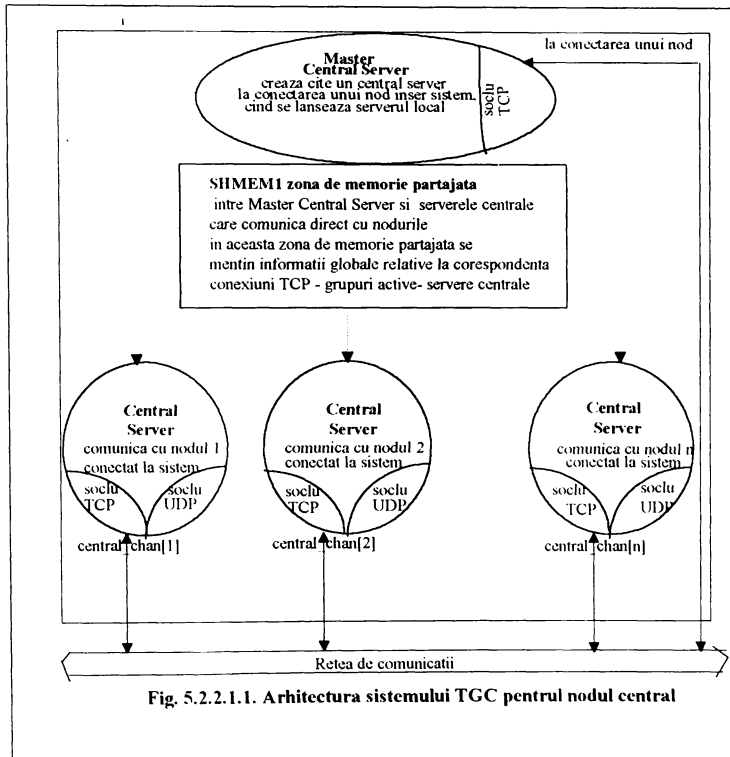


Fig. 5.2.2.1.1. Arhitectura sistemului TGC pentru nodul central

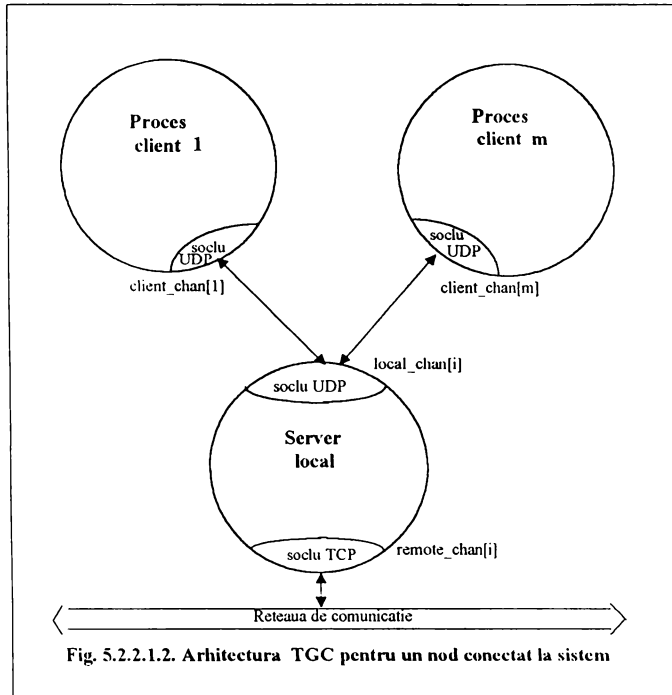
- (1) elimină posibilitatea sufocării unui nod, suprasolicitat într-o variantă centralizată;
- (2) asigură, într-un anumit grad, siguranță în funcționare: dacă unul dintre nodurile serverelor locale suferă o avarie, atunci serverul global închide conexiunea cu nodul respectiv și își actualizează corespunzător tablelele interne;
- (3) în sistemul distribuit, serverul global trebuie să transmită mesaje, numai către serverele locale; acestea, la rândul lor vor fi cele care transmit aceste mesaje către toate procesele client interne.

Sistemul TGC este compus din următoarele elemente:

- (1) *biblioteca de funcții* pusă la dispoziție programelor client; fiecare program client trebuie linkeditat cu aceasta bibliotecă. Ea implementează primitivile descrise în subcapitolul 5.2.2.3., pe care un client le poate apela. Aceste primitive sunt transformate în pachete corespunzătoare protocolului de comunicație și trimise serverului local;
- (2) *serverul local*, care se execută în fiecare nod TGC memorează o listă a proceselor client locale împreună cu apartenența lor la grupuri. Serverele locale asigură distribuția mesajelor la procesele de pe același nod (distribuție locală). Fiecare client trebuie să se conecteze la serverul său local utilizând un mecanism IPC (cozi de mesaje, FIFO);
- (3) *serverul global* care se execută pe un singur nod al sistemului. Rolul său este de a menține informații globale despre sesiunile multicast (grupurile active și identitatea nodurilor TGC) și de a distribui mesajele

între noduri. Mesajele către un grup multicast sunt trimise către un server local numai dacă pe acel nod grupul respectiv este activ (dacă există cel puțin un client conectat la acel grup).

Spațiul numelor alias se memorează în nivelul sistemului global; astfel, serverele locale nu tratează operațiile relative la alias-uri. Localizarea serverului global este transmisă serverelor locale fie ca variabile de mediu, fie ca parametri în linia de comandă la începutul execuției.



5.2.2.2. Protocolul de comunicație

TGC folosește pentru comunicația dintre serverele locale și clienți protocolul UDP, fiind adecvat pentru situația mai-mulți-la-unul din acest caz. Comunicația în rețea între serverele locale și serverele centrale se face prin protocolul TCP. Pentru fiecare nod care se conectează în sistem, serverul central master creează câte un fiu, un server central, cu rolul de a superviza comunicația cu serverul local din nodul respectiv.

O cerere adresată sistemului TGC are următoarea structură:

```

knd: kind. # tipul cererii
client_id: int # identificatorul clientului de la care s-a trimis cererea
dest_group_id: int # identificatorul grupului de la care s-a trimis cererea
source_group_id: int # identificatorul grupului sursa care a trimis mesajul sau cererea
group_name[1:MAX_LEN_GROUP]: char # nume de grup
len: int. # lungimea partii variabile a mesajului
mes[1:len]: char # mesajul
    
```

Răspunsurile returnate de serverul central au structura următoare:

```

knd: kind. # tipul cererii
error_code: int. # codul erorii sau zero pentru succes
client_id: int # identificatorul clientului de la care s-a trimis cererea
dest_group_id: int. # identificatorul grupului destinație
source_group_id: int. # identificatorul grupului sursa
len: int. # lungimea partii variabile a mesajului
mes[1:len]: char # mesajul
    
```

5.2.2.3. Descrierea serverelor locale și a serverului central

În cele ce urmează, se prezintă descrierea SR a serverelor locale și serverului central (master și slave). Pentru fiecare nod care se conectează la sistem, serverul master central creează câte un proces fiu cu rolul de a superviza comunicația cu acel nod.

Sunt prezentate mai întâi tipurile de mesaje definite pentru TGC și structura informației care circulă pe canalele de comunicație clienți - servere locale și servere locale - servere centrale..

```
type kind = enum (
    # tipuri de mesaje care circula in sistem
    INIT_TGC, CLOSE_TGC, CONNECTION_REQ, CREATE_GROUP, JOIN_GROUP,
    GET_ID_GROUP, LEAVE_GROUP, CREATE_ALIAS_GROUP,
    DESTROY_ALIAS_GROUP, SEND_TO_GROUP, RECEIVE_FROM_GROUP_NOW,
    RECEIVE_FROM_GROUP_INFINITE, NO_MORE_RECEIVE_FROM_GROUP,
    SEND_MESSAGE, GET_NO_OF_MESSAGES
)

type error = enum (
    # codificarea raspunsurilor la cererile trimise de clienti
    SUCCESS, MORE_MES, NO_WAIT_MESSAGES, ECLIENT_NOT_CONNECTED,
    ENO_AUTHORIZATION, ETOO_MANY_CLIENTS, ETOO_MANY_CONC_GROUPS,
    EFULL_GROUP, ECLIENT_ALREADY_CONNECTED,
    ECLIENT_NOT_IN_DEST_GROUP, EWRONG_SOURCE_GROUP,
    ECLIENT_NOT_IN_SOURCE_GROUP, EWRONG_DEST_GROUP,
    ENO_BUF, EGROUP_ALREADY_CREATED, ETOO_MANY_GROUPS,
    ETOO_MANY_HOSTS_IN_GROUP, EWRONG_ALIAS, EMES_TO_LONG, EBADF
)

chan client_chan[1..n, 1..MAX_CLI] (
    error_code: int,      # codul erorii sau zero pentru succes
    dest_group_id: int,   # identificatorul grupului destinatie
    source_group_id: int, # identificatorul grupului sursa
    len: int,             # lungimea partii variabile a mesajului
    mes[1:len]: char     # mesajul
)

chan local_chan[1..n] (
    # canal pentru comunicatie client -> server local
    kind: kind,          # tipul cererii
    client_id: int       # identificatorul clientului de la care s-a trimis cererea
    dest_group_id: int   # identificatorul grupului de la care s-a trimis cererea
    source_group_id: int # identificatorul grupului sursa care a trimis mesajul sau cererea
    group_name[1:MAX_LEN_GROUP]: char # nume de grup
    len: int,            # lungimea partii variabile a mesajului
    mes[1:len]: char    # mesajul
)

chan central_chan[1..n] (
    # canal pentru comunicatie server local -> server central
    kind: kind,          # tipul cererii
    client_id: int       # identificatorul clientului de la care s-a trimis cererea
    dest_group_id: int   # identificatorul grupului de la care s-a trimis cererea
    source_group_id: int # identificatorul grupului sursa care a trimis mesajul sau cererea
    group_name[1:MAX_LEN_GROUP]: char # nume de grup
    len: int,            # lungimea partii variabile a mesajului
    mes[1:len]: char    # mesajul
)

chan remote_chan[1..n] (
    # canal pentru comunicatie server central -> server local
    kind: kind,          # tipul cererii
    error_code: int,     # codul erorii sau zero pentru succes
    client_id: int       # identificatorul clientului de la care s-a trimis cererea
    dest_group_id: int   # identificatorul grupului destinatie
    source_group_id: int # identificatorul grupului sursa
    len: int,            # lungimea partii variabile a mesajului
    mes[1:len]: char    # mesajul
)
```

```
)  
chan send chan[1..n] (  
    knd: kind, # canal pentru comunicare server central → server central  
                # tipul cererii, numai SEND_MESSAGE  
    error_code: int, # codul erorii sau zero pentru succes  
    client_id: int, # identificadorul clientului de la care s-a trimis cererea  
    dest_group_id: int, # identificadorul grupului destinatie  
    source_group_id: int, # identificadorul grupului sursa  
    len: int, # lungimea partii variabile a mesajului  
    mes[1:len]: char # mesaj  
)  
chan conn_req chan (  
    knd: kind, # canal pentru conectarea unui server local la serverul central master  
                # tipul cererii, numai CONNECTION_REQ  
    sender: int # numarul nodului care se conecteaza  
)  
chan temp chan (  
    error_code: int, # canal pentru conectarea unui client la serverul local  
                # codul erorii  
    client_id: int # identificadorul clientului  
)
```

Serverele locale au următoarea structură:

LocalServer[1..n]:

```
const MAX_NAME_LEN = ... # lungimea maxima a numelor clientilor sau a grupurilor  
const MAX_MES_LEN = ... # lungimea maxima a unui mesaj  
const MAX_CLIENTS_IN_GROUP = ... # numarul maxim de clienti dintr-un grup  
const MAX_MES_FOR_A_CLIENT = ... # numarul maxim de mesaje care pot fi buferate pentru un client  
const MAX_CONC_GROUPS = ... # numarul maxim de grupuri la care poate apartine un client  
const MAX_GROUPS = ... # numarul maxim de grupuri  
const MAX_CLI = MAX_CLIENTS_IN_GROUP * MAX_GROUPS # numarul maxim de clienti  
const MAX_PROTO = ... # lungimea maxima a unui sir care indica protocolul folosit pentru autorizare clienti  
const MAX_STRING = ... # lungimea maxima a unei parole folosite pentru autorizarea clientilor  
const HEADER_LEN = ... # lungimea header-ului mesajelor
```

```
type CLIENT_CONN = record of  
    (pid: int, # identificadorul unui client  
     port: int, # identificadorul canalului rezervat dinamic unui client  
     name[1:MAX_NAME_LEN]: char # numele clientului dat la conectare  
    )  
# structura care mentine identitatea clientilor  
type MES_TO_GROUP = record of # structura care mentine informatii relative la un mesaj receptionat  
    (mes[1:MAX_MES_LEN]: char, # mesaj  
     counter: int, # nr proceselor care partajeaza mesajul  
     source_gid: int, # identificadorul grupului sursa  
     dest_gid: int # identificadorul grupului destinatie  
    )  
type GROUP_OF_CLIENTS = record of # structura care mentine componenta unui grup  
    (gid: int, # identificadorul de grup  
     clients_no: int, # numarul de clienti din acest grup  
     pid_tab[1:MAX_CLIENTS_IN_GROUP]: int # identicatori clienti  
    )  
type MES_TO_CLIENT = record of # buffer de mesaje pentru un client  
    (pid: int, # identificadorul clientului  
     pmes[MAX_MES_FOR_A_CLIENT]: pointer to MES_TO_GROUPS # indica mesajul  
    )  
type WCLIENTS_FOR_GROUP = record of # structura care mentine identicatorii proceselor blocate.  
    # in asteptare de mesaje de la un grup  
    (gid: int, # identificadorul de grup  
     wclients_no: int, # numarul de procese in asteptare de mesaje de la grupul gid  
     wpid_tab[MAX_CLIENTS_IN_GROUP]: int # identicatorii proceselor in asteptare
```

Capitolul 5: Proiectarea și implementarea unor mecanisme suport pentru dezvoltarea de aplicații distribuite

```
)
type CLIENTS_TO_GROUPS = record of # structura care mentine grupurile la care apartine un client
  (pid :int                                     # identificatorul clientului
   gid_no :int                                  # numarul de grupuri la care apartine un client
   gid_tab[MAX_CONC_GROUPS]: int               # identificatorii grupurilor
  )

var groups_of_clients_tab[1:MAX_GROUPS]: GROUP_OF_CLIENTS # mentine grupurile active din nod
var mes_to_client_tab[1:MAX_CLI]: MES_TO_CLIENT           # buffere de mesaje pentru clienti
var client_conn_tab[MAX_CLI]: CLIENT_CONN                # mentine clientii activi si conexiunile lor
var wclients_for_group_tab[1:MAX_GROUP]: WCLIENT_FOR_GROUP # mentine identitatea clientilor blocati
var client_to_groups_tab[1: MAX_CLI]: CLIENT_TO_GROUPS # mentine grupurile tuturor clientilor
var knd:int, error_code:int, client_id:in, dest_group_id:int, source_group_id:int # folosite pentru citire mesaje
var mes_len: int, mes[1:MAX_MES]: char # memoreaza lungimea si partea variabila a unui mesaj
var name[1: MAX_NAME_LEN]:char          # folosita pentru numele clientilor si grupurilor (alias-uri)
var auto_proto[1:MAX_PROTO]:char # folosita pentru citirea protocolului utilizat pentru autorizarea clientilor
var auto_string[1:MAX_STRING]: char     # folosita pentru citirea unei parole
var mes_to_groups: MES_TO_GROUP         # folosita la alocarea mesajelor receptionate
var host_name[1: MAX_NAME_LEN]         # folosita la memorarea numelui propriului nod
var pmessage: pointer to MES_TO_GROUP   # variabila de lucru
var message[1:MAX_MES]:char, source_gid: int # variabile de lucru

initializari : determinarea numelui propriului nod in host_name si adresei sale
send conn_req_chan(CONNECTION_REQ, host_name, i) #cerere de conexiune la serverul central
receive remote_chan[i](knd, error_code, client_id, dest_group_id, source_group_id, len , message)
if ( knd ≠ CONNECTION_REQ) or (error_code ≠ SUCCESS) → exit fi
do true →
  if not empty(local_chan[i]) →
    receive local_chan[i] (knd, client_id, dest_group_id, source_group_id, name, mes_len , mes)
    if knd = INIT_TGC or client_id ∈ client_conn id →
      if knd = INIT_TGC →
        determina protocolul de autorizare client. auto_proto, si parola, auto_string din mes
        valideaza autorizarea clientului si pozitioneaza error_code
        if error_code≠SUCCESS →
          send temp_chan(ENO_AUTHORIZATION,0)
        [] error_code = SUCCESS →
          stabileste un identificator pentru client in client_id,
          ca fiind numarul primei intrari libere din client_conn_tab
          if client_id > MAX_CLIENTS → send temp_chan(ETOO_MANY_CLIENTS, 0)
          [] client_id ≤ MAX_CLIENTS →
            creaza o intrare clients_conn_tab[client_id]
            creaza o intrare mes_to_client_tab[client_id]
            send temp_chan(SUCCESS, client_id)
        fi
      fi
    [] knd = JOIN_GROUP →
      if client_to_groups_tab[client_id].gid_no > MAX_CONC_GROUPS →
        send client_chan[client_id] (ETOO_MANY_CONC_GROUPS, dest_group_id, 0, 0, NULL)
      [] client_to_groups_tab[client_id].gid_no ≤ MAX_CONC_GROUPS
      if dest_group_id ∈ group_of_clients_tab →
        if client_id ∉ group_of_clients_tab[dest_group_id] →
          if group_of_clients_tab[dest_group_id].clients_no > MAX_CLIENTS_IN_GROUP →
            send client_chan[client_id](EFULL_GROUP, dest_group_id, 0, 0, NULL)
          [] group_of_clients_tab[dest_group_id].clients_no ≤ MAX_CLIENTS_IN_GROUP →
            inscreaza client_id in group_of_clients_tab[dest_group_id]
            inscreaza dest_group_id in clients_to_group_tab
```

Capitolul 5: Proiectarea și implementarea unor mecanisme suport pentru dezvoltarea de aplicații distribuite

```
        send client_chan[client_id](SUCCESS, dest_group_id, 0, 0, NULL)
    fi
    [] client_id ∈ group_of_clients_tab{dest_group_id} →
        send client_chan[client_id](ECLIENT_ALREADY_CONNECTED, dest_group_id, 0, 0, NULL)
    fi
    [] dest_group_id ∈ group_of_clients_tab →
        send central_chan[i](JOIN_GROUP, client_id, dest_group_id, 0, NULL, 0, NULL)
    fi
fi
[] kno = LEAVE_GROUP →
    if dest_group_id ∉ group_of_clients_tab →
        send client_chan[client_id](WRONG_DEST_GROUP, dest_group_id, 0, 0, NULL)
    [] dest_group_id ∈ group_of_clients_tab →
        if client_id ∈ group_of_clients_tab{dest_group_id}.pid_tab →
            sterge client_id din group_of_clients_tab{dest_group_id}.pid_tab
            fa mesajele i st mes_to_clients_tab[client_id].pmes[i]^dest_gid = dest_group_id →
                pmessage := mes_to_clients_tab[client_id].pmes[i]
                mes_to_clients_tab[client_id].pmes[i] = NULL;
                pmessage^.counter := pmessage^.counter - 1 # un mesaj in asteptare mai puțin
            if pmessage^.counter = 0 → dezaloca structura de la pmessage fi
        fa
        [] client_id ∉ group_of_clients_tab{dest_group_id}.pid_tab →
            send client_chan[client_id](ECLIENT_NOT_IN_DEST_GROUP, dest_group_id, 0, 0, NULL)
        fi
fi
[] kno = CREATE_ALIAS_GROUP or kno = DESTROY_ALIAS_GROUP or
    kno = GET_ID_GROUP or kno = CREATE_GROUP →
    send central_chan[i](kno, client_id, dest_group_id, 0, name, 0, NULL)
[] kno = SEND_TO_GROUP →
    if source_group_id ∉ group_of_clients_tab →
        send client_chan[client_id](EWRONG_SOURCE_GROUP, # se incearca o emisie dintr-un grup
            dest_group_id, source_group_id, 0, NULL) # sursa inexistent
    [] source_group_id ∈ group_of_clients_tab →
        if client_id ∉ group_of_clients_tab[source_group_id].pid_tab →
            send client_chan[client_id](ECLIENT_NOT_IN_SOURCE_GROUP, # se incearca o emisie
                dest_group_id, 0, 0, NULL) # in contul unui client neconectat la sursa

        [] client_id ∈ group_of_clients_tab[source_group_id].pid_tab →
            formeaza la message un mesaj din adresa proprie IP, client_conn_tab[client_id].name, mes
            send central_chan[i](SEND_TO_GROUP, client_id, dest_group_id, source_group_id, NULL,
                mes.len + HEADER_LEN, message)
        fi
    fi
[] kno = RECEIVE_FROM_GROUP INFINITE →
    if (source_group_id > MAX_GROUPS) →
        send client_chan[client_id](EWRONG_SOURCE_GROUP, # se incearca receptia de la un grup
            dest_group_id, source_gid, mes.len, message) # inexistent
    [] (source_group_id ≤ MAX_GROUPS) →
        if dest_group_id ∉ group_of_clients_tab →
            send client_chan[client_id](EWRONG_DEST_GROUP # se incearca o receptie in contul unui grup
                dest_group_id, source_group_id, 0, NULL) # care nu exista
        [] dest_group_id ∈ group_of_clients_tab →
            if client_id ∉ group_of_clients_tab{dest_group_id}.pid_tab →
                send client_chan[client_id](ECLIENT_NOT_IN_DEST_GROUP, # se incearca o receptie in
```

Capitolul 5: Proiectarea și implementarea unor mecanisme suport pentru dezvoltarea de aplicații distribuite

```
dest_group_id, 0, 0, NULL) # in contul unui client neconectat la grupul destinatie
[] client_id ∈ group_of_clients_tab[dest_group_id].pid_tab →
  if exista mesaje i in mes_to_client_tab[client_id] and
    mes_to_clients_tab[client_id].pmes[i]^dest_gid = dest_group_id and
    (source_group_id = ANY or
      (mes_to_clients_tab[client_id].pmes[i]^source_gid = source_group_id)) →
    extrage primul mesaj i in message si calculeaza lungimea sa in mes_len
    pmessage := mes_to_clients_tab[client_id].pmes[i]
    mes_to_clients_tab[client_id].pmes[i] = NULL;
    source_gid = pmessage^.source_gid
    pmessage^.counter := pmessage^.counter - 1 # un mesaj in asteptare mai putin
  if pmessage^.counter = 0 → deazaloca structura de la pmessage fi
  if (mai exista mesaje i in mes_to_client_tab[client_id]) and
    mes_to_clients_tab[client_id].pmes[i]^dest_gid = dest_group_id and
    ((source_group_id = ANY) or
      (mes_to_clients_tab[client_id].pmes[i]^source_gid = source_group_id)) →
    send client_chan[client_id](MORE_MES, dest_group_id, source_gid,
      mes_len, message) # se receptioneaza mesajul si o indicatie ca mai exista mesaje
  [] nu mai exista mesaje in asteptare de la source_group_id catre dest_group_id →
    send client_chan[client_id](SUCCESS, dest_group_id, source_gid, mes_len, message)
    #se receptioneaza mesajul si o indicatie ca nu mai exista mesaje
  fi
[]nu exista mesaje in asteptare de la source_group_id catre dest_group_id pentru client_id →
  if wclients_for_groups[source_group_id].wclients_no>MAX_CLIENTS_IN_GROUP →
    send client_chan[client_id](ENO_BUF, dest_group_id, source_gid, mes_len, message)
  [] wclients_for_groups[source_group_id].wclients_no ≤ MAX_CLIENTS_IN_GROUP →
    insereaza client_id in wclients_for_groups[source_group_id].wpid.tab
    si actualizeaza wclients_for_groups[source_group_id].wclients_no
    if source_group_id ∉ group_of_clients_tab →
      send central_chan[i](RECEIVE_FROM_GROUP_INFINITE, client_id,
        dest_group_id, source_group_id, NULL, 0, NULL)
        # pentru verificare source_group_id
    fi
  fi
fi
fi
fi
fi
[] knd = NO MORE RECEIVE FROM GROUP →
  sterge client_id wclients_for_groups[source_group_id].wpid.tab
  si actualizeaza wclients_for_groups[source_group_id].wclients_no
[] knd = RECEIVE FROM GROUP NOW →
  if (source_group_id > MAX_GROUPS) →
    send client_chan[client_id](EWRONG_SOURCE_GROUP, # se incearca receptia de la un grup
      dest_group_id, source_gid, mes_len, message) # inexistent
  [] (source_group_id ≤ MAX_GROUPS) →
    if dest_group_id ∉ group_of_clients_tab →
      send client_chan[client_id](EWRONG_DEST_GROUP #se incearca o receptie in contul unui grup
        dest_group_id, source_group_id, 0, NULL)# care nu exista
    [] dest_group_id ∈ group_of_clients_tab →
      if client_id ∉ group_of_clients_tab[dest_group_id].pid_tab →
        send client_chan[client_id](ECLIENT_NOT_IN_DEST_GROUP, # se incearca o receptie in
          dest_group_id, 0, 0, NULL) # in contul unui client neconectat la grupul destinatie
      [] client_id ∈ group_of_clients_tab[dest_group_id].pid_tab →
```



```

if exista mesaje i in mes to client tab[client id] and
    mes to clients tab[client id].pmes[i].dest_gid = dest_group id and
    (source_group id = ANY or
        (mes to clients tab[client id].pmes[i].source_gid = source_group id)) →
    extrage primul mesaj i in message și calculează lungimea sa in mes len
    pmessage := mes to clients tab[client id].pmes[i]
    mes to clients tab[client id].pmes[i] = NULL;
    source_gid = pmessage.source_gid
    pmessage.counter := pmessage.counter - 1 # un mesaj în așteptare mai puțin
    if pmessage.counter = 0 → deașază structura de la pmessage fi
if (mai exista mesaje i in mes to client tab[client id]) and
    mes to clients tab[client id].pmes[i].dest_gid = dest_group id and
    ((source_group id = ANY) or
        (mes to clients tab[client id].pmes[i].source_gid = source_group id)) →
    send client chan[client id](MORE_MES, dest_group id, source_gid,
        mes len, message) # se recepționează mesajul și o indicație că mai există mesaje
    [] nu mai există mesaje în așteptare de la source_group id către dest_group id →
    send client chan[client id](SUCCESS, dest_group id, source_gid, mes len, message)
    # se recepționează mesajul și o indicație că nu mai există mesaje
fi
    [] nu există mesaje în așteptare de la source_group id către dest_group id →
    send client chan[client id](NO_WAIT_MESSAGES, dest_group id, source_gid,
        0, NULL) # nu există mesaje pentru recepție
fi
fi
fi
fi
[] knd = GET_NO_OF_MESSAGES →
if (source_group id > MAX_GROUPS) →
    send client chan[client id](EWRONG_SOURCE_GROUP, # se încearcă recepția de la un grup
        dest_group id, source_gid, mes len, message) # inexistent
    [] (source_group id ≤ MAX_GROUPS) →
    if dest_group id ∉ group of clients tab →
        send client chan[client id](EWRONG_DEST_GROUP # se încearcă o recepție în conținutul unui grup
            dest_group id, source_group id, 0, NULL) # care nu există
        [] dest_group id ∈ group of clients tab →
            if client id ∉ group of clients tab[dest_group id].pid tab →
                send client chan[client id](ECLIENT_NOT_IN_DEST_GROUP, # se încearcă o recepție în
                    dest_group id, 0, 0, NULL) # în conținutul unui client neconectat la grupul destinat
            [] client id ∈ group of clients tab[dest_group id].pid tab →
                contorizează în k mesajele i cu
                mes to client[client id].pmes[i].source_gid = source_group id și
                mes to client[client id].pmes[i].dest_gid = dest_group id
                send client chan[client id](SUCCESS, dest_group id, source_group id, k, NULL)
            fi
        fi
    fi
    [] knd = CLOSE_TGC →
    șterge intrarea client id din client conn tab[client id]
    fa toate grupurile la care este conectat client id, client to groups[client id].gid.tab[i] →
    source_gid := client to groups[client id].gid.tab[i]
    șterge client id din group of clients tab[source_gid].gid tab[i]
    fa mesajele i și mes to clients tab[client id].pmes[i].dest_gid := source_gid →
    pmessage := mes to clients tab[client id].pmes[i]
    mes to clients tab[client id].pmes[i] = NULL;

```

apitolul 5: Proiectarea și implementarea unor mecanisme suport pentru dezvoltarea de aplicații distribuite

```

    pmessage^.counter := pmessage^.counter - 1      # un mesaj in asteptare mai putin
    if pmessage^.counter = 0 → dezaloca structura de la pmessage fi
fa
af
fi
fi
fi
if not empty(remote_chan[i]) →
    receive remote_chan[i] (knd, error_code, client_id, dest_group_id, source_group_id, mes_len, message)
    if knd = CREATE_GROUP, JOIN_GROUP →
        if error_code = SUCCESS →
            adauga o intrare in group_of_clients[dest_group_id] care contine client_id
            insereaza dest_group_id in clients_to_group_tab[client_id]
        || error_code ≠ SUCCESS →
            send_client_chan[client_id](error_code, dest_group_id, 0, 0, NULL)
        fi
    || knd = CREATE_ALIAS_GROUP or knd = DESTROY_ALIAS_GROUP
        or knd = GET_ID_GROUP or knd = SEND_TO_GROUP →
            send_client_chan[client_id](error_code, dest_group_id, 0, 0, NULL)
    || knd = RECEIVE_FROM_GROUP INFINITE →
        if error_code = SUCCESS →
            sterge_client_id din wclients_for_groups[source_group_id].wpid.tab
            send_client_chan[client_id](error_code, dest_group_id, source_group_id, 0, NULL)
        fi
    || knd = SEND_MESSAGE →
        if exista_client_id astfel incit (client_id ∈ dest_group_id) and
            (client_id ∉ wclients_for_group[source_group_id].wpid.tab) →
            se aloca o zona de memorie pentru o structura MES_TO_GROUP. pmessage
            pmessage^.mes := mes
            pmessage^.source_gid := source_group_id
            pmessage^.dest_gid := dest_group_id
            se contorizeaza in k toti membrii grupului dest_group care nu sunt in
                wclients_for_group[source_group_id].wpid.tab
            pmessage^.counter := k      fi
        fa client_id ∈ dest_group_id st client_id ∉ wclients_for_group[dest_group_id].wpid.tab
            insereaza pmessage in mes_to_client_tab[client_id].pmes
        fa
        fa (client_id ∈ dest_group_id) st (client_id ∈ wclients_for_group[source_group_id].wpid.tab)
            send_client_chan[client_id](SUCCESS, dest_group_id, source_group_id, mes_len, mes)
            #se receptioneaza mesajul si o indicatie ca nu mai exista mesaje
            sterge_client_id din wclients_for_groups[source_group_id].wpid.tab
        af
    fi
fi
od

```

Serverele centrale pot fi descrise astfel:

```

const MAX_HOSTS_FOR_A_GROUP # numarul maxim de noduri pe care poate fi deschis un grup
const MAX_HOSTS = ...      # numarul maxim de noduri care se pot conecta la sistem
const MAX_NAME_LEN = ...   # lungimea maxima a numelor clientilor sau a grupurilor
const MAX_ALIASES_FOR_A_GROUP = ... # numarul maxim de pseudonime pentru un grup
const MAX_HOSTS_NAME = ... # lungimea maxima a numelui unui nod
const MAX_ADDRESSE_LEN = ... # lungimea maxima a adreselor unui nod

```

```

type HOSTS_FOR_GROUP = record of # structura care mentine nodurile pe care s-a deschis grupul gid
    (gid: int, # identicatorul de grup
     hosts_no: int, # numarul de noduri pe care exista membri ai grupului gid
     port_no[MAX_HOSTS_FOR_A_GROUP]: int # porturile canalelor proceselor CentralServer
    ) # care controleaza comunicatia cu nodurile grupului
type GROUP_ALIASES = record of # structura care mentine pseudonimele grupului gid
    (gid: int,
     alias_no: int, # numarul de pseudonime acordate grupului gid
     alias_tab[MAX_ALIASES_FOR_A_GROUP][MAX_NAME_LEN]: char # pseudonime
    )
type HOSTS = record of # structura care mentine informatii despre noduri
    (name[MAX_HOSTS_NAME]: char,
     adresse[MAX_ADRESSE_LEN]: char,
     port: int; # porturile canalelor proceselor CentralServer
    ) # care controleaza comunicatia cu nodurile grupului

```

```

var hosts_for_group_tab[MAX_GROUPS]: HOSTS_FOR_GROUP #mentine repartitia grupurilor pe noduri
var group_aliases_tab[MAX_GROUPS]: GROUP_ALIASES # mentine pseudonimele grupului
var hosts_tab[MAX_HOSTS]: HOSTS # mentine informatii despre noduri

```

MasterCentralServer::

```

var knd: kind, host_name[MAX_NAME_LEN]: char # pentru citirea mesajelor
initializari
do true →
    receive conn_req_chan(knd, host_name, i) # cerere de conexiune primita de la un server local i
    if (knd ≠ CONNECTION_REQ)→
        send remote_chan[i](CONNECTION_REQ, BAD_REQUEST, 0, 0, 0, 0, NULL )
    [] if nu exista o intrare libera in hosts_tab →
        send remote_chan[i](CONNECTION_REQ, TOO_MANY_HOSTS, 0, 0, 0, 0, NULL )
    [] knd = CONNECTION_REQ and exista o intrare libera in hosts_tab →
        creaza si completeaza o intrare i in hosts_tab
        creaza un proces CentralServer si transmite ca parametru hosts_tab[i]
od

```

CentralServer::

```

var knd: int, error_code: int, client_id: int # folosite pentru mesaje
var dest_group_id: int, source_group_id: int, gid: int # folosite pentru mesaje
var name[1:MAX_NAME_LEN]:char,mes_len: int, message[1:MAX_MES]: char #variabile de lucru
do true →
    if not empty(central_chan[i]) →
        receive central_chan[i](knd, client_id, dest_group_id, source_group_id, name, mes_len, message)
        if knd = CREATE_GROUP →
            if exista un grup, gid pentru care name ∈ group_aliases_tab[gid]→
                send remote_chan[i](knd,EGROUP_ALREADY_CREATED, client_id, 0, 0, 0, NULL)
            [] nu exista un grup, gid pentru care name ∈ group_aliases_tab[gid]→
                if nu exista o intrare o intrare libera in hosts_for_group_tab →
                    send remote_chan[i](knd,TOO_MANY_GROUPS, client_id, 0, 0, 0, NULL)
                [] exista o intrare o intrare libera in hosts_for_group_tab →
                    creaza si completeaza o intrare gid in hosts_for_groups_tab si in group_aliases_tab
                    send remote_chan[i](CREATE_GROUP, SUCCESS, client_id, gid, 0, 0, NULL)
                fi
            fi
        [] knd = GET_ID_GROUP →
            if nu exista un grup, gid pentru care name ∈ group_aliases_tab[gid]→

```

Capitolul 5: Proiectarea și implementarea unor mecanisme suport pentru dezvoltarea de aplicații distribuite

```
    send remote_chan[i] (GET_ID_GROUP, EWRONG_ALIAS, client_id, 0, 0, 0, NULL)
  || exista un grup, gid pentru care name ∈ group_aliases_tab[gid] →
    send remote_chan[i] (GET_ID_GROUP, SUCCESS, client_id, gid, 0, 0, NULL)
  fi
|| knđ = CREATE_ALIAS_GROUP →
  if dest_group_id ∉ hosts_for_group_tab →
    send remote_chan[i] (knd, EWRONG_DEST_GROUP, client_id, dest_group_id, 0, 0, NULL)
  || dest_group_id ∈ hosts_for_group_tab →
    if exista un grup, gid pentru care name ∈ group_aliases_tab[gid] →
      send remote_chan[i] (CREATE_ALIAS_GROUP, EWRONG_ALIAS, client_id, 0, 0, 0, NULL)
    || nu exista un grup, gid pentru care name ∈ group_aliases_tab[gid] →
      if group_aliases_tab[dest_group_id].alias_no > MAX_ALIASES_FOR_A_GROUP →
        send remote_chan[i] (knd, ETOO_MANY_ALIASES, client_id, dest_group_id, 0, 0, NULL)
      || group_aliases_tab[dest_group_id].alias_no ≤ MAX_ALIASES_FOR_A_GROUP →
        adauga nume in group_aliases_tab[dest_group_id]
        send remote_chan[i] (CREATE_ALIAS_GROUP, SUCCESS, client_id, gid, 0, 0, NULL)
      fi
    fi
  fi
|| knđ = DESTROY_ALIAS_GROUP →
  if dest_group_id ∉ hosts_for_group_tab →
    send remote_chan[i] (knd, EWRONG_DEST_GROUP, client_id, dest_group_id, 0, 0, NULL)
  || dest_group_id ∈ hosts_for_group_tab →
    if name ∉ group_aliases_tab[dest_group_id] →
      send remote_chan[i] (DESTROY_ALIAS_GROUP, EWRONG_ALIAS, client_id, 0, 0, 0, NULL)
    || name ∈ group_aliases_tab[dest_group_id] →
      sterge nume din group_aliases_tab[dest_group_id]
      send remote_chan[i] (DESTROY_ALIAS_GROUP, SUCCESS, client_id, 0, 0, 0, NULL)
    fi
  fi
|| knđ = JOIN_GROUP →
  if dest_group_id ∉ hosts_for_group_tab →
    send remote_chan[i] (knd, EWRONG_DEST_GROUP, client_id, dest_group_id, 0, 0, NULL)
  || dest_group_id ∈ hosts_for_group_tab →
    if hosts_for_group_tab[dest_group_id].host_no ≥ MAX_HOSTS_FOR_A_GROUP →
      if portul nodului apelant nu exista in hosts_for_group_tab[dest_group_id].port_no →
        adauga portul nodului in intrarea grupului din hosts_for_group_tab[dest_group_id].port_no
        si actualizeaza hosts_for_group_tab[dest_group_id].host_no
      fi
      send remote_chan[i] (JOIN_GROUP, SUCCESS, client_id, dest_group_id, 0, 0, NULL)
    || hosts_for_group_tab[dest_group_id].host_no < MAX_HOSTS_FOR_A_GROUP →
      send remote_chan[i] (JOIN_GROUP, ETOO_MANY_HOSTS_IN_GROUP,
        client_id, dest_group_id, 0, 0, NULL)
    fi
  fi
|| knđ = LEAVE_GROUP →
  if dest_group_id ∈ hosts_for_group_tab →
    sterge din hosts_for_group_tab[dest_group_id].port_no portul nodului
    si actualizeaza hosts_for_group_tab[dest_group_id].host_no
    if hosts_for_group_tab[dest_group_id].host_no = 0 →
      sterge intrarile hosts_for_group_tab[dest_group_id] si group_aliases_tab[dest_group_id]
    fi
  fi
  send remote_chan[i] (LEAVE_GROUP, SUCCESS, client_id, dest_group_id, 0, 0, NULL)
|| knđ = RECEIVE_FROM_GROUP_INFINITE →
```

```

if source_group_id ∉ hosts_for_group_tab →
    send remote_chan[r] (knd, WRONG_SOURCE_GROUP, client_id, dest_group_id,
                        source_group_id, 0, NULL) # grup sursa inexistent
|| source_group_id ∈ hosts_for_group_tab →
    send remote_chan[r] (knd, SUCCESS, client_id, dest_group_id, source_group_id, 0, NULL)
fi
|| knd SEND_TO_GROUP →
if dest_group_id ∉ hosts_for_group_tab →
    send remote_chan[r] (knd, WRONG_DEST_GROUP, client_id, dest_group_id, 0, 0, NULL)
|| dest_group_id ∈ hosts_for_group_tab →
    send remote_chan[r] (knd, SUCCESS, client_id, dest_group_id, 0, 0, NULL)
    fa j = 1 to hosts_for_group_tab[dest_group_id].host_no →
        send send_chan[hosts_for_group_tab[dest_group_id].port_no[j]] (SEND_MESSAGE,
                                SUCCESS, client_id, dest_group_id, source_group_id, mes_len, message)
    af
fi
fi
if not empty(local_chan[r]) →
    receive send_chan[r] (knd, error_code, client_id, dest_group_id, source_group_id, name, mes_len, message)
    if knd SEND_MESSAGE →
        send remote_chan[r] (knd, SUCCESS, client_id, dest_group_id, source_group_id, mes_len, message)
fi
od

```

Funcțiile puse la dispoziția clienților pentru aplicații pot fi descrise astfel:

```

const MAX_NAME_LEN = # lungimea maxima a numelor clientilor sau a grupurilor
const MAX_PROTO = # lungimea maxima a unui sir care indica protocolul folosit pentru autorizare clienti
const MAX_STRING = # lungimea maxima a unei parole folosite pentru autorizarea clientilor
const HEADER_LEN = # lungimea header-ului mesajului
var client_id: int := 0 # identificadorul clientului returnat de serverul local la conectarea clientului
# trebuie transmis in continuare in corpul fiecărei cereri
var error_code: int, dest_group_id: int, source_group_id: int # folosite pentru mesaje
var mes_len: int, message[1:MAX_LEN]: char #variabile de lucru

```

```

InitTGC(user[MAX_NAME_LEN]: char, auto_proto[MAX_PROTO]: char, passwd[MAX_STRING]: char)::
    if user = NULL → return EBADF fi
    if client_id ≠ 0 →
        send local_chan[r] (CLOSE_TGC, client_id, 0, 0, NULL, 0, NULL)
        client_id = 0
    fi
    send local_chan[r] (INIT_TGC, 0, 0, 0, user, length(auto_proto + passwd), auto_proto + passwd)
    receive temp_chan (error_code, client_id)
    return error_code

```

```

CreateGroup(var pgroup: int, name[MAX_NAME_LEN]: char)::
    if client_id = 0 → return CLIENT_NOT_CONNECTED
    if name = NULL → return EBADF fi
    send local_chan[r] (CREATE_GROUP, client_id, 0, 0, name, 0, NULL)
    receive client_chan[client_id] (error_code, dest_group_id, source_group_id, mes_len, message)
    if error_code ≠ SUCCESS pgroup := dest_group_id
    return error_code

```

JoinGroup(group_id: int)::

```
if client_id = 0 → return ECLIENT_NOT_CONNECTED
send local_chan[i](JOIN_GROUP, client_id, group_id, 0, NULL, 0, NULL)
receive client_chan[client_id] (error_code, dest_group_id, source_group_id, mes_len, message)
return error_code
```

LeaveGroup(group_id: int)::

```
if client_id = 0 → return ECLIENT_NOT_CONNECTED
send local_chan[i](LEAVE_GROUP, client_id, group_id, 0, NULL, 0, NULL)
receive client_chan[client_id] (error_code, dest_group_id, source_group_id, mes_len, message)
return error_code
```

CreateAliasGroup(group: int, alias[MAX_NAME_LEN]: char)::

```
if client_id = 0 → return ECLIENT_NOT_CONNECTED
if alias = NULL → return EBADF fi
send local_chan[i](CREATE_ALIAS_GROUP, client_id, group_id, 0, name, 0, NULL)
receive client_chan[client_id] (error_code, dest_group_id, source_group_id, mes_len, message)
return error_code
```

DestroyAliasGroup(group: int, alias[MAX_NAME_LEN]: char)::

```
if client_id = 0 → return ECLIENT_NOT_CONNECTED
if alias = NULL → return EBADF fi
send local_chan[i](DESTROY_ALIAS_GROUP, client_id, group_id, 0, name, 0, NULL)
receive client_chan[client_id] (error_code, dest_group_id, source_group_id, mes_len, message)
return error_code
```

GetIdGroup(var pgroup: int, alias[MAX_NAME_LEN]: char)::

```
if client_id = 0 → return ECLIENT_NOT_CONNECTED
if alias = NULL → return EBADF fi
send local_chan[i](GET_ID_GROUP, client_id, 0, 0, name, 0, NULL)
receive client_chan[client_id] (error_code, dest_group_id, source_group_id, mes_len, message)
if error_code ≠ SUCCESS pgroup := dest_group_id
return error_code
```

**SendToGroup(dest_group_id: int, source_group_id: int, mes_len: int, message[MAX_MES]: char,
counter: int, pfilter: pointer to function)::**

```
if client_id = 0 → return ECLIENT_NOT_CONNECTED
if counter > MAX_MES → return EMES_TO_LONG
pfilter(message) # se codifica mesajul
mes_len := counter
send local_chan[i](SEND_TO_GROUP, client_id, dest_group_id, source_group_id, mes_len, message)
receive client_chan[client_id] (error_code, dest_group_id, source_group_id, mes_len, message)
return error_code
```

**ReceiveFromGroup(source_group_id: int, dest_group_id: int, var pcounter: int, var pmes[MAX_MES]:
char, time_limit: int, pfilter: pointer to function)::**

```
if client_id = 0 → return ECLIENT_NOT_CONNECTED
if time_limit = NOW → mode = RECEIVE_FROM_GROUP_NOW
[] time_limit = INFINITE → mode = RECEIVE_FROM_GROUP_INFINITE
fi
if time_limit = NOW or time_limit = INFINITE →
i := 0; error_code := MORE_MES
do (i < pcounter) and (error_code = MORE_MES or (error_code = SUCCESS and time_limit = INFINITE))
send local_chan[i](mode, client_id, dest_group_id, source_group_id, NULL, 0, NULL)
```

Capitolul 5: Proiectarea și implementarea unor mecanisme suport pentru dezvoltarea de aplicații distribuite

```
    receive client_chan[client_id] (error_code, dest_group_id, source_group_id, mes_len, message)
    if (error_code ≠ MORE_MES) or (error_code ≠ SUCCESS) →
        pcounter := i; return error_code
    fi
    pfilter(message + HEADER_LEN)
    pmes := pmes + mes_len; memoreaza message la pmes
    i := i + 1
od
[] time_limit ≠ NOW and time_limit ≠ INFINITE →
    i := 0; error_code := MORE_MES
    do (i < pcounter) and (error_code = MORE_MES or (error_code = SUCCESS))
        send local_chan[i](RECEIVE_FROM_GROUP_INFINITE, client_id, dest_group_id,
            source_group_id, NULL, 0, NULL)
        programeaza un semnal ALARM pentru n secunde, a carui rutina seteaza pe error_code = ALARM
        receive client_chan[client_id] (error_code, dest_group_id, source_group_id, mes_len, message)

        if error_code ≠ MORE_MES or error_code ≠ SUCCESS →
            pcounter := i; return error_code
        fi
        pfilter(message + HEADER_LEN)
        i := i + 1
    od
fi
pcounter := i;
return error_code
```

GetNoOfMessages(source_group_id: int, dest_group_id: int, var pcounter: int)::

```
if client_id = 0 → return ECLIENT_NOT_CONNECTED
send local_chan[i](GET_NO_OF_MESSAGES, client_id, dest_group_id, source_group_id, NULL, 0, NULL)
receive client_chan[client_id] (error_code, dest_group_id, source_group_id, mes_len, message)
if error_code ≠ SUCCESS → pcounter := 0; return error_code
[] pcounter := mes_len; return SUCCESS
fi
```

CloseTGC()::

```
if client_id = 0 → return ECLIENT_NOT_CONNECTED
send local_chan[i](CLOSE_TGC, 0, 0, 0, NULL, 0, NULL)
receive client_chan[client_id] (error_code, dest_group_id, source_group_id, mes_len, message)
return error_code
```


În descrierea de mai sus, nu a fost prezentată soluția adoptată pentru respectarea atomicității în transmitia mesajelor. Deoarece protocolul folosit în rețeaua de comunicații este protocolul TCP, singura prevedere ce trebuie luată este cea împotriva situației de "buffere ocupate"; aceasta înseamnă că, de fapt, trebuie asigurat controlul fluxului mesajelor la nivelul aplicației. Mai precis, controlul fluxului mesajelor trebuie asigurat între serverul central și serverele locale, la execuția unei operații *SEND_MESSAGE*. Soluția adoptată este următoarea: când, în urma primirii unui mesaj, un server local detectează că se ocupă toate bufferele sale și, în viitor, nu mai poate astfel primi nici un mesaj pentru un anumit grup, deoarece nu va mai avea unde să-l memoreze, informează serverul central despre acest lucru, transmițându-i un mesaj *FULL_BUFFERS*. Serverul central menține în tabloul *hosts_for_groups_tab[MAX_GROUPS]*, de tipul *HOSTS_FOR_GROUP* încă un cîmp, care arată dacă umplerea bufferelor pentru un nod este iminentă. Dacă, atunci când trebuie trimis un mesaj la grupul respectiv, se găsește pentru un nod în acest cîmp valoarea 1, serverul global returnează eroarea *ENO_BUFS*.

Tabloul *hosts_for_groups_tab[MAX_GROUPS]* se actualizează atunci, când, în nodul în cauză, procesele client recepționează mesaje de la acel grup, eliberînd astfel buffere; în această situație, serverul local transmite serverului global mesajul *FREE_BUFFER*.

5.2.3. Alte variante

Varianta prezentată la 5.2.2. se poate caracteriza ca o variantă semidistribuită: există câte un server local în fiecare nod și un server global unic, cu rol de proces coordonator, situat pe un nod fixat în rețea.

În această secțiune vor fi prezentate și alte soluții de proiectare a unui toolkit pentru comunicații de grup, menținînd ca valabile următoarele cerințe de proiectare: (1) sistemul TGC trebuie să se execute în mod utilizator. Singura excepție poate fi legată de faptul că modul broadcast în Unix poate fi o operație privilegiată. Este preferabil însă ca sistemul să se execute în mod utilizator, putînd fi astfel startat și de utilizatori obișnuiți. Totuși, unele componente, (eventual cele legate de stabilirea de parole pentru clienți), dacă este necesar, pot să fie prevăzute cu unele restricții de execuție. (2) Protocolul de comunicație trebuie să fie independente de hardware.

5.2.3.1. Variante complet distribuite

O variantă complet distribuită prezintă față de varianta semidistribuită avantajul eliminării serverului coordonator, dar și dezavantajul creșterii încărcării rețelei. Există două soluții: folosirea modului de lucru broadcast, prin utilizarea unui protocol adecvat, și folosirea unui protocol punct-la-punct.

Arhitectura generală a sistemului, în ambele cazuri, va fi aceeași; va exista cite un server general în fiecare nod. Acesta va menține aceleași structuri de date cu cele ale unui server local din varianta semidistribuită, dar va îngloba și structurile de date ale serverului central din varianta semidistribuită.

Prin urmare, va conține atît structurile de date necesare menținerii evidenței atașării clienților la grupuri dar și numele alias asociate tuturor grupurilor. În ceea ce privește structurile de date care mențin legătura dintre grupuri și nodurile pentru care există clienți atașați la aceste grupuri, ele sunt necesare numai în varianta în care nu se utilizează transmisie broadcast. Va fi necesar însă, în ambele variante introducerea mărcilor de timp în corpul mesajelor, pentru a ordona transmisiile în sistem.

Pentru tratarea transmisiei și recepționării de mesaje concurentă cu execuția operațiilor și pentru modularizarea implementării se separă protocolul de comunicație într-un proces separat, care livrează cererile de operații confirmate unui proces sever, ce tratează cererile. Arhitectura generală a sistemului este prezentată în fig. 5.2.3.1.

A. Varianta în care se folosește transmisia broadcast

Protocolul folosit va fi protocolul UDP.

Mecanismul confirmărilor va fi construit asemănător cu cel prezentat la 4.3.4.1. în cadrul implementării semafoarelor distribuite. Vor fi tratate de fiecare nod numai mesajele *complet confirmate*. Dacă o cerere nu este complet confirmată dar, în schimb, este urmată de una sau mai multe cereri complet confirmate, după citeva mesaje de probă cu nodul care nu a confirmat cererea (neconfirmate și ele) nodul respectiv se deconectează din sistem și nu va mai fi luat în considerare pentru contorizarea confirmărilor decît după o nouă cerere de conectare a sa.

Nu toate cererile de operații trebuie transmise în rețea. Operațiile care nu trebuie transmise broadcast în rețea sunt: *IniTGC()*, *GetIdGroup()*, *JoinGroup()*, *ReceiveFromGroup()*, *GetNoOfMessages()*, *CloseTGC()*; în plus, operația *LeaveGroup()* nu trebuie transmisă în rețea decît dacă generează, în urma execuției ei, ștergerea unui grup; evident, toate operațiile sunt transmise în rețea numai dacă nu generează, în urma execuției, un cod de eroare de tipul grup inexistent, client neatașat la grup, format eronat al parametrului etc.

Structura unui pachet care transmite o operație în rețea va fi:

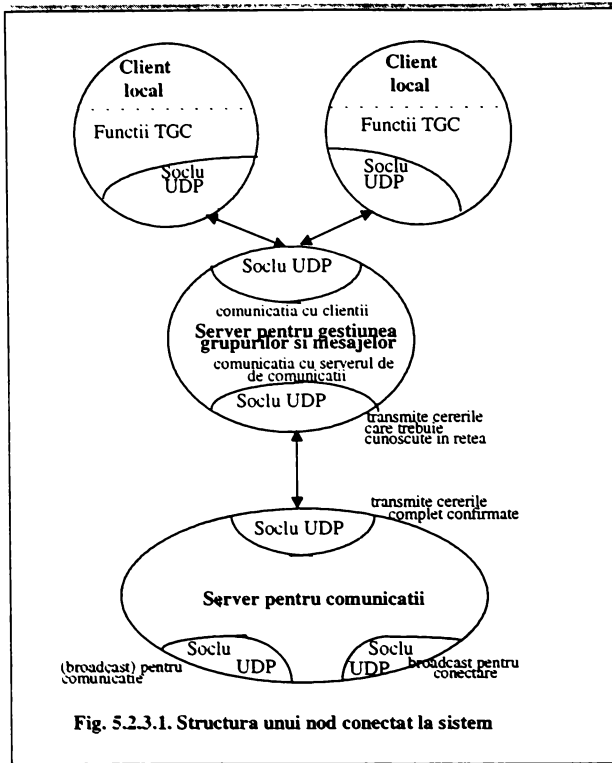


Fig. 5.2.3.1. Structura unui nod conectat la sistem

```

/* toate pachetele au un header de identificare */
LENGTHS kind ; /* tipul cererii */
ADDRESS sender; /* identificatorul emitatorului */
TS-TYPE time-stamp ; /* marca de timp */
LENGTHS duplicate; /* numarul de retransmisie al mesajului */
/* partea mesajii primita de la serverul pentru gestiunea grupurilor */
MESSAGE variable-part ;
    
```

Fiecare pachet are un header format din: (1) o zonă de dimensiune constantă, obligatorie pentru toate tipurile de cereri, care conține un octet pentru identificarea tipului cererii, identificatorul emițătorului și marca de timp și (2) mesajul propriu zis care conține și o zonă de dimensiune variabilă (în cazul operației *SendToGroup()* conține mesajul de transmis iar în cazul operațiilor *CreateGroup()*, *CreateAliasGroup()*, *DestroyAliasGroup()* numele de grup implicat în operație).

O atenție specială trebuie acordată operației *CreateGroup()*, care trebuie executată în regim de excludere reciprocă, pentru ca fiecare nod să poată alege un identificator de grup unic la crearea unui grup nou.

Atomicitatea transmisiilor se asigură prin mecanismul special al confirmărilor. Se presupune că mesajele și confirmările primite de un nod sunt identice cu cele primite de toate nodurile. Mesajele nu se pierd decât în patru situații: când un nod are bufferele ocupate, când un nod nu mai răspunde la nici un mesaj, când se constată lipsa unui mesaj sau a unei confirmări (în toate nodurile se va observa același lucru). Prima situație se poate rezolva transmiiind totuși un mesaj de confirmare dar cu un conținut care să indice eroarea *ENOBUFS*. Cea de-a doua situație poate fi detectată ușor, deoarece, dacă un nod nu mai este activ în sistem, toate mesajele următoare vor fi neconfirmate complet. În acest caz, se deconectează automat nodul respectiv din sistem. Pentru ultimele două situații, mesajele neconfirmate complet, urmate de altele confirmate complet sunt ignorate de toate nodurile. Pentru fiecare mesaj trimis în sistem, funcțiile TGC așteaptă să recepționeze un cod de eroare care arată modul în care s-a desfășurat execuția: cu succes sau cu o eroare. În ultimul caz, se reîncearcă transmisia de câteva ori înainte de a se un răspuns utilizatorului.

Descrierea serverului de comunicații în această variantă este prezentată la 5.2.4.1. iar a managerului de grupuri în anexa A.

ire)

B. Varianta în care nu se folosește transmisie broadcast

Protocolul de bază folosit în comunicație va fi protocolul TCP; s-a preferat acest protocol deoarece rezolvă automat problemele de retransmisie în cazul unor erori pe liniile de comunicație, oferind siguranță în funcționare. Se poate folosi însă și protocolul UDP, dacă se completează cu un mecanism de confirmări.

Structura pachetelor vehiculate în rețea este aceeași cu cea prezentată la punctul A. Spre deosebire de varianta A, fiecare nod va trebui să mențină structuri de date referitoare la disponerea grupurilor pe noduri în sistem.

Atomicitatea transmisiilor se rezolvă printr-un mecanism asemănător protocolului *two-phase commit* (vezi capitolul 7). Acest protocol presupune că orice mesaj se transmite în rețea în două faze:

- (1) În prima fază, *faza de pregătire*, se *încearcă* dacă transmisia poate reuși, transmițându-se mesajul la toate nodurile receptoare: acestea nu transmit mesaje la clienți, ci confirmă doar faptul că pot primi mesajul;
- (2) În cea de a doua fază, *faza de realizare*, nodul emițător colectează toate mesajele de confirmare și, în momentul în care s-au primit toate confirmările, transmite mesajul *commit* către toate nodurile receptoare, care, numai în acest moment, iau în considerare mesajul primit. Dacă nodul emițător nu a primit toate mesajele de confirmare, după câteva încercări nereușite trimite mesajul *abort* pentru ca mesajul trimis anterior să fie descărcat. În acest caz, deoarece nodul care nu se mai poate contacta poate bloca funcționarea sistemului în ansamblu, se decide eliminarea sa din sistem: toate grupurile care existau numai în acel nod sunt eliminate. Nodul respectiv nu va mai putea adresa cereri de operații sistemului decât după o cerere de conectare explicită. De remarcat că, prin utilizarea protocolului *two-phase commit* se asigură și cerințele prezentate la 3.5.2. referitoare la comunicațiile de grup.

5.2.3.2. Varianta în inel

O altă modalitate de a realiza transmisia broadcast este prin simularea ei. Simularea se poate realiza dacă se organizează serverele locale din fiecare nod, care recepționează cererile din rețea într-un inel logic în care se deplasează un token circulant. Evident că această organizare nu impune nici o restricție asupra topologiei rețelei locale utilizate, care poate fi de tip Ethernet.

Numai procesul care deține token-ul, la un moment dat, poate transmite un mesaj broadcast.

Soluția token-ului circulant asigură ordonarea mesajelor și un control al fluxului de mesaje în sistem. Arhitectura generală a sistemului, inclusiv structura token-ului circulant este prezentată în fig. 5.2.3.2. Câmpul *time_stamp* identifică cel mai mic număr de secvență neasociat unui mesaj din token; cînd un posesor transmite broadcast un mesaj, asociază mesajului acest număr de secvență și apoi incrementează valoarea acestui câmp. Descrierea serverului de comunicație este prezentată la 5.2.4.2. iar a managerului de comunicații în anexa A.

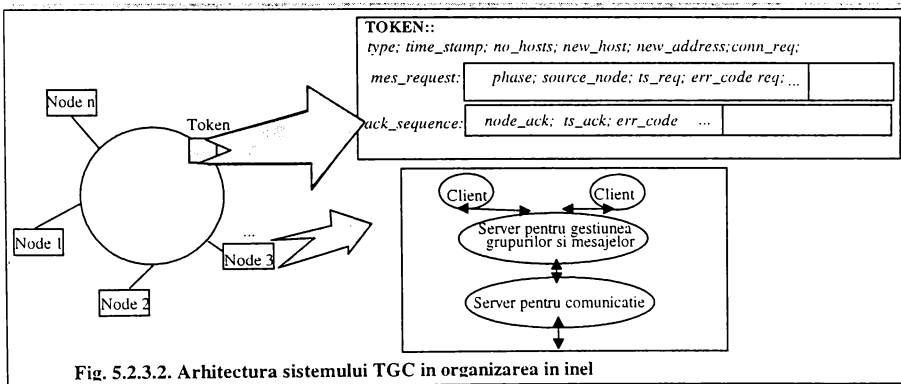


Fig. 5.2.3.2. Arhitectura sistemului TGC în organizarea în inel

5.2.4. Implementarea variantelor complet distribuite

Se va prezenta în acest subcapitol implementarea a 2 variante din cele 3 variante complet distribuite prezentate mai sus: varianta cu transmisie broadcast și varianta cu transmisie în inel.

5.2.4.1. Implementarea variantei cu transmisie broadcast

Se va prezenta implementarea serverului de comunicație din arhitectura prezentată în fig. 5.2.3.1. Managerul pentru gestiunea grupurilor și mesajelor (pe scurt GM-Group Manager) este același pentru ambele variante distribuite (varianta cu transmisie broadcast și varianta cu transmisie în inel) și este prezentat în anexa

A. Pentru comunicațiile din sistem, server de comunicație <-> rețea, server de comunicație <-> manager de grupuri și manager de grupuri <-> clienți se folosesc canalele reprezentate în figura 5.2.4.1.

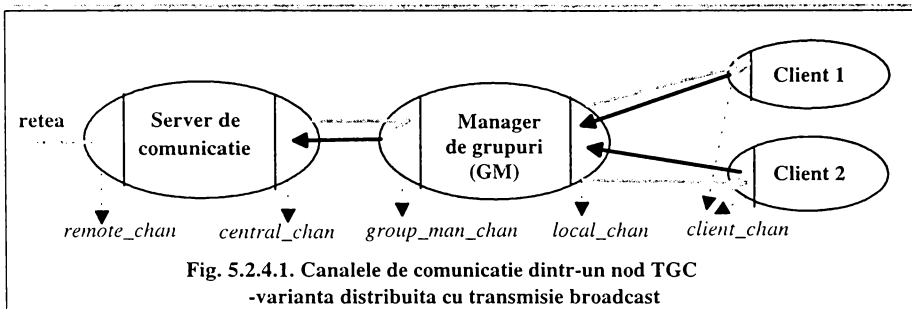


Fig. 5.2.4.1. Canalele de comunicație dintr-un nod TGC
-varianta distribuita cu transmisie broadcast

Serverul de comunicație poate fi descris astfel:

```
# Serverul poate detecta faptul ca un nod nu mai raspunde la mesaje.
# Daca detecteaza un nod care nu mai raspunde la mesaje, serverul de comunicatie il deconecteaza din sistem
# Structuri de date folosite de server in plus fata de cele descrise la 5.2.2.
```

```
const TIME_OUT_CONNECT = ...      # timpul maxim de asteptare la conectarea la sistem
const MAX_MESSAGES = ...         # numarul maxim de mesaje ce pot fi primite de toti clientii
const WAIT_FOR_CONN=...         # timpul dupa care se repeta o cerere de conectare
const MAX_DUP=...               # numarul maxim de retransmisii ale unui mesaj
const TIME_OUT_REQUEST=...      # timpul dupa care se repeta transmisia unei cereri

type kind = enum(CONN_REQ, CONN_REQ_ACK)
type error = enum ( ETOO_MANY_HOSTS , EBUSY, ELOST_A_HOST, ENO_BUFS, EMORE_CLIENTS)
type mes_state= enum (free, in_transit, busy)
type mes_to_group_manager = enum ( ACK, REQ)
type MES_TO_CLIENT = record of
    ( pid: int,                    # buffer de mesaje pentru un client
      mes_no[1:MAX_MES_FOR_A_CLIENT]: int # indica numarul intrarii in buferul de mesaje
      pmes[MAX_MES_FOR_A_CLIENT]: pointer to MES_TO_GROUPS # indica mesajul
      pmes_state[MAX_MES_FOR_A_CLIENT]: mes_state) # indica daca slotul este ocupat

type MESSAGE = record of
    ( knd: kind,                  # tipul unui mesaj primit in local_chan, central_chan
      client_id: int              # si transmis in rețea
      dest_group_id: int          # tipul cererii
      source_group_id: int        # identificatorul clientului de la care s-a trimis cererea
      group_name[1:MAX_LEN_GROUP]: char # identificatorul grupului de la care s-a trimis cererea
      len: int,                  # identificatorul grupului sursa care a trimis mesajul sau cererea
      mes[1:len]: char)          # nume de grup
                                  # lungimea partii variabile a mesajului
                                  # mesajul

type HOSTS: record of
    ( host_name[1:MAX_HOST_NAME]: char # tabela de host-uri
      address: int;
      id: int;
    )

type MES_BUF = record of
    ( state: mes_state            # structura unei intrari in bufferul de mesaje
      message: MES_TO_GROUP)     # starea unei intrari: libera, ocupata, in tranzit
                                  # cuprinde mesajul propriu zis de transmis

chan conn_chan = record of
    ( knd: kind                  # canalul pe care se transmite o cerere de conectare
      host_name[1:MAX_HOST_NAME] # tipul cererii CONN_REQ
      address:int                 # numele host-ului
                                  # adresa hostului

chan inout_chan [1:MAX_HOSTS] =
    ( knd: kind,                # canalul pe care se transmite raspunsul la o cerere de conectare
      # cererea de conectare
```

Capitolul 5: Proiectarea și implementarea unor mecanisme suport pentru dezvoltarea de aplicații distribuite

```
err_code:int          # codul de eroare
no_hosts: int         # numarul total de noduri
hosts_tab:HOSTS,     # tabela de noduri
no_groups: int,      # numarul de grupuri
group_aliases_tab[1:MAX_GROUPS]:GROUP_ALIASES)# tabela de grupuri
chan group_man_chan[1:MAX_HOSTS]= # canal pe care se transmit cererile din retea sau confirmarii
(type: int,          # tipul mesajului, cerere proprie,ACK sau din retea REQ
err_code: int,      # codul de eroare obtinut in nod in urma tratarii cererii
msg: MESSAGE)
chan remote_chan[1:MAX_HOSTS] # canal broadcast pe care se transmit cererile si confirmările in retea
(type: mes_to_group_manager # tipul mesajului, cerere proprie, ACK sau REQ, din retea
case type           # structura mesajului transmis in retea care contine o cerere
REQ:
sender: int,        # identificatorul emitatorului
ts:int,            # marca de timp a mesajului
dup: int,          # mesaj duplicat, numarul de retransmisie
mes: MESSAGE)     # mesajul propriu zis primit de la GM
ACK:               # structura mesajului transmis in retea care contine o confirmare
sender: int,        # identificatorul emitatorului
ts_req: int,       # marca de timp a mesajului confirmat
ts_ack: int,       # marca de timp a confirmării
err_code:int)     # codul de eroare
```

informatii partajate cu managerul de grupuri GM intr-o zona de memorie partajata

```
var no_groups: int          # numarul de grupuri
var group_aliases_tab[1:MAX_GROUPS]:GROUP_ALIASES # mentine denumirile grupurilor
var group_of_clients_tab[1:MAX_GROUPS]:GROUP_OF_CLIENTS # mentine inf. despre componenta gr-lor
var mes_to_clients_tab: [1: MAX_CLI]:MES_TO_CLIENT # mentine inf. despre mes. in asteptare pt.clienti
var messages_buf[1: MAX_MESSAGES]: MES_BUF # bufferul de mesaje
```

CommunicationServer()::

informatii locale

```
var hosts: int          # numarul de noduri conectate la sistem si de grupuri
var hosts_tab[1:MAX_HOSTS]:HOSTS: # tabela cu informatii despre noduri
var mes_queue: queue of (knd, int, int, int, MESSAGE) # coada de cereri: tipul, sender, ts, dup, mes
var ack_list: queue of (knd, int, int, int, int) # coada de confirmari: tipul cererii, sender, ts, err_code
var next_group_id, next_id: int # primul identificator liber: de grup, respectiv de nod
var my_host_name[MAX_HOST_NAME], host_name[MAX_HOST_NAME]:char # retine numele unui nod
var my_adresse, address: int # adresa proprie si adresa altui nod primita din retea
var lc:int:=0, dup: int :=0 # marca de timp locala si nr de retransmitere al unui msaj
var err_code, err_to: int # coduri de eroare
var mes, old_mes: MESSAGE # cerere primita de la managerul de grupuri sau din retea
var knd: kind, ts, ts_req, ts_ack, sender, my_id # folosite la receptia unui mesaj din retea
var ready: bool := true # conditioneaza emisia unui mesaj in retea
var excl: semaphore := 1 # pentru excludere reciproca in accesul la informatiile din mem. partajata
var start_GM: semaphore := 0 # pentru sincronizare server de comunicatii - manager de grupuri
```

initializari: my_host_name, my_adresse

```
broadcast conn_chan( CONN_REQ, my_host_name, my_adresse); # pt. conectare se transmite adresa
instaleaza handler-ul ALARM_function_connect pentru semnalul ALARM;
programeaza un semnal de alarma dupa un interval de TIME_OUT_CONNECT
err_code:=receive inout_chan[my_adresse](knd,err_code, hosts, host_tab, no_groups, group_aliases_tab);
if err_code = TIME_OUT_CONNECT→
    hosts = 1; initializeaza tabela hosts_tab pentru hosts = 1 # doar un singur nod in sistem
fi
fa j=1 to hosts -1 →
    receive inout_chan[my_adresse](knd,err_code, hosts, hosts_tab, no_groups, group_aliases_tab);
    insereaza, (daca nu exista) host_name, in tabela hosts_tab
af
```

```

if err_code = TOO_MANY_HOSTS → exit          # depasirea tabelii de host-uri
fi
instaleaza handler-ul ALARM_function_request pentru semnalul ALARM;
initializeaza hosts_tab, group_aliases_tab, group_of_clients_tab, mes_to_clients_tab
stabileste my_id ca numarul intrarii corespunzatoare nodului propriu in hosts_tab
stabileste next_group_id, next_id, ca urmatoarele intrari libere in hosts_tab, group_aliases_tab
Up(start_GM)          # managerul de grupuri poate sa transmita cereri

do true →          # bucla principala a serverului de comunicatii
  if not empty(conn_chan[my_address]) →
    receive conn_chan[my_address](knd, host_name, address)
    if knd = CONN_REQ →          # receptioneaza o cerere de conectare la sistem
      if hosts < MAX_HOSTS →
        insereaza, (daca nu exista) host_name, adresa in tabela hosts_tab
        actualizeaza hosts, next_id
      Down(excl)
      send inout_chan[address](CONN_REQ_ACK,SUCCESS,hosts,hosts_tab,no_groups,group_aliases_tab)
      Up(excl)
      [] hosts = MAX_HOSTS →
        Down(excl)
        send inout_chan[address](CONN_REQ_ACK,TOO_MANY_HOSTS,hosts,hosts_tab,
                                no_groups, group_aliases_tab)
      Up(excl)
    fi
  fi
fi
if not empty(central_chan[my_address]) →          # s-a primit o cerere de la managerul de grupuri
  receive central_chan[my_address](old_mes)
  if not ready → send group_man_chan[my_address](ACK,EBUSY, old_mes) fi
  dup := 0; ready := false;
  broadcast remote_chan(REQ, my_id, lc, dup, old_mes)
  lc := lc+1;
  programeaza un semnal de alarma care lanseaza handlerul ALARM_function_request dupa un interval
  de TIME_OUT_REQUEST          # pentru a detecta daca s-a confirmat complet
  # un mesaj sau trebuie retransmis
  # s-a primit un mesaj de la din retea
fi
if not empty(remote_chan[my_address]) →
  receive remote_chan[my_address](knd)
  lc := max(lc, ts+1); lc := lc+1;
  if knd = REQ →          # mesajul este o cerere
    receive remote_chan[my_address](sender, ts, dup, mes)
    if dup < MAX_DUP →          # nu s-a depasit numarul maxim de retransmisii
      insereaza in mes_queue (mes.knd, sender, ts, dup, mes)
      err_code := SUCCESS
    if mes.knd = LEAVE_GROUP →          # daca mai exista clienti in grup nu se sterge grupul
      Down(excl)
      if exista cel putin un client local din grupul mes.dest_group_id in group_of_clients_tab →
        err_code := EMORE_CLIENTS
      fi
      Up(excl)
    [] mes.knd = SEND_TO_GROUP →          # daca nu se poate bufera mesajul, eroare
      Down(excl)
    if exista clienti locali pentru mes.dest_group_id in group_of_clients_tab →
      aloca un slot k pentru mesaj in messages_buf[k] pe care il pune in starea in_transit
      if nu se poate aloca un slot k pentru mesaj → err_code := ENO_BUFS
      [] se poate aloca un slot k pentru mesaj →
        if nu se poate aloca un pointer in starea in_transit la message_buf[k] la un client din
        mes.dest_group_id parcurgind tabela mes_to_clients_tab →

```

Capitolul 5: Proiectarea si implementarea unor mecanisme suport pentru dezvoltarea de aplicatii distribuite

```

    err_code:=ENO_BUFS
    fi
    fi
    fi
    Up(excl)
fi
broadcast remote_chan(ACK, my_id, ts, lc, err_code) ; lc:=lc+1 # se transmite confirmarea la cerere
[] dup = MAX_DUP
    ready := true;
    gaseste nodul care nu a confirmat cererea cu marca de timp ts;
    elimina acel nod din hosts_tab
fa mesaje q din mes_queue neduplicate care nu au primit confirmarea de la acel nod→
    if q.sender=my_id → send group_man_chan[my_address](ACK, ELOST_A_HOST, q.mes)
    elimina mesajul q din mes_queue
    fi
af
fi
[] knd = ACK → # s-a primit o confirmare
receive remote_chan[my_address](sender, ts_req, ts_ack, err_code)
insereaza in ack_list: sender, ts_req, ts_ack, err_code
if cererea cu ts_req (coresp. confirmarii) este complet confirmata (s-au primit toate confirmarile)→
    elimina din coada mes_queue cererile incomplet confirmate cu marca de timp mai mica decit ts_req
    extrage din mes_queue in q mesajul complet confirmat
    if q.sender = my_id →
        ready := true;
        sterge programarea semnalului ALARM
    fi
    err_code := SUCCESS
    if q.knd = LEAVE_GROUP or q.knd = SEND_TO_GROUP →
        if exista o confirmare ack la ts_req cu ack.err_code ≠ SUCCESS → err_code:=ack.err_code fi
    fi
    if q.sender = my_id → send group_man_chan[my_address](ACK, err_code, q.mes) fi
    [] q.sender ≠ my_id →
        if err_code= SUCCESS →
            if q.knd = SEND_TO_GROUP
                Down(excl)
                if exista clienti locali pentru grupul q.mes.dest_group_id →
                    send group_man_chan[my_address](REQ, err_code, q.mes)
                    dealoca mesajul din buffer_mes, pointerul din mes_to_clients_tab iar starea lor devine free
                fi
                Up(excl)
                [] err_code ≠ SEND_TO_GROUP →
                    send group_man_chan[my_address](REQ, err_code, q.mes)
            fi
        fi
    fi
fi
fi
fi
od

ALARM_function_connect():: # functia apelata pentru un time_out la conectare
    err_code := TIME_OUT_CONNECT
    return
ALARM_function_request():: # functia apelata pentru un time_out la receptia tuturor confirmarilor
    if not ready →
        if dup ≤ DUP_MAX →
            dup := dup+1;
```



```

    send remote_chan[my_address](my_id, lc, dup, old_mes); # se transmite o cerere duplicat
    lc := lc + 1;
    programeaza un semnal de alarma care lanseaza handlerul ALARM_function_request
    dupa un interval de timp de TIME_OUT_REQUEST
fi
fi
return

```

5.2.4.2. Implementarea variantei cu transmisie în inel

Se va prezenta implementarea serverului de comunicație din arhitectura prezentată în fig. 5.2.3.2. Managerul pentru gestiunea grupurilor și mesajelor (pe scurt GM-Group Manager) este același pentru ambele variante distribuite (varianta cu transmisie broadcast și varianta cu transmisie în inel) și este prezentat în anexa A. Pentru comunicațiile din sistem, server de comunicație <-> rețea, server de comunicație <-> manager de grupuri și manager de grupuri <-> clienți se folosesc canalele reprezentate în figura 5.2.4.2.

Serverul de comunicație poate fi descris astfel:

```

# Serverul de comunicație este specific variantei cu transmisie in inel
# Serverul poate detecta faptul ca un nod nu mai raspunde la mesaje.

```

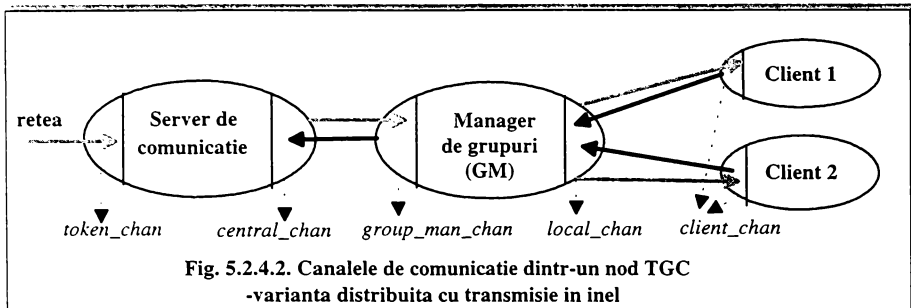


Fig. 5.2.4.2. Canalele de comunicație dintr-un nod TGC -varianta distribuita cu transmisie in inel

Dacă detectează un nod care nu mai răspunde la mesaje, serverul de comunicație îl deconectează din sistem

```

const MAX_MESSAGES = ... # numarul maxim de mesaje ce poate fi primit de toti clientii
const WAIT_FOR_CONN=... # timpul dupa care se repeta o cerere de conectare
const TIME_OUT_CONN= ... # timpul maxim de asteptare la conectarea la sistem
const TIME_OUT_TOKEN=... # timpul maxim cit se asteapta token-ul
const TIME_OUT_TEST =... # timpul maxim cit se asteapta un mesaj de test

type kind = enum(NAME_REQ.. NAME_REQ_ACK, CONN_REQ, CONN_REQ_ACK, NEW_TOKEN, TEST)
type mes_to_group_manager = enum ( ACK, REQ)
type mes_state=(free, in_transit, busy)
type error = enum ( ETOO_MANY_HOST , ETOO_MANY_REQ, ETOO_MANY_REQ_IN_NETW, EBUSY)
type MES_TO_CLIENT = record of # buffer de mesaje pentru un client
    ( pid: int, # identificadorul clientului
      mes_no[1:MAX_MES_FOR_A_CLIENT]: int # indica numarul intrarii in buferul de mesaje
      pmes[MAX_MES_FOR_A_CLIENT]: pointer to MES_TO_GROUPS # indica mesajul
      pmes_state[MAX_MES_FOR_A_CLIENT]: mes_state) # indica daca slotul este ocupat
type MESSAGE = record of
    ( # tipul unui mesaj primit in local_chan, central_chan
      # si transmis in rețea
      knd: kind, # tipul cererii
      client_id: int # identificadorul clientului de la care s-a trimis cererea
      dest_group_id: int # identificadorul grupului de la care s-a trimis cererea
      source_group_id: int # identificadorul grupului sursa care a trimis mesajul sau cererea
      group_name[1:MAX_LEN_GROUP]: char # nume de grup
      len: int, # lungimea partii variabile a mesajului
      mes[1:len]: char)# mesajul

```

Capitolul 5: Proiectarea si implementarea unor mecanisme suport pentru dezvoltarea de aplicatii distribuite

```
type ACK_TYPE = record of
    ( node_ack: int          # structura unei confirmari la o cerere in token
      ts_ack: int           # adresa proprie a nodului care confirma o cerere
      err_code: int)       # marca de timp a cererii confirmate
type MES_BUF = record of
    (state: mes_state      # codul de eroare obtinut la confirmarea unei cereri
      message: MES_TO_GROUP) # structura unei intrari in bufferul de mesaje
                                # starea unei intrari: libera, ocupata, in tranzit
                                # mesajul propriu zis transmis
type TOKEN = record of
    (type: int             # tipul mesajului TOKEN sau TEST
      time_stamp: int      # marca de timp a mesajelor
      no_hosts: int        # numarul de noduri
      max_req: int         # numarul de cereri din token
      max_ack: int         # numarul de confirmari din token
      new_host[1:MAX_HOST_NAME]: char # numele nodului nou conectat
      new_address: int      # adresa nodului nou conectat
      conn_req: int        # folosit la conectarea unui nod
      mes_request[1:MAX_HOSTS*2]: record of
          phase: int;
          source_node: int
          ts_req: int,
          err_code: int,
          req: MESSAGE)     # coada de cereri
      ack_sequence: [1:MAX_HOSTS*MAX_HOSTS*2]: ACK_TYPE # coada de confirmari
    )
chan token_chan[1..MAX_HOSTS](token: TOKEN) # canalul pe care se transmit mesaje in inelul logic
chan inout_chan [1..MAX_HOSTS] =           # canalul pe care se transmite o cerere de conectare si validarea la
    (knd: kind                             # cererea de conectare
     NAME_REQ_ACK, CONN_REQ.:
       no_hosts: int                       # numarul nodurilor
       host_name[1:MAX_HOST_NAME]         # numele host-ului
       address:int                          # adresa hostului
     CONN_REQ_ACK:
       err_code: int                       # codul de eroare
       token:TOKEN                         # token-ul primit la conectare
       no_groups: int,
       group_aliases_tab[1:MAX_GROUPS]:GROUP_ALIASES )
)
chan name_chan [1..MAX_HOSTS] =           # canalul pe care se transmite o cerere de deconectare a unui nod
    (                                     # din sistem sau un mesaj care anunta un nod care nu raspunde
     knd: kind                             # tipul cererii
     host_name[1:MAX_HOST_NAME]           # numele host-ului
     address:int)                          # adresa hostului
chan group_man_chan: [1..MAX_HOSTS] =     # canal pe care se transmit cereri din retea sau confirmari la GM
    (type: mes_to_group_manager           # tipul mesajului, cerere proprie, ACK sau REQ, din retea
     err_code: int                         # codul de eroare
     mesg: MESSAGE)
chan central_chan [1:MAX_HOSTS] =         # canal pe care se primesc cererile de la managerul de grupuri(GM)
    (mes:MESSAGE)

# informatii partajate cu managerul de grupuri intr-o zona de memorie partajata
var no_groups: int                        # numarul de grupuri
var group_aliases_tab[1:MAX_GROUPS]:GROUP_ALIASES # mentine denumirile grupurilor
var group_of_clients_tab[1:MAX_GROUPS]:GROUP_OF_CLIENTS # mentine inf. despre componenta gr-lor
var mes_to_clients_tab: [1: MAX_CLI]:MES_TO_CLIENT # mentine inf. despre mes. in asteptare pt.clienti
var messages_buf[1: MAX_MESSAGES]: MES_BUF # bufferul de mesaje
```

CommunicationServer()::

```
# informatii locale
var hosts: int # numarul de noduri conectate la sistem
var pending_req[1:MAX_REQUESTS]: MESSAGE # bufferul de cereri memorate pina la primirea token-ului
var hosts_tab[1:MAX_HOSTS]: record of ( # tabela de host-uri
    host_name[1:MAX_HOST_NAME]: char
    address: int)
var token:TOKEN # variabila care pastreaza token-ul
var my_host_name[MAX_HOST_NAME], host_name[MAX_HOST_NAME]:char # retine numele unui nod
var my_adresse, next_address, pred_address, address: int # adresa proprie si a nodului vecin
var new_conn: bool = false # indica faptul ca s-a primit o cerere de conectare
var knd: kind, err_code: int # tipul unei cereri si codul de eroare
var mes: MESSAGE # cerere primita de la GM
var excl: semaphore := 1 # pentru excludere reciproca in accesul la informatiile din mem. partajata
var start_GM: semaphore := 0 # pentru sincronizare server de comunicatii - manager de grupuri
```

```
initializari: my_host_name, my_address
again:broadcast name_chan( NAME_REQ, my_host_name, my_address); # cerere de conectare la sistem
instaleaza handler-ul ALARM_function_connect pentru semnalul ALARM;
programeaza un semnal de alarma dupa un interval de TIME_OUT_CONNECT
err_code := receive inout_chan[my_adresse](knd, hosts, host_name, address);
if err_code = TIME_OUT_CONNECT → # numai un singur nod in sistem
    hosts = 1;
    initializeaza tabela hosts_tab pentru hosts = 1
fi
fa j=1 to hosts -1 →
    receive inout_chan[my_address](knd)
    if knd = NAME_REQ_ACK → receive inout_chan[my_address]( hosts, host_name, address); fi
    insereaza, (daca nu exista) host_name, address in tabela hosts_tab (tabela sa fie ordonata dupa adrese IP)
fa
alege numele predecesorului in inel;
stabileste nodul succ in inel next_address si cel pred pred_address astfel incit tabela sa fie ordonata
send inout_chan[pred_adresse](CON_REQ,0, my_host_name, my_address);
receive inout_chan[my_address] (knd, err_code);
if err_code = EBUSY→
    sleep (WAIT_FOR_CONN); goto again
[] err_code = ETOO_MANY_HOST →
    exit;
[] receive inout_chan[my_adresse] (token, no_groups,group_aliases_tab);
fi
instaleaza handler-ul ALARM_function_token pentru semnalul ALARM;
initializeaza group_of_clients_tab, mes_to_clients_tab
initializeaza bufferul de mesaje messages_buf, bufferul de cereri pending_req
token.time_stamp := token.time_stamp+1
insereaza pentru toate cererile din mes_request din token o confirmare cu
    source_node:=my_address, ts_ack:=mes_request[j].ts, err_code:=SUCCESS
send token_chan [next_address](token)
programeaza un semnal de alarma (pentru cazul cind nu se receptioneaza token-ul)
    dupa un interval de TIME_OUT_TOKEN
Up(start_GM) # managerul de grupuri poate sa transmita cereri

do true → # bucla principala a serverului de comunicatii
    if not empty(name_chan[my_address]) →
        receive name_chan[my_address](knd,host_name, address)
        if knd = NAME_REQ → # receptioneaza o cerere de livrare a numelui si a adresei nodului
            send inout_chan[address](NAME_REQ_ACK,hosts, my_host_name, my_address)
        [] knd = NEW_TOKEN → # genereaza un nou token la pierderea token-ului
            actualizeaza hosts_tab, hosts:=hosts-1
            if my_host_name precede in inel pe host_name →
```

Capitolul 5: Proiectarea si implementarea unor mecanisme suport pentru dezvoltarea de aplicatii distribuite

```

    token.hosts:= token.hosts - 1
    elimina din token.ack_sequence toate confirmările date de host_name
    actualizeaza next_address la nodul urmator nodului host_name in hosts_tab
    send token_chan[next_address](token);
    programeaza un semnal de alarma (pentru cazul cind nu se receptioneaza token-ul)
    dupa un interval de TIME_OUT_TOKEN
fi
fi
fi
if not empty(inout_chan[my_address]) →
    receive inout_chan[my_address](knd) # receptioneaza o cerere de conectare in sistem
    if knd=CONN_REQ →
        receive inout_chan[my_address](temp,host_name, address); new_conn = true;
    fi
    if hosts = 1 →
        actualizeaza next_address, deci nodul succesori in inel ; hosts = hosts + 1
        insereaza, (daca nu exista) host_name, address in tabela hosts_tab astfel incit tabela sa fie ordonata
        in functie de adresele IP

        token.new_host:=NULL; token.new_address:=NULL;
        token.no_hosts:= hosts; token.time_stamp:= 1; token.conn_req:=0
        Down(excl)
        send inout_chan[next_address](CONN_REQ_ACK,SUCCESS,token,no_groups,group_aliases_tab)
        Up(excl)
    fi
fi
if not empty (token_chan[my_address]) # receptioneaza token-ul
    receive token_chan[my_address] (token); hosts:= token.no_hosts;
    fa mesaje j din token.mes_request in ordinea marilor de timp →
    if token.mes_request[j].phase = 2 →
        err_code :=token.mes_request[j].err_code
        if token.mes_request[j].source_node = myAdresse → # cerere locala
            send group_man_chan[my_address](ACK, err_code, token.mes_request[j].req)
            Down(excl)
            if err_code=SUCCESS and token.mes_request[j].req.knd= SEND_TO_GROUP
                and exista cel putin un client in group_of_clients_tab din grupul
                token.mes_request[j].req.dest_group_id→
                    dealoca intrarile marcate in_transit din mes_to_clients_tab si messages_buf
            fi
            Up(excl)
        [] token.mes_request[j].source_node ≠ myAdresse → # cerere din retea
        if err_code=SUCCESS →
            if token.mes_request[j].req.knd= SEND_TO_GROUP →
                Down(excl)
                if exista cel putin un client in group_of_clients_tab din grupul
                    token.mes_request[j].req.dest_group_id
                    send group_man_chan[my_address](REQ, err_code, token.mes_request[j].req)
                    dealoca intrarile marcate in_transit din mes_to_clients_tab si messages_buf
                fi
                Up(excl)
            [] token.mes_request[j].req.knd ≠ SEND_TO_GROUP →
                send group_man_chan[my_address](REQ, err_code, token.mes_request[j].req)
            fi
        fi
    fi
fi
[] token.mes_request[j].phase = 1 →
    aloca un slot l in token pentru confirmarea cererii token.mes_request[j]
    token.ack_sequence[j][l].node_ack := my_address;
    token.ack_sequence[j][l].ts_ack := token.mes_request[j].ts
```

```

if token.mes_request[j].req.knd = SEND_TO_GROUP →
  Down(excl)
  if exista clienti locali pentru token.mes_request[j].req.dest_group_id→
    aloca un slot k pentru mesaj in message_buf[k] pe care il pune in starea in_transit
  if nu se poate aloca un slot k pentru mesaj →
    token.ack_sequence[j][l].err_code := ENO_BUFS
    token.mes_request[j].err_code := ENO_BUFS
  [] se poate aloca un slot k pentru mesaj →
    if nu se poate aloca un pointer in starea in_transit la message_buf[k] la un client din
      token.mes_request[j].dest_group_id parcurgind tabela mes_to_clients_tab→
        token.ack_sequence[j][l].err_code := ENO_BUFS
        token.mes_request[j].err_code := ENO_BUFS
    fi
  fi
  Up(excl)
  [] token.mes_request[j].req.knd = LEAVE_GROUP →
  Down(excl)
  if exista cel putin un client in group_of_clients_tab din grupul
    token.mes_request[j].req.dest_group_id→
    token.ack_sequence[j][l].err_code := EMORE_CLIENTS
    token.mes_request[j].err_code := EMORE_CLIENTS
  fi
  Up(excl)
fi

fi
if token.mes_request[j].source_node = myAdresse → # cerere lansata de nodul respectiv
  token.mes_request[j].phase := token.mes_request[j].phase+1
  if token.mes_request[j].phase = 3 →
    dealoca mes_request[j] din token si ordoneaza mesajele in coada
    dealoca confirmarile pentru mes_request[j] din coada de confirmari din token
  fi
fi
af
if exista loc k in token.mes_request pentru urmatoarea cerere j din pending_req[j]→
  copiaza cererea pending_req[j] in token.mes_request; dealoca intrarea din pending_req[j]
  token.mes_request[k].err_code := SUCCESS; token.mes_request[k].phase := 1;
  token.mes_request[k].ts := token.time_stamp
  [] nu exista loc in token.mes_request pentru o cerere →
  send group_man_chan[my_address](ACK,ETOO_MANY_REQ_IN_NETW, pending_req[j])
  dealoca intrarea din pending_req[j]
fi
if (new_conn = true) →
  if (token.conn_req = 0)→
  if hosts = MAX_HOSTS →
    send inout_chan[address](CON_REQ_ACK, TOO_MANY_HOSTS);
    send token[next_address](token);
  [] hosts < MAX_HOSTS→
    actualizeaza next_address, deci nodul succesori in inel
    insereaza, (daca nu exista) host_name, address in tabela hosts_tab (tabela ordonata in fct de adrese IP)
    token.new_host := host_name; token.new_address := address;
    token.no_hosts := token.no_hosts + 1;
    token.time_stamp := token.time_stamp+1;
    token.conn_req := token.conn_req+1
    send inout_chan[next_address](CONN_REQ_ACK,SUCCESS,token,no_groups,group_aliases_tab)
  fi
  [] token.conn_req < token.no_hosts→
  send inout_chan[address](CONN_REQ_ACK, EBUSY)

```

Capitolul 5: Proiectarea și implementarea unor mecanisme suport pentru dezvoltarea de aplicații distribuite

```
insereaza, (daca nu exista) token.new_host, token.new_adress in tabela hosts_tab
hosts := token.no_hosts; new_conn:= false; token.conn_req:=token.conn_req+1
token.time_stamp:=token.time_stamp+1;
send token[next_address](token)
[] token.conn_req = token.no_hosts→
new_conn:=false; token.conn_req := 0
token.time_stamp:=token.time_stamp+1;
send token[next_address](token)
fi
[] (new_conn = false) →
token.time_stamp:=token.time_stamp+1;
if token.conn_req > 0→
insereaza, (daca nu exista) host_name, address in tabela hosts_tab;
token.conn_req:=token.conn_req+1
fi
send token[next_address](token);
fi
fi
if not empty (central_chan[my_address])→ # cerere de la managerul de grupuri
receive central_chan[my_address](mes);
if hosts = 1 → send group_man_chan[my_address](ACK, SUCCESS, mes)
[] if exista loc in pending_request → # cererea se bufereaza pentru a fi transmisa la primirea token-ului
aloca urmatorul slot si memoreaza mes
[] nu exista loc in pending_request →
send group_man_chan[my_address](ACK, ETOO_MANY_REQ, mes)
# daca nu exista loc cerere este refuzata
fi
fi
od

ALARM_function_connect():: # functia apelata pentru un time_out la conectare
err_code := TIME_OUT_CONNECT
return

ALARM_function_token():: # functia apelata pentru un time_out la primirea token-ului
token.type := TEST
send token[next_address](token)
programeaza un semnal de alarma care lanseaza handlerul ALARM_function_disconnect
dupa un interval de timp de TIME_OUT_REQUEST
receive token[my_address](token)
return

ALARM_function_disconnect():: # functia apelata pentru un time_out la primirea token-ului de test
broadcast name_chan(NEW_TOKEN, nume nod care precede pe my_host_name,
adresa nod care precede pe my_address)
return
```

5.2.5. Comparație între variantele propuse și concluzii

În urma analizei celor patru variante propuse, varianta 1- varianta semidistribuită, varianta a 2 a - varianta distribuită cu mesaje broadcast, varianta a 3 a - varianta distribuită fără utilizare broadcast, varianta a 4 a - varianta în inel, se pot enunța următoarele concluzii:

(1) În toate soluțiile se consideră că se cunoaște numărul participanților în sistem pentru a se putea lua decizia dacă un mesaj a fost recepționat sau nu de *toți* membrii unui grup. În variantele 1 și 3 trebuie cunoscută și identitatea participanților (nodurilor). În varianta a 4 a trebuie cunoscută identitatea vecinilor din inel.

(2) Conectarea unui nod nou la sistem se poate face numai după confirmarea cererii de conectarea a acestuia de către toate nodurile în variantele 2, 3, după confirmarea serverului central în varianta 1 și după confirmarea de către nodurile vecine în inel în varianta a 4 a. În variantele 2, 3 și 4, la conectare, nodul

conectat trebuie să-și construiască și tabela care menține legătura *group_id* - *alias*-uri, eventual și disponerea grupurilor pe noduri (varianta a 3 a).

(3) Implementarea Unix și Windows a descrierii variantei 1 prezentate la 5.2.2.3. este directă: pentru serverele locale se folosește o implementare cu *select*, pentru serverul central master se folosește în Unix *fork+exec* și o zonă de memorie partajată între serverul central master și toate serverele centrale fii iar în Windows serverele centrale se implementează cu *fire*.

(4) Atomicitatea se obține folosind pentru comunicația în rețea protocolul TCP și fie folosind direct un protocol *two-phase-commit* (varianta a 3 a), fie un protocol echivalent (în varianta a 4 a se simulează *two-phase-commit* în inelul logic al procesoarelor). În varianta a 2 a se folosește protocolul UDP pentru a putea beneficia de transmisia broadcast iar atomicitatea se obține datorită faptului că fiecare nod "cunoaște" toate confirmările sosite pentru un mesaj.

(5) În cazul în care un nod este deconectat de la sistem fără a fi emis o cerere de deconectare, și deci refuză să mai primească mesaje, în toate variantele se forțează automat, pentru nodul respectiv, o cerere de deconectare de la sistem. În varianta a 4 a trebuie reconstruit inelul logic al procesoarelor.

(6) În ceea ce privește o comparație a sistemului TGC proiectat, cu alte sisteme de comunicații de grup se pot face următoarele observații: sistemul Totem [MMA96] folosește soluția organizării procesoarelor în inel, sistemele MCL și S:Net Linda Kernel folosesc transmisie broadcast, iar sistemul Vartlaap [LRM93] o soluție semidistribuită. Recent, unele sisteme de comunicații de grup au fost suplimentate cu facilități pentru toleranța la defecte. Sistemul TGC, deși suportă deconectarea din sistem a unui nod nu este prevăzut cu o completă toleranță la defecte, în special în prima variantă în care serverul central este unic. Toleranța la defecte în prima variantă impune menținerea unei copii a serverului central (mecanismul de "mirroring" al serverului central).

(7) În toate variantele, identificatorii de grup, globali în întreg sistemul, sunt refolosiți după distrugerea unui grup, la părăsirea acestuia de către ultimul membru. Se apreciază că, la o extindere ulterioară a sistemului TGC se poate lua în considerare observația ca identificatorii de grup accordați noilor grupuri să aibă valori nefolosite (asemănător *pid*-urilor din Unix).

(8) În toate variantele alegerea unui identificator pentru un grup trebuie făcută în excludere reciprocă. Prin urmare, execuția operației *Create_Group* trebuie făcută în regim de excludere reciprocă folosind unul din algoritmi prezentați și implementați în subcapitolele 4.3.2., 4.3.3., 4.3.4. În prima variantă, această problemă se pune doar local deoarece serverele centrale fii aleg identificatorii de grup; accesul la identificatori se face controlat de un semafor local.

(9) În ceea ce privește primitivile sistemului TGC acestea pot fi extinse astfel: (a) primitiva *SendToGroup()* poate fi extinsă astfel încât să admită ca destinație și grupul *ALL* (în sensul generic de toate grupurile). (b) poate fi introdusă o nouă primitivă care să găsească toate procesele componente ale unui grup, la un moment dat. (c) un grup poate să rămână în folosință chiar dacă nici un proces nu mai este conectat la el; aceasta implică introducerea unei noi primitive pentru distrugerea explicită a unui grup.

(10) Succesiunea de cereri adresate de un client sistemului TGC, în oricare din cele trei variante pentru conectarea la un grup este următoarea:

Clienți[*i*:1..*n*]:

```
var done: bool. code: int #variabile de lucru
var MyGroupId: int = 0 # identificador de grup
var MyGroupAlias[1: MAX_NAME_LEN]: char = ... # numele grupului
var MyName[1: MAX_NAME_LEN]: char = ... #numele procesului
```

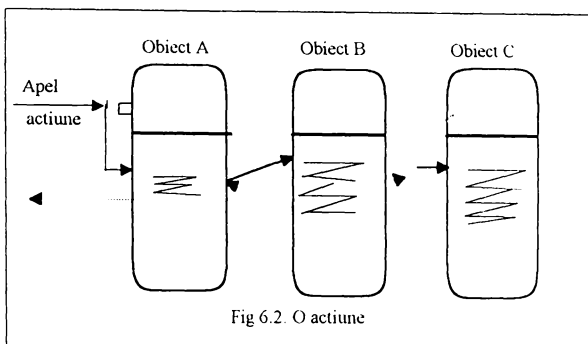
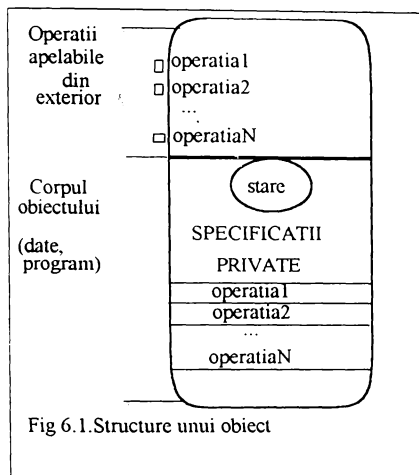
```
done := false;
IniTGC(MyName, NULL, NULL)
do not done →
  code := GetIdGroup(MyGroupAlias, MyGroupId)
  if code = SUCCESS →
    code := JoinToGroup(MyGroupId);
    if code = SUCCESS → done := true fi
  [] code = SUCCESS →
    code := CreateGroup(MyGroupAlias, MyGroupId)
    if code = SUCCESS → done := true fi
  fi
od
```


Capitolul 6.

SISTEME DE PROGRAMARE DISTRIBUITĂ BAZATE PE OBIECTE (SPDBO). CRITERII DE PROIECTARE

6.1. Concepte fundamentale

Un *obiect* este definit ca o entitate care încapsulează informații de stare (sau date) și un set de operații sau proceduri care manipulează aceste date într-un mod controlat, astfel încât pot fi tratate, colectiv, ca o singură unitate (fig 6.1).



Singura modalitate în care poate fi modificată sau examinată starea unui obiect este prin apelul operațiilor facute publice, deci accesibile din exterior. Aceasta crează o interfață exterioară pentru obiect, prin care se specifică operațiile accesibile asupra obiectului pastrand însă implementarea lor și informațiile de stare inaccesibile în exterior.

O operație asupra unui obiect poate apela alte operații, posibil chiar

asupra altor obiecte; aceste operații la rândul lor pot apela operații asupra altor obiecte s.a.m.d. O înlanțuire a unor astfel de apeluri de operații se numește *acțiune* (fig 6.2).

Limbajele de programare pot încuraja lucrul cu obiecte, permițând dezvoltarea unor programe constând din obiecte autonome care interacționează: astfel de limbaje pot fi clasificate [CC91] în *limbaje bazate pe tipul obiect*, pentru care tipul obiect este o caracteristică a limbajului: Ada, Modula-2 și *limbaje orientate pe obiect* (care implementează și conceptul de moștenire): C++, SR și Smalltalk. Conform principiilor urmărite de aceste limbaje, un program trebuie să aibă următoarele cinci caracteristici [CC91]:

1. *să fie abstract*. Concepțiile majore de proiectare trebuie să fie separate de detaliile de implementare.
2. *să fie structurat*. Un program complex trebuie să poată fi descompus în mai multe componente de dimensiuni adecvate și între care să poată fi stabilite relații bine definite.
3. *să fie modular*. Structura internă a fiecărei componente trebuie să fie astfel concepută încât să nu depindă de structura altei componente.

4. *să fie concis*. Codul programului trebuie să fie cât mai clar.

5. *să fie verificabil*. Programul trebuie să fie ușor de testat și depanat.

Un *sistem de programare bazat pe obiecte (SPBO)* poate fi definit ca un mediu de calcul furnizat de un limbaj de programare bazat pe obiecte și de un sistem de operare sau o extensie a sa care implementează conceptul de obiect. Aceasta permite obiectelor să fie manipulate, gestionate și utilizate eficient, dar și partajarea obiectelor de către mai mulți utilizatori: ultima caracteristică nu este asigurată de limbajele de programare. Sistemul de operare dintr-un SPBO furnizează facilități în ceea ce privește mai ales interacțiunea obiectelor și gestiunea resurselor.

Un *sistem de programare distribuit bazat pe obiecte (SPDBO)* reunește caracteristicile unui sistem de programare bazat pe obiect ca și cele ale unui sistem de calcul distribuit. Aceste caracteristici sunt [CC91]:

1. *Distribuire*. Un SPDBO unește nodurile de lucru independente ale unei rețele locale de calculatoare, posibil chiar eterogene astfel încât să furnizeze un sistem de calcul descentralizat.

2. *Transparență*. Un SPDBO trebuie să ascundă caracteristică de distribuție a mediului de calcul față de utilizatori. De exemplu, localizarea obiectelor în rețea trebuie să fie transparentă pentru utilizatori și pe de altă parte trebuie să furnizeze acces uniform la toate obiectele sistemului, fie că aceste obiecte sunt active în memorie sau inactive în memoria secundară.

3. **Integritatea datelor.** Într-un SPDBO, înaintea execuției oricărei operații asupra sa, un obiect trebuie să fie într-o stare validă. Un obiect este într-o *stare validă* dacă starea sa a rezultat în urma execuției cu succes a unei operații. Dacă o operație nu se execută cu succes, atunci sistemul trebuie să anuleze toate modificările survenite în cadrul acelei operații.

4. **Toleranța la defecte.** Apariția unui defect al unui nod sau al unui obiect trebuie să reprezinte pentru SPDBO numai un defect parțial; restul sistemului trebuie să fie capabil să continue prelucrările, eventual cu inconvenientul lipsei unui serviciu.

5. **Recuperabilitate.** În urma apariției unui defect al unui nod și înlăturării lui, un SPDBO trebuie să fie capabil să recupereze obiectele permanente de pe acel nod.

6. **Autonomia obiectelor.** Un SPDBO poate permite proprietarului unui obiect să specifice clienții cărora li se permite să execute operații asupra acelui obiect.

7. **Concurența la nivelul programelor.** Un SPDBO trebuie să furnizeze programelor posibilitatea ca operațiile asupra diferitelor obiecte să se execute concurrent, în procesoare diferite. (fig 6.3.)

8. **Concurența la nivelul obiectelor.** Asupra unui obiect pot să se execute concurrent mai multe operații, evident cu respectarea restricțiilor de consistență ale obiectului.

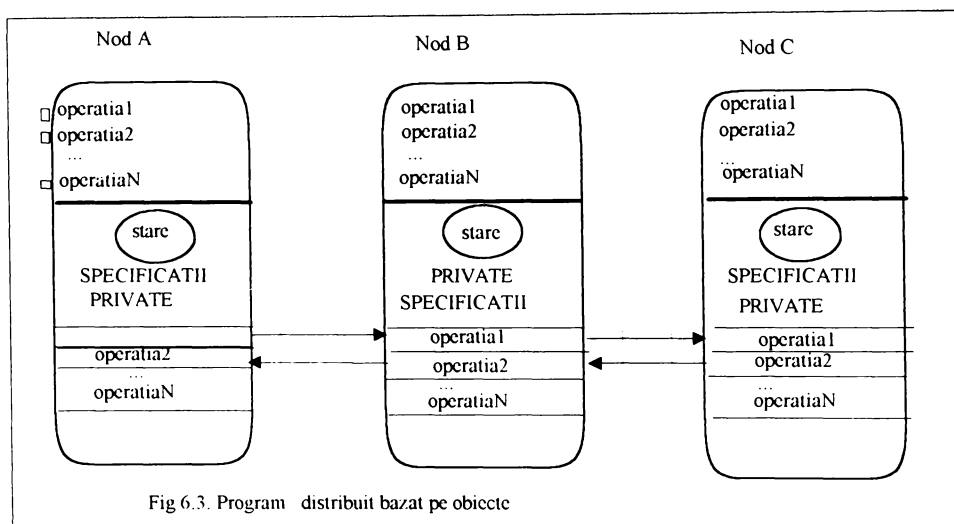


Fig 6.3. Program distribuit bazat pe obiecte

6.2. Structura obiectelor

Dimensiunea obiectelor și cantitatea de prelucrări executate asupra lor caracterizează **granulitatea obiectelor**. În cel mai simplu caz un SPDBO permite numai obiecte complexe, de dimensiune mare, cu operații conținând un număr mare de instrucțiuni și cu relativ puține interacțiuni cu alte obiecte. Exemple de astfel de obiecte sunt un fișier sau o baza de date monoutilizator. Caracteristic acestor obiecte este faptul că sunt rezidente în spațiul propriu de adrese.

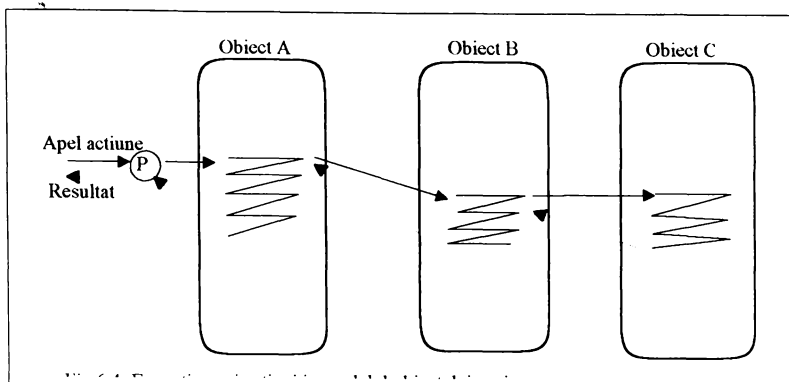
Pentru un control mai eficient al datelor, un SPDBO poate permite și obiecte de complexitate medie, cu exponenți tipici listele înlanțuite și componentele unei baze de date multiutilizator. Aceste obiecte pot fi rezidente în spațiul de adrese al obiectelor complexe. Aceasta permite un grad mai ridicat de concurență, deoarece sincronizarea accesului la date se face la nivelul obiectelor de complexitate medie, dar și reducerea cantității de date copiate în memoria secundară atunci când o acțiune este desavârșită (*commits*). Obiectele de complexitate mică sunt caracterizate prin dimensiune mică, operații compuse dintr-un număr mic de instrucțiuni și număr mare de interacțiuni cu alte obiecte. Exemple specifice sunt tipurile de date standard furnizate de limbajele de programare: numere întregi, valori logice, numere complexe. Un astfel de obiect este rezident în spațiul de adresă al unui obiect de complexitate medie sau mare.

Structura obiectelor mai este decisă și de relațiile dintre procese și obiecte. Procesele pot fi implementate separat și asociate temporar obiectelor pe care le examinează sau modifică sau pot fi asociate permanent obiectelor respective. Aceste două implementări corespund celor două modele cunoscute: *modelul obiectului pasiv* și *modelul obiectului activ*.

În *modelul obiectului pasiv* procesele și obiectele sunt entități complet separate. Un proces nu este restricționat la un singur obiect, dar numai într-un singur proces sunt înglobate toate operațiile dintr-o acțiune.

În consecință, un proces se poate executa asupra mai multor obiecte pe toată perioada existenței sale. Când dintr-un proces se apelează o operație asupra unui obiect, execuția procesului se comută, migrează în obiectul

respectiv. După terminarea execuției operației, procesul este returnat la operația apelantă; un exemplu este prezentat în figura 6.4.:



În *modelul obiectului activ* se crează unul sau mai multe procese servere care se asignează obiectului și care implementează operațiile ce se pot executa asupra obiectului. Fiecare proces este deci restricționat la un obiect particular; atunci când este distrus obiectul, sunt distruse și procesele server asociate lui. Interacțiunea dintre un proces server și un proces client se desfășoară în modelul obiectului activ astfel [CC91] [DG89]:

- (1) Procesul client prezintă cererea de execuție a unei operații împreună cu o listă de argumente;
- (2) Procesul server acceptă cererea, apoi localizează și execută operația specificată;
- (3) Dacă în cursul execuției operației se apelează o operație asupra unui alt obiect, procesul server, devenit acum proces client prezintă cererea și așteaptă rezultatul.
- (4) Când execuția operației se termină, procesul server returnează rezultatul clientului.

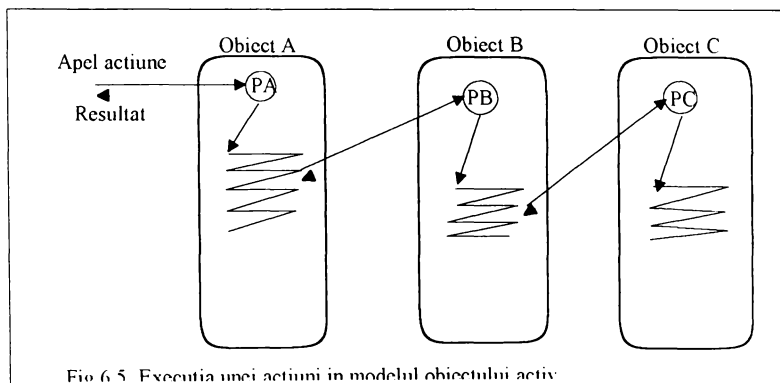


Fig 6.5. Execuția unei acțiuni în modelul obiectului activ

În acest model al obiectului activ, în execuția unei singure acțiuni pot fi implicate mai multe procese (fig.6.5.): În acest model, al obiectului activ, într-o singură acțiune pot fi implicate mai multe procese.

În *varianta statică a modelului activ*, la crearea, fiecărui obiect i se asociază un număr fix de procese server. Orice cerere de execuție a unei operații asupra unui obiect este direcționată către un proces server liber. Dacă toate procesele server sunt ocupate, cererile sunt înlantuite într-o coadă de așteptare, urmînd a fi servite în momentul cînd serverele devin disponibile. Gradul de multiprogramare obținut prin această variantă este, evident, limitat de numărul de procese server create pentru un obiect (vezi 3.1.2.).

În *varianta dinamică a obiectului activ*, se crează câte un proces server pentru fiecare obiect, atunci cînd se recepționează o cerere de execuție a unei operații asupra aceluși obiect. La terminarea execuției, procesul respectiv este distrus. Astfel, în această variantă nu se înlantuite cererile de execuție a operațiilor asupra obiectelor, însă cu dezavantajul unui timp adițional necesar pentru crearea și distrugerea dinamică a proceselor (vezi 3.1.2.).

Modelul obiectului activ introduce posibilitatea apariției interblocărilor. O interblocare poate apare ori de cîte ori nu există procese server suficiente pentru a executa cererile de operații recepționate pentru obiectul respectiv. De exemplu, în varianta statică o singură acțiune poate apela același obiect de mai multe ori decît numărul de procese server asociate obiectului (o astfel de situație poate fi un apel recursiv); în momentul în care nu mai există procese server disponibile apare interblocarea. În varianta dinamică, interblocarea este mai puțin o problemă; ea poate apare totuși dacă se permite ca o acțiune să aibă un număr oarecare de apeluri de operații.

6.3. Gestionarea obiectelor

Acest subcapitol descrie pe scurt caracteristicile unui SPDBO în ceea ce privește permanentizarea efectelor unei acțiuni asupra obiectelor, sincronizarea execuțiilor concurente de apeluri asupra unui obiect, protejarea obiectelor împotriva clienților neautorizați și recuperarea obiectelor în cazuri de defect.

6.3.1. Gestionarea acțiunilor

Acțiunile trebuie să aibă următoarele trei proprietăți ([CC91]):

1. **Seriabilitate** Mai multe acțiuni care se execută concurrent trebuie planificate astfel încât efectul total să fie același ca și cum ele s-ar fi executat serial.
2. **Atomicitate** O acțiune este fie executată fără nici o eroare în decursul vreunei operații, fie nu are nici un efect.
3. **Permanența** Efectul unei acțiuni executată cu succes nu poate fi pierdut, cu excepția cazului erorilor catastrofale.

O acțiune poate genera un număr oarecare de apeluri de operații asupra mai multor obiecte. Ea se consideră executată cu succes numai dacă toate operațiile au fost executate fără erori. O acțiune terminată cu succes este denumită *desăvârșită* ("committed"); o acțiune care nu poate fi terminată cu succes datorită unei erori aparute într-o operație se spune că s-a terminat prin abort. Dacă o acțiune se termină prin *commit*, modificările efectuate asupra obiectelor sunt făcute permanente prin înregistrarea lor în memoria secundară; în caz contrar toate aceste modificări sunt inversate, astfel încât obiectul respectiv rămâne în starea în care era înaintea acțiunii.

Cel mai cunoscut mecanism pentru desăvârșirea unei acțiuni este protocolul *two-phase commit*: [Tan88][Gos91][Cri88][Lam81](vezi 7.3.):

- (1) O cerere "*prepare*" este transmisă tuturor obiectelor afectate;
- (2) Un obiect care a recepționat o cerere "*prepare*" scrie modificările asupra obiectului în memoria secundară și returnează un mesaj de confirmare;
- (3) Dacă toate obiectele implicate în acțiune au returnat mesaje de confirmare, se transmite un mesaj "*commit*" către toate obiectele. În caz contrar, se transmite un mesaj "*abort*" către toate obiectele și dialogul este terminat.
- (4) Un obiect care a primit mesajul "*commit*" permanentizează modificările asupra obiectelor, eliberează structurile de control (eventual lock-uri, de exemplu) și returnează un al doilea mesaj de confirmare. Un obiect care a primit mesajul "*abort*" inversează modificările și eliberează structurile de control.
- (5) O cerere *commit* neconfirmată este repetată pînă la primirea confirmării. Cînd toate confirmările au fost primite, protocolul *two-phase commit* se consideră încheiat.

Inițierea protocolului de desăvârșire a unei acțiuni (sau de execuție *commit*) poate fi în responsabilitatea sistemului SPDBO sau a utilizatorilor. Cea mai simplă soluție este așa numitul *mechanismul cererii (request)*, în care utilizatorul este cel care decide cînd trebuie inițiată procedura de "*commit*". Un utilizator nu este obligat să termine fiecare acțiune prin *commit*, acest lucru poate îmbunătăți performanțele sistemului, dar se pierde caracteristica de atomicitate și permanență.

Pentru a garanta toate proprietățile unei acțiuni, sistemul SPDBO trebuie să fie cel responsabil pentru gestiunea acțiunilor și desăvârșirea lor cînd se pot încheia cu succes.

În *mechanismul tranzacției*, cînd toate operațiile unei acțiuni au fost executate cu succes, sistemul este cel care realizează procedura de *commit* pentru a permanentiza toate modificările. Dacă cel puțin o operație se termină prin eroare, tranzacția trebuie terminată prin abort.

O extensie a mecanismului tranzacțiilor este mecanismul *tranzacțiilor incubate*. Conform acestui mecanism, un singur nivel exterior de acțiuni crează și gestionează mai multe subnivele de *subacțiuni*. Eșuarea unei subacțiuni nu conduce obligatoriu la esuarea întregii acțiuni; acțiunea poate alege orice decizie în urma esuării unei subacțiuni. De exemplu poate alege repetarea execuției subacțiunii, poate termina prin abort întreaga acțiune sau poate ignora pur și simplu terminarea cronată a subacțiunii. Permanentizarea modificărilor făcute de subacțiuni se face numai cînd întreaga acțiune execută "*commit*". Dacă acțiunea se termină cu abort toate modificările făcute de subacțiuni sunt inversate. Mecanismul tranzacțiilor incubate are avantajul unui control mai detaliat al unei acțiuni permițînd esuarea subacțiunilor, dar are dezavantajul unei reprezentări mai complicate în memorie a obiectelor pentru a permite inversarea efectelor unei subacțiuni eșuate.

6.3.2. Sincronizarea acțiunilor

O importantă funcție a unui SPDBO este interdicția conflictelor și interferențelor acțiunilor la nivelul aceluiași obiect. O acțiune nu trebuie să poată observa sau modifica starea unui obiect care a fost parțial modificat de o altă acțiune care nu este terminată. În caz contrar, se poate ajunge la situația de *abort în cascadă*: orice acțiune care a observat starea parțial modificată a unui obiect de către o altă acțiune care mai târziu a fost terminată cu abort trebuie la rîndul ei terminată cu abort. Pentru a elimina situațiile de abort în cascadă și, pe de altă parte, pentru a păstra consistența obiectelor este necesară o *schema de sincronizare*.

Schemele de sincronizare cunoscute se clasifică în scheme *pesimiste* și *optimiste*. În cadrul schemelor de sincronizare *pesimiste*, sistemul trebuie să ia măsurile potrivite pentru a preveni conflictele de acces. O acțiune care apelează o operație asupra unui obiect asupra caruia se execută o altă acțiune trebuie suspendată pînă la terminarea

acesteia sau terminată anormal. Cel mai cunoscut mecanism de sincronizare dintre cele înglobate în schemele pesimiste este cele bazate pe lock-uri Read/Write (vezi 7.4.4.); de asemenea mai sunt utilizate semafoare (eventual semafoare distribuite, vezi 4.3.4.1.) și mecanisme bazate pe mărcile de timp (vezi 4.3.1.2.).

În cadrul schemelor de sincronizare *optimiste*, nu se ia nici o măsură la nivelul unui obiect în cadrul caruia se execută concurent mai multe acțiuni. În schimb, la comiterea unei acțiuni se testează condițiile de seriabilitate astfel încât informațiile de stare luate în considerare de acțiune să nu intre în conflict cu modificările realizate de orice altă acțiune care s-a terminat. Dacă informațiile observate de o acțiune nu au fost alterate de o altă acțiune care s-a terminat prin *commit*, atunci acțiunea poate fi *desavîrsită (commit)*; în caz contrar, modificările realizate de acțiunea respectivă, bazate pe informațiile observate anterior vor fi incorecte și în consecință acțiunea trebuie terminată prin *abort*.

Schemele de sincronizare optimiste permit cel mai mare grad de concurență posibil la nivelul unui obiect; acțiunile nu sunt niciodată puse în așteptare, ca în schemele pesimiste. Problema majoră a schemelor optimiste este însă faptul că acțiuni executate fără erori pot fi terminate prin abort; în plus trebuie menținute în memorie mai multe copii ale obiectului, iar în memoria secundară înregistrate toate modificările realizate de aceste acțiuni.

O schema pesimistă evită regia de sistem sporită datorată reexecuției cererilor și inversării efectuării modificărilor unei acțiuni pe seama reducerii concurenței. De exemplu, două acțiuni care examinează și modifică părți distincte ale unui aceluiasi obiect nu pot fi executate concurent, desi nu există nici un conflict de informații. O schema pesimistă, pe de altă parte, evită regia de sistem datorată întârzierii execuției operațiilor asupra unui obiect pe seama execuției unor operații de inversare a efectului modificărilor ("*undo*") sau a instalării modificărilor ("*redo*").

De aceea, o schema pesimistă se comporta mai bine decît o schema optimistă cînd conflictele sunt frecvente pe cînd, în cazul contrar, cînd conflictele sunt mai puțin frecvente, o schema optimistă se comportă mai bine decît o alta pesimistă.

6.3.3. Siguranța obiectelor

Un SPDBO trebuie să fie capabil să detecteze și recupereze obiectele în cazul de defect al unui nod. Aceasta este o caracteristică importantă deoarece cu cît numărul de componente într-un sistem crește cu atît crește și probabilitatea apariției unui defect. Există două metode generale de creștere a fiabilității unui SPDBO: o prima metoda propune recuperarea unui obiect imediat după detectarea și înlăturarea defectului, limitînd la maximum timpul în care obiectul rămîne inactiv; a doua metodă propune replicarea obiectelor pe mai multe noduri.

Metoda recuperării obiectelor are două variante: varianta recuperării *roll-back* și varianta recuperării *roll-forward*.

În varianta recuperării *roll-back* un obiect defect este restaurat la ultima sa stare consistentă înregistrată în memoria secundară de către procedura *commit*. Toate acțiunile în curs de deslășurare în momentul apariției defectului vor fi fără efect (vor fi pierdute).

În varianta recuperării *roll-forward* un obiect defect, *precum și toate obiectele cu care acesta a lucrat* sunt restaurate la ultima lor stare consistentă înregistrată în memoria secundară de procedura *commit*. Toate procesele și acțiunile în progres în momentul apariției defectului sunt restartate și executate complet. Deoarece în aceasta variantă de recuperare, sistemul trebuie să apară după recuperare ca și cînd nu ar fi apărut nici un defect, varianta este mai complexă decît cea a recuperării *roll-back*; în plus are dezavantajul că dacă eroarea este software ea va apare din nou cînd obiectul este restaurat.

Metoda obiectelor replicate permite existența în mai multor copii ale unui obiect pe mai multe noduri. Aceasta permite unui SPDBO să tolereze defecte la un anumit număr de noduri menținînd totuși funcționalitatea sistemului în ansamblu. Metoda ridică probleme legate de menținerea consistenței informațiilor între copiele unui obiect și de sincronizarea activității clienților.

Cea mai simplă strategie este de a permite numai *obiectelor imuabile* să fie replicate. (Obiectele imuabile sunt cele care pot fi numai examinate de clienți și deci a căror stare nu poate fi modificată. Acest tip de obiecte pot fi replicate fără probleme de menținerea consistenței stării și a sincronizării accesului.)

O alternativă la strategia de mai sus este strategia *copiei primare*. O copie a obiectului este desemnată drept copia primară, în timp ce restul copiilor sunt ordonate și menținute în stații separate drept copii secundare. Cererile de citire pot fi direcționate către orice copie, în schimb cererile de scriere nu pot fi direcționate decît către copia primară care propagă apoi toate modificările și către copiele secundare. Există două variante ale strategiei copiei primare: *varianta copiei primare statice* și *varianta copiei primare dinamice*. În varianta copiei primare statice, dacă serviciul copiei primare devine indisponibil, cererile de scriere nu sunt servite decît la recuperarea copiei primare. În varianta copiei primare dinamice, dacă serviciul copiei primare devine indisponibil, sistemul desemnează una din copiele secundare drept copie primară, putîndu-se astfel trata în continuare cererile de modificare ale obiectului; problema care apare este ca identitatea copiei primare precum și resursele asociate ei precum porturi și lock-uri să poată fi asociate copiei secundare desemnate ca primare. O altă strategie este *strategia copiilor egale* în care nu se desemnează un obiect primar sau secundar; fiecare copie a obiectului se considera egală cu oricare alta. Atît cererile de citire cît și cele de scriere pot fi servite de oricare copie; în acest caz însă este necesară cooperarea unora dintre copii în scopul servirii corecte a acestor cereri. Există mai multe variante la aceasta strategie: *sistemul votului, al copiilor disponibile sau sistemul regenerării*.

6.4. Interacțiunea obiectelor

Un SPDBO este responsabil și pentru gestionarea interacțiunii obiectelor. Atunci când o acțiune apelează o operație asupra unui obiect, sistemul trebuie să localizeze obiectul respectiv, să transmită cererea către acesta și eventual să returneze răspunsul.

6.4.1. Localizarea obiectelor

Un SPDBO trebuie să furnizeze utilizatorilor sai caracteristica de *transparență* în localizarea obiectelor, astfel încât un client să nu fie obligat să cunoască adresa fizică a nodului în care un obiect este rezident; la orice cerere de execuție a unei operații sistemul trebuie să identifice ce obiect este implicat și unde este rezident. Consecința este că sistemul trebuie să asigneze câte un identificator fiecărui obiect, iar acești identificatori trebuie să fie unici. Mecanismul de localizare poate chiar să fie atât de flexibil încât să permită migrarea obiectelor de la un nod la altul.

O primă soluție a problemelor de mai sus este *codificarea localizării obiectului* chiar în identificatorul obiectului. În momentul apelului, sistemul examinează câmpul respectiv din identificatorul obiectului și decodifică adresa nodului în care obiectul este rezident. Soluția este foarte simplă, eficientă, dar și restrictivă: odată ce un obiect a fost asignat unui nod el nu mai poate fi mutat, deoarece identificatorul obiectului trebuie modificat în acest caz. În consecință, pe toată durata de viață a sa, un obiect este fixat pe un nod.

O altă soluție este strategia *serverului de nume distribuit*. Conform acestei strategii, sistemul crează un grup de servere de nume ale obiectelor, menținute în general, (dar nu și totdeauna necesar) pe toate nodurile. Aceste servere cooperează astfel încât, per total ele conțin toate informațiile necesare pentru localizarea oricărui obiect din sistem. Există două variante ale acestei strategii. În prima variantă, un server de nume menține informații complete pentru localizarea oricărui obiect astfel încât orice server poate deservi orice cerere de localizare. În cea de a doua variantă, un server de nume menține informații parțiale pentru localizare; dacă o cerere de localizare nu poate fi rezolvată de un server de nume se delegă altul s.a.m.d. În aceasta variantă, cel puțin un server trebuie înștiințat atunci când un obiect este creat.

O altă strategie este strategia *cache/broadcast*. În fiecare nod se menține o memorie cache în care se memorează localizarea ultimelor "n" obiecte referite. La orice referire a unei cereri de operații asupra unui obiect se examinează această memorie cache; dacă localizarea obiectului nu este găsită în aceasta listă, se trimite un mesaj broadcast în rețea pentru a se determina nodul care conține obiectul respectiv. Fiecare nod care primește acest mesaj, caută în propria sa memorie cache pentru a vedea dacă deține acest obiect; în caz afirmativ trimite un mesaj de confirmare la nodul emițător, care își actualizează memoria cache relativ la acest obiect. Strategia este foarte eficientă, deoarece este de așteptat ca localizarea obiectelor să fie găsită în majoritatea cazurilor în memoria cache locală. Este, deasemenea, flexibilă deoarece permite migrarea unui obiect, fără a fi nevoie de transmisie de mesaje la alte noduri sau la un server de nume. Dezavantajul este faptul că trebuie transmise mesaje broadcast, ceea ce înseamnă o încărcare mare a rețelei; având drept efect conturbarea funcționării mai multor noduri deși unul singur este implicat în dialog.

6.4.2. Apelul de operații

Atunci când un client apelează o operație asupra unui obiect, sistemul SPDBO este cel care trebuie să transmită cererea către serverul obiectului respectiv și să returneze răspunsul înapoi la apelant. Există două strategii pentru gestionarea acestor cereri: *strategia transmiterii de mesaje* și *strategia apelului direct*.

Prima strategie, cea a *transmiterii de mesaje* este, în mod tipic folosită în sistemele care implementează modelul obiectului activ. Când un client apelează o operație asupra unui obiect, parametrii cererii sunt împachetați în mesaj; mesajul este transmis procesului server sau la un port asociat obiectului respectiv. Procesul server recepționează mesajul, despachetează parametrii din mesaj, execută operația și împachetează rezultatul într-un alt mesaj pe care-l transmite apelantului. Legătura client-server este dinamică și se realizează la fiecare apel. Un dezavantaj este complexitatea dialogului, pentru apeluri în același nod.

Cea de-a doua strategie, a *apelului direct* este proprie sistemelor care implementează modelul obiectului pasiv, în care un singur proces este responsabil pentru execuția tuturor operațiilor asociate unei acțiuni. Consecința este ca un proces poate migra în funcție de operație de la un obiect la alt obiect.

Atunci când un proces apelează o operație asupra unui obiect rezident pe aceeași stivă, execută următoarele patru etape (figura 6.6.) [CC91]:

- (1) Starea procesului și obiectele apelate se înregistrează în zona stivă a procesului. Sistemul poate proteja aceste zone împotriva examinării sau alterării ei de alte apeluri.
- (2) Parametrii apelului sunt adăugați în stivă.
- (3) Obiectul apelat este încărcat în memorie și se execută un apel de procedură pentru operația cerută.
- (4) La terminarea operației, rezultatele sunt returnate clientului și procesul este în starea avută înaintea apelului.

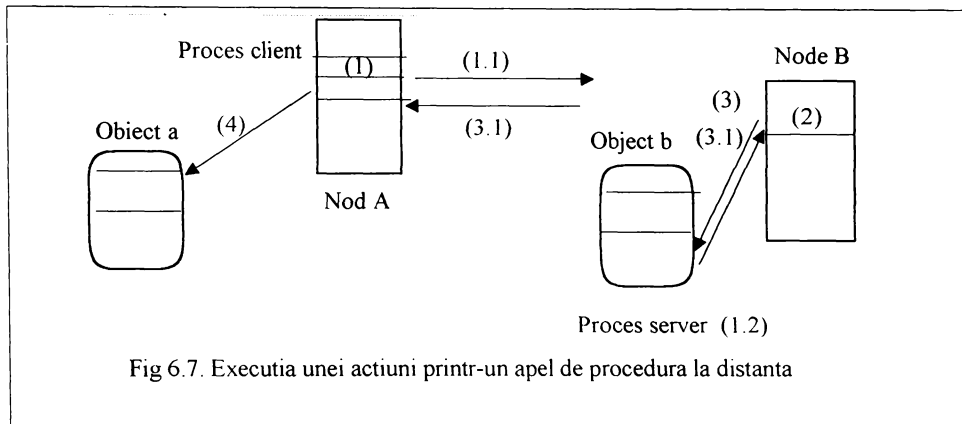
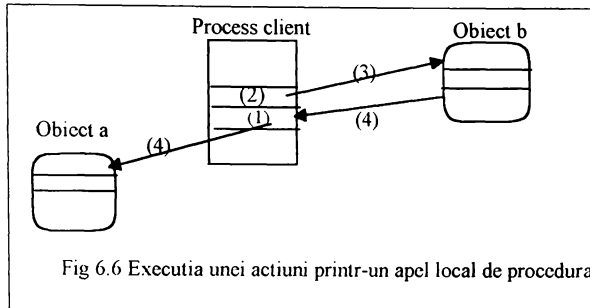
Atunci când un proces execută o acțiune asupra unui obiect rezident în alt nod, se execută în plus următoarele trei etape, primele două după (1), iar a treia după (3) (fig 6.7.):

- (1.1) Se formează un mesaj cuprinzând parametrii cererii care se trimite nodului în care obiectul este rezident

(1.2) Nodul care primește mesajul crează un proces care se execută în continuarea procesului curent

(3.1) Când operația a fost executată, se formează un mesaj conținând rezultatele care se transmite nodului clientului. Procesul din nodul de la distanță este apoi distrus.

Apelul direct apare astfel asemănător unui apel de procedură la distanță. Metoda are dezavantajul creerii și distrugerii de noi procese în nodurile în care obiectele sunt rezidente.



6.4.3. Detectarea eșecurilor în apelarea operațiilor

Eșecurile în apelul unei operații se clasifică în *eșecuri permanente* și *eșecuri tranzitorii*. Un *eșec permanent* se definește ca un eșec apărut înainte de începerea execuției operației. Cea mai cunoscută situație de acest gen este cea în care obiectul apelat nu poate fi localizat. Aceste tipuri de eșecuri sunt ușor de detectat și tratat.

Un *eșec tranzitoriu* este un eșec care apare în timp ce apelul operației se execută. Un astfel de eșec apare deci după acceptarea unei operații, dar înainte ca modificările făcute să devină permanente prin desăvârșirea acțiunii (înainte de terminarea prin "commit"). Aceste eșecuri sunt mai greu de detectat și de tratat deoarece pot avea cauze diverse. De exemplu, un eșec al unui apel de operație se poate datora eșecului obiectului client, obiectului server sau liniilor de comunicație. Un SPDBO trebuie să prevadă mecanisme pentru detectarea eșecurilor atât la nivelul obiectului client cât și la cel al serverului precum și pentru recuperare după astfel de eșecuri.

Sunt cunoscute mai multe metode pentru detectarea unor astfel de eșecuri: bazate pe mecanismul time-out, pe teste ale obiectului sau ale apelurilor. Dacă un eșec tranzitoriu nu este detectat de client, clientul și acțiunea corespunzătoare pot fi blocați la nesfârșit. În consecință, un client trebuie să poată detecta eșecul unui apel de operație și să poată iniția o procedură de recuperare a eșecului. Procedura de recuperare trebuie să elibereze resursele deținute de client și să termine acțiunea cu eroare, eventual să repete apelul operației.

Dacă un eșec tranzitoriu nu este detectat de obiectul server, resurse ale sistemului pot fi irosite nedefinit. Un obiect startat în urma unui apel de operație, ale cărei rezultate însă nu mai sunt preluate, este denumit *orfan*; el poate fi generat datorită unui eșec al unui obiect client, datorită unei tranzacții terminate cu abort, sau datorită unei erori ale rețelei și irosește resurse precum lock-uri și timp procesor. În consecință, trebuie eliminat cât mai repede posibil (vezi 3.4.5. pentru metode de eliminare a orfanilor).

6.5. Reprzentarea obiectelor în memorie

Un SPDBO trebuie să furnizeze mecanisme pentru gestionarea resurselor precum memoria principală, dispozitive de memorie externă, noduri ale rețelei. Specific SPDBO este modul de reprezentare în memoria internă și secundară a obiectelor, precum și modul de transfer al obiectelor între cele două, ca și cum se asignează obiectele procesoarelor

Obiectele a căror reprezentare este pierdută în urma defectării nodului se numesc *obiecte volatile*. Obiectele volatile sunt temporare, sunt rezidente numai în memoria internă și deci ușor de manipulat. Obiectele a căror reprezentare se menține după defectarea nodului, în limitele unei mari probabilități, se numesc *obiecte persistente*. Reprezentarea obiectelor persistente trebuie făcută obligatoriu și în memoria secundară; una sau chiar mai multe copii ale acestui obiect se poate găsi și în memoria internă. Acțiunile modifică versiunea din memorie a unui obiect persistent. Versiunea din memoria secundară este modificată numai după ce acțiunea se termină prin *commit*, asigurându-se astfel menținerea unei versiuni consistente și stabile a obiectului, însă cu prețul regiei de sistem sporite prin realizarea operațiilor de copiere în memoria secundară.

Un obiect persistent rezident atât în memoria internă cât și în memoria secundară este numit *obiect activ*. Un obiect persistent având reprezentarea numai în memoria secundară este numit *obiect inactiv*. La orice apel al unei operații asupra unui obiect inactiv se copiaza o copie a obiectului în memoria internă; eventual se pot crea și noi procese pentru gestionarea acestei copii volatile. Dacă procedura de *commit* relativă la obiect se termină cu succes, obiectul este copiat în memoria secundară (actualizat), iar memoria internă ocupată de copia sa dealocată. Aceasta salvează automat a obiectelor active persistente în memoria secundară și încărcarea lor din memoria secundară trebuie realizată transparent de SPDBO. Acesta trebuie să ascundă faptul că se folosește memorie secundară pentru memorarea obiectelor, în aparență fiind vorba numai de obiecte întotdeauna disponibile. Modul în care un obiect este reprezentat în memorie depinde de schema de sincronizare folosită și de modul de gestionare a acțiunilor prevazute în SPDBO. Schema de sincronizare influențează numărul de versiuni ale obiectului, iar modul de gestionare al acțiunilor influențează reprezentarea fiecărei versiuni.

Dacă sistemul SPDBO utilizează o schemă de sincronizare pesimistă, se menține în memorie, în general o singură versiune a obiectului. Toate acțiunile care apelează operații asupra unui același obiect modifică aceeași versiune volatilă. Dacă schema de sincronizare este optimistă se mențin în memorie noi versiuni ale fiecărui obiect creat. Fiecărei acțiuni care apelează un obiect i se asignează propria versiune volatilă a obiectului asupra căreia operează toate modificările. Astfel se permite ca multiple acțiuni să apeleze același obiect simultan fără interferențe.

Reprezentarea fiecărui obiect în memorie depinde de prezența sau absența facilității de tranzații încuibate. Dacă nu se utilizează tranzații încuibate, este suficientă reprezentarea obiectului cu o singură copie în memorie. Dacă tranzația eșuează, este suficientă dezalocarea copieii respective. În cazul când sunt permise tranzații încuibate, lucrurile se complică deoarece mai multe subacțiuni pot modifica același obiect. Modificările asupra obiectului nu sunt permanentizate decât cînd acțiunea de la nivelul cel mai superior execută *commit* și, pe de altă parte, fiecare subacțiune poate fi terminată cu *commit* sau poate eșua. Dacă o subacțiune modifică un obiect și apoi se termină prin abort, modificările trebuie invalidate ("undone") astfel încît obiectul să poată fi restaurat la starea de dinaintea execuției subacțiunii. În consecință, este necesar un mecanism suplimentar pentru a se putea inversa efectul modificărilor unei subacțiuni. Sunt cunoscute două procedee: utilizarea unui jurnal undo/redo sau menținerea unei copii "hotstandby", o versiune a obiectului creat de subacțiune.

Un sistem SPDBO trebuie să mențină în memoria secundară suficiente informații pentru ca un obiect persistent să poată fi restaurat la o stare consistentă, ca de exemplu după un defect al nodului în care obiectul este rezident. Un sistem SPDBO poate înregistra întreaga stare a unui obiect în momentul execuției *commit* sau poate înregistra numai modificările efectuate asupra obiectului, relative la o stare înregistrată anterior.

În prima metoda, denumită *metoda punctelor de reluare (checkpoint scheme)* se înregistrează în memoria secundară întreaga stare a obiectului modificat, cînd acțiunea respectivă execută *commit*. În faza de pre-*commit* a acțiunii versiunea modificată este scrisă în memoria secundară, fără interferență cu versiunea veche a obiectului. În faza de *commit*, versiunea modificată înlocuiește versiunea veche care este apoi dealocată. Dacă se execută însă abort, se dealocă versiunea modificată în timp ce versiunea veche rămîne nemodificată.

Un avantaj al acestei metode este faptul că menține în memoria secundară numai o singură copie a obiectului. Dezavantajul este că și în cazul unor puține modificări asupra unui obiect avînd informații de stare de dimensiuni mari trebuie înregistrate integral în memoria secundară.

O variantă a acestei metode este *metoda istoriei punctelor de reluare*. Metoda propune reprezentarea unui obiect persistent în memoria secundară ca o colecție ordonată de puncte de reluare. În momentul în care trebuie creată o nouă versiune a unui obiect, vechea versiune nu este distrusă; noua versiune se adaugă secvenței de puncte de reluare. În afara de dezavantajul spațiului additional necesar în memoria secundară, metoda mai are și dezavantajul problemei determinării momentului în care vechile puncte de reluare nu mai sînt necesare și, de aceea ștersc.

O metoda complet diferită este cea a *jurnalelor log*. Metoda propune ca numai modificările executate asupra unui obiect persistent să fie înregistrate într-un jurnal menținut în memoria secundară (vezi 7.3.4.). Cantitatea de informație înregistrată în memoria secundară depinde de schema particulară de înregistrare în jurnal folosită. Evident că, față de metoda precedentă, metoda jurnalelor log folosește spațiul memoriei secundare mai eficient. O consecință a memorării numai a modificărilor este faptul că, cu cît jurnalul log crește, cu atît crește și regia de sistem necesară menținerii lui ca și pentru tratarea recuperărilor.

Capitolul 7

CONTROLUL CONCURENȚEI, INTERBLOCĂRI SI RECUPERAREA OBIECTELOR

Acest capitol tratează în principal al doilea aspect al sincronizării în sisteme distribuite, prezentat în introducerea capitolului 4. Alături însă de problemele de control al comunicației și de tranzacțiile atomice sunt tratate și alte probleme legate de acestea, precum recuperarea tranzacțiilor și controlul interblocărilor.

7.1 Introducere

Se apreciază că proiectarea aplicațiilor distribuite fiabile poate fi obținută prin două metode generale, conceptual diferite: (1) prima metodă utilizează conceptul de *tranzacție atomică*; (2) a doua metodă utilizează copii multiple ale obiectelor necesare aplicațiilor, soluționând corect problema actualizării copiilor ([Gos91], [Tan93], [Tan96], [CC91]). Problemele în legătură cu controlul concurenței sunt legate în sistemele distribuite de problemele de comunicație; s-a aratat în capitolul 4 că sincronizarea în sistemele distribuite nu se poate realiza fără comunicație.

Controlul concurenței este și parte componentă a protocolului OSI CCR (*Commitment, Concurrency and Recovery*) care guvernează desfășurarea tranzacțiilor atomice. Controlul concurenței se consideră un aspect comun al sistemelor de operare distribuite și bazelor de date distribuite. Conceptul de tranzacție a apărut inițial în bazele de date, ca un instrument în menținerea consistenței datelor; consistența trebuia menținută atât în cazuri de avarii ale unor noduri, cât și în cazul acceselor concurente. Bazele de date nu sunt însă singurele aplicații în care trebuie menținută consistența datelor, în situațiile critice de mai sus; majoritatea sistemelor de fișiere pentru rețele de calculatoare de astăzi oferă facilități pentru blocarea la nivel de înregistrare fizică, blocare la nivel de înregistrare logică sau la nivel de fișier și, în special facilități pentru utilizarea tranzacțiilor atomice. Un exemplu este sistemul de gestiune fișiere Novell Netware care pune la dispoziția utilizatorilor API-uri de sincronizare (*Synchronization Services*) între care figurează apeluri pentru manipularea lock-urilor și pentru lucrul cu tranzacții atomice.

Un alt exemplu este sistemul distribuit de fișiere QuickSilver Distributed File System (QDFS) ([CPS93]). Serverele pentru fișiere ale acestui sistem furnizează accesul la fișiere într-un spațiu de nume structurat ca un arbore și prezintă o interfață similară cu cea din UNIX, cu excepția faptului că nu permite accesul aleatoriu la fișiere. Spre deosebire de UNIX, toate operațiile implementate în DFS sunt atomice. Aceasta înseamnă că, de exemplu, un client poate șterge, rescrie un fișier sau redenumi un catalog și apoi decide dacă modificările vor fi permanentizate sau nu. QDFS garantează că, fie toate modificările vor fi permanentizate, fie niciuna nu va fi realizată; aceasta garanție se face chiar pentru situații de avarii în rețea (defecțiuni ale nodurilor, ale liniilor de comunicație).

Tranzacțiile atomice au fost incluse de asemenea în Argus (limbajul CLU) ([Lis88], [LCJ87], [Lis90], [CC91][TCM94]), în sistemul de operare distribuit Clouds ([For88][CC91], [TCM94], [Tan96]), în TABS și Camelot ([Spe90], [SDE85], [SBD86], [CC91], [TCM94]), în Eden ([ABN85], [CC91]), în Arjuna ([CC91]), Encina (Transarc Corporation, 1991, [CC91]), Tuxedo (UNIX System Laboratories, 1991, [CC91]). Se apreciază că, la ora actuală, deși suportul tranzacțiilor atomice migrează de la bazele de date la sistemele de operare, acest lucru nu este încă unanim acceptat ([CPS93], [KS89]).

Pe de altă parte, trebuie precizat că modelul tranzacțiilor clasice, adoptat în bazele de date nu satisface cerințele unor aplicații complexe, precum aplicații CAD/CAM, medii de dezvoltare de software (SDT-uri), gestiunea rețelelor, informatică medicală, aplicații care sunt orientate pe date masive, în sensul că generează și manipulează volume mari de date. Este de dorit ca pentru aceste tipuri de aplicații să existe facilități similare cu cele tradiționale în bazele de date precum adăugarea, ștergerea, regăsirea datelor, alături de menținerea consistenței informațiilor stocate, cu observația că mecanismul tradițional al tranzacțiilor nu este adecvat. Astfel, aplicații precum gestiunea rețelelor și informatică medicală pot necesita prelucrări în timp real; altele, precum proiectare CAD/CAM sau medii de dezvoltare software necesită cooperare între mai mulți utilizatori. Controlul concurenței în manieră convențională nu este aplicabil pentru aceste cazuri, pentru care se utilizează termenii de *tranzacții lungi* și *tranzacții cooperante*. Unii cercetători susțin chiar că termenul de tranzacție în sensul clasic își pierde semnificația pentru tranzacțiile lungi și cele cooperante, din moment ce acestea nu mai conservă proprietățile de *atomicitate* și *seriabilitate*. Pe parcursul acestui capitol, în secțiunea 7.4.7 în special, se prezintă și cerințele specifice și mai ales soluții legate de controlul adecvat al concurenței pentru astfel de aplicații, denumite *aplicații complexe*.

7.2 Modelul tranzacțional

În cele ce urmează, se va utiliza noțiunea de *sistem de obiecte*.

Definiție.

Se numește *sistem de obiecte* o colecție de obiecte (de exemplu înregistrări sau fișiere) văzută ca o resursă unică, asupra căreia se execută operațiile de bază, denumite *acțiuni*, *read* și *write*.

Prelucrările asupra sistemului de obiecte pot fi văzute ca secvențe de acțiuni aplicate unui subset de obiecte ale sistemului de obiecte.

7.2.1 Consistența obiectelor. Noțiunea de tranzacție

Fiecare sistem de obiecte trebuie să satisfacă un set de restricții denumite *constrângeri de consistență*. Exemple clasice sunt cel ale unei baze de date bancare, la care o constrângere este aceea că balanța fiecărui cont trebuie să fie egală cu suma depunerilor minus suma extragerilor și cel al unui sistem de rezervare de locuri la care o constrângere este faptul că fiecare loc poate fi rezervat numai de un singur pasager. Starea bazei de obiecte este *consistentă* dacă sunt respectate toate constrângerile de consistență. Uneori, constrângerile de consistență pot fi foarte complexe. Pentru a trata și rezolva problemele de consistență a fost introdus conceptul de *tranzacție*.

Definiție.

O *tranzacție* este o secvență de acțiuni, care constituie unitatea logică de lucru cu sistemul de obiecte și care transferă sistemul de obiecte dintr-o stare consistentă în altă stare consistentă.

În bazele de date tradiționale, tranzacțiile servesc pentru trei scopuri: (1) sunt unități logice care grupează acțiunile unei prelucrări; (2) sunt unități atomice a căror execuție conserva consistența bazei de date; (3) sunt unități de recuperare.

Pentru a crește performanțele, este indicat să se permită execuția concurentă a tranzacțiilor. Totuși, intermixarea arbitrară a acțiunilor tranzacțiilor concurente are deseori ca rezultat o stare inconsistentă; în consecință, este necesar controlul acceselor concurente, sau pe scurt *controlul concurenței*. Există două situații de anomalie care generează stările de inconsistență: *anomalia actualizărilor pierdute*, și *anomalia regasirilor inconsistente*. Exemplul din fig. 7.2.1.1 ilustrează aceste situații.

Dacă tranzacțiile T_1 și T_2 , din acest exemplu ar fi executate concurent, prin intermixarea acțiunilor s-ar putea obține o balanță finală a contului B de 500000 sau 200000; aceasta este *anomalia actualizărilor pierdute*. Dacă tranzacția T_3 citește conturile A și B în timpul execuției tranzacției T_1 , atunci afișarea rezultatelor poate conduce la situația în care o sumă de bani transferată dintr-un cont în altul dispăre; aceasta este situația de *regasire inconsistentă* (A este debitor dar B nu a fost creditat). Deși T_1 și T_2 , individual, păstrează consistența datelor, ele prezintă însă următoarele caracteristici: (1) *inconsistența temporară*: după prima operație *write* fie a lui T_1 , fie a lui T_2 , sistemul de obiecte este inconsistent; (2) *conflict*: dacă tranzacția T_2 se execută între prima și a doua operație *write* a lui T_1 , atunci suma finală a balanței conturilor va fi de 1 600 000 și nu de 1 400 000, contrazicând cerințele de consistență.

O tranzacție poate fi terminată *normal* sau *anormal*. În primul caz se spune că este *realizată* sau *comisă*; în cel de-al doilea caz că este *abortată*. Tranzacția însăși poate cere să fie abortată, ca urmare, de exemplu, a unor date de intrare eronate; pe de altă parte, o tranzacție poate fi însă abortată și de sistem, ca urmare a violării unor constrângeri de consistență.

Pentru a se respecta constrângerile de consistență impuse unui sistem de obiecte, fără a se impune condiții excesiv de restrictive asupra prelucrărilor concurente ale aplicațiilor, o tranzacție trebuie să aibă patru proprietăți cunoscute sub numele de *proprietăți ACID* în literatură ([CPS93], [BHG87], [SDE85]): (1) *Proprietatea de atomicitate* este proprietatea conform căreia o tranzacție fie că se realizează, caz în care toate actualizările asupra obiectelor se permanentizează, fie este abortată, caz în care

Balanta initiala a conturilor
A 800 000
B 400 000
C 200 000
Tranzactii
T1 :
BeginTransaction
read cont A
read cont B
write balanta A - 100 000 in contul A
write balanta B + 100 000 in contul B
EndTransaction
T2 :
BeginTransaction
read cont B
read cont C
write balanta B - 200 000 in contul B
write balanta C + 200 000 in contul C
EndTransaction
T3 :
BeginTransaction
read cont A
read cont B
print balanta A
print balanta B
EndTransaction

Fig. 7.2.1.1 Tranzactii concurente

efectul ei este nul asupra obiectelor. (2) Proprietatea de *consistență* este proprietatea conform căreia o tranzacție transformă un sistem de obiecte dintr-o stare consistentă în alta. (3) Proprietatea de *izolare* sau de *scriabilitate* este proprietatea care garantează că efectul tranzațiilor care se execută concurrent este echivalent cu efectul execuției lor în mod serial, într-o ordine oarecare. (4) Proprietatea de *durabilitate* sau de *persistență* a efectului este proprietatea conform căreia efectele unei tranzații realizate rămân valabile și după o defecțiune a sistemului.

O tranzacție care execută operații asupra unor obiecte dispersate într-un sistem distribuit se numește *tranzacție distribuită*. Distribuirea obiectelor pe mai multe noduri crește performanța sistemului legată de concurență și incurajează utilizarea copiilor multiple pentru un obiect. Pentru a crește și mai mult gradul de concurență, tranzațiile pot fi *încuibate* ([Spe84]); tranzațiile încuibate se caracterizează prin următoarele elemente: (1) o tranzație de nivel superior poate iniția mai multe tranzații încuibate care se execută în paralel; acestea, la randul lor, pot crea alte tranzații, s.a.m.d. În modelul general, o tranzație părinte își suspendă activitatea până când fiii săi se termină (normal sau anormal); (2) o tranzație fii poate obține un *lock* (folosit de mecanismul de concurență, vezi subcapitolul 7.4.4.) de la o tranzație ascendentă; (3) când o tranzație se termină cu *commit*, toate *lock*-urile sale, inclusiv cele moștenite, sunt returnate tranzației părinte; astfel, numai la terminarea tranzației celei mai exterioare, se eliberează *lock*-urile. Analog, efectul tranzațiilor încuibate se permanentizează numai când tranzația cea mai exterioară se termină (cu *commit*); (4) dacă are loc un incident hardware, ce generează abortarea unei tranzații interioare, atunci toate *lock*-urile deținute de acea tranzație sunt returnate tranzației părinte, care este anunțată de terminarea anormală. Tranzacția părinte poate continua sau se poate aborta ca însăși.

Se apreciază că tranzațiile atomice, utilizate împreună cu mecanismul RPC, cu tehnici de autentificare și protecție, și cu un suport lingvistic adecvat, pot constitui o bază importantă pentru construcția de programe sigure (Gos91).

7.2.2 Modelul de lucru pentru sistemul distribuit de obiecte

Se prezintă în această secțiune modelul de lucru al sistemului distribuit de obiecte: *modelul obiect-entitate*.

Colecțiile de obiecte sunt împrăștiate pe mai multe noduri; aceste noduri sunt interconectate într-o rețea și identificate în mod unic printr-un număr. Nodurile comunică între ele prin mesaje, livrate cu o întârziere oarecare în cadrul rețelei; se presupune însă că toate mesajele sunt livrate, și, în plus, sunt livrate în ordinea în care au fost emise. Sistemul distribuit de obiecte constă într-un set de *entități* (de exemplu înregistrări sau fișiere), care primesc denumiri unice și servesc ca unitate indivizibilă pentru acces. O entitate se materializează în unul sau mai multe obiecte de date, fiecare identificat prin perechea <nume_entitate> <identificator_nod>. Replicarea unei entități are loc atunci când o entitate este reprezentată de mai multe obiecte identice. Replicarea entităților poate conduce la o îmbunătățire a performanțelor dacă costul memorării și menținerii obiectelor copii este mai mic decât costul de acces; de asemenea, replicarea poate mări viteza de acces la obiecte și siguranța de funcționare a sistemului, în sensul că dacă un nod devine inaccesibil, accesul este posibil la un nod care deține o altă copie a obiectului.

În figura 7.2.2.1 se arată corespondența între entități și obiecte.

O entitate E_i , replicată în nodurile S_1, \dots, S_k este reprezentată prin obiectele $\langle E_i, S_1 \rangle, \dots, \langle E_i, S_k \rangle$. Fiecare tranzație T este supervizată printr-un manager de tranzație, al cărui scop principal este de a supraveghea prelucrările distribuite. Se definesc patru primitive la interfața tranzație-manager: *read*(F_i)

care returnează valoarea entității E_i , *write*($E_i, \text{valoare}$) care setează valoarea entității E_i la valoarea argumentului și primitivele *BeginTransaction* și *EndTransaction* care delimitează o tranzație. Fiecare cerere de citire a unei entități este translatată de managerul de tranzații într-o acțiune *read* de forma: *read*($T, \text{read obiect} \langle E_i, S \rangle, V$) asupra obiectului corespunzător; analog, fiecare cerere de scriere a unei entități este translatată într-un set de acțiuni *write* de forma: *write*($T, \text{write obiect} \langle E_i, S \rangle, V$), cite o acțiune *write* pentru fiecare obiect care constituie o copie (dacă se utilizează replicarea entităților). În plus, pentru a

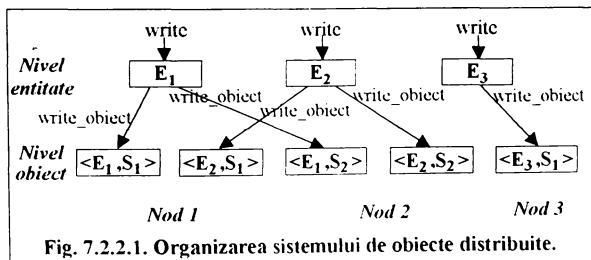


Fig. 7.2.2.1. Organizarea sistemului de obiecte distribuite.

comite tranzația, managerul de tranzații generează cite o primitivă *commit* de forma ($T, \text{commit obiect} \langle E_i, S \rangle$), pentru fiecare nod S la care s-a executat o acțiune.

Managerii de obiecte sunt procesele care gestionează obiectele și execută asupra obiectelor acțiunile cerute de managerul de tranzații. Ca răspuns la o acțiune *read*, de exemplu, managerul de obiecte returnează managerului de tranzații valoarea obiectului. O acțiune *commit obiect* are ca efect permanentizarea efectelor acțiunilor *write* dintr-o tranzație asupra aceluși obiect. La un moment dat, pentru o tranzație pot fi în lucru mai mulți manageri de obiecte, supravegheați de managerul de tranzație la care s-a deschis tranzația cu primitivă *BeginTransaction*. Execuția cererilor la nivelul unui manager de obiecte, pentru același obiect, se face secvențial.

În modelul prezentat se va considera că fiecare entitate fie nu este replicată, fie este complet replicată (în toate nodurile). Se consideră că această restricție nu este o restricție conceptuală majoră, deoarece o entitate parțial replicată poate fi tratată ca o grupare de entități distincte punându-se o constrângere globală de consistență: aceea că toate entitățile distincte să conțină aceleași valori. Se impune ca o cerere de scriere pentru o entitate să genereze actualizarea fiecărei copii a unei entități, în timp ce o cerere de citire să acceseze o singură copie. Cu aceasta manieră de a trata replicarea, modelul poate fi generalizat la modelul cu voturi al lui Gifford (vezi secțiunea 1.3.2.2). Pe

parcursul acestui capitol, dacă nu se va specifica explicit, se va considera că o entitate nu este replicată, ceea ce implică mapare unu la unu entitate-obiect.

Există mai multe modele pentru descrierea și reprezentarea unei tranzacții într-un sistem distribuit de obiecte. Primul model este acela în care se creează câte un *procces agent* în fiecare nod implicat în execuția unei tranzacții; procesul agent obține accesul la obiectul respectiv și va elibera obiectul după îndeplinirea cererii și, eventual, va returna un semnal de sfârșit acces apelantului sau. În acest model (Menasc și Muntz) se consideră că procesele agent se creează când este necesar acces la un obiect situat la distanță și se distrug când accesul a fost realizat. Un al doilea model (Gray) consideră că procesele de la distanță se creează la începutul tranzacției și rămân în viață până la terminarea ei. Un al treilea model (Rosenkrantz) consideră că o tranzacție originară dintr-un nod migrează dintr-un nod în altul, după cum se generează acțiunile asupra obiectelor. În cele ce urmează în acest capitol se va considera primul model.

Există două modele de arhitectură în ceea ce privește managerii de tranzacții și managerii de obiecte: *modelul integrat* și *modelul client-server*. În *modelul integrat* (fig.7.2.2.2), toate nodurile sunt dotate cu managerii de tranzacție eventual, cu managerii de obiecte dacă există sisteme de obiecte pe acele noduri. Modelul poate fi folosit în rețele locale dar și în rețele de largă răspândire geografică. În *modelul client-server* (fig.7.2.2.3), gestiunea tranzacțiilor (clienți) este separată de gestiunea obiectelor (servere). Acest model este utilizat numai în rețelele locale.

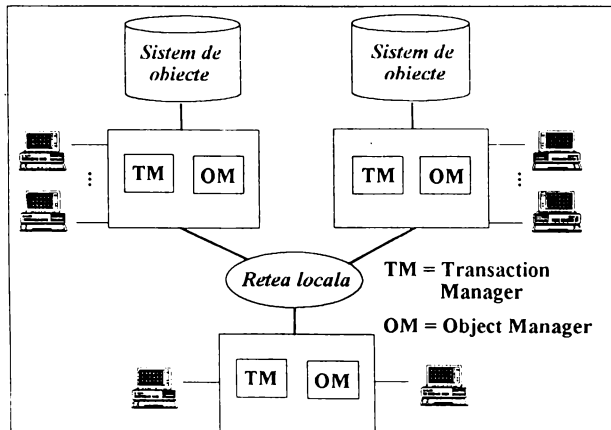


Fig. 7.2.2.2. Modelul integrat.

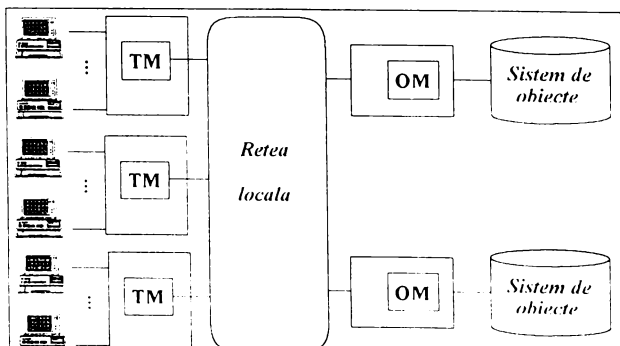


Fig. 7.2.2.3. Modelul client-server.

7.2.3 Cerințele modelului sistemului distribuit de obiecte

Un sistem distribuit de obiecte trebuie să îndeplinească următoarele *cerințe*: (1) *Integrarea strictă a datelor*. Pentru regăsirea și gestiunea datelor, sistemul trebuie să implementeze o interfață uniformă pentru toate instrumentele de acces la date. (2) *Orientarea pe aplicație*. Sistemul trebuie să organizeze datele în structuri capabile să capteze cât mai mult din semantica specifică aplicației. (3) *Integritatea datelor*. Sistemul trebuie să asigure consistența și recuperabilitatea datelor, pentru a garanta astfel satisfacerea constrîngerilor de consistență cerute de aplicație. (4) *Facilizarea accesului*. Sistemul trebuie să asigure o interfață (eventual un limbaj de interogare) pentru a accesa un set de obiecte la un anumit moment. (5) *Independența datelor*. Sistemul trebuie să ascundă structura internă a datelor interfeței de acces, astfel încît dacă structura se modifică, impactul asupra aplicației care utilizează datele să fie cît mai mic.

- Detaliile referitoare la distribuire, replicare și concurență trebuie să fie transparente pentru utilizator. Există patru cerințe de transparență care se aplică unui sistem distribuit de obiecte: (1) *Transparența localizării*. Deși datele sunt distribuite pe mai multe noduri și pot migra de pe un nod pe altul, programele nu trebuie să cunoască localizarea lor. (2) *Transparența replicării*. Programele de aplicație trebuie să trateze fiecare entitate ca un singur element, fără a cunoaște dacă entitatea este sau nu replicată. (3) *Transparența concurenței*. Programele percep propriile lor tranzacții ca fiind singura activitate din sistem. (4) *Transparența defecțiunilor*. Acțiunile tranzacțiilor se materializează fie în totalitate, fie deloc. După realizarea unei tranzacții, efectele ei trebuie să persiste chiar și în cazul unui defect hardware sau a unei erori software.

7.3 Realizarea atomicității tranzacțiilor în sistemele distribuite

S-a aratat în subcapitolul 7.2.1 că conceptul de *tranzacție atomică*, caracterizat prin cele cinci proprietăți cunoscute sub numele de ACID este conceptul fundamental în structurarea sistemelor software tolerante la defecte; importanța sa constă în faptul că asigură detectarea erorilor și recuperarea obiectelor. Din punct de vedere al atomicității, tranzacțiile se caracterizează în următoarele patru moduri: (1) procesele care execută o acțiune atomică nu au cunoștință despre existența altor procese neimplicate în acțiunea atomică și vice versa; (2) procesele care execută o acțiune atomică nu comunică cu alte procese în timp ce se execută acțiunea; (3) procesele care execută o acțiune atomică nu pot detecta o modificare de stare a obiectelor, alta decît cea datorată lor, și nu pot instala o modificare de stare decît cînd acțiunea atomică s-a terminat; (4) o acțiune atomică este indivizibilă și instantanee în sensul că efectul său asupra obiectelor este același, ca și cînd nu ar fi executată concurent cu alte acțiuni.

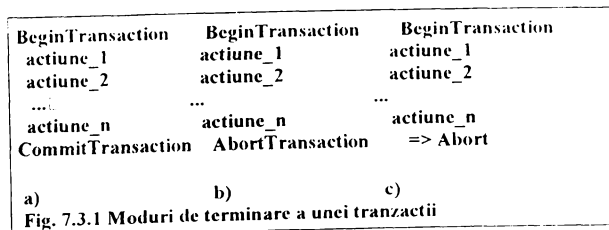
Dacă în timpul execuției unei tranzacții atomice apare o eroare care împiedică terminarea normală a ei, starea tuturor obiectelor modificate de această tranzacție trebuie restaurată la starea de dinaintea începerii tranzacției; se spune atunci că efectul actualizărilor este anulat și că tranzacția a fost recuperată. În acest sens, o tranzacție ale cărei actualizări pot fi realizate fără a anula efectele altor tranzacții se spune că este *recuperabilă*.

Recuperabilitatea și atomicitatea tranzacțiilor pot fi realizate dacă se impun următoarele două cerințe implementării tranzacțiilor: (1) obiectele actualizate nu vor fi eliberate decît la terminarea tranzacției; (2) starea inițială a tuturor obiectelor actualizate de tranzacție poate fi reconstruită. Prima cerința poate fi satisfăcută prin utilizarea unui protocol adecvat pentru comiterea tranzacției, iar cea de-a doua poate fi satisfăcută prin utilizarea unui mecanism adecvat pentru recuperare.

7.3.1 Protocele pentru comiterea tranzacțiilor

Figura 7.3.1. arată trei moduri posibile de desfășurare a unei tranzacții. Dacă tranzacția a ajuns în final, se execută o acțiune *commit* pentru a permanentiza modificările asupra obiectelor (7.3.1.a). Dacă tranzacția a detectat condiții pentru care nu poate fi încheiată normal (erori ale datelor de intrare, de exemplu), atunci ea se poate aborta (termină anormal) ca însăși (7.3.1.b). O tranzacție poate fi abortată și de către sistem datorită unor cauze externe: time-out, rezolvarea unei interblocări, violarea protecțiilor, erori pe liniile de comunicație, avarii ale nodului (fig.7.3.1.c.).

Deoarece într-o tranzacție distribuită participa mai multe noduri, este nevoie de un protocol pentru a se obține consensul tuturor participanților pentru comiterea tranzacțiilor. Se apreciază că problema comiterii tranzacțiilor distribuite este asemănătoare cu problema încercării închiderii unei conexiuni într-un sistem cu linii de



comunicații nesigure; aceasta problemă a fost comparată cu vestita problemă a generalilor ([MOO87]). De exemplu, presupunem că o sumă de bani este transferată dintr-un cont în altul din București în Timișoara; ambele conturi, contul debitor din București și contul creditor din Timișoara trebuie, fie amândouă procesate (terminate cu *commit*), fie amândouă rejectate. Prin renunțarea la condiția ca protocolul să fie finit ca lungime (ceea ce problema generalilor demonstrează că nu este posibil) poate fi găsită o soluție practică la realizarea atomicității tranzațiilor distribuite.

În acest subcapitol va fi descris *protocolul comiterii tranzațiilor distribuite în două faze (two phase commit)*. Protocolul presupune existența unui coordonator al comiterii tranzației, care supervizează comunicația cu toți agenții de la distanță ai unei tranzații. În cadrul acestui protocol, fiecărui participant trebuie să i se permită să aborteze o tranzație în orice moment: acest lucru apare din necesitatea ca fiecare participant să-și poată controla accesul la obiectele dorite. De exemplu, multe mecanisme de control al concurenței (vezi subcapitolul 7.4.) blochează accesul la obiectele accesate de o tranzație activă. Dacă o tranzație blochează accesul la un obiect pentru o lungă perioadă de timp, un participant trebuie să poată aborta aceasta tranzație. Din acest punct de vedere însă, protocoalele de comitere a tranzațiilor distribuite au o fereastră de vulnerabilitate, în cadrul căreia participantii nu pot aborta unilateral o tranzație; pentru protocolul comiterii în două faze, aceasta fereastră în cursul căreia nu li se permite participanților să aborteze tranzație este mică. Protocolul se execută după ce toate actualizările asupra obiectelor s-au încheiat, la execuția comenzii *CommitTransaction*. În acest moment unul dintre participanți devine coordonatorul comiterii tranzației, în mod obișnuit cel din nodul de la care tranzația a fost inițiată sau cel dintr-un nod foarte sigur, de exemplu un server central.

Protocolul se desfășoară, la nivelul coordonatorului, astfel:

Coordonator Faza nr.1:

1. Se trimite un mesaj *PrepareToCommit* tuturor agenților participanți la tranzație și se așteaptă un mesaj de răspuns de la fiecare. Într-un fișier jurnal local, memorat pe un mediu fizic sigur (*stable storage*) se scrie un articol *PREPARED*.
- 2.a. Dacă un participant trimite ca mesaj de răspuns *Aborted* sau nu răspunde într-un interval de timp fixat, atunci tranzația este anulată: coordonatorul trimite un mesaj *AbortTransaction* la toți participanții. Concomitent coordonatorul urmărește desfășurarea tranzației înregistrând în fișierul jurnal local un articol *ABORT*, așteaptă mesajele de confirmare de la fiecare participant (*AckAbort*) și tranzația se termină.
- 2.b. Dacă toți participanții trimit ca mesaj de răspuns *Prepared*, se trece la faza nr.2.

Coordonator Faza nr.2:

1. Coordonatorul scrie un articol *COMMIT* în fișierul jurnal și trimite un mesaj *CommitTransaction* către toți participanții după care așteaptă un mesaj de răspuns *AckCommit*, de la fiecare participant.
2. La primirea tuturor mesajelor *AckCommit*, se scrie un articol *COMPLETE* în fișierul jurnal și tranzația se termină.

La nivelul participanților, protocolul se desfășoară astfel:

Participant Faza nr.1:

1. Se așteaptă primirea mesajului *PrepareToCommit* de la coordonator.
2. Se blochează obiectele locale implicate în acțiunea atomică, astfel încât să nu mai poată fi accesate de către alt cererător: totodată starea inițială a datelor se înregistrează pe un mediu fizic sigur (*stable storage*).
3. Dacă pasul 2 a fost realizat cu succes, se scrie un articol *PREPARED* într-un fișier jurnal local și se trimite un mesaj *Prepared* coordonatorului; în caz contrar se trimite un mesaj *Aborted* coordonatorului.

Participant Faza nr.2:

1. Se așteaptă un mesaj de la coordonator.
- 2.a. Dacă mesajul primit este *AbortTransaction*, se anulează efectul tranzației, se eliberează *lock*-urile și se trimite un mesaj *AckAbort* coordonatorului.
- 2.b. Dacă mesajul primit este *CommitTransaction*, se eliberează *lock*-urile, se permanentizează actualizările obiectelor, și se trimite un mesaj *AckCommit* coordonatorului și se înregistrează în jurnalul local un articol *COMMITTED*.

În eventualitatea avariei și apoi a restartării unui nod se realizează următoarele acțiuni: (1) Dacă coordonatorul execută *commit* și apoi nodul său suferă o avarie, înainte de a scrie articolul *COMMIT* în fișierul jurnal (pasul 1 al fazei nr.2) atunci trebuie trimis un mesaj *AbortTransaction* fiecărui participant. (2) Dacă coordonatorul scrie articolul *COMMIT* și nodul său suferă o avarie înainte de scrierea unui articol *COMPLETE* (pasul 2 al fazei nr.2), atunci trebuie trimis un mesaj *CommitTransaction* fiecărui participant. (3) Dacă coordonatorul suferă o avarie după scrierea articolului *COMPLETE*, atunci la restartare tranzația se consideră încheiată și poate fi ignorată.

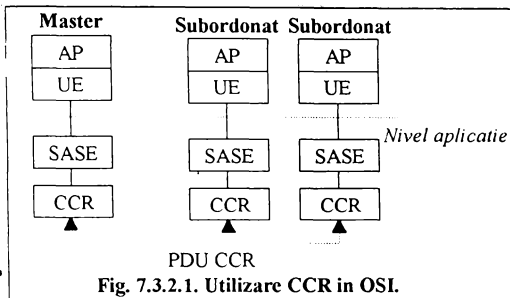
Protocolul se numește în *două faze* deoarece dacă unul dintre participanți nu poate comite acțiunea sa locală sau nodul său suferă o avarie în timpul fazei nr.1, atunci, la restartare, participantul respectiv (mai precis un *manager* de recuperare din acel nod) va întreba coordonatorul dacă trebuie să anuleze tranzația neterminată sau să o execute încă o dată.

Deoarece restartările pot avea loc de mai multe ori, operațiile de anulare a operațiilor locale și de refacere operațiilor locale (numite *undo* și respectiv *redo*) trebuie să fie *idempotent*. Dacă coordonatorul suferă o avarie înainte de execuția unui *commit* înainte de primirea unui mesaj *CommitTransaction*. Pe de altă parte, în cazul în care coordonatorul, la restartare, găsește articolul *COMMIT* în fișierul jurnal, atunci coordonatorul transmite *COMMIT* tuturor participanților. Efectul final este că fie toți participanții execută *commit*, fie niciunul.

Protocolul *two-phase commit* este adecvat în sistemele în care se dispune de transmisia broadcast a mesajelor și în care se cere un grad mare de paralelism. Numărul de mesaje pentru n participanți este $4n$; acest număr poate fi redus la $2n+2$ dacă sistemul dispune de transmisie broadcast sau la $2n$ dacă mesajele sunt transmise circular, de la un nod la altul; ultimul caz este adecvat pentru sisteme cu următoarele caracteristici: (1) nu este disponibilă transmisia broadcast și transmisia unui mesaj implică un cost mare; (2) necesitățile de concurență sunt reduse; (3) identitatea participanților la tranzație este de natura statică și cunoscută.

7.3.2 Protocolul CCR în OSI

Facilitățile de comitere a tranzațiilor, concurență și recuperare a tranzațiilor au fost standardizate de ISO în DIS 8649/3 și DIS 8650/3. Aceste standarde definesc serviciile, respectiv protocolul, pentru componentă OSI, *Commitment, Concurrency and Recovery (CCR)*; serviciile acestei componente CASE sunt disponibile aplicațiilor fie direct, fie prin intermediul altor componente SASE înrudite (fig.7.3.2.1).



CCR se bazează pe conceptul de tranzație atomică și pe protocolul *two-phase commit* mai sus descris (fig.7.3.2.2). În CCR, inițiatorul unei tranzații este denumit *master* iar agenții participanți *subordonați (slaves)*, înainte de inițierea unei tranzații trebuie să existe o asociere între aplicație (sau SASE) cu rolul de master și aplicațiile subordonate (fig.7.3.2.3.a). O tranzație atomică este inițiată de un master printr-o primitivă *C_Begin.request* având ca parametru un identificator al tranzației (numele aplicației urmat de un sufix unic) și timpul maxim de așteptare până la transmiterea unei cereri *RollBack* (fig.7.3.2.3). La primirea unei primitive *C_Begin.indication*, fiecare subordonat

crează o copie a obiectelor locale și startează mecanismul de blocare folosit; în practică primitivă *C_Begin* are efectul stabilirii la nivelul sesiunii a unor puncte de sincronizare pentru aceasta asociere. În continuare, se transmit și se execută operațiile specifice aplicației, utilizând primitivele SASE. În fig.7.3.2.3 se arată primitivele CCR pentru cele două cazuri posibile de terminare a unei tranzații: normal (b) sau anormal (c). În primul caz, (b), toți subordonații returnează succes la cererea de *commit*, utilizând serviciul *C_Ready*; în cel de-al doilea caz, (c), unul sau mai mulți subordonați returnează un răspuns negativ masterului, utilizând serviciul *C_Refuse*. În practică, serviciul *C_Prepare* este opțional și, dacă nu este utilizat, subordonații indică disponibilitatea lor printr-un *C_Ready* sau *C_Refuse* imediat după terminarea operațiilor aplicațiilor. În plus, există și serviciul *C_Restart* care poate fi folosit atât de un master cât și de un subordonat, pentru a aduce obiectele implicate într-o tranzație la o stare anterioară; poate fi folosit însă de un subordonat

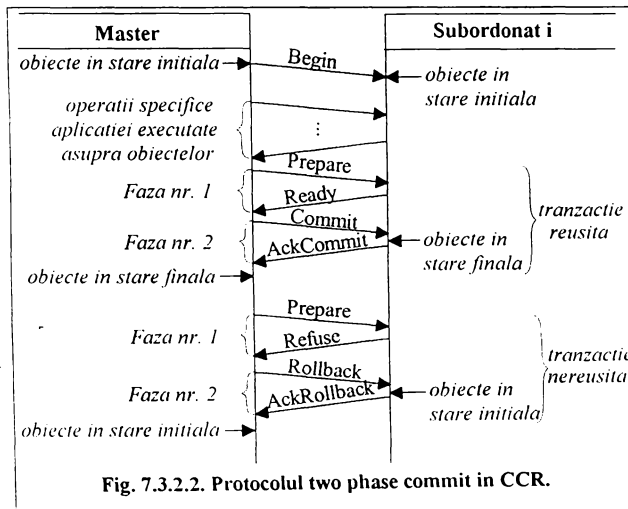
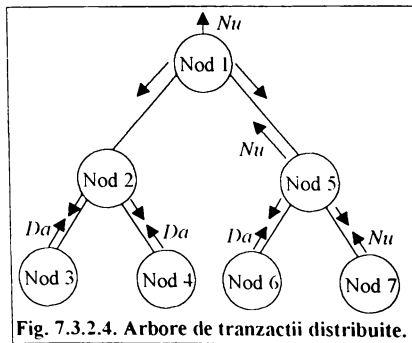
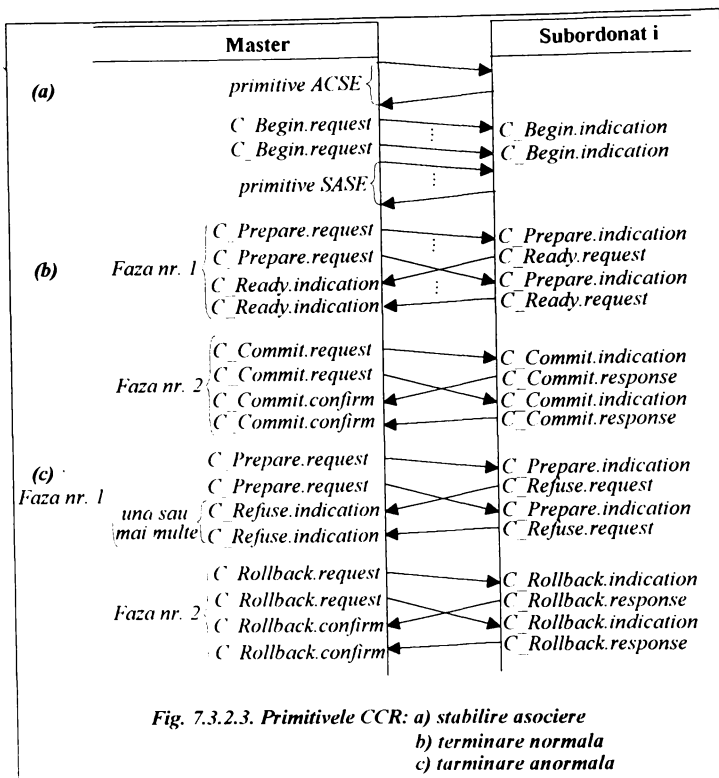


Fig. 7.3.2.2. Protocolul two phase commit in CCR.

numai dacă acesta nu și-a dat acordul pentru *commit*. Un master determină starea curentă a unui subordonat executând primitivă *C_Restart.request* cu un parametru care indică punctul de reluare astfel: (1) ACTION, semnifică restartarea acțiunii atomice la început; (2) COMMIT, semnifică restartarea acțiunii atomice după ultimul punct COMMIT; (3) ROLLBACK, semnifică restartarea acțiunii atomice după ultimul punct ROLLBACK. Similar, punctul de reluare dintr-o primitivă de confirmare se indică astfel: (1) DONE, semnifică faptul că subordonatul a realizat ultima cerere de *rollback* sau *commit*. (2) REFUSE, semnifică faptul că subordonatul nu a realizat ultima cerere și de aceea masterul trebuie să execute o cerere *rollback*; (3) ACTION, semnifică faptul că subordonatul trebuie să restarteze acțiunea de la început; (4) RETRY_LATER, semnifică faptul că subordonatul nu poate continua cu o restartare în acest moment și că masterul trebuie să încerce mai târziu.

Protocolul CCR poate genera un arbore de tranzații distribuite, în care un subordonat se comportă ca un master pentru o altă tranzație atomică descrisă prin subarboarele respective. Un parametru, identificatorul de ramură, se utilizează pentru a specifica o ramură din aceasta ierarhie. Dacă unul sau mai mulți subordonați refuză o propunere *commit*, refuzul se propagă pînă la masterul arborelui; acesta execută atunci o cerere *rollback*, care se propagă în jos la toți subordonații (fig.7.3.2.4).



7.3.3. Implementarea tranzacțiilor atomice

Cînd o tranzacție este în desfășurare, actualizările sale asupra obiectelor trebuie făcute într-o manieră reversibilă. Rezultă că participanții trebuie să mențină două stări ale obiectelor implicate într-o tranzacție; cea *inițială*, premergătoare tranzacției, care trebuie restaurată în cazul unui incident hardware sau a unei abortări a tranzacției și cea *intermediară*, inconsistentă, modificată pe măsură ce tranzacția evoluează. Cînd un participant recepționează din partea coordonatorului cererea de comitere a tranzacției, starea intermediară este considerată stabilă și înlocuiește starea inițială. Pe de altă parte, mecanismele de control al concurenței pot impune uneori lucrul cu o copie a obiectelor și nu cu obiectele propriu-zise (vezi subcapitolul 7.4.6.). Există două implementări mai cunoscute ale tranzacțiilor atomice care respectă aceste cerințe: prima implementare este cunoscută sub numele de *tehnica paginării umbră*; cea de-a doua implementare este cunoscută sub numele de *tehnica jurnalelor undo/redo*.

În cadrul primei tehnici, tranzacția nu operează direct asupra obiectelor, ci asupra unor copii ale obiectelor situate într-o zonă de lucru a tranzacției. Pînă în momentul comiterii tranzacției toate operațiile asupra obiectelor se execută asupra copiilor obiectelor din zona de lucru și nu asupra obiectelor propriu-zise. Dacă, de exemplu, obiectele

sunt fișiere, atunci, înainte de a actualiza un bloc al unui fișier, se obține un bloc liber din lista de blocuri libere, în care se copiază vechiul bloc al fișierului și se actualizează această copie; concomitent, se copiază în zona de lucru și se actualizează și harta blocurilor alocate fișierului. Astfel, în timp ce tranzacția respectivă lucrează cu o noua hartă de alocare a fișierelor, celelalte aplicații percep fișierul ca neschimbat. Noile blocuri alocate fișierului se numesc *pagini umbra*. Fig.7.3.3.1 ilustrează această tehnică. Se observă că se utilizează și un fișier jurnal, *jurnal de intenții*, care înregistrează toate blocurile umbră introduse de o tranzacție. La comiterea tranzacției, harta fișierului se actualizează, iar spațiul ocupat de vechile blocuri ale fișierului se eliberează. Dacă are loc un incident hardware în timpul comiterii tranzacției, atunci, la repunerea în funcțiune, având înregistrat jurnalul de intenții (pe un mediu fizic sigur-*stable storage*) orice participant va actualiza hartile de fișier

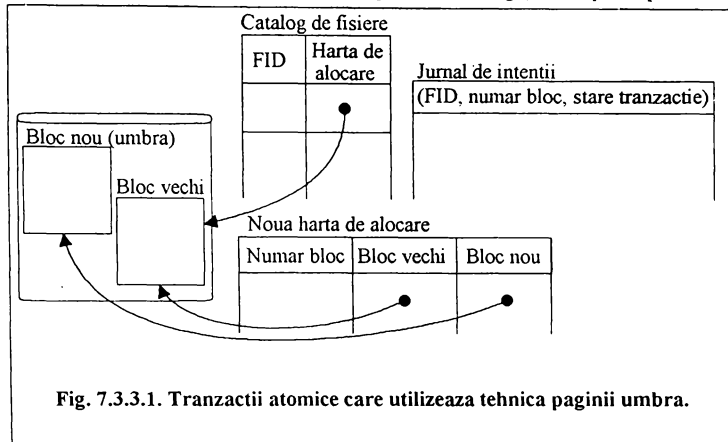


Fig. 7.3.3.1. Tranzacții atomice care utilizează tehnica paginii umbra.

conform jurnalului de intenții.

O altă tehnică de implementare a tranzacțiilor atomice este *tehnica jurnalelor undo/redo* numită uneori *tehnica listelor de intenții*. Aceasta elimină redundanța tehnicii paginii-umbra plecând de la două premise: (1) în general, majoritatea tranzacțiilor nu esuează; (2) regia de sistem sporită în cazul tehnicii paginii-umbra se datorează copiilor de blocuri de disc cu care se lucrează, ocolind blocurile alocate. Tehnica jurnalelor undo/redo permite efectuarea actualizărilor direct asupra obiectelor; dar, înainte de modificarea unui obiect, valoarea veche și noua a acestuia se memorează într-un fișier jurnal; o intrare în acest fișier va conține deci identificatorul tranzacției, identificatorul obiectului, vechea și noua valoare. Articolele fișierului jurnal pot fi repertorizate într-un fișier hartă a jurnalului (fig.7.3.3.2).

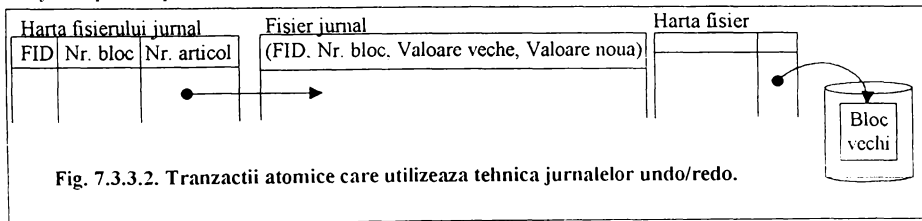


Fig. 7.3.3.2. Tranzacții atomice care utilizează tehnica jurnalelor undo/redo.

În afară de regia de sistem mai redusă, tehnica jurnalelor undo/redo mai are și următoarele avantaje față de tehnica paginii-umbra: (1) menține alocarea inițială a blocurilor unui fișier, păstrând eventual contiguitatea blocurilor, chiar și în cazul unei actualizări într-o ordine aleatoare a blocurilor, ceea ce nu este adevărat în tehnica paginii-umbra; (2) poate fi construită ca un nivel deasupra sistemului de gestiune al fișierelor, deoarece nu implică acces la blocuri și la hartile de alocare ale fișierelor.

7.3.4. Recuperarea obiectelor

Recuperarea obiectelor se definește ca fiind tehnica prin care se restaurează o stare consistentă anterioară a unui sistem de obiecte, în cazul unei defecțiuni hardware sau a unei erori software.

Soluția tradițională în bazele de date centralizate este crearea periodică a unor puncte de reluare și recuperarea prin revenirea la cel mai recent punct de reluare. Într-un sistem distribuit însă această soluție se complică datorită cerinței de respectare a consistenței datelor în sistemul distribuit și necesității de recuperare în urma oricăror combinații de erori ale liniilor de comunicații, ale nodurilor, erori software, incluzând chiar erorile din timpul realizării recuperării.

Se lucrează de aceea cu un *sistem de recuperare al obiectului*, bazat pe un *fișier jurnal (log)*; acesta conține o istorie a tuturor acțiunilor recuperabile (al căror efect se poate anula) și implementează *abortarea unei tranzacții*. Recuperarea obiectelor este strâns legată de modul de implementare al tranzacțiilor prezentat în subcapitolul anterior.

7.3.4.1. Protocolul Write-Ahead Logging

Un fișier *log* este un fișier secvențial, potențial nelimitat ca lungime, care poate fi modificat numai prin adăugarea de articole. În practică, articolele care trebuie scrise în *log* sunt în mod obișnuit bufferate în memoria principală, în timp ce pe disc se rezervă o dimensiune fixă. Când dimensiunea devine prea mare pentru ca articolele să poată începe în această zonă, fie se salvează o parte necritică a conținutului într-o zonă off-line, fie se trunchiază periodic fișierul astfel încât să nu se piardă informația critică. Scrierea într-un fișier *log* se face într-o manieră atomică, utilizând tehnica *stable storage*. În mod obișnuit, tehnica *stable storage*, fie că utilizează instrucțiuni de I/O la nivel fizic (operațiile de bază furnizate de hardware), fie că utilizează apeluri ale sistemului de operare pentru scriere (dar în sistemul de operare UNIX, de exemplu, acestea trebuie urmate și de apeluri *fsync()* pentru a asigura scrierea efectivă pe disc) realizează două copii ale datelor. Ca atare, dacă în timpul unei operații de scriere apare un incident hardware, cel puțin una dintre copii rămâne validă, și astfel și cealaltă poate fi corectată (sau creată din nou).

În general, scrierea articolelor în fișierele *log* se face urmărindu-se două scopuri: (1) pentru a anula efectul unor tranzații abortate; (2) pentru a restabili efectul unor tranzații care au executat *commit*, dar ale căror actualizări au fost întrerupte sau distruse de un incident hardware.

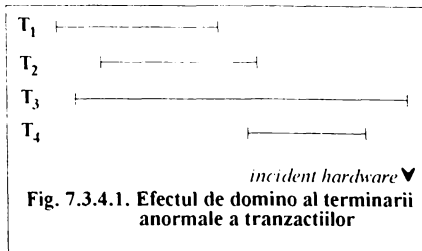
Protocolul *Write-Ahead Logging* (WAL) este un protocol care definește momentele în care articolele *log*, bufferate în memorie, trebuie scrise pe disc, cu scopul de a asigura atomicitatea tranzațiilor. Informațiile pentru anularea tranzațiilor (*undo* information) trebuie scrise pe disc înainte de modificarea versiunii persistente a unui obiect printr-o operație a unei tranzații. Informațiile pentru restabilirea efectului unei tranzații (*redo* information) trebuie scrise pe disc înainte ca tranzația să execute *commit*. Acest protocol oferă avantaje unanim recunoscute față de alte tehnici pentru recuperare. Când o tranzație accesează date ale unui obiect (definit de utilizator), se crează în memorie o copie a acelui obiect; dacă obiectul este modificat, copia respectivă se rescrie pe disc (versiunea persistentă a obiectului). Deoarece informațiile necesare pentru anularea actualizărilor tranzațiilor sunt scrise în *log*, versiunea persistentă a obiectului poate fi oricând înlocuită pe disc cu o copie actualizată a obiectului, chiar dacă aceasta reflectă actualizări ale unei tranzații încă necomise. Pentru alte mecanisme (ca de exemplu tehnica paginii umbră, prezentată în subcapitolul precedent) alocatorul de memorie trebuie să rezerve un nou spațiu pe disc pentru noua versiune actualizată a obiectului; versiunea nemodificată a obiectului trebuie reținută pentru ca tranzația să poată fi abortată. Astfel, deoarece obiectele actualizate își mută poziția pe disc, strategiile pentru plasarea grupată a obiectelor înrudite (*clustering strategies*) sunt dificil de implementat (dezavantajul (1) al tehnicii paginii-umbră din subcapitolul precedent). Analog, deoarece fișierul *log* conține și informații pentru restabilirea efectului unei tranzații, obiectele actualizate nu trebuie obligatoriu scrise pe disc înainte de comiterea unei tranzații. Obiectele frecvent actualizate pot rămâne în memoria principală atâta timp cât este necesar, fără a fi rescrise pe disc (obiectul din memorie joacă astfel rolul de cache).

În plus, în afară de articolele utilizate pentru informațiile *redo* și *undo*, în fișierele *log* se mai înregistrează și alte tipuri de informații: articole *checkpoint*, care conțin informații pentru a limita căutarea în fișierul *log* în cursul unei recuperări, articole care descriu modificarea stării tranzațiilor în decursul unui protocol *two-phase commit* și, eventual, informații de stare necesare mecanismului de control al concurenței care trebuie reținute și după incidentele hardware.

7.3.4.2. Tehnicile Value-Logging și Operation Logging

Tehnica *Value-Logging* este cea mai simplă tehnică pentru recuperarea obiectelor, utilizând protocolul WAL, când se actualizează un obiect, sistemul de recuperare scrie în fișierul *log* valoarea obiectului neactualizat și valoarea obiectului actualizat; trebuie făcută observația că se poate scrie în fișierul *log* și numai un cîmp al valorii obiectului (cîmpul care se actualizează), și nu toate cîmpurile, pentru economie de spațiu. Dacă tranzația va executa *commit*, și apare un incident hardware, sistemul de recuperare trebuie să restaureze valoarea obiectului la cea de dinaintea operațiilor de actualizare, pe care o găsește în fișierul *log*. Dacă tranzația execută *commit*, și apare un incident hardware care împiedică actualizarea versiunii persistente a obiectului, sistemul de recuperare trebuie să restaureze valoarea obiectului la cea finală de după toate operațiile de actualizare, pe care o găsește în fișierul *log*. Este evident că operațiile făcute de sistemul de recuperare a obiectelor, în cazul acestei tehnici, pentru a anula și reinstala efectul unei tranzații sunt *idempotente*; execuția lor de mai multe ori, în cazul unor incidente în cursul procesului de recuperare are același efect.

Mecanismul de control al concurenței, care coordonează accesul la obiecte a mai multor tranzații, trebuie să se asigure că dacă un obiect este actualizat de o tranzație, el nu va mai fi preluat de nici o altă tranzație pînă la comiterea primei tranzații. În caz contrar, dacă prima tranzație ar fi abortată după cel puțin o actualizare și deci efectul ei anulat, atunci ar trebui abortate toate tranzațiile care au accesat acel obiect, chiar și cele care între timp au fost comise. Figura 7.3.4.1. ilustrează această situație. Înainte de incidentul hardware, T_1 actualizează obiectul O_1 și deblochează accesul la acest obiect; în continuare T_4 citește obiectul O_1 și actualizează obiectul O_2 , apoi se termină normal cu *commit*. După repararea incidentului, T_1 poate fi abortată și apoi restartată, dar T_4 este deja terminată și nu poate fi abortată, chiar dacă a realizat, în mod eronat, actualizarea obiectului O_2 în funcție de valoarea lui O_1 . Acest efect se numește *efect domino* sau *abort în cascada al tranzațiilor*.



În multe situații, totuși, tehnica value-logging nu poate fi folosită datorită limitării severe a performanțelor. De exemplu, dacă se înregistrează mai multe articole ale unui fișier într-o pagină și se consideră o pagină ca un obiect, atunci când se actualizează un obiect din aceea pagină de către o tranzacție, nici o altă tranzacție nu va mai putea utiliza nu numai acel articol, ci orice alt articol din aceea pagină, pînă la comiterea primei tranzacții; aceasta restricționează într-un mod acceptabil concurența.

Tehnica Operation Logging poate fi utilizată pentru a diminua restricțiile cu privire la actualizările efectuate de tranzacții neterminate. În cadrul acestei tehnici, sistemul de recuperare scrie în fișierul log, o descriere a operației care a provocat actualizarea unui obiect și nu valoarea obiectului de dinainte și de după actualizare. Afița timp cît pentru fiecare operație există și o operație compensatorie, care să anuleze efectul operației, o tranzacție poate fi întotdeauna abortată. Operația poate fi și reexecută, cînd un incident hardware apare după *commit*. Datorită faptului că reprezentarea unei operații în fișierul log se face mult mai compact decît reprezentarea valorii obiectelor, tehnica Operation logging prezintă avantajul economiei de spațiu în log și al unui timp mai mic de scriere în log.

Se dau următoarele exemple de operații care sunt adecvate pentru aceasta tehnică: operații de incrementare și decrementare a unui contor, de inscriere și ștergere a unui nod într-un arbore B care poate memora, de exemplu, maparea numelor de fișiere din cataloage în identificatori de fișiere sau relațiile de indexare dintr-o bază de date.

Deși tehnica operation logging oferă anumite avantaje, totuși introduce un nou gen de probleme. Multe din operațiile care actualizează obiectele *nu sunt idempotente*; dacă în cursul unui proces de recuperare se repetă aceste operații, rezultatele obținute vor fi eronate. Cea mai cunoscută soluție la acest tip de problemă este de a asocia fiecărui obiect un număr de *secvență log*, cu valori strict crescătoare, care va reflecta starea obiectului scris pe disc. Înainte de a repeta o actualizare, se compară acest număr din log cu numărul asociat obiectului; dacă numărul asociat obiectului este mai mare sau egal cu cel din log, actualizarea descrisă în log a fost deja aplicată și nu mai trebuie reinstalată.

7.3.4.3. Tehnicile Redo-Only Logging și Undo-Only Logging

În unele situații poate fi avantajos să se întîrzie actualizările unei tranzacții pînă în momentul comiterii ei; acest lucru asigură faptul că operațiile tranzacției nu vor trebui să fie niciodată anulate și de aceea vechile valori ale obiectelor nu se salvează în log. Aceasta tehnică, numită *tehnica redo-only logging* este avantajoasă, de exemplu, în cazul unei operații de ștergere a conținutului unui fișier sau a unui articol al unui fișier. Operațiile de ștergere se reexecută simplu, dar efectul lor se anulează greu; informațiile de *undo* pentru aceste operații sunt de dimensiuni mari. Dacă, în schimb, actualizările sunt întîrziate pînă în momentul comiterii tranzacției, astfel încît informațiile de *undo* nu trebuie scrise în log, sistemul trebuie să garanteze că, odată ce tranzacția execută *commit*, orice actualizări întîrziate sunt realizate, chiar dacă intervine un incident hardware. Tehnica generală este de a crea o *lista de intenții* care să cuprindă actualizările în așteptare și de a forța scrierea acestor liste în log, înainte de faza *prepare to commit*. Există sisteme de fișiere a căror tehnică de recuperare se bazează în întregime pe tehnica redo-only logging: sistemele de fișiere Alpine ([BKT85]) și Oregon ([Hag87]) interzic rescrierea obiectelor pe disc pînă în momentul în care tranzacția execută *commit*.

O altă tehnică utilă este de a scrie în fișierul log numai informațiile *undo* și de a forța scrierea pe disc a obiectelor actualizate înainte de *commit*. Aceasta tehnică, numită *tehnica undo-only logging* asigură faptul că actualizările nu vor trebui niciodată reexecutate și, de aceea, noile valori ale obiectelor nu se scriu în log. Tehnica este utilă pentru cazul în care se crează fișiere noi, într-un sistem de fișiere.

Trebuie specificat, în final, că, în general, un sistem de recuperare poate folosi o combinație de tehnici de recuperare prezentate în subcapitolele 7.3.4.2. și 7.3.4.3., deci o tehnică hibridă. În acest sens exemple sugestive sunt sistemul de fișiere QuickSilver Distributed File System (QDFS) și sistemul de gestiune de baze de date Starbust ([CPS93]). Astfel, aceste sisteme folosesc atît tehnica *value-logging* (în Starbust, pentru înregistrarea conținutului bazei de date, în QDFS pentru înregistrarea hărții de alocare a fișierelor), cît și tehnicile *operation-logging* (pentru inserarea sau ștergerea cheilor într-un arbore B), *redo-only logging* (pentru ștergerea fișierelor în DFS și a relațiilor și fișierelor de relații în Starbust) și *undo-only logging* (pentru crearea fișierelor noi care conțin articole ale bazei de date).

7.3.4.4. Structura sistemului de recuperare

Un sistem de recuperare al obiectelor în cadrul modelului tranzacțional prezentat în subcapitolul 7.2.2. are structura generală prezentată în fig. 7.3.4.4.

El este compus din două elemente principale:

1. Un *manager de recuperare* al cărui rol este de a urmări evoluția tranzacțiilor, și de a lua deciziile adecvate conform tehnicii de recuperare alese, în timpul execuției acțiunilor tranzacției, ca și la execuția comenzilor *commit*, *abort*.
2. Un *manager pentru log*, utilizat de managerul de recuperare și de alte componente ale sistemului distribuit pentru a scrie articole în log; scrierea articolelor în fișierul log este cerută fie de managerul de recuperare, fie de sistem.

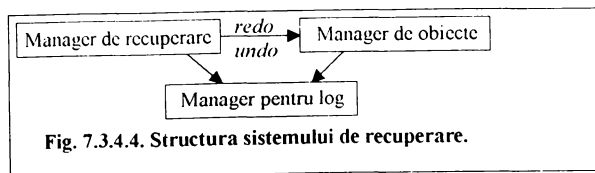


Fig. 7.3.4.4. Structura sistemului de recuperare.

7.4. Controlul concurenței

În scopul menținerii consistenței datelor distribuite, se impun anumite constrângeri în structura tranzacțiilor sau în maniera de realizare a concurenței.

7.4.1. Cerințele mecanismului de control al concurenței

Kohler ([Koh81]) a arătat că, dacă într-un sistem distribuit nu apare o defecțiune hardware, *toate mecanismele de control al concurenței trebuie să asigure următoarele:*

- (1) consistența obiectelor;
- (2) realizarea fiecărei tranzacții atomice într-un timp finit.

În plus, un *bun* mecanism de control al concurenței se caracterizează astfel:

- (1) permite execuția paralela a acțiunilor, în scopul satisfacerii cerințelor de performanță ale sistemelor;
- (2) implică prelucrări minore;
- (3) se execută cu performanțe satisfăcătoare chiar și în rețele în care întârzierile de transmisie sunt semnificative;
- (4) impune puține constrângeri în structura tranzacțiilor.

7.4.2 Serializarea tranzacțiilor

În continuarea acestui capitol se va folosi următoarea definiție formală pentru o tranzacție.

Definiție:

O *tranzacție* T_i este o pereche $(T_i, <_i)$, în care $T_i = \{T_{ij}; 1 \leq j \leq m\}$ este setul de acțiuni aparținând tranzacției i iar $<_i$ este o relație de precedență care ordonează setul T_i .

Pentru a se preciza acțiunea ce se execută asupra sistemului de obiecte și obiectul asupra căruia se execută (îr modelul prezentat la 7.2.2) se vor utiliza următoarele notații:

citire: $\langle T_i, \text{read_object}_i(X, S_k), V \rangle$

scriere: $\langle T_i, \text{write_object}_i(X, S_k), V \rangle$

în care prin T_i se precizează tranzacția, prin *read_obiect_i*, respectiv *write_obiect_i*, se precizează acțiunea, prin X obiectul, prin S_k nodul pe care se afla obiectul iar prin V valoarea rezultată prin citire, respectiv implicată în scriere.

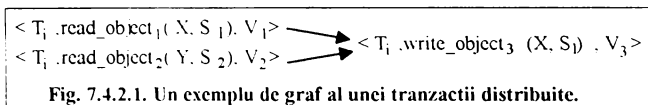


Fig. 7.4.2.1. Un exemplu de graf al unei tranzacții distribuite.

O tranzacție poate fi reprezentată printr-un graf $G = (N, A)$ în care N este un set de noduri corespunzând acțiunilor tranzacției T_i iar $A = \{(T_i, T_{ik}); T_{in} <_i T_{ik}\}$ este un set de arce indicind ordinea de execuție a acțiunilor. În figura 7.4.2.1 este prezentat un exemplu de graf

al unei tranzacții simple. Acest graf reprezintă tranzacția T_i compusă din acțiuni asupra a trei entități (sau obiecte, deoarece entitățile nu sunt replicate). Acțiunea *write* trebuie executată după terminarea ambelor acțiuni de citire; nu există restricții asupra ordinii de execuție a celor două citiri.

Pot fi utilizate două modele pentru tranzacțiile distribuite: *modelul tranzacțiilor seriale* și *modelul tranzacțiilor concurențe*. În modelul tranzacției seriale se presupune că relația de precedență $<_i$ ordonează total seturile T_i ; în modelul tranzacției concurențe, se presupune că relația $<_i$ ordonează parțial seturile T_i . În figurile 7.4.2.2 și 7.4.2.3 se prezintă două exemple de grafuri pentru cele două modele.

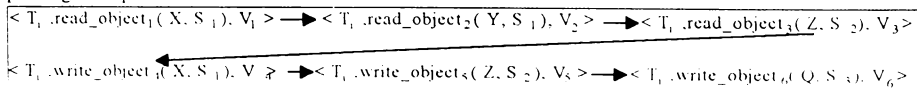


Fig. 7.4.2.2. Un exemplu de graf al unei tranzacții distribuite seriale.

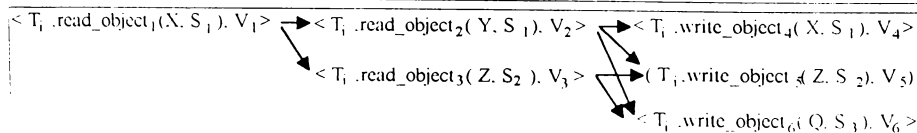


Fig. 7.4.2.3. Un exemplu de graf al unei tranzacții distribuite concurențe.

S-a arătat în subcapitolul 7.2 că principala cauză pentru situațiile de violare a consistenței sistemelor la obiecte este execuția concurență a unui set de acțiuni provenind de la tranzacții diferite. Controlul concurenței trebuie să acționeze prin intermedierea corectă a acestor acțiuni; controlul concurenței va consta în planificarea acțiunilor tranzacțiilor. Altfel spus, ordonarea acțiunilor unui set de tranzacții este rezultatul aplicării *algoritmului de control al concurenței* setului de tranzacții.

Definiție.

O *planificare a unui set de tranzacții* τ este o mulțime parțial ordonată $(T(\tau), <_s)$ în care $T(\tau) = U_i T_i \cup T_0 \cup T_r$ este mulțimea tuturor acțiunilor executate de tranzacții, completată cu încă două tranzacții, o *tranzacție inițială* T_0 care scrie (stabilește) starea inițială a sistemului distribuit de obiecte și o *tranzacție finală* T_r care citește starea finală iar $<_s$ este o relație de ordine parțială peste $T(\tau)$ pentru care sunt îndeplinite următoarele condiții:

$$(1) <_s \ U_i <_r$$

(2) pentru orice pereche de acțiuni T_{ij} din $T(\tau)$, T_{ki} din $T(\tau)$, cu $i < k$, care cer acces la același obiect, dintre care cel puțin una este de scriere, este adevărată relația:

$$T_{ij} <_s T_{ki} \text{ sau } T_{ij} <_s T_{ki}$$

În contrast cu sistemele centralizate, într-un sistem distribuit este imposibil să se determine la un moment dat starea globală a întregului sistem sau să se determine aprioric ordinea de execuție a acțiunilor în prezența întârzierilor de transmisie. De aceea, relația $<_s$, care desemnează ordinea de execuție a acțiunilor într-o planificare $S(\tau)$ este o relație de ordine parțială. Condiția (1) impune ca relația $<_s$ să respecte ordonarea acțiunilor stipulată în cadrul tranzacțiilor. Condiția (2) impune ca acțiunile *read* și *write* sau *write* și *write* asupra aceluiași element al unui obiect să fie executate serial.

O planificare $S(\tau)$ poate fi reprezentată printr-un *graf al planificării*.

Definiție.

Un *graf al unei planificări* $SG(S(\tau))$ este un graf orientat (N, A) ale cărui noduri corespund tuturor acțiunilor $T(\tau)$ și ale cărui arce sunt toate perechile (T_{ij}, T_{ki}) astfel încât $T_{ij} <_s T_{ki}$.

Planificările pot fi *seriale* sau *concurente*. O planificare $S(\tau)$ este *serială* dacă pentru orice două tranzacții, toate acțiunile unei tranzacții se execută înaintea tuturor acțiunilor celei de-a doua tranzacții; în caz contrar, cînd acțiunile unor tranzacții diferite se execută în paralel (sau intercalat), planificarea este *concurentă*. Cu alte cuvinte, o planificare conținînd tranzacțiile T_1, T_2, \dots, T_m este serială (și consistentă) dacă pentru orice $i=1, \dots, m-1$ tranzacția T_i se execută complet înainte ca tranzacția T_{i+1} să-și înceapă execuția.

O planificare a unui set de tranzacții τ executate în N noduri într-un sistem distribuit poate fi reprezentată printr-un set de planificări locale $S(\tau) = \{S_1, S_2, \dots, S_k\}$ care descriu execuția tranzacțiilor la sistemele locale de obiecte.

Formal, o planificare locală S_k a unui set de tranzacții τ în nodul S_k este o mulțime parțial ordonată $(T^k(\tau), <_k)$ în care $T^k(\tau) = \{T_{ij} : T_{ij} \text{ se execută în nodul } S_k\} \in T(\tau)$. Se consideră în continuare problema corectitudinii planificarilor tranzacțiilor în contextul unor *informații sintactice* asupra tranzacțiilor; aceasta înseamnă că algoritmul de control al concurenței are la dispoziție numai informații despre identitatea elementelor obiectelor la care se cere acces (modelul *controlului sintactic al concurenței*). Un criteriu firesc de corectitudine pentru planificările concurente este *criteriul de serializare al tranzacțiilor*. Înainte de a enunța acest criteriu se mai introduc următoarele definiții:

Definiții.

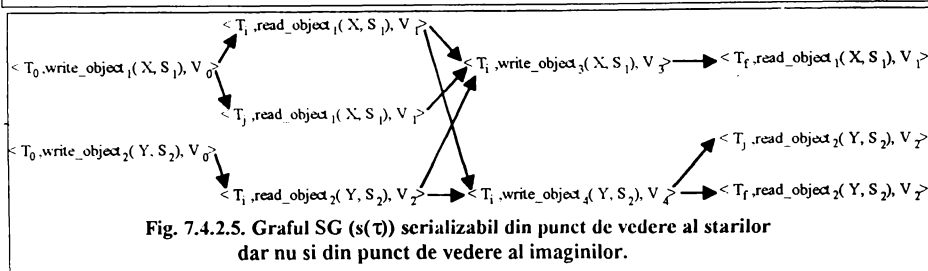
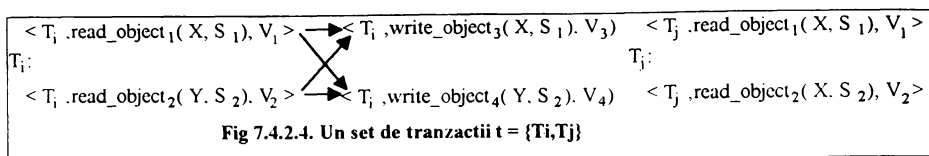
- (1) Se numește *stare* a unui sistem de obiecte mulțimea de valori asignate elementelor tuturor obiectelor din sistemul distribuit de obiecte.
- (2) Se numește *imaginea* sistemului de obiecte percepută de o tranzacție T_i într-o planificare $S(\tau)$ setul de valori ale datelor citite de acțiunile tranzacției.
- (3) Două planificări $S(\tau) = (T(\tau), <_s)$ și $S'(\tau) = (T(\tau), <_{s'})$ sunt *echivalente din punct de vedere al imaginilor* dacă imaginea percepută de orice tranzacție $T_i \in S(\tau)$ este identică cu imaginea percepută de aceeași tranzacție în planificarea $S'(\tau)$.
- (4) Două planificări $S(\tau) = (T(\tau), <_s)$ și $S'(\tau) = (T(\tau), <_{s'})$ sunt *echivalente din punct de vedere al stărilor*, dacă pentru orice prelucrări ale tranzacțiilor conținute în τ , stările finale obținute de S și S' sunt identice.
- (5) Două planificări $S(\tau) = (T(\tau), <_s)$ și $S'(\tau) = (T(\tau), <_{s'})$ sunt *echivalente* dacă ele sunt echivalente ca imagine și ca stări.

Prin urmare, două planificări S_1 și S_2 sunt echivalente dacă: (1) setul de tranzacții din S_1 și S_2 este identic; (2) pentru orice obiect X din S_1 , dacă tranzacția T_1 execută o citire a obiectului X și valoarea lui X citită de T_1 este modificată de T_2 , aceeași situație trebuie să apară și în S_2 (sincronizare *read-write*); (3) Pentru fiecare obiect X din S_1 , dacă tranzacția T_1 execută o scriere a obiectului X , înainte ca T_2 să execute și ea o scriere a lui X , aceeași situație trebuie să apară și în S_2 (sincronizare *write-write*).

Criteriu de serializare

O planificare $S(\tau)$ este *corectă* dacă este echivalentă cu o planificare serială a setului de tranzacții.

O planificare echivalentă cu o planificare serială se numește *planificare serializabilă*; mulțimea tuturor planificarilor serializabile ale unui set de tranzacții τ se va nota cu SR . În figura 7.4.2.5. se arată un graf al



tranzacțiilor din figura 7.4.2.4. serializabil din punct de vedere al starilor, dar neserializabil din punct de vedere al imaginilor deoarece tranzacția T_i citește valoarea obiectului X , care este inconsistentă.

Pentru a verifica serializabilitatea unei planificări concurente, $S(\tau)$, trebuie examinată ordinea acceselor la date al acțiunilor tranzacțiilor. Verificarea se face utilizând un *graf de serializare*.

Definiție:

Graful de serializare al planificării $S(\tau)$, notat $SRG(S(\tau))=(N,A)$ este un graf orientat, a cărui mulțime de noduri N corespunde setului de tranzacții extins cu două tranzacții ipotetice T_o și T_p , care stabilesc starea inițială, respectiv starea finală și a cărui mulțime de arce se definește astfel: dacă există un obiect X și o acțiune $write_object_i(X)$ T_i și o acțiune $read_object_k(X)$ T_k astfel încât a doua acțiune citește valoarea scrisă de prima, atunci:

- (1) $(T_i, T_k) \in A$
- (2) dacă $T_i \neq T_o$, $T_k \neq T_p$ și există o tranzacție $T_p \neq T_o$, care scrie pe X , atunci fie $(T_p, T_k) \in A$ fie $(T_k, T_p) \in A$.
- (3) dacă $T_i \neq T_o$, atunci $(T_o, T_i) \in A$.
- (4) dacă $T_i = T_o$, $T_k \neq T_p$ și există o tranzacție $T_p \neq T_o$, care scrie pe X , atunci $(T_k, T_p) \in A$.
- (5) dacă $T_k = T_p$ și există o tranzacție T_p , care scrie pe X atunci $(T_p, T_i) \in A$.

Prin condiția (1) se impune ca, dacă într-o planificare $S()$ o tranzacție T_k citește valoarea lui X scrisă de o tranzacție T_i , atunci T_i trebuie să precedă T_k în orice planificare serială echivalentă; precedența este reprezentată prin arcul (T_i, T_k) .

Condiția (2) impune ca, dacă într-o planificare $S(\tau)$, o tranzacție T_k citește o valoare scrisă de T_i și o altă tranzacție T_p scrie X , atunci T_p trebuie fie să precedă pe T_i , fie să urmeze lui T_i în orice planificare serială echivalentă. Condițiile (3),(4),(5) se referă la tranzacția inițială și la cea finală. Dacă T_i este tranzacția inițială, atunci nici o tranzacție T_p nu o poate precede: de asemenea, dacă T_k este tranzacția finală, nici o tranzacție T_p nu poate să o urmeze.

Teorema:

- (1) planificare $S(\tau)$ este *serializabila* dacă și numai dacă poate fi construit un graf de serializare $SRG(S(\tau))$ aciclic.

Deoarece verificarea existenței unui astfel de graf implică construcția tuturor grafurilor de serializare, o problemă NP completă în practică, soluțiile propuse pentru controlul concurenței utilizează un alt criteriu de corectitudine, mai restrictiv, cel al *serializabilității D*. Pentru a defini serializabilitatea D , trebuie definite în prealabil următoarele elemente:

Definiție:

- (1) Două acțiuni T_u și T_v sunt în conflict în accesul la obiectul X , conflict desemnat de $T_u \rightarrow^c_X T_v$, dacă ambele accesează același obiect și cel puțin o acțiune este de scriere.
- (2) Două tranzacții sunt în conflict de acces la date, conflict specificat prin $T_i \rightarrow^c T_j$ dacă există un obiect X și acțiunile $T_u \in T_i$ și $T_v \in T_j$ astfel încât $T_u \rightarrow^c_X T_v$.
- (3) Acțiunea T_u preceda acțiunea T_v în contextul planificării $S(\tau)$ în raport cu obiectul X , $T_u \rightarrow^c_X T_v$, dacă $T_u \rightarrow^c T_j$ și $T_j \rightarrow T_v$.
- (4) Tranzacția T_i precedă tranzacția T_k în planificarea $S(\tau)$, $T_i \rightarrow T_k$, dacă există un obiect X și o acțiune $T_u \in T_i$ și $T_v \in T_k$ astfel încât $T_u \rightarrow^c_X T_v$.

În aceste condiții, relația de precedență \rightarrow ordonează parțial setul de tranzacții.

Teorema:

- (1) planificare $S(\tau)$ este *serializabila D* dacă și numai dacă relația sa de precedență este aciclică.

Verificarea corectitudinii unei planificări $S(\tau)$ se reduce la a testa dacă graful relației de precedență \rightarrow este aciclic. Acest graf, denumit *graful de serializare D al unei planificări $S(\tau)$* , $DSRG(S(\tau))$ se definește astfel:

Definiție:

Graful de serializare D al unei planificări $S(\tau)$, $DSRG(S(\tau))$, este un graf orientat $DSRG(S(\tau))=(V,A)$ a cărei mulțime de noduri corespunde setului de tranzații extinse cu tranzația inițială și cu una finală, T_0 și T_f și a cărei mulțime de arce este $A=\{(T_i, T_j) : T_i \rightarrow T_j, T_i, T_j \in \tau \cup \{T_0, T_f\}\}$.

Ordinea tranzațiilor dintr-o planificare, implicată de relația \rightarrow se numește *ordine de serializare*; semnificația ei este următoarea: dacă tranzațiile dintr-un set ar fi executate serial, în ordinea de serializare, atunci planificarea obținută în acest mod ar fi echivalentă cu planificarea $S(\tau)$. În fig. 7.4.2.8. se arată un exemplu de graf de serializare D, pentru setul de tranzații din fig. 7.4.2.6., cu graful planificării în fig. 7.4.2.7. Ordinea de serializare a tranzațiilor este T_0, T_2, T_1, T_3, T_f . Setul tuturor planificărilor unui set de tranzații care satisface criteriul de serializare D, DSR, este inclus în SR.

Criteriul de serializare D este mai restrictiv decât criteriul de serializare prezentat anterior. Exemplul din figurile 7.4.2.9, 7.4.2.10, 7.4.2.11. demonstrează acest lucru. Se consideră planificarea din fig. 7.4.2.9. Graful de serializare construit pentru această planificare în figura 7.4.2.10. este aciclic, prin urmare planificarea este serializabilă. Setul de arce se determină astfel: arcul (T_1, T_2) trebuie să existe în $SRG(S(\tau))$, deoarece $\langle T_2, read_obiect_1(X, S_1), V_1 \rangle$ citește valoarea lui X scrisa de $\langle T_1, write_obiect_1(X, S_1), V_1 \rangle$; arcul (T_2, T_1) trebuie să existe în $SRG(S(\tau))$, deoarece $\langle T_1, read_obiect_1(X, S_1), V_1 \rangle$ citește valoarea lui X scrisa de $\langle T_2, write_obiect_1(X, S_1), V_2 \rangle$. Datorită componenței planificării $S(\tau)$ trebuie inclus unul din următoarele două arce: (T_3, T_1) sau (T_3, T_2) ; fie (T_3, T_1) arcul inclus. În sfârșit, conform condițiilor (3), (4), (5) din definiția grafului de serializare mai trebuie incluse în $SRG(S(\tau))$ și următoarele arce: (T_0, T_1) , (T_0, T_2) , (T_3, T_2) . Planificarea $S(\tau)$ este serializabilă, dar nu și serializabilă D, deoarece graful $DSRG(S(\tau))$ (fig. 7.4.2.11.) conține un ciclu.

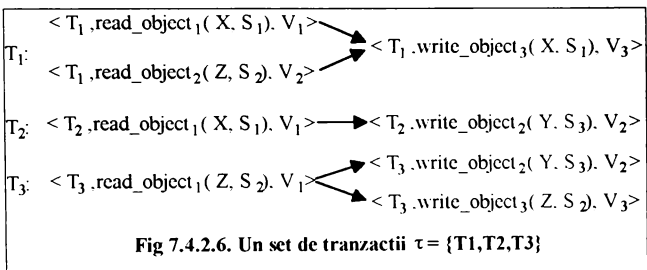


Fig. 7.4.2.6. Un set de tranzații $\tau = \{T1, T2, T3\}$

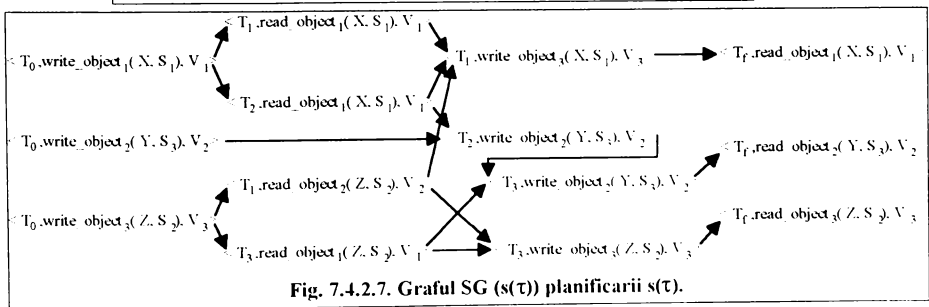


Fig. 7.4.2.7. Graful SG (s(tau)) planificării s(tau).

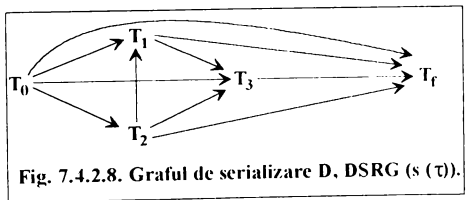
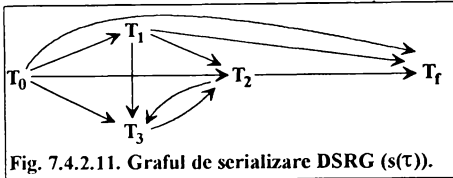
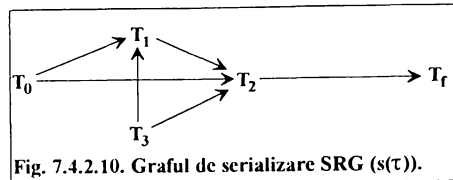
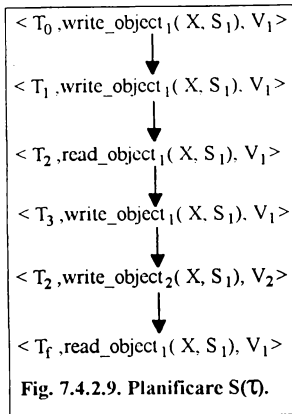


Fig. 7.4.2.8. Graful de serializare D, DSRG (s(tau)).



7.4.3. Mecanisme tradiționale pentru controlul concurenței

Problema consistenței datelor în aplicațiile tradiționale de baze de date a fost redusă la testarea serializabilității planificărilor; acest lucru se datorează faptului că nu se cunosc constrîngerile de consistență: fiecare acțiune într-o tranzație se abstractizează printr-o operație *read* sau *write*. Problema obținerii serializabilității poate fi descompusă în două subprobleme: *sincronizarea read-write* și *sincronizarea write-write*, notate pe scurt sincronizările *rw* și *ww* ([BG81]). Sincronizarea *rw* se referă la serializarea tranzațiilor în așa fel încît fiecare operație *read* citește aceeași valoare ca și cînd tranzațiile ar fi executate serial; sincronizarea *ww* se referă la serializarea tranzațiilor în așa fel încît ultima operație *write* a fiecărei tranzații lasă datele în aceeași stare ca și execuția serială. Respectarea serializărilor *rw* și *ww* generează o prelucrare consistentă. Astfel, deși sistemul de gestiune al bazei de date distribuite nu dispune de informații asupra constrîngerilor de consistență specifice aplicației, garantează consistența, permițînd numai execuții serializabile ale tranzațiilor concurente. Conceptul de serializabilitate al tranzațiilor este propriu mecanismelor tradiționale pentru controlul concurenței. Dacă sunt disponibile informații semantice despre tranzații și operațiile lor pot fi permise și planificări neserializabile, dar totuși consistente (vezi secțiunea 7.4.7.2).

Mecanismele tradiționale de rezolvare ale problemelor controlului concurenței se împart în următoarele categorii: mecanisme bazate pe blocarea cu *lock*-uri, mecanisme bazate pe mărcile de timp și mecanisme pentru controlul optimist al concurenței. Unele mecanisme adaugă și alte concepte precum granularitatea blocărilor și tranzații încuibate.

7.4.4. Mecanisme bazate pe blocarea cu *lock*-uri

Cele mai utilizate mecanisme pentru controlul concurenței sunt cele bazate pe *lock*-uri, acestea reprezentînd o facilitate specială pentru controlul accesului la obiectele partajate. Utilizînd *lock*-uri, tranzațiile pot bloca obiecte în scopul asigurării accesului exclusiv, pe durata unei stări temporare inconsistente. Aceasta implică faptul că alte tranzații care necesită accesul la obiectele blocate, fie vor fi puse în așteptare, fie vor fi abortate, fie vor întrerupe tranzația care deține obiectul blocat. Un obiect blocat și modificat de o tranzație abortată sau întreruptă poate fi deblocat numai după ce starea obiectului este restaurată la cea de dinaintea începerii tranzației.

Există trei tipuri de acțiuni relative la *lock*-uri: (1) $\langle T, \text{lock}_X \langle E, S \rangle \rangle$ specifică o cerere a tranzației *T* pentru un *lock* exclusiv la obiectul $\langle E, S \rangle$; cererea este satisfăcută numai dacă nici o altă tranzație nu deține un *lock* de orice tip asupra obiectului; (2) $\langle T, \text{lock}_S \langle I, S \rangle \rangle$ specifică o cerere a tranzației *T* pentru un *lock* partajat la obiectul $\langle E, S \rangle$; cererea este satisfăcută numai dacă nici o altă tranzație nu deține un *lock* exclusiv asupra obiectului; (3) $\langle T, \text{unlock} \langle E, S \rangle \rangle$ specifică eliberarea *lock*-ului deținut asupra obiectului $\langle E, S \rangle$.

Două *lock*-uri sunt *in conflict* dacă implică același obiect și cel puțin unul dintre ele este un *lock* exclusiv. Una tranzație *T* se permite să obțină un *lock* asupra unui obiect, de mai multe ori în moduri diferite: modul exclusiv de blocare este menținut, însă, pînă la o cerere *unlock*.

Avantajele principale ale metodelor bazate pe *lock*-uri sunt: (1) ordinea de serializare nu este stabilită a priori, ci depinde de progresul execuțiilor tranzațiilor; (2) sunt disponibili algoritmi ierarhici eficienți. Dezavantajele sunt legate de (1) accesul la obiecte este restrictiv, datorită *lock*-urilor; gradul de concurență este mai mic decît într-o serializare D; (2) pot apărea interblocări.

7.4.4.1. Mecanismul blocării în două faze (2PL)

Mecanismul blocării în două faze, introdus în [EGT78] este astăzi unanim acceptat ca soluție standard pentru controlul concurenței în bazele de date convenționale.

Mecanismul impune următoarele constrângeri ale manierei în care tranzacțiile sunt structurate: (1) tranzacția este *corect structurată*, adică citește un obiect numai cît timp deține un *lock* exclusiv sau partajat asupra obiectului și scrie un obiect numai dacă deține un *lock* exclusiv asupra obiectului; (2) tranzacția este *în două faze*, adică nu execută o acțiune *lock* după ce a executat o acțiune *unlock*. Prima fază a tranzacției, denumită *fază de acumulare*, începe cu prima acțiune și se termină la prima acțiune *unlock*. Cea de-a doua fază, denumită *fază de restrangere*, începe cu prima acțiune *unlock* și continuă pînă la ultima acțiune.

Execuția concurentă a unui set de tranzacții generează execuția unei secvențe de acțiuni intermixate asupra obiectelor din sistem. Fiecare planificare dintr-un nod constituie o planificare locală, iar controlul structurii ei este aspectul fundamental al mecanismului de control al concurenței. O planificare este *legală* dacă în cadrul ei tranzacțiile nu contin *lock*-uri în conflict. Următoarele două teoreme stabilesc relația dintre structura tranzacțiilor și structura planificărilor locale și planificărilor globale.

Teorema 1: (pentru un sistem centralizat)

Dacă $\{T_1, T_2, \dots, T_k\}$ este un set de tranzacții corect structurate și în două faze, atunci orice planificare Q_0 este consistentă.

Teorema 2: (pentru un sistem distribuit)

Fie Q_1, \dots, Q_n un set de planificări locale ale tranzacțiilor $\{T_1, T_2, \dots, T_n\}$. Atunci există o planificare Q astfel încît: (a) secvențele de acțiuni din T_i și Q_i sunt subsecvențe ale lui Q_i ; (b) dacă fiecare planificare Q_i este legală, atunci și Q este legală.

Demonstrație:

Deoarece aceasta teoremă stabilește că execuția acțiunilor în planificările locale poate fi considerată ca execuția unei singure planificări într-un nod centralizat se propune următoarea demonstrație. Se consideră că fiecare nod j este dotat cu un ceas (imaginar) implementat printr-un simplu contor, a cărui valoare este desemnată prin C_j . Presupunem că fiecare manager de tranzacție (TM) menține de asemenea un ceas desemnat prin C_T . Dacă TM este în nodul i , atunci $C_T = C_i$. Cînd MT adaugă o cerere pentru o acțiune la obiectul din nodul j , incrementează mai întîi C_i și apoi trimite un mesaj: ("REQUEST", C_i , T , A , $\langle 0 \rangle$, V) la nodul j . La primirea mesajului, nodul j actualizează ceasul sau $C_j = \max(C_j, C_i)$. Cînd nodul j execută acțiunea A asupra obiectului O pentru tranzacția T , incrementează C_j și trimite un mesaj de confirmare la TM, ("ACK", C_j , T , A , $\langle 0 \rangle$, V). La primirea mesajului de confirmare, TM stabilește: $C_i = \max(C_i, C_j)$ și $C_T = C_i$. Pentru a forma o planificare globală echivalentă, se pot combina planificările locale, sortind acțiunile în ordine crescătoare după $(C_i, \text{număr nod})$. Aceasta planificare globală satisface proprietățile a) și b) și teorema este demonstrată.

Un exemplu de planificare legală și serială a tranzacțiilor din fig. 7.2.1.1. este prezentat în fig. 7.4.4.1.1.

Implementarea unui mecanism de control al concurenței bazate pe *lock*-uri într-un sistem distribuit poate fi făcută dacă se menține în fiecare nod un manager al *lock*-urilor, căruia se adresează acțiunile relative la *lock*-uri pentru obiectele de pe acel nod. În fig. 7.4.4.1.2. se arată un exemplu în acest sens: în nodul 1 sunt două obiecte A și B, iar în nodul 2 obiectul C. O tranzacție T va adresa acțiunile *lock* pentru obiectele A și B managerului de *lock*-uri local, iar cererile *lock* pentru obiectul C managerului de *lock*-uri de pe nodul 2. Managerul de *lock*-uri poate fi parte integrantă a managerului de obiecte.

Se apreciază că algoritmul 2PL este optimal, în sensul că asigură un grad mai mare de concurență decît alte metode dacă: (1) obiectele sistemului distribuit sunt independente și reprezintă același nivel de abstractizare; (2) tranzacțiile sunt dependente de date.

```

< T1. lock_X (A. S) >
< T1. read_object (A.S). BalantaA >
< T1. lock_X (B. S) >
< T1. read_object (B. S). BalantaB >
< T1. write_object (A. S). BalantaA-100 000 >
< T1. write_object (B. S). BalantaB+100 000 >
< T1. unlock (A. S) >
< T1. unlock (B. S) >
< T2. lock_X (B. S) >
< T2. read_object (B.S). BalantaB >
< T2. lock_X (C. S) >
< T2. read_object (C.S). BalantaC >
< T2. write_object (B. S). BalantaB-200 000 >
< T2. write_object (C. S). BalantaC+200 000 >
< T2. unlock (B. S) >
< T2. unlock (C. S) >
< T3. lock_S (A. S) >
< T3. read_object (A. S). BalantaA >
< T3. lock_S (B. S) >
< T3. read_object (B. S). BalantaB >
print BalantaA
print BalantaB
< T3. unlock (A. S) >
< T3. unlock (A. S) >
    
```

Fig. 7.4.4.1.1. O planificare legala si seriala

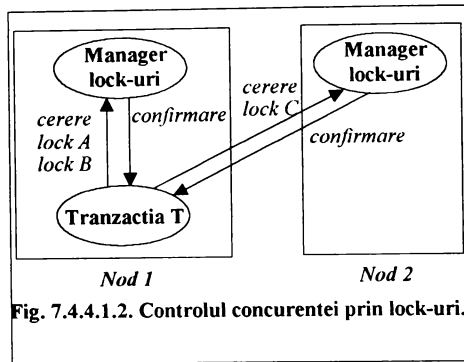


Fig. 7.4.4.1.2. Controlul concurenței prin lock-urii.

7.4.4.2 Algoritmi de blocare ierarhici

Algoritmii de blocare ierarhici pun problema granularității datelor blocate: aceasta constă în a selecta dimensiunea optimală a unității de blocare. Problema generează acceptarea unui compromis între maximizarea gradului de concurență și minimizarea regiei de sistem necesară pentru blocare. Pentru a maximiza gradul de concurență pentru tranzații mici, sunt adecvate multe unități de blocare mici; pentru a minimiza regia de sistem pentru tranzații mari, sunt adecvate puține unități de blocare dar mari. Conceptul de granularitate a fost propus pentru a încapsula ideea unor unități de blocare ierarhice. O ierarhie de granularități poate fi reprezentată printr-un graf aciclic ale cărui noduri sunt granularitățile și ale cărui arce conectează granularitățile de nivel înalt cu granularitățile de nivel jos; un exemplu este prezentat în figura 7.4.4.2.1.

Algoritmii de blocare ierarhici utilizează două tipuri de moduri de blocare: *modul de bază* (prezentat anterior la 7.4.4.1) și *modul intențiilor de blocare*. Două moduri de blocare sunt *compatibile* dacă permit accesul concurrent la același obiect. Pentru a bloca un obiect de la un anumit nivel, trebuie obținută mai întâi permisiunea pentru intenția de acces de la toate componentele ascendente de pe nivelele superioare; de exemplu, dacă un obiect A are deja asociat *lock-ul* de intenție I, atunci o cerere de blocare a lui A în modul S (shared) sau X (exclusiv) va fi nesatisfăcută deoarece s-ar putea ca A să aibă o componentă de nivel mai jos deja blocată în modul X; astfel, modurile X, S și I sunt incompatibile unul cu altul.

Pentru a permite însă un grad mai mare de concurență, s-au definit trei moduri de intenții de acces la un obiect: *intenția de acces în mod partajat* (desemnată prin IS), *intenția de acces în mod exclusiv* (desemnată prin IX) și *accesul mixt* (desemnat prin SIX). Cele cinci moduri de blocare pot fi rezumate astfel: (1) modul IS asupra unui obiect specifică faptul că tranzația dorește să citească componentele descendente ale obiectului (permite blocarea partajată a acestora în modurile S sau IS); (2) modul IX asupra unui obiect specifică faptul că tranzația dorește să scrie și/sau să citească componentele descendente ale obiectului; (3) modul S asupra unui obiect specifică faptul că tranzației i se permite accesul în mod partajat la obiectul respectiv și implicit la toate componentele descendente; (4) modul SIX asupra unui obiect specifică faptul că tranzația va deține modul S și IX asupra obiectului (poate citi obiectul) și că are dreptul să scrie și/sau citească descendenții săi; (5) modul X asupra unui obiect specifică faptul că tranzația va avea acces exclusiv asupra obiectului și, implicit, acces exclusiv și asupra componentelor descendente. În tabelul din fig. 7.4.4.2 se prezintă compatibilitatea dintre cele 5 moduri.

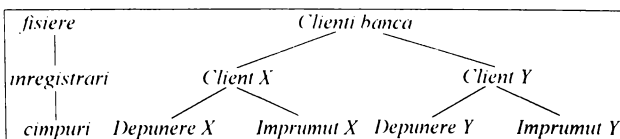


Fig. 7.4.4.2.1. O ierarhie de lock-urii.

mod de blocare curent	mod de blocare cerut	IS	IX	S	SIX	X
IS		DA	DA	DA	DA	NU
IX		DA	DA	NU	NU	NU
S		DA	NU	DA	NU	NU
SIX		DA	NU	NU	NU	NU
X		NU	NU	NU	NU	NU

Fig. 7.4.4.2.2 Compatibilitatea modurilor de blocare pentru algoritmi de blocare ierarhici.

Pentru a obține corect accesul la o componentă obiect, o tranzacție trebuie să se conformeze următorului algoritm:

Algoritm:

- (1) Pentru a obține accesul S sau IS la o componentă a unui obiect, o tranzacție trebuie să dețină un lock IS sau IX asupra tuturor componentelor predecesoare;
- (2) Pentru a obține accesul X, SIX sau IX la o componentă a unui obiect, o tranzacție trebuie să dețină un lock SIX sau IX asupra componentelor predecesoare;
- (3) Toate lock-urile trebuie eliberate fie la sfârșitul execuției tranzacției, fie în ordinea frunze către rădăcină în timpul execuției tranzacției.

În figura 7.4.4.2.3. se prezintă ca exemplu două tranzacții care actualizează cîmpurile "Depunere Y", utilizînd corect lock-uri ierarhice.

Se poate demonstra că orice planificare concurentă S(T) a tranzacțiilor T, care respectă algoritmul de mai sus, este serializabilă; în figura 7.4.4.2.4 se prezintă un exemplu de planificare serializabilă a tranzacțiilor din figura 7.4.4.2.3. (în exemplu s-au omis operațiile de scriere și citire).

Algoritmul aratat prezintă însă pericolul apariției interblocării. Un exemplu de planificare a tranzacțiilor T₁ și T₂ în care apare interblocare, este prezentat în figura 7.4.4.2.5.

```

T1:
<T1.lock_IX. "Clienti_banca">
<T1.lock_X. "Client X">
<T1.lock_S. "Client Y">
<T1.read "Depunere X">
<T1.read "Depunere Y">
<T1.read "Imprumut Y">
<T1.write "Imprumut X">
<T1.unlock "ClientY">
<T1.unlock "Client X">
<T1.unlock "Clienti_banca">

T2:
<T2.lock_IX. "Clienti_banca">
<T2.lock_IX. "Client Y">
<T2.lock_X. "Depunere Y">
<T2.read "Depunere Y">
<T2.lock_IS. "Client X">
<T2.lock_S. "Depunere X">
<T2.read "Depunere X">
<T2.write "Depunere Y">
T2.unlock "Depunere X">
<T2.unlock "Client X">
<T2.unlock "Depunere Y">
<T2.unlock "Client Y">
<T2.unlock "Clienti_banca">
Fig. 7.4.4.2.3. Tranzacții care utilizează lock-uri ierarhice
    
```

```

S(t):
<T1. lock_IX. "Clienti_banca">
<T2. lock_IX. "Clienti_banca">
<T1. lock_X. "Client X">
<T1. lock_S. "Client Y">
<T1. unlock. "Client_Y">
<T2. lock_IX. "Client_Y">
<T2. lock_X. "Depunere Y">
<T1. unlock. "Client_X">
<T2. lock_IS. "Client_X">
<T2. lock_S. "Depunere X">
<T1. unlock. "Clienti_banca">
<T2. unlock. "Depunere X">
<T2. unlock. "Client_X">
<T2. unlock. "Depunere Y">
<T2. unlock. "Client_Y">
<T2. unlock. "Clienti_banca">

Fig.7.4.4.2.4 Planificare serializabilă a tranzacțiilor T1 și T2
    
```

```

S(t):
<T1. lock_IX. "Clienti_banca">
<T2. lock_IX. "Clienti_banca">
<T1. lock_X. "Client_X">
<T2. lock_IX. "Client_Y">
<T2. lock_X. "Depunere Y">
<T1. lock_S. "Client_Y">
<T2. lock_IS. "Client_X">
... interblocare
Fig.7.4.4.2.5. Planificare ce generează o interblocare
    
```

7.4.4.3 Grafuri aciclice de lock-uri

Conceptul de lock-uri ierarhice poate fi extins și la un set de obiecte organizate într-un graf orientat aciclic, în care componentă a unui obiect poate avea mai mult decât un părinte; un exemplu este prezentat în fig. 7.4.4.3. În acest model, se consideră următoarele extensii asupra lock-urilor: (1) O tranzație deține implicit un lock S, dacă a obținut, implicit sau explicit un lock S, SIX sau X la cel puțin unul din părinții componente; aceasta implică faptul că tranzația a obținut explicit un acces S, SIX sau X la unul dintre predecesori. (2) O tranzație deține implicit un lock X asupra tuturor părinților; aceasta înseamnă că pentru o componentă în fiecare cale care conduce la componentă C s-a obținut explicit modul X și pentru toți predecesorii săi s-a obținut modul IX sau SIX. Pentru a obține accesul, fiecare tranzație trebuie să respecte următorul protocol: (1) Pentru a obține un lock S sau IS, o tranzație trebuie să dețină un lock IS (sau un lock mai puternic) la cel puțin unul dintre părinți și la toți predecesorii săi dintr-o cale către rădăcina. (2) Pentru a obține un lock IX, SIX sau X la o componentă a unui obiect, o tranzație trebuie să dețină un lock IX (sau un lock mai puternic) pentru toți părinții săi. Aceasta implică faptul că tranzația va deține un lock IX (sau unul mai puternic) pentru toți predecesorii acestor componente și nici o altă tranzație nu va putea deține un lock asupra unei componente predecesoare care este incompatibil cu IX. (3) cind o tranzație se termină, lock-urile pot fi eliberate în orice ordine; în caz contrar, ele trebuie eliberate în ordinea frunza-rădăcină. În figura 7.4.4.3.b) se prezintă un exemplu de obținere a unui acces S la înregistrările dintr-un fișier F, utilizind indexul I.

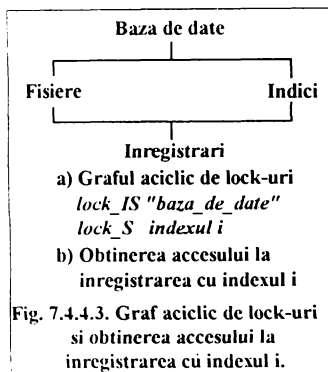


Fig. 7.4.4.3. Graf aciclic de lock-uri și obținerea accesului la înregistrarea cu indexul i.

7.4.4.3.b) se prezintă un exemplu de obținere a unui acces S la înregistrările dintr-un fișier F, utilizind indexul I.

7.4.5 Mecanisme bazate pe marca de timp

Una dintre problemele mecanismelor de control al concurenței este posibilitatea apariției interblocărilor. Problema interblocărilor poate fi rezolvată dacă se asociază fiecărei tranzații o marcă de timp. Mecanismele de control al concurenței din această secțiune se bazează toate pe marca de timp.

7.4.5.1 Concepte de bază

Marca de timp asociată unei tranzații, $TS(T_i)$ este un număr obținut prin concatenarea valorii ceasului (logic sau fizic) cu identificatorul nodului de la care a fost inițiată tranzația. Pentru fiecare obiect al sistemului distribuit de obiecte s-asociază două mărci de timp: *marca de timp pentru citire* $TS_R(x)$ și *marca de timp pentru scriere* $TS_W(x)$; acestea se setează la valoarea mărcii de timp asociată ultimei tranzații care a citit, respectiv a scris obiectul. Mecanismele de control al concurenței bazate pe marca de timp presupun că ordinea de serializare a unei planificări concurente a tranzațiilor este strict determinată prin valorile mărcilor de timp ale tranzațiilor. Algoritmii de control bazați pe marca de timp trebuie să verifice că această ordine nu este violată în timpul execuției concurente a tranzațiilor (datorită întârzierilor de comunicație, de exemplu). Ordinea de serializare hotărâtă a priori este principalul dezavantaj al metodelor bazate pe marca de timp, deoarece orice violare a acestei ordini conduce la abortări multiple ale tranzațiilor ale căror cereri de acces sosesc prea târziu; de aceea au fost propuși diverși algoritmi al căror scop era minimizarea efectelor acestui dezavantaj.

Algoritmii bazați pe marca de timp pot fi grupați în trei categorii de bază: determiniști, nedeterminiști și semi-determiniști.

Algoritmii determiniști includ algoritmi în care managerii de obiecte permit accesul la obiecte numai dacă acesta nu provoacă restartarea unei tranzații; cu alte cuvinte, o cerere de acces la un obiect de către o tranzație este buferată de un manager de obiecte pînă nu mai există tranzații cu mărci de timp mai mici care cer acces la același obiect. În această categorie se includ *algoritmii conservativi* ([CP84], [BG81], [CG87]).

Algoritmii nedeterminiști se consideră descendenți direcți ai procedurii *two-phase commit*. Cînd o tranzație T_i cer un acces la un obiect în conflict cu accesul la același obiect garantat unei alte tranzații T_j , $TS(T_i) > TS(T_j)$, accesul la obiectul T_i depinde de starea tranzației T_j . Dacă cererea tranzației T_j a fost acceptată, atunci cererea lui T_i va fi rejectată; în caz contrar, cererea lui T_i este acceptată, provocînd rejectarea cererii lui T_j .

Algoritmii semi-determiniști includ algoritmi care utilizează o coadă de tranzații, în scopul minimizării numărului de abortări. Fiecărui obiect i se asociază o coadă de tranzații care așteaptă accesul la un obiect. Un manager de obiecte, după execuția cererii de acces curente, preia tranzația cu cea mai mică marcă de timp din coadă. În cazul în care coada asociată unui obiect este vidă în momentul sosirii unei cereri, atunci cererea este acceptată imediat cu condiția ca tranzația de la care provine cererea să aibă o marcă de timp mai mare decât tranzația care a lansat ultima cerere de acces executată.

Avantajele metodelor bazate pe mărcile de timp sunt: (1) O tranzație poate actualiza un obiect chiar imediat după ce a citit-o; actualizările pot fi executate concurrent atîta timp cît ordinea mărcilor de timp nu este violată. Astfel, gradul de concurență este îmbunătățit. (2) Interblocările sunt eliminate. Dezavantajul principal este faptul că ordinea de serializare hotărâtă a priori conduce la abortarea și restartarea tranzațiilor care încearcă să acceseze obiecte în ordine inversă a mărcilor de timp. Minimizarea restartărilor la infinit este problema crucială.

7.4.5.2 Algoritmul fundamental pentru controlul concurenței pe baza mărcii de timp

Algoritmul fundamental în controlul concurenței pe baza mărcii de timp poate fi descris, într-o primă fază, ignorând respectarea atomicității, astfel:

Algoritm:

1. Managerul de tranzații asociază o marcă de timp unică fiecărei noi tranzații T .
2. Managerul de tranzații asociază marca de timp tuturor acțiunilor trimise pentru execuție managerului de obiecte, în cadrul acelei tranzații.
3. Fiecare manager de obiecte folosește un planificator pe bază de marcă de timp care planifică acțiunile locale corespunzător ordinii mărcilor de timp. Pentru fiecare obiect O , planificatorul înregistrează cea mai mare marcă de timp a tuturor acțiunilor *read_object* terminate și a tuturor acțiunilor *write_object* terminate. Aceste mărci de timp vor fi notate cu TS_R și TS_w .
4. Pentru o acțiune *read_object*, planificatorul compară marca sa de timp TS cu TS_w . Dacă $TS < TS_w$ planificatorul rejectează acțiunea iar managerul de tranzații abortează tranzația; în caz contrar, planificatorul plasează acțiunea într-o coadă FIFO cu acțiuni în așteptare la obiectul O și setează TS_R la $\max(TS_R, TS)$.
5. Pentru o acțiune *write_object*, planificatorul compară marca sa de timp TS cu maximum dintre TS_R și TS_w . Planificatorul rejectează acțiunea dacă TS este mai mică; în caz contrar, plasează acțiunea în coadă și setează TS_w la TS .
6. Dacă este necesar să se aborteze o tranzație managerul de tranzații îi va asigna o nouă marcă de timp (mai mare) și o va restarta.

Algoritmul de bază introduce însă trei probleme specifice: *restartări ciclice*, *restartări necesare* și *restartări la infinit*. Pentru a prezenta *restartările ciclice* se consideră următoarea situație: două tranzații T_1 și T_2 fiecare constind într-o acțiune *read_object(S,O)* urmată de o acțiune *write_object(S,O)* relative la același obiect O sunt executate de următoarea planificare:

$$\dots \langle T_1, \text{read_object}(S,O), V_1 \rangle \dots \langle T_2, \text{read_object}(S,O), V_1 \rangle \dots \\ \langle T_1, \text{write_object}(S,O), V_2 \rangle \dots \langle T_2, \text{write_object}(S,O), V_2 \rangle \dots$$

Dacă marca de timp a lui T_1 este mai mică decât cea a lui T_2 , atunci $\langle T_1, \text{write_object}(S,O), V_2 \rangle$ va fi rejectată iar T_1 restartată cu o marcă de timp mai mare. Dacă imediat T_1 execută o nouă acțiune *read_object(S,O)* atunci $\langle T_1, \text{write_object}(S,O), V_2 \rangle$ va fi rejectată iar T_2 restartată. Apare astfel o secvență de restartări ale tranzațiilor T_1 și T_2 .

Restartările necesare apar dacă secvența de execuție a acțiunilor este în ordinea inversă a mărcilor de timp. Un exemplu în acest sens îl constituie cazul unei tranzații T_2 , a cărei marcă de timp este mai mare decât cea a lui T_1 , dar toate acțiunile lui T_2 sunt terminate complet înaintea oricărei acțiuni a lui T_1 ; T_1 va fi rejectată chiar dacă planificarea serială în care T_2 este urmată de T_1 este perfect acceptabilă.

Restartările la infinit apar cînd o tranzație mai veche, la fiecare restartare, intră în conflict cu o tranzație mai nouă: aceasta este o problema de aminare la infinit.

Pentru a asigura atomicitatea tranzațiilor este necesar să se integreze procedura de mai sus cu procedura *two-phase commit* care permite abortarea sau realizarea tuturor operațiilor unei tranzații (în principal, contează operațiile *write*). În acest caz, procedura *two-phase commit* se realizează pe baza utilizării unor operații *prewrite*, care sunt executate de tranzații în locul operațiilor *write*. Operațiile *prewrite* sunt bufărate în managerii de obiecte și nu se aplică obiectelor locale: bufărea unei operații înseamnă înregistrarea ei într-un bufer pentru o execuție ulterioară. Numai cînd toate operațiile *prewrite* ale unei tranzații sunt acceptate, tranzația poate fi realizată și efectuate și operațiile *write*. Dacă toate operațiile *prewrite* ale unei tranzații au fost acceptate, operațiile *write* corespundente nu pot fi rejectate.

Scopul procedurii *two-phase commit* este de a evita actualizările parțiale provocate, de exemplu de avaria unui nod. În plus, procedura *two-phase commit* previne așa numitele *abortări în cascadă* sau *efectul domino*; acestea au loc cînd abortarea unei tranzații parțial comisă provoacă abortarea altor tranzații parțial comise care au citit valori scrise de tranzația abortată, efectul putîndu-se propaga în continuare la alte tranzații.

În descrierea algoritmului fundamental pentru controlul concurenței pe baza mărcii de timp, se vor folosi următoarele variabile: TS_R - marca de timp maximă a acțiunilor *read-object* executate, TS_w - marca de timp maximă a acțiunilor *write-object*, \min_TS_R - marca de timp minimă a oricărei acțiuni *read-object* bufărate, \min_TS_w = marca de timp minimă a oricărei acțiuni *write-object* bufărate, \min_TS_p - marca de timp minimă a oricărei acțiuni *prewrite* bufărate. Se presupune ca: (1) toate aceste variabile, ca și mărcile de timp generate de managerii de tranzații sunt inițializate cu zero; (2) \min_TS_p , \min_TS_w și \min_TS_R rețin valorile cînd bufărele corespundente devin vide; (3) fiecare acțiune *read-object*, *write-object* și *prewrite* este o operație indivizibilă raportată la planificatorul local.

În aceste condiții, algoritmul poate fi descris astfel:

```

<Ti, read-object (X, S), V> ::
  if TS(Ti) < TSw →
    <rejectionea cererea. abortează tranzacția și restartează tranzacția cu o nouă marca de timp>
  [] TS(Ti) ≥ TSw →
    if TS(Ti) > min_TSr(X) →
      <plasează cererea de citire în buffer>
    [] TS(Ti) ≤ min_TSr(X) →
      <read_X>
      TSr(X) := max(TS(Ti), TSr(X))
    fi
  fi

<Ti, write-object (X, S), V> ::
  if TS(Ti) > min_TSr(X) or TS(Ti) > min_TSw(X) →
    plasează cererea de scriere în buffer
  [] TS(Ti) ≤ min_TSr(X) and TS(Ti) ≤ min_TSw(X) →
    <realizează_write (Ti, X)>
    TSw(X) := TS(Ti)
  fi

<Ti, prewrite (X, S), V> ::
  if TS(Ti) < min_TSr(X) or TS(Ti) < min_TSw(X) →
    <rejectionea cererea. abortează și restartează Ti tranzacția cu o nouă marca de timp >
  [] TS(Ti) ≥ min_TSr(X) and TS(Ti) ≥ min_TSw(X) →
    <plasează cererea prewrite în buffer >
  fi
  
```

Procedura realizează write(T_i , X) realizează actualizarea obiectului și înlătură din buffer cererea prewrite. Dacă valoarea min_TS_r(X) crește după această ștergere, atunci inițiază o procedură care testează buffer-ul și determină dacă noua valoare min_TS_r(X) nu deblochează o cerere read sau write a altei tranzacții, fie ea T_j, pentru care este adevărată acum condiția TS(T_j) < min_TS_r(X). Deblocarea acestor cereri poate provoca actualizarea din nou a variabilelor min_TS_r(X) și min_TS_w(X), și, în consecință, testarea din nou a buffer-ului.

7.4.5.3. Optimizări ale algoritmului de bază

În cele ce urmează sunt descrise trei tehnici de optimizare ale algoritmului de bază ([BK91]):

1. Regula lui Thomas

Dacă două cereri write-object consecutive sunt în conflict și cererea mai nouă are marca de timp mai mică decât cealaltă, atunci cea mai nouă poate fi ignorată, eliminând astfel acțiunile write-object necesare.

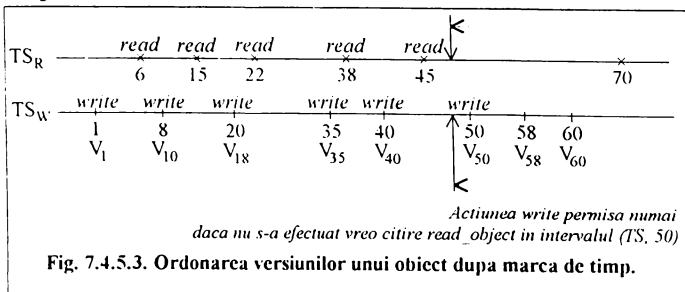
2. Folosirea mai multor versiuni ale unui obiect

Prin folosirea mai multor versiuni pentru un obiect, cu costul irosirii de spațiu, se evită rejectarea acțiunilor read-object. În acest caz un obiect O constă într-o secvență de versiuni, fiecare de forma <TS_w, valoare>, care înregistrează mărcile de timp și valorile conținute în urma execuției cererilor write-object <O,S> și o secvență de mărci de timp TS_r pentru toate acțiunile read-object <O,S>. Crearea unor copii multiple ale unui obiect este arătată în fig. 7.4.5.3. Pentru a se reduce cerințele de memorie, versiunile mai vechi ale unui obiect sunt descărcate.

Implementarea acestei tehnici se face ținând cont de următoarele reguli:

- (1) O acțiune read-object cu marca de timp TS se execută prin citirea versiunii obiectului O cu cea mai mare marcă de timp mai mică decât TS și adăugând marca TS setului de mărci TS_r.
- (2) O acțiune write-object cu marca de timp TS se execută dacă nu există TS_r în intervalul semideschis [TS, cea mai mică TS_w > TS). În caz contrar se creează o nouă versiune a obiectului O cu marca de timp TS.

3. Algoritmi Wait-Die și Wound-Wait



Algoritmi Wait-Die și Wound-Wait se utilizează în conjuncție și cu prevenirea interblocării, fiind prezenți de aceea timp pe larg în secțiunea 7.5.2. Pe scurt, protocolul wait-die, o tranzacție T_j poate aștepta după T_i numai dacă T_j este mai veche (marca de timp mai mică); altfel este abortată și apoi restartată. În protocolul wound-die, o tranzacție T_j poate aștepta după T_i numai dacă T_j este mai nouă (cu o marcă de timp mai mare)

7.4.5.4. Algoritmul conservativ de ordonare a mărcilor de timp

Algoritmul conservativ de ordonare al mărcilor de timp a fost introdus pentru a reduce numărul abortărilor din algoritmul de bază (CP84, [CG87]). În cadrul acestui algoritmului o cerere de acces este acceptată numai dacă nu provoacă restartarea altei tranzații din sistem. În scopul eliminării restartărilor, operațiile tranzațiilor mai recente sunt bufferate pînă la terminarea execuțiilor unor tranzații mai vechi, cu care sunt în conflict. Aceasta presupune că fiecare nod comunică cu celelalte în scopul de a culege informații despre setul de tranzații curent procesate.

Fiecare manager de obiecte are asociat un set de $2*N$ cozi, care conțin cererile de acces de la tranzațiile inițiate de la N manageri de tranzații: TM_1, \dots, TM_N . Fiecare manager de obiecte conține o coadă $RQueue(TM_i)$ de cereri de acces *read* și o coadă $WQueue(TM_i)$ de cereri de acces *write* pentru fiecare modul TM_i . Fiecare cerere de acces sosită la un modul OM_i de la nodul TM_i este plasată în $RQueue(TM_i)$ sau $WQueue(TM_i)$. Se presupune că toate cererile de acces sosesc de la modulul TM_i la modulul OM_i în ordinea mărcilor de timp; pentru aceasta sunt necesare două condiții: (1) ordinea mărcilor de timp nu trebuie să fie modificată de mediul de comunicație; (2) fiecare modul TM trebuie să trimită cererile de acces în ordinea mărcilor de timp. Aceste cerințe pot fi îndeplinite în două moduri. Primul este de a executa tranzațiile în ordinea strict serială, la fiecare modul TM . Aceasta micșorează gradul de concurență al sistemului la N . Un al doilea mod este de a executa toate tranzațiile dintr-un modul TM , în două etape. Toate cererile de citire se generează în prima etapă și toate cererile de scriere în etapa a doua; aceasta asigură execuția concurrentă la nivelul unui modul TM_i .

Se prezintă în continuare procedurile de citire și scriere pentru algoritmul conservativ. Pentru simplificare, se presupune că, atunci cînd un modul OM_k primește o cerere de scriere sau citire, fiecare coadă $RQueue(TM_i)$ și $WQueue(TM_i)$, $i=1,2,\dots,N$ conține cel puțin o cerere de acces la obiecte; deoarece s-a presupus că toată cererile de acces sosesc la OM_k în ordinea mărcilor de timp, $RQueue$ și $WQueue$ conțin întotdeauna cele mai vechi cereri de acces netratate. Cu alte cuvinte, un modul OM_k nu poate primi o cerere de acces de la un manager de tranzații TM_i astfel încît marca de timp a acestei tranzații să fie mai mică decît marca de timp a tranzațiilor în așteptare în $RQueue(TM_i)$ și $WQueue(TM_i)$. Procedurile de citire și scriere pentru algoritmul conservativ sunt următoarele:

```

<Ti, read_object(X,S),V>::
test = true
fa TMk=TM1, TMN →
fa Tj ∈ WQueue(TMk) →
if TS(Tj) < TS(Ti) → test = false fi
af
if test → <read X>
[] not test → <plaseaza cererea
read de la Ti in RQueue(TMk)>
fi

```

```

<Ti, write_object(X,S),V>::
test = true
fa TMk=TM1, TMN →
fa Tj ∈ WQueue(TMk) →
if TS(Tj) < TS(Ti) → test = false fi
fa
fa Tj ∈ RQueue(TMk) → if TS(Tj) < TS(Ti) → test = false fi
af
if test → < read (X) >
[] not test → <plaseaza cererea write de la Ti in WQueue(TMk)>
fi

```

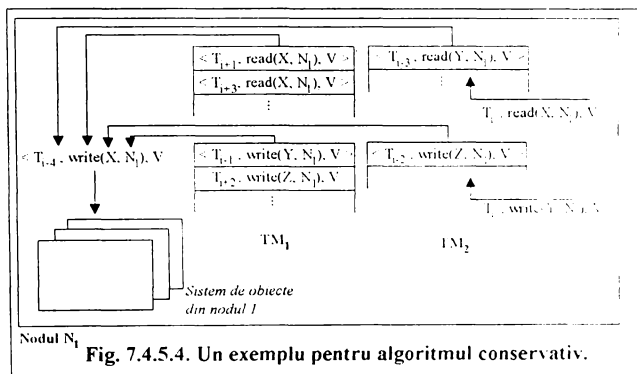


Fig. 7.4.5.4. Un exemplu pentru algoritmul conservativ.

În algoritmul de mai sus s-a presupus că atunci cînd un modul OM_k primește o cerere de acces la date, nici una din cele $2*N$ cozi ale sale, corespunzătoare modulelor TM_1, \dots, TM_N nu este vidă; în caz contrar, ultima cerere dintr-o coadă ar trebui întârziată pînă la sosirea unei noi cereri. Problema concilor vide poate fi însă soluționată și astfel, fiecare modul TM , transmite periodic acțiuni nule fiecărui modul OM ; o acțiune nulă nu specifică nici o operație asupra obiectelor, dar, fiind însoțită de marca de timp, rolul său este de a furniza ultima marcă de timp generată de un modul TM . O acțiune nulă poate fi transmisă periodic sau la cererea directă a unui modul OM .

În fig. 7.4.5.4. se prezintă un exemplu de utilizare al algoritmului conservativ. Sistemul este format din două noduri N_1 și N_2 , fiecare conținând câte un modul OM și TM ; în figură se consideră execuția cererilor de acces la obiectul din N_2 . Se presupune că mărcile de timp ale tranzațiilor sunt egale cu indicii lor. Figura arată starea cozilor $RQueue$ și $WQueue$ ale modulului OM_1 cînd se primește o cerere de citire $\langle T_{1,1}, read(X, N_1), V \rangle$ în momentul prelucrării cererii $\langle T_{1,2}, write(X, N_1), V \rangle$. Conform algoritmului conservativ, planificarea pentru execuție a cererilor este următoarea: $\dots \langle T_{1,2}, write(X, N_1), V \rangle, \langle T_{1,1}, write(Y, N_1), V \rangle, \langle T_{1,2}, write(Z, N_1), V \rangle, \langle T_{1,1}, write(Y, N_1), V \rangle, \langle T_{1,1}, read(X, N_1), V \rangle$.

7.4.6 Controlul optimist al concurenței

În multe aplicații, mecanismele bazate pe *lock*-uri restring concurența și generează o regie de sistem sporită. În ([BG81][Gos91]) se arată că mecanismele bazate pe *lock*-uri prezintă următoarele dezavantaje: (1) introduc o regie de sistem sporită nenecesară pentru tranzațiile read-only care nu afectează integritatea datelor; (2) nu există mecanisme bazate pe *lock*-uri care să ofere un grad înalt de concurență; (3) interzice eliberarea *lock*-urilor până la sfârșitul tranzației, deși nu este absolut necesar; aceasta se face în scopul evitării abortărilor în cascadă, micșorează gradul de concurență; (4) fiind se lucrează cu obiecte mari, aflate în memoria secundară, blocarea obiectelor accesate frecvent (referite ca noduri suprasolicitate) pe timpul cât se așteaptă pentru accesul la memoria secundară poate provoca o descreștere semnificativă a concurenței; (5) este posibil ca, în unele cazuri, tranzațiile să nu se suprapună (în accesul la aceleași date); în aceste cazuri, blocările sunt nenecesare. Ele sunt utile numai în cazurile cele mai dezavantajoase, când tranzațiile accesează aceleași date.

Pentru a evita aceste dezavantaje, Kung și Robinson au prezentat pentru prima oară un algoritim pentru controlul optimist al concurenței.

7.4.6.1. Concepte de baza

Ideea de bază a mecanismelor ce oferă un control optimist al concurenței este presupunerea că, în accesul la date, conflictele apar destul de rar în execuția tranzațiilor concurente. De aceea nu se abortează sau întârzie nici o tranzație în timpul desfășurării ei, ca în cazul mecanismelor bazate pe *lock*-uri sau pe marca de timp; în schimb, în momentul realizării tranzației se execută un test de validare pentru a se detecta eventualele conflicte care au apărut și care ar fi putut afecta consistența sistemului. Dacă rezultatul testului este pozitiv, atunci tranzația va fi abortată și restartată.

Mecanismele care oferă un control optimist al concurenței sunt adecvate pentru prelucrările în care conflictele la date apar rar.

Algoritmii pentru controlul optimist al concurenței se pot grupa în două categorii. Prima categorie ([BG81], [KR81], [CG87]) cuprinde algoritmii care construiesc graful de precedență (vezi sect. 7.4.2) și îl testează pentru determinarea aciclității. Graful relației de precedență este actualizat de câte ori se acceptă o cerere de scriere sau citire. După cum s-a aratat în secțiunea 7.4.2., absența ciclurilor în graful relației de precedență este suficientă pentru ca o planificare a unui set de tranzații să fie serializabilă. Problema crucială în sistemele distribuite este însă proiectarea unui mecanism eficient pentru menținerea grafului relației de precedență. O soluție a acestei probleme este menținerea unei copii a grafului relației de precedență în fiecare nod și actualizarea ei la fiecare operație *read* sau *write*; aceasta soluție implică însă o circulație mare de mesaje în sistem și de aceea se utilizează mai mult cea de a doua categorie de algoritmi.

A doua categorie de algoritmi ([KR81], [CO82], [CG87], [Gos91], [BK91]) execută o tranzație în trei faze: o fază de citire, o fază de validare și o fază de scriere. Pe tot parcursul fazei de citire, tranzația se execută local, în zona de lucru proprie. În faza de validare, se testează dacă, prin execuția operațiilor de actualizare pregătite de tranzație în zona de lucru nu se violează consistența obiectelor. Dacă testul nu este validat, urmează faza de scriere în care are loc scrierea obiectelor actualizate din zona de lucru în sistemul de obiecte; dacă testul nu este validat, tranzația este abortată și eventual, restartată, ca o tranzație nouă.

Controlul optimist al concurenței prezintă următoarele avantaje: (1) execuția tranzațiilor care nu periclitează consistența datelor nu este întârziată; (2) nu pot apare interblocări și restartări ciclice, dar și următoarele dezavantaje: (1) trebuie colectate informații despre tranzații, în vederea realizării fazei de validare; (2) este posibilă restartarea infinită.

7.4.6.2. Algoritmi pentru controlul optimist al concurenței

În aceasta secțiune vor fi prezentați doi algoritmi pentru controlul optimist al concurenței în sisteme centralizate, după care aceștia se vor extinde la variantele pentru sisteme distribuite ([MOO87], [CG87]).

Execuția unei tranzații se desfășoară în trei faze: faza de citire, faza de validare și faza de scriere. Fazele de citire și de scriere au fost explicate în secțiunea precedentă; centrul de greutate în aceasta secțiune va cădea pe faza de validare. Voi fi prezentați doi algoritmi: algoritmul de validare serială și de validare paralela.

În cadrul algoritmului de validare serială, fiecărei tranzații T_i i se aștează un identificator unic $TC(T_i)$ care nu este neaparat o marcă de timp, ci un număr generat de un numărator global (GTC). Identificatorii se asignează tranzațiilor după execuția cu succes a fazei de validare. Fie $TC_{start}(T_i)$ cel mai mare identificator de tranzație generat în momentul inițierii tranzației T_i și $TC_{finish}(T_i)$ cel mai mare identificator de tranzație generat în momentul începerii fazei de validare a lui T_i . Tranzațiile ale căror identificatori sunt: $TC_{start}(T_i)+1, TC_{start}(T_i)+2, \dots, TC_{finish}(T_i)$ constituie setul de tranzații a căror fază de scriere s-a executat concurrent cu faza de citire a lui T_i ; acest set se va desemna prin $T_{set}(T_i)$. Cu alte cuvinte, $T_{set}(T_i)$ se definește astfel: $T_{set}(T_i) = \{T_j | TC(T_j) \in \{TC_{start}(T_i)+1, \dots, TC_{finish}(T_i)\}\}$

Setul de date citite de tranzația T_i va fi desemnat prin $O_c(T_i)$, iar cel de date scrise de tranzația T_i va fi desemnat prin $O_w(T_i)$. Cu aceste condiții, algoritmul de validare serială poate fi descris astfel:

ValidareSeriala(T_i) ::

```

TCfinal( $T_i$ ) = GTC; test = true
fa  $T_j \in T_{sf}(T_i) \rightarrow$ 
    if  $O_r(T_i) \cap O_w(T_j) \neq \emptyset \rightarrow$  test = false fi
af
if test  $\rightarrow$  GTC = GTC+1
    TS( $T_i$ ) = GTC
    < write phase >
[] not test  $\rightarrow$  < abort  $T_i$  >
fi
    
```

Pentru a se putea garanta serializabilitatea, procedurile ValidareSeriala() trebuie executate în regim de excludere reciproca (sunt secțiuni critica). Ca atare execuția unei faze de validare și de scriere a unei tranzații blochează parțial întregul sistem deoarece nici o altă tranzație nu mai poate intra în aceste faze în acest timp, în schimb, fazele de citire ale altor tranzații pot fi executate.

Se observă că acest algoritm impune actualizarea serială a obiectelor, în ordinea identificatorilor de tranzații, de aceea algoritmul se numește de *validare seriala*.

Tranzațiilor read-only nu li se asociază identificatori deoarece ele nu iau parte în validarea tranzațiilor următoare din sistem; pentru aceste tranzații $TC_{start}(T_i) = TC_{final}(T_i)$.

Algoritmul de *validare paralela* permite unei faze de validare a unei tranzații să se desfășoare concurrent cu fazele de scriere ale altor tranzații. În acest algoritm se presupune că tranzației i se asociază identificatorul la sfârșitul fazei de scriere. Ca și în algoritmul de validare serială, la inițierea unei tranzații și la sfârșitul fazei de citire se citește conorul GTC pentru a determina valorile $TC_{start}(T_i)$ și $TC_{final}(T_i)$. Tranzațiile cu identificatorii $TC_{start}(T_i)+1, TC_{start}(T_i)+2, \dots, TC_{final}(T_i)$ constituie setul $T_w(T_i)$ al tranzațiilor care au executat fazele de scriere concurrent cu faza de citire a lui T_i . Tranzațiile continute în $T_w(T_i)$ și-au terminat faza de scriere când T_i își începe faza de scriere. În plus, sistemul pastrează informații despre setul de tranzații care au terminat faza de citire dar nu au terminat faza de scriere. Acest set desemnat prin T_{activ} conține tranzațiile ale căror faze de scriere și validare au fost executate concurrent cu aceleași faze ale lui T_i . Algoritmul de validare paralela poate fi descris astfel:

ValidareParalela(T_i) ::

```

TCfinal( $T_i$ ) = GTC
TlocActiv( $T_i$ ) = Tactiv #se realizeaza o copie locala a lui Tactiv
Tactiv( $T_i$ ) = Tactiv  $\cup$  { $T_i$ }
test = true
fa  $T_j \in T_w(T_i) \rightarrow$ 
    if  $O_r(T_i) \cap O_w(T_j) \neq \emptyset \rightarrow$  test = false fi
af
fa  $T_j \in T_{locActiv}(T_i) \rightarrow$ 
    if  $(O_r(T_i) \cup O_w(T_j)) \cap O_w(T_i) \neq \emptyset \rightarrow$  test = false fi
if test  $\rightarrow$  <write phase>
    GTC = GTC+1
    TS( $T_i$ ) = GTC
    Tactiv = Tactiv - { $T_i$ }
[] not test  $\rightarrow$  Tactiv = Tactiv - { $T_i$ }
    <abort  $T_i$ >
fi
    
```

În acest algoritm numai operațiile asupra lui GTC și T_{activ} se execută în excludere reciprocă.

Un dezavantaj al acestui algoritm este acela că o tranzație T_i poate fi abortată în mod nencesar datorită unui conflict cu o tranzație invalidată $T_k \in T_{locActiv}$ care este de asemenea abortată. În exemplul următor se prezintă o asemenea situație. Se consideră planificarea concurrentă a tranzațiilor T_1, T_2, T_3 și T_4 din fig. 7.4.6.2. și se presupune că tranzațiile acționează asupra următoarelor elemente:

```

O1( $T_1$ ) = { $O_1, O_2$ }, Ow( $T_1$ ) = { $O_1, O_2$ }
Or( $T_2$ ) = { $O_3, O_4$ }, Ow( $T_2$ ) = { $O_2, O_3$ }
Or( $T_3$ ) = { $O_2, O_5$ }, Ow( $T_3$ ) = { $O_1, O_6$ }
Or( $T_4$ ) = { $O_4, O_6$ }, Ow( $T_4$ ) = { $O_6$ }
    
```

Validarea tranzațiilor se desfășoară astfel: (1) T_1 este validată fără condiții deoarece atât $T_w(T_1)$ cât și $T_{locActiv}$ sunt vide; (2) pentru validarea lui T_2 se ia în considerare doar $T_w(T_2) = \{T_1\}$; $T_{locActiv}(T_2)$ este vidă; T_2 este validată deoarece $O_r(T_2)$ și $O_w(T_1)$ sunt disjuncte; (3) pentru validarea lui T_3 se ia în considerare $T_w(T_3) = \{T_1, T_2\}$; $T_{locActiv}(T_3)$

este vida. Deoarece $O_r(T_1) \cap O_w(T_1) = \{O_2\}$ și $O_r(T_1) \cap O_w(T_2) = \{O_2, O_3\}$, tranzacția T_1 este abortată; (4) pentru validarea lui T_4 se ia în considerare $T_{is}(T_4) = \{T_1, T_2\}$ și $T_{loc,av}(T_4) = \{T_3\}$. Tranzacția T_4 este abortată deoarece :

$$(O_r(T_4) \cup O_w(T_4)) \cap O_w(T_3) = \{O_6\}$$

Totuși, tranzacția T_4 este abortată inutil, deoarece conflictul apare cu o tranzacție care va fi invalidată.

Pentru extinderea algoritmilor de validare serială și paralela la sisteme distribuite se consideră următoarele: (1) se presupune că fiecare tranzacție primește un identificator unic global și acest identificator se asociază tuturor subtranzacțiilor lui T_i : $T_i^{S_1}, \dots, T_i^{S_n}$; (2) validarea tranzacțiilor se desfășoară pe două nivele: local și global. Nivelul de *validare locală* implică acceptarea fiecărei subtranzacții $T_i^{S_k}$, local, în nodul S_k , conform procedurilor de validare prezentate mai sus. Nivelul de *validare globală* implică acceptarea tranzacției distribuite T_i , după acceptarea tuturor subtranzacțiilor $T_i^{S_1}, \dots, T_i^{S_n}$ (*two-phase commit*).

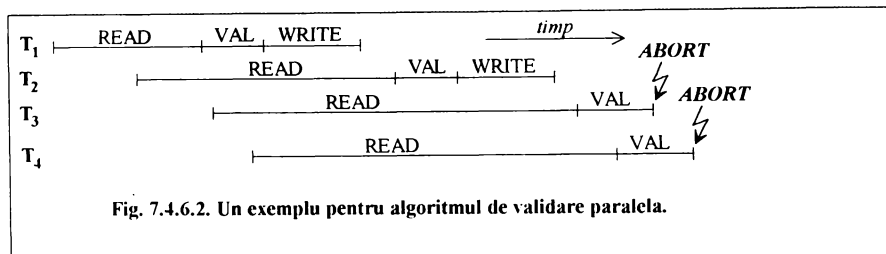


Fig. 7.4.6.2. Un exemplu pentru algoritmul de validare paralela.

7.4.7 Extinderea mecanismelor tradiționale pentru controlul concurenței

Mecanismele tradiționale de control al concurenței impun serializarea execuției tranzacțiilor, având în vedere numai operațiile de citire și scriere, datorită lipsei de informații relative la operațiile specifice aplicației. Aceasta lipsă de informație conduce la imposibilitatea de a specifica sau verifica constrângeri semantice asupra datelor; însă în condițiile existenței unor astfel de informații, se poate construi o planificare neserializabilă dar consistentă. Prin egalizarea noțiunii de consistență cu noțiunea de serializabilitate se micșorează gradul de concurență.

7.4.7.1 Cerințe ale controlului concurenței în aplicațiile avansate

Se apreciază că aplicațiile avansate impun următoarele noi cerințe mecanismelor de control al concurenței ([YEK87], [BK91]):

(1) Furnizarea unui suport adecvat pentru tranzacțiile lungi. În aplicații lungi, precum compilări de module sursă sau proiectări de circuite, operațiile asupra obiectelor sunt operații de lungă durată. Dacă aceste operații ar fi încapsulate în tranzacții, aceste tranzacții vor fi, de asemenea, tranzacții lungi, pentru care nu se accepta blocarea unuia până la comiterea altuia.

(2) Furnizarea unui suport care să ofere control utilizatorului. Mecanismul de control al concurenței trebuie să furnizeze acest suport în scopul controlului de către utilizator a unor aplicații de natură interactivă și nedeterministă; utilizatorul trebuie să poată iniția o tranzacție, execută operații în cadrul ei, în mod interactiv, să o poată restructura dinamic și comite sau aborta în orice moment. Natura nedeterministă a tranzacțiilor implică faptul că mecanismul de control al concurenței nu va putea determina dacă se violează constrângerile de consistență decât în momentul execuției și validării operațiilor asupra datelor. Acest lucru poate avea ca urmare abortarea unor tranzacții în urma unor ore de muncă; de aceea utilizatorul trebuie să poată inversa explicit efectul unor operații, pentru a refăce consistența datelor;

(3) Furnizarea unui suport pentru cooperare sinergetică. Cooperarea între utilizatori în timpul execuției unor aplicații precum sisteme CAD/CAM, SDE, poate avea implicații asupra mecanismului de control al concurenței. În cadrul acestor medii de proiectare, este probabil ca utilizatorii să dorească schimbul de informații între ei la un moment dat (să partajeze anumite informații). Activitățile a doi sau mai mulți utilizatori care lucrează cu aceleași obiecte partajate, pot să nu fie serializabile. Utilizatorii pot să-și transmită controlul asupra obiectelor într-un mod care să poată fi realizat printr-o planificare serializabilă sau să modifice concurent două părți ale unui obiect, cu intenția de a le integra într-o nouă versiune a obiectului. Acest mod de partajare și schimb de informații a fost denumit *interacțiune sinergetică* ([YEK87], [BK91]).

În plus față de aceste cerințe, multe aplicații avansate necesită suport pentru obiecte complexe. De exemplu, într-un sistem de proiectare software, obiectele pot fi organizate ierarhic (proiectele sunt constituite din module, modulele din funcții). Tehnicile pentru gestionarea obiectelor complexe sunt proprii programării orientate pe obiecte; totuși complexitatea structurii și dimensiunea obiectelor poate avea influențe puternice asupra mecanismului de control al concurenței.

7.4.7.2 Extinderea mecanismelor bazate pe serializabilitate și renunțarea la serializabilitate

Pentru rezolvarea problemelor legate de tranzațiile lungi au fost propuse două soluții: extinderea mecanismelor bazate pe serializabilitate, care generează tot planificări serializabile și mecanisme care generează planificări neserializabile (renunțarea la serializabilitate). Ambele soluții, pentru creșterea gradului de concurență, utilizează informații semantice specifice aplicațiilor.

Extinderea mecanismelor bazate pe serializabilitate.

În mecanismele tradiționale, toate operațiile asupra obiectelor sunt abstractizate în două operații: citire și scriere. Această abstractizare este necesară pentru proiectarea unor mecanisme de control al concurenței cu caracter general, care să nu depindă de specificul aplicațiilor, de exemplu mecanismul blocării în două faze. Performanțele acestui mecanism pot fi însă inacceptabile în aplicațiile avansate, deoarece obiectele pot rămâne blocate mult timp după utilizarea lor, blocând în acest fel alte tranzații. Mecanismele de control optimist al concurenței oferă de asemenea performanțe slabe pentru aplicațiile avansate, deoarece când rata conflictelor crește, care este în general cazul tipic al aplicațiilor avansate, crește numărul de abortări și restartări ale tranzațiilor. O rezolvare a problemelor introduse de tranzațiile lungi este utilizarea unor informații semantice despre tranzații și operații pentru a extinde mecanismele tradiționale. În secțiunea 7.4.8. se prezintă modalitatea de extindere a mecanismelor bazate pe lockuri și a celor de control optimist al concurenței în vederea satisfacerii cerințelor tranzațiilor lungi.

Renunțarea la serializabilitate.

O altă soluție pentru rezolvarea cerințelor tranzațiilor lungi este renunțarea la planificările serializabile și folosirea unor planificări care se bazează pe informațiile semantice relative fie la obiecte, fie la operațiile specifice aplicației. Mecanismele bazate pe informații semantice se împart în două mari grupuri ([SZ89]): primul grup definește și folosește proprietăți relative la concurență ale obiectelor ([SGA87]), iar cel de-al doilea utilizează informații semantice relative la tranzații. Algoritmii din primul grup impun constrângeri în intercalarea operațiilor tranzațiilor pe baza conflictelor dintre operațiile, definite asupra unor obiecte de anumite tipuri. Controlul concurenței în accesul la obiecte ale căror tipuri și operații se cunosc va fi prezentat în secțiunea 7.4.9. În secțiunea 7.4.8. vor fi comentați algoritmi din al doilea grup.

7.4.8. Mecanisme pentru controlul concurenței în aplicații avansate

În această secțiune va fi prezentată extinderea mecanismelor bazate pe lock-uri și a celor de control optimist al concurenței la aplicații complexe.

7.4.8.1 Mecanismul blocării altruiste

Un important gen de informații care poate fi folosit pentru a crește concurența este declararea obiectelor ca libere în momentul terminării utilizării lor; aceste obiecte vor putea fi folosite de alte tranzații. Un mecanism care folosește această tehnică este *mechanismul blocării altruiste* ([SGA87], [BK91]). Acest mecanism utilizează informații despre așa numitele *șabloane de acces* la obiecte ale unei tranzații pentru a decide ce obiecte pot fi eliberate. Se folosesc două tipuri de informații: *șabloane de acces negativ*, care descriu obiectele care nu vor fi accesate de tranzație și *șabloanele de acces pozitiv* care descriu obiectele ce vor fi accesate de tranzație și ordinea de acces. Împreună, aceste șabloane permit tranzațiilor lungi să elibereze obiectele după utilizarea lor. Setul de obiecte care au fost blocate și apoi eliberate de o tranzație lungă se numește *urna tranzației* ("wake"). Eliberarea unui obiect este o operație *unlock* condiționată deoarece permite altor tranzații să folosească obiectul numai dacă respectă cele două condiții de mai jos, care asigură serializabilitatea.

O planificare în *două faze cu eliberare* ("two phase with release") se definește ca o planificare care respecta condițiile:

Condiția 1

Două tranzații nu pot deține lock-uri asupra aceleiași date, simultan, cu excepția cazului când una dintre tranzații a eliberat lock-ul înainte ca cealaltă să-l obțină. În acest caz, lock-ul deținut de ultima se spune că se afla în *urna* primei tranzații.

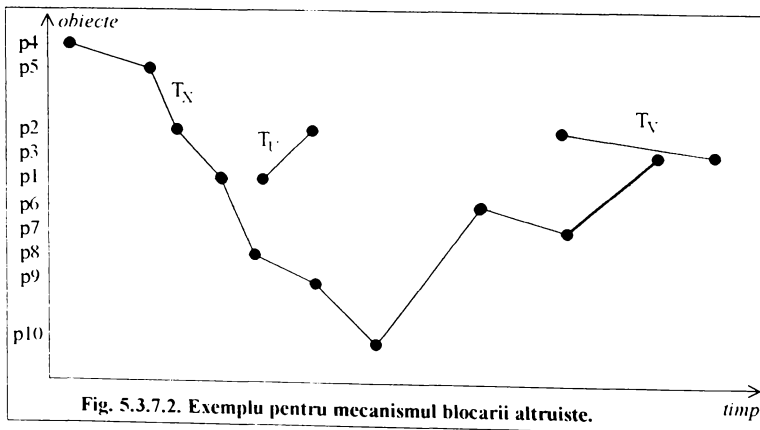
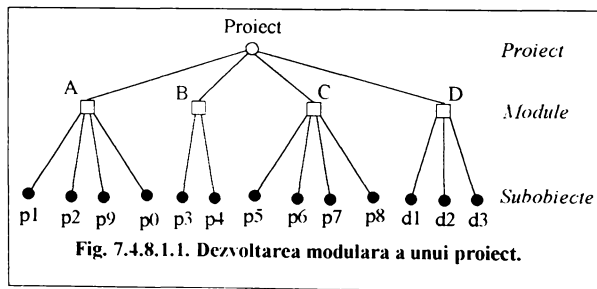
Condiția 2

Dacă o tranzație este în *urna* altei tranzații atunci ea trebuie să fie complet în *urna* acelei tranzații. Aceasta înseamnă că dacă tranzația T_2 obține un lock asupra unui obiect eliberat de tranzația T_1 , atunci orice obiect accesat de T_1 și T_2 care este deținut la momentul curent de T_2 , trebuie să fi fost eliberat de T_1 înainte de a fi fost obținut de T_2 .

Aceste două condiții de protocol garantează serializabilitatea tranzațiilor fără a altera structura lor. Protocolul presupune că tranzațiile sunt *programate* și nu *controlate* de utilizator. Exemplul următor ilustrează acest mecanism. Considerăm aplicația din figura 7.4.8.1.1 în care se presupune că o echipa de proiectanți dezvoltă modular un sistem de circuite electronice. Fiecare modul din proiect conține descrierea unui anumit număr de plăci (subobiecte). Presupunem că utilizatorul X, fiind un membru nou în cadrul proiectului compus din proiectanții U și V dorește să se familiarizeze cu toate schemele. În această situație, X va iniția o tranzație nouă, T_x , care va accesa toate subobiectele, pe rând; tranzația T_x va fi o tranzație lungă în decursul căreia X va citi fiecare subiect, va adăuga câteva comentarii și va elibera subiectul. În același timp, utilizatorul U va iniția o altă tranzație (scurtă),

T_U , care va accesa doar două subobiecte p_1 și p_2 din modulul A. Presupunem că T_X a accesat și eliberat deja p_2 , în momentul curent fiind în citirea lui p_1 . T_U trebuie să aștepte pînă cînd T_X termină utilizarea lui p_1 ; în acest moment T_U poate accesa p_1 , găsindu-se în *urma* lui T_X . T_U poate intra în *urma* lui T_X deoarece toate obiectele la care T_U dorește acces (p_1 și p_2) sunt în *urma* lui T_X . După terminarea lucrului cu p_1 , T_U poate accesa imediat p_2 , fără a întârzia, deoarece p_2 a fost eliberat de T_X . Presupunem acum că V inițiază o altă tranzație (scurta), T_V , în care dorește să acceseze p_2 și p_3 , care nu este în *urma* lui T_X . T_V poate accesa p_2 după ce T_U termină tranzația, dar, apoi trebuie să aștepte pînă fie T_X eliberează pe p_3 (cu alte cuvinte p_3 intră în *urma* lui T_X) sau pînă T_X se termină. Dacă T_X nu accesează niciodată p_3 , T_V este forțată să aștepte terminarea lui T_X . Pentru a crește concurența în aceste de situații, în [SGA87] s-a propus un mecanism pentru extinderea dinamică a *urmei* unei tranzații lungi, pentru a permite tranzațiilor scurte care sunt deja în *urma* unei tranzații lungi să-și continue execuția. Acest mecanism utilizează informații despre obiectele care *nu* vor fi accesate de o tranzație, în scopul de a introduce aceste obiecte automat în *urma* tranzației. Continuînd în acest mod exemplul, p_3 va fi adăugat automat în *urma* lui T_X , ceea ce va permite tranzației T_V să acceseze p_3 fără întârziere. Figura 7.4.8.1.2. prezintă acest exemplu. Fiecare tranzație este reprezentată printr-o linie frîntă, care unește subobiectele accesate în cadrul tranzației. Tranzația lungă T_X accesează subobiectele $p_4, p_5, p_2, p_1, p_3, p_9, p_{10}, p_6$ și p_7 ; p_3 nu este accesat de T_X (linia groasă din figură nu aparține tranzației lungi). Tranzația T_U este complet în *urma* lui T_X , deoarece fiecare subobiect accesat de T_U (p_1 și p_2) a fost accesat înainte de T_X . În schimb, T_V nu se află în *urma* lui T_X . Pentru a se permite execuția lui T_V , se expandează *urma* lui T_X adăugându-se p_3 (în figura acest lucru se simbolizează cu linia groasă). Planificarea celor trei tranzații este serializabila și echivalentă cu execuția serială T_X, T_U și T_V .

Avantajul principal al mecanismului blocării altruiste este posibilitatea de a utiliza eliberarea obiectului înainte de terminarea tranzațiilor. Mecanismul generează planificări serializabile iar în cazul absenței informațiilor se reduce la blocarea standard în două faze. Prezintă însă un important dezavantaj: nu evită abortările în cascadă, cînd o tranzație lungă este abortată, toate tranzațiile scurte aflate în *urma* ei trebuie abortate, *chiar dacă sunt deja terminate*.



7.4.8.2 Mecanismul validării instantanee

Mecanismul blocării altruiste are la baza mecanismul tradițional de blocare în două faze (2PL) și ca atare moștenește unele dezavantaje ale acestuia (de exemplu, poate provoca interblocări). O alternativă este folosirea mecanismului de control optimist al concurenței, ca punct de plecare. Mecanismele tradiționale de validare a tranzațiilor prezintă însă și ele dezavantaje precum: (1) necesitatea de a aborta o tranzație datorită unui conflict în faza de validare, chiar dacă acest conflict nu conduce la o violare a criteriului de serializabilitate, (2) irosirea timpului datorită descoperirii conflictelor în faza de validare, (3) riscul restartărilor infinite, iar în cazul validării paralele, a restartărilor ciclice. Dezavantajul de tipul (1), și ca urmare și (2) provine din definiția prea largă a noțiunii de conflict. În [PSU86] se arată că riscul restartării inutile a tranzațiilor poate fi micșorat dacă se împart conflictele de acces la date în două categorii: *conflicte grave* în cazul cărora tranzațiile trebuie abortate și apoi restartate și *conflicte aparente* în cazul cărora tranzațiile își pot continua execuția. Tot în [PSU86] se propune mecanismul validării instantanee a tranzațiilor care utilizează distincția între cele două tipuri de conflicte. Pentru a ilustra mecanismul validării instantanee se consideră exemplul din figura 7.4.8.2. Presupunem planificarea tranzațiilor T_x , T_U , T_V din fig.7.3.7.3. Conform mecanismului tradițional de validare, T_U și T_V ar trebui să fie abortate datorită conflictelor cu tranzația T_x :

$$O_w(T_x) \cap O_r(T_U) = \{p_1\}$$

$$O_w(T_x) \cap O_r(T_V) = \{p_2\}$$

Cu toate acestea, conflictul dintre T_U și T_x este diferit față de conflictul dintre T_V și T_x . Tranzația T_U citește subiectul p_1 , care este apoi modificat de T_x ; astfel $T_U >_{p_1} T_x$. Deoarece în faza de scriere, T_U poate accesa iarăși pe p_1 , se poate spune că și $T_x >_{p_1} T_U$ și deci se formează un ciclu în graful relației de precedență. Astfel, execuția concurentă (din figură) a tranzațiilor T_x și T_U violează criteriul de serializabilitate și de aceea T_U trebuie abortată. Conflictul dintre T_V și T_x , dimpotrivă, nu violează criteriul de serializabilitate. Tranzația T_V citește subiectul p_2 după ce acesta a fost actualizat de T_x ; astfel, $T_x >_{p_2} T_V$. Deoarece algoritmul de validare actualizează subiectele în conformitate cu ordinea identificatorilor de tranzații, pentru orice subiect p , $p \in \{O_w(T_x) \cap O_w(T_V)\}$ relația $T_x >_p T_V$ este adevărată. Astfel, nu este necesară abortarea tranzației T_V .

O altă constatare care se poate face în urma analizei acestui exemplu este aceea că acest conflict dintre T_x și T_U poate fi detectat mai devreme decât în faza de validare a lui T_U : în momentul în care T_x se termină (momentul snapshot din figură). Detectarea mai timpurie a conflictelor poate mări eficacitatea unui algoritim de validare.

Prima modificare a mecanismului tradițional de validare a tranzațiilor constă în adăugarea în setul $O_r(T_i)$ al unei tranzații T_i a identificatorilor tranzațiilor care și-au terminat faza de scriere în timpul fazei de citire a lui T_i (identificatori *end-of-transaction*); acești identificatori vor fi desemnați prin EOT_i . Setul $O_r(T_i)$ poate fi ordonat în ordinea cronologică a execuției operațiilor de citire. De aceea, pentru a valida tranzația T_i în raport cu o tranzație $T_j \in T_{sf}(T_i)$ este suficient să se considere numai un subset al lui $O_r(T_i)$ de la începutul fazei de citire a lui T_i până la apariția lui EOT_j . Fie, de exemplu setul $O_r(T)$ al unei tranzații T ordonat astfel: $O_r(T) = \{a, b, EOT_j, c, EOT_k, d, e, EOT_l\}$. Potrivit mecanismului tradițional de validare a tranzației T în raport cu $T_{sf}(T) = \{T_j, T_k, T_l\}$ sunt necesare următoarele teste:

$$O_w(T_i) \cap \{a, b, c, d, e\}$$

$$O_w(T_j) \cap \{a, b, c, d, e\}$$

$$O_w(T_k) \cap \{a, b, c, d, e\}$$

Prin modificarea mecanismului de validare, sunt necesare testele:

$$O_w(T_i) \cap \{a, b\}$$

$$O_w(T_j) \cap \{a, b, c\}$$

$$O_w(T_k) \cap \{a, b, c, d, e\}$$

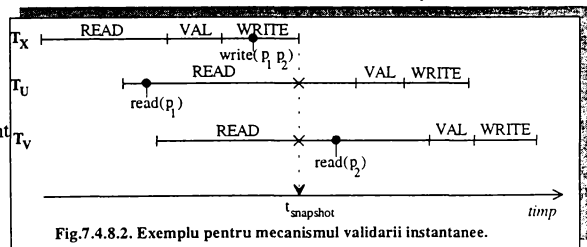


Fig.7.4.8.2. Exemplu pentru mecanismul validării instantanee.

A doua modificare a algoritmului constă în substituirea procedurii de validare (serială sau paralelă), executată o singură dată pe durata de viață a unei tranzații printr-un *număr de validări instantanee*. Validările instantanee se execută la sfârșitul fazei de scriere a fiecărei tranzații. Fie T_i o tranzație în momentul terminării fazei de scriere; în acest moment, fiecare tranzație T_j aflată în faza de citire trebuie validată în raport cu T_i . Setul $O_r(T_i)$ conținând datele citite de T_j până la sfârșitul fazei de scriere a lui T_i se compară cu $O_w(T_i)$. Dacă intersecția lor este nevidă, atunci T_j trebuie abortată deoarece a citit date care au fost actualizate de T_i . Fazele de scriere și validare se execută în regim de excludere reciprocă.

În concluzie, mecanismul de validare instantanee, față de mecanismul tradițional de validare are următoarele avantaje: (1) numărul de tranzații care se abortează este mai mic; (2) conflictele de acces la date sunt detectate mai devreme, prevenind irosirea unui timp suplimentar, până la validarea tranzației; (3) setul de date citite de o tranzație care se folosesc pentru validare este mai mic și astfel detectarea conflictelor se face mai rapid.

7.4.9. Mecanisme pentru controlul concurenței bazate pe informații semantice. Sincronizarea accesului la tipuri de date abstracte

Operațiile care se execută asupra tipurilor abstracte partajate trebuie să respecte următoarele cerințe [Gos91][CG87][SP84]: (1). Între apelul ultimei operații și execuția ei, obiectul poate să nu fie într-o stare consistentă. Nu se permite însă altor operații să observe sau să manipuleze obiectul în acest interval; (2). Când un grup de tranzații apelează operații ce se execută asupra unui obiect partajat, aceste operații trebuie astfel intercalate încât să mențină seriabilitatea tranzațiilor sau să respecte o *proprietate de ordine* (mai puțin restrictivă decât seriabilitatea) în grupul tranzațiilor; (3). Operațiile efectuate asupra obiectelor partajate trebuie să nu poată observa informații care se pot modifica dacă o tranzație va fi invalidată. Acest lucru necesită întârzierea execuției unor operații aparținând unor tranzații, pînă ce o altă tranzație este terminată, cu succes sau fără succes.

Metodele de sincronizare a accesului la tipuri abstracte partajate prezentate în acest capitol se concentrează asupra respectării cerințelor 2 și 3 de mai sus, tipice sistemelor bazate pe tranzații și reliefează importanța informațiilor semantice în extinderea concurenței.

7.4.9.1. Relația de dependență

S-a văzut că intercalarea operațiilor unui grup de tranzații este descrisă de planificări. O *planificare* poate fi definită ca o secvență de perechi <tranzație, operație> care semnifică ordinea în care sînt intercalate operațiile componente unor tranzații (vezi 7.4.2.).

Pentru a analiza tranzații care conțin operații asupra unor tipuri abstracte arbitrare, planificatoarele trebuie să aibă caracterizate explicit operațiile. De exemplu, o planificare poate conține operații pentru a introduce un element într-o coadă, pentru a insera o intrare într-un director, pentru a insera un nod într-un arbore B. Deoarece tipurile abstracte pot fi implementate într-o manieră arbitrară sînt necesare doua tipuri distincte de planificări: ([SP84]): (1) *planificările abstracte* caracterizează proprietățile de consistență a unui tip abstract, indiferent de implementare; (2) *planificările de apel* caracterizează proprietățile de concurență ale unei implementări.

Operațiile apelate într-o planificare abstractă reprezintă operații furnizate de o implementare ideală a tipului abstract. Această implementare este ideală în sensul că nu este permis să se modifice ordinea operațiilor executate asupra unui obiect. Astfel, planificările abstracte conțin operațiile în ordinea în care ele afectează obiectele partajate. De exemplu, să considerăm tipul abstract *listă FIFO* (o coadă) și operațiile de introducere a unui element la sfîrșitul cozii, QEnq, respectiv de extragere a unui element de la începutul cozii, QDeq. Presupunem că inițial coada Q este vidă. Atunci, planificarea abstractă:

t1: <T₁, QEnq(Q,S),X>

t2: <T₂, QEnq(Q,S),Y>

t3: <T₃, QDeq(Q,S),V>

descrie exact starea cozii la fiecare moment t1, t2, t3. Astfel, se poate deduce că a 3-a operație extrage în variabila V elementul X, iar în coadă rămîne numai elementul Y. Proprietățile de consistență ale unui tip pot fi specificate sub forma unui set de planificări abstracte. Acestea sînt denumite *planificări consistente* și o descriere a acestui set este o parte a specificării tipului.

Implementările reale pot reordona operațiile executate asupra unui obiect pentru a mări concurența, fără a modifica proprietățile de consistență. Să considerăm, de exemplu, o implementare a unui tip coadă în care elementele destinate a fi introduse în coadă sînt mai întii depozitate într-un buffer local al tranzației și introduse în bloc la sfîrșitul tranzației. Această implementare permite oricărui număr de tranzații să apeleze QEnq simultan, cu restricția serializării corecte cînd sînt implementate mai multe cozi. Prin execuția operațiilor de introducere în bloc, această implementare reordonează efectiv operațiile QEnq. Planificările de apel intercalează operațiile în ordinea apelurilor implementărilor reale care poate fi alta decît cea a implementării abstracte. De exemplu, utilizînd implementarea reală de mai sus a unui tip coadă, poate rezulta următoarea planificare de apel:

<T₂, QEnq(Q,S), Y>

<T₁, QEnq(Q,S), X>

<T₃, QDeq(Q,S), V>

Dacă T₁ se termină normal înaintea lui T₂ implementarea reală a reordonat cele două operații de introducere.

Correspondența între planificările de apel și planificările abstracte este de tipul mai multe la unul; fiecare planificare de apel implementează exact o planificare abstractă, dar o planificare abstractă poate fi implementată prin mai multe planificatoare de apel. Strategia de sincronizare utilizată de o implementare determină un set de planificări de apel denumite *planificări legale*, care sînt permise de implementare.

Prin examinarea unui planificator abstract, este posibil să se determine ce dependențe există între tranzațiile

implicate. Notăția $D:T_i:X \rightarrow_o T_j:Y$ va fi utilizată pentru a reprezenta dependența D formată cînd tranzacția T_i execută operația X și ulterior, tranzacția T_j execută operația Y asupra aceluiași obiect O . Setul de perechi ordonate $\{(T_i, T_j)\}$ pentru care există X, Y și O astfel încît $D:T_i:X \rightarrow_o T_j:Y$ formează o relație, notată $<_D$. Dacă $T_i <_D T_j$, se spune că T_i precede pe T_j și T_j depinde de T_i , prin dependența D .

Dacă se presupune că toate operațiile asupra obiectelor sînt alese arbitrar și nu se iau în considerare informații semantice utile pentru implementarea mecanismului de concurență, singura cerință este ca fiecare tranzacție, privită individual, să transforme obiectul dintr-o stare inițială consistentă într-o stare finală consistentă. În aceste condiții, numai planificatoarele abstracte serializabile pot garanta că păstrează corectitudinea tranzațiilor individuale. Urmînd presupunerea că toate operațiile sînt nedistingibile, numai o singură posibilă dependență D poate fi definită astfel: $T_1 <_D T_2$, dacă T_1 execută orice operație asupra unui obiect, înainte ca T_2 să execute orice operație asupra acestui.

Definiție.

Fie $<_D^*$ închiderea tranzitivă a relației $<_D$. O planificare este *ordonabilă* referitor la $\{D\}$ dacă $<_D^*$ este o relație de ordine. Cu alte cuvinte, nu există cicluri de forma

$$T_1 <_D T_2 <_D \dots <_D T_n <_D T_1.$$

În general, o planificare este ordonabilă, referitor la S , unde S este un set de relații de dependență, dacă toate relațiile din S au o închidere tranzitivă de ordine parțială. Relațiile din S sînt denumite *relații prescise*.

Se poate demonstra că ordonabilitatea referitoare la $\{<_D\}$ este echivalentă cu seriabilitatea.

Luînd în considerare informații semantice despre operațiile aplicate asupra obiectelor în tranzații, este posibil să se identifice relații de dependență pentru care este permisă formarea de cicluri. Să considerăm, de exemplu, mecanismul tradițional *read/write* de control al concurenței în accesul la un obiect. Un asemenea sistem implementează două tipuri de operații asupra obiectelor: *read_object*(R) și *write_object*(W). Există relații de dependență posibile între două tranzații care accesează același obiect:

1. $D_1:T_i:R \rightarrow_o T_j:R$; T_i citește un obiect, ulterior citit de T_j
2. $D_2:T_i:R \rightarrow_o T_j:W$; T_j citește un obiect, ulterior modificat de T_j
3. $D_3:T_i:W \rightarrow_o T_j:R$; T_i modifică un obiect, ulterior citit de T_j
4. $D_4:T_i:W \rightarrow_o T_j:W$; T_i modifică un obiect, ulterior modificat de T_j

Mecanismul anterior, nedistingînd aceste dependențe, previne formarea de cicluri în relația de dependență $<_D$, reuniunea tuturor patru de mai sus. În contrast, mecanismul de control al concurenței *Read/Write* ține însă cont de faptul că dependențele $R \rightarrow R$ nu influențează comportarea sistemului. Astfel, fiind dată o pereche de tranzații T_1 și T_2 , o planificare abstractă în care atît T_1 cît și T_2 execută o operație *Read* asupra unui obiect partajat, semantica operațiilor *Read* asigură faptul că nici T_1 , nici T_2 și nici altă tranzație din planificare nu determină dacă $T_1 <_{D_1} T_2$ sau $T_2 <_{D_1} T_1$.

În cazul *read/write*, condiția necesară de seriabilitate poate fi reformulată, în termenii relațiilor de dependență astfel: o planificare este serializabilă dacă este ordonabilă referitor la $\{D_2 \cup D_3 \cup D_4\}$. Permiînd operații de citire multiple, mecanismul *read/write*, permite formarea de cicluri în dependența $<_{D_1}$ și în relații compuse care includ $<_{D_1}$ prevenind însă formarea de cicluri în relații compuse din reuniuni $<_{D_2}$, $<_{D_3}$ și $<_{D_4}$. De exemplu, să considerăm următoarele planificări, cu efecte identice în starea sistemului:

- P1:: $\langle T_1, R(O,S),A \rangle$
 $\langle T_2, R(O,S),B \rangle$
 $\langle T_1, W(O,S),C \rangle$
- P2:: $\langle T_2, R(O,S),B \rangle$
 $\langle T_1, R(O,S),A \rangle$
 $\langle T_1, W(O,S),C \rangle$

În prima planificare, $T_1 <_{D_1} T_2$ și $T_2 <_{D_2} T_1$. Există deci un ciclu în relația $<_{D_1 \cup D_2}$, dar nu există cicluri în relațiile $<_{D_2 \cup D_3 \cup D_4}$. În a doua planificare, primele două operații sînt inversate și nu se formează astfel nici un ciclu.

Pe de altă parte, următoarele două planificări nu au efecte identice:

- P1:: $\langle T_1, R(O,S),A \rangle$
 $\langle T_2, W(O,S),B \rangle$
 $\langle T_1, W(O,S),C \rangle$
- P2:: $\langle T_2, W(O,S),B \rangle$
 $\langle T_1, R(O,S),A \rangle$
 $\langle T_1, W(O,S),C \rangle$

În acest caz, prima planificare nu este serializabilă deoarece $T_1 <_{D_2} T_2$ și $T_2 <_{D_4} T_1$, deci se formează astfel un ciclu în relația $<_{D_2 \cup D_4}$. Într-adevăr, T_1 observă O înainte de a fi modificat de T_2 , dar starea finală a lui O reflectă modificarea W din T_1 și nu din T_2 , implicând o execuție a lui T_1 înainte de T_2 . A doua planificare nu are cicluri și deci este serializabilă.

În concluzie, se poate afirma că **ordonabilitatea cu referire la un set de relații de dependență prescise furnizează o modalitate de caracterizare a planificatoarelor consistente**. Pentru un control al concurenței care forțează serializabilitatea fără informații semantice, setul de relații prescise conține $<_D$, care este echivalentă cu reuniunea tuturor relațiilor de dependență posibile. Pentru un mecanism Read/Write, setul conține relația compusă $<_{R \rightarrow W \cup W \rightarrow R \cup W \rightarrow W}$. Când se iau în considerare informații specifice tipului, relațiile prescise sînt de asemenea specifice tipului abstract. În unele cazuri aceste proprietăți sînt echivalente cu serializabilitatea, dar, în general, un grad mai mare de concurență se obține numai cînd sînt garantate proprietăți mai puțin restrictive.

În afară de utilitatea lor în ceea ce privește definirea consistenței, relațiile de dependență sînt utile și în detectarea situațiilor **de abort în cascadă**. Un abort în cascadă este posibil atunci cînd se formează o dependență între două tranzații, prima fiind neterminată. Un abort al primei tranzații poate genera un abort în cascadă a celor ce depind de ea. Condițiile în care un abort în cascadă apare efectiv depind de relația de dependență implicată și de proprietățile acelui obiect. Dacă considerăm cele patru relații de dependență prezentate în mecanismul de control al concurenței Read/Write atunci: (1) relația de dependență $R \rightarrow R$ nu va cauza niciodată abort în cascadă; (2) analog, relațiile $R \rightarrow W$ și $W \rightarrow W$ nu pot genera abort în cauză, deoarece, în ambele cazuri, a doua tranzație nu depinde de datele modificate de prima tranzație; (3) În absența unor informații semantice suplimentare, pentru relația $W \rightarrow R$ se poate presupune că un abort al primei tranzații o va afecta pe a doua, care trebuie și ea terminată prin abort.

După detectarea dependențelor care pot genera abort în cascadă, trebuie controlată formarea lor. Transpusă în termenii planificărilor abstracte, afirmația devine:

Teoremă.

Considerînd operațiile care formează dependența, acestea nu trebuie să se intercaleze în două tranzații, dacă prima tranzație nu este terminată înaintea celei de-a doua.

Trebuie făcută observația că nu toate planificatoarele consistente sînt și fără abort în cascadă și invers. Mecanismele de sincronizare trebuie să restrîngă setul de planificări de apel legale la planificările abstracte aflate la intersecția celor consistente și celor fără cascadă.

7.4.9.2. Specificarea tipurilor abstracte

Metodele tradiționale pentru sincronizarea accesului la o instanțiere a unui tip de date abstract partajat sînt destinate să asigure corectitudinea individuală a operațiilor asupra unui obiect, cu alte cuvinte să prevină anomaliile datorate concurenței. Într-un sistem bazat pe tranzații, suplimentar, pe lîngă corectitudinea individuală a operațiilor, trebuie asigurată pastrarea proprietăților de consistență la nivelul întregii tranzații, care poate conține un număr arbitrar de operații.

Acest subcapitol, prin compararea unui tip listă serializabilă cu un tip semilistă care restricționează cerințele de consistență, arată că dacă se exploatează la maximum informațiile semantice pentru operații, gradul de concurență poate fi crescut prin renunțarea la serializabilitate și înlocuirea ei cu ordonabilitatea referitoare la un set de relații de dependență.

A. Directoare

Ca prim exemplu, se consideră tipul de date *director*, creat în scopul menținerii unei corespondențe între un șir alfanumeric (cheie), și pointeri pentru obiecte arbitrare. Operațiile uzuale furnizate sînt:

1. *InsertDir(dir, str, ptr)*: se inserează obiectul cu cheia *str* în directorul *dir*. Se returnează OK sau cheie duplicată.
2. *DeleteDir(dir, str)*: se șterge obiectul cu cheia *str* din catalogul *dir*. Se returnează OK sau negăsit.
3. *LookupDir(dir, str)*: se caută cheia *str* în *dir*. Se returnează *ptr* sau negăsit.
4. *Dumpdir(dir)*: se returnează un vector de perechi $\langle str, ptr \rangle$ cu conținutul complet al directorului *dir*.

Intercalarea acestor operații aplate de tranzații executate concurrent trebuie controlată cu mare grijă. Tipul director trebuie să aibă următoarele proprietăți:

1. Dacă o tranzație ia în considerare o intrare cu cheia "S" dintr-un director, toate operațiile următoare ale aceleiași tranzații trebuie să poată accesa această intrare cu cheia "S" (numai dacă, evident, una dintre aceste operații nu o șterge). Nici unei alte tranzații nu îi va fi permis, însă, să șteargă intrarea "S" sau să modifice pointerul care îi corespunde înainte ca tranzația care o accesează să se termine.

2. Dacă o tranzație constată că un director *nu* conține o intrare cu cheia "S", nici unei alte tranzații nu i se poate permite să creeze o intrare cu această cheie, pînă la terminarea primei tranzații.

3. Dacă o tranzație creează o intrare nouă într-un director, nici unei alte tranzații nu i se poate permite să observe

această intrare pînă cînd prima tranzacție nu se termină cu succes.

4. Dacă o tranzacție șterge o intrare dintr-un director D și o înscriează într-un alt director D', nu trebuie să fie posibil ca o altă tranzacție să observe intrarea respectivă fie în D fie în D'.

O modalitate de respecta aceste proprietăți este prin modelarea operațiilor *InsertDir* și *DeleteDir* ca o operație de citire urmată de o operație de scriere și a operațiilor *Lookup* și *DumpDir* ca o operație de citire. Tipul director poate atunci fi specificat utilizînd relațiile de dependență *read/write* discutate în subcapitolul anterior.

Aceste informații semantice sînt însă insuficiente și restrîng concurența în mod necesar. Să presupunem, de exemplu, că implementăm tipul director utilizînd mecanismul standard de blocare *read/write* în două faze. Considerînd operația *LookupDir(dir, "PSDR")* care, încercînd să obțină un *lock* S, ea va fi blocată de o altă tranzacție care a executat *DeleteDir(dir, "PNL")* și care deține un *lock* X asupra obiectului director. Însă operația *LookupDir(dir, "PSDR")* nu depinde în nici un caz de execuția operației *DeleteDir(dir, "PNL")* care poate fi mai tîrziu și abortată, astfel încît blocarea tranzacției nu este necesară. De aceea, tranzacția care deține *lock*-ul X, poate fi, din acest punct de vedere o tranzacție oricît de lungă. Această situație este un exemplu în care gradul de concurență este redus nemotivat; cauza este lipsa de informații semantice despre tipul director. Prin utilizarea de informații semantice, blocajul nemotivat poate fi înlăturat. Următorul set de reguli de sincronizare care utilizează informații mai detaliate despre operațiile relative la un director pot fi folosite pentru a specifica modul în care multiple tranzacții pot fi intercalate fără a viola seriabilitatea și fără a permite abort în cascadă:

1. Operațiile *InsertDir* sau *DeleteDir* care specifică chei diferite pot fi executate în orice ordine.
2. Operațiile *InsertDir*, *DeleteDir*, sau *LookupDir* care specifică aceeași cheie ca alte operații *InsertDir* sau *DeleteDir* executate anterior în tranzații neterminate, trebuie blocate pînă la terminarea acelor tranzații.
3. Operațiile *LookupDir* pot fi executate în orice ordine.
4. Operațiile *InsertDir* care specifică aceeași cheie ca și operații anterioare *LookupDir* nereușite executate de alte tranzații neterminate trebuie blocate pînă acele tranzații sînt terminate.
5. Operațiile *DeleteDir* care specifică aceeași cheie ca și operații anterioare *LookupDir* reușite executate de alte tranzații neterminate trebuie blocate pînă acele tranzații sînt terminate.
6. Operațiile *DumpDir* trebuie blocate dacă au fost executate operații *InsertDir* sau *DeleteDir* de alte tranzații neterminate.
7. Operațiile *InsertDir* sau *DeleteDir* trebuie blocate dacă operații *DumpDir* au fost executate în alte tranzații neterminate.

Aceste specificații pot fi precizate utilizînd relații de dependență și ordonabilitatea. Prima etapă este identificarea dependențelor pentru acest tip. Ca și regulile de sincronizare ele se bazează pe operațiile definite pentru acel tip. Pentru a păstra numărul de dependențe la minimum operațiile pentru tipul director pot fi împărțite în trei grupuri:

1. Acele operații care modifică o intrare particulară într-un director, identificat de cheia α . Operațiile *InsertDir* și *DeleteDir* terminate cu succes, sînt în această clasă. Ele sînt denumite operații *Modify* (M).
2. Acele operații care observă prezența, absența sau conținutul unei intrări particulare într-un director, identificată printr-o cheie α . Operațiile *LookupDir* ca și operațiile *InsertDir* și *DeleteDir* nereușite sînt în această clasă. Această clasă este clasa operațiilor *Lookup* (L).

3. Acele operații care observă proprietăți ale directorului care nu pot fi izolate într-o intrare individualizată. Singura operație din această clasa este *DumpDir* (dintre cele definite; mai pot fi date și alte exemple, astfel poate fi definită o operație care returnează numărul de intrări din director). Această clasă este clasa operațiilor *Dump* (D).

Suplimentar definiții operațiilor ca atare, regulile de sincronizare iau în considerare și argumentele furnizate de operații. În următorul set complet de dependențe, simbolurile α și α' reprezintă chei distincte în operații cu directoare:

1. $D_1: T_i: M(\alpha) \rightarrow T_j: M(\alpha')$; T_i modifică o intrare cu cheia α și T_j modifică ulterior o intrare cu o cheie diferită α' .
2. $D_2: T_i: M(\alpha) \rightarrow T_j: M(\alpha)$; T_i modifică o intrare cu cheia α și T_j modifică ulterior aceeași intrare.
3. $D_3: T_i: M(\alpha) \rightarrow T_j: L(\alpha')$; T_i modifică o intrare cu cheia α și T_j observă ulterior o intrare cu o cheie diferită α' .
4. $D_4: T_i: M(\alpha) \rightarrow T_j: L(\alpha)$; T_i modifică o intrare cu cheia α și T_j observă ulterior aceeași intrare.
5. $D_5: T_i: L(\alpha) \rightarrow T_j: L(\alpha')$; T_i observă o intrare cu cheia α și T_j observă ulterior o intrare cu o cheie diferită α' .
6. $D_6: T_i: L(\alpha) \rightarrow T_j: L(\alpha)$; T_i observă o intrare cu cheia α și T_j observă ulterior aceeași intrare cu cheia α .
7. $D_7: T_i: L(\alpha) \rightarrow T_j: M(\alpha')$; T_i observă o intrare cu cheia α și T_j modifică ulterior o altă intrare cu cheia α' .
8. $D_8: T_i: L(\alpha) \rightarrow T_j: M(\alpha)$; T_i observă o intrare cu cheia α și T_j modifică ulterior aceeași intrare.
9. $D_9: T_i: D \rightarrow T_j: M(\alpha)$; T_i observă întregul conținut al directorului și T_j modifică ulterior o intrare cu cheia α .
10. $D_{10}: T_i: D \rightarrow T_j: L(\alpha)$; T_i observă întregul conținut al directorului și T_j observă o intrare cu cheia α .
11. $D_{11}: T_i: M(\alpha) \rightarrow T_j: D$; T_i modifică o intrare cu cheia α și T_j observă ulterior întregul conținut al directorului.
12. $D_{12}: T_i: L(\alpha) \rightarrow T_j: D$; T_i observă o intrare cu cheia α și T_j observă ulterior întregul conținut al directorului.
13. $D_{13}: T_i: D \rightarrow T_j: D$; T_i observă întregul conținut al directorului și T_j ulterior observă din nou întregul conținut al directorului.

Ca și dependența $R \rightarrow R$, multe din aceste dependențe nu trebuie să afecteze tranzațiile ulterioare. De aceea, ele nu trebuie luate în considerare la definirea condițiilor necesare pentru consistență. Dependențele următoare sînt cele care nu

trebuie luate în considerare:

1. Acelea pentru care nici o operație din dependență nu modifică obiectul (directorul): $D_6, D_{10}, D_{12}, D_{13}$. Acestea sînt analoge cu dependența $R \rightarrow R$.
 2. Acelea pentru care cele două operații din dependență referă două chei diferite: D_1, D_3, D_5, D_7 .
- Rezultă, în concluzie, că un abstract pentru tipul director este consistent dacă este ordonabil referitor la $\{<_{D_2 \cup D_4} \cup D_8 \cup D_9 \cup D_{11}\}$.

B. Liste FIFO

Specificații similare cu cele prezentate mai sus există și pentru alte tipuri de date.

Listele FIFO sînt un exemplu mai des întîlnit în literatură ([AC93][CG87][SP84][Mus94d]).

Luăm în considerare două operații asupra listei FIFO:

1. $QEnq(queue, ptr)$: adaugă o intrare la sfîrșitul cozii.
2. $QDeq(que)$: citește și elimină o intrare de la începutul cozii. Dacă lista este vidă, operația este blocată pînă cînd lista devine nevidă.

În scopul serializării tranzacțiilor care conțin operații asupra listelor FIFO, numeroase proprietăți trebuie garantate.

De exemplu:

1. Dacă o tranzacție adaugă cîteva intrări într-o listă FIFO, aceste intrări trebuie să apară împreună în aceeași ordine la sfîrșitul listei.
2. Orice intrări adăugate într-o listă coadă în timpul unei tranzacții nu pot fi observate de alte tranzacții decît după ce prima tranzacție se termină cu succes.
3. Dacă două tranzacții creează fiecare intrări în două liste, ordinea relativă a intrărilor create de cele două tranzacții trebuie să se regăsească și în cele două liste.

Aceste proprietăți, prea restrictive în ceea ce privește concurența, pot fi modificate pentru a permite execuția concurrentă a mai multor operații. De exemplu, dacă operațiile $QEnq$ din mai multe tranzacții sînt intercalate, atunci intrările create de fiecare tranzacție nu vor apare în bloc la începutul listei FIFO.

Analog dacă se permite execuția unei operații $QDeq$ într-o tranzacție urmată de o altă tranzacție neterminată care execută deasemenea o operație $QDeq$, un abort al unei tranzacții va interschimba pozițiile celor două intrări.

Setul de planificatoare abstracte consistente pentru liste FIFO serializabile poate fi specificat în termeni de dependență și ordonabilitate. La fel ca și în cazul directoarelor, este necesar să se distingă individual elemente în lista FIFO. Se presupune că, în momentul execuției operației $QEnq$ fiecare element are asociat un identificator unic. În continuare, simbolurile α și α' se vor utiliza pentru a simboliza identificatori distincți pentru elemente diferite, iar operațiile $QEnq$ și $QDeq$ vor fi prescurtate cu E respectiv D. Setul complet de dependențe pentru liste FIFO este:

1. $D_1: T_1: E(\alpha) \rightarrow_Q T_1: E(\alpha')$; T_1 introduce un element α' în coada Q după ce T_1 a introdus anterior elementul α .
2. $D_2: T_1: E(\alpha) \rightarrow_Q T_1: D(\alpha')$; T_1 elimină elementul α' după ce T_1 a introdus elementul α .
3. $D_3: T_1: E(\alpha) \rightarrow_Q T_1: D(\alpha)$; T_1 elimină elementul α care a fost introdus anterior de T_1 .
4. $D_4: T_1: D(\alpha) \rightarrow_Q T_1: E(\alpha')$; T_1 introduce elementul α' după ce T_1 a eliminat elementul α .
5. $D_5: T_1: D(\alpha) \rightarrow_Q T_1: D(\alpha')$; T_1 elimină elementul α' după ce T_1 a eliminat elementul α .

Folosind schema de sincronizare *read/write*, $QEnq$ poate fi modelată ca o operație de scriere în timp ce $QDeq$ poate fi modelată ca o operație de citire urmată de o operație de scriere. Atunci, conform schemei de sincronizare *read/write*, trebuie prevenită formarea de cicluri cuprinzînd relația de dependență $<_{R \rightarrow W \cup W \rightarrow R \cup W \rightarrow W}$. Totuși, această soluție este restrictivă într-un grad mai mare decît necesar. Considerăm, de exemplu dependența D_2 , care se formează atunci cînd o tranzacție elimină un element după ce o altă tranzacție a introdus anterior un alt element sau cazul D_4 , care este de fapt cazul invers lui D_2 . Concurența poate fi mărită dacă nu se includ în mecanismul de control aceste dependențe. Prin urmare, pentru crearea unei liste FIFO stricte trebuie garantată ordonabilitatea cu referire la relația $<_{D_1 \cup D_3 \cup D_5}$ a planificărilor abstracte, dar pot fi admise cicluri cu dependențe din relația D_2 sau D_4 . De exemplu, considerăm următoarea planificare, în care operează două tranzacții într-o listă FIFO care conține inițial elementele $\{A, B\}$:

- t1: $<T_1, QEnq(Q, S), X>$
 t2: $<T_2, QDeq(Q, S), A>$ returnează X
 t3: $<T_1, QEnq(Q, S), Y>$

La pasul t2 al planificării, se formează o dependență $T1 <_{D_2} T2$, iar la pasul t3, dependența $T2 <_{D_4} T1$. Astfel se formează un ciclu în relația compusă $<_{D_2 \cup D_4}$, dar planificarea este ordonabilă în sensul de mai sus.

Relațiile de dependență pot fi folosite pentru a planificări și din punct de vedere al situației de abort în cascadă. Relația de dependență $<_{D_1}$ este similară dependenței $W \rightarrow W$ și deci nu poate provoca abort în cascadă. Relațiile $<_{D_3}$ și $<_{D_5}$

sînt similare dependențelor $W \rightarrow R$. Într-o dependență D3, informația este transferată în listă (sub forma unui element α) între tranzații; este evident că o astfel de dependență poate provoca abort în cascadă. O dependență D5 poate, similar, să cauzeze abort în cascadă deoarece eliminarea unui element de către prima tranzație afectează ce element este primit de a doua tranzație.

Deși această definiție pentru consistența unei liste FIFO este o schemă *read/write* îmbunătățită, totuși restrînge în mod nenecesar concurența. Ea permite la cel mult două tranzații să acceseze concurrent o listă FIFO, una executînd operații *QEnq*, iar cealaltă operații *QDeq*. Spre deosebire de director, o listă FIFO trebuie să păstreze și ordinea elementelor pe care le conține. Un sistem bazat pe tranzații serializabile garantează faptul că tranzațiile pot fi plasate în orice ordine; prin forțarea unei anumite ordini a elementelor, tipuri de date precum liste FIFO sau LIFO restricționează concurența.

C. Semiliste FIFO

Este posibil să se restrîngă cerințele de integritate pe care programatorul le poate impune asupra unui tip, garantîndu-se consistența prin proprietăți de consistență mai puțin ferme decît seriabilitatea. Acest lucru permite mărirea gradului de concurență și este demonstrat în continuare printr-un tip special de listă FIFO.

Unul din modulele obișnuite de lucru cu o listă FIFO este folosirea ei ca buffer între un producător și un consumator. În mod uzual, ordinea exactă a intrărilor din listă nu este importantă. Ceea ce este esențial este ca intrările introduse la un capăt al cozii să ajungă garantat în celălalt capăt într-un timp aproximativ egal cu alte intrări create în aproximativ același timp. Tipul de listă FIFO cu această proprietate este numit *semilistă FIFO* (WFIFO). Intrările nu pot rămîne la nesfîrșit în listă dacă:

1. Tranzația care le introduce se termină într-un interval de timp finit.
2. Tranzația care le elimină se termină într-un interval de timp finit.
3. Numai un număr finit de tranzații elimină o intrare și apoi se termină cu abort.

Operațiile permise într-o semilistă sînt aceleași cu cele definite pentru liste FIFO adică:

1. *WQEnq(queue, ptr)*: adaugă o intrare la sfîrșitul listei.

2. *WQDeq(queue)*: citește și elimină o intrare de la începutul listei. Dacă lista este vidă, operația este blocată pînă cînd lista devine nevidă.

Regulile de sincronizare sînt însă mai puțin restrictive decît cele pentru liste FIFO. Ele sînt:

1. Operațiile *WQEnq* pot fi executate chiar dacă o altă tranzație neterminată a executat alte operații *WQEnq*.
2. Operațiile *WQDeq* pot fi executate chiar dacă o altă tranzație neterminată a executat *WQDeq*.
3. O operație *WQDeq* se aplică intrării de la începutul listei WFIFO; tranzația care a referit această intrare trebuie să se fi terminat.

În exemplele care urmează, o listă WFIFO va reprezentată printr-o secvență de litere, cea mai din stînga reprezentînd începutul listei. Literele mici italice simbolizează intrări pentru care operațiile corespondente *WQEnq* aparțin unor tranzații încă neterminate. Literele mari simbolizează intrări care nu au fost eliminate și pentru care tranzațiile care le-au introdus au fost terminate. Literele mari italice simbolizează intrări care au fost eliminate prin operații *WQDeq* aparținînd unor tranzații încă neterminate. Cifrele asociate intrărilor simbolizate prin litere italice indică numărul tranzației care a executat acea operație.

Presupunem că inițial lista este vidă. Dacă tranzația T_1 execută *WQEnq(WQ,c)* atunci lista conține:

$\{a^1\}$

Conform regulii 1, o altă tranzație T_2 este liberă să creeze o altă intrare:

$\{a^1, b^2\}$

În continuare, T_1 poate să adauge o nouă intrare:

$\{a^1, b^2, c^1\}$

Dacă ambele T_1 și T_2 se termină normal, starea finală devine:

$\{A, B, C\}$

Presupunem acum că T_3 execută *WQDeq*. Conform regulii 2, o altă tranzație T_4 este liberă pentru a elimina două elemente:

$\{A^3, B^4, C^1\}$

Dacă în acest moment T_4 se termină normal iar T_3 execută *abort*, starea finală a listei devine:

$\{A\}$

Se observă că, în acest caz, elementele A și C s-au inversat, chiar dacă au fost inițial inserate de aceeași tranzație. Pe de altă parte, dacă T_3 nu ar fi executat *abort* ci *commit*, deci în condiții normale, comportarea este asemănătoare unei liste FIFO.

Un alt exemplu ilustrează situația cînd o tranzație neterminată normal ajunge la începutul listei WFIFO.

Presupunem că starea inițială este:

$\{a^2, b^1\}$

Dacă T_6 execută *commit*, dar T_5 rămîne încă neterminată starea devine:

$\{a^5, B\}$

Dacă T_7 elimină un element în acest moment, va fi returnat B, starea devenind:

$\{a^5\}$

după ce T_7 execută *commit*. Pe de altă parte, dacă T_5 execută *commit* după T_6 , dar înainte de eliminarea din T_7 , va fi returnat A. Din nou se constată că lista WFIFO încearcă să emuleze o listă FIFO obișnuită; în acest caz se permite ca intrări din lista WFIFO, aparținând unor tranzații neterminate să fie rezervate pînă la terminarea tranzațiilor cu *commit*.

Proprietățile de consistență ale unei liste WFIFO pot fi descrise succint prin relații de dependență. Dependențele sînt aceleași ca și pentru tipul listă FIFO. Totuși, spre deosebire de lista FIFO strictă care pretinde ca planificatoarele abstracte să fie ordonabile ca referire la relația $\{<_{D1} \cup D3 \cup D5\}$, lista WFIFO permite formarea de cicluri cu toate relațiile de dependență exceptînd relația $<_{D3}$. Permițînd formarea de cicluri în relația $<_{D1}$, se permite astfel intercalarea creerii de intrări de către multiple tranzații. Similar, permițînd formarea de cicluri în relația $<_{D5}$, se permite intercalarea ștergerilor de intrări de către mai multe tranzații.

Pentru a concluziona comparația dintre listele FIFO și WFIFO, trebuie evidențiat faptul că în listele WFIFO au fost sacrificate două proprietăți ale listelor FIFO. În primul rînd, nu se respectă ordinea strictă a intrărilor listei. În al doilea rînd, tranzațiile care operează asupra listelor WFIFO nu mai sînt serializabile. În schimb, sînt menținute alte proprietăți esențiale. De exemplu, o listă WFIFO nu va amîna nedefinit "consumul" unei intrări față de "producerea" ei. În plus, lista WFIFO forțează ordonarea acelor tranzații care operează printr-un element comun al listei. Acest lucru se asigură prin ordonabilitatea față de $\{<_{D3}\}$. Toate aceste modificări măresc într-un grad mare concurența, fără a compromite utilitatea specifică a tipului de date.

Pînă acum s-au considerat doar două operații posibile în lucrul cu liste de tip coadă (FIFO sau WFIFO). Există însă și alte operații, dintre care unele pot restrînge drastic gradul de concurență. Un astfel de exemplu este operația *QEmpty(q)*. Această operație restrînge mult gradul de concurență deoarece informația care trebuie returnată nu poate fi izolată într-o intrare particulară a listei. De aceea, ca și operația *DumpDir*, operația *QEmpty(q)* blochează toate operațiile ulterioare asupra cozii care modifică informația observată. Analog, această operație va fi întotdeauna blocată ori de cîte ori există un dubiu în rezultatul său (în cazuri de inserări și eliminări de elemente în tranzații neterminate). Nu întotdeauna serializarea tranzațiilor care apelează operația de mai sus este necesară. Să presupunem, de exemplu, că o operație *QEmpty* returnează fals, deci există cel puțin o intrare în listă. Dacă următoarea operație din tranzație este *QDeq*, poate surveni o eroare "listă vidă", dacă tranzația nu este serializată. Dacă, dimpotrivă, *QEmpty* returnează true, nu este nevoie de serializarea tranzației. Rezultă că o tranzație trebuie sau nu serializată în funcție de rezultatul apelului operației *QEmpty*.

D. Tranzații care conțin operații asupra mai multor tipuri

Considerațiile prezentate în subcapitolul anterior relativ la specificațiile asupra unor tipuri abstracte pot fi extinse pentru a descrie modul în care proprietățile de consistență a mai multor tipuri interacționează în cazul în care într-o tranzație sînt implicate mai multe tipuri.

Trebuie generalizată definiția de consistență aplicată asupra planificărilor abstracte conținînd operații asupra unor tipuri multiple. Presupunem că specificațiile proprietăților de consistență pentru tipul Y_1 garantează ordonabilitatea față de $\{<_{D1}\}$, pentru tipul Y_2 garantează ordonabilitatea față de $\{<_{D2}\}$, etc.

Definiție.

*Setul de planificatoare abstracte consistente față de tipurile Y_1, Y_2, \dots, Y_n se definește ca setul de planificatoare ordonabile față de $\{<_{D1} \cup D2 \cup \dots \cup Dn\}$, deci față de reuniunea dependențelor prescise pentru fiecare tip. Un set de tipuri care satisface aceste proprietăți este denumit un *set de tipuri cooperative*.*

Exemplul va demonstra insuficiența ordonabilității față de dependențele prescise pentru fiecare tip în parte. Considerăm următoarea planificare conținînd tranzații care execută operații asupra tipurilor discutate anterior (directoare și liste FIFO).

```
<T1, QEnq(Q,T), Y>
<T2, QEnq(Q,T), Z>
<T2, InsertDir(D,S), "PDSR">
<T1, DeleteDir(D,S), "CDR">
```

Fie $<_{Dir}$ notația pentru relația $<_{D2 \cup D4 \cup D8 \cup D9 \cup D11}$ (definită anterior pentru tipul director) și $<_Q$ notația pentru $<_{D1 \cup D3 \cup D5}$ (definită anterior pentru tipul listă FIFO). Deși planificarea este ordonabilă față de relația $\{<_{Dir}, <_Q\}$ nu este ordonabilă față de relația $\{<_{Dir} \cup <_Q\}$. Pentru a realiza ordonabilitatea, tipurile director și listă FIFO trebuie să coopereze în prevenirea de cicluri în relația compusă.

7.4.9.3. Metode de sincronizare

Pentru demonstrarea corectitudinii implementării unui tip trebuie efectuate următoarele etape: (1). caracterizarea setului de planificări de apel legale, adică acele planificări permise de mecanismul de sincronizare utilizat în implementare; (2). realizarea unei corespondențe între planificarea de apel și cea abstractă; (3). demonstrarea corectitudinii corespondenței planificării de apel-planificării abstracte.

Cînd nu se utilizează reordonarea operațiilor într-o implementare a unei planificări de apel (cazul ideal) atunci planificarea de apel este echivalentă cu planificarea abstractă; în acest caz etapa a doua și a treia nu mai sînt necesare.

A. Lock-uri specifice tipurilor abstracte

Mecanismul de sincronizare propus în acest subcapitol se bazează pe mecanismul tradițional de blocare pe bază de lock prezentat în subcapitolul 7.4.4. Există mai multe variante ale acestui mecanism, toate bazîndu-se pe același principiu: înainte de a opera asupra unui obiect, orice tranzație trebuie să obțină un lock pentru obiect; alte tranzații nu vor putea accesa obiectul decît după eliberarea lock-ului obținut de tranzația respectivă.

Mecanismul de blocare restricționează formarea dependențelor între tranzații, restrîngînd setul de planificatoare de apel legale. Ori de cîte ori o tranzație este forțată de către o alta să aștepte la un lock, formarea unei dependențe între cele două tranzații este întîrziată pînă cînd prima tranzație eliberează lock-ul. Protocolul de blocare în două faze, cel mai cunoscut, impune ca nici o tranzație care eliberează un lock să nu mai pretindă un alt lock. Acest lucru are ca efect conversia eventualelor cicluri de dependență în interblocări. Acestea pot fi detectate și pentru că nu sînt permise astfel de dependențe, fiecare tranzație poate fi terminată cu abort fără a afecta altă tranzație. S-ar putea afirma că protocolul blocării în două faze aplică o politică conservatoare; într-adevăr el întîrzie formarea de orice dependențe dintr-o relație prescrisă, nu numai cele care conduc la cicluri. Totuși, acesta nu este un dezavantaj propriu-zis, deoarece formarea unor astfel de dependențe cu transfer de informație trebuie întîrziată datorită pericolului unui abort în cascadă.

Dezavantajul cel mai mare al mecanismelor lock tradiționale este acela că utilizează minimum de informații semantice asupra obiectelor manipulate. Cele mai simple scheme lock utilizează un singur tip de *lock*; schemele *lock_S/lock_X* utilizează minimum de informație, dar sînt totuși insuficiente pentru a ține cont de avantajele dependențelor specifice tipurilor.

Două observații pot fi făcute în ceea ce privește dependențele specifice tipurilor. În primul rînd, ele specifică modul în care operațiile specifice tipurilor, incluse în tranzații, pot fi intercalate. Analog, schema de blocare necesită specificarea claselor *lock*-urilor care corespund operațiilor asupra acestui tip. În al doilea rînd, dependențele reflectă datele furnizate ca argumente de către operații. De aceea, un obiect *lock* în schema de blocare generalizată constă în două părți: tipul specific clasei lock-ului și date specifice instanțierii lock-ului. Notăția *{LockClass (date)}* este utilizată pentru a reprezenta o instanțiere a unui lock.

După definirea claselor *lock*-urilor pentru un tip definit trebuie furnizată o funcție booleană care să specifice dacă un nou *lock* (de un anumit tip) este disponibil în funcție de *lock*-urile deja deținute de obiect. Pentru similitudine cu notațiile utilizate în 7.4.4., această funcție va fi reprezentată printr-o *tabelă de compatibilitate* a *lock*-ului. Numai *lock*-uri cerute de alte tranzații trebuie verificate pentru compatibilitate. O cerere a unui *lock* este compatibilă întotdeauna cu *lock*-urile deținute de aceeași tranzație.

B. Blocarea cu lock-uri în cazul directoarelor

Pentru ilustrarea mecanismului de blocare în cazul tipului director se consideră o implementare ideală (care nu reordonează operațiile) a tipului specificat în paragraful 7.4.9.2. Se presupune, de asemenea, că se asigură sincronizarea la nivelul operațiilor asupra tipului director (printr-un mecanism oarecare de excludere reciprocă) deoarece mecanismul de blocare se utilizează numai pentru intercalarea corectă a operațiilor asupra directoarelor din diverse tranzații. Mecanismul de blocare și de excludere reciprocă nu pot fi totuși complet independente, deoarece excluderea reciprocă trebuie eliminată cînd se așteaptă un lock.

Deoarece corespondența între planificările de apel și planificările abstracte este imediată pentru implementarea ideală, etapa a doua nu este necesară. Interesează numai prima și a treia etapă: caracterizarea setului de planificări legale și comparația acestora cu setul de planificări consistente.

Se lucrează cu aceleași tipuri de operații asupra directoarelor, care pot fi împărțite în trei grupuri:

1. operații de *modificare*, care alterează o intrare din director, identificată prin cheia (string) α .
 2. operații de *căutare*, care detectează prezența, absența sau conținutul unei intrări particulare din director, identificată prin cheia α .
 3. operații de *parcursere*, care observă proprietățile unui director, care nu pot fi izolate într-o intrare particulară.
- Corespunzător acestor grupuri, pot fi definite trei clase de lock:

1. {*ModifyDir*(α)}; pentru a indica faptul că o tranzacție incompletă a inserat sau șters o intrare cu cheia α .
2. {*LookupDir*(α)}; pentru a indica faptul că o tranzacție incompletă observă o intrare cu cheia α .
3. {*DumpDir*}; pentru a indica faptul că o tranzacție incompletă a executat o operație *DumpDir* de listare a întregului director.

Tabela de compatibilitate pentru lock-uri poate fi reprezentată astfel:

Lock Reținut				
Cerere lock	Nimic	ModifyDir (α)	LookupDir (α)	DumpDir
ModifyDir (α)	DA	Nu	Nu	Nu
ModifyDir (α')	DA	DA	DA	Nu
LookupDir (α)	DA	Nu	DA	DA
LookupDir (α')	DA	DA	DA	DA
DumpDir	DA	Nu	DA	DA

Fiecare intrare din această tabelă reflectă specificul relațiilor de dependență pentru directoare. Intrările compatibile reprezintă relațiile de dependență pentru care este permisă formarea de cicluri: de exemplu, pentru linia 2, coloana 2 valoarea este OK deoarece sînt permise cicluri în relația de dependență $\langle M(\alpha) \rightarrow M(\alpha) \rangle$.

Protocolul utilizat de operațiile asupra directoarelor pentru alocarea și eliberarea lock-urilor este:

1. Operațiile *InsertDir* și *DeleteDir* care specifică cheia α obțin un lock {*ModifyDir*(α)}. Dacă operația reușește, lock-ul este reținut alocat pînă la sfîrșitul tranzației. Dacă operația nu reușește, lock-ul este convertit într-un lock {*Lookup*(α)}, care este reținut pînă la sfîrșitul tranzației. Această conversie reprezintă o violare a protocolului strict 2PL și care demonstrează cum se poate îmbunătăți gradul de concurență pentru un protocol specific unui tip.
2. Operațiile *Lookup* care specifică o cheie α obțin un lock {*Lookup*(α)} pentru director pînă la sfîrșitul tranzației.
3. Operațiile *DumpDir* obțin un lock {*DumpDir*} asupra directorului pînă la sfîrșitul tranzației.

C. Blocarea cu lock-uri în cazul listelor

Se prezintă în acest subcapitol specificul folosirii lock-urilor pentru liste FIFO, liste WFIFO, și liste cu buffer local.

1. Liste FIFO

Ca și în exemplul anterior, se presupune o implementare ideală a tipului FIFO inclusiv respectîndu-se excluderea reciprocă la nivelul operațiilor. Deoarece în cazul acestui tip de date există două operații, *QDeq* și *QEnq*, vor fi definite două clase de lock: {*QEnq*(α)} și {*QDeq*(α)}. Ca și în cazul directoarelor, lock-urile trebuie să identifice intrarea particulară din listă. Deoarece intrările unei liste FIFO nu se identifică printr-o cheie, se va presupune în continuare că la momentul efectuării unei operații *QEnq*, fiecare element din listă contribuie la identificarea ei printr-un identificator unic. Acești identificatori corespund celor folosiți în relațiile de dependență. Astfel, un lock {*QEnq*(α)} indică faptul că un element α a fost inserat în listă de o tranzacție incompletă. Analog, un lock {*QDeq*(α)} indică faptul că un element α a fost eliminat din listă. Protocolul folosit în execuția operațiilor asupra unor liste FIFO este:

1. Operațiile *QEnq* trebuie să obțină un lock {*QEnq*(α)}, în care α este identificatorul intrării ce trebuie adăugată. Acest lock este reținut alocat pînă la sfîrșitul tranzației.
2. Operațiile *QDeq* trebuie să obțină un lock {*QDeq*(α)}, în care α este intrarea de la începutul listei. Acest lock este reținut alocat pînă la sfîrșitul tranzației.

Tabelul următor arată compatibilitatea lock-urilor pentru liste FIFO:

Lock Reținut			
Cerere Lock	Nimic	QEnq (α)	QDeq (α)
QEnq (α)	DA	-	-
QEnq (α')	DA	NU	DA
QDeq (α)	DA	NU	-
QDeq (α')	DA	DA	NU

Simbolurile α și α' reprezintă identificatorii a două elemente diferite. Incompatibilitatea dintre operațiile $\{QDeq(\alpha)\}$ și $\{QEnq(\alpha)\}$ asigură interzicerea eliminării intrărilor aparținând unor tranzații care nu au executat încă *commit*, înlăturându-se astfel situații de abort în cascadă.

2. Liste WFIFO

Se folosește aceeași implementare idealizată a unei liste FIFO și aceleași operații $QEnq$ și $QDeq$. Diferența față de tipul precedent constă în tabela de compatibilitate a lock-urilor:

Lock cerut	Nimic	Lock reținut	
		WQEnq (α)	WQDeq (α)
WQEnq (α)	DA	-	-
WQEnq (α')	DA	DA	DA
WQDeq (α)	DA	NU	-
WQDeq (α')	DA	DA	DA

Posibilitatea intercalării operațiilor $WQEnq$ de către diferite tranzații se reflectă în tabelă prin definirea lock-urilor $\{WQEnq(\alpha)\}$ și $\{WQEnq(\alpha')\}$ drept compatibile, ca și a lock-urilor $\{WQDeq(\alpha)\}$ și $\{WQDeq(\alpha')\}$. Singura restricție rămâne linia 3, coloana 2, în scopul prevenirii eliminării intrărilor din tranzații încă neterminate (care nu au executat *commit*). Aceasta previne formarea de cicluri în relația de dependență prescrisă $<_{E(\alpha) \rightarrow D(\alpha)}$ și previne formarea de abort în cascadă.

Protocolul de blocare pentru operațiile asupra unei liste WFIFO este același ca și în cazul unei liste FIFO. Singura diferență este aceea că dacă operația $WQDeq$ nu poate obține un lock $\{WQDeq(\alpha)\}$ pentru elementul de la începutul cozii, nu blochează tranzația ci caută în continuarea listei până la obținerea unui lock $\{WQDeq(\alpha')\}$ asupra altui element. Aceasta reflectă proprietatea listelor WFIFO de a permite eliminarea de elemente "din adâncul" listei, în cazul în care elementul de început este implicat într-o tranzație ce nu a executat *commit*.

3. Liste FIFO cu buffer local

Considerăm acum o implementare neidealizată a unei liste FIFO, aceea prezentată anterior, în care elementele de inserat în listă se depozitează într-un buffer local care se inserează în listă la terminarea tranzației prin *commit*. Pentru a distinge între cele două implementări se va nota implementarea folosind un buffer ca listă *CFIFO*. Pentru listele CFIFO se definesc aceleași operații $CQEnq$ și $CQDeq$, dar, în plus, și o operație $CQCommit(cq, tid)$ apelată de managerul tranzației pentru a salva bufferul în momentul când tranzația execută *commit*.

Lista CFIFO poate fi implementată utilizând o listă FIFO globală, nebufferată, *GlobalQueue*, ca o componentă a implementării și liste locale tranzațiilor asupra cărora se acționează prin operații $CQxxx$, în regim de excludere reciprocă, dar fără mecanism de blocare, listele nefiind partajate. Operația $CQCommit(cq, tid)$ apelează $QEnq(GlobalQueue, ptr)$ pentru a plasa fiecare din intrările listei locale în lista FIFO globală.

Primul apel $QEnq$ trebuie să poată obține un lock $\{QEnq(\alpha_i)\}$ asupra listei FIFO globale; alocarea acestui lock permite unei tranzații să obțină acces exclusiv la lista globală, deoarece lock-ul $\{QEnq(\alpha_i)\}$ necesar pentru fiecare intrare particulară de inserat este compatibil cu cel precedent.

Implementarea CFIFO permite planificări de apel în care operațiile $CQEnq$ din diferite tranzații pot fi intercalate. Amînarea operațiilor $QEnq$ asupra listei FIFO globale până la execuția operației *commit* transformă planificările de apel în planificări abstracte fără intercalare de operații.

Dacă lista CFIFO globală nu este vidă, prima dintre intrări trebuie eliminată și returnată. Pentru a preveni cicluri în relația $<_{E(\alpha) \rightarrow D(\alpha)}$ este necesară blocarea, dacă intrarea este reținută de o tranzație care nu a executat încă *commit*. Dacă lista *GlobalQueue* este vidă, orice intrări din bufferul local trebuie transferate în *GlobalQueue* printr-un proces identic cu cel din $CQCommit$. Prima intrare din bufferul local de intrări al tranzației poate apoi fi eliminată din *GlobalQueue* și returnată. Deplasarea intrărilor buferate în bloc în *GlobalQueue* interzice formarea de cicluri în relațiile prescrise $<_{E(\alpha) \rightarrow E(\alpha)}$ și $<_{D(\alpha) \rightarrow D(\alpha')}$. Un al treilea caz apare atunci când și lista *GlobalQueue* și bufferul sînt vide; în acest caz operația trebuie blocată pînă cînd o altă tranzație inserează o intrare în *GlobalQueue* și execută *commit*.

Operațiile descrise nu sînt însă suficiente pentru a implementa $CQDeq$. Dificultatea constă în aceea că $QEnq(GlobalQueue)$ va fi blocată dacă lista globală este vidă, lucru care trebuie evitat pentru a beneficia de avantajele bufferului local. De aceea, mai trebuie introdusă o operație, $QDeqNE$ necesară în cazul prezentat mai sus. Dacă se poate elimina o intrare din lista globală, $QDeqNE$ alocă un lock $\{QDeq(\alpha_i)\}$ asemenea cu $QDeq$. În caz contrar și dacă există intrări în bufferul local, acestea sînt deplasate în lista globală, după care se apelează operația standard $QDeq$ pentru a elimina și returna prima intrare din lista globală. În acest caz, tranzația reține alocat un lock $\{QDeq(\alpha_i)\}$ asupra listei pînă la terminare.

Capitolul 8

CONTROLUL INTERBLOCĂRII OBIECTELOR

Problemele care se pun în legătură cu controlul interblocării în sistemele distribuite sunt asemănătoare cu cele pentru sistemele centralizate. În sistemele distribuite, însă, datorită dispersiei informațiilor pe mai multe noduri, interblocările sunt și mai greu de prevenit, evitat și detectat.

Se pot distinge două tipuri de interblocări în sistemele distribuite: interblocările rezultate ca urmare a utilizării resurselor în acces exclusiv și interblocările de comunicare, care apar când un set de procese așteaptă circular primirea unui mesaj. În secțiunile următoare se va trata primul caz de interblocări, considerând că, pe de o parte, al doilea caz poate fi inclus în primul dacă se consideră canalele de comunicație ca resurse și, pe de altă parte, se ține cont de aprecierea ([Tan93], [Tan96], [Gos91]) că o astfel de comunicație circulaa apare destul de rar în aplicații (paradigma de interacțiune a proceselor cea mai des întâlnită este cea client-server).

8.1. Metode de control al interblocării

În principiu, toate metodele de control al interblocării cunoscute în sistemele centralizate sunt valabile și pentru sistemele distribuite. Astfel, metoda prevenirii interblocării prin ordonarea resurselor poate fi utilizată prin definirea unei relații de ordine peste evenimentele sistemului. Se asignează tuturor resurselor sistemului distribuit numere unice iar un proces poate cere o resursă i , de pe orice procesor, dacă și numai dacă nu deține o resursă cu număr mai mare sau egal cu i . Similar, algoritmul bancherului poate fi utilizat într-un sistem distribuit dacă se desemnează unul dintre procesele din sistem ca proces central care menține toate informațiile necesare algoritmului; fiecare cerere de acces la resursă va fi adresată acestui proces. Totuși, în afară de faptul că, fiind centralizată, metoda poate conduce la sufocarea procesului central, ea prezintă același dezavantaj ca și în sistemele centralizate: impune cunoașterea în avans a resurselor utilizate, ceea ce, în sistemele distribuite este arareori posibil. Ca atare, în sistemele distribuite aceste metode nu se folosesc. Se folosesc, în schimb, metodele de prevenire a interblocării care se bazează pe tranzații atomice și pe ordonarea tranzațiilor prin asocierea de mărci de timp și metode de detecție a interblocării.

Utilizarea tranzațiilor atomice facilitează procedura de distrugere și reluare a proceselor implicate într-o interblocare; în cazul tranzațiilor atomice, acestea pot fi readuse în starea de înaintea începerii tranzației, pe când reluarea unor procese simple de la o stare anterioară necesită mai multe complicații.

8.2. Prevenirea interblocărilor

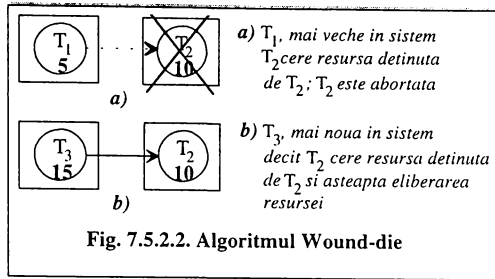
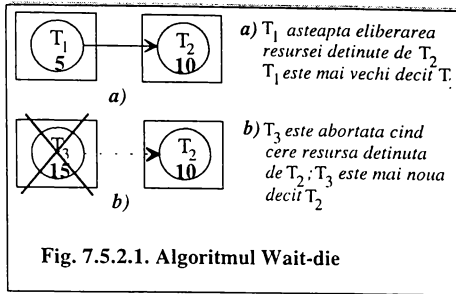
În sistemele distribuite în care sunt implementate tranzații atomice și timp global sunt folosiți în practică doi algoritmi pentru prevenirea interblocărilor. Ambii se bazează pe ideea asignării fiecărei tranzații, în momentul în care începe, a unei mărci de timp. Fiecărei tranzații i se mai asignează și o prioritate, care poate fi numărul de proces; aceasta urmează să fie folosită în cazul în care două tranzații au aceeași marcă de timp.

Cei doi algoritmi sunt:

(1) algoritmul *wait_die* care se bazează pe o tehnică nepreemptivă. Când o tranzație T_i cere accesul la o resursă deținută la momentul respectiv de T_j , T_i va aștepta alocarea resursei numai dacă are o marcă de timp mai mică decât T_j (aceasta înseamnă că T_i este o tranzație mai veche în sistem decât T_j); în caz contrar T_i este abortată. În figura 7.5.2.1 se prezintă un exemplu de funcționare al acestui algoritmu.

(2) algoritmul *wound_die* care se bazează pe o tehnică preemptivă. Când o tranzație T_i cere accesul la resursă deținută la momentul respectiv de T_j , T_i va aștepta alocarea resursei numai dacă are o marca de timp mai mare decât T_j (aceasta înseamnă că T_i este o tranzație mai nouă în sistem decât T_j); în caz contrar T_j este abortată, iar T_i continuată (fig 7.5.2.2).

Ambii algoritmi evită aminarea la infinit a tranzațiilor, cu condiția ca, la reluarea tranzațiilor după abortare, să se pastreze aceeași marcă de timp. Deoarece mărcile de timp au valori crescătoare, o tranzație abortată va deveni, la un moment, cea mai veche din sistem și va avea cea mai mică marca de timp din sistem. Există diferențe semnificative în ceea ce privește modul de operare al algoritmilor în acest sens. Astfel: (1) În algoritmul *wait_die*, o tranzație mai veche așteaptă ca o tranzație mai nouă să elibereze resurse, pe când în algoritmul *wound_die* o tranzație mai veche nu așteaptă niciodată la o tranzație mai nouă. (2) În algoritmul *wait_die*, o tranzație mai nouă este abortată când cere accesul la o resursă deținută de o tranzație mai veche. În mod normal, programul va încerca din nou realizarea tranzației și, din același motiv, ea poate fi abortată de mai multe ori; în contrast, în algoritmul *wound_die*, o tranzație nouă va aștepta eliberarea resursei de către tranzația mai veche. În schimb, algoritmul *wound_die* abortează o tranzație nouă când o alta mai veche cere aceeași resursă. Ca atare, se poate aprecia că algoritmul *wound_die* necesită, în general, mai puține operații abort.



8.3. Detecția interblocărilor în sisteme distribuite

Prevenirea interblocării poate genera abortarea unor tranzații care uneori nu generează interblocare. Pentru a înlătura acest dezavantaj, se folosesc *algoritmi de detecție* a interblocării. Detecția interblocării în sisteme distribuite constituie un subiect al multor cercetări, propunându-se diverși algoritmi, pentru a căror certitudine s-au făcut demonstrații, dar pentru care însă, ulterior s-au găsit și contraexemple. Corectitudinea unui algoritm de detecție depinde de două condiții: (1) *condiția de progres*: fiecare interblocare trebuie detectată; (2) *condiția de siguranță*: dacă s-a detectat o interblocare, atunci ea trebuie să existe într-adevăr. În secțiunea 7.5.3.3 se arată că, datorită întârzierilor în transmisia mesajelor, într-un sistem distribuit pot fi detectate interblocări false.

8.3.1 Modelul de lucru

Teoria detecției interblocărilor poate fi aplicată direct sistemelor de obiecte distribuite. Ca atare, în această secțiune va fi descris modelul de lucru folosit în secțiunile următoare și aplicat sistemelor de obiecte distribuite. Acest model constă într-un set de N noduri, S_1, S_2, \dots, S_N conectate într-o rețea; se consideră transmisia mesajelor ca sigură iar nodurile complet conectate. În fiecare nod există un sistem de obiecte, centralizat, component al sistemului de obiecte distribuite. Există M tranzații T_1, T_2, \dots, T_M în execuție în contul sistemelor de obiecte distribuite. O tranzație adresează cererile de execuție a unei operații asupra obiectului local (fișiere, înregistrări, bază de date locală) unui *manager central al tranzațiilor* sau *controler*. Există câte un controler C_i pentru fiecare nod S_i . O cerere de acces la obiect poate consta, de exemplu, într-o operație de scriere, citire sau blocare a obiectului. O tranzație este blocată din momentul în care adresează o cerere managerului central și pînă în momentul în care controlerul oferă accesul la obiect (cînd tranzația devine activă). Cererea se poate referi la un obiect local sau la un obiect situat la distanță, caz în care tranzația devine distribuită.

O tranzație distribuită T_i se implementează prin mai mulți agenți ai tranzației, t_{ij} , câte unul pentru fiecare nod S_j . Cînd o tranzație cere acces la un obiect situat la distanță, se creează un agent care are rolul de a realiza cererea la obiect, accesul la obiect și eliberarea obiectului; agentul este distrus cînd cererea la distanță a fost îndeplinită. Pentru simplitate, agenții tranzațiilor se vor putea nota și cu un singur indice, t_i semnificînd deci al i -lea agent.

În cazul în care agentul t_{ij} cere accesul la un obiect nelocal, gestionat de controlerul C_m , controlerul C_j transmite cererea la agentul t_{im} , via controlerul C_m . Cînd t_{im} obține accesul la obiect de la C_m trimite un mesaj la t_{ij} , anunțîndu-l că accesul a fost obținut. Cînd agentul tranzației T_i nu mai are nevoie de accesul la obiectul controlat de C_m , îi comunică agentului t_{im} care este responsabil pentru eliberarea resursei la C_m . Astfel toate cererile între noduri sunt transformate în cereri între agenți ai aceleiași tranzații.

Modelul matematic al cererilor de acces la obiecte ale tranzațiilor este un graf, *graful dependentelor tranzației* (TWFG, *transaction wait_for graph*). Nodurile acestui graf se asociază proceselor agent ale tranzațiilor. Arcele reprezintă relații de blocare între procesele agent. Un nod care are un arc ce iese din acel nod corespunde unui proces agent blocat. Mai precis, într-un graf TWFG există un arc de la agentul t_{ij} la agentul t_{im} dacă controlerul C_j are o cerere de la t_{ij} referitoare la o resursă deținută de t_{im} ; un astfel de arc se numește *arc intracontroler*. Un arc de la agentul t_{ij} la agentul t_{im} există dacă t_{ij} așteaptă un mesaj de permisiune acces la obiect de la t_{im} ; un astfel de arc se numește *arc intercontroler*.

O structură ciclică în acest graf indică o interblocare. Definiția precisă a termenului de *structură ciclică* depinde de modelul de interblocare considerat (vezi secțiunea 7.5.3.2.2). Un prim exemplu este prezentat în fig. 7.5.3.1. Există 5 tranzații, implementate prin 7 procese agent. Prezența ciclului $t_{11} \rightarrow t_{31} \rightarrow t_{32} \rightarrow t_{22} \rightarrow t_{21}$ indică o interblocare dacă tranzațiile T_2 și T_3 sunt necesare ambele obiecte (1 și 2) pentru a continua prelucrările; dacă, însă îi este necesar numai un singur obiect dintre cele două, nu apare interblocare.

Principala problemă în sistemele distribuite este construcția grafului global TWFG. O serie de algoritmi propun menținerea în fiecare nod a grafurilor TWFG locale. Pentru a se demonstra că nu apare interblocare, trebuie arătat că graful global obținut prin reuniunea tuturor grafurilor TWFG locale nu are structuri ciclice.

8.3.2 Clase de algoritmi pentru detecția interblocării

Există numeroase metode de a organiza graful TWFG într-un sistem distribuit. Clasificarea algoritmilor de detecție a interblocării se face în funcție de modul în care se construiește graful TWFG.

8.3.2.1. Clasificarea algoritmilor

Detecția interblocării poate fi realizată utilizând algoritmi cu caracter *centralizat, ierarhic și distribuit*. Algoritmii *centralizați* construiesc graful global TWFG prin reuniunea grafurilor locale. Acesta este construit și menținut de către un singur proces: *procesul coordonator* al detecției interblocărilor. Algoritmii centralizați prezintă dezavantajele oricărei metode centralizate: vulnerabilitatea procesului central la avarii și traficul ridicat de mesaje. Unui detector centralizat îi este necesar un interval de timp considerabil pentru colectarea și asamblarea grafurilor locale TWFG; ca atare un alt dezavantaj este și faptul că o interblocare rămâne nedetectată un interval de timp relativ mare.

O alternativă la menținerea tuturor informațiilor necesare detecției interblocării într-un singur proces este de a distribui aceste informații între mai multe procese și de a impune o organizare ierarhică a facilităților de detecție, obținându-se o altă categorie de algoritmi, *algoritmii ierarhici*.

O a treia categorie de algoritmi sunt *algoritmii complet distribuiți*, în cazul acestora problemele principale sunt menținerea consistenței informațiilor, pentru a nu se detecta interblocări false și pentru nu se detecta aceeași interblocare în mai multe sisteme. Costul unui algoritm distribuit se masoară în numărul de mesaje schimbate între nodurile implicate în interblocare și timpul de calcul.

8.3.2.2 Modele de interblocare

În funcție de aplicație, se pot permite diverse tipuri de cereri de operații asupra obiectelor. De exemplu, o tranzație poate cere o combinație de operații asupra mai multor obiecte (obiectul a și obiectul b sau obiectul c) sau poate cere accesul numai la un singur obiect. Aceasta distincție permite o clasificare a algoritmilor distribuți de detecție, conform complexității cererilor de acces pe care aceștia le permit.

Cel mai simplu model este cel în care o tranzație poate avea în așteptare la un moment doar o singură cerere de acces la un obiect. Aceasta înseamnă că ordinul de ieșire al graficului TWFG (numărul maxim de arce ce pot ieși dintr-un nod) este 1. A detecta interblocarea pentru un astfel de model înseamnă a găsi un ciclu în graful global TWFG.

Modelul AND permite unei tranzații să ceară la un moment dat accesul la un set de obiecte. O tranzație va fi blocată dacă nu primește accesul la toate obiectele; de aceea modelul se numește modelul AND. Si în acest model, detecția interblocării înseamnă a găsi un ciclu în graful global TWFG.

O definiție mai riguroasă a interblocării în modelul AND a fost dată de Chandy și Misra ([CM82]). Un proces t_i este dependent de un proces agent t_j dacă există o secvență $seq = t_1, t_2, \dots, t_m, t_i$ de procese agent blocate și fiecare proces din seq , cu excepția primului, deține o resursă la care așteaptă procesul anterior din seq ; se spune că t_i este dependent local de t_j dacă toate procesele agent din seq aparțin aceluiași controler. În acest caz, t_i este interblocat dacă este dependent de el însuși sau de un agent care este dependent de el însuși; în ambele cazuri interblocarea apare numai dacă există un ciclu de procese agent blocate.

Algoritmii de detecție a interblocării pentru modelul AND stabilesc că există interblocare dacă și numai dacă există cicluri de procese blocate. Aceasta condiție nu impune însă, că dacă un agent t_i este interblocat el va fi și detectat ca interblocat. De fapt, dacă t_i este interblocat, dar nu face parte dintr-un ciclu de procese blocate, t_i nu va fi niciodată declarat interblocat. Un exemplu este agentul t_{31} din figura 7.5.3.1.

Interblocarea în modelul cu un singur obiect poate fi definită asemănător, punându-se în plus restricția ca un proces agent să poată avea în așteptare numai o singură cerere (un singur arc care iese din nodul procesului agent). De aici rezultă că modelul AND este o generalizare a modelului cu un singur obiect. Pentru modelul AND au fost propuși și discutați prin exemple și contraexemplu mulți algoritmi de detecție ([Gos91], [CM82][CMH83], [GS80], [MM79]).

Modelul OR permite unei tranzații să ceară un singur obiect dintr-un set de obiecte; un exemplu adecvat este o cerere de citire dintr-un sistem de obiecte replicate: citirea poate avea loc la orice copie disponibilă. În modelul OR, apariția unui ciclu este insuficientă pentru detecția interblocării. În același exemplu din fig.1, dacă se consideră cereri de tip OR, tranzația T_2 nu este blocată, deoarece t_{11} nu are arce de ieșire și după ce t_{11} eliberează obiectul, T_2 poate continua. Definiția interblocării pentru modelul OR a fost dată în [CM82]. În acest model se definește un proces agent blocat, un proces care are o cerere OR în așteptare; fiecărui proces blocat i se asociază un set de procese, denumit *set dependent*, un proces blocat își reia execuția dacă primește permisiunea de acces la un obiect de la un proces din setul său dependent. Un set de procese S este *interblocat* dacă: (1) toate procesele din S sunt blocate; (2) setul dependent al fiecărui proces din S este un subset al lui S; (3) nu există în tranzit mesaje de permisie acces la obiect între procesele din S. Un proces este *interblocat* dacă aparține unui set interblocat.

În termeni TWFG, o interblocare în modelul OR este indicată de prezența unui nod *knot*; prin definiție un nod este de tip *knot* dacă oricare ar fi w cu proprietatea că w poate fi atins din v atunci și v poate fi atins din w . Astfel, detecția interblocării în modelul OR se reduce la a găsi nodurile *knot* din graful TWFG. În secțiunea 8.3.5.2 se prezintă un algoritm de detecție pentru acest model.

Modelul AND_OR este o generalizare a celor două modele anterioare. În modelul AND_OR se poate specifica orice combinație *and* și *or* de cereri de acces la obiect; de exemplu, (a *and* (b *or* c)) *or* d. Nu există o structură din teoria grafurilor care să descrie situația de interblocare în modelul TWFG. În principiu, detectarea

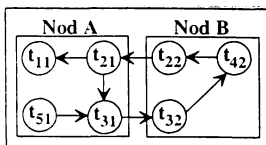


Fig. 7.5.3.1. Graf TWFG.

interblocării în modelul AND_OR se face prin aplicarea repetată a testului pentru modelul OR, exploatând astfel faptul că interblocarea este o proprietate stabilă (adică nu dispăre de la sine).

Modelul (n^k) permite unei tranzații să ceară K resurse dintr-un set de n resurse disponibile. Este o generalizare a modelului AND_OR; este adevărată și reciprocă, adică modelul AND poate fi reprezentat în modelul (n^k) ca (n^k), iar modelul OR ca (1^n). O definiție și un algoritm de detecție a interblocării pentru acest model au fost date în [BT84].

Modelul fără restricții nu impune nici o structurare a cererilor de acces la obiecte; singura condiție impusă este stabilitatea interblocării. O teorie generală a acestui model poate fi găsită în [CL85], [CM86].

8.3.2.3 Clase de algoritmi distribuiți pentru detecția interblocării

Există patru mari clase de algoritmi distribuiți pentru detecția interblocării: *algoritmi care transmit caile, algoritmi care urmăresc arcele, algoritmi cu difuzie și algoritmi care stabilesc starea globală*.

Algoritmii care transmit caile (path-pushing) constituie prima clasă de algoritmi proiectați pentru detectarea interblocării. Acești algoritmi construiesc explicit graful global TWFG. Ideea de bază este de a construi o formă simplificată a grafului TWFG în fiecare nod. În acest scop, fiecare nod transmite graful său local la un anumit număr de vecini de fiecare dată când se execută testul de interblocare. După actualizarea structurilor de date locale, de fiecare nod, grafurile locale TWFG sunt transmise mai departe, procedura repetându-se pînă unul dintre noduri are o imagine suficient de completă pentru a stabili existența sau absența interblocării. Deoarece caracteristica principală a acestor algoritmi este transmiterea căilor din graful TWFG ei sunt înglobați sub denumirea de algoritmi *path-pushing*.

Algoritmii care urmăresc arcele (edge-chasing) detectează prezența unui ciclu într-un graf TWFG prin propagarea unor mesaje de probă în lungul arcelor grafului. Mesajele de probă sunt distincte de mesajele cerere a unui obiect sau mesajele de permisiune acces la obiect. Când un inițiator al unei astfel de proceduri recepționează un mesaj de probă, înseamnă că există un ciclu în acel graf. Un exemplu de algoritm din această clasă este prezentat în [CM82] și descris în secțiunea 8.3.5.1.

Algoritmii cu difuzie au ca idee de bază activarea unei prelucrări cu difuzie, de exemplu, de un manager de tranzații care suspectează o interblocare. Caracteristic prelucrării este faptul că graful global TWFG este reflectat implicit în structura prelucrării; graful global nu este construit explicit. Prelucrarea cu difuzie se împărtășește prin transmiterea de mesaje *query* și se restrînge prin recepția de mesaje *reply*; în momentul când prelucrarea cu difuzie a ajuns înapoi la inițiator, algoritmul se termină și inițiatorul stabilește o interblocare. Un exemplu de algoritm din această clasă este prezentat în [CMH83] și descris în secțiunea 8.3.5.2.

Algoritmii care stabilesc starea globală se bazează pe rezultatele obținute de Chandy și Lamport în [CL85]. Noțiunea cheie este starea globală consistentă a sistemului care poate fi determinată fără suspendarea temporară ("înghețarea") a prelucrărilor aplicației. Chandy și Lamport arată cum se poate obține o stare globală consistentă a unui sistem distribuit prin propagarea unor marcatori de-a lungul canalelor sistemului; o astfel de stare este denumită o vedere (snapshot) a sistemului și poate fi examinată off-line în ceea ce privește interblocările.

8.3.3 Algoritmi centralizați pentru detecția interblocării

În algoritmii centralizați pentru detecția interblocării, graful global TWFG se construiește prin reuniunea tuturor grafulor TWFG locale. El este menținut de către un proces central, denumit *coordonatorul detecției interblocării*. Datorită întârzierilor mesajelor în sistemele distribuite, se disting două tipuri de grafuri TWFG: *graful real* care descrie starea reală a sistemului distribuit, din orice moment, stare care însă nu poate fi cunoscută și *graful construit* care este o aproximare a grafului real, realizat de coordonator în timpul execuției algoritmului. Evident, graful construit trebuie generat astfel încît atunci cînd algoritmul de detecție este executat, rezultatele raportate să fie corecte în sensul condițiilor prezentate în secțiunea 8.3.3. Construcția grafului poate avea loc în următoarele trei moduri: (1) ori de cîte ori un nou arc este adăugat sau eliminat dintr-un graf local; (2) periodic, după un număr de modificări în graf; (3) ori de cîte ori coordonatorul execută algoritmul de detectare a unui ciclu în graf.

Considerăm modelul cererii de acces la un singur obiect. cînd se execută algoritmul de detectare a unui ciclu în graf, coordonatorul analizează graful global; dacă se găsește un ciclu, se selectează o tranzație victimă care se abortează. Există două situații care se pot solda în mod eronat cu abortarea tranzațiilor:

(1) detectarea ciclurilor false care pot apare în graful TWFG global. Exemplul din figura 7.5.3.3.1 ilustrează acest lucru. Ciclul fals apare datorită recepției mesajului de inserare a unui arc înainte de recepția celui de ștergere. Detectarea ciclurilor false de acest tip poate fi ocolită dacă se folosește o marcă de timp pentru fiecare mesaj și ordonarea evenimentelor folosind algoritmul lui Lamport. Astfel, cînd coordonatorul recepționează un mesaj cu marca de timp TS care conduce la o situație de interblocare, poate trimite un mesaj broadcast care conține marca de timp către toate nodurile, pentru că acestea, dacă au cumva mesaje cu mărci de timp mai mici, să le trimită imediat. Metoda necesită însă o circulație excesivă de mesaje.

(2) abortări necesare ale tranzațiilor pot apare și dacă există într-adevar o interblocare și a fost aleasă o tranzație victimă, dar în același timp una dintre tranzații este abortată din alte motive (de exemplu și-a depășit timpul maxim de existență).

Algoritmul centralizat următor [SP88] detectează toate interblocările fără a detecta și interblocări false. Cînd o tranzație T_i din nodul A are un obiect deținut de T_j în nodul B , se trimite către coordonator un mesaj cu marca de timp TS_a și se inserează arcul (T_i, T_j, TS_a) în graful local al lui A . Arcul (T_i, T_j, TS_a) se inserează și în graful local al lui B numai dacă T_j a recepționat mesajul care conține cererea și nu poate elibera resursa. O cerere de la T_i la T_j în același nod este tratată în aceeași manieră dar nu se mai asociază o marcă de timp arcului (T_i, T_j) . Algoritmul de detecție al interblocării este următorul: (1) Coordonatorul trimite un mesaj de inițiere la fiecare nod din sistem; (2) la primirea acestui mesaj, toate nodurile trimit grafurile lor locale coordonatorului. Fiecare dintre aceste noduri conține graful

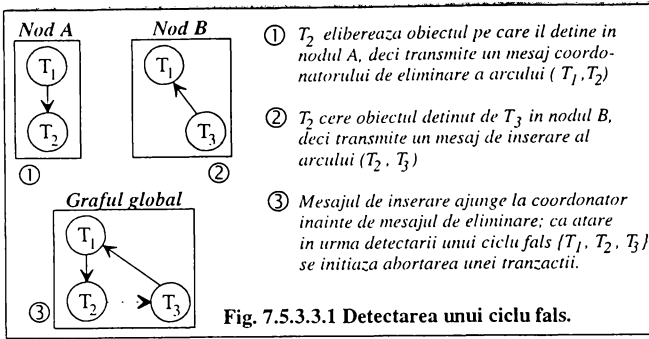


Fig. 7.5.3.3.1 Detectarea unui ciclu fals.

real care reprezintă o stare instanțelor a nodului, dar aceasta nu este sincronizată în raport cu celelalte noduri ale sistemului; (3) când coordonatorul a primit toate răspunsurile, construiește grafurile globale astfel: (a) Grafurile globale conțin câte un nod pentru fiecare tranzație din sistem; (b) Grafurile globale au un arc (T_i, T_j) dacă și numai dacă: (i) există un arc (T_i, T_j) într-unul din grafurile locale; (ii) un arc (T_i, T_j, TS_k) apare în mai mult decât un graf local. Dacă există un ciclu în grafurile globale astfel construite, atunci sistemul este interblocat; dacă

nu există un astfel de ciclu, atunci sistemul nu este interblocat în momentul începerii execuției algoritmului.

8.3.4. Algoritm ierarhic pentru detecția interblocării

Ca și în cazul algoritmului centralizat, fiecare nod menține grafurile sale locale. În contrast însă cu algoritmi centralizați, grafurile globale TWGF este distribuit într-un anumit număr de noduri coordonatoare. Acestea sunt organizate într-un arbore, în care fiecare frunză conține grafurile locale al unui singur nod din sistemul distribuit; celelalte noduri (care nu sunt frunze) ale arborelui conțin

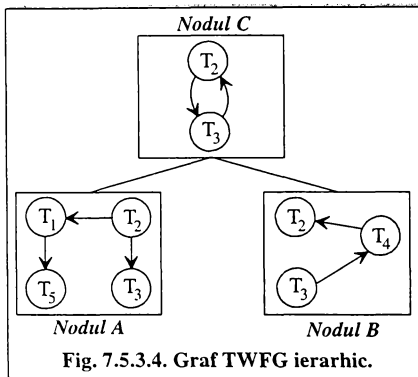


Fig. 7.5.3.4. Graf TWGF ierarhic.

informații relevante ale grafurilor coordonatorilor din subarborii respectivi. Fie trei noduri coordonatoare A, B și C, C fiind cel mai tinar strămoș comun al lui B și C. Presupunem că nodul T_1 apare în grafurile locale ale lui A și B; atunci T_1 trebuie să apară în următoarele grafuri: (1) coordonatorului C; (2) fiecărui coordonator din calea de la C la A; (3) fiecărui coordonator din calea de la C la B. În plus, dacă T_1 și T_2 apar în grafurile globale ale coordonatorului D și există o cale de la T_1 la T_2 în grafurile locale ale lui D, atunci trebuie să existe un arc (T_1, T_2) în grafurile globale ale lui D. Dacă există un ciclu în oricare dintre aceste grafuri, atunci există interblocare. Pentru a ilustra acest algoritm se consideră exemplul din fig. 7.5.3.4. Deoarece T_2 și T_3 apar în A și B, atunci apar și în C; deoarece există în A o cale de la T_2 la T_5 , atunci se include în C un arc (T_2, T_5) . Analog, deoarece există o cale de la T_3 la T_4 în B, arc (T_3, T_4) se include în C. Se observă că grafurile globale conțin un ciclu, ceea ce implică prezența interblocării.

8.3.5. Algoritmi distribuți pentru detecția interblocării

În această secțiune vor fi prezentați doi algoritmi distribuți: un algoritm din clasa algoritmilor edge-chasing pentru detecția interblocării în modele AND și un algoritm din clasa algoritmilor cu difuzie pentru detecția interblocării în modele OR.

8.3.5.1. Algoritm edge-chasing pentru detecția interblocării în modele AND

Algoritm următor a fost propus de Chandy și Misra în [CM82] și [CMH83] pentru detecția interblocării în modele AND. Se presupune că pentru detecția interblocării fiecare controler execută o copie a algoritmului. Pentru a determina când un agent al unei tranzații în așteptarea accesului la un obiect este interblocat, controlerul său inițiază transmisia unor mesaje de probă. Inițierea mesajelor de probă poate fi făcută pentru mai multe tranzații.

Un mesaj de probă constă într-un triplet (i, j, k) în care i specifică inițiatorul, iar j și k că mesajul a fost transmis pe arcul intercontroler (t_j, t_k) . Un controler trimite un mesaj de probă (i, j, k) dacă sunt adevărate următoarele condiții: (1) t_j este în așteptare la o resursă; (2) t_j așteaptă la t_k ; (3) t_j este dependent de t_k . Mesajele de probă recepționate de un controler pot fi așteptate sau ignorate. Un mesaj de probă (i, j, k) este așteptat dacă t_k este în așteptare la o resursă. Este evident că, dacă un controler al unui agent t_j primește mesajul de probă (i, j, i) , oricare ar fi j , rezultă că t_j este interblocat. (fig. 7.5.5.1.)

Algoritm este descris în continuare:

inițierea mesajelor de probă de către controlerul agentului de tranzație t_i :

if t_i este dependent local de el însuși →
declara t_i interblocat

[] t_i nu este dependent local de el însuși →

fa arcele de ieșire (t_i, t_k) ale lui t_i → send (i, j, k)

af

fi

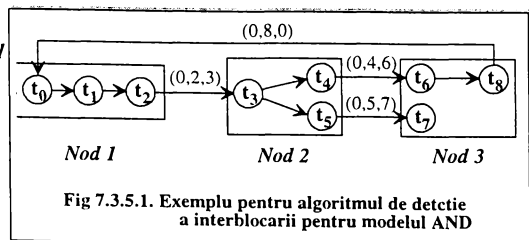


Fig 7.5.5.1. Exemplu pentru algoritmul de detecție a interblocării pentru modelul AND

controlder care recepționează un mesaj de probă::

```

if mesajul de proba (i, j, k) este acceptat →
    if k=i → declara ti interblocat
    [] fa arcele de iesire (tp, tq) ale lui tk → send (i, p, q) af
fi

```

Pentru a întrerupe o interblocare există mai multe soluții. O soluție este ca procesul care a inițiat mesajele de probă să se autodistrugă. Pentru algoritmul prezentat mai sus, această soluție nu este convenabilă, deoarece algoritmul poate fi executat simultan de mai multe procese, deci se pot autodistrugere mai multe procese. O altă soluție este ca fiecare proces să-și adauge identitatea în cadrul mesajului de probă, astfel încât la returnarea mesajului la inițiator lista să fie completă. Inițiatorul poate alege pentru distrugere procesul cu numărul cel mai mare, trimittându-i un mesaj să se autodistrugă. Chiar dacă algoritmul de detecție a interblocării a fost inițiat din mai multe noduri, pentru autodistrugere se va alege același proces.

8.3.5.2 Algoritm cu difuzie pentru detecția interblocării în modele OR

Algoritmul următor a fost propus de Chandy, Misra și Haas în [CMH83] pentru detecția interblocării în modele OR. Un proces agent poate determina cînd este interblocat prin inițierea unei prelucrări cu difuzie. Mai multe procese pot iniția prelucrări cu difuzie în același timp. Mesajele trimise în cadrul unei prelucrări cu difuzie au formele: *query* (*i, j, k*) și *reply* (*i, j, k*), unde *i* specifică procesul inițiator p_i , *j* specifică emittorul p_j , iar *k* receptorul p_k . Există cel mult un mesaj de forma *query* (*i, j, k*) și cel mult o replică *reply* (*i, j, k*) la mesajul *query* (*i, j, k*). Un mesaj inițiator trimite mesaj *query* la toate procesele din setul său dependent (vezi fig. 7.5.3.2.2). Ideea de bază este ca un proces blocat, la primirea unui mesaj *query*, trebuie să propage mesajul în întregul său set dependent, dacă nu a făcut-o deja. Dacă există o secvență de procese blocate p_1, \dots, p_i astfel încît fiecare proces din secvență, cu excepția primului, este în setul dependent al procesului anterior din secvență, un mesaj *query* inițiat de p_i se va propaga la p_j . La primirea unui mesaj *query* sau *reply*, un proces p_k , dacă este activ, îl ignoră; dacă este blocat în așteptare la un obiect, există mai multe posibilități. Dacă p_k a primit primul mesaj *query* de la un inițiator *i* (primul mesaj *query* se numește *de angajare*) el propagă la toate procesele din setul său dependent și memorează numărul de mesaje într-o variabilă locală, fie ea *num[i]*; o altă variabilă locală *wait[i]* specifică faptul că p_k a fost continuu blocat de la primirea mesajului *query*. Dacă p_k primește alte mesaje *query*, raspunde cu *reply* dacă *wait[i]* are valoarea adevarat; în caz contrar, dacă *wait[i]* este fals, ignoră mesajul *query*. cînd p_k primește un mesaj *reply* și *wait[i]* este adevarat, decrementează *num[i]*. Un proces p_k trebuie să raspunda la un mesaj de angajare numai cînd a recepționat un mesaj *reply* pentru fiecare cerere pe care a propagat-o, deci dacă *num[i]=0*. Cînd un proces p_k inițiază o prelucrare cu difuzie, el trimite un mesaj *query* (*k, k, j*) la fiecare proces din setul său dependent și setează *num[k]* la numărul de mesaje trimise. Dacă inițiatorul p_k primește răspunsuri *reply* la toate mesajele *query* trimise, atunci inițiatorul este interblocat. Algoritmul este în continuare:

inițierea unei prelucrări cu difuzie de catre un inițiator p_i ::

```

fa  $p_j$  din setul dependent S al lui  $p_k$  → send_query (k, k, j) af
num[k]:=|S|;
wait[k]:=true;

```

un proces blocat p_k la primirea unui mesaj query (i, j, k) executa ::

```

if mesajul query este de angajare →
    fa  $p_m$  din setul S al lui  $p_k$  → send_query (i, k, m) af
    num[i]:=|S|; wait[i]:=true;
[] mesajul query nu este de angajare →
    if wait[i] = true → send_reply (i, k, j) la  $p_j$  fi
fi

```

un proces blocat p_k la primirea unui mesaj reply (i, k, j) executa ::

```

if wait[i]=true →
    num[i]:=num[i]-1
    if num[i] = 0 →
        if i = k → declara interblocare pentru  $p_k$ 
        [] i ≠ k → send_reply (i, k, j) la procesul care a trimis mesajul de angajare la  $p_k$ 
fi
fi

```

Pentru a garanta că fiecare interblocare va fi detectată, trebuie eliminată restricția ca prelucrarea cu difuzie să poată fi inițiată numai o singură dată de un proces; ca atare, ori de cîte ori se va bloca, un proces va putea iniția o prelucrare cu difuzie. Pentru a se distinge prelucrările cu difuzie inițiate de același proces p_i , mesajele *query* și *reply* vor conține un parametru în plus care să specifice numărul de secvență al prelucrării cu difuzie inițiată de p_i .

Capitolul 9.

PROIECTAREA UNUI SUPORT SOFTWARE PENTRU PROGRAMARE DISTRIBUITĂ BAZATĂ PE OBIECTE ȘI PE TRANZACȚII ATOMICE

Există actualmente un interes crescut în metodologii și tehnici cu caracter general care simplifică construcția sistemelor software distribuite robuste și eficiente. Acest capitol prezintă o astfel de metodologie bazată pe tranzacții atomice și include proiectarea unui suport software pentru programarea distribuită bazat pe obiecte (DOSTP). Tranzacțiile atomice micșorează dificultatea scrierii unor programe de aplicații distribuite prin simplificarea tratării defectiunilor nodurilor și prin rezolvarea problemelor de concurență.

9.1. Cerințe de proiectare

Înainte de a prezenta arhitectura sistemului DOSTP se enumeră cerințele de proiectare impuse:

- (1) simplitate de apel ale funcțiilor sistemului DOSTP la nivelul aplicațiilor utilizator;
- (2) modularitate și extensibilitate, în sensul că pot fi introduși în sistem noi manageri de obiecte printr-un efort de programare minim; practic, activitatea de programare a unui manager pentru un anumit tip de obiecte se concentrează numai pe structura obiectelor și a cererilor primite relative la tipul de obiecte respective, pe tratarea cererilor și rezolvarea problemelor de concurență în accesul la obiecte și pe recuperarea obiectelor în urma defectării nodului pe care acestea au fost rezidente. În ceea ce privește conectarea la sistemul tranzacțional al unui manager de obiecte, sistemul DOSTP trebuie să pună la dispoziție modalitățile și funcțiile necesare.
- (3) sistemul trebuie să ofere posibilitatea execuției în paralel a mai multor operații asupra unor obiecte diferite;
- (4) sistemul trebuie să ofere posibilitatea recuperării obiectelor în urma defectării unuia sau mai multor noduri implicate într-o tranzacție atomică;
- (5) sistemul trebuie să se execute în mod utilizator;
- (6) sistemul trebuie să ofere o comunicație sigură;
- (7) obiectele trebuie să fie portabile (cu modificări minime, absolut necesare) atât în rețele UNIX cât și în rețele Microsoft Windows (Win32);
- (8) tranzacțiile să nu poată fi amânate la infinit. Se cere deci evitarea situației de amânare la infinit sau a restartărilor la infinit.

În urma analizei cerințelor de mai sus s-a optat ca DOSTP să ofere primitive pentru: (1) conectarea și deconectarea aplicațiilor utilizator la sistem; (2) inițierea unei tranzacții atomice; (3) execuția într-o manieră asincronă (returnarea controlului se face imediat după apelul operației) de operații asupra unor obiecte, ai căror manageri sunt conectați la sistem; (4) execuția, într-o manieră sincronă (returnarea controlului se face după execuția completă a operației) de operații asupra unor obiecte, ai căror manageri sunt conectați la sistem; (5) așteptarea terminării execuției tuturor operațiilor asincrone lansate anterior, în cadrul tranzacției atomice curente; (6) terminarea normală a unei tranzacții atomice, deci permanentizarea operațiilor executate; (7) terminarea anormală a unei tranzacții atomice.

Se prezintă, în continuare, o descriere în detaliu a primitivelor furnizate de DOSTP:

(1) **int InitDOSTP (char *name)**

Scop:

Permite conectarea unei aplicații utilizator la sistemul DOSTP.

Parametri:

name - numele aplicației utilizator

Valoare returnată:

Funcția returnează SUCCESS în cazul în care operația a reușit sau un cod de eroare în caz contrar. Codurile de eroare care pot fi obținute sunt:

EBADF

-numele aplicației este un șir nul

EDOSTP_INITIALIZE_DONE

-aplicația era conectată la DOSTP

ECREATE_APL_FIFO

-nu se poate crea canalul de comunicație cu aplicația

(2) **int BeginTransaction (void)**

Scop:

Permite specificarea începerii unei noi tranzacții atomice.

Valoare returnată:

Funcția returnează SUCCES în cazul în care operația a reușit sau cod de eroare în caz contrar. Codurile de eroare care pot fi obținute sunt:

EDOSTP_INITIALIZE_EXPECTED - nu s-a apelat *InitDOSTP()*

EBEGIN_TRANS_DONE - aplicația are tranzacția deschisă

EOPEN_APL_FIFO - nu se poate deschide canalul de comunicație cu aplicația

EOPEN_TM_FIFO - nu se poate deschide canalul de comunicație cu managerul de tranzacții

EWRITE_TM_FIFO_FAIL - nu se poate comunica cu managerul de tranzacții

ERead_APL_FIFO_FAIL - nu se poate citi canalul aplicației

ETM_CREATE_SHM - man. de tranzacții nu a reușit crearea zonei de memorie partajate

EAPL_ACCESS_SHM_FAIL - nu se poate realiza atașarea la zona de memorie partajată

ETM_CREATE_SEM - man. de tranzacții nu a reușit crearea semafoarelor

ETM_OPEN_TM_FIFO - man. de tranzacții nu a reușit deschiderea canalului de comunicație

ECS_TIME_OUT - eroare time-out la comunicatia cu serverul de comunicatii

(3) **int SyncRequest(char *objectName, void *sourceAddr, void *destAddr, int timeout)**

Scop:

Permite specificarea execuției unei operații sincrone asupra unui obiect *objectName*. Aplicația este pusă în așteptare pînă în momentul terminării execuției operației. Utilizatorul trebuie să rezerve o zonă de memorie în care să transmită parametri operației, a cărei adresă o transmite ca parametru și o altă zonă de memorie în care sistemul returnează răspunsul la cererea sincronă, a cărei adresă se transmite de asemenea ca parametru. Structura celor două zone de memorie depinde de tipul obiectului. Ea trebuie cunoscută numai de managerul obiectelor de acest tip și de aplicație. Pentru sistemul tranzacțional DOSTP structura acestor zone de memorie nu are nici o importanță: atît cererea cît și răspunsul sunt tratate ca un șir de octeți.

Parametri:

sourceAddr - adresa de început a zonei de memorie la care apelantul a depus cererea de execuție a operației sincrone

destAddr - adresa de început a zonei de memorie la care sistemul returnează răspunsul la cererea sincronă

objectName - numele obiectului implicat în tranzacție

timeout - valoarea maximă pentru timpul de așteptare al răspunsului

Valoare returnată:

Funcția returnează valoarea SUCCES în cazul în care cererea a fost tratată de managerul de obiecte, ceea ce nu înseamnă că operația s-a executat cu succes, sau un cod de eroare în caz contrar. Codurile returnate pot avea următoarele valori:

EBADF - numele obiectului este un șir nul

EDOSTP_INITIALIZE_EXPECTED - nu s-a apelat *InitDOSTP()*

EINVALID_TRANS - nu s-a executat *BeginTransaction()*.

ENO_BUF - nu mai există loc în buferele sistemului pentru aceasta cerere și pentru răspuns

ETIME_OUT - s-a depășit timpul maxim de așteptare specificat în parametrul de intrare

ETOO_MANY_OBJECTS - s-a depășit numărul maxim de obiecte permise într-o tranzacție

EOBJ_NOT_FOUND - obiect inexistent

(4) **int AsyncRequest (char *nameObject, void *sourceAddr, void*destAddr, int *pErrorCode, int timeout)**

Scop:

Permite specificarea execuției unei operații asupra unui obiect *objectName*. Controlul execuției este redat aplicației imediat după apelul funcției. Ca și în cazul funcției precedente, *SyncRequest*, utilizatorul trebuie să rezerve două zone de memorie pentru cererea asincronă, respectiv pentru răspuns, ale căror adrese de început se dau ca parametri de intrare. Formatul acestor zone trebuie cunoscut numai de managerul de obiecte. Avînd în vedere că răspunsul la cerere se obține la *destAddr* după ce controlul a fost returnat aplicației, la *pErrorCode* se livrează o prima indicație despre desfășurarea execuției, în sensul ca cererea a fost sau nu acceptată de sistem.

Parametri:

nameObject - numele obiectului

sourceAddr -adresa de început a zonei de memorie la care apelantul a depus cererea de execuție a operației sincrone
destAddr -adresa de început a zonei de memorie la care sistemul returnează răspunsul la cererea asincronă
pErrorCode -adresa la care sistemul va returna codul de eroare obținut în urma tratării cererii de către sistemul tranzacțional. Codul erorii obținute în urma tratării cererii de către managerul de obiecte se codifică în răspunsul livrat la *destAddr*.
timeOut -aloarea maxima pentru timpul de așteptare al răspunsului. Dacă aplicația nu primește răspuns în acest interval de timp la adresa *pErrorCode* se va returna ETIME_OUT.

Valoare returnată:

Funcția returnează SUCCES în cazul în care cererea a fost acceptată de sistem sau unul dintre următoarele coduri de eroare:

EBADF -numele obiectului este un șir nul
EDOSTP_INITIALIZE_EXPECTED -nu s-a apelat *InitDOSTP()*
EINVALID_TRANS -nu s-a apelat *BeginTransaction()*
ENO_BUF -nu mai există loc în buferele sistemului pentru această cerere și pentru răspuns
La *pErrorCode* se pot returna următoarele valori:
ETIME_OUT -s-a depășit timpul maxim de așteptare specificat în parametrul de intrare.
ETOO_MANY_OBJECTS - s-a depășit numărul maxim de obiecte permise într-o tranzacție
EOBJ_NOT_FOUND - obiect inexistent

(5) **int Synchronize (void)**

Scop:

Permite așteptarea terminării tuturor operațiilor asincrone lansate anterior.

Valoare returnată:

Funcția returnează SUCCES sau următorul cod de eroare atunci când nu s-a inițiat nici o tranzacție dar se cere execuția ei:

EDOSTP_INITIALIZE_EXPECTED - nu s-a apelat *InitDOSTP()*
EINVALID_TRANS -nu s-a executat *BeginTransaction()*.

(6) **int EndTransaction (void)**

Scop:

Permite terminarea normală a unei tranzacții atomice lansate anterior. Funcția comandă sistemului DOSTP execuția celor două faze ale procedurii *two-phase commit*. Înainte de a se încerca desavârșirea tranzacției aplicația trebuie să încerce mai întâi terminarea tuturor operațiilor asincrone lansate, printr-un apel *Synchronize()*.

Valoare returnată:

Funcția returnează SUCCES în cazul în care desavârșirea tranzacției a reușit sau următoarele coduri de eroare:

EDOSTP_INITIALIZE_EXPECTED - nu s-a apelat *InitDOSTP()*
EINVALID_TRANS - nu s-a apelat *BeginTransaction()*
ESYNCRONIZE_EXPECTED - trebuie apelată funcția *Synchronize()* și apoi *EndTransaction()*
ETIME_OUT - desavârșirea tranzacției nu a reușit în intervalul de timp maxim
ECOMMIT_FAIL - eroare în desavârșirea tranzacției (obiectele nu au putut realiza în totalitate operațiile cerute)

(7) **int AbortTransaction (void)**

Scop:

Permite terminarea anormală sau anularea unei tranzacții atomice lansate anterior. Funcția trebuie apelată obligatoriu după orice operație asupra unui obiect nereușită.

Valoare returnată:

Funcția returnează SUCCESS sau valorile:

EDOSTP_INITIALIZE_EXPECTED -nu s-a apelat *InitDOSTP()*
EINVALID_TRANS -nu s-a apelat *BeginTransaction()*

(8) **int CloseDOSTP (void)**

Scop:

Funcția se apelează pentru a deconecta aplicația de la sistemul TCP.

Valoare returnată:

Funcția returnează SUCCESS sau valorile:

EDOSTP_INITIALIZE_EXPECTED -nu s-a apelat *InitDOSTP()*

9.2. Modelul sistemelor DOSTP

În contextul celor prezentate în capitolele 6, 7 și 8 se apreciază că:

(1) Modelul de arhitectură al sistemului DOSTP este *modelul integrat* (vezi 7.2.2). În cadrul acestui model toate nodurile (în care se inițiază tranzațiile) sunt dotate cu manageri de tranzații și cu manageri de obiecte, dacă există obiecte pe acele noduri.

(2) Pentru realizarea tranzațiilor distribuite se folosește modelul în care se crează câte un *proces agent* în fiecare nod implicat în execuția unei tranzații (vezi 7.2.2); procesul agent va obține accesul la obiectul de pe acel nod și va realiza dialogul manager de tranzație - manager de obiect.

(3) DOSTP utilizează *varianta dinamică a obiectului activ* (vezi 6.3) în care, atunci când se recepționează o cerere de implicare a unui obiect într-o tranzație, se crează câte un proces server pentru fiecare obiect. Toate operațiile cerute de acea tranzație pentru obiectul respectiv vor fi tratate numai de acel proces server. Avantajul acestei metode este faptul că nu limitează numărul de tranzații ce pot fi tratate simultan la nivelul unui obiect (decât, evident, de numărul maxim de procese concurente permise de sistemul de operare, ceea ce, în mod normal, este puțin probabil de atins).

(4) Localizarea obiectelor se face folosind soluția *serverelor de nume distribuite* (vezi 6.4.1). Sistemul crează câte un server de nume pe fiecare nod. Există un manager central, care cunoaște toate nodurile componente ale sistemului DOSTP și la care se conectează toate nodurile la intrarea în sistem. Acest manager are și rolul de a actualiza informațiile din serverele de nume, la conectarea sau eliminarea unui manager de obiecte în sistem, prin transmisia de mesaje către toate serverele de nume.

(5) Operațiile asupra obiectelor sunt apelate folosind *transmisia de mesaje* (vezi 6.4.2), strategie tipică pentru sistemele în care se implementează modelul obiectului activ. Când o aplicație apelează o operație asupra unui obiect, parametrii cererii sunt împachetați într-un mesaj; mesajul este transmis apoi procesului manager al obiectului, după ce obiectul a fost localizat folosind serverul de nume. Procesul manager al obiectului recepționează mesajul, despachetează parametrii din mesaj, execută operația și împachetează rezultatul într-un alt mesaj pe care-l transmite apelantului. Pentru apelurile la obiecte situate în același nod cu aplicația care a inițiat tranzația, în DOSTP s-a optat ca mesajele să fie transmise via o zonă de memorie comună, partajată de aplicație și de managerul de obiecte.

(6) În DOSTP se consideră că se lucrează cu cereri de tratare a unor operații a căror execuție nu depășește un timp prestabilit, astfel încât răspunsul la o cerere se consideră și o confirmare a recepționării cererii (vezi 3.1.3).

(7) În ceea ce privește gestionarea operațiilor asupra obiectelor (vezi 6.3.1) DOSTP folosește *mecanismul tranzațiilor atomice*. Cind toate operațiile grupate în cadrul unei tranzații au fost executate cu succes, sistemul este cel care realizează procedura de permanentizare a modificărilor. Dacă cel puțin o operație se termină prin eroare, tranzația trebuie terminată prin abort.

(8) Protocolul ales pentru terminarea tranzațiilor distribuite este protocolul *two-phase commit* (vezi 7.3.1). Protocolul presupune existența unui coordonator al comiterii tranzației; coordonatorul unei tranzații distribuite este întotdeauna în DOSTP managerul de tranzații de pe nodul de la care s-a inițiat tranzația respectivă. Coordonatorul tranzației supravezează comunicația cu toți agenții de la distanță ai unei tranzații.

(9) Schema de sincronizare folosită în DOSTP pentru sincronizarea operațiilor la nivelul aceluiași obiect, în cadrul unui manager de obiecte, se poate încadra în *schemele optimiste* prezentate pe scurt în 6.3.2 și mai în detaliu în 7.4.6. Schema de sincronizare folosită este proprie managerului de obiecte. Pot fi realizați mai mulți manageri de obiecte pentru același tip de obiecte, fiecare folosind propria sa schemă de sincronizare, de o manieră optimistă sau pesimistă. Pentru managerii de date realizați până în momentul de față s-a folosit o schema de sincronizare optimistă. Se apreciază că o direcție de extindere a acestui sistem poate consta în testarea mai multor scheme de sincronizare pentru același tip de obiect. Se poate utiliza chiar și o schemă hibridă care funcționează astfel: monitorizându-se frecvența de conflicte ale cererilor de operații asupra unui obiect, atunci când frecvența este mai mică decât limita maximă, se utilizează varianta optimistă, iar când conflictele sunt mai dese (frecvența mare) se utilizează varianta pesimistă. Schemele de sincronizare optimiste permit cel mai mare grad de concurență posibil la nivelul unui obiect; operațiile nu sunt puse niciodată în așteptare. Prețul plătit pentru această performanță este faptul că acțiuni executate fără erori pot fi terminate anormal (la *commit*) deoarece la desăvârșirea tranzațiilor se testează condițiile de serabilitate astfel încât informațiile de stare luate în considerare de o operație să nu intre în conflict cu modificările realizate de orice altă operație care s-a terminat (a fost acceptată în faza *prepare-to-commit*); în plus, trebuie menținute în memorie mai multe copii ale obiectului.

(10) În ceea ce privește reprezentarea obiectelor în memorie (vezi 6.5.) se precizează faptul că DOSTP lucrează cu o copie persistentă a obiectului, memorată într-un fișier care, la încărcarea obiectului, este mapat în memorie. Fiecare manager de obiect menține, local unei tranzații, toate modificările făcute asupra obiectului de către acea tranzație care se permanentizează în faza de *commit*.

(11) Pentru siguranța obiectelor, DOSTP folosește metoda recuperării obiectelor, varianta *roll-back* (vezi 6.3.3 și 7.3.4). Se lucrează cu un sistem de recuperare al obiectelor bazat pe *fisiere jurnal (log)* care folosește protocolul *Write-Ahead Logging* (vezi 7.3.4.1) și tehnicile *Redo-only-logging* și *Value-logging* (vezi 7.3.4.2 și 7.3.4.3). Tehnica *Redo-only-logging* este avantajoasă ținând cont că se folosește o schemă de sincronizare optimistă; actualizările sunt întârziate până în momentul comiterii tranzacției, astfel încât informațiile de *undo* nu trebuie scrise în log. Sistemul garantează că, odată ce tranzacția execută *commit*, orice actualizări întârziate sunt realizate, chiar dacă intervine un incident hardware. Tehnica constă în crearea unor *liste de intenții* care să cuprindă actualizările în așteptare și de a forța scrierea acestor liste în log, înainte de faza *prepare-to-commit*.

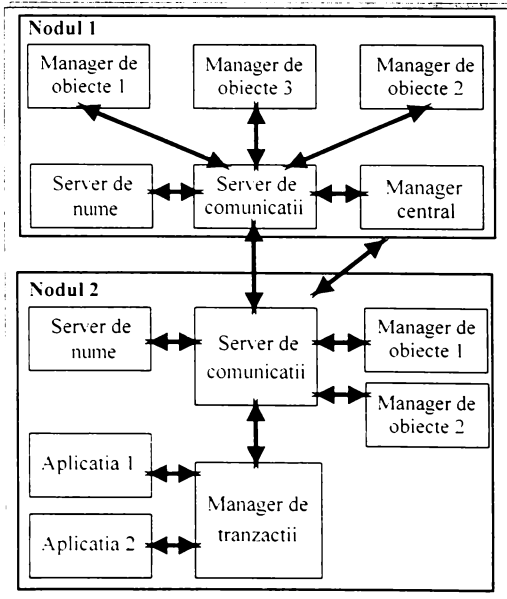
În ceea ce privește structura sistemului de recuperare (vezi 7.3.4.4), elementele sale sunt încorporate la nivelul managerilor de obiecte și la nivelul managerilor de tranzacții.

(12) În ceea ce privește controlul interblocării obiectelor se fac două observații (vezi 7.4.6.1): (1) Deoarece se folosește o schemă de sincronizare optimistă nu pot apărea situații de interblocări și restartări ciclice. (2) Pot apărea însă situații de restartare infinită. Prin urmare, pentru respectarea cerinței 8 de proiectare impusă sistemului DOSTP, care interzice apariția restartărilor la infinit, s-a optat pentru implementarea unui mecanism *Wound-die* (vezi 8.2). Acest mecanism evită restartările la infinit ale tranzacțiilor, dar impune asocierea unor mărci de timp tranzacțiilor și reluarea tranzacțiilor după abortare cu aceeași marcă de timp.

Deoarece tranzacțiile pot fi inițiate de pe orice nod din sistem, pentru alocarea unei mărci de timp unice unei tranzacții, care de fapt este o prioritate asociată unei tranzacții- cu cât prioritatea are o valoare mai mică, tranzacția este mai prioritară, sau, altfel spus, are o marcă de timp mai mică, deci este mai veche în sistem și deci mai prioritară)- pot fi folosite două soluții: o soluție complet distribuită care presupune însă circulație de mesaje pentru alegerea mărcii de timp (sau priorității) în regim de excludere reciprocă (poate fi folosit unul din algoritmi prezentati la 4.3.2, 4.3.3, 4.3.4) sau o soluție cu proces coordonator. În DOSTP s-a optat, pentru a simplifica circulația de mesaje, pentru soluția unui proces coordonator. În DOSTP, pe un nod fixat există un *manager central* cu dublu rol: (1) primul rol a fost prezentat la punctul (2); el este legat de conectarea nodurilor la sistemul DOSTP și de conectarea sau eliminarea unor manageri de obiecte; (2) al doilea rol al managerului central constă în alocarea de priorități tranzacțiilor, la inițierea lor. La recepționarea unei cereri *BeginTransaction()* fiecare manager de tranzacție contactează managerul central pentru a-i livra prioritatea tranzacției. Când o tranzacție este reluată, prioritatea tranzacției nu se modifică; astfel, după câteva încercări nereușite, devine cea mai veche din sistem și poate fi terminată normal.

9.3. Arhitectura generală a sistemului DOSTP

Componentele sistemului DOSTP dintr-un nod în care se inițiază tranzacții și în care există și manageri de obiecte și dintr-un nod în care există numai obiecte sunt prezentate în figura 9.3.1. Managerul central se poate găsi pe orice nod al sistemului, însă fixat și cunoscut. Dacă într-un nod nu se implementează obiecte, atunci managerii de obiecte pot lipsi; pe de altă parte,



într-un nod pot fi rezidenți mai mulți manageri de obiecte, pentru același tip de obiecte sau pentru tipuri diferite. Analog, managerul de tranzacții poate lipsi dintr-un nod, dacă de pe acel nod nu se intenționează să se execute aplicații care folosesc sistemul DOSTP. Pe un nod există (eventual) un singur manager al tranzacțiilor, un singur server de nume, un singur server pentru comunicații, dar pot exista mai mulți manageri de obiecte și mai multe aplicații.

Pe scurt, funcțiile componentelor DOSTP sunt:

1. *Managerii de obiecte* stochează informațiile de stare ale obiectelor, realizează mecanismul de sincronizare specific tipului implementat și execută operațiile asupra obiectelor. Ei cooperează cu managerii de recuperare pentru a implementa protocolul pentru recuperarea tranzacțiilor. Managerii de obiecte sunt proiectați folosind *tehnica programării orientate pe obiecte* astfel încât integrarea unui manager de obiecte în sistemul tranzacțional să

nu depindă de tipul obiectului. Managerii de obiecte sunt singurele componente care cunosc structura cererilor aplicațiilor și structura obiectelor.

2. *Serverele de nume* (există câte un server în fiecare nod conectat la DOSTP) realizează localizarea obiectelor în rețea; ele mențin structuri de date care memorează poziția obiectelor pe noduri.

3. *Serverele de comunicație* transmit mesajele inter-noduri, generează procesele agent în cadrul tranzațiilor distribuite, și detectează defecțiunile altor noduri cu care inițiază comunicația.

4. *Managerii de tranzații* coordonează inițierea, desăvârșirea și terminarea anormală a unei tranzații.

5. *Sistemul de recuperare* are trei funcții majore: (a) înainte de transmisia confirmării mesajului *prepare-to-commit*, înscrie listele de intenții în fisierul log; formatul înregistrărilor depinde de tipul obiectelor, de aceea s-a considerat oportun ca această funcție să fie înglobată în managerii de obiecte; (b) o altă funcție este recuperarea obiectelor în urma defectării nodului respectiv și apoi a repunerii în funcțiune a lui; și această funcție depinde de structura obiectelor și a fost înglobată în managerii de obiecte; (c) a treia funcție constă în luarea deciziei de a executa *commit* sau *abort* în urma analizei fișierelor de log după o defecțiune a unui nod coordonator al unei tranzații urmate de repunerea lui în funcțiune. Aceste fișiere de log care conțin informații despre starea tranzațiilor și obiectelor implicate în tranzații se mențin la nivelul managerului de tranzații, în calitatea sa de coordonator al tranzațiilor distribuite, și, în consecință, funcția (c) este înglobată la nivelul managerului de tranzații.

7. *Procesele de aplicație* inițiază tranzațiile și apelează operațiile implementate în serverele de date.

8. *Managerul central* este un singur proces în tot sistemul, pe un nod fixat. El are două funcții: (a) ține evidența nodurilor conectate la sistem și transmite serverelor de nume, la conectarea lor, localizarea obiectelor din sistem; (b) livrează prioritățile pentru tranzații.

9.3.1. Procesele utilizator ale sistemului DOSTP

Procesele utilizator (sau de aplicație) ale sistemului DOSTP constituie procesele de nivelul cel mai înalt în utilizarea tranzațiilor, procesele care inițiază tranzațiile. Aplicațiile apelează operații asupra obiectelor prin trimiterea de mesaje și recepționarea rezultatelor de la managerii de date.

Spre deosebire de managerii de date, aplicațiile nu memorează obiecte în corpul lor și de aceea nu acceptă mesaje de la alte servere de date. Aplicațiile nu sînt afectate de mecanismele de recuperare și deci nu sunt restartate după o defecțiune de nod.

Pentru a putea apela operații asupra obiectelor, aplicațiile trebuie să cunoască structura cererilor recunoscute de manageri pentru acele tipuri de obiecte. Pe de alta parte, aplicațiile trebuie să realizeze conectarea cu sistemul DOSTP. Aceste cerințe sunt îndeplinite prin linkeditarea aplicațiilor cu un set de funcții puse la dispoziție aplicațiilor de sistemul DOSTP. Pe un nod al sistemului în care se execută aplicații, trebuie să existe un manager de tranzații; acesta poate deservi un număr practic nelimitat de aplicații locale.

9.3.2. Managerul tranzațiilor

Managerul tranzațiilor este un proces a cărui funcție este de a coordona inițierea și terminarea unei tranzații. Informațiile despre starea tranzațiilor sunt furnizate prin două clase de mesaje recepționate de managerul tranzațiilor de la alte componente STAD.

Prima clasă de mesaje declanșează inițierea și terminarea unei tranzații. A două clasă informează managerul de tranzații ce noduri la distanță și ce manageri locali sînt participanți într-o tranzație dată.

Procesul care a inițiat o tranzație trebuie să fie obligatoriu și cel care încearcă să o termine cu *commit*. Pentru terminarea unei tranzații, managerul tranzațiilor utilizează o variantă de protocol *two-phase commit*. Protocolul este bazat pe construcția unui arbore; fiecare nod este coordonator de tranzație pentru nodurile sale fii. Un nod A este părintele unui alt nod B dacă și numai dacă nodul A a fost primul nod care a apelat o operație în nodul B în decursul aceleiași tranzații. Informațiile referitoare la relațiile unui nod cu alte noduri ascendente sau descendente direct în arborile tranzațiilor sînt păstrate în managerul de tranzații.

9.3.3. Serverele de comunicații

Serverele de comunicații sînt procese rezidente în fiecare nod, care acționează ca un tampon între emițătorul unui mesaj și receptorul sau, realizînd localizarea transparentă a proceselor. Perechea de procese corespondente, servere ale comunicațiilor între cele două noduri asigură astfel gestiunea comunicațiilor.

Serverele de comunicație furnizează și funcția de monitorizare a rețelei. Un server de comunicație detectează defecțiunile de comunicație și de noduri în următoarele situații:

1. *pierderea conexiunii*. Indică o defecțiune a nodului cu care se comunică în mod direct. Este detectată de managerul comunicațiilor, fie datorită lipsei mesajului de confirmare, fie datorită nerecunoașterii unei conexiuni deja folosite și confirmate anterior.

2. *erori de tip time-out*. Serverul comunicațiilor furnizează proceselor emițătoare (din nodul său) un port pentru a fi adresat atunci cînd procesul receptor este la distanță. El nu cunoaște noțiunea de tranzație dar

este responsabil de construcția arborelui tranzacției; înregistrarea identității nodurilor fii este menținută însă la nivelul managerului de tranzacții.

Serverele de comunicații sunt singurele procese dintr-un nod DOSTP care comunică cu alte noduri. Un server de comunicație dintr-un nod este contactat de managerul de tranzacții, de serverul de nume și de managerii de obiecte pentru orice schimb de mesaje cu alte procese de pe alte noduri.

9.3.4. Serverele de nume

Sistemul DOSTP prevede un serviciu de gestiune a numelor permanente pentru obiecte. Procesele servere de nume se execută în fiecare nod și cooperează cu serverul central pentru gestiunea numelor permanente. Serverele de nume livrează managerilor de comunicații locali, informații despre localizarea obiectelor din rețea, în momentul în care managerii de tranzacție primesc o cerere de operație asupra unui obiect nou implicat într-o tranzacție.

9.3.5. Managerul central

Managerul central este un proces unic în sistemul DOSTP. El este contactat în trei situații:

(1) la conectarea unui nod în sistem, la începutul execuției serverului de comunicații. Serverul de comunicații își anunță prezența în sistem, după care primește de la managerul central o listă cu toate obiectele din sistem pe care o transmite serverului de nume din același nod.

(2) la introducerea sau eliminarea unui nou manager de obiecte în sistem;

(3) la inițierea unei tranzacții când managerul de tranzacții, prin intermediul serverului de comunicații cere o prioritate pentru noua tranzacție.

Pentru acordarea identificatorilor și priorităților tranzacțiilor se folosește un fișier local care reține valoarea următoare de acordat identificadorului, respectiv priorității unei tranzacții.

9.3.6. Managerii de obiecte

Managerii de date încapsulează două tipuri de date, obiectele și informația de sincronizare asociată și patru tipuri de funcții:

(1) funcțiile care implementează operațiile care acționează asupra obiectelor;

(2) funcțiile care implementează mecanismul pentru sincronizarea tranzacțiilor în accesul la obiecte;

(3) funcțiile care implementează dialogul cu managerul de comunicații;

(4) funcțiile care implementează restaurarea obiectelor după o terminare anormală sau defect al unui nod.

Funcțiile din categoriile (1), (2) și (4) depind de tipul obiectului gestionat de manager; cele din categoria (3) nu depind de tipul obiectului. Având în vedere aceasta observație o consecință este imediată: folosind tehnici de programare specifice programării orientate pe obiecte trebuie să se separe în clase separate funcțiile de tipul (3) și funcțiile de tipurile (1), (2) și (4); în acest fel, pentru proiectarea unor noi manageri de obiecte, vor fi rescrise numai funcțiile de tipul (1), (2) și (4).

Fiecare manager de obiecte acceptă mesaje care invocă execuția unor operații și returnează rezultatul execuției în alte mesaje.

După un defect al unui nod, se startează noi procese manageri de date pentru a accesa obiectele protejate. În aceste situații, starea obiectelor partajabile se restaurează; la restaurare, contribuie funcțiile din categoria (4).

Managerii de date conțin și implementarea mecanismului ales pentru sincronizarea tranzacțiilor în DOSTP; acesta poate folosi o schema pesimistă, optimistă sau o schema hibridă. Managerii de obiecte vor conține relațiile de compatibilitate pentru lock-uri, tabelele care descriu care lock-uri sunt deținute de tranzacții și identitatea acestor tranzacții și codul care manipulează aceste structuri de date.

Se apreciază că managerii de obiecte constituie o abstractizare a unor tipuri de date puse la dispoziția aplicațiilor.

9.3.7. Sistemul de recuperare

Sistemul de recuperare nu este implementat în DOSTP ca proces distinct; el se compune din module plasate în managerii de obiecte și în managerii de tranzacții precum și din fișiere jurnal (sau fișiere de log) plasate analog lângă managerul de tranzacție aferent sau lângă managerul de obiecte.

Sistemul de recuperare se constituie din:

(1) codul procedurilor pentru recuperarea obiectelor, proceduri lansate în execuție după defectarea și apoi repunerea în funcțiune a unui nod și plasate în managerii de obiecte;

(2) codul funcțiilor pentru scrierea articolelor în fișierele log, conform tehnicilor *Redo-only-logging* și *Value-logging* plasate în managerii de obiecte;

(3) codul funcțiilor pentru analiza fișierelor log în scopul determinării stării tranzacțiilor în momentul apariției defecțiunii unui nod, plasat în managerii de tranzacție.

(4) fișierele de log care rețin informații despre tranzacțiile active și despre obiectele implicate în aceste tranzacții; aceste fișiere sunt plasate lângă managerul de tranzacție, coordonator al tranzacției respective.

(5) fișierele de log care mențin informațiile cu care trebuie actualizată copia permanentă a obiectelor în momentul execuției fazei *commit*.

9.4. Proiectarea elementelor componente ale sistemului DOSTP

În subcapitolul de față se prezintă modul de proiectare al elementelor componente ale sistemului DOSTP, insistându-se asupra structurilor de date utilizate și a circulației mesajelor între componente.

9.4.1. Structuri de date ale sistemului

Înainte de a se prezenta structurile de date ale sistemului este necesar să se arate pe scurt modul de desfășurare al unei tranzacții în cadrul sistemului DOSTP și modul de comunicare între elementele sistemului DOSTP.

Există patru momente distincte în funcționarea sistemului DOSTP:

(M1) Momentul conectării unui nod la sistem, când se lansează în execuție serverul de comunicații și serverul de nume. Dacă pe acel nod urmează să se execute aplicații se lansează în execuție și managerul de tranzacții. Serverul de comunicații contactează managerul central și primește o listă cu localizarea tuturor obiectelor introduse în sistem. Din acest moment, aplicațiile pot să lucreze cu tranzacții.

(M2) Momentul introducerii sau eliminării unui obiect din sistem, când managerul de obiecte trimite un mesaj adecvat managerului central care anunță toate nodurile din sistem (via servere de comunicații - servere de nume) despre modificarea operată în lista de obiecte.

(M3) Desfășurarea unei tranzacții are loc astfel:

(3.1) Aplicația lansează un apel *BeginTransaction()*, a cărui urmare este transmiterea unui mesaj ALP-BEGIN-TRANSACTION la managerul de tranzacții.

(3.2) Managerul de tranzacții (TM) care așteaptă în buclă continuă recepționarea de mesaje pe canalul *chan-TM*, primește mesajul și creează un proces fiu ce creează o zonă de memorie partajată și un set de semafoare pentru accesul la această zonă și trimite identificatorul zonei și semafoarelor către TM; acesta va comunica cu aplicația și cu serverul de comunicații *numai prin zona de memorie partajată* în cadrul tranzacției pe care o va gestiona.

(3.3) Procesul fiu transmite un mesaj TM-BEGIN-TRANSACTION către serverul de comunicații (CS) care, la rândul său, creează un fiu al cărui scop va fi tratarea tranzacției respective. Procesul fiu va contacta managerul central pentru obținerea priorității tranzacției.

(3.4) După obținerea priorității, procesul fiu al CS-ului, confirmă începerea tranzacției trimițând un mesaj ACK-TM-BEGIN-TRANSACTION fiului TM-ului.

(3.5) Fiul TM-ului confirmă începerea tranzacției trimițând aplicației un mesaj ACK-TM-BEGIN-TRANSACTION.

(3.6) Toate cererile de efectuare a unor operații asupra obiectelor sunt trimise numai prin memoria partajată.

(3.7) Prima cerere dintr-o tranzacție către un obiect este defalcată la nivelul serverului de comunicații în două cereri: o primă cerere transmisă într-un canal creat static de către obiect (în care așteaptă cereri de implicare) *chan-OB* și cererea propriu-zisă transmisă în memoria partajată.

(3.8) Pentru obiectele aflate pe alte noduri, se angajează un dialog între serverele de comunicație: cel de pe nodul aplicației și cel de pe nodul obiectului. Este sarcina serverului de comunicație de pe nodul obiectului să creeze o zonă de memorie partajată cu managerul obiectului când se implică obiectul respectiv în tranzacție.

(3.9) După primirea unei cereri de implicare a unui obiect într-o tranzacție, managerul obiectului creează un proces fiu care va trata toate cererile pentru acel obiect provenind din tranzacția respectivă.

(3.10) La încheierea tranzacției, aplicația transmite (tot prin memoria partajată) mesajul END-TRANSACTION sau ABORT-TRANSACTION.

(3.11) Procesul fiu creat de TM pentru supervizarea tranzacției este cel care asigură înfaptuirea celor două faze ale protocolului *two-phase commit*.

(3.12) După terminarea tranzacției, procesele fii create de TM, CS, OBM se distrug.

(M4) După defectarea unui nod și repunerii lui în funcțiune trebuie luate în considerare două situații:

(4.1) nodul în cauză este cel pe care se afla coordonatorul de tranzacție

(4.2) nodul în cauză este un nod pe care se aflau numai obiecte.

După prezentarea sumară de mai sus a modului de comunicație între componentele sistemului DOSTP în cele patru momente distincte din funcționare, rezulta câteva concluzii în ceea ce privește structurile de date ale sistemului:

(1) procesele componente ale sistemului vor avea două tipuri de canale de comunicație: (A) *canale statice*, fixate din momentul compilării, pe care se transmit mesajele de la inițierea unei tranzacții (3.1, 3.2, 3.3), implicarea unui nou obiect într-o tranzacție (3.7, 3.8), conectarea unui nod la sistem (M1), inserarea sau

eliminarea unui obiect din sistem (M2) și (B) *canale dinamice*, create în timpul execuției, în momentul când procesele fii ale serverelor de comunicație de pe noduri inițiază dialogul în cadrul unei tranzacții.

(2) Într-un nod, pe timpul desfășurării unei tranzacții, mesajele se transmit printr-o zonă de memorie partajată la care accesul se face controlat de semafoare, după modelul *producator-consumator*. Astfel, dacă, de exemplu, un proces fiu al lui TM care supravezează tranzacția T1 (pe scurt, TM-T1) trebuie să transmită un mesaj procesului fiu al lui CS, care supravezează T1 (CS-T1), atunci el execută o operație Up() asupra semaforului asociat lui CS-T1, deblocând procesul CS-T1 care, eventual, aștepta la semafor. În cele ce urmează, se prezintă structurile de date globale ale sistemului:

```
enum kind = (APL_BEGIN_TRANSACTION, END_TRANSACTION, REQUEST, FIRST_REQUEST,
TRANSACTION_ACCEPTED, TRANSACTION_ENDED, ACK_APL_BEGIN_TRANSACTION,
CS_END, APL_END, ABORT, COMMIT, ACK_ABORT, ACK_COMMIT, CONNECT, DISCONNECT,
RECORD, DELETE, GET_TRANS_ID, GET_OBJECT, ACK_TM_BEGIN_TRANSACTION,
TM_BEGIN_TRANSACTION, TM_END, BEGIN_TRANSACTION, TRANSACTION_INTERRUPTED,
INTERRUPT, TRANSACTION_NOT_ACCEPTED, TRANS_ID_ACK, ACK_RECORD, OBJ_FOUND,
UNKNOWN, ACK_REQUEST, ACK_FIRST_REQUEST, RECORD_END)

enum errors = (EBADF, EDOSTP_INITIALIZE_DONE, ECREATE_APL_FIFO,
EDOSTP_INITIALIZE_EXPECTED, EBEGIN_TRANS_DONE, EOPEN_APL_FIFO,
EAPL_OPEN_TM_FIFO, EWRITE_TM_FIFO_FAIL, EREAD_APL_FIFO_FAIL,
EAPL_ACCESS_SHM_FAIL, EINVAL_TRANS, ENO_BUF, ETIME_OUT, NETW_ERROR,
ESYNCRONIZE_EXPECTED, ECOMMIT_FAIL, ETM_CREATE_SHM, ETM_ACCESS_SHM_FAIL,
ETM_CREATE_SEM, ETM_OPEN_TM_FIFO, ETOO_MANY_OBJECTS, EOBJ_NOT_FOUND)

const MAX_SLOTS = ... # numarul maxim de pachete din zona de memorie partajata
const MAX_OBJ_LEN = ... # dimensiunea maxima a numelui unui obiect
const MAX_APL_LEN = ... # dimensiunea maxima a numelui unei aplicatii
const S_FIRST_UNUSED = 6
const SEM_NO = ... # numarul de semafoare alocate obiectelor dintr-o tranzactie
const MAX_TRY = ... # nr. max. de incercari pentru implicarea unui nod la distanta
const MAX_OBJ_NO = SEM_NO - S_FIRST_UNUSED # nr maxim de obiecte dintr-o tranzactie
const TIME_OUT_APL_TM_BT = ... # valori de time-out
const TIME_OUT_APL_TM_SREQ = ...
const TIME_OUT_APL_TM_AREQ = ...
const TIME_OUT_TM_CS_BTACK = ...
const TIME_OUT_REQ = ...
const TIME_OUT_RES = ...

type SHMEM = record of # structura zonei de memorie partajata
(sbtTable[1:MAX_SLOTS]: PACKET, # zona de pachete
freePackList[1:MAX_SLOTS]: int, # lista de pachete libere si ocupate
freePackListIdx: int, # index în lista de pachete; indica, la dreapta pachete
# libere, iar la stanga pachete ocupate
aplList[1:MAX_SLOTS]: int, # coada de pachete trimise operajiei
aplListIdx: int, # index în coada de pachete trimise aplicatiei;
transmanList[1:MAX_SLOTS]: int, # coada de pachete trimise managerului de tranzactii (fiu)
transManListIdx: int, # indice în coada de pachete de mai sus
commServList[1:MAX_SLOTS]: int, # coada de pachete trimise serverului de comunicatii fiu
commServListIdx: int, # indice în coada de pachete de mai sus
objectList[1:MAX_OBJ][1:MAX_SLOTS]: int, # cozile de pachete trimise managerului de obiecte
objectListIdx: int # indice în cozile de mai sus
)

type PACKET = record of # structura unui pachet din zona de memorie partajata
(type: int, # tipul pachetului
errorCode: int, # codul erorii
transId: int, # identificatorul tranzactiei în care este implicat obiectul
priority: int, # prioritatea tranzactiei
seqNo: int, # numar de utilizare al slotului
aplSlot: int, # numarul slotului utilizat de aplicatie
aplSlotUsage: int, # numarul slotului aplicatiei
objId: int, # utilizat intern de TM pentru a identifica obiectele
```

Partea a II a: Criterii de proiectare a sistemelor de programare distribuită bazate pe obiecte.

```
objName [1:MAX_OBJ_LEN]: char, # numele unui obiect
text : REQUEST # cererea propriu-zisa, structura ei nu are importanta decat
#pentru aplicatie si managerul de obiecte
)
type BT_PACKET = record of # structura unui pachet transmis la initierea unei tranzactii
(type: int, # tipul mesajului
transId: int, # identificatorul tranzactiei
pri: int, # prioritatea tranzactiei
aplName[1:MAX_APL_LEN]: char # numele aplicatiei
)
type BT_ACK_PACKET = record of # structura unui pachet transmis la initierea unei tranzactii
(type: int, # tipul mesajului
errCode: int, # codul de eroare
transId: int, # identificatorul tranzactiei
pri: int, # prioritatea tranzactiei
shmlId, semId: int, # identificatori pentru memoria partajata si pentru semafoare
objName[1:MAX_OBJ_LEN]: char, # numele obiectului (folosit pentru recuperare obiectelor dintr-o tr.)
aplName[1:MAX_APL_NAME]: char # numele aplicatiei
)
type RECOVER = record of # mesaje folosite de obiecte la recuperare
(type: int, # tipul mesajului
transNo: int, # identificatorul tranzactiei
objName[1:MAX_OBJ_LEN]: char, # numele obiectului (folosit pentru recuperare obiectelor dintr-o tr.)
)
typedef CD_PACKET = record of # mesaje folosite introducerea sau eliminarea unui nod din sistem
(type: int # CONNECT, DISCONNECT sau GET_OBJECT
host: ADDR # adresa nodului asociat
port: int # portul asociat la CS
)
typedef RD_OBJECT = record of # mesaje folosite introducerea sau eliminarea unui obiect din sist.
(type: int # RECORD sau DELETE
host: ADDR # adresa nodului asociat
port: int # portul asociat la CS obiectului
pri: int # insa folosit si la implicarea unui nod
aplName[1:MAX_APL_NAME]: char # numele aplicatiei
objName[1:MAX_OBJ_LEN]: char, # numele obiectului
)
type BT_OBJECT = record of # structura unui pachet transmis la implicarea unui obiect
(type: int, # tipul mesajului
errCode: int, # codul de eroare
transId: int, # identificatorul tranzactiei
pri: int, # prioritatea tranzactiei
shmlId, semId: int, # identificatori pentru memoria partajata si pentru semafoare
semNo: int # numarul semaforului din set rezervat pentru obiect
aplName[1:MAX_APL_NAME]: char # numele aplicatiei
)
typedef ACK_ID_TRANS = record of # mesaje folosite pentru determinarea prioritatii si
(type: int # identificatorului unei prioritati
tranId: int # identificator
pri: int, # numele obiectului
)
# canale statice
chan chanf_APL(BT_ACK_PACKET) # canal static pentru mesajele destinate aplicatiei
chan chanf_TM(BT_PACKET) # canal static pentru mesajele destinate managerului de tranzactii
chan chanf_CS(kind; # canal static pentru mesaje la serverul de comunicatii:
case TM_BEGIN_TRANSACTION: BT_ACK_PACKET # TM_BEGIN_TRANSACTION de la TM
case RECORD, DELETE: RD_PACKET) # RECORD, DELETE de la managerii de obiecte locale
chan chans_CS() # canal static pentru mesajele destinate serverului de comunicatii
chan chanu_CS() # canal static pentru mesajele destinate serverului de comunicatii
chan chans_CM() # canal static pentru mesajele destinate managerului central
```



```
chan chanu_NS() # canal static pentru mesajele destinate serverului de nume
chan chanu_OBJ(BT_OBJECT) # canal static cu mesaje la initierea și terminarea tranzacțiilor
# canale dinamice pentru comunicare prin memorie partajată
chan chan_S_APL (PACKET) #canal dinamic pentru mesajele destinate aplicației
chan chan_S_TM (PACKET) # canal dinamic pentru mesajele destinate man. de tranzacții fiu
chan chan_S_CS (PACKET) #canal dinamic pentru mesajele destinate serv. de comunicații fiu
# canale dinamice tip socluri unix
chan chan_U_TM (RECOVER) # canal dinamic pentru mesajele destinate man. de tranzacții fiu
chan chan_U_CS (BT_ACK_PACKET) #canal dinamic pentru mesajele destinate serv. de comunicații fiu
chan chan_U_NS (CD_PACKET)
#canale dinamice tip soclu inet
chan chan_K_CS(kind; # canal dinamic (soclu inet) pentru mesaje la serverul de comunicații:
  case NETWORK_BEGIN_TRANSACTION: BT_ACK_PACKET # NETWORK de la TM
  case RECORD_DELETE: RD_PACKET) # RECORD, DELETE de la managerii de obiecte locale
chan chan_K_CM(kind; # canal static (soclu inet) pentru mesaje la serverul de comunicații:
  case NETWORK_BEGIN_TRANSACTION: BT_ACK_PACKET # NETWORK de la TM
  case RECORD_DELETE: RD_PACKET) # RECORD, DELETE de la managerii de obiecte locale
chan chan_KO_CS (BT_OBJECT), chan_KT_CS (BT_OBJECT) # pt comunicare CS-CS
```

```
var slotsTab[1: MAX_SLOTS]: PACKET # tabela de sloturi pentru pachetele comunicate prin mem. partajată
var S_SHM, S_APL, S_TM, S_CS, S_Child, S_Free_Slot: semaphore # semafoare pentru acces la mem. part.
```

9.4.2. Integrarea proceselor de aplicație în sistemul DOSTP

Setul de funcții prezentat la 9.1 realizează integrarea proceselor de aplicație în sistemul DOSTP, integrare transparentă pentru utilizatori. Aceste funcții sunt incluse într-o bibliotecă cu care se linkedează aplicațiile. Ele pot fi descrise astfel:

```
const MAX_TAB = MAX_SLOTS - 2 # numărul maxim de cereri asincrone care pot fi emise de aplicație
const SYNC_RES_SLOT = MAX_TAB # slot rezervat pentru cererile sincrone
type REQUEST_PACKET = record of # cerere de efectuare a unei operații asupra unui obiect
  (empty : bool # pachetul este sau nu este folosit
  seqNo : int # numărul cererii în cadrul tranzacției
  errorCode : pointer to int # memorează adresa la care va fi returnat codul de eroare
  resAdr : pointer to char # memorează adresa la care va fi returnat răspunsul la cerere
  slotUsage : int # contor pentru folosirea pachetului
  timeout : time # timpul după care nu se mai așteaptă răspunsul
  )
var reqTab[1:MAX_TAB]: REQUEST_PACKET # tabela de cereri asincrone
var validDOSTP: bool = FALSE # indică dacă aplicația s-a conectat la DOSTP
var validTrans: bool = FALSE # indică dacă o tranzacție este activă
var alreadyTrans: bool = FALSE # indică dacă aplicația a inițiat o tranzacție
var oldTrSucc: bool = true # # o nouă prioritate (de la managerul central) pentru o tranzacție nouă
var oldPri = 0 # prioritatea ultimei tranzacții executate
var asyncReq = 0: int # numărul cererii în cadrul tranzacției
var seqNo: int # numerele de secvență ale pachetelor
var pkt: PACKET # pentru trimiterea pachetelor în memoria partajată
var errorCode : int # reține codurile de eroare returnate de funcții
var errReqCode: int # reține codul de eroare pentru operații
var slot: int # pentru a reține numărului locației pachetului în memoria partajată
```

```
InitDOSTP[1:MAX_APL_LEN]: char) ::
  if name = sir nul → return EBADF
  if validDOSTP → return EDOSTP_INITIALIZE_DONE
  generează numele aplicației în aplName
  creează canalul static de comunicare cu numele aplName.FIFO și obține identitatea lui în chanf_APL
  if chanf_APL nu a putut fi creat → return ECREATE_APL_FIFO
  validDOSTP := true;
```

Partea a II a: Criterii de proiectare a sistemelor de programare distribuită bazate pe obiecte.

return SUCCES

BeginTransaction()::

```
var btPckt: BT_PACKET # pachet trimis la initierea unei tranzactii
var btAckPckt: BT_ACK_PACKET # pachet trimis pentru confirmarea initierii unei tranzactii
if validDOSTP = FALSE return EDOSTP_INITIALIZE_EXPECTED
if alreadyTrans → return EBEGIN_TRANS_DONE
fa i = 0 to MAX_TAB + 1 → reqTab [i].empty := 1; reqTab[i].slotUsage := 0; af # initializari
genereaza numele managerului de tranzactii in TMName.FIFO
deschide canalul chanf_TM (cu numele TMName.FIFO) pentru comunicatia aplicatie -> TM
if chanf_TM nu a putut fi deschis → return EOPEN_TM_FIFO
deschide canalul chanf_APL pentru comunicatia TM -> aplicatie
if chanf_APL nu a putut fi deschis → return EOPEN_APL_FIFO
btPckt.type := APL_BEGIN_TRANSACTION ;
if oldTrSucc = FALSE → btPckt.pr := -1 # prioritate noua pentru tranzactie
[] oldTrSucc = true → btPckt.pr := oldPri # daca se reia o tranzactie se retine vechea prioritate
fi
if ((errCode := send chanf_TM(btPckt)) ≠ SUCCESS) → return EWRITE_TM_FIFO_FAIL
programeaza un semnal de alarma TIME_OUT_APL_TM pentru receptia pe canalul chanf_APL
if ((errCode := receive chanf_APL(btAckPckt) ≠ SUCCESS) return EREAD_APL_FIFO_FAIL
inchide canalele chanf_APL, chanf_TM
if btAckPckt.errCode ≠ SUCCESS → return btAckPckt.errCode
[] btAckPckt.errCode = SUCCESS →
    oldPri := btAckPckt.pri; semId := btAckPckt.semId; shmId := btAckPckt.shmId
    if nu se poate realiza atasarea la memoria partajata → return EACCESS_SHM_FAIL
    [] se poate realiza atasarea la memoria partajata →
        validTrans := true; alreadyTrans := true; asyncReq := 0; # initializari
        return SUCCESS
fi
fi
```

**SyncRequest (objectName [1: MAX_OBJ_NAME], sourceAddr: pointer to char,
destAddr: pointer to char, timeOut: int)::**

```
if validDOSTP = FALSE → return EDOSTP_INITIALIZE_EXPECTED
if validTrans = FALSE → return EINVALID_TRANS
if objectName = sir nul → return EBADF
if asyncReq > MAX_TAB → return ENO_BUF
if objectName ≠ sir nul → pckt.type := REQUEST; pckt.objName := objectName fi
pckt.seqNo := seqNo; pckt.aplSlot := SYNC_RES_SLOT; pckt.text := sourceAddr
pckt.aplSlot Usage := reqTab [SYNC_RES_SLOT].SlotUsage;
se completeaza o cerere reqTab [SYNC_RES_SLOT] cu valorile: .resAdr := destAddr, .seqNo := seqNo,
.perrorCode := errReqCode, .empty := FALSE, .timeOut := timeOut
send chan_S_TM(slot)
programeaza un semnal de alarma pentru receptia pe chan_S_APL pentru timeOut sec
again: errCode := receive chan_S_APL(slot);
if errCode = SUCCES →
    if slotsTab[slot].aplSlot = SYNC_RES_SLOT →
        if reqTab [SYNC_RES_SLOT].slotUsage = slotsTab[slot].aplSlotUsage →
            i = slotsTab[slot].aplSlot
            completeaza la reqTab [i].resAdr si reqTab [i].perrorCode
            cu slotsTab[slot].text, slotsTab[slot].errCode # s-a rezolvat raspunsul
            elibereaza slot # pentru cererea sincrona
            reqTab [i].slotUsage := reqTab [i].slotUsage + 1
        [] reqTab [SYNC_RES_SLOT].slotUsage ≠ slotsTab[slot].aplSlotUsage →
            elibereaza slot # pachet sincron intirziat; se ignora
            goto again
        goto again
```

```

    fi
[] slotsTab[slot], aplSlot ≠ SYNC_RES_SLOT →
    i = slotsTab[slot], aplSlot
    if reqTab[i].empty = FALSE and reqTab[i].slotUsage = slotsTab[slot].aplSlotUsage →
        completeaza la reqTab[i].resAdr si reqTab[i].perrorCode
        cu slotsTab[slot].text, slotsTab[slot].errCode # s-a rezolvat raspunsul
        asyncReq := asyncReq - 1 # pentru cererea asincrona
        reqTab[i].slotUsage := reqTab[i].slotUsage + 1
    fi
    elibereaza slot
    goto again
fi
[] errCode = TIME_OUT →
    reqTab[SYNC_RES_SLOT].SlotUsage = reqTab[SYNC_RES_SLOT].SlotUsage + 1
    seqNo := seqNo + 1; reqTab[SYNC_RES_SLOT].empty := true
    reqTab[SYNC_RES_SLOT].perrorCode := TIME_OUT
    return TIME_OUT
fi
seqNo := seqNo + 1
if reqTab[SYNC_RES_SLOT].perrorCode SUCCESS → return COMMIT_ERROR
[] return SUCCESS

```

AsyncRequest (objectName [1: MAX_OBJ_NAME]: char, sourceAddr: pointer to char,
destAddr: pointer to char, perrCode: pointer to int, timeOut: int)::

```

if validDOSTP = FALSE → return EDOSTP_INITIALIZE_EXPECTED
if validTrans = FALSE → return EINVAL_TRANS
if objectName = sir nul → return EBADF
if asyncReq > MAX_TAB → return ENO_BUF
if nu exista o intrare libera i in reqTab → return ENO_BUF
if objectName ≠ sir nul → pkt.type := REQUEST; pkt.objName := objectName fi
pkt.seqNo := seqNo; pkt.aplSlot := i; pkt.text := sourceAddr;
pkt.aplSlotUsage := reqTab[i].SlotUsage;
se completeaza o cerere reqTab[i] cu valorile:
.resAdr := destAddr, .seqNo := seqNo, .perrorCode := errReqCode, .empty := 1, .timeOut := timeOut
send chan_S_TM(slot)
seqNo := seqNo + 1; asyncReq := asyncReq + 1
return SUCCESS

```

Synchronize (::):

```

if validDOSTP = FALSE → return EDOSTP_INITIALIZE_EXPECTED
if validTrans = FALSE → return EINVAL_TRANS
if asyncReq = 0 → return SUCCESS
determina in maxTimeOut maximum dintre toate valorile reqTab[i].timeOut
if maxTimeOut > timpul curent →
    programeza un semnal de alarma peste maxTimeOut - timpul curent
    again: errCode := receive chan_S_APL(slot);
    if errCode ≠ TIME_OUT →
        i = slotsTab[slot], aplSlot
        if reqTab[i].empty = FALSE and reqTab[i].slotUsage = slotsTab[slot].aplSlotUsage →
            completeaza la reqTab[i].resAdr si reqTab[i].perrorCode
            cu slotsTab[slot].text, slotsTab[slot].errCode # s-a rezolvat raspunsul
            asyncReq := asyncReq - 1 # pentru cererea asincrona
            reqTab[i].slotUsage := reqTab[i].slotUsage + 1
        fi
        elibereaza slot
        goto again

```

```
fi
fi
fa i:=1 to MAX_TAB st reqTab[i].empty = FALSE
    reqTab[i].empty = true ; reqTab [i].slotUsage:= reqTab [i].slotUsage+1
    reqTab [i].perrorCode := TIME_OUT; asyncReq := asyncReq - 1
af
return SUCCESS
```

EndTransaction()::

```
if validDOSTP = FALSE → return EDOSTP_INITIALIZE_EXPECTED
if validTrans = FALSE → return EINVALID_TRANS
if asyncReq > 0 → return ESINCRONIZE_EXPECTED
pckt.type := END_TRANSACTION
errCode = SyncRequest_S(NULL, pckt, pckt, TIME_OUT_APL_TM_SREQ);
if errCode = ETIME_OUT oldTrSucc := FALSE
[] errCode ≠ ETIME_OUT oldTrSucc := true
fi
if errCode ≠ SUCCESS →
    if errCode = ETIME_OUT → return ETIME_OUT
    [] errCode ≠ ETIME_OUT → errCode = ECOMMIT_FAIL
pckt.type := APL_END
errCode = AsyncRequest_S(NULL, pckt, pckt, errCode, TIME_OUT_APL_TM_AREQ);
alreadyTrans := FALSE;
if errCode ≠ SUCCESS →
    if errCode = ETIME_OUT → return ETIME_OUT
    [] errCode ≠ ETIME_OUT → errCode = ECOMMIT_FAIL
return SUCCESS
```

AbortTransaction()::

```
if validDOSTP = FALSE → return EDOSTP_INITIALIZE_EXPECTED
if validTrans = FALSE → return EINVALID_TRANS
pckt.type := ABORT_TRANSACTION
errCode = AsyncRequest_S(NULL, pckt, pckt, errCode, TIME_OUT_APL_TM_AREQ);
alreadyTrans := FALSE; oldTrSucc := FALSE;
inchide canalele chanf_APL, chanf_TM
return SUCCESS
```

CloseDOSTP()::

```
if validDOSTP = FALSE → return EDOSTP_INITIALIZE_EXPECTED
inchide canalele chanf_APL, chanf_TM; sterge chanf_APL
validDOSTP := FALSE
return SUCCESS
```

9.4.3. Proiectarea managerilor de tranzactii

Aplicațiile trimit cererile de operații obiectelor prin intermediul managerilor de tranzacții. Un manager de tranzacție creează câte un proces fiu, TM-T_i pentru fiecare tranzacție *i* pe care o inițiază o aplicație. Un manager de tranzacție se constituie deci în mai multe procese: procesul principal TM și procesele fiu, care supraveghează tranzacțiile. Procesul TM așteaptă într-o buclă continuă recepționarea mesajelor de inițiere ale tranzacțiilor de la aplicații, și recepționarea mesajelor de acceptare a unei tranzacții sau terminare a unei tranzacții de la procesele fiu TM-T_i. Procesul TM ține evidența tranzacțiilor pastrind liste interne cu

correspondența identificator tranzacție - nume aplicație; listele sunt actualizate la începerea sau terminarea unei tranzacții. Identificatorul de tranzacție, ca și prioritatea tranzacțiilor, se obțin prin contactarea managerului central. Pe baza priorității unei tranzacții, managerii de obiecte decid, în cazul unui conflict de acces, care dintre tranzacțiile în conflict se va termina cu succes.

Procesul TM recepționează mesaje într-un singur canal, *chan-TM*, creat static, la lansarea sa în execuție. Cererile de operații asupra obiectelor sunt recepționate de procesele *TM-Ti* în zona de memorie partajată creată de *TM-Ti* la lansarea în execuție; mecanismul de recepție a unui mesaj din zona de memorie partajată va fi reprezentat mai jos, în descrierea managerilor de tranzacție, prin recepționarea unui mesaj din canalul guvernat de semaforul asociat procesului, *chan-S-TM [Ti]*. Procesele *TM-Ti* constau într-o buclă principală în care se recepționează mesaje din zona de memorie partajată (cereri de operații de la aplicație sau răspunsuri de la cereri anterioare, răspunsuri trimise de serverul de comunicație sau de managerii de obiecte locale).

Trebuie precizat că, la lansarea în execuție, într-o fază de inițializare, înainte de a intra în bucla de recepționare a mesajelor, managerul de tranzacții (TM) analizează fișierele sale de log pentru a detecta dacă execuția sa are loc după o terminare anormală, survenită în urma defectării nodului. La nivelul managerului de tranzacții, TM, există două fișiere log: primul fișier (TMREC1) memorează identificatorii și starea tranzacțiilor, iar cel de-al doilea, (TMREC2) memorează pentru fiecare tranzacție numele obiectelor implicate în tranzacție și starea obiectului în tranzacție. O tranzacție se poate găsi în fișierul TMREC1 în una din următoarele patru stări: (1) tranzacția nu a ajuns la *EndTransaction()* sau *AbortTransaction()* (WORKING); (2) s-a executat faza *prepare to commit*, dar nu și *commit* pentru tranzacție (PREPARE-TO-COMMIT) și deci nu se cunosc răspunsurile managerilor de obiecte implicați în tranzacție; (3) s-a executat faza COMMIT, în urma primirii răspunsurilor afirmative de la toți managerii de obiecte (COMMIT); (4) s-a executat faza ABORT (ABORT). Stările posibile ale obiectelor din tranzacții, înregistrate în fișierul TMREC2 sunt: (1) obiectul recepționează cereri și le tratează (BUSY); (2) obiectul a răspuns afirmativ la un mesaj PREPARE-TO-COMMIT (OK); (3) obiectul a răspuns negativ la un mesaj PREPARE-TO-COMMIT, constatând un conflict cu o altă tranzacție mai prioritară, FAIL.

După ce managerii de obiecte au răspuns cu ACK-ABORT sau ACK-COMMIT intrările corespunzătoare tranzacției respective din TMREC1 și TMREC2 se invalidează (se modifică numărul tranzacției la zero, număr invalid, lasând celelalte informații neschimbate).

În faza de inițializare, în urma analizei fișierului TMREC1, managerul de tranzacții TM creează câte un proces pentru fiecare tranzacție a cărei intrare este neștearsă; procesul fiu, pe baza informațiilor din TMREC2 și în funcție de starea tranzacției va asigura recuperarea tranzacției astfel:

(1) Pentru o tranzacție găsită în starea WORKING se trimite la toți managerii de obiecte (găsiți în TMREC2) comanda ABORT.

(2) Pentru o tranzacție găsită în starea PREPARE trebuie verificate răspunsurile managerilor de obiecte implicați în tranzacție (din fișierul TMREC2); dacă toți managerii au răspuns OK, atunci se trimite comanda COMMIT; dacă cel puțin un manager a răspuns cu FAIL atunci se trimite ABORT.

(3) Pentru o tranzacție găsită în starea COMMIT se transmite COMMIT.

(4) Pentru o tranzacție găsită în starea ABORT se transmite ABORT.

Dacă, pe de altă parte, în timpul desfășurării tranzacției apare o defecțiune a unui nod care conține un obiect implicat în tranzacție, atunci serverul de comunicații detectează acest lucru și sesizează managerul de comunicații. Acesta va consulta fișierele de recuperare și va decide acțiunea de executat astfel: pentru tranzacțiile în starea COMMIT se transmite COMMIT, pentru tranzacțiile în orice stare diferită de COMMIT se transmite ABORT. Managerii de tranzacții pot fi descriși astfel:

```
enum TRANS_STATE = (ABORT, WORKING, PREPARE_TO_COMMIT, COMMIT) # stările unei tranzacții
enum OBJ_RESPONSE = (BUSY, OK, FAIL) # răspunsurile unui obiect la prepare-to-commit
type TMREC1_ART = record of # structura unui articol din fișierul TMREC1
    (transNo: int,
     transState: TRANS_STATE,
     transName[1: MAX_APL_LEN]: char
    )
type TMREC2_ART = record of # structura unui articol din fișierul TMREC2
    (transNo: int,
     objResp: OBJECT_RESPONSE,
     objName[1: MAX_OBJ_LEN]: char
    )
type OBJECTS = record of
    (free: bool, # obiect nefolosit cind free este true
     errCode: int,
     objTransId: int,
     ackPrToCom, ackCom, ackAbort: int
     objName[1: MAX_OBJ_LEN]: char
    )
```

```
)  
TransactionManager::  
var TMREC1: file of TMREC1_ART # fisierele de recuperare  
var TMREC2: file of TMREC2_ART  
var mes: BT_PACKET # mesaj primit pe canalul chanf_TM  
var S_Rec1, S_Rec2: semaphore := (1,1) # semafoare pentru acces exclusiv la fisierele log  
var semId: int # identificatorul setului de semafoare  
var mes: BT_ACK_PACKET # pentru transmitia mesajelor pe canale statice  
var semIdR: semaphore # folosit pentru contorizarea fiilor care au efectuat recuperarea obiectelor  
var shmlD: int # identificator al zonei de memorie partajata  
var objTab: OBJECTS[1:MAX_OBJ_NO] # tabela cu informatii despre obiectele implicate in tranzactie  
var objNo: int # numarul de obiecte din tranzactie  
var previous_abort: int  
  
# faza de initializare  
programeaza tratarea semnalului SIGCHLD  
if nu se pot deschide semafoarele S_Rec1, S_Rec2 identificate prin semId pentru controlul accesului  
la fisierele de recuperare → Panic()  
fi  
if nu se poate crea si deschide canalul static pentru receptia mesajelor chanf_TM → Panic() fi  
if nu se pot deschide fisierele de recuperare → Panic() fi  
# faza recuperare dupa repunerea in functiune a managerului de tranzactii  
do not eof (TMREC1)  
if exista in fisierul TMREC1 o tranzactie neinvalidata →  
lanseaza in executie un proces fiu (fork+execlp) TransactionManagerServer caruia i se transmite:  
transName, semId, transNo, transState, 0  
fi  
od  
# faza de tratare a cererilor de la aplicatii  
do true →  
receive chanf_TM(mes)  
if mes.type = APL_BEGIN_TRANSACTION →  
lanseaza in executie un fiu (fork+execlp) TransactionManagerServer caruia i se transmite:  
mes.aplName, semId, 0, 0, mes.pri,  
|| mes.type = TRANSACTION_ACCEPTED →  
insereza tranzactia in lista de tranzactii active; pentru o tranzactie se memoreaza:  
mes.aplName, mes.pri, mes.transId,  
|| mes.type = TRANSACTION_ENDED → elimina tranzactia din lista de tranzactii active  
fi  
od  
  
TransactionMangerServer(aplName[MAX_APL_LEN]:char, semId: int, transNo: int, transState, pri: int)::  
var waiting: bool # folosita pentru a indica validitatea tranzactiei  
var childNo: int # contorizeaza numarul de procese fii create pentru recuperare  
var recover: RECOVER # pentru mesajele primite de la obiecte la recuperare  
var rec: int, errCode: int # variabila de lucru  
var pckt, pcktr: PACKET  
var transId: int; priority: int, objId: int, status: int # variabile de lucru  
var seqNoEnd: int, aplSlotUsage: int, aplSlot: int # variabile de lucru  
var ackPrNo: int, csFlag: int, aplFlag: int, resp: int # variabile de lucru  
  
if nu se poate determina din fisierul de configurare adresa managerului central → Panic() fi  
programeaza tratarea semnalelor SIGCHLD si SIGALRM  
if nu se poate deschide canalul de comunicatie cu CS chanf_CS → Panic() fi  
if transNo ≠ 0 → # faza de recuperare a tranzactiilor ramase neterminate de la ultima executie  
if nu se poate deschide semaforul semIdR initializat la 0 → Panic() fi  
if transState = ABORT or transState = WORKING →  
modifica starea tranzactiei din fisierul TMREC1 la ABORT  
if nu se poate deschide fisierul TMREC2 → Panic() fi
```

```

    childNo := 0
    do cit timp exista un obiect objName implicat in tranzactia transNo →
        childNo := childNo + 1
        TransactionRecovery(ABORT, transNo, pri, objName, semIdR)
    od
    do cit timp valoarea semaforului semIdR ≠ childNo → od # se asteapta terminarea tuturor fiilor
    invalideaza in TMREC1 tranzactia transId
[] transState = PREPARE_TO_COMMIT →
    if nu se poate deschide fisierul TMREC2 → Panic() fi
    if toate obiectele tranzactiei transId din TMREC2 au raspuns OK →
        rec:= COMMIT; modifica starea tranzactiei in TMREC1 la COMMIT
    [] nu toate obiectele tranzactiei transId din TMREC2 au raspuns OK →
        rec := ABORT; modifica starea tranzactiei in TMREC1 la ABORT
    childNo := 0
    do cit timp exista un obiect objName implicat in tranzactia transNo →
        childNo := childNo + 1
        TransactionRecovery(rec, transNo, pri, objName, semIdR)
    od
    do cit timp valoarea semaforului semIdR ≠ childNo → od # se asteapta terminarea tuturor fiilor
    invalideaza in TMREC1 tranzactia transId
[] transState = COMMIT →
    childNo := 0
    do cit timp exista un obiect objName implicat in tranzactia transNo →
        childNo := childNo + 1
        TransactionRecovery(COMMIT, transNo, pri, objName, semIdR)
    od
    do cit timp valoarea semaforului semIdR ≠ childNo → od # se asteapta terminarea tuturor fiilor
    invalideaza in TMREC1 tranzactia transId
fi
exit
fi
# faza de initializare
if nu se poate deschide canalul static pentru comunicatia cu aplicatia chanf_APL → Panic() fi
if nu se poate crea o zona de mem. partajata cu id. shmId → FatalError(ETM_CREATE_SHM) fi
if eroare la atasarea la zona de mem. partajata la adr. pShm → FatalError(ETM_ACCESS_SHM_FAIL) fi
if nu se poate crea setul de sem. semId pt. acc. la mem. partajata → FatalError(ETM_CREATE_SEM) fi
se initializeaza tabela de sloturi
if nu se poate deschide canalul de comunicatie chanf_TM → FatalError(ETM_OPEN_TM_FIFO) fi
mes.type := TM_BEGIN_TRANSACTION; mes.errCode := SUCCESS; mes.aplName := aplName
mes.shmId := shmId; mes.semId := semId; send chanf_CS(mes)
programeaza un semnal de alarma pentru receptia pe chan_S_MT la TIME_OUT_TM_CS_BTACK
errCode := receive chan_S_MT(slot)
if errCode = TIME_OUT → FatalError(ECS_TIME_OUT) fi
if slotsTab[slot].errCode ≠ SUCCESS → FatalError(slotsTab[slot].errCode) fi
mes.type := ACK_APL_BEGIN_TRANSACTION; send chanf_APL(mes) # raspuns la aplicatie
transId := slotsTab[slot].transId; priority:= slotsTab[slot].priority
mes.type := TRANSACTION_ACCEPTED; mes.transId:= transId; mes.pri:= priority
send chanf_TM(mes)
scrie in fisierul TMREC1 un articol pentru tranzactia initiata, avind cimpurile:
    .transNo:= transId; transState:= WORKING; .transName:= aplName;
fa i= 1 to MAX_OBJ_NO objTab.free := true # initializeaza tabela de obiecte
waiting = true
# faza de receptionare a cererilor de operatii si a raspunsurilor
do waiting = true →
receive chan_S_TM(slot); pcktr:= slotTab[slot]
if pcktr.type = REQUEST→
    if pcktr.objName exista in pozitia i in objTab[i] objId := i →
    [] objId := -1
    fi

```



```

if objId ≠ -1 → # nu este prima cerere pentru acest obiect in tranzactie
    slotTab[slot].transId:= transId; slotTab[slot].objId:= objId; send chan_S_CS(slot)
[] objId = -1 → # este prima cerere pentru acest obiect in tranzactie
    if nu mai exista loc in tabela de obiecte objTab →
        slotTab[slot].errCode:= ETOO_MANY_OBJECTS; send chan_S_APL(slot)
    [] inereaza obiectul in objTab intr-o pozitie libera objId
        objNo := objNo +1; slotTab[slot].:= FIRST_REQUEST
        slotTab[slot].transId:= transId; slotTab[slot].objId:= objId; send chan_S_CS(slot)
        Down(S_Rec2)
        scrie in fisierul TMREC2 un articol pentru noul obiect cu cimpurile:
            .objName:= slotTab[slot].objName, .transId:= transId, resp:= BUSY
        Up(S_Rec2)
    fi
fi
[] pcktr.type = ACK_FIRST_REQUEST→
    send chan_S_APL(slot);
    if slotTab[slot].errCode = EOBJ_NOT_FOUND →
        sterge din objTab intrarea obiectului slotTab[slot].objName; objNo := objNo - 1
        Down(S_Rec2)
        sterge din fisierul TMREC2 articolul obiectului slotTab[slot].objName
        Up(S_Rec2)
        alocu un nou slot slot1 in slotsTab; slotTab[slot1].type:= ABORT
        AbortTransaction(slot1); # forteaza terminarea anormala a tranzactiei
[] pcktr.type = ACK_REQUEST→
    send chan_S_APL(slot);
[] pcktr.type = ABORT_TRANSACTION→
    AbortTransaction(slot);
[] pcktr.type = ACK_ABORT_TRANSACTION→
    objTab[pcktr.objId].ackAbort := 1
    sterge din objTab intrarea obiectului slotTab[slot].objName; objNo := objNo - 1
    Down(S_Rec2);sterge din TMREC2 articolul obiectului slotTab[slot].objName; Up(S_Rec2)
    elibereaza slot
    if objNo = 0 →
        aloca slot1; slotTab[slot1].type := END_TRANSACTION; send chan_S_CS(slot1)
        again:errCode := receive chan_S_TM(slot)
        if errCode = SUCCESS and slotsTab[slot].type ≠ CS_END →
            elibereaza slot; goto again:
        [] errCode = TIME_OUT→
            aloca slot; slotTab[slot].type := ABORT_TRANSACTION
            AbortTransaction(slot)
    fi
    End()
[] Up(S_Child)
fi
[] pcktr.type = END_TRANSACTION→
    seqNoEnd:= pcktr.seqNo; aplSlotUsage:= pcktr.aplSlotUsage; aplSlot:= pcktr.aplSlot;
    if objNo ≠ 0 →
        if freePackListIdx ≠ 1 → End() fi #in mod normal trebuie sa existe numai pachetul End
        elibereaza slot;
        # faza prepare_to_commit
        initializeaza pentru toate obiectele objTab cimpurile .ackPrToCom :=0, .errCode := 0
        initializeaza semaforul S_Child la MAX_SLOTS - 1; ackPrNo := 0;
        se creaza un proces fiu care transmite mesajul la toate obiectele; el executa:
            Down(S_Rec1);
            modifica in TMREC1 starea tranz. .transState:= ABORT;
            Up(S_Rec1)
            fa obiectele i st objTab[i].free = false →
                Down(S_Child); aloca slot1; slotTab[slot1].type := PREPARE_TO_COMMIT

```

```

        slotTab[slot1].objId:= i;slotTab[slot1].objTransId:= transId;slotTab[slot1].errCode=0
slotTab[slot1].objName:= objectTab[i]; send chan_S_CS (slot1)
    af
        exit
[] objNo = 0 →
    status:= 0
    fa obiectele i st objTab[i].free = false →
        if objTab[i].errCode ≠ SUCESS → status := objTab[i].errCode fi
    af
        slotTab[slot].errCode := status; slotTab[slot].seqNo= seqNoEnd;
        slotTab[slot].aplSlotUsage:= aplSlotUsage; slotTab[slot].aplSlot:= aplSlot;
        send chan_S_APL(slot)
        aloca slot1; slotTab[slot1].type := END_TRANSACTION; send chan_S_CS (slot1)
        csFlag := 0;
        do csFlag=0 →
            receive chan_S_TM(slot)
            if slotTab[slot].type = CS_END → csFlag=1
            [] slotTab[slot].type = APL_END → aplFlag=1
            elibereaza slot

        od

        do aplFlag=0 →
            receive chan_S_TM(slot)
            if slotTab[slot].type = APL_END → aplFlag=1
            elibereaza slot

        od
    End()
[] pcktr.type = ACK_PREPARE_TO_COMMIT →
    objTab[pcktr.objId].ackPrToCom := 1; resp:= OK;
    if pcktr.errCode ≠ SUCESS → objTab[pcktr.objId].errCode := pcktr.errCode
    [] objTab[pcktr.objId].errCode := text.errCode
    fi
    if objTab[pcktr.objId].errCode ≠ SUCESS → resp := FAIL fi
    ackPrNo := ackPrNo +1
Down(S_Rec2);modifica in fis TMREC2 articolul pentru obiect la cimpul .resp:= resp; Up(S_Rec2)
    if resp ≠ OK → modifica in fisierul TMREC1 articolul pentru obiect in cimpul .tranState:= ABORT fi
    elibereaza slot
    if ackPrNo = objNo →
        if toate obiectele au raspuns cu OK → status := OK
        [] nu toate obiectele au raspuns cu OK → status := FAIL
        if status ≠ OK →
            rezerva slot
            slotTab[slot].errCode := status; slotTab[slot].seqNo= seqNoEnd;
            slotTab[slot].aplSlotUsage:= aplSlotUsage; slotTab[slot].aplSlot:= aplSlot;
            send chan_S_APL(slot)
            aloca slot1; slotTab[slot1].type := ABORT_TRANSACTION;
            AbortTransaction(slot1)
        [] if status = OK → # faza de commit
            initializeaza pentru toate obiectele objTab cimpurile .ackPrToCom :=0, .errCode := 0
            initializeaza semaforul S_Child la MAX_SLOTS - 1
            se creaza un proces fiu care transmite mesajul la toate obiectele; el executa:
            Down(S_Rec1);modifica in TMREC1 starea tranz. .tranState:=ABORT; Up(S_Rec1)
            fa obiectele i st objTab[i].free = false →
                Down(S_Child); aloca slot1; slotTab[slot1].type := COMMIT
                slotTab[slot1].objId:= i;slotTab[slot1].objTransId:= transId;slotTab[slot1].errCode=0
                slotTab[slot1].objName:= objectTab[i]; send chan_S_CS (slot1)
            af
                exit
        fi
    fi

```

```

[] Up(S_Child)
[] pcktr.type = ACK_COMMIT →
  objNo := objNo - 1; elibereaza slot
  if objNo = 0 →
    if toate obiectele au raspuns cu OK → status := OK
    [] nu toate obiectele au raspuns cu OK → status := FAIL
    fa toate obiectele din objTab st objTab.free = false →
      elibereaza intrarea
      Down(S_Rec2); sterge din TMREC2 articolul obiectului Up(S_Rec2)
    af
    rezerva slot; slotTab[slot].errCode := status; slotTab[slot].seqNo= seqNoEnd;
    slotTab[slot].aplSlotUsage:= aplSlotUsage; slotTab[slot].aplSlot:= aplSlot;
    send chan_S_APL(slot)
    aloca slot1; slotTab[slot1].type := END_TRANSACTION;
    send chan_S_CS(slot)
    csFlag := 0;
    do csFlag=0 →
      receive chan_S_TM(slot)
      if slotTab[slot].type = CS_END → csFlag=1
      [] slotTab[slot].type = APL_END → aplFlag=1
      elibereaza slot
    od
    do aplFlag=0 →
      receive chan_S_TM(slot)
      if slotTab[slot].type = APL_END → aplFlag=1
      elibereaza slot
    od
  End()
[] Up(S_Child)
[] pcktr.type = APL_END → aplFlag := 1
[] pcktr.type = CS_END → aplFlag := 1 waiting := false
[] pcktr.type = INTERRUPT →
  status := starea tranzactiei citita din TMREC1
  if status ≠ COMMIT → pcktr.type := INTERRUPT_ABORT
  [] status = COMMIT → pcktr.type := INTERRUPT_COMMIT
  fi
  AbortTransaction(slot)
fi
od
AbortTransaction(slot: int)::
var type: int
  type := slotTab[slot].type; elibereaza slot
  if type = INTERRUPT_COMMIT → transState:= COMMIT
  [] transState:= ABORT;
  fi
  Down(S_Rec1);
  modifica in fisierul TMREC1 articolul pentru tranzactie la starea transState
  Up(S_Rec1);
  programeaza un semnal de alarma peste TM_ABORT_TIME
  if previous_abort = 0
    previous_abort := previous_abort + 1
  if objNo = 0 →
    aloca slot1; slotTab[slot1].type := END_TRANSACTION; send chan_S_CS(slot1)
    again:errCode := receive chan_S_TM(slot)
    if errCode = SUCCESS and slotsTab[slot].type ≠ CS_END →
      elibereaza slot; goto again:
    [] errCode = TIME_OUT →

```

```

        alocă slot; slotTab[slot].type := ABORT_TRANSACTION
        AbortTransaction(slot)
    fi
    End()
[] objNo ≠ 0 →
    initializează toate cimpurile .errCode și .ackAbort din objTab la 0
    initializează semaforul S_Child la MAX_SLOTS - 1
    se creează un proces fiu care transmite mesajul la toate obiectele; el execută:
        fa obiectele i st objTab[i].free = false →
            Down(S_Child); alocă slot1; slotTab[slot1].type := type;
            slotTab[slot1].objId:= i; slotTab[slot1].objTransId:= transId;
            slotTab[slot1].objName:= objectTab[i]; send chan_S_CS (slot1)
        af
        exit
    fi
fi
End():
    invalidează în fisierul TMREC1 articolul pentru tranzacția transId
    se închid semafoarele semId, se eliberează zona de memorie partajată pShm
    mes.type := TRANSACTION_ENDED; send chanf_TM(mes)
    exit
FatalError(errCode: int)::
    mes.type := ACK_APL_BEGIN_TRANSACTION; mes.errCode := errCode
    send chanf_APL(mes)
    se închid semafoarele semId, se eliberează zona de memorie partajată pShm
    exit
Recovery(mes.type: int, transNo: int, pri: int, objName: [1:MAX_OBJ_LEN]: char, semIdR: int)::
    receive chan_U_TM (recover)
    if recover.type = ACK_ABORT or recover.type = ACK_COMMIT
        Down(S_Rec2); invalidează în TMREC2 articolul corespunzător lui transNo, objName; Up(S_Rec2)
    fi
    Up(semIdR)
    închide chan_U_CS, chan_U_TM
    exit()
TransactionRecovery(type: int, transNo: int, pri: int, objName: [1:MAX_OBJ_LEN]: char, semIdR: int)::
    generează mes de tip BT_ACK_PACKET cu cimpurile:
        .type := type, .transNo := transNo, .objName := objName, .aplName := aplName
    if nu se poate obține o conexiune chan_U_CS, chan_U_TM de la canalul chan_U_CS → Panic() fi
    send chan_U_CS (mes)
    lansează în execuție un fiu care execută Recovery(mes.type, transNo, pri, objName, semIdR)
    închide chan_U_CS, chan_U_TM
    return

```

9.4.4. Proiectarea serverelor de comunicații

În fiecare nod conectat la sistemul DOSTP se execută câte un server de comunicații. Acesta are două scopuri principale: (1) *asigura destășurarea tranzacțiilor*. La inițierea unei tranzacții, managerul de tranzacții contactează serverul de comunicații; acesta creează un proces server fiu, care va gestiona exclusiv comunicațiile tranzacției respective. Pentru fiecare obiect de la distanță implicat în aceasta tranzacție, procesul server fiu va crea un alt proces care va recepționa mesajele de la obiectul de la distanță. Serverul de comunicație creat pentru a deservi o tranzacție se folosește de serviciile serverului de nume pentru a determina localizarea obiectelor în sistem. Dacă obiectul este local, transferul de mesaje se execută prin zona de memorie partajată. (2) *asigura notificarea sistemului în cazul pierderii legăturii cu un nod participant la o tranzacție*. Acțiunile unui server de comunicații care detectează pierderea legăturii cu un nod participant la o tranzacție depind de poziția managerului de tranzacții care a contactat serverul de comunicații. Dacă acesta se află pe același nod, atunci serverul de comunicații anunță managerul de tranzacții că s-a pierdut legătura cu un nod implicat în tranzacție și așteaptă decizia lui (procesul fiu al serverului de comunicații care dialoga cu nodul respectiv este terminat); în caz contrar, serverul de comunicații anunță managerul de obiect că s-a pierdut legătura cu nodul coordonator și se așteaptă decizia coordonatorului, la repunerea sa în funcțiune. (vezi 10.2.)

Descrierea SR a serverului de comunicații este:

Partea a II a: Criterii de proiectare a sistemelor de programare distribuită bazate pe obiecte.

```
const MAX_PROC = ... # numarul maxim de obiecte de la distanta care pot fi implicate
const LOCAL = 1; const DISTANCE = 0 # indicatori pentru localizarea tranzactiilor si a obiectelor
typedef OBJECT = record of # informatii despre obiectele implicate intr-o tranzactie
    (sem: int # numarul semaforului din setul semId asociat unei SHM
     local: int # obiect local sau nu
     chan: int # canalul (soclul) asociat pentru obiectele de la distanta
     name[1:MAX_OBJ_LEN]: char # numele obiectului
    )

CommunicationServer():
var semList[1:SEM_NO]: int # numerele semafoarelor din setul semId asociate accesului la zona partajata
var semListIdx: int := FIRST_UNUSED # index in lista semList
var objTab[1:MAX_OBJ_NO]: OBJECT # retine localizarea obiectelor dintr-o tranzactie
var semIdR: semaphore # semafor folosit la recuperare
var mes: CD_PACKET # pachet trimis la conectare CS la sistem
var typ: int # primeste tipul mesajelor
var pckRD: RD_PACKET # pentru mesaje de la managerii de la obiecte
var pckBTAck: BT_ACK_PACKET # pentru mesaje de la managerul central
var addrHost: ADDR # adresa nodului
var transErr: bool := false # indica faptul ca nu mai este valida nici o cerere pt. tranzactie
var aplName[1:MAX_APL_NAME]: char # variabile pentru memorarea informatiilor caracteristice tranz.
var shmId, semId: int # variabile de lucru care memoreaza identif. pentru SHM
var pidTab[1: MAX_PROC]: int # identificatorii proceselor fii create

# faza de initializare
if nu se poate determina din fisierul de configurare canalul managerului central → Panic() fi
programeaza tratarea semnalului SIGCHLD
fa i = 1 to SEM_NO semList[i] := i # initializarea numerelor semafoarelor din set
fa i = 1 to MAX_OBJ_NO → objTab[i].local := 0; objTab[i].chan := -1 fi # initializarea tabelii de obiecte
if nu se poate crea si deschide canalul de comunicatie propriu chanf_CS → Panic() fi
if nu se poate crea si deschide canalul de comunicatie propriu chans_CS → Panic() fi
if nu se poate crea si deschide canalul de comunicatie propriu chanu_CS → Panic() fi
if nu se poate crea sem. semIdR pentru recuperare, cu valoarea initiala 0 → Panic() fi
if nu obtine adresa nodului curent la addrHost → Panic() fi
# faza de conectare la serverul central
completeaza un mesaj mes cu cimpurile .type=CONNECT, .host= addrHost, .port = portul fixat pentru CS
send chans_CM(mes)
# receptie si tratare cereri locale pe chanf_CS
if not empty (chanf_CS) → receive chanf_CS(typ);
    if typ = RECORD or typ = DELETE →
        receive chanf_CS(pckRD);
        creaza un proces fiu care executa:
            obtine o conexiune chan_U_CS, chan_U_NS de la chanu_NS; send chan_U_NS(pckRD);
            obtine o conexiune chan_K_CM, chan_K_CM de la chans_CM; send chan_K_CM(pckRD); exit
    [] typ = TM_BEGIN_TRANSACTION →
        receive chanf_CS(pckBTAck);
        creaza un proces fiu care executa: aplName := pckBTAck.aplName;
        instaleaza o rutina pentru SIGALRM; Services(LOCAL, pckBTAck); exit
    fi
fi
# receptie si tratare cereri de la distanta pe chans_CS
if not empty (chans_CS) →
    stabileste o conexiune chan_K, chan_K_CS; receive chan_K_CS(typ);
    if typ = RECORD or typ = DELETE →
        receive chan_K_CS(pckRD);
        creaza un proces fiu care executa:
            obtine o conexiune chan_U_CS, chan_U_NS de la chanu_NS; send chan_U_NS(pckRD); exit
    [] typ = RECORD_END →
        receive chan_K_CS(pckRD);
        creaza un proces fiu care executa:
```

```

    obtine o conexiune chan_U_CS, chan_U_NS de la chanu_NS; send chan_U_NS(pckRD);
    receive chan_U_NS(pckRD); Up (semldR); exit
[] typ = NETWORK_BEGIN_TRANSACTION →
    receive chan_K_CS(pckBTack);
    creaza un proces fiu care executa: aplName := pckBTack.aplName;
    instaleaza o rutina pentru SIGALRM; Services(DISTANCE, pckBTack); exit
[] typ = ABORT or typ = COMMIT →
    receive chan_K_CS(pckBTack);
    creaza un proces fiu care executa:
        Down (semldR); # asteapta sa se termine inregistrarea obiectelor;
        Up (semldR); Recover(pckBTack, chan_K); exit
fi
fi
# receptie si tratare cereri locale pe chanu_CS
if not empty (chanu_CS) →
    stabileste o conexiune chan_U_CS, chan_U_TM; receive chan_U_CS(typ);
    if typ = ABORT or typ = COMMIT →
        receive chan_U_CS(pckBTack);
        creaza un proces fiu care executa:
            Down (semldR); # asteapta sa se termine inregistrarea obiectelor;
            Up (semldR); Recover(pckBTack, chan_U_TM); exit
    fi
fi
Recover(data: BT_ACK_PACKET, chan_R: chan):
var mes: RD_OBJECT; dt: BT_OBJECT
    copiaza cimpurile comune de la data la dt
    completeaza mes cu valorile .type:= GET_OBJECT, .objName:= data.objName
    obtine o conexiune chan_U_CS, chan_U_NS de la chanu_NS;
    send chan_U_NS(mes); receive chan_U_CS(mes);
    if mes.type ≠ UNKNOWN →
        if addrHost = mes.addr → # obiect local
            obtine o conexiune chan_U_CS, chan_U_OBJ de la chanu_OBJ pentru mes.objName;
            send chan_U_OBJ(dt); receive chan_U_CS(dt); send chan_R(dt);
        [] addrHost ≠ mes.addr → # obiect la distanta
            incerca pina se reuseste stabilirea unei conexiuni chan_KT_CS, chan_KO_CS
                de la CS-ul din nodul mes.addr
            send chan_KT_CS(dt); receive chan_KO_CS(dt); send chan_R(dt);
        fi
    fi
fi
return
Services(typ:int, data: BT_ACK_PACKET):
var pA: BT_ACK_PACKET; pckt: PACKET; aut: RD_PACKET; # variabile de lucru
var errCode: int, slot: int, res: ACK_ID_TRANS, waiting:bool # variabile de lucru
var bto: BT_OBJECT;
    if typ = LOCAL → # cerere de la o aplicatie locala
        copiaza data la pA; aplName:= data.aplName; semld:= data.semld; shmlid:= data.shmlid;
        if nu se poate realiza atasarea la memoria partajata → return
        aut.type:=GET_TRANS_ID; aut.pri:= pA.pri;
        obtine o conexiune chan_K_CS, chan_K_CM de la chans_CM; send chan_K_CM(aut);
        if (errCode:= receive chan_K_CM(res)=TIME_OUT→
            pckt.type := INTERRUPT; aloca slot pt pckt; send chan_S_TM(slot); exit;fi
        [] errCode=SUCCES →
            pckt.type :=ACK_TM_BEGIN_TRANSACTION; pckt.pri:= res.pri; pckt.errCode := SUCCES
            pckt.transId := res.trans.Id; aloca slot pt pckt; send chan_S_TM(slot)
        fi
        waiting:= true
        do waiting = true →
            receive chan_S_CS(slot);
            TreatMessage(slot, LOCAL); # trateaza un mesaj de la o aplicatie locala

```

```
od
[] typ = DISTANCE →
if nu se poate crea si atasa la o zona de mem. partajata shmld sau semaf. semlD return
completeaza cimpurile bto cu: .type:=BEGIN_TRANSACTION, .transId:=data.transId,
.shmld:=data.shmld, .semlD:=data.semlD, .semNo:=FIRST_UNUSED, .queueNo:=0
.aplName:= data.aplName
obține o conexiune chan_U_CS, chan_U_OBJ de la chanu_OBJ pentru data.objName;
send chanu_U_OBJ(bto); waiting:= true
do waiting = true →
instaleaza un handler pentru SIGALRM
programeaza un semnal de alarma pentru TIME_OUT_REQ
if (errCode:= receive chan_K_CS(pckt)= SUCCESS)→
TreatMessage(slot, DISTANCE);
if slotsTab[slot].type = ACK_ABORT or slotsTab[slot].type = ACK_COMMIT→
waiting = true; elibereaza zona de memorie partajata
fi
[] errCode ≠ SUCCESS →
genereaza un slot cu slotsTab[slot]=INTERRUPT;
send chan_S_OBJ(slot)
again: receive chan_S_CS(pckt);i
if pckt.type ≠ ACK_INTERRUPT→ goto again
elibereaza zona de memorie partajata si semafoarele; exit
fi
od
fi
return
TreatMessage(slot: int, where:int)::
var bto: BT_OBJECT; mes: RD_OBJECT ; wp:PACKET
pckt:= ref(slotsTab[slot]);
if where=LOCAL → # mesaj de TM local
if pckt.type = FIRST_REQUEST → # intii determina localizarea obiectului
completeaza mes cu valorile .type:= GET_OBJECT, .objName:= pckt.objName
obține o conexiune chan_U_CS, chan_U_NS de la chanu_NS;
send chan_U_NS(mes); receive chan_U_CS(mes);
if mes.type = UNKNOWN →
transErr:= true; objTab[pckt.objId].name:= ""
completeaza slot cu cimpurile .type:= ACK_FIRST_REQUEST
.errCode:= EOBJ_NOT_FOUND; transId:= -1; send chan_S_TM(slot); return
[] if mes.type = OBJ_FOUND →
objTab[pckt.objId].name:= pckt.objName; # obiect local
if addrHost = mes.host →
completeaza cimpurile bto cu:
.type:=BEGIN_TRANSACTION, .transId:=pckt.transId,
.shmld:=shmld, .semlD:=semlD,
.semNo:= semList[semListIdx];semListIdx:=semListIdx+1
.queueNo:=semNo-FIRST_UNUSED; .aplName:= aplName
obține o conexiune chan_U_CS, chan_U_OBJ
de la chanu_OBJ pentru mes.objName;
send chanu_U_OBJ(bto); waiting:= true
send chan_S_OBJ[semNo](pckt);
objTab[pckt.objId].sem:= bto.semNo; objTab[pckt.objId].local:=1
[] addrHost ≠ mes.host # obiect la distanta
objTab[pckt.objId].local:=0;
if dupa MAX_TRY nu se poate obține o conex. cu chans_CS de adr. mes.addr→
transErr:= true; objTab[pckt.objId].name:= ""
completeaza slot cu cimpurile .type:= ACK_FIRST_REQUEST
.errCode:= EOBJ_NOT_FOUND; transId:= -1; send chan_S_TM(slot);
return
[] se poate obține o conex. chan_KT_CS,chan_KO_CS
objTab[pckt.objId].socket:= chan_KD_CS
```



```

mes.type:= NTEWORK_BEGIN_TRANSACTION;
mes.transId := pkt.transId; mes.pri:= pkt.pri;
errCode := send chan_KO_CS (bto)
if errCode ≠ SUCCESS →
    pkt.type:= INTERRUPT; send chan_S_TM(slot);
    return
[] errCode = SUCCESS
    creaza un proces fiu pid pentru obiectul la distanta din tranz. locala
    inchide toate canalele dinamice cu exceptia chan_KD_CS
    instaleaza o functie pentru tratarea semnalului SIGALRM
    do true
        programeaza un semnal de alarma pentru
        errCode:= receive chan_KT_CS(wp);
        if errCode=SUCCES →
            aloca slot; slotsTab[slot]:=wppkt;
            send chan_S_TM(slot);
            if w.type=ACK_ABORT or w.type = ACK_COMMIT→
                exit
            [] errCode ≠SUCCES →
                pkt.type:= INTERRUPT; send chan_S_TM(slot);
            inchide canalele; exit
        fi
    od
    pidTab[pkt.objId]:=pid
    fi
fi
fi
[] pkt.type = REQUEST or pkt.type = PREPARE_TO_COMMIT or pkt.type = COMMIT→
    if transErr=false → Transmit(slot, false)
[] pkt.type = ABORT→ Transmit(slot, false)
[] pkt.type = INTERRUPT_ABORT → pkt:= ABORT;Transmit(slot, true)
[] pkt.type = INTERRUPT_COMMIT → pkt:= COMMIT;Transmit(slot,true)
[] pkt.type = ACK_FIRST_REQUEST →
    if objTab[pkt.objId].local ≠0 →
        if pkt.transId=-1→semList[semListIdx]:= objTab[pkt.objId].sem; semListIdx:=semListIdx-1
        fi
    fi
    send chan_S_MT(pkt)
[] pkt.type = ACK_REQUEST or pkt.type = PREPARE_TO_COMMIT
    or pkt.type =ACK_COMMIT or pkt.type =ACK_ABORT →
    send chan_S_MT(pkt)
[] pkt.type = END_TRANSACTION →
    pkt.type := CS_END; send chan_S_MT(pkt); waiting := false;
fi
[] where = DISTANCE →
    if pkt.type = FIRST_REQUEST or pkt.type=PREPARE_TO_COMMIT or
    pkt.type = ABORT or pkt.type = COMMIT→
        send chan_S_OBJ(); receive chan_S_CS(slot)
        errCode:=send chan_KT_CS (slotsTab[slot])
        if errCode ≠ SUCCESS →
            rezerva un slot; slotsTab[slot].type:= INTERRUPT; send chan_S_OBJ(slot);
            again: receive chan_S_CS(slot);
            if slotsTab[slot].type ≠ INTERRUPT → goto again fi
            elibereaza zona de memorie partajata
        fi
    fi
return
Transmit(slot:int, abort: bool):: # se apeleaza pentru a transmite comenzile locale tranzactiei la distanta
var bto: BT_OBJECT, mes: RD_OBJECT, sem: int, # variabile de lucru
    pkt:= ref(slotsTab[slot]);

```

```
if objTab[pckt.objId].objName = Ø → # obiect cu nume nul
  if pckt.type = ABORT →
    genereaza un slot cu .type:=ACK_ABORT, .objName:=pckt.nameObj; .objId:=pckt.objId
    send chan_S_TM(slot)
    return
  fi
fi
if objTab[pckt.objId].local = 0 → # obiect local, pe acelasi nod cu tranzactia
  sem:= objTab[pckt.objId].sem; send chan_S_OBJ(pckt) #transmite pachetul la ob. local prin SHM
[] if abort → # situatia cind nodul obiectului nu mai raspunde
  inchide canalul vechi objTab[pckt.objId].chan
  distruge eventualul proces care mai asteapta la receive pt. ca CS-ul pereche sa forteze INTERRUPT
  identificatorul procesului este memorat in pidTab[pckt.objId]
  bto.type:=pckt.type; bto.transId:=pckt.transId;
  creaza un proces fiu care executa:
    inchide toate canalele din objTab[].chan
    completeaza mes cu valorile .type:= GET_OBJECT, .objName:= pckt.objName
    obtine o conexiune chan_U_CS, chan_U_NS de la chanu_NS; # intii determina
    send chan_U_NS(mes); receive chan_U_CS(mes); # localizarea obiectului
    instaleaza cite handler pentru tratare SIG_PIPE
    incearca periodic, la infinit, obtinerea unei conexiuni
      chan_KT_CS, chan_KO_CS, cu CS-ul mes.host
    if conexiunea s-a putut stabili →
      send chan_KO_CS(bto); receive chan_KT_CS(bto);
      pckt.type:=bto.type; send chan_S_TM(slot);
    fi
  fi
[] abort = false →
  creaza un proces fiu care executa: # este necesar un fiu pentru cazul cind nodul obi. nu mai rasp.
  send objTab[pckt.objId].chan(pckt) # se transmite SIGPIPE si este distrus numai acest fiu
fi
return
```

9.4.5. Proiectarea serverelor de nume

Rolul serverelor de nume este de a menține o listă cu identitatea obiectelor conectate în sistem și localizarea lor. Un server de nume este contactat de un server de comunicații în trei situații:

(1) când se apelează prima operație asupra unui obiect nou implicat în tranzacție, pentru ca serverul de comunicații să cunoască unde va plasa cererea (ca și următoarele pentru același obiect): local, în memoria partajată sau contactînd serverul de comunicații de pe nodul obiectului;

(2) la conectarea sau eliminarea unor obiecte în sistem, cînd managerul central transmite modificarea efectuată;

(3) la conectarea nodului în sistem, cînd managerul central transmite identitatea și localizarea tuturor obiectelor integrate pînă în acel moment în sistem.

Managerul central transmite informații serverelor de nume via serverele de comunicații.

Descrierea SR a serverelor de nume este:

```
NameServer():
var pcktRD, obj: RD_PACKET
var res:ACK_ID_TRANS, mes:BT_OBJECT
var hostAddr: ADDR
#faza de initializare
initializeaza lista obiectelor conectate ca vida
determina adresa nodului curent la hostAddr
#faza de primire si tratare cereri
do true →
  accepta o conexiune chan_U_S, chan_U_NM pe canalul chanu_NS
  receive chan_U_(typ)
  if typ = RECORD →
    receive chan_U_NS(pcktRD); inscrie noul obiect in lista obiectelor conectate
    if pcktRD.host = hostAddr →
      mes.type := ACK_RECORD; send chan_U_S(mes);
```

```
[] typ = DELETE →
    receive chan_U_NS(pktCD); sterge obiectul din lista obiectelor conectate
>[] typ = RECORD_END →
    receive chan_U_NS(pktCD); send chan_U_S(pktCD);
>[] typ = GET_OBJECT →
    receive chan_U_S(pktRD);
    if pktRD.objName exista in lista de obiecte obj →
        pktRD.type:= OBJ_FOUND; pktRD.host:=obj.host;pktRD.port:=obj.port;
    [] pktRD nu exista in lista de obiecte obj →
        pktRD.type:= UNKNOWN;
    send chan_U_S(pktRD);
    fi
fi
od
```

9.4.6. Proiectarea managerului central

Managerul central a fost introdus în vederea realizării următoarelor două scopuri: (1) ținerea (centralizată) a evidenței nodurilor conectate la sistem; (2) livrarea unui număr de prioritate unic tranzacțiilor inițiate în mediul distribuit.

Descrierea SR a managerului central este:

CentralManager():

var pktRD, obj: RD_PACKET

var pktCD, node: CD_PACKET

var res:ACK_ID_TRANS

#faza de initializare

if nu se poate deschide fisierul *TrIdPr* cu val. curenta de acordat id. de tranzactii si prioritatii → *Panic()*
initializeaza lista de noduri conectate si a obiectelor conectate ca vide

#faza de primire si tratare cereri

do true →

accepta o conexiune *chan_K1_CS, chan_K1_CM* pe canalul *chans_CM*

receive *chan_K_CM(typ)*

if *typ = CONNECT* →

receive *chan_K_CM(pktCD)*; inscrie noul nod in lista nodurilor conectate

fa *obj* din lista de obiecte introduse in sistem →

completeaza la *pktRD* cimpurile *.host:= obj.host, .type:=RECORD,*
.objNam:=obj.objName, .port:=obj.port

send *chan_K1_CS(pktRD)*;

af

pkt.type:=RECORD_END ; **send** *chan_K1_CS(pktRD)*;

[] *typ = DISCONNECT* →

receive *chan_K_CM(pktCD)*; sterge nodul din lista nodurilor conectate

[] *typ = RECORD* →

receive *chan_K_CM(pktRD)*; inscrie noul obiect in lista obiectelor conectate

fa *node* din lista de noduri introduse in sistem →

obtine o conexiune *chan_K2_CS, chan_K2_CM* de la canalul *chans_CS*
de pe nodul *node.host*

send *chan_K2_CS(pktRD)*;

af

inchide canalele *chan_K2_CS, chan_K2_CM*

[] *typ = DELETE* →

receive *chan_K_CM(pktRD)*; sterge obiectul in lista obiectelor conectate

fa *node* din lista de noduri introduse in sistem →

obtine o conexiune *chan_K2_CS, chan_K2_CM* de la canalul *chans_CS*
de pe nodul *node.host*

send *chan_K2_CS(pktRD)*;

af

inchide canalele *chan_K2_CS, chan_K2_CM*

[] *typ = GET_TRANS_ID* →

receive *chan_K_CM(pktRD)*;

citeste si apoi incrementeaza *transId* din fisierul local *TrIdPri*;

```

rescrie valoarea incremenat in fisier : res.transId:=transId;
if pktRD.pri = -1 →
    citeste si apoi incrementeaza pri din fisierul local TrIdPri;
    rescrie valoarea incremenat in fisier : res.pri:=pri;
[] pktRD.pri ≠ -1 → res.pri:=pktRD.pri
res.type:= TRANS_ID_ACK; send chan_K_CS(res)
fi
od

```

9.4.7. Proiectarea managerilor de obiecte

S-a arătat la 9.3.6. că managerii de obiecte încapsulează două tipuri de date (obiectele și informația de sincronizare) și patru tipuri de funcții care, în principal, tratează lucrul cu obiecte, sincronizarea accesului la obiecte, recuperarea obiectelor după defecte și integrarea obiectelor în sistem. Până în momentul de față, în DOSTP au fost implementați manageri de obiecte care folosesc o schemă de sincronizare optimistă și o tehnică de recuperare redo-only (adekvată pentru schemele optimiste).

Implementarea managerilor se face folosind tehnica programării pe obiecte (vezi 10.3.).

În cele ce urmează, se prezintă descrierea SR a unui manager de obiecte, fără a se lua în considerare modul de implementare obiectual.

ObjectManager::

```

type ACLOG_TYPE = record of
    (transId: int,
    pri: int,
    op: int,
    logPC: int)
type TSFILE_TYPE = record of
    (transId: int,
    pri: int)
var: serverName[1:MAX_APL_LEN], hostName[1:MAX_APL_LEN];
var: aut: RD_PACKET
var mes:BT_OBJECT, pkt:PACKET
var ts: int, pri:int, op:int, logPC:int
var semId:int, semNo:int, queueNo:int, shmId:int
var waiting: bool:=true, conn:int:=false

```

= structura articolelor din fisierul cu tranzactiile care
 = au executat prepare: identificatorul tranzactiei
 = prioritatea tranzactiei
 = operatia de executat dupa prepare: commit sau abort
 = deplasamentul in fisierul PCLOG
 = structura articolelor din fisierul TSFILE care retine
 = tranzactiile terminate: identificatorul tranzactiei
 = prioritatea tranzactiei
 = memoreaza nume
 = folosit la conectarea obiectului
 = mesaje receptionate si trimise
 = informatii din ACLOG pentru o tranzactie
 = informatii necesare pentru accesul la SHM
 = indicatori

faza de initializare

```

if nu se poate crea si deschide canalul de comunicatie propriea shm (obj) → Paru(x) fi
obține numele nodului curent la hostName si numele obiectului serverName
completeaza aut cu .type:=RECORD, .objName:=serverName, .host:=hostName
if nu se poate crea si deschide canalul de comunicatie cu chanf_CS → Paru(x) fi
send chanf_CS(aut);
if nu se poate deschide fisierul ACLOG in modul append → Paru(x) fi
if nu se poate deschide fisierul PCLOG in modul append → Paru(x) fi
if nu se poate deschide fisierul TSFILE in modul append → Paru(x) fi
seteaza modul de lucru fara buferare pentru fisierele ACLOG, PCLOG, TSFILE
if nu se poate deschide si mapa in memorie fisierul cu copia obiectului → Paru(x) fi
if nu se poate deschide si mapa in memorie fisierul cu informatiile de stare ale tranzactiilor → Paru(x) fi
if nu se poate crea semaforul semObj pentru accesul la copia partajata a obiectului in memorie → Paru(x) fi
if nu se poate crea semaforul semTs pentru accesul la fisierul de tranzactii TSFILE → Paru(x) fi

```

faza de recuperare initiala a obiectelor

```

pozitionare la inceputul fisierului ACLOG
fa articolele fisierului ACLOG →
    citeste cimpurile unui articol ACLOG_TYPE la ts, pri, op, logPC
    if ts>0 Recovery ts, pri, op, logPC; fi #realizeaza commit sau abort specific obiectului respectiv

```

af

```

creaza o alta copie pentru fisierul PCLOG avind inregistrările invaldate eliminate
# reface efectele efecte tranzactiilor ramase inca nerevenimate
pozitionare la inceputul fisierului PCLOG
fa articolele fisierului PCLOG →

```

citeste cimpurile unui articol la *ts, pri, op, ind, val*
completeaza *oldTrans[oldTransIdx]*: *.transNo:=ts, .pri:= pri, logPC:= pozitia curenta din PCLOG*
seteaza lock-urile pentru operatia *op, ind* respectiva

af
faza de asteptare a unui mesaj pentru implicarea intr-o intr-o tranzactie sau de terminare tranzactie
do true →

stabileste o conexiune *chan_U_OBJ, chan_U_CS* la canalul *chanu_OBJ*

receive *chan_U_OBJ (mes)*;

if *conn = false* →

if *mes.type = ACK_RECORD and mes.errCode = SUCCESS* → *conn:= true*;

[*conn = true* →

if *mes.type = TRANSACTION_ACCEPTED* →

insereaza tranzactia in lista de tranzactii active

[*mes.type = TRANSACTION_ACCEPTED* →

insereaza tranzactia in lista de tranzactii active

[*mes.type = TRANSACTION_ENDED* →

sterge tranzactia in lista de tranzactii active

[*mes.type = TRANSACTION_INTERRUPTED* →

sterge tranzactia in lista de tranzactii active

[*mes.type = BEGIN_TRANSACTION* →

ts:= mes.transId;

creaza un proces fiu care executa:

ts:= mes.transId; pr:=mes.pri, semid:=mes.semId; shmId:=mes.shmId;

queueNo:=mes.queueNo; semNo:= mes.semNo;

if nu se poate realiza atasarea la zona de memorie partajata →

mes.type:= TRANSACTION_NOT_ACCEPTED; send chanu_OBJ(mes); exit

fi

waiting:= true

mes.type:= TRANSACTION_ACCEPTED; send chanu_OBJ(mes)

do *waiting = true* →

receive *chan_S_OBJ(pkt)*;

if *pckt.type = FIRST_REQUEST*→

trateaza cererea seteaza *pckt.errCode*;

pckt.type := ACK_FIRST_REQUEST; send chan_S_CS(pkt)

[*pckt.type = REQUEST* →

trateaza cererea seteaza *pckt.errCode*;

pckt.type := ACK_REQUEST; send chan_S_CS(pkt)

[*pckt.type =PREPARE_TO COMMIT* →

trateaza cererea si seteaza *pckt.errCode*;

pckt.type := ACK_PREPARE_TO_COMMIT; send chan_S_CS(pkt)

[*pckt.type = COMMIT* →

trateaza cererea si seteaza *pckt.errCode*;

pckt.type := ACK_COMMIT or send chan_S_CS(pkt);

mes.type:= TRANSACTION_ENDED; send chanu_OBJ(mes); waiting := false

[*pckt.type = ABORT* →

trateaza cererea si seteaza *pckt.errCode*;

pckt.type := ACK_ABORT send chan_S_CS(pkt);

mes.type:= TRANSACTION_ENDED; send chanu_OBJ(mes); waiting := false

[*pckt.type = INTERRUPT* →

trateaza cererea si seteaza *pckt.errCode; send chan_S_CS(pkt)*;

mes.type:= TRANSACTION_INTERRUPTED; send chanu_OBJ(mes);

pckt.type:= ACK_INTERRUPT; send chan_S_CS(pkt); waiting := false;

fi

elibereaza memoria partajata

exit

[*mes.type = ABORT or mes.type = COMMIT* →

creaza un proces fiu care executa:

mes.transId:=trans.Id; op:=mes.type;

if *transId* exista in fisierul *TSFILE* → *errCode := RECOVERY_OK* # tranzactie terminata

```

[] transId nu exista in fisierul TSFILE
do cit timp transId exista exista in tabloul tranzactiilor in curs de recuperare→sleep(TIME)
  if starea tranzactiei din oldTrans = WORKING →
    mes.errCode := RECOVERING # recuperare deja in lucru pentru alt proces
  [] starea tranzactiei din oldTrans = DONE →
    mes.errCode := DONE # recuperare rezolvata
  [] starea tranzactiei din oldTrans = UNCHANGED →
    inregistreaza in oldTrans starea tranzactiei transId ca RECOVERING
    Recovery(ts, pri, op, logPC); inregistreaza in oldTrans starea tr. transId ca DONE
    mes.errCode := RECOVERY_OK
  [] tranzactia nu exista in oldTrans →
    mes.errCode := TRANS_NOT_FOUND
  fi
  if errCode = RECOVERY_OK→
    if op = ABORT mes.type := ACK_ABORT
    [] op = COMMIT mes.type := ACK_COMMIT
    fi
  fi
  send chan_U_CS_OBJ(mes)
  exit
od
fi
od
```

Capitolul 10.

IMPLEMENTAREA SISTEMULUI DOSTP ÎN REȚELE UNIX ȘI MICROSOFT WINDOWS (WIN32)

Capitolul de față arată cum se implementează sistemul DOSTP, a cărei proiectare a fost prezentată în capitolul 9, în rețele Unix și Microsoft Windows (Win 32).

10.1. Particularități de implementare

Subcapitolul 10.1 se referă la implementarea tuturor componentelor sistemului DOSTP, cu excepția managerilor de obiecte, a căror implementare este prezentată (pentru un tip de date particular) în subcapitolul 10.3.

10.1.1. Particularități de implementare în rețele Unix

Implementarea sistemului DOSTP în rețelele Unix s-a făcut urmărind îndeaproape proiectarea prezentată în descriere SR în capitolul 9. Se fac însă următoarele observații:

(1) Canalele de comunicație între serverele de comunicație ale nodurilor, pe de o parte, și, pe de altă parte între serverele de comunicație și serverul central au fost implementate prin socluri Inet. Pentru recepționarea și emisia mesajelor din/pe aceste canale au fost scrise două funcții speciale care, apelând (eventual de mai multe ori) funcțiile de sistem *read()* și *write()* asigură citirea/scrierea întregului mesaj (un anumit număr de octeți). Acest lucru este necesar ținând cont de specificul soclurilor TCP (fluxuri de octeți).

(2) Canalele de comunicație statice server de nume -> server de comunicație, aplicație -> manager de tranzații, manager de tranzații -> server de comunicații și manager de obiecte -> server de comunicație s-au implementat prin canale FIFO, iar canalele de comunicație server de comunicație -> manageri de obiecte, server de comunicație -> server de nume s-au implementat prin socluri Unix.

(3) Comunicația locală dintre aplicație <-> manager de tranzații, manager de tranzații <-> server de comunicații, server de comunicații <-> manager de obiecte se realizează printr-o zonă de memorie partajată, a cărei creare se face în managerul de tranzații astfel:

```
if ((shmld=shmget(IPC_PRIVATE, DIM_SHM, IPC_CREAT | S_IRUSR | S_IWUSR)) < 0)
    Panic ("TM->Error : Cannot creat shared memory");
if (pShm = (SHM*) shmat (shmld, 0, 0)) = (SHM*)-1)
    Panic ("TM->Error : Cannot attach to shared memory\n");
```

(4) Pentru accesul la zona de memorie partajată se folosesc semafoarele *S_APL* (de aplicație), *S_TM* (de managerul de tranzații), *S_CS* (de serverul de comunicații), în scopul recepționării pachetelor (emițătorul execută Up(semafor), iar receptorul Down(semafor)). În plus, mai sunt prevazute semafoarele: *S_SHM* pentru accesul exclusiv la lista de pachete libere din zona de memorie partajată (în momentul alocării sau eliberării unui pachet) și *S_FREE_SLOT* inițializat la numărul total de pachete disponibile în zona de memorie partajată, decrementat (Down(semafor)) la fiecare alocare a unui pachet și incremental la fiecare eliberare de pachet (astfel încât în situația că nu mai sunt pachete libere în zona, apelantul este blocat).

Crearea și inițializarea semafoarelor se face imediat după crearea zonei de memorie partajată astfel:

```
if ((semld=semget(IPC_PRIVATE, SEM_NO, IPC_CREAT | S_IRUSR | S_IWUSR)) < 0)
    Panic ("TM->Error : Cannot create semaphores");
if (semctl (semld, S_SHM, SETVAL, 1) == -1 ||
    semctl (semld, S_APL, SETVAL, 0) == -1 ||
    semctl (semld, S_TM, SETVAL, 0) == -1 ||
    semctl (semld, S_CS, SETVAL, 0) == -1 ||
    semctl (semld, S_FREE_SLOT, SETVAL, MAX_SLOTS) == -1)
)
    Panic ("TM -> Cannot initialize semaphores\n");
```

(5) Comunicația între două procese prin zona partajată a fost reprezentată în descrierea SR din subcapitolul 9.4, sub forma unor operații de recepție și emisie. Implementarea comunicației prin zona de

Partea a II a: Criterii de proiectare a sistemelor de programare distribuită bazate pe obiecte.

memorie partajată, de exemplu, în cazul managerului de tranzacții care poate trimite pachete atât aplicației cât și serverului de comunicații a fost realizată prin următoarea funcție:

```
void SendPacket (int dest, int slot) {
    WaitOnSemaphore (semId, S_SHM);
    switch (dest)
    {
    case APL:
        aplList[aplListIdx] = slot; /*se depune numarul pachetului care va fi transmis */
        aplListIdx=(aplListIdx+1) % MAX_SLOTS; /*înca un pachet depus în coada */
        signalSemaphore(semId,S_APL); /*se anunța aplicația */
        break;
    case CS:
        commServList[commServListIdx]=slot;
        commServListIdx=(commServListIdx+1) % MAX_SLOTS;
        SignalSemaphore(semId, S_CS);
        break;
    default:
        Debug("No such destination\n");
    }
    SignalSemaphore(semId, S_SHM);
}
```

(6) Rezervarea și dealocarea unui pachet din zona de memorie partajată a fost implementată astfel:

```
void FreeSlot(int slot) {
    WaitOnSemaphore(semId, S_SHM); /* acces exclusiv la zona partajată */
    freeSlotList[--freeSlotListIdx]=slot; /*pachetele eliberate se depun la stînga indexului */
    SignalSemaphore(semId, S_SHM); /* s-a eliberat zona */
    SignalSemaphore(semId, S_SLOT_FREE); /* s-a mai eliberat un pachet */
}
int GetSlot(void) {
    int slot;
    WaitOnSemaphore(semId, S_SLOT_FREE);/*se așteaptă eliberarea unui pachet */
    WaitOnSemaphore(semId, S_SHM); /* acces exclusiv la zona partajată */
    slot=freeSlotList[freeSlotListIdx++]; /*se ocupă un pachet */
    SignalSemaphore(semId, S_SHM); /* se eliberează zona */
    return slot;
}
```

(7) Managerii de tranzacții și serverele de comunicații dedicați unei tranzacții au fost implementați în procese separate create de managerul de tranzacții, respectiv serverul de comunicații. Aceștia există numai pe parcursul desfășurării unei tranzacții. Pentru a nu se genera procese *zombie* la terminarea proceselor fii, a trebuit captat și tratat semnalul SIGCHLD.

(8) Managerul de tranzacții, serverul de comunicații, serverul de nume și managerul central au fost realizați ca programe XWindows. Aceștia prezintă o interfață simplă care oferă și facilități pentru vizualizarea tranzacțiilor în lucru (managerul de tranzacții), a obiectelor conectate la sistem (serverul de nume și managerul central) și a nodurilor conectate la sistem (managerul central). Receptionarea mesajelor pe canalele statice se face asincron, în cadrul unor funcții declarate *XInputCallbackProc*.

(9) Pentru ca semnalele Unix, care sosesc asincron, să nu întrerupă un apel sistem *semop* ce punea în așteptare un proces la un semafor, implementarea unei funcții *WaitOnSemaphore()* trebuie făcută avîndu-se în vedere acest lucru, astfel:

```
void WaitOnSemaphore(int semId, int sem) {
    struct sembuf semBuf;
    int code;
    semBuf.sem_num=sem;
    semBuf.sem_op=-1;
    semBuf.sem_flg=0;
    again: errno=0;
    code=semop(semId, &semBuf, 1);
    if (code==-1) {
        if (errno==EINTR) goto again;
        if (errno==EIDRM || errno==EINVAL)
            Panic(" Access to semaphores lost\n");
    }
}
```

(10) Pentru detectarea pierderii legăturii cu un nod se captează și tratează semnalele SIGALARM și SIGPIPE.

10.1.2. Particularități de implementare în rețele Microsoft Windows (Win 32)

Pentru varianta Windows au fost implementate componentele DOSTP care trebuie să gestioneze un nod ce conține obiecte: serverul de comunicație, serverul de nume și managerii de obiecte (pentru managerii de obiecte vezi subcapitolul 10.3).

Ca și în varianta Unix, implementarea serverului de comunicație și a serverului de nume s-a făcut urmărind proiectarea prezentată la 9.4. Trebuie făcute însă următoarele observații:

(1) Canalele de comunicație server de comunicație <-> manageri de obiecte, server de comunicație <-> server de nume au fost implementate folosind numai socketuri Inet. Nu se mai pot folosi, ca în varianta Unix socketuri locale (specifice numai pentru Unix) sau canale FIFO.

(2) Comunicația locală, dintre serverul de comunicație și managerul de obiecte se face tot printr-o zonă de memorie partajată. Crearea ei se face astfel:

```
if (( hMmf = CreateFileMapping ((HANDLE)0xffffffff, NULL, PAGE_READWRITE, 0, sizeof(structSHM), szBuff)
    ==NULL)
    Panic ("CS->Error : Cannot creat shared memory"); );
if ((*pShm = (struct SHM*)MapViewOfFile( hMmf, FILE_MAP_WRITE, 0, 0, sizeof(struct SHM))) == NULL )
    Panic ("CS->Error : Cannot map shared memory");
```

(3) Pentru accesul la zona de memorie partajată se folosesc analog semafoarele descrise la punctul (4) de la 10.1.1 dar asocierea lor, în mod unic, la o zonă de memorie partajată se face asigurându-le un nume unic în sistem astfel:

```
printf( szBuff, "SHM%s%d", objName, transId);
if ( NULL == (*phSemSHM = CreateSemaphore( NULL, 1, 1, szBuff))
    Panic ("CS->Error : Cannot create SHM semaphore");

printf( szBuff, "CS%s%d", objName, transId);
if ( NULL == (*phSemCS = CreateSemaphore( NULL, 0, NR_MAX_SLOTURI, szBuff))
    Panic ("CS->Error : Cannot create CS semaphore");

printf( szBuff, "SLOT%s%d", objName, transId);
if ( NULL == (*phSemSlot = CreateSemaphore( NULL, NR_MAX_SLOTURI, NR_MAX_SLOTURI, szBuff) )
    Panic ("CS->Error : Cannot create SLOT semaphore");
}
```

(4) Rezervarea și dealocarea unui pachet din zona de memorie partajată a fost implementată astfel:

```
void FreeSlot( int slot, HANDLE hSemSHM, HANDLE hSemSlot, struct SHM *pShm)
    WaitForSingleObject( hSemSHM, INFINITE );
    freeSlotList[freeSlotListIdx] = slot;
    ReleaseSemaphore( hSemSHM, 1, NULL);
    ReleaseSemaphore( hSemSlot, 1, NULL);
}

int GetSlot( HANDLE hSemSHM, HANDLE hSemSlot, struct SHM *pShm, DWORD pid){
    int slot;
    WaitForSingleObject( hSemSlot, INFINITE);
    WaitForSingleObject( hSemSHM, INFINITE);
    slot = freeSlotList[freeSlotListIdx++];
    ReleaseSemaphore( hSemSHM, 1, NULL);
    return slot;
}
```

(5) Comunicația între două procese prin zona de memorie partajată, se realizează în varianta Windows folosind următoarea funcție SendPacket():

```
void SendPachet(PACHET*pPachet,HANDLE hSemSHM,HANDLE hSem,HANDLE hSemSlot,structSHM*pShm)
{
    int slotNo;
    slotNo = GetSlot( hSemSHM, hSemSlot, pShm, pid);
    WaitForSingleObject( hSemSHM, INFINITE);
    objQueueList[0][objQueueListIdx[0]] = nrSlot;
    objQueueListIdx[0] = objQueueListIdx[0](++indexOb[0]) % MAX_SLOTS;
    ReleaseSemaphore( hSemSHM, 1, NULL );
    memcpy( &slotsTab[ slotNo], pPachet, sizeof( slotsTab[slotNo]));
    ReleaseSemaphore( hSem, 1, NULL);
}
}
```

(6) Serverul de comunicații principal și serverele de comunicație care gestionează câte o tranzație sunt implementate în *fire* separate. Fiecare fir este creat și lansat în execuție de procesul programului principal la primirea în soclul INET a unui mesaj *RECORD*, *DELETE*, *END_RECORD*, *ABORT*, *COMMIT*, *NETWORK_BEGIN_TRANSACTION*, când se stabilește o nouă conexiune. Firul care se creează trebuie să-și copieze cererea primită și soclul noii conexiuni, după care procesul programului principal poate trata altă cerere de acceptare a unei conexiuni. Pentru sincronizarea proces principal-fir pe parcursul copierii acestor informații se folosește un semafor special.

(7) Detectarea pierderii legăturii cu nodul coordonator după stabilirea unei conexiuni se realizează comparând valoarea returnată de apelurile *recv()* și *send()* folosite pentru comunicație prin socluri cu valoarea *SOCKET_ERROR*.

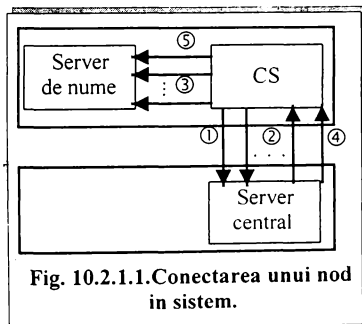
10. 2. Exemple de desfășurare a comunicațiilor în sistem

Se prezintă în continuare câteva exemple de desfășurare a comunicațiilor în sistem. Sunt luate în considerare toate cele șase situații distincte care pot să apară în funcționarea sistemului DOSTP: (1) conectarea unui nod la sistemul DOSTP; (2) integrarea sau eliminarea unui obiect din sistem; (3) inițierea unei noi tranzații; (4) implicarea unui obiect local și a unui obiect de la distanță într-o tranzație; (5) realizarea protocolului *two_phase commit*; (6) recuperare după o defecțiune a unui nod și a reintrării sale în sistem.

Pentru a se deosebi mesajele transmise în canalele create static și asociate managerilor de tranzații, serverelor de comunicații, de cele create dinamic, în toate figurile din această secțiune (10.2.) transmiterea mesajelor din prima categorie va fi reprezentată cu linii mai îngroșate.

10.2.1. Conectarea nodurilor și a managerilor de obiecte la sistemul tranzacțional

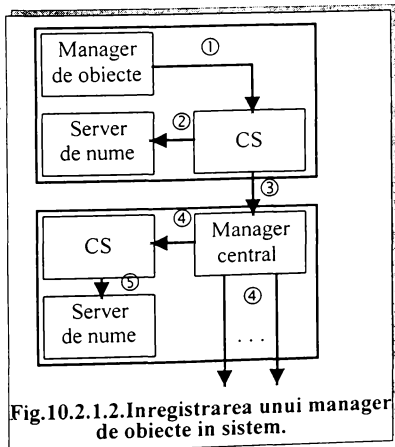
Conectarea unui nod la sistemul DOSTP este inițiată de serverul de comunicații. Într-un nod care se conectează la sistemul DOSTP trebuie să se lanseze în execuție (în această ordine) serverul de comunicații, serverul de nume, managerul de tranzații și managerii de obiecte. Înaintea, însă, a conectării oricărui nod trebuie să se lanseze managerul central pe nodul fixat.



Circulația mesajelor în sistem în cazul conectării unui nod este următoarea (fig. 10.2.1.1):

- (1) Serverul de comunicații trimite mesajul *CONNECT* managerului central, pe canalul fixat static al acestuia.
- (2) Managerul central recepționează mesajul *CONNECT* și transmite serverului de comunicații (apelant) prin mesaje de tipul *RECORD*, lista cu identificatorii și localizarea obiectelor conectate la sistem.
- (3) Serverul de comunicații transmite mai departe mesajele *RECORD* la serverul de nume local.
- (4) După transmisia întregii liste a obiectelor, managerul central transmite mesajul *END_RECORD*.
- (5) După recepția mesajului *END_RECORD*, serverul de comunicații transmite mesajul mai departe la serverul de nume.

În cazul introducerii sau eliminării unui manager de obiecte la sistem, circulația mesajelor este următoarea: (fig. 10.2.1.2.)



- (1) Managerul de obiecte scrie în canalul (fixat static) al serverului de comunicație local (soclu TCP) mesajul pentru conectarea sau eliminarea sa din sistem: *RECORD*, respectiv *DELETE*.
- (2) Serverul de comunicații trimite mesajul primit (*RECORD* sau *DELETE*) la serverul de nume local.
- (3) Serverul de comunicații trimite mesajul la managerul central.
- (4) Managerul central transmite la serverele de comunicație de pe toate nodurile conectate la sistem mesajul de introducere sau eliminare a unui obiect.
- (5) Serverele de comunicații (cu excepția celui de pe nodul în care se afla managerul de obiect) transmit mesajul primit la serverele lor de nume locale.

10.2.2. Inițierea unei tranzacții

Tranzacțiile se inițiază de către aplicații prin apelul funcției *BeginTransaction()*, care trimite apoi mesajul *BEGIN_TRANSACTION* managerului de tranzacții.

Ațiunile întreprinse de sistemul DOSTP la primirea unui mesaj *BEGIN_TRANSACTION* sunt următoarele (fig. 10.2.2.):

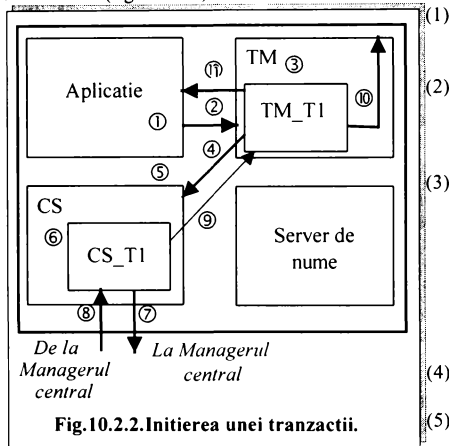


Fig.10.2.2. Inițierea unei tranzacții.

- (1) Aplicația trimite mesajul *BEGIN_TRANSACTION* în canalul fixat static *chanf_MT* al managerului de tranzacții.
- (2) Managerul de tranzacții, care așteaptă în buclă continuă primirea de mesaje în *chanf_TM* recepționează mesajul *APL_BEGIN_TRANSACTION*.
- (3) Managerul de tranzacții (TM) creează un proces fiu, *TM_T1* care va deservi toate cererile trimise de tranzația T1 în memoria partajată; TM transmite mesajul *APL_BEGIN_TRANSACTION* fiului nou creat; *TM_T1* va crea și se va atașa la o zonă de memorie partajată în care se vor transmite, local, toate mesajele pentru acea tranzație.
- (4) *TM_T1* transmite serverului de comunicații mesajul *TM_BEGIN_TRANSACTION*.
- (5) Serverul de comunicații (CS) recepționează mesajul *TM_BEGIN_TRANSACTION*.
- (6) Serverul de comunicații, CS, creează un fiu *CS_T1* care va gestiona exclusiv comunicațiile pentru tranzația respectivă T1.
- (7) *CS_T1* transmite un mesaj *GET_TRANS_ID* la managerul central, pentru a obține un identificator și o prioritate pentru tranzație.
- (8) Managerul central recepționează mesajul *GET_TRANS_ID*, alocă un identificator și o prioritate pentru tranzație și răspunde apelantului transmițându-i identificatorul și prioritatea tranzației.
- (9) *CS_T1* recepționează mesajul transmis de managerul central și transmite managerului *TM_T1* mesajul *ACK_TM_BEGIN_TRANSACTION*.
- (10) *TM_T1* transmite managerului de tranzacții TM mesajul *TRANSACTION_ACCEPTED*.
- (11) *TM_T1* transmite aplicației *ACK_APL_BEGIN_TRANSACTION*.

10.2.3. Exemple de implicare a obiectelor în tranzacții

Există două situații distincte de implicare a obiectelor în tranzacții: (1) obiectul pentru care s-a cerut operația (prima operație pentru acel obiect) este local; (2) obiectul este situat la distanță.

Ațiunile întreprinse de sistem pentru implicarea unui obiect local într-o tranzație T1 sunt următoarele (fig. 10.2.3.1):

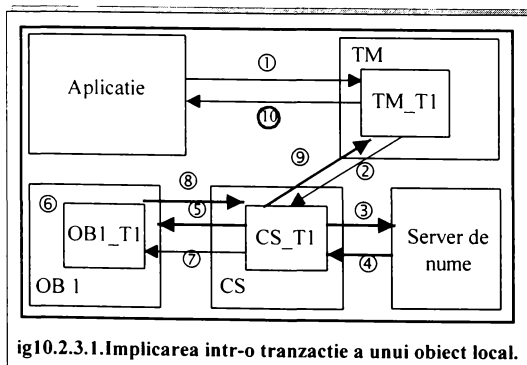


Fig.10.2.3.1. Implicarea într-o tranzație a unui obiect local.

- (1) aplicația trimite mesajul *REQUEST* (în zona de memorie partajată), mesaj care cuprinde operația cerută.

- (2) managerul de tranzacții TM_T1 recepționează cererea și, în urma analizei tabelelor sale, constatînd că obiectul nu a mai fost apelat în cadrul tranzacției T1, modifică tipul mesajului primit la *FIRST_REQUEST* și îl trimite mai departe la CS_T1.
- (3) CS_T1 recepționează mesajul pentru a localiza obiectul, contactează serverul de nume trimișîndu-i un mesaj de tipul *GET_OBJECT*.
- (4) serverul de nume recepționează mesajul și trimite răspunsul la CS_T1.
- (5) CS_T1 trimite mesajul *BEGIN_TRANSACTION* în canalul fixat static al managerului obiectului.
- (6) managerul obiectului crează un fiu care se va ocupa de tranzacția respectivă și care se atașează la zona de memorie partajată (OB_T1) și așteaptă primirea de mesaje în zona de memorie partajată.
- (7) CS_T1 trimite în zona de memorie partajată pentru OB_T1 cererea primită în mesajul *FIRST_REQUEST* de la TM_T1.
- (8) OB_T1 primește mesajul și trimite către CS_T1 răspunsul la cerere într-un mesaj de tipul *ACK_FIRST_REQUEST*.
- (9) CS_T1 primește mesajul și transmite răspunsul la TM_T1.
- (10) TM_T1 transmite răspunsul cererii la aplicație.

În scopul implicării unui obiect de la distanță, acțiunile întreprinse sunt următoarele:

- (1) aplicația trimite mesajul *REQUEST* (în zona de memorie partajată), mesaj care cuprinde o operație asupra obiectului OB1.
- (2) managerul de tranzacții TM_T1 recepționează cererea și, în urma analizei tabelelor sale, constatînd că obiectul nu a mai fost apelat în cadrul tranzacției T1, modifică tipul mesajului primit la *FIRST_REQUEST* și îl trimite mai departe la CS_T1.
- (3) CS_T1 transmite mesajul *GET_OBJECT* la serverul de nume, pentru a localiza obiectul.
- (4) serverul de nume trimite răspunsul la CS_T1.
- (5) CS_T1 începe implicarea obiectului de la distanță în tranzacție prin transmisia unui mesaj *NETWORK_BEGIN_TRANSACTION*.
- (6) CS_T1 crează un proces fiu CS_T1_OB1 pentru a recepționa mesaje de la nodul pe care se află obiectul.
- (7) serverul de comunicații de pe nodul pe care se află obiectul, în urma recepționării mesajului de la (5) crează un fiu CS_OB1_T1 care va gestiona comunicația cu obiectul OB1 în cadrul tranzacției de la distanță T1 și crează o zona de memorie partajată.
- (8) CS_OB1_T1 (din nodul obiectului) transmite un mesaj *BEGIN_TRANSACTION* în canalul static al managerului obiectului cu informații despre zona de memorie partajată.
- (9) managerul obiectului crează un fiu OB1_T1 care se va ocupa de tranzacție; acesta se atașează la memoria partajată și transmite managerului obiectului OB1 un mesaj *TRANSACTION_ACCEPTED*.
- (10) CS_T1 transmite la CS_OB1_T1 cererea, într-un mesaj *FIRST_REQUEST*.
- (11) CS_OB1_T1 transmite cererea la OB1_T1.
- (12) OB1_T1 execută operația și transmite mesajul *ACK_FIRST_REQUEST* către CS_OB1_T1.
- (13) CS_OB1_T1 transmite răspunsul la CS_T1_OB1.
- (14) CS_T1_OB1 transmite la TM_T1 răspunsul la cerere.
- (15) TM_T1 transmite răspunsul la aplicație.

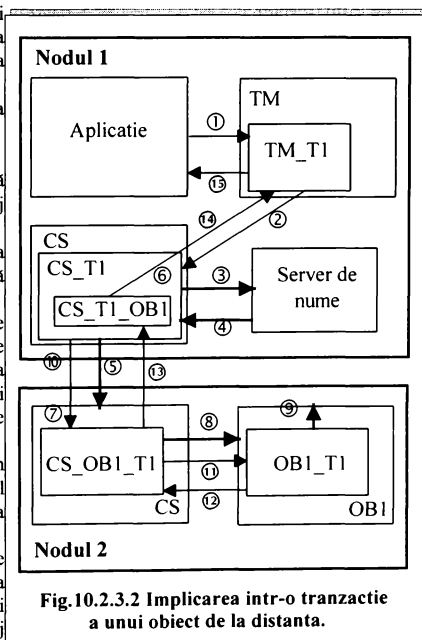


Fig.10.2.3.2 Implicarea într-o tranzacție a unui obiect de la distanță.

10.2.4. Exemple de desfășurare a protocolului two-phase commit

În cazul protocolului *two-phase commit* există trei aspecte care trebuie luate în considerare:

- (1) acțiunile întreprinse de sistem în faza *prepare_to_commit*.
- (2) acțiunile întreprinse de sistem în faza *commit*.
- (3) acțiunile întreprinse de sistem în faza *abort*.

Pentru primul caz, acțiunile sunt reprezentate în fig. 10.2.4.1., în situația unei tranzacții care folosește două obiecte: OB1 local și OB2 la distanță.

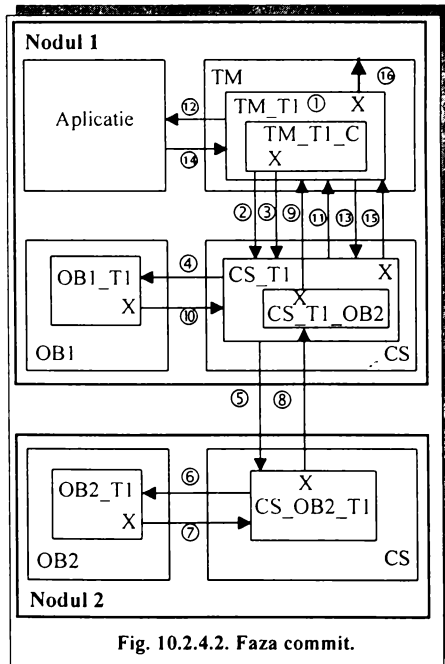
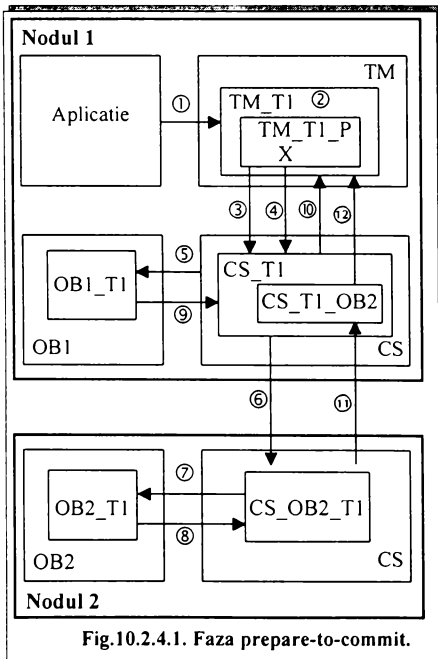
Acțiunile sunt următoarele:

- (1) aplicația transmite, în urma execuției funcției *EndTransmission()*, mesajul *END_TRANSACTION*.
- (2) procesul *TM_T1*, care controlează tranzacția și recepționează mesajul *END_TRANSACTION*, crează un fiu *TM_T1_P*; acesta identifică obiectele participante la tranzacție: *OB1* și *OB2*.
- (3) *TM_T1_P* trimite către *CS_T1* mesajul *PREPARE_TO_COMMIT* (pentru *OB1*).
- (4) *TM_T1_P* trimite către *CS_T1* mesajul *PREPARE_TO_COMMIT* (pentru *OB2*) și *TM_T1_P* se termină (*exit()*).
- (5) *CS_T1* transmite mesajul *PREPARE_TO_COMMIT* către *OB1_T1*.
- (6) *CS_T1* transmite mesajul *PREPARE_TO_COMMIT* către *CS_OB2_T1*.
- (7) *CS_OB2_T1* transmite la *OB2_T1* mesajul primit, *PREPARE_TO_COMMIT*.
- (8) *OB2_T1* transmite *ACK_PREPARE_TO_COMMIT* la *CS_OB2_T1*.
- (9) *OB1_T1* transmite la *CS_T1* mesajul *ACK_PREPARE_TO_COMMIT* (*OB1*).
- (10) *CS_T1* transmite la *TM_T1* mesajul primit de la *OB1_T1*.
- (11) *CS_OB2_T1* transmite la *CS_T1* *OB2* mesajul primit de la *OB2_T1*.
- (12) *CS_T1* transmite la *TM_T1* mesajul *ACK_PREPARE_TO_COMMIT* (*OB2*).

În momentul recepționării confirmărilor la mesajele *PREPARE_TO_COMMIT* de la toate obiectele implicate în tranzacție, managerul tranzacțiilor decide dacă tranzacția se poate termina prin *commit* sau trebuie terminată cu *abort*.

În cazul când tranzacția se poate termina cu *commit*, acțiunile și fluxul de mesaje din sistem sunt următoarele (fig. 10.2.4.2):

- (1) *TM_T1* crează un proces fiu *TM_T1_C*.
- (2) *TM_T1_C* transmite mesajul *COMMIT* la *CS_T1* (pentru obiectul *OB1*).
- (3) *TM_T1_C* transmite mesajul *COMMIT* la *CS_T1* (pentru obiectul *OB2*) și apoi se termină.
- (4) *CS_T1* transmite mesajul *COMMIT* la *OB1_T1*.
- (5) *CS_T1* transmite mesajul *COMMIT* la *CS_OB2_T1*.



- (6) CS_OB2_T1 transmite mesajul *COMMIT* la OB2_T1.
- (7) OB2_T1 tratează cererea *COMMIT*, transmite mesajul *ACK_COMMIT* la CS_OB2_T1 și apoi se termină.
- (8) CS_OB2_T1 transmite mai departe mesajul *ACK_COMMIT* la CS_T1_OB2 și se termină.
- (9) CS_T1_OB2 recepționează mesajul *ACK_COMMIT*, îl transmite mai departe la TM_T1 și se termină.
- (10) OB1_T1 tratează cererea *COMMIT*, transmite mesajul *ACK_COMMIT* către CS_T1 și apoi se termină.
- (11) CS_T1 transmite către TM_T1 mesajul recepționat, *ACK_COMMIT*.
- (12) TM_T1 transmite la aplicație codul de eroare obținut în urma recepționării tuturor confirmărilor la *ACK_COMMIT*.
- (13) TM_T1 transmite apoi la CS_T1 mesajul *TM_END*.
- (14) după recepționarea răspunsului la *COMMIT*, aplicația trimite la TM_T1 mesajul *APL_END*.
- (15) CS_T1, după recepționarea mesajului *TM_END* transmite la TM_T1 mesajul *CS_END*.
- (16) TM_T1 transmite mesajul *TRANSACTION_ENDED* către TM, și apoi se termină.

Terminarea unei tranzacții prin abort generează în sistem o circulație de mesaje asemănătoare cu cea de mai sus (fig. 10.2.4.3) cu excepția faptului că se transmit mesaje *ABORT* (2, 3, 4, 5, 6, 7) și *ACK_ABORT* (8, 9, 10, 11, 12) în locul mesajelor *COMMIT* și *ACK_COMMIT* și a faptului că răspunsul către aplicație se transmite imediat.

10.2.5. Exemple de tratare a situațiilor de recuperare

Există două situații de pierdere a contactului cu un nod implicat într-o tranzacție, care trebuie tratate în mod diferit: (1) când nodul cu care s-a pierdut contactul conține un obiect la distanță implicat într-o tranzacție; (2) când nodul cu care s-a pierdut contactul conține managerul de tranzacții coordonator al unei tranzacții. Există însă și o a treia situație în care trebuie raportată o eroare de tip "manager de obiecte inexistent", situație în care un nod care deține unul sau mai mulți manageri de obiecte nu mai răspunde la mesaje, dar care nu a fost implicat într-o tranzacție. În momentul în care un obiect de pe acest nod este implicat într-o tranzacție T1 în urma transmiterii mesajului *GET_OBJECT*, serverul de comunicații va primi ca răspuns un cod de eroare "manager de obiect inexistent". În acest caz se poziționează în *TM_T1* un indicator *transErr* astfel încât să poată fi refuzate toate cererile de operații următoare din cadrul acelei tranzacții.

A acțiunile și fluxul de mesaje desfășurate de sistemul DOSTP într-o situație de tipul (2), când nodul pe care se afla coordonatorul unei tranzacții T1 nu mai răspunde la mesaje, și apoi este repus în funcțiune sunt următoarele (fig. 10.2.5.1):

(1) când nodul pe care se afla coordonatorul tranzacției nu mai este integrat în sistem, nodurile pe care se află obiectele implicate în tranzacții sesizează acest lucru prin faptul că, fie se primește un semnal de time-out în urma unei operații de recepție mesaj, fie se primește semnalul *SIGPIPE* (Unix) fie se returnează eroare la recepția sau emisia unui mesaj.

- (2) CS_OB2_T1 transmite la OB2_T1 mesajul *INTERRUPT*.
- (3) OB2_T1, după recepția mesajului *INTERRUPT*, transmite mesajul *ACK_INTERRUPT* și se termină.
- (4) CS_OB2_T1, după recepția mesajului *ACK_INTERRUPT* distruge zona de memorie partajată și se termină.
- (5) Se repune în funcțiune nodul coordonator

⇒

(A). Managerul de tranzacții T1, în urma analizei fișierului *TMRECI*, gasește tranzacția T1 neterminată și creează un proces fiu, *TM_T1* care, în funcție de starea tranzacției, realizează următoarele:

- (1) Dacă starea tranzacției se găsește a fi *ABORT* sau *WORKING*, atunci *TM_T1* va forța execuția unei faze *ABORT* pentru toate obiectele implicate în tranzacție. Va crea un proces fiu care va trimite mesaje *ABORT* la toate obiectele și va aștepta răspunsuri de la acestea. Deci, mai întâi trimite mesajul *ABORT* (pentru obiectul 1) la CS.
- (2) *TM_T1* trimite mesajul *ABORT* (pentru obiectul OB2) la CS.
- (3) CS, după recepția mesajului *ABORT* pentru obiectul OB1, creează un proces fiu CS_OB1.
- (4) CS, după recepția mesajului *ABORT* pentru obiectul OB2, creează un proces fiu CS_OB2.
- (5) CS_OB1 trimite un mesaj *GET_OBJECT* la serverul de nume, pentru a localiza obiectul.
- (6) CS_OB2 trimite un mesaj *GET_OBJECT* la serverul de nume, pentru a localiza obiectul.
- (7) Serverul de nume răspunde transmițând localizarea OB1.
- (8) Serverul de nume răspunde transmițând localizarea OB2.
- (9) CS_OB1 transmite mesajul *ABORT* la OB1.
- (10) OB1 transmite la CS_OB1 confirmarea *ACK_ABORT*.
- (11) CS_OB1 transmite la *TM_T1* răspunsul *ACK_ABORT* și se termină.
- (12) CS_OB2 transmite mesajul *ABORT* la serverul de comunicații CS de pe nodul obiectului OB2.
- (13) CS primește mesajul *ABORT* și creează un proces fiu CS_OB2_R.
- (14) CS_OB2_R transmite mesajul la OB2 și așteaptă confirmarea.

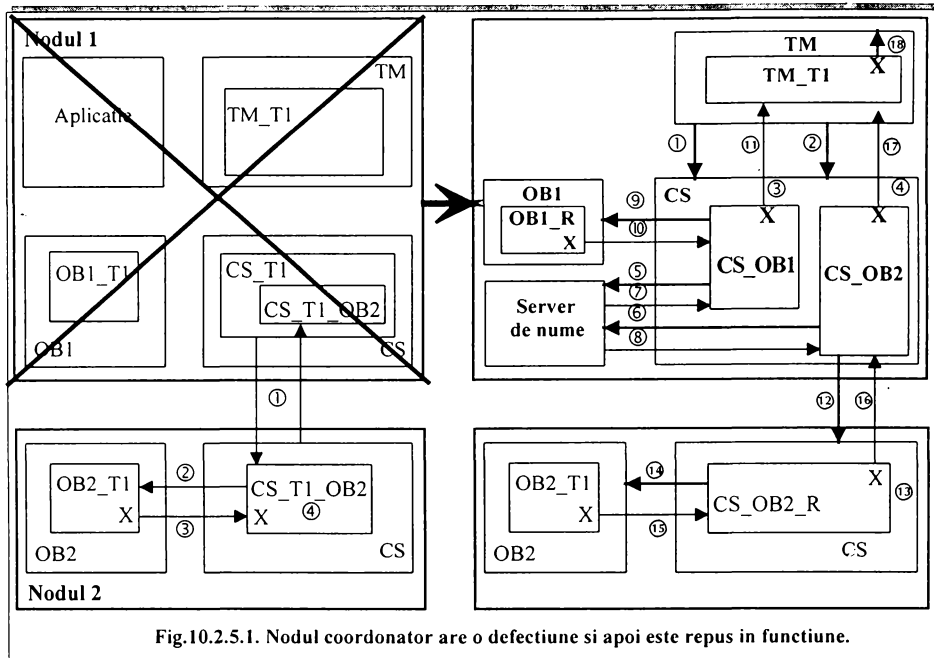


Fig.10.2.5.1. Nodul coordonator are o defectiune si apoi este repus in functiune.

- (15) OB2 transmite ACK_ABORT la CS_OB2_R.
- (16) CS_OB2_R transmite confirmarea la CS_OB2 și se termină.
- (17) CS_OB2 transmite confirmarea la TM_T1 și se termină.
- (18) TM_T1 primește și confirmarea de la obiectul 2 și se termină.

B. Dacă starea tranzacției este *PREPARE_TO_COMMIT*, atunci, dacă în urma analizei fișierului TMREC2 se găsește starea tuturor obiectelor OK, se forțează faza *commit* iar circulația de mesaje este asemănătoare cu cea de mai sus 1-18, dar se lucrează cu mesaje *COMMIT* și confirmări pentru *COMMIT*. Dacă cel puțin un obiect este găsit în TMREC2 într-o stare diferită de OK atunci se forțează faza *abort* iar circulația mesajelor este aceeași ca mai sus 1-18.

C. Dacă starea tranzacției este *COMMIT* atunci se forțează faza *commit*.

Acțiunile și circulația mesajelor în DOSTP în cazul în care nodul pe care se află un obiect nu mai poate fi contactat sunt (fig. 10.2.5.2.):

- (1) Când nodul pe care se află un obiect nu mai este conectat în sistem (fie nodul are o defectiune, fie linia de comunicație cu acel nod este defectă) nodul coordonator sesizează acest lucru fie obținând erori la apelul unei operații de recepție sau emisie, fie primind un semnal SIGPIPE (Unix) fie un semnal de time-out.
- (2) CS_T1_OB2 transmite un mesaj *INTERRUPT* la TM_T1 (în zona de memorie partajată) și se termină.
- (3) TM_T1 analizează fișierul TMREC1 și în funcție de starea tranzacției decide execuția unei faze *commit* sau *abort* și crează un proces fiu TM_T1_R.
- (4) TM_T1_R trimite un mesaj *INTERRUPT_ABORT* (sau *INTERRUPT_COMMIT*) la CS_T1 pentru obiectul OB1.
- (5) TM_T1_R trimite un mesaj *INTERRUPT_ABORT* (*INTERRUPT_COMMIT*) la CS_T1 pentru obiectul OB2 și se termină.
- (6) CS_T1 trimite un mesaj *ABORT* (sau *COMMIT*) la OB_T1.
- (7) OB_T1 confirmă mesajul prin *ACK_ABORT* la CS_T1.
- (8) CS_T1 transmite *ACK_ABORT* la TM_T1.
- (9) TM_T1 distruge procesul care s-ar mai putea afla în execuție pentru a recepționa mesajele de la obiectul OB2 (CS_OB2_T1) și închide vechiul canal și crează un nou proces fiu CS_T1_R.
- (10) CS_T1_R transmite un mesaj *GET_OBJECT* la serverul de nume pentru a determina localizarea obiectului OB2.
- (11) Serverul de nume răspunde transmițând localizarea lui OB2.
- (12) CS_T1_R încearcă într-o buclă continuă să realizeze o conexiune cu serverul de comunicație.
- (13) Se repune în funcțiune nodul care conține obiectul OB2.

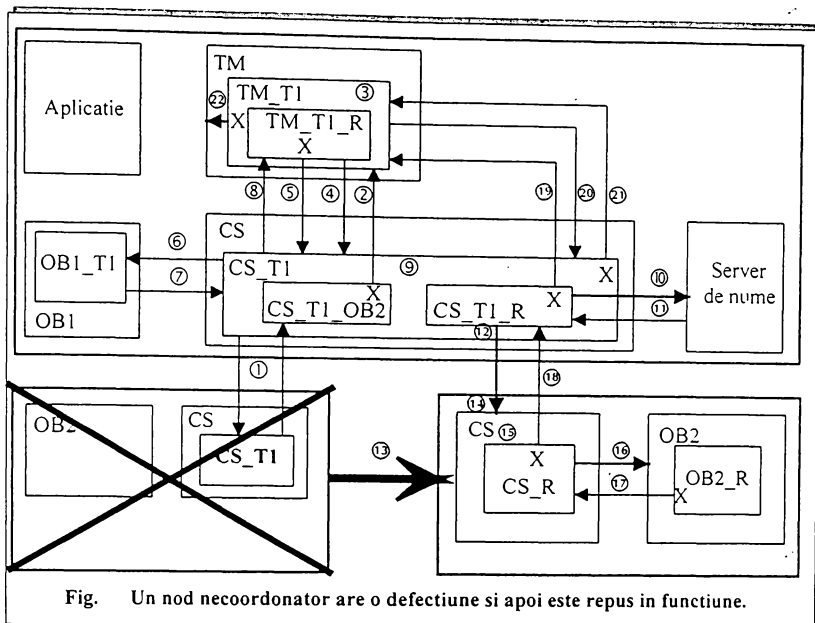


Fig. Un nod necoordinator are o defectiune si apoi este repus in functiune.

⇒

- (14) Serverul de comunicații al nodului repus în funcțiune, CS, primește mesajul *ABORT* (sau *COMMIT*).
- (15) CS crează un proces fiu CS_R care execută recuperarea.
- (16) CS_R transmite obiectului OB2 mesajul *ABORT* (sau *COMMIT*).
- (17) OB2 confirmă mesajul primit de la CS_R cu *ACK_ABORT* (sau *ACK_COMMIT*).
- (18) CS_R transmite la TM_T1_R confirmarea primită și se termină.
- (19) CS_T1_R transmite confirmarea la TM_T1 și se termină.
- (20) TM_T1 transmite la CS_T1 mesajul *TM_END*.
- (21) CS_T1 transmite la TM_T1 mesajul *CS_END* și se termină.
- (22) TM_T1 transmite managerului de tranzacții *TRANSACTION_ENDED* și se termină.

10.3. Exemplu de implementare al unui manager de obiecte

Subcapitolul următor prezintă modul de implementare, folosind programarea orientată pe obiecte, al unui manager de obiecte.

10.3.1. Descrierea managerului de obiecte

Scopul principal în realizarea managerilor de obiecte este ascunderea detaliilor de conectare cu sistemul tranzacțional. Managerii de obiecte sunt realizați astfel încât ei comunică numai cu serverul de comunicație local, fără a avea cunoștință dacă serviciile le-au fost solicitate local sau de la distanță.

Orice manager de obiecte care urmează să se integreze în sistemul DOSTP trebuie să implementeze următoarele cinci seturi de funcții:

- (1) *funcția de conectare la sistem* care constă în transmiterea unui mesaj *RECORD* către managerul central, prin intermediul serverului de comunicații local
- (2) *funcția de eliminare din sistem* care constă în transmiterea unui mesaj *DELETE* către managerul central, prin intermediul serverului de comunicații local, în momentul când utilizatorul dorește eliminarea obiectului din sistem
- (3) *funcția de lansare* a unei noi tranzacții, la solicitarea unui client; această funcție trebuie apelată la recepționarea mesajului *BEGIN_TRANSACTION* în canalul static (soclu) al managerului de obiecte. Rolul său este de a genera un proces fiu (în Unix) sau un fir (în Windows) care să se atașeze la zona de memorie partajată, să returneze confirmarea *ACK_FIRST_REQUEST* și să deservască apoi cererile de operații asupra obiectului respectiv, recepționate sub formă de pachete în zona de memorie partajată.
- (4) *funcția de recepționare/transmisie a pachetelor* din zona din/în zona de memorie partajată.
- (5) *funcțiile necesare modului de lucru tranzacțional*. Acestea cuprind funcțiile care tratează mesajele *REQUEST*, *PREPARE_TO_COMMIT*, *COMMIT* și *ABORT* primite ca pachete în zona de memorie

partajată în timpul desfășurării normale a unei tranzacții și mesajul *INTERRUPT* primit ca urmare a pierderii legăturii cu nodul coordonator precum și funcțiile care tratează mesajele *COMMIT* și *ABORT* primite de la nodul coordonator în canalul fixat al managerului de obiecte în urma unei defecțiuni și apoi a repunerii în funcțiune a nodului (coordonator sau necoordonator).

10.3.1.1. Organizarea managerului de obiecte

Deoarece funcțiile din seturile 1-4 trebuie să le asigure orice manager de obiect, o primă consecință este că acestea, o dată scrise, realizarea unui nou manager de obiecte poate fi făcută fără a se cunoaște detaliile despre implementarea sistemului tranzacțional; trebuie implementate numai funcțiile din setul 5. A doua consecință este că un nou manager de obiecte poate fi testat separat, înainte de integrarea sa în sistemul tranzacțional. Soluția implementării obiectuale este utilă din ambele puncte de vedere.

Managerul de obiecte care s-a realizat până în momentul de față în sisteme Unix și Windows gestionează vectori cu elemente întregi.

După analiza seturilor de funcții de la 10.3.1. s-a optat pentru implementarea managerului de obiecte a următoarelor clase:

- (1) o clasă de bază *CObject*, clasă abstractă, de la care, prin derivare să poată fi implementat orice tip de obiect. Această clasă trebuie să cuprindă funcțiile (virtuale) pe care un manager trebuie întotdeauna să le implementeze. Structura ei este următoarea:

```
class CObject :public CDebug, public CError {
public:
    CObjectManager(int debug=1, FILE *f =stdout); #constructor
    virtual int TreatRequest(PACKET *)=0;
    virtual int PrepareToCommit(void)=0;
    virtual int Commit(void)=0;
    virtual int Abort(void)=0;
    virtual int Interrupt(void)=0;
    virtual int Continue(int, int)=0;
    virtual int ObjectInit(void)=0;
    virtual int Finish(void)=0;
    virtual void SetTransNo(int, int)=0;
    virtual GetInitData(char *)=0;
    virtual void Startup(void)=0;
    virtual void Cleanup(void)=0;
};
class CError {
public:
    int m_errCode=0;
    virtual void Error(int)=0;
};
```

Pentru a ușura depanarea programelor, clasa *CObject* mai moștenește încă două clase (vezi anexa D): clasa *CDebug*, al cărei rol este de a fixa fișierul în care se scriu mesajele pentru depanare (implicit stdout), de a determina identificatorul procesului curent (Unix) sau firului curent (Windows) care prefixează mesajele de depanare și de a fixa modul de lucru (cu sau fără depanare), și clasa *CError* care menține codul de eroare al ultimei operații și o metodă virtuală *Error()* care se particularizează la derivare.

Din definiția clasei *CObject* rezultă că orice manager de obiecte trebuie să implementeze funcțiile modului de lucru tranzacțional (*TreatRequest*, *PrepareToCommit*, *Commit*, *Abort*) și funcțiile care asigură recuperarea (*Interrupt* și *Continue*). În plus, au mai fost introduse funcții pentru: execuția unor operații de inițializare la crearea obiectului (*ObjectInit*), execuția unor operații de dealocare resurse la ieșirea din sistem (*Finish*), execuția unor operații la lansarea/ terminarea unei tranzacții (*Cleanup*, *Startup*), memorarea identificatorului și priorității unor tranzacții (*SetTransNo*) și preluarea unor date inițiale (*GetInitData*).

- (2) o clasă *CObjectManager* care să asigure integrarea unui obiect în sistemul tranzacțional, implementând funcțiile din seturile 1, 2, 3 de la 10.3.1. Fiecarei instanțe de tip *CObject* i se asociază un obiect de tip *CObjectManager* care realizează legătura obiect-sistem tranzacțional.
- (3) o clasă *CDispatcher* care asigură implementarea funcțiilor din setul 4 de la 10.3.1.

10.3.1.2. Concurența la nivelul managerului de obiecte

În continuare, mecanismele de concurență și de recuperare, care depind de tipul obiectelor implementate vor fi descrise pentru managerul de vectori implementat în DOSTP. Clasa care implementează comportarea acestui obiect, obținută prin derivarea clasei *CObject* a fost denumită *CVector* și este prezentată în anexa D. Asupra unui obiect de tip vector de elemente întregi pot fi realizate două operații: *Read* și *Write*. Structura unei astfel de cereri și a răspunsului la o cerere este următoarea:

```
typedef struct {
    int type;           /* tipul operației READ sau WRITE */
    int errorCode;     /* codul de eroare returnat în răspuns */
    int position;      /* poziția în vector */
    int value;         /* valoarea returnată la READ sau valoarea de scris la WRITE */
} VEC_REQUEST;
```

Vectorul este un *obiect persistent*: el este reprezentat în memoria externă printr-un fișier, al cărui nume este parametru al constructorului clasei. Pentru implementarea *listelor de intenții* se folosesc două fișiere: un fișier care urmărește execuția fazelor tranzațiilor (*prepare-to-commit, commit sau abort*) și un altul pentru memorarea operațiilor executate de o tranzație și care trebuie permanentizate în faza *commit*. Numele acestor două fișiere se primesc tot ca parametri ai constructorului clasei *CVector*. În plus, pentru a evita încercările de recuperare a unor tranzații deja terminate, se mai lucrează cu un fișier care menține identitatea tranzațiilor încheiate.

La nivelul unui obiect, pot fi în lucru mai multe tranzații, deservite fiecare de cite un proces separat (în Windows un fir separat); deoarece aceste procese trebuie să partajeze obiectul, soluția adoptată este următoarea: programul principal prezintă o interfață grafică (XWindows sau Windows) care furnizează butoane pentru crearea unui obiect (vector), eliminarea unui obiect și afișarea tranzațiilor active pentru un obiect; în momentul creării unui obiect, în funcția de tip *callback*, apelată în programul principal, se crează un obiect al clasei *CVector*. Constructorul acestui obiect mapează în memorie fișierul corespunzător obiectului persistent, deschide în modul *append* fișierele *log* și alocă memorie pentru informațiile de stare. Din punct de vedere al concurenței, ultima operație prezintă importanță. Informațiile de stare ale obiectului vor fi partajate de toate procesele (sau firele) care deservesc tranzații pentru acel obiect. Ca atare, ele vor fi plasate într-o zonă de memorie care trebuie să poată fi partajată de aceste procese (sau fire) iar accesul la zonă să fie realizat în regim de excludere reciprocă, deci controlat prin semafoare.

Procesele care vor deservi tranzații pentru obiectul creat vor fi create de programul principal. În consecință, în Unix ele vor moșteni codul programului inițial, deci vor cunoaște adresa zonei de memorie partajată în care sunt depuse informațiile de stare, la care trebuie însă să se ataseze pentru a avea acces. În Windows, lucrându-se cu fire și mecanismul firelor nepunând la dispoziție conceptul de moștenire al codului ca în Unix, este necesară o copie explicită a obiectului *CVector* din programul principal în firele care deservesc tranzații, pentru a se lucra cu o altă copie a obiectului, dar cu aceeași zonă de memorie pentru informațiile de stare. De aceea, în anexa D, în definiția clasei *CObject*, varianta Windows, apare o metodă *CopyObject* în plus față de varianta Unix. Structura și semnificația informațiilor de stare este următoarea:

```
struct STATE_INF {
    TRANS_INF vTransId[MAX_TRANS_NO]; /* identificatorii si prioritatile tranzactiilor active */
    short vBlockedTransId[MAX_TRANS_NO]; /*contine TRUE pentru tr. abortate prematur de o tr. mai prioritara */
    short vReadPos[MAX_TRANS_NO][MAX_OBJ]; /* pentru fiecare tranzactie se retin pozitiile citite */
    short vWritePos[MAX_TRANS_NO][MAX_OBJ]; /* pentru fiecare tranzactie se retin pozitiile scrise */
    short vBlockedWPos[MAX_OBJ]; /* pozitia lock-urilor pentru scriere (write-lock-uri) */
    short vBlockedRPos[MAX_OBJ]; /* pozitia lock-urilor pentru citire (read-lock-uri) */
    OLD_TRANS vOldTrans[MAX_OLD_TRANS]; /* tr. in starea OK_PC, dar care nu au primit Commit sau Abort */
    int vOldTransIdx; /* indice in vOldTrans */
};

typedef struct TRANS_INF {
    int transId; /* identificatorul tranzacției */
    int priority; /* prioritatea tranzacției */
} TRANS_INF;

typedef struct OLD_TRANS {
    int transNo; /* identificatorul tranzacției */
    int priority; /* prioritatea tranzacției */
    long inPC; /*poziția în fișierul PCLOG la care încep înregistrările acestei tranzacții */
    short state; /*starea tranzacției:UNCHANGED,RECOVERED,RECOVERING,DONE,NOT_FOUND */
} OLD_TRANS; /* starea tranzacției este vazuta in raport de procesul de recuperare */
```

Procesele (firele) care deservește tranzațiile sunt create de obiectul *CObjectManager* atașat unui obiect *CVector*; ele moștenesc (sau copiază) obiectul *CVector*, generând astfel *obiectele de lucru* (care partajează informațiile de stare) și gestionează informații locale unei tranzații precum starea curentă a tranzației relativă la *prepare-to-commit* (*OK_PC*, *FAIL_PC*, *WORKING_PC*), operațiile Read și Write efectuate (*m_vRPos*, *m_vWPos*), poziția în fișierul PCLOG a înregistrărilor corespunzătoare fazei *prepare-to-commit*.

Se va arăta în continuare modul cum au fost implementate funcțiile modului de lucru tranzacțional (declarate virtuale în clasa de bază) în clasa *CVector*, punându-se accentul pe aspectele de concurență de la nivelul obiectului:

(1) Funcția *ReadRequest* se apelează când se primește o cerere *Read* sau *Write* (pachet de tip *REQUEST*). O cerere de *Read* se tratează astfel: (a) se testează dacă cererea a survenit *după* primirea unui mesaj *PREPARE TO COMMIT* returnându-se în acest caz codul de eroare *EWROG_STATUS*; (b) se testează dacă poziția din cerere nu depășește dimensiunea prestabilită pentru vector, returnându-se codul de eroare *EOVERFLOW*; (c) se testează dacă tranzația nu a fost blocată de o altă mai prioritară, returnându-se codul *ELOW_PRIORITY_ABORT* (se testează *m_pTrStatus -> vBlockedTransld[m_transldx]*); (d) se testează dacă pe poziția solicitată nu există un *write-lock* (o altă tranzație a ajuns în faza *prepare_to_commit* și va scrie această poziție) inspectându-se, pentru poziția solicitată, vectorul *m_pTrStatus -> vBlockedWPos* și returnându-se în acest caz *ELOCK_RECORD*; (e) dacă poziția nu a mai fost citită în această tranzație, se memorează faptul că s-a făcut o citire, alocându-se o nouă intrare în vectorul local *m_vRPos[rldx]* și în zona partajată; (f) se completează în pachetul de răspuns valoarea citită, luată fie direct din zona de memorie în care este mapat obiectul (dacă tranzația nu a scris încă această poziție) fie din *m_vWPos[index].val*.

O cerere *Write* se tratează analog cu o cerere *Read* la punctele (a), (b), (c); la punctul (d) apare o diferență, în sensul că dacă poziția solicitată are un *read-lock*, operația de scriere nu eșuează deoarece scrierea efectivă se face la *commit* și dacă până atunci tranzația care l-a citit se termină (șansele fiind mari, deoarece ea este în faza *prepare_to_commit*) totul este permis. Scrierea se face în copia locală a obiectului (*m_vWPos*) și se marchează și în zona partajată actualizarea acestei poziții.

În ambele operații, în caz de eșec, se marchează tranzația în *m_pTrStatus -> vBlockedTransld[m_transldx]* ca abortată și se deblochează articolele scrise/citite de această tranzație; totodată se face observația că accesul la zona de memorie partajată se face în regim de excludere reciprocă controlat de semaforul *m_semObjId* (aceiași pentru toate *obiectele de lucru*).

(2) Funcția *PrepareToCommit()* se apelează când se primește un pachet de tip *PREPARE TO COMMIT*. Ea realizează prima fază a protocolului *two-phase commit* astfel: (a) se testează dacă starea tranzației este *WORKING_PC*; numai pentru această stare se execută *prepare_to_commit*, altfel se returnează *EWROG_STATUS*; (b) dacă tranzația a fost marcată ca abortată, se deblochează articolele scrise/citite de tranzația respectivă și se returnează un cod de eroare *ELOW_PRIORITY_ABORTED*; (c) dacă cel puțin o poziție din cele citite în cadrul acestei tranzații are un *write-lock*, sau invers, cel puțin o poziție scrisă are un *read-lock* se procedează ca la (b); (d) în acest moment, este sigur că tranzația nu este în conflict cu alte tranzații. Se verifică dacă cel puțin o poziție pentru care tranzația trebuie să obțină un *write-lock* nu a fost cumva citită de alta tranzație *mai prioritară* (nu și invers, o poziție pentru care trebuie un *read-lock* dacă nu cumva a fost citită de alta tranzație mai prioritară, deoarece această situație este permisă); în caz afirmativ, starea tranzației devine *FAIL_PC* și se returnează *ETRY_AGAIN* (c) se obțin lock-urile (se incrementează pozițiile respective în *m_ptrStatus -> vBlockedWPos*, *m_pTrStatus -> vBlockedRPos* și se marchează ca blocate (abort) tranzațiile care au citit pozițiile pe care tranzația curentă le scrie; (f) se completează fișierul PCLOG și se completează deplasamentul său în *m_logPos*, se modifică starea tranzației la *OK_PC* și se returnează *SUCCESS*. Întreaga execuție a funcției *PrepareToCommit()* se execută sub controlul semaforului *m_semObjId* astfel încât actualizarea lock-ului și a fișierului se face în regim de excludere reciprocă.

(3) Funcția *Commit()* se apelează când se primește un pachet de tip *COMMIT*. Ea realizează faza *commit* a protocolului *two phase commit*, partea de server (vezi 7.3.1) astfel: (a) se testează dacă starea tranzației este diferită de *OK_PC* returnându-se *EWROG_STATUS* în acest caz; (b) se înscrie în fișierul ACLLOG o înregistrare care să specifice identificatorul tranzației, starea *COMMIT*, prioritatea ei, poziția listei de intenții în fișierul PCLOG; (c) se actualizează fișierul obiectului conform informațiilor obiectului de lucru local; (d) se înscrie în fișierul tranzațiilor terminate identificatorul acestei tranzații; (e) se marchează pentru ștergere înregistrările acestei tranzații din PCLOG și ACLLOG; (f) se eliberează lock-urile și intrarea corespunzătoare din *m_pTrStatus -> vTransld*.

(4) Funcția *Abort()* se apelează când se primește un pachet de tip *ABORT*; ea realizează faza *abort* a protocolului *two phase commit*, partea de server (vezi 7.3.1) astfel: (a) dacă starea tranzației este *OK_PC* se înscrie în fișierul ACLLOG o înregistrare care să specifice identificatorul tranzației, starea *ABORT*, prioritatea ei și poziția listei de intenții în fișierul PCLOG; (b) se eliberează lock-urile (dacă starea tranzației era *PC_OK*), se demarhează articolele citite/scrise și se eliberează intrarea corespunzătoare din *m_pTrStatus -> vTransld*.

10.3.1.2. Recuperarea obiectelor

La nivelul obiectelor, există două aspecte ale recuperării obiectelor (vezi și 10.2.5):

- (1) primul aspect apare când se lansează în execuție un manager al unui obiect, după ce nodul obiectului pierduse contactul cu nodul coordonator al unei tranzații în care era implicat obiectul. În această situație managerul de obiecte analizează fișierele de log și ia fie decizia de *abort*, fie de *commit*, fie de așteptare a deciziei coordonatorului;
- (2) Al doilea aspect apare când managerul de obiect primește un pachet *INTERRUPT* de la serverul de comunicație ca urmare a pierderii legăturii cu coordonatorul unei tranzații.

Primul aspect este tratat de funcția *ObjectInit()* a clasei *CObject* și deci și a managerului *CVector*. Funcția *ObjectInit* se apelează imediat după crearea unui obiect și are trei scopuri:

- (1) termină tranzațiile care au recepționat de la coordonator comanda *Abort* sau *Commit*, dar, datorită unei defecțiuni, nu au putut-o realiza;
- (2) păstrează în fișierele de log numai tranzațiile cu înregistrări valide, ștergând tranzațiile pentru care s-a realizat faza *Commit* sau *Abort*;
- (3) înregistrează în lista *m_pTrStatus* -> *vOldTrans* informații relative la tranzațiile care nu au putut fi terminate (sunt în starea *OK_PC*, dar nu s-a înregistrat în *ACI.LOG* comanda dată de coordonatorul tranzației) și face blocările de obiecte necesare; pentru fiecare articol scris de o tranzație, obține un *write_lock* și pentru fiecare articol citit obține un *read_lock*.

Al doilea aspect este tratat de funcția *Interrupt()* și *Continu()* ale clasei *CObject* (și deci și ale managerului *CVector*). Funcția *Interrupt()* este apelată atunci când serverul de comunicații detectează întreruperea conexiunii cu coordonatorul tranzației (și aplicația se află la distanță). Ea realizează următoarele operații în regim de excludere reciprocă:

- (1) demarchează articolele scrise/citite de această tranzație;
- (2) dacă tranzația este în starea *OK_PC*, atunci se înscrie această tranzație în vectorul tranzațiilor neterminate (*m_pTrStatus* -> *vOldTrans*);
- (3) eliberează intrarea alocată acestei tranzații în *m_pTrStatus* -> *vTransId*.

Funcția *Continu()* se apelează când, în urma repunerii în funcțiune a nodului coordonator se primește comanda *Abort* sau *Commit* în canalul static. În acest caz, acțiunile care trebuie realizate depind de starea tranzației neterminate și sunt următoarele:

- (1) dacă tranzația este găsită în fișierul tranzațiilor încheiate, *m_tf*, atunci ea a fost terminată și se transmite confirmarea fără nici o altă acțiune;
- (2) dacă tranzația este găsită în tabelul *m_pTrStatus* -> *vTransId* atunci tranzația este încă activă; această situație poate apare când procesul care tratează tranzația nu a tratat încă pachetul *INTERRUPT*, dar coordonatorul a fost repus foarte repede în funcțiune. În acest caz, se așteaptă pînă tranzația nu mai apare ca activă, dar este, în schimb, înscrisă în *m_pTrStatus* -> *vOldTrans* și se iau următoarele măsuri: (2.1) dacă starea tranzației memorată în *m_pTrStatus* -> *vOldTrans* este *RECOVERING*, atunci rezultă că mai este în lucru un proces (sau fir) care are în lucru recuperarea tranzației respective și, ca atare, procesul curent se va termina, lasînd ca răspunsul să fie dat de procesul în lucru; (2.2) dacă starea tranzației este *UNCHANGED*, atunci procesul curent recuperează tranzația; (2.3) dacă starea tranzației este *DONE* înseamnă că deja un proces a executat recuperarea și se transmite confirmarea; (2.4) dacă tranzația nu este găsită în nici o stare din cele de mai sus, înseamnă că tranzația a fost abortată și se transmite confirmarea.

Pentru recuperarea propriu-zisă a unei tranzații în starea *OK_PC* a fost construită o clasă specială *CRecovery* care furnizează următoarele funcții:

- (1) *ResumeCommit()* pentru actualizarea versiunii persistente a obiectului conform actualizărilor memorate în fișierul log *PCLLOG* și actualizarea apoi a fișierelor log;
- (2) *ResumeAbort()* pentru actualizarea fișierelor de log;
- (3) *Commit()* pentru execuția tuturor acțiunilor funcției *Commit* din *CVector*, dar, pe baza informațiilor de actualizare găsite în fișierul *PCLLOG*;
- (4) *Abort()* pentru execuția tuturor acțiunilor funcției *Abort* din *CVector*.

Funcțiile (1) și (2) se apelează când trebuie continuată comiterea/abortarea unei tranzații (care nu a putut fi terminată din cauza unei defecțiuni, primul aspect al recuperării), în funcția *ObjectInit()*, după crearea unui obiect *CRecovery* caruia i se transmit ca parametri (constructorului său): fișierele de log (numele și pointerii la structura lor de tipul *FILF*), pointerul la structura de stare a obiectului (*m_pTrStatus*) deplasamentul articolelor corespunzătoare tranzației în fișierul *PCLLOG*, prioritatea și identificatorul tranzației, semaforul pentru controlul accesului exclusiv.

Funcțiile (3) și (4) se apelează când se primește comanda *Abort* sau *Commit* în canalul static al managerului de obiecte (al doilea aspect al recuperării, în funcția *Continu()*).

10.3.2. Conectarea obiectelor în sistemul DOSTP

S-a aratat la 10.3.1.1 ca, pentru conectarea obiectelor în sistem, au fost prevăzute două clase: *CObjectManager* și *CDispatcher*. În acest subcapitol vor fi prezentate aceste două clase și modul de integrare automată a obiectelor în sistemul tranzacțional.

10.3.2.1. Descrierea clasei CObjectManager

Clasa *CObjectManager* (a cărei definiție este prezentată în anexa D) a fost introdusă în vederea următoarelor două scopuri: (1) să realizeze comunicarea prin canalul static dintre un manager de obiecte cu serverul de comunicații (canalul static este implementat sub forma de soclu *Unix* pentru versiunea Unix și soclu *Inet* pentru versiunea Windows); (2) să pună la dispoziție funcții de vizualizare a tranzacțiilor active. În vederea atingerii acestor scopuri au fost introduși membri precum porturile și soclurile pentru comunicare, pointer la instanța *CObject* asociată, o listă de tranzacții care folosesc instanța *CObject* și alți membri necesari și specifici implementării Unix (widget label, quit, box) sau Windows (semafoare folosite la copierea informațiilor transmise firelor create la inițierea unei noi tranzacții).

În ceea ce privește funcțiile puse la dispoziție de această clasă se enumeră: (1) funcția *ConnectObjectToSystem()* care realizează transmiterea mesajului *RECORD* la managerul central (via server de comunicații) și punerea soclului în așteptare de recepție mesaje; (2) *DisconnectObjectFromSystem()* care realizează transmiterea mesajului *DELETE* la managerul central (via serverul de comunicații), închiderea soclurilor de comunicare, apelul *Finish()* al instanței obiectului asociat *CObject* pentru dealocarea resurselor (eliberarea semafoarelor pentru accesul exclusiv la informațiile de stare și pentru accesul exclusiv la vectorul tranzacțiilor neterminate *vOldTrans*); (3) funcția *GetNextRequest()* care realizează: (3.1) acceptarea unei noi conexiuni și preluarea primei cereri din soclu (în mod normal de tip *BEGIN_TRANSACTION* sau de tip *COMMIT/ABORT*); (3.2) dacă cererea este de tip *BEGIN_TRANSACTION*, se crează un proces (sau fir) care declară un obiect de tip *CDispatcher* care va realiza toată conversația cu serverul de comunicații prin zona de memorie partajată; (3.3) dacă cererea este de tip *ABORT/COMMIT* (trimisă de un manager de tranzacții care vrea să recupereze o tranzacție) se crează un proces (sau fir) care declară un obiect de tip *CDispatcher* și se apelează metoda *Continue* a acestuia; aceasta metodă, după câteva operații de rutină, va apela în cele din urmă metoda *Continue* a instanței *CObject* asociate.

10.3.2.2. Descrierea clasei CDispatcher

Clasa *CDispatcher* (a cărei definiție este prezentată în anexa D) a fost introdusă pentru a implementa funcțiile de preluare/transmitere a pachetelor din zona de memorie partajată. Se crează câte o instanță a acestei clase la fiecare inițiere a unei tranzacții (la primirea unui mesaj *BEGIN_TRANSACTION*) sau la fiecare cerere *ABORT/COMMIT* pentru o tranzacție nedeterminată. Instanțele se crează în cadrul unui proces (sau fir) separat care deservește deslășurarea normală a tranzacției respectiv recuperarea tranzacției.

În cadrul acestei clase sunt prevăzute funcții pentru: (1) atașarea și detașarea la zona de memorie partajată prin care are dialogul cu serverul de comunicații; (2) blocare pînă la recepționarea unui pachet din zona de memorie partajată; (3) preluarea unui pachet din zona de memorie partajată, apelul funcției de tratare a cererii din instanța *CObject* atașată și returnarea răspunsului la cerere tot în zona de memorie partajată; (4) recuperarea unei tranzacții folosind funcția *Continue* din instanța *CObject* asociată. Spre deosebire de primele 3 clase de funcții, pentru recuperarea unei tranzacții, dialogul cu serverul de comunicații se face printr-un canal static (soclu *Unix* sau *Inet*) și nu printr-un canal creat în mod dinamic (zona de memorie partajată).

10.3.2.3. Integrarea automată a obiectelor în sistemul tranzacțional

Pentru a ușura scrierea codului pentru integrarea obiectelor în sistemul tranzacțional s-au generat patru macrodefiniții. Acestea sunt definite în anexa D și au următoarele scopuri:

(1) Macrodefiniția *INCLUDE_SYSTEM*, inclusă în secțiunea "public" a definiției unei clase (ea este inclusă de fapt în definiția clasei *CVector* : public *INCLUDE_SYSTEM* (*CVector*)) are efectul de a introduce automat în definiția clasei a încă citorva membri. Pentru varianta Unix acești membri sunt: un pointer la instanța *CObjectManager* asociată obiectului (creată dinamic), o funcție *friend init* care primește ca parametrii contextul aplicației XWindows obținut în programul principal și identificatorul widget-ului în care instanța *CObjectManager* atașată obiectului va face afișările. Pentru varianta Windows acești membri sunt: portul pe care managerul așteaptă cereri, semaforul folosit la crearea firelor în care acționează instanța *CDispatcher*, și două funcții *friend* : *Init()* și *CreateThreads()*.

(2) Macrodefiniția *DEFINE_SYSTEM* asigură definirea funcțiilor prietene de mai sus.

(3) Macrodefiniția *START_SYSTEM* asigură conectarea unui obiect la sistemul tranzacțional. Ea constă numai în apelul funcției *Init()* definită prin *DEFINE_SYSTEM*.

Partea a II a: Criterii de proiectare a sistemelor de programare distribuită bazate pe obiecte.

(4) Macrodefiniția *STOP_SYSTEM* asigură deconectarea unui obiect de la sistemul tranzacțional. Ea constă numai în apelul metodei *DisconnectObjectFromSystem()* a instanței *CObjectManager*. Folosind aceste definiții, un manager de obiecte (programul principal) poate fi scris astfel:

```
void Quit(Widget w,XtPointer client_data,XtPointer call_data);
void Add(Widget w,XtPointer client_data,XtPointer call_data);

CVector v1 ("v1","data1","PCLOG1","ACLOG1","Trans","Status"),*v2: // se creaza un obiect static

XtAppContext app_cont;
Widget shell ,pane,box,quit,add;
FILE *f1,*f2;

void main (int argc,char** argv)
{
    shell = XtVaAppInitialize( &app_cont, "XBox", NULL, 0, &argc,argv, NULL,NULL );
    pane = XtVaCreateManagedWidget("pane",panedWidgetClass, shell, NULL);
    box = XtVaCreateManagedWidget("box", boxWidgetClass,pane,NULL);
    quit = XtVaCreateManagedWidget("quit",commandWidgetClass,box,NULL);
    add = XtVaCreateManagedWidget("adv2",commandWidgetClass,box,NULL);

    f1=fopen("debug1","w+"); // deschiderea unui fisier pentru depanare
    setbuffer(f1,0,0); // mod de lucru fara buferare
    v1.Init(1,f1); // stabilirea lui f1 ca fisier in care se vor inscrie mesajele de depanare
    v1.ObjectInit(); // faza de recuperare initiala pentru tranzactiile care pot fi terminate

    START_SYSTEM (&v1,app_cont,pane); // conectarea la sistem

    XtAddCallback(quit,XtNcallback,Quit,NULL); // functie apelata la stergerea unui obiect
    XtAddCallback(add,XtNcallback,Add,NULL); // functie apelata la adaugarea unui obiect de acelasi tip
    XtRealizeWidget(shell);
    XtAppMainLoop(app_cont);
}

void Add(Widget w,XtPointer client_data,XtPointer call_data)
{
    v2=new CVector("v2","data2","PCLOG2","ACLOG2","Trans","Status");

    f2=fopen("debug2","w+"); // deschiderea unui fisier pentru depanare
    setbuffer(f2,0,0); // mod de lucru fara buferare
    v2->Init(1,f2); // stabilirea lui f1 ca fisier in care se vor inscrie mesajele de depanare
    v2->ObjectInit(); // faza de recuperare initiala pentru tranzactiile care pot fi terminate

    START_SYSTEM(v2,app_cont,pane); // conectarea la sistem

    XtRemoveCallback(add,XtNcallback,Add,NULL);
    XtDestroyWidget(add);
}

void Quit(Widget w,XtPointer client_data,XtPointer call_data)
{
    if (v2) STOP_SYSTEM(v2); // deconectarea de la sistem a obiectului v1
    v2=&v1;

    STOP_SYSTEM(v2); // deconectarea de la sistem a obiectului v2
    exit(0);
}
```

10.3.3. Particularități ale implementării în sistemele Unix și Windows

Diferențele între implementarea unui manager de obiecte în sistemul Unix și în sistemul Windows provin în principal din faptul că în sistemul Unix se lucrează cu procese, care moștesc codul procesului părinte, și care, după creare, nu mai au acces la acesta, iar în Windows se lucrează cu fire. Cele trei clase de bază prezentate la 10.3.2 sunt aproape identice în ambele implementări. Diferă însă lucrul cu semafoare și cu zonele de memorie partajată. Se enumeră în continuare câteva diferențe între cele două implementări:

- (1) mesajele *BEGIN_TRANSACTION* conțin numai pentru versiunea Unix identificatori pentru semafoare și zona de memorie partajată la care instanța *CDispatcher* trebuie să se atașeze.
- (2) comunicația între serverul de comunicații și managerul de obiecte, la nivelul mesajelor *BEGIN_TRANSACTION*, *ABORT*, *COMMIT* se face printr-un soclu *Inet* în versiunea Windows și printr-un soclu *Unix* în versiunea Unix.
- (3) instanțierile clasei *CDispatcher* se fac în procese separate în versiunea Unix și în fire separate în versiunea Windows.
- (4) În versiunea Windows firele nu moștesc codul părintelui și de aceea, în clasa *CObject* trebuie introdusă o metodă *CopyObject* care să asigure copierea instanței *CObject* de pe firul principal în firele similare, create la inițierea unei tranzacții.
- (5) În versiunea Unix, deoarece se crează procese care moștesc descriptorii de fișiere ai procesului părinte, aceștia trebuie închiși în procesele fii.
- (6) dimensiunea articolelor din fișierelor text PCLOG și ACLOG este mai mică cu un octet în fișierele Unix.
- (7) În Unix se lucrează cu funcția *fsync()* pentru siguranța scrierii pe disc, iar în Windows cu *commit()*.
- (8) macrodefinițiile pentru integrarea automată în sistem a obiectelor diferă.
- (9) fișierele folosite sunt setate fără buferare (pentru ca operațiile *fprintf* de scriere să se facă direct pe fișiere) cu funcția *setbuf()* în Windows și cu funcția *setbuffer()* în Unix.
- (10) maparea fișierului obiectului (versiunea persistentă) se face în Unix cu apelul:

```
m_pData = (int *)mmap(0 ,size , PROT_READ|PROT_WRITE ,MAP_FILE | MAP_SHARED ,m_df,0);  
if ( (int)m_pData==-1) Panic("ObjMan: mmap object error")
```

iar în Windows cu apelurile:

```
m_df=CreateFile(m_dataFile,GENERIC_READ|GENERIC_WRITE,FILE_SHARE_READ|FILE_SHARE_WRITE,  
NULL, OPEN_EXISTING, FILE_ATTRIBUTE_NORMAL, NULL);  
if (m_df == INVALID_HANDLE_VALUE)  
Panic("ObjMan: CreatFile object error")  
if((m_hData=CreateFileMapping(m_df, NULL,PAGE_READWRITE,0,sizeof(int)*MAX_OB,m_dataFile))==NULL)  
Panic("ObjMan: MappingFile object error")  
if ( (m_pData = (int *)MapViewOfFile (m_hData, FILE_MAP_WRITE, 0, 0,sizeof(int)*MAX_OB ) ) == NULL )  
Panic("ObjMan: MapViewOfFile object error")
```

- (11) În programul principal din versiunea Unix trebuie captat semnalul SIGCHLD, pentru a se executa apelul *wait()* în momentul în care se termină un proces fiu (se evită apariția proceselor *zombie*).

10.3.4. Exemplu de aplicație care apelează managerii de obiecte de tipul proiectat

Orice aplicație care folosește managerii de obiecte de tipul proiectat trebuie să folosească modelul prezentat la 9.4.2. Acesta înscamnă că pentru lucrul cu sistemul tranzacțional secvența de instrucțiuni care trebuie executată trebuie să cuprindă:

```
/* conectarea la sistemul tranzacțional DOSTP */
(1) ...
/* initierea unei tranzacții */
(2) ...
/* apeluri de operații sincrone sau asincrone asupra unor obiecte locale sau la distanță */
/* recepționarea rezultatelor */
(3) ...
/* terminarea operațiilor asincrone lansate */
(4) ...
/* încheierea tranzacției cu commit sau abort
(5) ...
/* deconectarea de la sistem */
(6) ...
```

Etapele 1-4 pot fi realizate, într-o aplicație foarte simplă astfel:

```
(1) if (errorCode = InitDOSTP ("Test") { /* încercare de conectare la sistem */
    printf ("Init error = %d\n", errorCode); /* încercare nereușită */
    close DOSTP(); /* în acest caz, deconectare de la sistem */
}

(2) if (errorCode = BeginTransaction () { /* încercare de inițiere a unei tranzacții */
    printf ("BeginTransaction error = %d\n", errorCode); /* încercare nereușită */
    close DOSTP(); /* în acest caz, deconectare de la sistem */
}

(3) VEC_REQUEST vrequest; /* cerere de execuție a unei operații */
int status;
vrequest.position = position;
vrequest.value = value;
errorCode = AsyncRequest ("v1", &vrequest, &request, &status, TIME_OUT)
if (errorCode)
    printf ("AsyncRequest error = %d\n", errorCode);
.....
errorCode = SyncRequest ("v2", &vrequest, &vrequest, TIME_OUT)
if (errorCode)
    printf ("SyncRequest error = %d\n", errorCode);

(4) errorCode = Synchronize();
if (errorCode)
    printf ("Synchronize error = %d\n", errorCode);

(5) again:
errorCode = EndTransmission(); /* se încerca terminarea tranzacției */
if (errorCode = ETRY_AGAIN) { /* se așteaptă 1 secundă */
    sleep(1);
    if (retry++ < MAX_RETRY) /* se reîncearcă commit */
        goto again;
    else
        CloseDOSTP(); /* dacă nu se reușește după MAX_RETRY încercări, deconectare */
/* de la sistem */
}

(6) CloseDOSTP();
```

10.4. Concluzii și direcții de dezvoltare ulterioară

În urma analizei proiectării și implementării sistemului DOSTP pot fi enunțate următoarele concluzii:

(1) Sistemul DOSTP constituie un suport software pentru dezvoltarea de aplicații distribuite. Arhitectura sa nu este complet distribuită deoarece se prevede un manager central pe un nod fixat în sistem. Introducerea managerului central este un compromis acceptat din motive de performanță; dacă nu s-ar fi folosit managerul central, atunci, pentru obținerea, într-o manieră complet distribuită, a priorității și identificatorilor tranzacțiilor, unici în întreg sistemul, ar fi trebuit folosită o metodă prezentată la 4.3.2, 4.3.3 sau 4.3.4, ceea ce implică însă o circulație mare de mesaje.

(2) Sistemul a fost proiectat astfel încât să ofere o eficiență maximă. Comunicațiile locale între componentele sistemului, în timpul desfășurării unei tranzacții, când au loc multe transferuri de date se fac prin intermediul unei zone de memorie partajate. Metoda oferă rapiditate dar are dezavantajul că necesită mai multe resurse auxiliare (semafoare pentru controlul accesului în regim de excludere reciprocă, canale pentru stabilirea legăturii) și un protocol relativ lent de conectare (se creează mai întâi zona de memorie partajată, semafoarele pentru controlul accesului, apoi se transmit identificatorii acestora prin canale FIFO). În schimb, după realizarea conectării, transferul de date decurge foarte rapid, compensându-se pierderea de timp de la început.

(3) La implementarea versiunii Unix și a versiunii Windows au trebuit să fie rezolvate diferențele între funcțiile ce lucrează cu soclurile, cu memoria partajată și cu semafoarele. În Unix se lucrează cu seturi de semafoare, fiecare set fiind caracterizat de un identificator propriu; dacă un proces posedă acest identificator, atunci el poate folosi setul de semafoare. În Windows nu există seturi de semafoare; semafoarele identificându-se prin nume; dacă un proces cunoaște numele respectiv, atunci el poate folosi semaforul respectiv. Analog stau lucrurile și cu zonele de memorie partajată: în Windows, zonele de memorie partajată (eventual cu fișiere mapate în memorie) se identifică prin nume iar în Unix printr-un identificator.

(4) Diferențele majore între versiunile Unix și Windows se datorează facilităților diferite în lucrul în procese și fire oferite de Unix respectiv Windows; acestea nu au afectat logica programelor dar au introdus necesitatea unor operații de copiere și sincronizare în plus.

(5) În Windows nu se pot folosi soclurile pentru comunicarea locală (socluri Unix), de aceea, în versiunea Windows, procesele de pe același nod comunică prin socluri TCP. Un soclu Unix este caracterizat printr-un nume de "port", care respectă regulile de nume ale fișierelor, spre deosebire de un soclu de tip Inet (TCP) care este caracterizat printr-un nume de port. În consecință, obiectelor de pe Windows li se asociază pentru a fi contactate un șir de caractere drept port iar celor de pe Unix un număr.

(6) Pe parcursul realizării testelor, au fost sesizate aspecte nespecificate ale apelurilor sistem: dacă, de exemplu, un server de comunicație de pe un nod Unix execută apelul *write()* și serverul de comunicație pereche Windows are o defecțiune sau este terminat forțat, atunci serverul de comunicație primește semnalul SIGPIPE, sesizând apariția defecțiunii; dacă nodul care are defecțiunea este nod Unix atunci apelul *write()* reusește, cu toate că partenerul de comunicație a dispărut (dispariția va fi sesizată la următorul apel *read()*). Apelul funcției sistem *read()*, în aceleași condiții sesizează dispariția nodului pereche, returnând valoarea -1.

(7) Pentru situațiile de recuperare nu se mai folosește memoria partajată, având în vedere că se transmit doar câteva comenzi și confirmări (traficul de mesaje este mic).

(8) Pe sisteme uniprocessor diferențele între timpii de execuție ale aplicațiilor care efectuează cereri asincrone către obiecte locale și între cei ai aplicațiilor care efectuează aceleași cereri sincron sunt mici. Avantajul prelucrării asincrone apare în cazurile în care aplicațiile execută cereri asincrone asupra mai multor obiecte de la distanță (prelucrările se fac paralel).

(9) Prin maniera de conectare automată a managerilor de obiecte la sistem (folosind clasele definite la 10.3.2.1 și 10.3.2.2 și macrodefinițiile de la 10.3.2.3) managerii de obiecte pot fi implementați și testați separat de sistemul DOSTP.

(10) Pentru a pune la dispoziția componentelor sistemului anumite informații precum localizarea managerului central, direcțiile pentru fișiere temporare se folosește un fișier de configurație. Folosind acest fișier, pot fi făcute modificări fără a fi necesară recompilarea surselor.

(11) O posibilă extindere a sistemului DOSTP ar consta în introducerea drepturilor de acces la sistemul DOSTP. Acest lucru poate fi realizat specificându-se o parolă în funcția de conectare la sistem *InitDOSTP()*.

(12) În varianta actuală a sistemului DOSTP se presupune că nodul pe care se afla managerul central nu se defectează. O extindere a sistemului ar putea fi făcută și în acest sens; astfel, se pot folosi fișiere log la nivelul managerului central care conțin informații despre nodurile conectate la sistem și obiectele introduse în sistem, și care se pot citi la relansarea managerului central, și eventual adăugat un protocol de interogare al acestora, dacă mai sunt în funcțiune.

Capitolul 11

Concluzii

11.1 Concluzii generale

Scopul acestei teze, intitulată "Contribuții la proiectarea sistemelor de operare distribuite" este de a-și aduce contribuția, într-o manieră originală, la tratarea unor probleme și aspecte ale proiectării sistemelor de operare distribuite, un domeniu relativ nou pentru informatica din țara noastră. Autorul și-a propus ca în această lucrare să abordeze atât aspecte teoretice ale acestui domeniu, considerând că domeniul respectiv se află încă în faza de identificare și încercare de soluționare a problemelor sale fundamentale, dar și aspecte practice constând în implementarea unor API-uri destinate să asiste utilizatorii în dezvoltarea de aplicații în rețele locale Unix, Microsoft Windows (Win32) și Novell NetWare.

Lucrarea este împărțită în două mari părți: prima parte, care cuprinde primele cinci capitole, are ca obiective principale prezentarea problemelor specifice unui sistem de operare distribuit, accentuându-se problemele de comunicare și sincronizare, paradigmele interacțiunii dintre procese precum și legătura sistem de operare distribuit-aplicație distribuită. În cadrul primei părți se prezintă modul de proiectare și implementare a două subsisteme software. Scopul acestor subsisteme software este: (1) controlul accesului concurrent la o bază de date în rețele Unix, Microsoft Windows (Win32) și Novell NetWare; (2) furnizarea unor API-uri pentru comunicații de grup în rețele Unix, Microsoft Windows (Win32) și Novell NetWare.

A doua parte a tezei, care cuprinde capitolele 6-10, reprezintă un studiu sinteză asupra modalităților și criteriilor de proiectare a sistemelor de programare distribuită bazată pe obiecte (SPDBO) precum și proiectarea și implementarea unui suport software bazat pe modelul tranzațiilor atomice pentru SPDBO în rețele Unix și Microsoft Windows (Win32).

În urma parcurgerii celor 10 capitole ale lucrării, pot fi enumerate următoarele concluzii:

Partea I: SISTEME DISTRIBUITE

Capitolul 1: Sisteme de operare distribuite. Suportul oferit prelucrărilor distribuite: probleme specifice, modele de proiectare.

Sistemul de operare asigură un mediu de execuție pentru programele aplicative, asumându-și în contul acestora toate sarcinile necesare, începând de la activități de gestiune, partajare și control al accesului la diverse resurse, până la asigurarea unor interfețe uniforme și a unor niveluri diferențiate de abstractizare/virtualizare.

În capitolul 1 se abordează și problema interacțiunii prelucrare distribuită-sistem de operare, prezentându-se, prin sistemul TDL, o modalitate de conectare a programelor în rețele Unix și Windows. În mod curent, utilizatorii unui mediu eterogen dezvoltă programe în manieră tradițională: se scriu procese separate ca programe în diferite limbaje (Pascal, C) implementate pe anumite procesoare și apoi printr-un program separat se încarcă aceste programe pentru execuție. Deseori, aceste programe sunt scrise astfel încât înglobează în ele specificații privind disponerea lor pe procesoare sau chiar privind structura rețelei. Acest mod de lucru face dificilă sau chiar imposibilă utilizarea aplicației într-un mediu diferit, complică o dezvoltare ulterioară a aplicației care ar schimba structura programelor și suportul pentru siguranță la defecte (prin restartarea unei noi copii a programului pe un procesor diferit). O manieră nouă de abordare a acestor probleme este de a separa activitatea de programare a programelor de activitatea de cooperare a lor în codul aplicației. Vor exista astfel două nivele de programare în cadrul unei aplicații: *un nivel de programare al programelor componente ale unei aplicații și un nivel de programare a aplicației care utilizează programele dezvoltate anterior*. O aplicație distribuită în cadrul sistemului proiectat, TDL, constă într-un set de taskuri care comunică prin cozi de mesaje. . Descrierile taskurilor și declarațiile de tipuri specifică modul în care taskurile sunt lansate în execuție și executate ca procese concurente, tipul datelor prin care comunică procesele, precum și cozile de date necesare proceselor producătoare și consumatoare. Există trei faze distincte în dezvoltarea unei aplicații distribuite utilizând limbajul TDL: crearea unei biblioteci de taskuri, crearea unei aplicații utilizând biblioteca de taskuri și execuția aplicației. În mediul de execuție există trei componente active: *taskurile de aplicații, serverul și planificatorul*. În orice nod al rețelei în care se execută un task, se activează și un exemplar de server.

O concluzie a subcapitolelor 1.3 și 1.4 este aceea că există o tendință în dezvoltarea sistemelor de operare distribuite de a fi structurate ca un set de servere de sistem situate deasupra unui micro-kernel minimal (Amoeba, Chorus, Mach, Plan 9, Topaz, V-System); acest mod de structurare a componentelor sistemului de operare oferă suport pentru modularitate, portabilitate, scalabilitate, dar și adaptare la noi arhitecturi

hardware în condiții de modificare rapidă a cerințelor pieței. În cadrul acestei arhitecturi, micro-kernelul furnizează câteva funcții de bază, independente de particularitățile sistemului de operare, legate de alocarea procesorului și gestiunea memoriei precum și facilitarea de comunicare între procese (Inter-Process Communication-IPC). Astfel, se oferă severelor de sistem posibilitatea de a interacționa independent de mediul în care sunt executate, mediu care poate fi multiprocesor, multicomputer sau configurație de rețea.

Serverele de sistem furnizează, pe baza primitivelor puse la dispoziție de micro-kernel, un set de funcții specifice sistemului de operare, care realizează gestiunea resurselor fizice și logice (fișiere, dispozitive de I/O, linii de comunicație); un set de servere de sistem este uneori denumit *subsistem*. Subsistemele, pe de o parte, separă funcțiile sistemului de operare în seturi de servicii furnizate de servere autonome, și, pe de altă parte, oferă seturi de API-uri programelor de aplicație.

Pentru asigurarea unui mediu de programare bazat pe standardul Unix, un subsistem special poate fi rezervat pentru implementarea apelurilor sistem Unix. Serviciile Unix pot fi astfel extinse la medii distribuite, multiprocesor sau chiar timp real, într-o manieră transparentă pentru programele aplicative.

Arhitectura bazată pe un micro-kernel oferă, pe de altă parte și suport adecvat pentru dezvoltări ulterioare ale sistemului de operare. Omogenitatea interfețelor serverelor de sistem, datorată folosirii aceluiași facilități IPC, oferite de micro-kernel, permite programatorilor de sistem să dezvolte noi servere pe care să le integreze ulterior în sistem. De exemplu, noi strategii de gestiune a fișierelor sau facilități de toleranță la defecte pot fi înglobate în noi servere, care întâi vor fi testate ca programe utilizator obișnuite, fără a afecta astfel funcționalitatea sistemului în ansamblu și integrate ulterior, după testarea lor, în sistem.

O altă concluzie a capitolului întâi este legată de instrumentele ce asistă proiectantul în dezvoltarea sistemului de operare distribuit; aceste instrumente implică și un anumit grad de abstractizare. Astfel, în proiectarea sistemelor de operare distribuite există două modele: *modelul obiect* și *modelul proces*. Conceptul de obiect folosit în modelul obiect dispune de proprietățile de încapsulare și abstractizare. După cum obiectele active permit execuția paralelă a mai multor procese componente sau a unui singur proces, modelul obiect se divide în *modelul obiecte multiactive* și *modelul obiecte active secvențiale*.

Capitolul 2. Comunicații în sisteme distribuite. Protocoale.

Arhitectura bazată pe un micro-kernel minimal furnizează, pe de o parte, o soluție de rezolvare într-o manieră structurată și modulară a problemei creșterii continue a complexității sistemelor de operare, dar, pe de altă parte, furnizează și avantajul unui cadru general pentru proiectarea, dezvoltarea și integrarea sistemelor "deschise". Termenul de arhitectură deschisă este tratat în capitolul 2 în sensul accepțiunii OSI: un sistem *deschis* este în accepțiunea ISO/OSI un sistem cu o arhitectură standardizată public, care poate să fie conectat și deci să comunice cu orice alt sistem deschis care utilizează același model de arhitectură. La baza conceptului de arhitectură deschisă stau două principii esențiale: (1) fiecare nivel valorifică serviciile asigurate de nivelele inferioare, culminând cu nivelul aplicație, căruii îi sunt furnizate cât mai multe servicii pentru proiectarea de aplicații distribuite; (2) este asigurată independența fiecărui nivel, prin definirea serviciilor care trebuie asigurate nivelului următor, indiferent de modul de realizare al acestor servicii.

Există două direcții în cadrul cercetărilor asupra protocoalelor nivelului transport: (1) proiectarea și implementarea unor protocoale specifice unui sistem de operare distribuit (*protocoale cu caracter restrâns*); (2) reexaminarea proiectării și implementării *protocoalelor cu caracter universal*. Pentru ambele tipuri de protocoale există argumente pro și contra. Astfel, *protocoalele cu caracter universal* oferă avantaj precum portabilitate, independența aplicațiilor, dar și dezavantaje care provin din ignorarea aspectelor specifice ce apar în sistemele distribuite; aceste dezavantaje constau într-o complexitate excesivă (protocoalele de transport la nivelul ISO/OSI sunt un exemplu în acest sens) cu repercursiuni asupra performanțelor, orientarea preponderentă pe variante de orientate pe conexiune, lipsa unui suport pentru implementarea unor facilități speciale în rețelele locale precum partajarea resurselor multicasting. *Protocoalele cu caracter restrâns* au ca principal avantaj performanțe mai bune, (întirzieri mici, printr-o execuție eficientă) dar și dezavantajul aplicabilității numai într-un anumit tip de rețea și/sau arhitectură de aplicații.

Există o serie de cerințe pe care trebuie să le respecte un nou protocol Internet, rezolvate în totalitate de IPv6, motiv pentru care acesta a fost desemnat ca standard al noii generații de protocoale Internet: (1) *compatibilitate cu IPv4*; (2) *facilități extinse în ceea ce privește posibilitățile de adresare și de rutare*. (3) *noi funcționalități* precum servicii real-time, autentificare, asigurarea securității, integrității și confidențialității datelor.

Capitolul 3. Paradigme ale interacțiunii dintre procese în sisteme distribuite.

Într-un sistem distribuit, o aplicație poate fi structurată în mai multe procese repartizate, pe mai multe procesoare, dispuse pe mai multe noduri într-o rețea. Cele mai cunoscute clase de primitive utilizate pentru comunicația între procese sunt: (1) *comunicarea prin mesaje*, care se asociază modelului client/server și la care fluxul de mesaje poate fi unidirecțional sau bidirecțional; (2) *apel de procedură la distanță (RPC)*, care oferă un mecanism puternic tipizat pentru transferul bidirecțional al informației; (3) *tranzacții*, care oferă un mecanism complex pentru executarea și sincronizarea de operații asupra unor obiecte partajate, asigurând recuperarea informației în cazul apariției de erori hardware și software; (4) *comunicare prin conducte (pipes)*

ca o generalizare a modelului conductelor locale introdus de sistemele Unix. Aceste clase de primitive permit implementarea a patru tipuri de procese într-un program distribuit: *filtre, clienți, serveri și perechi*.

Capitolul 4 prezintă problemele care apar la implementarea primitivelor din clasele 1,2 și 4. Modelul client-server este abordat atât sub aspectul concurenței la nivelul serverelor și clienților, al protocolului utilizat, orientat sau neorientat pe conexiune, cât și sub aspectul cerințelor impuse realizării unor sisteme client-server.

Comunicarea prin mesaje este asociată cu modelul client/server și constituie o extensie a comunicării locale interprocese prin mesaje; este o formă de comunicare în care procesele utilizator generează explicit mesajele transmise prin subsistemul de comunicație și utilizează explicit mecanismele de livrare și recepție ale mesajelor. Referitor la semantica primitivelor pentru comunicare, se separă primitivele în primitive *cu și fără blocare, sigure și nesigure, cu și fără buferare*; în capitolul 4 se arată, că, pentru primitiva *send*, această clasificare este generată de faptul că un program apelant al unei operații *send* poate obține controlul în următoarele 4 moduri: (1) semantica *no-wait-send* sau *asynchronous-send*- procesul emitator așteaptă doar preluarea mesajului de către subsistemul de comunicație; (2) semantica *synchronization send*- procesul emitator așteaptă până când nodul destinație primește mesajul; (3) semantica *rendez-vous* - procesul emitator așteaptă până când procesul receptor primește mesajul; (4) semantica *remote-invocation send* sau *remote procedure call*- procesul așteaptă până când primește un răspuns de la procesul receptor, în urma tratării cererii.

Modul de lucru cu buferare la nivelul programului utilizator, pentru primitiva *receive*, este exemplificat prin aplicația prezentată în anexa C4, în care mesajele primite de un server aflat pe un calculator dintr-o rețea Novell sunt memorate într-o coadă de mesaje, lansându-se o rutină *EventServiceRoutine*.

Modelul de comunicație bazat pe *apel de procedură la distanță (RPC=Remote Procedure Call)* permite unui proces să apeleze și să execute o procedură de la distanță ca și pe o procedură locală. Pe timpul execuției procedurii de la distanță, apelantul este blocat, până când se recepționează rezultatul. Prin RPC se permite astfel extinderea apelului de proceduri, folosit în limbajele de programare la mediile distribuite. Modelul RPC este atractiv pentru că oferă un mecanism de comunicație în cadrul programelor distribuite, care se poate caracteriza astfel: (1) este simplu și familiar (datorită similarității cu apelurile de proceduri locale); (2) este general (apelul de procedură este folosit în toate aplicațiile); (3) poate fi implementat eficient. Modelul RPC simplifică construcția programelor distribuite, permițând ignorarea de către acestea a detaliilor de comunicație și a erorilor de transmisie. Sunt necesare cinci elemente pentru a se obține efectul execuției unei proceduri locale într-un mediu distribuit: *procesul client, codul stub pentru client, mecanismul de transport, codul stub pentru server și procesul server*. Modelul RPC poate fi interpretat și ca o extindere a primitivelor sigure, cu blocare, structurate: *Send, GetRequest, SendResponse*. Spre deosebire de aceste primitive, datorită simplității și totodată generalității sale, modelul RPC poate fi însă extins la mai mulți clienți și la mai multe servere, ca și la medii eterogene.

Exemplificările pentru modelul RPC sunt realizate în capitolul 3 pentru DCE RPC, Microsoft RPC și Sun RPC; se arată cum se realizează legătura clienților la servere, concurența multifir la nivelul serverelor, dezvoltarea aplicațiilor distribuite.

Capitolul 3 abordează și celelalte două tipuri de interacțiuni între procese în sisteme distribuite: comunicațiile de grup și comunicația la distanță prin *pipe*. Un *grup multicast* este un set de procese care reprezintă destinația aceleiași mesaj. *Comunicațiile de grup (sau multicast)* se referă la comunicații de tipul unul-la-mai-multi; grupurile sunt entități dinamice; pot fi create și distruse grupuri, iar componența grupurilor poate fi modificată. În cadrul unui grup, procesele pot fi egale sau poate fi creată o ierarhie între procese; de exemplu, un proces poate îndeplini rolul de coordonator, în timp ce celelalte pot îndeplini aceeași sarcină de lucru. Primul model prezintă o simetrie; dacă unul dintre procese dispare, grupul devine mai mic, dar continuă să lucreze. În cel de-al doilea, dispariția coordonatorului pune în pericol întregul grup. În schimb, în primul model, luarea unei decizii este o operație complicată, care implică un proces de votare. Pentru *gestiunea grupurilor* (crearea, stergerea grupurilor, modificarea componenței grupurilor) există două soluții: (1) o soluție *centralizată*, în cadrul căreia un server de grupuri primește și tratează toate cererile care se referă la componența grupurilor; (2) o soluție *distribuită*, în cadrul căreia un proces transmite un mesaj tuturor proceselor din toate grupurile la intrarea sa într-un grup, ca și la parasirea unui grup. Grupurile pot fi *închise*, caz în care numai membrii unui grup pot transmite mesaje grupurilor sau *deschise*, caz în care orice proces din sistem poate transmite mesaje la orice grup.

Comunicarea între procese prin conducte (*pipes*) este un mecanism care permite transferul datelor între procese utilizând strategia FIFO, concomitent cu sincronizarea proceselor implicate în transfer. Ea poate fi aplicată și modelului client server prezentat anterior deși nu îi este caracteristică. Microsoft Windows (Win32) furnizează pentru comunicația la distanță *conducte cu nume*. În subcapitolul 3.3.2 se arată cum pot fi proiectate procesele cu rol de server într-o comunicație prin conducte astfel încât să funcționeze multifir; soluția este să se creeze mai multe instanțe ale unei conducte, iar fiecărei instanțe să i se asocieze un fir în cadrul căruia se așteaptă cereri de conectare și se deservesc cereri de servicii.

O concluzie care rezultă din cele prezentate în capitolul 3 ar fi următoarea: comunicarea prin mesaje, care este considerată o caracteristică definitorie a arhitecturii bazate pe un micro-kernel este cea mai naturală cale de a structura sistemele în care componentele sunt distribuite pe un set de procesoare slab cuplate; ea oferă o izolare foarte clară între componentele sistemului, impunând reguli de comunicație explicită între

acestea, dar, în același timp, oferă și o soluție flexibilă pentru a asambla ierarhic componentele distribuite, în entități globale, situate pe nivele superioare. Comunicația între componente poate îmbrăca forma unor protocoale în care se folosesc primitivele *send/receive simple*, care asigură transferul de date, sau forma unor mesaje *send/reply*, care pot fi combinate într-o formă de RPC, pentru a deservi cât mai bine comunicațiile de tip client-server.

Capitolul 4. Sincronizarea în sisteme de operare distribuite.

Comunicația între procese asigură nu numai transferul propriu zis al datelor ci și sincronizarea execuției proceselor. Capitolul 4 este dedicat în întregime problemelor de sincronizare în sisteme distribuite. Datorită distribuției resurselor, întâzierilor de transmisie, lipsa unei memorii comune și deci lipsa unor informații globale, algoritmi și metodele pentru sincronizare din sistemele centralizate nu pot fi utilizate în sistemele distribuite. Există două clase mari de mecanisme de sincronizare în sistemele distribuite : **centralizate și distribuite**. Unele dintre aceste mecanisme se bazează pe conceptul de *excludere reciprocă* care presupune o ordonare a unor secvențe de cod executate de procesele concurente; aceste mecanisme pot utiliza diferite tehnici precum: *un proces central coordonator* sau un *token circulant*. Alte mecanisme se bazează pe înregistrarea evenimentelor semnificative în cursul unor prelucrări asincrone; acestea utilizează obiecte precum *contori de evenimente și secvențiatori*, sau *semafoare distribuite*.

Capitolul 4 prezintă următorii algoritmi și mecanisme pentru sincronizare: (1) *algoritmi distribuți pentru excludere reciprocă*: algoritmul lui Lamport, algoritmul lui Ricard și Agrawala, algoritmul lui Rymond pentru excludere reciprocă în accesul la k resurse identice. Caracteristica comună a acestor algoritmi este faptul că se bazează pe mărcile de timp asociate cererilor de intrare în secțiunea critică și pe schimb de mesaje; (2) *algoritmi care utilizează un token circulant* pentru excludere reciprocă: algoritmul pentru sisteme structurate în inel și algoritmul lui Chandy, pentru sisteme nestructurate în inel precum și o soluție proprie a autorului pentru implementarea mecanismului de *secvențiator circulant*; (3) *algoritmi broadcast* care folosesc facilitatea transmisiei broadcast. S-au prezentat: soluția Andrews de implementare a *semafoarelor distribuite* precum și soluții proprii ale autorului pentru implementarea *contorilor de evenimente distribuți și a secvențiatorilor distribuți*. Prezentarea tuturor algoritmilor s-a făcut în mod unitar, utilizându-se pentru descrierea lor, descrierea SR introdusă de Andrew [And91]. Toți algoritmi prezentați au fost implementați în rețele locale Unix, Microsoft Windows (Win32) și Novell NetWare; în anexele B1-B20 se prezintă și codul sursă al implementărilor (din lipsă de spațiu, se prezintă pentru fiecare algoritm numai o singură variantă: Unix, Microsoft Windows sau Novell NetWare). În anexa B20 se arată cum pot fi utilizați contorii de evenimente distribuți și secvențiatorii distribuți pentru a obține excluderea reciprocă.

O serie de algoritmi și mecanisme de sincronizare prezentați în acest capitol vor fi folosiți în capitolele următoare pentru realizarea sincronizării în diferite aplicații. Astfel, pentru controlul accesului concurrent la o bază de date în rețele locale Unix, Microsoft Windows (Win32) se folosesc semafoare distribuite; pentru implementarea unui toolkit pentru comunicații de grup, în scopul obținerii accesului exclusiv în momentul alegerii unui identificator pentru un nou grup, se folosește o metodă cu token circulant în varianta în inel a implementării și o metodă broadcast în varianta broadcast a implementării.

Capitolul 5. Proiectarea și implementarea unor mecanisme suport pentru dezvoltarea de aplicații în rețele Unix, Microsoft Windows (Win32) și Novell NetWare.

Prima parte a lucrării se încheie cu un capitol cu caracter pur aplicativ, în care se prezintă proiectarea și implementarea a două toolkit-uri.

Subcapitolul 5.1. este dedicat în întregime *proiectării și implementării unui toolkit pentru controlul accesului concurrent la o bază de date*. Sunt prezentate 3 soluții cu semafoare și o soluție bazată pe un proces coordonator. Implementările au fost realizate în rețele locale Unix, Microsoft Windows (Win32) și Novell NetWare. În finalul capitolului 5.1, în paragraful 5.1.4. se prezintă o comparație a celor 4 soluții propuse și citeva concluzii rezultate în urma implementărilor; dintre care se enumeră: (1) Toate soluțiile prezentate necesită implementarea unor variabile de stare; acestea pot fi implementate fie separat, în cadrul unui proces coordonator, fie pot fi simulate ca "variabile comune" prin asimilarea lor cu valorile unor semafoare distribuite; (2) Prin faptul că mențin o listă de procese, prima soluție cu semafoare și soluția procesului central pot planifica procesele pentru execuție după diverși algoritmi; (3) Pe de altă parte, a doua și a treia soluție cu semafoare, gestionând liste, sunt mai rapide; (4) Implementarea soluției bazată pe un proces coordonator se face având în vedere raportul dintre timpul necesar tratării unei cereri și rata de sosire a cererilor, pentru a elimina posibilitatea apariției situației de "sufocare" a serverului.

Subcapitolul 5.2. este dedicat *proiectării și implementării unui subsistem software* (în sensul descris mai sus, acela de set de servere de sistem) pentru comunicații de grup, care, la nivelul utilizator, se materializează într-un toolkit. Sunt propuse 4 variante: o variantă *semidistribuită*, două variante complet distribuite, care folosesc *transmisie broadcast*, respectiv protocolul *two phase commit* și o variantă *în inel*. În toate soluțiile se consideră că se cunoaște numărul participanților în sistem pentru a se putea lua decizia dacă un mesaj a fost recepționat sau nu de *toți* membrii unui grup. În variantele 1 și 3 trebuie cunoscută și identitatea participanților (nodurilor). În varianta a 4 a trebuie cunoscută identitatea vecinilor din inel.

În proiectarea și implementarea celor două toolkit-uri din capitolul 5 s-au folosit concepte, algoritmi și diverse observații prezentate în toate capitolele anterioare; astfel au fost folosite protocoalele TCP, UDP și IPX prezentate în capitolul 2, s-a ținut seama de observațiile prezentate în subcapitolul 3.1.2. relativ la concurența la nivelul serverelor, s-au folosit semafoarele distribuite prezentate la 4.3.4., algoritmi pentru excludere reciprocă prezentați la 4.2. și s-au folosit ambele variante de realizare a serverelor: *multiproces* sau *multifir* (capitolul 1 și 3). Testele au fost realizate utilizând sistemele de operare Linux 2.0, DEC Alpha OSF 1.3, Microsoft Windows NT 4.0 și Windows 95 precum și Novell NetWare 3.11.

Partea a II-a: CRITERII DE PROIECTARE A SISTEMELOR DE PROGRAMARE DISTRIBUITĂ BAZATĂ PE OBIECTE (SPDBO). PROIECTAREA UNUI SUPORT SOFTWARE BAZAT PE MODELUL TRANZACȚIILOR ATOMICE PENTRU SPDBO ÎN REȚELE UNIX ȘI WINDOWS.

Capitolul 6. Sisteme de programare distribuită orientată pe obiecte (SPDBO). Criterii de proiectare

Cea de-a doua parte a lucrării cuprinde suportul teoretic necesar pentru proiectarea sistemelor de programare distribuită bazată pe obiecte precum și proiectarea și implementarea în rețele locale a unui suport software pentru lucrul cu tranzații atomice în aplicații distribuite.

Abordarea acestei teme s-a făcut ca urmare a constatării că, în ultimii ani, se conturează tot mai accentuat tendința ca, serviciile de gestiune ale tranzațiilor, tradițional implementate în sistemele de gestiune a bazelor de date să fie implementate în sistemele de operare. Aceasta permite ca serviciile pentru controlul concurenței și pentru lucrul cu tranzații atomice să fie disponibile oricărui program de aplicație, nu numai clienților unui manager al unei baze de date. Exemple sunt sistemul de operare Novell NetWare care pune la dispoziție funcții pentru lucrul cu tranzații atomice și pentru manipularea lock-urilor asupra fișierelor (*Synchronization Services*) și sistemul distribuit de fișiere Quick Silver Distributed System(CPS93), care furnizează o interfață de acces similară cu cea a sistemului Unix, cu observația că toate operațiile sunt realizate într-o manieră atomică; aceasta înseamnă că, de exemplu, un client poate șterge, rescrie un fișier sau redenumi un director și apoi decide dacă modificările vor fi permanentizate sau nu. QDFS garantează că, fie toate modificările vor fi permanentizate, fie niciuna nu va fi realizată; aceasta garanție se face chiar pentru situații de avarii în rețea (defecțiuni ale nodurilor, ale liniilor de comunicație).

Partea a doua a acestei lucrări abordează aspectul tranzațiilor atomice din perspectiva execuției lor în cadrul unui sistem distribuit de obiecte. Un *obiect* este definit ca o entitate care încapsulează informații de stare (sau date) și un set de operații sau proceduri care manipulează aceste date într-un mod controlat, astfel încât pot fi tratate, colectiv, ca o singură unitate. În capitolul 6 se tratează aspecte relative la componența și structura obiectelor, ale gestionării obiectelor, interacțiunii dintre obiecte și al reprezentării obiectelor în memorie. O concluzie care se poate enunța în urma parcurgerii capitolului 6 este aceea că proiectarea sistemelor de programare distribuită este o sarcină complexă și dificilă, care trebuie să țină seama de o multitudine de caracteristici și funcții necesare aplicațiilor, unele specifice numai unei clase de aplicații. Abordarea obiectuală simplifică însă proiectarea, ținând seama că obiectele servesc ca unități de protecție, sincronizare, recuperare și mobilitate.

Capitolul 7. Controlul concurenței și recuperarea obiectelor

Capitolul 7 este axat în principal pe două probleme: realizarea atomicității tranzațiilor efectuate asupra obiectelor și controlul concurenței acțiunilor la nivelul obiectelor.

În legătură cu *realizarea atomicității tranzațiilor*, este prezentat mai întâi modelul tranzațional, apoi protocolul *two-phase commit*, cu implementarea sa OSI/ISO, protocolul *write-ahead logging*, alături de tehnicile cunoscute pentru recuperarea obiectelor într-un sistem tranzațional: *value-logging*, *redo-only* și *undo-only logging*. În finalul subcapitolului 7.3., al cărui subiect este realizarea atomicității tranzațiilor se prezintă structura unui sistem de recuperare de obiecte; el este alcătuit din: (1) un *manager de recuperare* al cărui rol este de a urmări evoluția tranzațiilor, și de a lua deciziile adecvate conform tehnicii de recuperare alese, în timpul execuției acțiunilor tranzației, ca și la execuția comenzilor *commit*, *abort*; (2) un *manager pentru log*, utilizat de managerul de recuperare și de alte componente ale sistemului distribuit pentru a scrie articole în log; scrierea articolelor în fișierul log este cerută fie de managerul de recuperare, fie de sistem.

În legătură cu *mecanismele de control ale concurenței* se arată că acestea trebuie să asigure următoarele: (1) consistența obiectelor; (2) realizarea fiecărei tranzații atomice într-un timp finit. În plus, un *bun* mecanism de control al concurenței se caracterizează astfel: (1) permite execuția paralelă a acțiunilor, în scopul satisfacerii cerințelor de performanță ale sistemelor; (2) implică prelucrări minore; (3) se execută cu performanțe satisfăcătoare chiar și în rețele în care întârzierile de transmisie sunt semnificative; (4) impune puține constrângeri în structura tranzațiilor. Soluția tradițională de rezolvare a problemelor de concurență în accesul la un sistem de obiecte este prin serializarea *tranzațiilor*. Problema consistenței datelor în aplicațiile tradiționale de baze de date a fost redusă la testarea serializabilității planificărilor; acest lucru se datorează faptului că nu se cunosc constrângerile de consistență: fiecare acțiune într-o tranzație se abstractizează printr-o operație *read* sau *write*. Problema obținerii serializabilității poate fi descompusă în două subprobleme: *sincronizarea read-write* și *sincronizarea write-write*, notate pe scurt sincronizările *rw* și *ww*. Sincronizarea *rw* se referă la serializarea tranzațiilor în așa fel încât fiecare operație *read* citește aceeași valoare ca și când

tranzacțiile ar fi executate serial; sincronizarea *rw* se referă la serializarea tranzacțiilor în așa fel încât ultima operație *write* a fiecărei tranzacții lasă datele în aceeași stare ca și execuția serială. Respectarea serializărilor *rw* și *rw* generează o prelucrare consistentă. Astfel, deși sistemul de gestiune al bazei de date nu dispune de informații asupra constrângerilor de consistență specifice aplicației, garantează consistența, permițând numai execuții serializabile ale tranzacțiilor concurente. În capitolul 7 se clasifică mecanismele tradiționale de rezolvare a controlului concurenței în: mecanisme bazate pe blocarea cu *lock*-uri, mecanisme bazate pe *mărcile de timp* și mecanisme pentru *controlul optimist* al concurenței. Conceptul de serializabilitate al tranzacțiilor este propriu mecanismelor tradiționale pentru controlul concurenței. Dacă sunt disponibile informații semantice despre tranzacții și despre operațiile lor pot fi permise și planificări neserializabile, dar totuși consistente. O soluție pentru rezolvarea cerințelor *tranzacțiilor lungi* este renunțarea la planificările serializabile și folosirea unor planificări care se bazează pe informațiile semantice relative fie la obiecte, fie la operațiile specifice aplicației. În subcapitolele 7.4.7. și 7.4.8. se prezintă două posibilități de extindere a mecanismelor tradiționale pentru controlul concurenței: *mecanismul blocării altruiste* și *mecanismul validării instanțelor*.

Capitolul 8. Controlul interblocării obiectelor

Problemele care se pun în legătură cu controlul interblocării în sistemele distribuite sunt asemănătoare cu cele pentru sistemele centralizate. În sistemele distribuite, însă, datorită dispersiei informațiilor pe mai multe noduri, interblocările sunt și mai greu de prevenit, evitat și detectat.

Utilizarea tranzacțiilor atomice facilitează procedura de distrugere și reluare a proceselor implicate într-o interblocare; în cazul tranzacțiilor atomice, acestea pot fi readuse în starea de înaintea începerii tranzacției, pe cînd reluarea unor procese simple de la o stare anterioară necesită mai multe complicații.

Pentru a detecta și rezolva eficient interblocarea, este necesară utilizarea unor informații suplimentare despre aplicații. Aceste informații se referă la tipul cererii de alocare a resurselor. În consecință, algoritmi de detecție a interblocărilor se pot clasifica în funcție de modelele folosite pentru cererile de alocare: *modelul AND*, *modelul OR*, *modelul AND-OR*, *modelul (\wedge^*)* și *modelul fără nici o restricție*.

Interblocările pot fi detectate și rezolvate utilizând algoritmi *centralizați*, *ierarhici sau distribuiți*. Algoritmi de detecție și rezolvare a interblocărilor sunt elaborați avînd la bază următoarele concepte: transmiterea căilor grafului dependențelor tranzacției, urmărirea unui mesaj de probă de-a lungul arcelor grafului dependențelor tranzacției, activarea unei prelucrări cu difuzie și stabilirea stării globale consistente a sistemului.

În capitolul 8 se prezintă din fiecare clasă de algoritmi (centralizați, ierarhici sau distribuiți) cite cel puțin un algoritm exemplu. Pentru clasa algoritmilor distribuiți se prezintă un algoritm *edge-chasing* pentru detecția interblocării în modele AND și un algoritm *cu difuzie* pentru detecția interblocării în modele OR.

În sistemele distribuite în care sunt implementate tranzacții atomice și timp global sunt folosiți în practică doi algoritmi pentru *prevenirea interblocărilor*. Ambii se bazează pe ideea asignării fiecărei tranzacții a unei mărci de timp, în momentul în care începe. Fiecărei tranzacții i se mai asignează și o prioritate, care poate fi numărul de proces; aceasta urmează să fie folosită în cazul în care două tranzacții au aceeași marcă de timp. Cei doi algoritmi sunt: (1) algoritmul *wait_die* care se bazează pe o tehnică nepreemptivă; (2) algoritmul *wound_die* care se bazează pe o tehnică preemptivă. Ambii algoritmi evită amînarea la infinit a tranzacțiilor, cu condiția ca, la reluarea tranzacțiilor după abortare, să se păstreze aceeași marcă de timp.

Capitolul 9. Proiectarea unui suport software pentru programare distribuită bazată pe obiecte și pe trazații atomice (DOSTP)

Există actualmente un interes crescut în metodologii și tehnici cu caracter general care simplifică construcția sistemelor software distribuite robuste și eficiente. Acest capitol prezintă o astfel de metodologie bazată pe tranzacții atomice și include proiectarea unui suport software pentru programarea distribuită bazat pe obiecte (DOSTP). Tranzacțiile atomice micșorează dificultatea scrierii unor programe de aplicații distribuite prin simplificarea tratării defecțiunilor nodurilor și prin rezolvarea problemelor de concurență. Tranzacțiile atomice sunt, la ora actuală o paradigmă caracteristică programării distribuite; ele au fost incluse într-o serie de sisteme și chiar limbaje orientate pe obiecte.

Sistemul DOSTP a fost proiectat astfel încît să respecte următoarele cerințe: (1) simplitate de apel ale funcțiilor sistemului DOSTP la nivelul aplicațiilor utilizator; (2) modularitate și extensibilitate, în sensul ca pot fi introduși în sistem noi manageri de obiecte printr-un efort de programare minim; (3) sistemul trebuie să ofere posibilitatea execuției în paralel a mai multor operații asupra unor obiecte diferite; (4) sistemul trebuie să ofere posibilitatea recuperării obiectelor în urma defecțiunii unuia sau mai multor noduri implicate într-o tranzacție atomică; (5) sistemul trebuie să se execute în mod utilizator; (6) sistemul trebuie să ofere o comunicație sigură; (7) obiectele trebuie să fie portabile (cu modificări minime, absolut necesare) atît în rețele UNIX cat și în rețele Microsoft Windows (Win32); (8) tranzacțiile să nu poată fi amînate la infinit sau restartate la infinit.

În ceea ce privește maniera de proiectare a sistemului DOSTP se fac următoarele observații: (1) DOSTP folosește *modelul integrat* ca arhitectură generală; (2) pentru execuția tranzacțiilor atomice se crează

cite un proces agent pentru fiecare obiect de la distanță implicat într-o tranzacție; (3) DOSTP utilizează *varianta dinamică a obiectelor active*, iar localizarea obiectelor se face folosind soluția *serverelor de nume distribuite*; (4) concurența la nivelul obiectelor se rezolvă printr-o *schemă optimistă*; (5) pentru recuperarea obiectelor, DOSTP utilizează tehnicile *value-logging și redo-only logging*; (6) pentru evitarea restartărilor la infinit ale tranzațiilor, în DOSTP s-a implementat un mecanism *wound-die*.

În capitolul 9 se prezintă în detalii arhitectura sistemului DOSTP și proiectarea elementelor sale componente: serverele de comunicații, managerii de tranzații, managerii de obiecte, serverele de nume, managerul central, procesele utilizator care realizează tranzațiile. Sistemul de recuperare al obiectelor nu este implementat în DOSTP ca proces distinct; el se compune din module plasate în managerii de obiecte și în managerii de tranzații și din fișierele lor *log* aferente.

Capitolul 10. Implementarea sistemului DOSTP în rețele Unix și Microsoft Windows (Win32)

Capitolul 10 este dedicat în întregime descrierii implementării sistemului DOSTP în rețele locale UNIX și rețele Microsoft Windows (Win32). Descrierea este făcută comparativ, deoarece există deosebiri între cele două implementări, generate în special de faptul că în Unix s-au implementat servere multiproces iar în Windows servere multifir și de faptul că există diferențe în lucrul cu semafoare și memoria partajată între cele două sisteme. Un element important, care trebuie subliniat, este proiectarea obiectuală a managerilor de obiecte. Scopul principal în realizarea managerilor de obiecte este ascunderea detaliilor de conectare cu sistemul tranzațional. Managerii de obiecte sunt realizați astfel încât ei comunică numai cu serverul de comunicație local, fără a avea cunoștință dacă serviciile le-au fost solicitate local sau de la distanță. Orice manager de obiecte care urmează să se integreze în sistemul DOSTP trebuie să implementeze următoarele cinci seturi de funcții: (1) *funcția de conectare la sistem*; (2) *funcția de eliminare din sistem*; (3) *funcția de lansare* a unei noi tranzații, la solicitarea unui client; (4) *funcția de recepționare/ transmisie* a pachetelor din zona din/in zona de memorie partajată. (5) *funcțiile necesare modului de lucru tranzațional*. Deoarece funcțiile din seturile 1-4 trebuie să le asigure orice manager de obiect, o primă consecință este că acestea, o dată scrise, realizarea unui nou manager de obiecte poate fi făcută fără a se cunoaște detalii despre implementarea sistemului tranzațional; trebuie implementate numai funcțiile din setul 5. A doua consecință este că un nou manager de obiecte poate fi testat separat, înainte de integrarea sa în sistemul tranzațional. Soluția implementării obiectuale este utilă din ambele puncte de vedere.

Se prezintă, în final, concluziile autorului în legătură cu aspectele critice și caracteristicile fundamentale ale unei noi generații de sisteme de operare distribuite, așa după cum reies din parcurgerea acestei lucrări:

(1) Se apreciază că trebuie să se ofere facilități pentru integrarea în sisteme deschise, pentru multiprocesare, și distribuție masivă;

(2) Un element decisiv îl constituie optimizarea comunicațiilor interprocese și interprocesoare și distribuția, transparentă pentru utilizatori, a resurselor;

(3) O cerință importantă este portabilitatea pe diverse arhitecturi hardware, de la sisteme uniprocessor și multiprocessor (puternic sau slab cuplate) pînă la medii puternic distribuite, de la spații de adrese liniare, segmentate, pînă la arhitecturi cu memorie virtuală;

(4) Este necesar să se ofere suport pentru configurarea dinamică a serverelor de sistem dependente de hardware;

(5) Este necesară posibilitatea distribuirii serviciilor sistemului de operare ;

(6) Este necesară respectarea compatibilității cu interfețele standard ale sistemelor deschise;

(7) Trebuie să se ofere un mediu adecvat pentru dezvoltarea, depanarea, extinderea și integrarea altor aplicațiilor cit și a noilor componente ale sistemului de operare; în ultimul caz trebuie să fie posibilă testarea componentelor în mod de lucru utilizator urmat, ulterior, după depanare, de migrarea în mod sistem;

(8) Arhitectura bazată pe un micro-kernel minimal permite o construcție modulară a sistemului de operare, într-o manieră în care componentele de sistem pot fi dezvoltate și asamblate în diferite moduri respectînd platformele hardware actuale, pe de o parte și cerințele aplicațiilor pe de altă parte, făcînd posibilă cooperarea și interacționarea elementelor de sistem prin intermediul liniilor de comunicație de mare performanță;

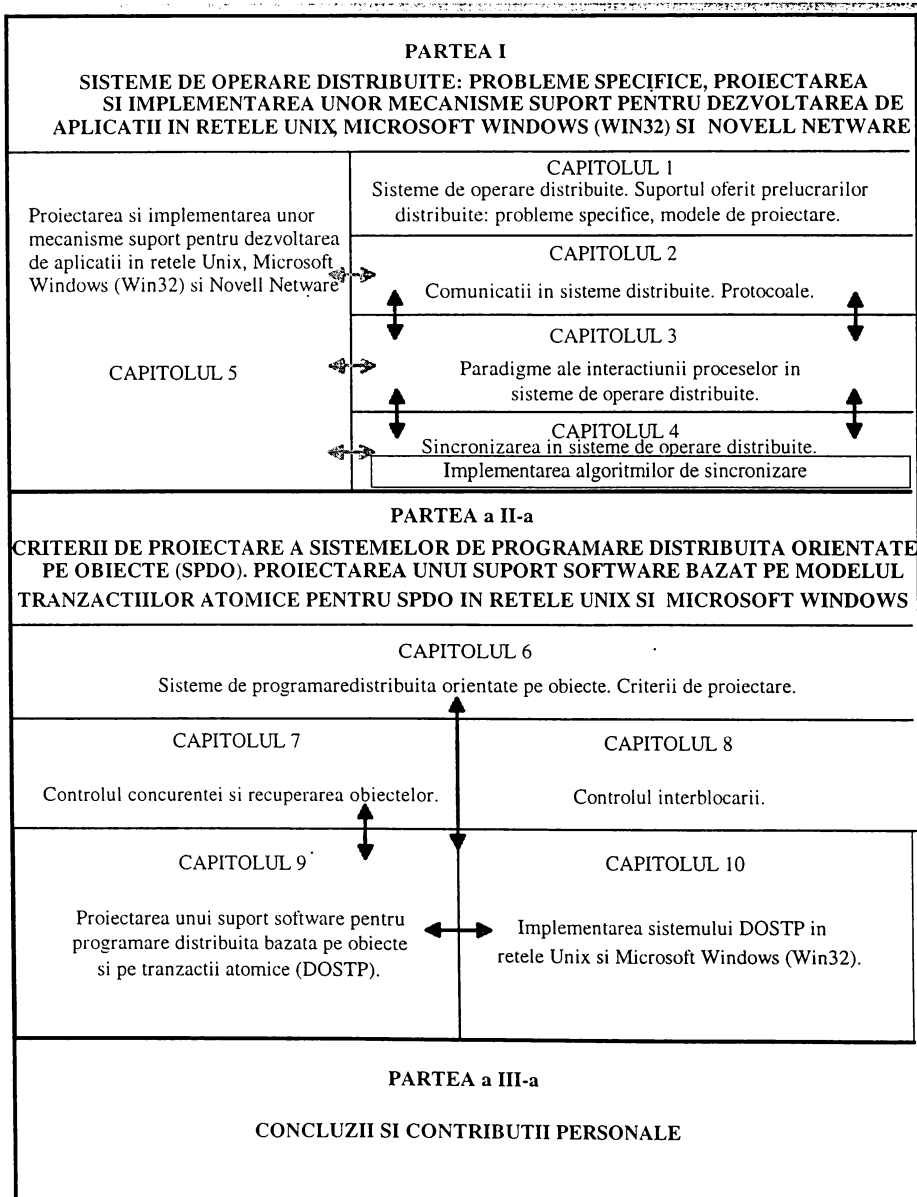
(9) O importanță mare o prezintă infrastructura oferită pentru lucrul cu fire, actualmente standardizat de POSIX. Se apreciază însă că este nevoie de un suport mai bun pentru lucrul cu fire în ceea ce privește limbajele de programare, depanarea programelor și modificarea bibliotecilor în biblioteci multifir.

(10) În arhitecturile bazate pe obiecte, pentru a oferi suport adecvat aplicațiilor, este necesară implementarea obiectelor partajate și persistente, dotate cu controlul concurenței și posibilități de recuperare. Aceste cerințe sunt realizate prin utilizarea conceptului de tranzație atomică, implementat prin protocolul *two-phase commit*, tehnicile de recuperare a obiectelor și mecanisme de control al concurenței (*optimiste sau pesimiste*).

(11) Conceptul de tranzație distribuită, caracterizat de proprietățile *ACID (atomicitate, consistență, integritate și durabilitate)* poate constitui un suport software de mare ajutor pentru aplicații;

(12) Sistemele distribuite cu facilități de toleranță la defecte necesită în plus față de serviciile obișnuite tip *datagram unicast*, *servicii de comunicații de grup*, fie la nivel de procese, fie la nivel de obiecte.

Interdependența problemelor abordate în această lucrare și plasarea subsistemelor software dezvoltate de autor în raport cu noțiunile teoretice abordate sunt prezentate în figura următoare.



11.2. Contribuții originale, valoare aplicativă și direcții de dezvoltare

Pornind de la dublul scop al acestei lucrări - acela de a prezenta aspecte teoretice ale proiectării sistemelor de operare distribuite și de a proiecta și implementa în rețele locale Unix, Microsoft Windows(Win32) și Novell NetWare subsisteme software - autorul a urmărit realizarea următoarelor obiective:

(A) selecția și sistematizarea unui mare volum de cunoștințe din domeniul sistemelor de operare distribuite, atât din lucrările congreselor de specialitate din ultimii 15 ani cât și din lucrările publicate în reviste de specialitate;

(B) elaborarea unor variante originale pentru: (a) proiectarea unui toolkit pentru comunicații de grup; (b) proiectarea unor mecanisme și algoritmi pentru rezolvarea problemelor de concurență la nivelul obiectelor și pentru recuperarea obiectelor într-un sistem de proiectare distribuită bazat pe obiecte și tranzații atomice, și găsirea unor soluții de conectare automată a managerilor de obiecte la acest sistem; (c) proiectarea unor mecanisme pentru controlul accesului concurrent la o bază de date.

(C) elaborarea unor variante originale de implementare în rețele locale a unor concepte precum contori de evenimente distribuiți, secvențiatori distribuiți, și implementarea în rețele locale Unix, Microsoft Windows(Win32) și Novell NetWare a mai multor algoritmi cunoscuți pentru sincronizarea proceselor în sisteme distribuite (algoritmi bazați pe marca de timp, algoritmi bazați pe un token circulant și algoritmi bazați pe transmisie broadcast)

(D) implementarea în rețele locale Unix, Microsoft Windows(Win32) și Novell NetWare a 3 subsisteme software: (1) pentru comunicații de grup; (2) pentru controlul accesului concurrent la o bază de date; (3) pentru proiectarea unui suport pentru un sistem de proiectare distribuită bazat pe obiecte și tranzații atomice.

Având în vedere scopul și obiectivele acestei lucrări se pot scoate în evidență următoarele contribuții originale ale acestei teze:

(1) Prezentarea comparativă a noțiunilor de sistem de operare distribuit și sistem de operare pentru rețele de calculatoare, prin prezentarea totodată a elementelor caracteristice fiecăruia.

(2) Expunerea unitară și sistematică a problemelor referitoare la gestiunea proceselor și gestiunea fișierelor în sisteme de operare distribuite.

(3) Abordarea și prezentarea legăturii sistem de operare - prelucrare distribuită din 2 puncte de vedere: interconectarea programelor în rețele locale, pe baza unei descrieri a aplicației (sistemul TDL, capitolul 1) și punerea la dispoziția utilizatorilor, într-un sistem de operare cunoscut (Unix) a unor API-uri pentru dezvoltarea de aplicații distribuite (DCE OSF).

(4) Prezentarea comparativă a 2 modele de proiectare a sistemelor de operare distribuite: *modelul proces* și *modelul obiect*, cu accentuarea submodelelor *obiecte active secvențiale* sau *obiecte multiactive* ultimul fiind folosit de mai multe ori, pe parcursul acestei lucrări, pentru implementarea unor subsisteme.

(5) Sistematizarea cunoștințelor legate de specificațiile ISO/OSI în raport cu comunicațiile în sisteme distribuite, alături de relieful rolului principal deținut de nivelul transport în furnizarea de servicii pentru sistemele de operare.

(6) Exemplificarea protocoalelor prezentate în capitolul 2, protocolul IP - accentuându-se importanța noii versiuni IPv6-, protocolul TCP, Novell IPX și a comunicării prin socluri în Unix și Microsoft Windows prin aplicații prezentate în anexele B și C.

(7) Sistematizarea și prezentarea unitară a aspectelor specifice ale suportului de comunicație pentru sisteme de operare distribuite, cu realizarea unei comparații între protocoalele cu caracter universal și cele cu caracter restrâns.

(8) Prezentarea *modelului client-server* sub multiple aspecte: acela al concurenței la nivelul serverelor, al protocoalelor specifice și al sistemelor client-server.

(9) Expunerea completă a comunicării prin mesaje în sisteme distribuite, cu exemplificarea utilizării primitivelor *receive* fără blocare și fără buferare (rezolvarea problemei buferării mesajelor la nivelul aplicației) în aplicația din anexa C4.

(10) Prezentarea comunicației prin conducte (*pipes*) la distanță, cu exemplificarea serverelor multifir în Microsoft Windows NT 4.0.

(11) Prezentarea unitară a *modelului RPC*, cu exemplificări pentru DCE OSF, Microsoft RPC sub multiple aspecte: transmiterea parametrilor și rezultatelor, legarea clienților la servere, concurență, semantica RPC în prezența defectelor, protocoale RPC, extensii ale mecanismului de bază, aspecte critice în proiectare.

(12) Expunerea unitară a algoritmilor pentru sincronizarea proceselor în sisteme distribuite, utilizând descrierea SR introdusă de G. Andrews.

(13) Prezentarea și implementarea în rețele locale Unix, Microsoft Windows (Win32) și Novell NetWare a mai multor algoritmi pentru excludere reciprocă atât din clasa algoritmilor bazați pe marca de timp

cît și din clasa algoritmilor bazați pe transmiterea unui token circulant. Implementările algoritmilor sunt prezentate în anexele B1-B6.

(14) Elaborarea unor variante originale pentru implementarea unor mecanisme pentru sincronizarea proceselor în sisteme distribuite, precum contori de eveniment distribuiți, secvențiatori distribuiți, secvențiator circulant.

(15) Implementarea în rețele locale Unix, Microsoft Windows (Win32) și Novell NetWare a unor mecanisme pentru sincronizare precum: semafoarele distribuite, contori de eveniment distribuiți, secvențiatori distribuiți, cu exemplificarea folosirii lor pentru obținerea excluderii reciproce.

(16) Elaborarea mai multor variante pentru controlul accesului concurrent la o bază de date în rețele locale de calculatoare și comparația acestor variante.

(17) Implementarea în rețele locale Unix, Microsoft Windows (Win32) și Novell NetWare a 4 soluții pentru controlul accesului concurrent la o bază de date în rețele locale: 3 soluții bazate pe semafoare și o soluție bazată pe un proces coordonator. Implementările sunt prezentate în anexa C.

(18) Elaborarea a 4 variante originale pentru implementarea comunicațiilor de grup: o variantă semidistribuită, o variantă cu transmisie broadcast, o variantă care folosește protocolul *two-phase commit* și o variantă cu transmisie în inel.

(19) Implementarea în rețele Unix a 3 variante pentru comunicațiile de grup: varianta semidistribuită (cap 5.2.2.), varianta cu transmisie broadcast și varianta cu transmisie în inel (anexa A).

(20) Elaborarea unui studiu teoretic privind proiectarea sistemelor de programare distribuită bazată pe obiecte și pe tranzacții atomice, care tratează problemele de concurență la nivelul obiectelor, recuperarea obiectelor și realizarea atomicității tranzacțiilor.

(21) Proiectarea unui suport software pentru programarea distribuită bazat pe obiecte și pe tranzacții atomice (DOSTP). Sunt prezentate detaliile de implementare ale fiecărei componente: serverele de comunicații, managerii de tranzacții, managerii de obiecte, serverele de nume, managerul central, procesele utilizator care realizează tranzacțiile.. Se prezintă, de asemenea, circulația mesajelor în sistem în toate momentele distincte care pot apărea în funcționarea sistemului: conectarea unui nod la sistem, inițierea unei tranzacții, implicarea unui obiect local sau de la distanță într-o tranzacție, încheierea unei tranzacții prin *commit* sau *abort*, recuperarea obiectelor după repunerea în funcțiune a unui nod coordonator sau neecordonator al unei tranzacții. Soluția aleasă pentru implementare se bazează, în principal, pe două observații: (1) dezvoltarea unui sistem de programare distribuită poate fi simplificată prin utilizarea obiectelor, care devin entități autonome ce servesc ca unități pentru sincronizare, protecție, recuperabilitate, securitate și mobilitate. (2) gruparea acțiunilor asupra obiectelor în tranzacții și extinderea mecanismului tranzacțiilor la tipuri abstracte devine un instrument puternic în proiectarea unui astfel de sistem.

(22) Elaborarea unei metode de proiectare obiectuală a managerilor de obiecte atît pentru conectarea lor automată la sistemul tranzacțional cît și pentru dezvoltarea și testarea lor autonomă.

(23) Implementarea sistemului DOSTP în rețele locale Unix și Microsoft Windows.

Aspectele prezentate în lucrare sunt rezultatele cercetărilor autorului din ultimii 7 ani, fiind valorificate în cadrul a 18 lucrări publicate la conferințe științifice internaționale și în reviste de specialitate dintre care se remarcă articolul "The Design and the Implementation of a DOSTP- A Distributed Object Oriented System Based On Transactions Processing", publicat în revista *Facta Universitatis, series: Electronics and Energetics*, vol 9, no 2, 1996, pp 149-172, Nis, Yugoslavia, articolul "Object Oriented System for Distributed Transactions during the Development of a Product within a group of Companies", publicat la *1-st International Conference on Engineering Design and Automation*, Bangkok, Thailand, Martie 18-21, 1997 și raportul de cercetare "Designing a Distributed Object Oriented System", publicat și susținut în programul TEMPUS prin două prelegeri la Technological Educational Institute of Piraeus, Faculty of engineering, Department of Computer Systems, Athens, Greece, mai 1996. Totodată, unele aspecte ale lucrării au putut fi valorificate într-un contract de cercetare științifică încheiat cu Ministerul Învățămîntului în anul 1994 (grant no 3006/14c/B9/1994) și în activitatea didactică desfășurată cu studenții secției de calculatoare, Facultatea de Automatică, Calculatoare și Electronică, Universitatea din Craiova.

Se remarcă caracterul aplicativ al rezultatelor acestei teze, pe viitor conturîndu-se noi perspective de valorificare, ținînd cont de interesul crescînd care se manifestă față de sistemele distribuite.

ANEXA A

Implementarea unui toolkit pentru comunicații de grup Managerul de grupuri pentru variantele distribuite

Serverul pentru tratarea cererilor(managerul de grupuri) inglobeaza functiile serverului central si ale celui
local din varianta semidistribuita prezentata la 5.2.2. si este identic pentru varianta cu transmisie broadcast
si pentru varianta cu transmisie in inel

```
const MAX_NAME_LEN = ... # lungimea maxima a numelor clientilor sau a grupurilor
const MAX_MES_LEN = ... # lungimea maxima a unui mesaj
const MAX_CLIENTS_IN_GROUP = ... # numarul maxim de clienti dintr-un grup
const MAX_MES_FOR_A_CLIENT = ... # numarul maxim de mesaje care pot fi buferate pentru un client
const MAX_CONC_GROUPS = ... # numarul maxim de grupuri la care poate apartine un client
const MAX_GROUPS = .. # numarul maxim de grupuri
const MAX_CLI = MAX_CLIENTS_IN_GROUP * MAX_GROUPS # numarul maxim de clienti
const MAX_PROTO = ... # lungimea maxima a unui sir care indica protocolul folosit pentru autorizare clienti
const MAX_STRING = ... # lungimea maxima a unei parole folosita pentru autorizarea clientilor
const HEADER_LEN = ... # lungimea header-ului mesajelor
const MAX_ALIASES_FOR_A_GROUP = ... # numarul maxim de pseudonime pentru un grup
```

```
type kind = enum ( # tipuri de mesaje care circula in sistem
    INIT_TGC, CLOSE_TGC, CREATE_GROUP, JOIN_GROUP, GET_ID_GROUP,
    LEAVE_GROUP, CREATE_ALIAS_GROUP, DESTROY_ALIAS_GROUP,
    SEND_TO_GROUP, RECEIVE_FROM_GROUP_NOW, GET_NO_OF_MESSAGES
    RECEIVE_FROM_GROUP_INFINITE, NO_MORE_RECEIVE_FROM_GROUP)
```

```
type error = enum ( # codificarea raspunsurilor la cererile trimise de clienti
    SUCCESS, MORE_MES, NO_WAIT_MESSAGES, ECLIENT_NOT_CONNECTED,
    ENO_AUTHORIZATION, ETOO_MANY_CLIENTS, ETOO_MANY_CONC_GROUPS,
    EFULL_GROUP, ECLIENT_ALREADY_CONNECTED,
    ECLIENT_NOT_IN_DEST_GROUP, EWRONG_SOURCE_GROUP,
    ECLIENT_NOT_IN_SOURCE_GROUP, EWRONG_DEST_GROUP,
    ENO_BUF, EGROUP_ALREADY_CREATED, ETOO_MANY_GROUPS,
    EWRONG_ALIAS, EMES_TO_LONG, EBADF)
```

```
type mes_state = enum (free, in_transit, busy) # starile unui slot in buferul de mesaje
type mes_to_group_manager = enum (ACK, REQ) # tipul unui mesaj transmis server comunicatie->manager
# de gestiune grupuri
```

```
type CLIENT_CONN = record of
    ( pid: int, # identificadorul unui client
      port: int, # identificadorul canalului rezervat dinamic unui client
      name[1:MAX_NAME_LEN]:char # numele clientului dat la conectare
    ) # structura care mentine identitatea clientilor
```

```
type MES_TO_GROUP = record of # structura care mentine informatii relative la un mesaj receptionat
    ( mes[1:MAX_MES_LEN]: char, # mesaj
      counter: int, # nr proceselor care partajeaza mesajul
      source_gid: int, # identificadorul grupului sursa
      dest_gid: int # identificadorul grupului destinatie
    )
```

```
type GROUP_OF_CLIENTS = record of # structura care mentine componenta unui grup
    (gid: int, # identificadorul de grup
      clients_no: int, # numarul de clienti din acest grup
      pid_tab[1:MAX_CLIENTS_IN_GROUP]: int # identicatori clienti
    )
```

```
type MES_TO_CLIENT = record of # buffer de mesaje pentru un client
    ( pid: int, # identificadorul clientului
```

```

mes_no[1:MAX_MES_FOR_A_CLIENT]: int # indica numarul intrarii in buferul de mesaje
pmes[1:MAX_MES_FOR_A_CLIENT]: pointer to MES_TO_GROUPS # indica mesajul
pmes_state[MAX_MES_FOR_A_CLIENT]: mes_state) # indica daca slotul este ocupat
)
type WCLIENTS_FOR_GROUP = record of # structura care mentine identificatorii proceselor blocate,
# in asteptare de mesaje de la un grup
(gid: int, # identificatorul de grup
wclients_no: int, # numarul de procese in asteptare de mesaje de la grupul gid
wpid_tab[1:MAX_CLIENTS_IN_GROUP]: int, # identificatorii proceselor in asteptare
)
type CLIENTS_TO_GROUPS = record of # structura care mentine grupurile la care apartine un client
(pid :int # identificatorul clientului
gid_no: int # numarul de grupuri la care apartine un client
gid_tab[1:MAX_CONC_GROUPS]: int # identificatorii grupurilor
)
type MES_BUF = record of # structura unei intrari in bufferul de mesaje
(state: mes_state # starea unei intrari: libera, ocupata, in tranzit
message: MES_TO_GROUP) # mesajul propriu zis transmis
type GROUP_ALIASES = record of # structura care mentine pseudonimele grupului gid
(gid: int,
alias_no: int, # numarul de pseudonime acordate grupului gid
alias_tab[1:MAX_ALIASES_FOR_A_GROUP][1:MAX_NAME_LEN]: char # pseudonime
chan client_chan[1..MAX_CL] ( # canal pentru comunicatie server local → client
error_code: int , # codul erorii sau zero pentru succes
dest_group_id: int, # identificatorul grupului destinatie
source_group_id: int, # identificatorul grupului sursa
len: int, # lungimea partii variabile a mesajului
mes[1:len]: char ) # mesajul
chan local_chan[1..n] ( → # canal pentru comunicatie client → server local
knd: kind, # tipul cererii
client_id: int # identificatorul clientului de la care s-a trimis cererea
dest_group_id: int # identificatorul grupului de la care s-a trimis cererea
source_group_id: int # identificatorul grupului sursa care a trimis mesajul sau cererea
group_name[1:MAX_LEN_GROUP]: char # nume de grup
len: int, # lungimea partii variabile a mesajului
mes[1:len]: char ) # mesajul
chan central_chan[1..n]( # canal pentru comunicatie server local → server central
knd: kind, # tipul cererii
client_id: int # identificatorul clientului de la care s-a trimis cererea
dest_group_id: int # identificatorul grupului de la care s-a trimis cererea
source_group_id: int # identificatorul grupului sursa care a trimis mesajul sau cererea
group_name[1:MAX_LEN_GROUP]: char # nume de grup
len: int, # lungimea partii variabile a mesajului
mes[1:len]: char ) # mesajul
chan group_man_chan[1..n]( # canal pentru comunicatie server comunicatii → manager grupuri
type: mes_to_group_manager , # ACK sau REQ
err_code: int, # codul de eroare
knd: kind, # tipul cererii
client_id: int # identificatorul clientului de la care s-a trimis cererea
dest_group_id: int # identificatorul grupului de la care s-a trimis cererea
source_group_id: int # identificatorul grupului sursa care a trimis mesajul sau cererea
group_name[1:MAX_LEN_GROUP]: char # nume de grup
len: int, # lungimea partii variabile a mesajului
mes[1:len]: char ) # mesajul
chan temp_chan ( # canal pentru conectarea unui client la serverul local
error_code: int, # codul erorii
client_id: int) # identificatorul clientului

```

```
# informatii partajate cu managerul de grupuri GM intr-o zona de memorie partajata
var no_groups: int # numarul de grupuri
var group_aliases_tab[1:MAX_GROUPS]:GROUP_ALIASES # mentine denumirile grupurilor
var group_of_clients_tab[1:MAX_GROUPS]:GROUP_OF_CLIENTS # mentine inf. despre componenta gr-lor
var mes_to_clients_tab: [1: MAX_CLI]:MES_TO_CLIENT # mentine inf. despre mes. in asteptare pt.clienti
var messages_buf[1: MAX_MESSAGES]: MES_BUF # bufferul de mesaje

# alte structuri de date
var client_conn_tab[MAX_CLI]: CLIENT_CONN # mentine clientii activi si conexiunile lor
var wclients_for_group_tab[1:MAX_GROUP]:WCLIENT_FOR_GROUP# mentine identitatea clientilor blocati
var client_to_groups_tab[1: MAX_CLI]: CLIENT_TO_GROUPS # mentine grupurile tuturor clientilor
var knd:int, error_code:int,client_id:int, dest_group_id:int,source_group_id:int # folosite pentru citire mesaje
var mes_len: int, mes[1:MAX_MES]: char # memoreaza lungimea si partea variabila a unui mesaj
var name[1: MAX_NAME_LEN]:char # folosita pentru numele clientilor si grupurilor (alias-uri)
var auto_proto[1:MAX_PROTO]:char # folosita pentru citirea protocolului utilizat pentru autorizarea clientilor
var auto_string[1:MAX_STRING]: char # folosita pentru citirea unei parole
var message[1:MAX_MES]:char, source_gid: int # variabile de lucru
var pmessage: pointer to MES_TO_GROUP #variabila de lucru
var excl: semaphore := 1 # folosit pentru accesul la memoria partajata
var start_GM: semaphore:=0 # pentru sincronizarea server comunicati - manager de grupuri
```

initializari

```
Down(start_GM)
```

```
do true →
```

```
if not empty(local_chan[i]) →
```

```
receive local_chan[i] (knd, client_id, dest_group_id, source_group_id, name, mes_len , mes)
```

```
if knd = INIT_TGC and client_id ∉ clients_conn_tab →
```

```
if knd = INIT_TGC →
```

```
determina protocolul de autorizare client, auto_proto, si parola, auto_string din mes
```

```
valideaza autorizarea clientului si pozitioneaza error_code
```

```
if error_code ≠ SUCCESS →
```

```
send temp_chan(ENO_AUTHORIZATION,0)
```

```
[] error_code = SUCCESS →
```

```
stabileste un identificator pentru client in client_id,
```

```
ca fiind numarul primei intrari libere din client_conn_tab
```

```
if client_id > MAX_CLIENTS → send temp_chan(ETOO_MANY_CLIENTS, 0)
```

```
[] client_id ≤ MAX_CLIENTS →
```

```
creaza o intrare clients_conn_tab[client_id]
```

```
Down(excl);creaza o intrare mes_to_client_tab[client_id]; Up(excl)
```

```
send temp_chan (SUCCESS, client_id)
```

```
fi
```

```
fi
```

```
[] knd = JOIN_GROUP →
```

```
if client_to_groups_tab[client_id].gid_no > MAX_CONC_GROUPS →
```

```
send client_chan[client_id] (ETOO_MANY_CONC_GROUPS, dest_group_id, 0, 0, NULL)
```

```
[]client_to_groups_tab[client_id].gid_no ≤ MAX_CONC_GROUPS
```

```
if dest_group_id ∈ group_of_clients_tab →
```

```
if client_id ∉ group_of_clients_tab[dest_group_id] →
```

```
if group_of_clients_tab[dest_group_id].clients_no > MAX_CLIENTS_IN_GROUP →
```

```
send client_chan[client_id](EFULL_GROUP, dest_group_id, 0, 0, NULL)
```

```
[] group_of_clients_tab[dest_group_id].clients_no ≤ MAX_CLIENTS_IN_GROUP →
```

```
Down(excl);insereaza client_id in group_of_clients_tab[dest_group_id];Up(excl)
```

```
insereaza dest_group_id in clients_to_group_tab
```

```
send client_chan[client_id](SUCCESS, dest_group_id, 0, 0, NULL)
```

```
fi
```

```
[] client_id ∈ group_of_clients_tab[dest_group_id] →
```

```
send client_chan[client_id](ECLIENT_ALREADY_CONNECTED, dest_group_id, 0, 0, NULL)
```

```
fi
```

```
    [] dest_group_id ∈ group_of_clients_tab →
        send client_chan[client_id](EWRONG_DEST_GROUP, dest_group_id, 0, 0, NULL)
    fi
fi
[] knd = LEAVE_GROUP →
    if dest_group_id ∈ group_of_clients_tab →
        send client_chan[client_id](WRONG_DEST_GROUP, dest_group_id, 0, 0, NULL)
    [] dest_group_id ∈ group_of_clients_tab →
        if client_id ∈ group_of_clients_tab[dest_group_id].pid_tab →
            Down(excl); sterge client_id din group_of_clients_tab[dest_group_id].pid_tab;
            fa mesajele i st mes_to_clients_tab[client_id].pmes[i]^dest_gid = dest_group_id →
                pmessage := mes_to_clients_tab[client_id].pmes[i]
                mes_to_clients_tab[client_id].pmes[i] = NULL;
                pmessage^.counter := pmessage^.counter - 1 # un mesaj in asteptare mai putin
            if pmessage^.counter = 0 →
                deazaloca structura de la pmessage si intrarea
                    mes_to_clients_tab[client_id].mes_no[i] din message_buf devine free
            fi
        af
        Up(excl)
        send client_chan[client_id](SUCCESS, dest_group_id, 0, 0, NULL)
        if group_of_clients_tab[dest_group_id].clients_no = 0 →
            send central_chan[i](LEAVE_GROUP, client_id, dest_group_id, 0, NULL, 0, NULL)
        fi
    [] client_id ∈ group_of_clients_tab[dest_group_id].pid_tab →
        send client_chan[client_id](ECLIENT_NOT_IN_DEST_GROUP, dest_group_id, 0, 0, NULL)
    fi
fi
[] knd = CREATE_GROUP →
    if exista un grup, gid pentru care name ∈ group_aliases_tab[gid]→
        send client_chan[client_id](EGROUP_ALREADY_CREATED, 0, 0, 0, NULL)
    [] nu exista un grup, gid pentru care name ∈ group_aliases_tab[gid]→
        if nu exista o intrare o intrare libera in group_of_clients_tab →
            send client_chan[client_id](ETOO_MANY_GROUPS, 0, 0, 0, NULL)
        [] send central_chan[i](knd, client_id, 0, 0, name, 0, NULL)
    fi
fi
[] knd = GET_ID_GROUP →
    if nu exista un grup, gid pentru care name ∈ group_aliases_tab[gid]→
        send client_chan[client_id] (EWRONG_ALIAS, 0, 0, 0, NULL)
    [] exista un grup, gid pentru care name ∈ group_aliases_tab[gid]→
        send client_chan[client_id] (SUCCESS, gid, 0, 0, NULL)
    fi
[] knd = CREATE_ALIAS_GROUP →
    if dest_group_id ∈ group_of_clients_tab →
        send client_chan[client_id] (EWRONG_DEST_GROUP, dest_group_id, 0, 0, NULL)
    [] dest_group_id ∈ group_of_clients_tab →
        if exista un grup, gid pentru care name ∈ group_aliases_tab[gid]→
            send remote_chan[i] (CREATE_ALIAS_GROUP, EWRONG_ALIAS, client_id, 0,0,0,NULL)
        [] nu exista un grup, gid pentru care name ∈ group_aliases_tab[gid]→
            if group_aliases_tab[dest_group_id].alias_no > MAX_ALIASES_FOR_A_GROUP →
                send remote_chan[i] (knd, ETOO_MANY_ALIASES, client_id, dest_group_id, 0,0,0,NULL)
            [] group_aliases_tab[dest_group_id].alias_no ≤ MAX_ALIASES_FOR_A_GROUP →
                send central_chan[i] (knd, client_id, gid, 0, name, 0, NULL)
            fi
        fi
    fi
fi
```

```
[] knd = DESTROY_ALIAS_GROUP →
    if dest_group_id ∉ group_of_clients_tab →
        send client_chan[client_id](EWRONG_DEST_GROUP, dest_group_id, 0, 0, NULL)
    [] dest_group_id ∈ group_of_clients_tab →
        if name ∉ group_aliases_tab[dest_group_id] →
            send client_chan[client_id](EWRONG_ALIAS, dest_group_id, 0, 0, NULL)
        [] name ∈ group_aliases_tab[dest_group_id] →
            send central_chan[i](knd, client_id, dest_group_id, 0, name, 0, NULL)
        fi
    fi
[] knd = SEND_TO_GROUP →
    if source_group_id ∉ group_of_clients_tab →
        send client_chan[client_id](EWRONG_SOURCE_GROUP, # se incearca o emisie dintr-un grup
            dest_group_id, source_group_id, 0, NULL) # sursa inexistent
    [] source_group_id ∈ group_of_clients_tab →
        if client_id ∉ group_of_clients_tab[source_group_id].pid_tab →
            send client_chan[client_id](ECLIENT_NOT_IN_SOURCE_GROUP, # se incearca o emisie
                dest_group_id, 0, 0, NULL) # in contul unui client neconectat la sursa
        [] client_id ∈ group_of_clients_tab[source_group_id].pid_tab →
            if dest_group_id ∉ group_of_clients_tab →
                send client_chan[client_id](EWRONG_DEST_GROP, # se incearca o emisie la un grup destinatie
                    dest_group_id, source_group_id, 0, NULL) # inexistent
            [] formeaza la message un mesaj din adresa proprie IP, client_conn_tab[client_id].name, mes
                send central_chan[i](knd, client_id, dest_group_id, source_group_id, NULL,
                    mes_len + HEADER_LEN, message)
            fi
        fi
    fi
[] knd = RECEIVE_FROM_GROUP_INFINITE →
    if (source_group_id > MAX_GROUPS) or source_group_id ∉ group_of_clients_tab →
        send client_chan[client_id](EWRONG_SOURCE_GROUP, # se incearca receptia de la un grup
            dest_group_id, source_gid, mes_len, message) # inexistent
    [] (source_group_id ∈ MAX_GROUPS) →
        if dest_group_id ∉ group_of_clients_tab →
            send client_chan[client_id](EWRONG_DEST_GROUP, # se incearca o receptie in contul unui grup
                dest_group_id, source_group_id, 0, NULL) # care nu exista
        [] dest_group_id ∈ group_of_clients_tab →
            if client_id ∉ group_of_clients_tab[dest_group_id].pid_tab →
                send client_chan[client_id](ECLIENT_NOT_IN_DEST_GROUP, # se incearca o receptie
                    dest_group_id, 0, 0, NULL) # in contul unui client neconectat la grupul destinatie
            [] client_id ∈ group_of_clients_tab[dest_group_id].pid_tab →
                if exista mesaje i in mes_to_client_tab[client_id] and
                    mes_to_clients_tab[client_id].pmes[i]^dest_gid = dest_group_id and
                    (source_group_id = ANY or
                    (mes_to_clients_tab[client_id].pmes[i]^source_gid = source_group_id)) →
                    Down(excl);
                    extrage primul mesaj i in message si calculeaza lungimea sa in mes_len
                    pmessage := mes_to_clients_tab[client_id].pmes[i]
                    mes_to_clients_tab[client_id].pmes[i] = NULL;
                    source_gid = pmessage^.source_gid
                    pmessage^.counter := pmessage^.counter - 1 # un mesaj in asteptare mai putin
                    if pmessage^.counter = 0 → dezaloca structura de la pmessage si intrarea
                        mes_to_clients_tab[client_id].mes_no[i] din message_buf devine free
                    fi
                fi
                Up(excl);
                if (mai exista mesaje i in mes_to_client_tab[client_id]) and
                    mes_to_clients_tab[client_id].pmes[i]^dest_gid = dest_group_id and
                    ((source_group_id = ANY) or
```

```

        ( mes_to_clients_tab[client_id].pmes[i]^source_gid = source_group_id )) →
    send client_chan[client_id](MORE_MES, dest_group_id, source_gid,
        mes_len, message) # se receptioneaza mesajul si o indicatie ca mai exista mesaje
[] nu mai exista mesaje in asteptare de la source_group_id catre dest_group_id →
    send client_chan[client_id](SUCCESS, dest_group_id, source_gid, mes_len,message)
    #se receptioneaza mesajul si o indicatie ca nu mai exista mesaje
fi
[]nu exista mesaje in asteptare de la source_group_id catre dest_group_id pentru client_id →
    if wclients_for_groups[source_group_id].wclients_no>MAX_CLIENTS_IN_GROUP →
        send client_chan[client_id](ENO_BUF, dest_group_id, source_gid, mes_len,message)
    [] wclients_for_groups[source_group_id].wclients_no ≤ MAX_CLIENTS_IN_GROUP →
        insereaza client_id in wclients_for_groups[source_group_id].wpid.tab
        si actualizeaza wclients_for_groups[source_group_id].wclients_no
    fi
fi
fi
[] knid = NO_MORE_RECEIVE_FROM_GROUP →
    sterge client_id din wclients_for_groups[source_group_id].wpid.tab
    si actualizeaza wclients_for_groups[source_group_id].wclients_no
[] knid = RECEIVE_FROM_GROUP_NOW →
    if source_group_id > MAX_GROUPS or source_group_id ∉ MAX_GROUPS →
        send client_chan[client_id](EWRONG_SOURCE_GROUP, # se incearca receptia de la un grup
            dest_group_id, source_gid, mes_len,message) # inexistent
    [] (source_group_id ∈ MAX_GROUPS) →
        if dest_group_id ∉ group_of_clients_tab →
            send client_chan[client_id](EWRONG_DEST_GROUP #se incearca o receptie in contul unui grup
                dest_group_id, source_group_id,0, NULL)# care nu exista
        [] dest_group_id ∈ group_of_clients_tab →
            if client_id ∉ group_of_clients_tab[dest_group_id].pid_tab →
                send client_chan[client_id](ECLIENT_NOT_IN_DEST_GROUP, # se incearca o receptie in
                    dest_group_id, 0, 0, NULL) # in contul unui client neconectat la grupul destinatie
            [] client_id ∈ group_of_clients_tab[dest_group_id].pid_tab →
                if exista mesaje i in mes_to_client_tab[client_id] and
                    mes_to_clients_tab[client_id].pmes[i]^dest_gid = dest_group_id and
                    (source_group_id = ANY or
                    ( mes_to_clients_tab[client_id].pmes[i]^source_gid = source_group_id )) →
                    extrage primul mesaj i in message si calculeaza lungimea sa in mes_len
                    Down(excl);
                    pmessage := mes_to_clients_tab[client_id].pmes[i]
                    mes_to_clients_tab[client_id].pmes[i] = NULL;
                    source_gid = pmessage^.source_gid
                    pmessage^.counter :=pmessage^.counter - 1 # un mesaj in asteptare mai putin
                    if pmessage^.counter = 0 → deazaloca structura de la pmessage si intrarea
                        mes_to_clients_tab[client_id].mes_no[i] din message_buf devine free
                fi
                Up(excl);
                if (mai exista mesaje i in mes_to_client_tab[client_id]) and
                    mes_to_clients_tab[client_id].pmes[i]^dest_gid = dest_group_id and
                    ((source_group_id = ANY) or
                    ( mes_to_clients_tab[client_id].pmes[i]^source_gid = source_group_id )) →
                    send client_chan[client_id](MORE_MES, dest_group_id, source_gid,
                        mes_len, message) # se receptioneaza mesajul si o indicatie ca mai exista mesaje
                [] nu mai exista mesaje in asteptare de la source_group_id catre dest_group_id →
                    send client_chan[client_id](SUCCESS, dest_group_id, source_gid, mes_len,message)
                    #se receptioneaza mesajul si o indicatie ca nu mai exista mesaje
                fi
            [] nu exista mesaje in asteptare de la source_group_id catre dest_group_id →

```

```

        send client_chan[client_id](NO_WAIT_MESSAGES, dest_group_id, source_gid,
                                0,NULL) #nu exista mesaje pentru receptie
    fi
fi
fi
fi
[] kno = GET_NO_OF_MESSAGES →
if source_group_id > MAX_GROUPS or source_group_id ∉ MAX_GROUPS →
    send client_chan[client_id](EWRONG_SOURCE_GROUP, # se incearca receptia de la un grup
                                dest_group_id, source_gid, mes_len,message) # inexistent
[] (source_group_id ∈ MAX_GROUPS) →
    if dest_group_id ∉ group_of_clients_tab →
        send client_chan[client_id](EWRONG_DEST_GROUP #se incearca o receptie in contul unui grup
                                dest_group_id, source_group_id,0, NULL)# care nu exista
    [] dest_group_id ∈ group_of_clients_tab →
        if client_id ∉ group_of_clients_tab[dest_group_id].pid_tab →
            send client_chan[client_id](ECLIENT_NOT_IN_DEST_GROUP, # se incearca o receptie in
                                dest_group_id, 0, 0, NULL) # in contul unui client neconectat la grupul destinatie
        [] client_id ∈ group_of_clients_tab[dest_group_id].pid_tab →
            contorizeaza in k mesajele i cu
                mes_to_client[client_id].pmes[i]^ .source_gid = source_group_id si
                mes_to_client[client_id].pmes[i]^ .dest_gid = dest_group_id
            send client_chan[client_id](SUCCESS, dest_group_id, source_group_id, k, NULL)
        fi
    fi
fi
[] kno = CLOSE_TGC →
sterge intrarea client_id din client_conn_tab[client_id]
Down(excl);
fa toate grupurile la care este conecat client_id, client_to_groups[client_id].gid.tab[i] →
    source_gid := client_to_groups[client_id].gid.tab[i]
    sterge client_id din group_of_clients_tab[source_gid].gid_tab[i]
    fa mesajele i st mes_to_clients_tab[client_id].pmes[i]^ .dest_gid =source_gid →
        pmessage := mes_to_clients_tab[client_id].pmes[i]
        mes_to_clients_tab[client_id].pmes[i] = NULL;
        pmessage^.counter :=pmessage^.counter - 1 # un mesaj in asteptare mai putin
        if pmessage^.counter = 0 → deazaloca structura de la pmessage si intrarea
            mes_to_clients_tab[client_id].mes_no[i] din message_buf devine free
        fi
    fa
af
Up(excl);
fi
fi
if not empty(group_man_chan[i]) →
    receive group_man_chan [i](type, error_code, kno,
                                client_id, dest_group_id, source_group_id, mes_len, message)

    if kno = CREATE_GROUP →
        if error_code = SUCCESS →
            Down(excl); adauga o intrare in group_of_clients[dest_group_id]; Up(excl)
        fi
        if type = ACK →
            Down(excl)
            send client_chan[client_id](error_code, dest_group_id,0,0, NULL)
            insereaza dest_group_id in clients_to_group_tab[client_id]
            insereaza client_id in intrarea group_of_clients[dest_group_id]
            Up(excl)
        fi
    fi
fi

```



```
fi
[] kncl = CREATE_ALIAS_GROUP
  if error_code = SUCCESS → adauga nume in group_aliases_tab[dest_group_tab] fi
  if type = ACK → send client_chan[client_id](error_code, dest_group_id,0,0, NULL) fi
[] kncl = DESTROY_ALIAS_GROUP
  if error_code = SUCCESS → sterge nume din group_aliases_tab[dest_group_tab] fi
  if type = ACK → send client_chan[client_id](error_code, dest_group_id,0,0, NULL) fi
[] kncl = LEAVE_GROUP →
  if error_code = SUCCESS →
    Down(excl);sterge intrarea dest_group_id din group_of_clients_tab ; Up(excl);
  fi
[] kncl = SEND_TO_GROUP →
  if type = ACK → send client_chan[client_id](error_code, dest_group_id,0,0, NULL)
  if error_code= SUCCESS →
    if exista client_id astfel incit (client_id ∈ dest_group_id) and
      (client_id ∉ wclients_for_group[source_group_id].wpid_tab) →
      Down(excl);
      se alocă o zona de memorie pentru o structura MES_TO_GROUP, in intrarea j din message_buf
      message_buf[j].state = busy; message_buf[j].message.mes := mes;
      message_buf[j].message.source_gid := source_group_id;
      message_buf[j].message.dest_gid := dest_group_id;
      se contorizeaza in k toti membrii grupului dest_group care nu sunt in
      wclients_for_group[source_group_id].wpid_tab
      message_buf[j].message.counter := k
      Up(excl)
    fi
    Down(excl);
    fa client_id ∈ dest_group_id st client_id ∉ wclients_for_group[dest_group_id].wpid_tab
    insereaza message_buf[j].message in
      mes_to_client_tab[client_id].pmes si j in mes_to_client_tab[client_id].mes_no
    af
    Up(excl);
    fa (client_id ∈ dest_group_id) st (client_id ∈ wclients_for_group[source_group_id].wpid_tab)
    send client_chan[client_id](SUCCESS, dest_group_id, source_group_id, mes_len,mes)
      #se receptioneaza mesajul si o indicatie ca nu mai exista mesaje
    sterge client_id din wclients_for_groups[source_group_id].wpid.tab
    af
  fi
fi
fi
od
```

Anexa B

Sincronizare in sisteme distribuite

B1. Exemple de programe client care, pentru realizarea excluderii reciproce în rețele locale, folosesc un server care implementează unul din algoritmiți prezentați în capitolul 4, la 4.3.2 și 4.3.3

```

.....
Program client XWindows
Foloseste functiile puse la dispozitie pentru un client
de catre fiecare implementare a algoritmiților :

int ExclAccessInitialize(void)
void ExclAccessOpen(void)
void ExclAccessClose(void)
void ExclAccessEnd();

..... */
#include <stdio.h>
#include <X11/Intrinsic.h>
#include <X11/StringDefs.h>
#include <X11/Shell.h>
#include <X11/Xaw/List.h>
#include <X11/Xaw/Command.h>
#include <X11/Xaw/Pane.h>
#include <X11/Xaw/Box.h>
#include <X11/Xaw/Dialog.h>
#include <X11/Xaw/AsciiText.h>
#include <string.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <arpa/inet.h>
#include <netdb.h>
#include "exclu.h"
#define MAX_MES_LENGTH 100

Widget shell,box,quit,info,pane,req,rel,label;

void Quit(Widget w,XtPointer client_data,XtPointer call_data);
void Info(Widget w,XtPointer client_data,XtPointer call_data);
void Request(Widget w,XtPointer client_data,XtPointer call_data);
void Release(Widget w,XtPointer client_data,XtPointer call_data);

main(argc,argv)
int argc;
char **argv;
{
    XtAppContext app_cont;
    shell = XtVaAppInitialize(&app_cont, "XBox",
        NULL,0,&argc,argv, NULL,NULL);
    pane = XtVaCreateManagedWidget ("pane",
        paneWidgetClass, shell, NULL);
    box = XtVaCreateManagedWidget ("box",
        boxWidgetClass, pane, NULL);
    info = XtVaCreateManagedWidget("info",
        commandWidgetClass, box, NULL);
    req= XtVaCreateManagedWidget("request",
        commandWidgetClass, box, NULL);
    rel = XtVaCreateManagedWidget("release",
        commandWidgetClass, box, NULL);
    quit = XtVaCreateManagedWidget("quit",
        commandWidgetClass, box, NULL);

    label = XtVaCreateManagedWidget("label",
        labelWidgetClass, pane,
        XtNlabel,"n Mutual Exclusion",
        XtNwidth,400, NULL);

    XtAddCallback(quit,XtNcallback,Quit,0);
    XtAddCallback(info,XtNcallback,Info,0);
    XtAddCallback(req,XtNcallback,Request,0);
    XtAddCallback(rel,XtNcallback,Release,0);
    if (ExclAccessInitialize())
        exit(1);
    XtRealizeWidget(shell);
    XtAppMainLoop(app_cont);
}

void Info(w,client_data,call_data)
/* apelata la actionarea butonului info */
/* afiseaza adresa nodului */
Widget w;
XtPointer client_data,call_data;
{
    char msg[MAX_MES_LENGTH];
    char *adr;
    if ((adr = GetServerAddress()) != NULL){
        sprintf(msg,"Host address= %s\n", adr);
        XtVaSetValues(label,XtNlabel,msg,NULL);
    }
    else {
        sprintf(msg," Unable to find the address\n");
        XtVaSetValues(label,XtNlabel,msg,NULL);
    }
}

void Request(w,client_data,call_data)
/* apelata la actionarea butonului request */
/* realizeaza intrarea intr-o sectiune critica */
Widget w;
XtPointer client_data,call_data;
{
    char msg[MAX_MES_LENGTH];
    strcpy(msg," Waiting ...");
    XtVaSetValues(label,XtNlabel,msg,NULL);
    ExclAccessOpen();
    strcpy(msg," In critical section ... ");
    XtVaSetValues(label,XtNlabel,msg,NULL);
}

void Release(w,client_data,call_data)
/* apelata la actionarea butonului release */
/* realizeaza iesirea dintr-o sectiune critica */
Widget w;
XtPointer client_data,call_data;
{
    char msg[MAX_MES_LENGTH];
    ExclAccessClose();
    strcpy(msg," Done request...");
    XtVaSetValues(label,XtNlabel,msg,NULL);
}

void Quit(w,client_data,call_data)
/* apelata la actionarea butonului quit */
/* asigura deconectarea de la server */
Widget w;
XtPointer client_data,call_data;
{
    ExclAccessEnd();
    exit(0);
}

```

Anexa B: Sincronizare in sisteme distribuite

```

/* .....
Acelasi client dar in varianta Microsoft Windows
..... */

#include <stdio.h>
#include <stdlib.h>
#include <winsock.h>
#include <windows.h>
#include "resource.h"
#include "generic.h"
#include "exclu.h"
#define MAX_STR 100

LRESULT CALLBACK WndProc( HWND hWnd, UINT uMsg,
WPARAM wParam, LPARAM lParam );
HINSTANCE hInst;
LPCTSTR lpszAppName = "MyApp";
LPCTSTR lpszTitle = "Mutual Exclusion";

int WINAPI WinMain( HINSTANCE hInstance,
HINSTANCE hPrevInstance, LPTSTR lpCmdLine, int
nCmdShow)
{
    MSG msg;
    HWND hWnd;
    WNDCLASS wc;
    WSADATA wSData;
    UNREFERENCED_PARAMETER( lpCmdLine );
    UNREFERENCED_PARAMETER( hPrevInstance );
    if(GetVersion() & 0x80000000 && (GetVersion() & 0xFF) == 3){
        MessageBox( NULL,
            "This application cannot run on Windows 3.1.\n"
            "This application will now terminate.",
            "WinApp",
            MB_OK | MB_ICONSTOP | MB_SETFOREGROUND );
        return( 1 );
    }
    hInst = hInstance;
    wc.style = CS_HREDRAW | CS_VREDRAW;
    wc.lpszWndProc = (WNDPROC)WndProc;
    wc.cbClsExtra = 0;
    wc.cbWndExtra = 0;
    wc.hInstance = hInstance;
    wc.hIcon = LoadIcon( hInstance, lpszAppName );
    wc.hCursor = LoadCursor(NULL, IDC_ARROW);
    wc.hbrBackground = (HBRUSH)(COLOR_WINDOW+1);
    wc.lpszMenuName = MAKEINTRESOURCE(IDR_MENU1);
    wc.lpszClassName = lpszAppName;
    RegisterClass( &wc );
    hWnd = CreateWindow( lpszAppName,
        lpszTitle, WS_OVERLAPPEDWINDOW,
        CW_USEDEFAULT, 0, CW_USEDEFAULT, 0,
        NULL, NULL, hInstance, NULL );
    if ( !hWnd )
        return( FALSE );
    ShowWindow( hWnd, nCmdShow );
    WSASStartup(MAKEWORD(1,1), &wSData );
    if (ExclAccessInitialize()){
        MessageBox( NULL,
            "This application will now terminate .", "WinApp",
            MB_OK | MB_ICONSTOP | MB_SETFOREGROUND );
        exit(1);
    }
    while( GetMessage( &msg, NULL, 0, 0 ) )
        DispatchMessage( &msg );
    return( msg.wParam );
}

LRESULT CALLBACK WndProc(HWND hWnd, UINT uMsg,
WPARAM wParam, LPARAM lParam)
{
    HDC hdc;
    char str[MAX_STR], *adr;
    switch (uMsg)
    {
        case WM_CREATE:
            return DefWindowProc( hWnd, uMsg, wParam, lParam );
            break;
        case WM_DESTROY:
            PostQuitMessage( 0 );
            break;
        case WM_COMMAND:
            switch ( LOWORD( wParam ) )
            {
                case ID_INFO:
                    SetWindowText( hWnd, "Server Address" );
                    if ((adr=GetMyAddress()) != NULL){
                        hdc = GetDC(hWnd);
                        TextOut(hdc,1,1,"",20);
                        sprintf(str,"%s",adr);
                        TextOut(hdc,1,1,str,strlen(str));
                        ReleaseDC(hWnd,hdc);
                    }
                    break;
                case ID_REQUEST:
                    SetWindowText( hWnd, "Waiting ... " );
                    ExclAccessOpen();
                    SetWindowText( hWnd, "In critical section " );
                    break;
                case ID_RELEASE:
                    ExclAccessClose();
                    SetWindowText( hWnd, "Done!! " );
                    break;
                case ID_EXIT:
                    DestroyWindow( hWnd );
                    break;
            }
            break;
        default:
            return DefWindowProc( hWnd, uMsg, wParam, lParam );
    }
    return 0;
}

.....
Funciile puse la dispozitia unui program client de catre
toti algoritmi care realizeaza excluderea reciproca, prezenati
in capitolul 4, la 4.3.2 si 4.3.3 -versiunea Unix
..... */

#include <stdio.h>
#include <string.h>
#include <signal.h>
#include <unistd.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <arpa/inet.h>
#include <netdb.h>
#include <sys/ipc.h>
#include <sys/sem.h>
#include <fcntl.h>
#include <sys/stat.h>
#include "exclu.h"

#define BROADCAST_ADDR "193.226.37.255"
#define LOCAL_PORT 6561
#define FIFO_NAME "enterChan"
#define MAX_HOST_NAME 50
#define MAX_MES_FIFO 6
#define LEN 400

int GetServerAddress(void);
int enterChan; /* handler pipe */
int lChan;
struct sockaddr_in sad;

int ExclAccessInitialize (void)
{
    int errCode;
    if((enterChan= open(FIFO_NAME, O_RDWR))<0)
        return ERROR4;
    if ((errCode=GetServerAddress()) != 0)
        return errCode;
}

```

Anexa B: Sincronizare in sisteme distribuite

```

bzero((char *)&sad, sizeof(sad));
sad.sin_family = AF_INET;
sad.sin_addr.s_addr = inet_addr(MyServerAddress);
sad.sin_port = htons(LOCAL_PORT);
if(!IChan = socket(AF_INET, SOCK_DGRAM, 0)) < 0)
    return ERRORS;
return 0;
}

void ExclAccessOpen(void)
{
char info[LEN];
sprintf(info,"%s","REQUEST");
sendto(IChan,info,LEN,0,(struct sockaddr*)&sad,sizeof(sad));
read(enterChan,info,MAX_MES_FIFO);
}

void ExclAccessClose()
{
char info[LEN];
sprintf(info,"%s","RELEASE");
sendto(IChan,info,LEN,0,(struct sockaddr*)&sad,sizeof(sad));
}

char * GetMyAddress(void)
{
int errCode;
if ((errCode = GetServerAddress()) != 0)
    return NULL;
else
    return MyServerAddress;
}

int GetServerAddress(void)
/*determina adresa serverului propriu*/
{
struct in_addr;
char HostName[MAX_HOST_NAME];
struct hostent *host;

if (gethostname(HostName,MAX_HOST_NAME) != 0)
    return ERROR1;
if((host = gethostbyname(HostName)) == NULL)
    return ERROR2;
if(host != NULL)
    adr.s_addr = ((struct in_addr *)
(host->h_addr_list[0]))->s_addr;
if(strcmp(MyServerAddress,inet_ntoa(adr)) == NULL)
    return ERROR3;
return 0;
}

void ExclAccessEnd(void)
{
unlink(FIFO_NAME);
close(IChan);
}

/* ..... */
Funciile puse la dispozitia unui program client de catre
toti algoritmi care realizeaza excluderea reciproca, prezentati
in capitolul 4, la 4.3.2 si 4.3.3 -versiunea Windows

/* ..... */
#include <winsock.h>
#include <windows.h>
#include "client.h"
#include "exclwerr.h"

#define BROADCAST_ADDR "193.226.37.255"
#define LOCAL_PORT 6561
#define MAX_HOST_NAME 50
#define LEN 400

int GetServerAddress(void);
HANDLE enterChan;
int IChan;

struct sockaddr_in sad;

int ExclAccessInitialize(void)
{
int errCode;
enterChan=OpenSemaphore(SYNCHRONIZE,TRUE,
"LocalSem");

if (enterChan==NULL)
    return ERROR4;
if ((errCode = GetServerAddress()) != 0)
    return errCode;
memset((void FAR *)&sad, 0, sizeof(sad));
sad.sin_family = AF_INET;
sad.sin_addr.s_addr = inet_addr(MyServerAddress);
sad.sin_port = htons(LOCAL_PORT);
if(!IChan = socket(AF_INET, SOCK_DGRAM, 0)) < 0)
    return ERRORS;
return 0;
}

void ExclAccessOpen(void)
{
char info[LEN];
sprintf(info,"%s","REQUEST");
sendto(IChan,info,LEN,0,(struct sockaddr*)&sad,sizeof(sad));
WaitForSingleObject(enterChan,INFINITE);
}

void ExclAccessClose()
{
char info[LEN];
sprintf(info,"%s","RELEASE");
sendto(IChan,info,LEN,0,(struct sockaddr*)&sad,sizeof(sad));
}

void ExclAccessEnd(void)
{
closesocket(IChan);
}

char * GetMyAddress(void)
{
int errCode;
if ((errCode = GetServerAddress()) != 0)
    return NULL;
else
    return MyServerAddress;
}

int GetServerAddress(void)
/*determina adresa serverului propriu*/
{
struct in_addr adr;
char HostName[MAX_HOST_NAME];
struct hostent *host;

if (gethostname(HostName,MAX_HOST_NAME) != 0)
    return ERROR1;
if((host = gethostbyname(HostName)) == NULL)
    return ERROR2;
if(host != NULL)
    adr.s_addr = ((struct in_addr *) (host->h_addr_list[0]))->s_addr;
if(strcmp(MyServerAddress,inet_ntoa(adr)) == NULL)
    return ERROR3;
return 0;
}

Program demonstrativ pentru toti algoritmi de excludere
reciproca prezentati in capitolul 4, la 4.3.2 si 4.3.3 -versiunea
Novell NetWare
- Programul utilizeaza functiile :
int ExclAcclInitialize(void);
int SetUsers(char users[MAX_CON][USER_LEN],
int hosts, int noObjects);

void ExclAccEnd(void);
void OpenExclAccess(void);
void CloseExclAccess(void);
puse la dispozitie de implementarea fiecarui algoritim

/* ..... */

```

```
#include <stdio.h>
#include <conio.h>
#include <stdlib.h>
#include "exc.h"

void GetUsers (char users[MAX_CONJ][USER_LEN], int *hosts);
void GetObjects (int *noOb);
void ExitOnError (int val);
extern code;

void main(void)
{
    char c;
    char users[MAX_CONJ][USER_LEN];
    int hosts, noOb;
    int value, answer;
    clrscr();
    printf("Demonstration program for exclusive access:\n\n");
    (1) Lamport algorithm\n(2) Ricard & Agrawala algorithm\n
    (3) Ray algorithm\n(4) Token ring algorithm\n(5) Token\
    algorithm\n\n
    Choose an algorithm:");
    scanf("%d", &answer);
    do{
        GetUsers (users, &hosts);
        switch (answer){
            case 1:
            case 2:
            case 4:
            case 5:
                value = SetUsers(users,hosts,1);
                break;
            case 3:
                GetObjects(&noOb);
                value = SetUsers(users,hosts,noOb);
                break;
        }
        if (value) printf ("User %d not available! Try again.\n",value);
    }while (value);
    if ((value = ExclAccessInitialize()) != 0) ExitOnError(value);
    do{
        printf("Continue (Yes,No)?");
        c = getch();
        putchar("\n");
        if (c == 'n' || c == 'N'){
            ExclAccessEnd();
            printf("\nEnd demonstration.\n");
            exit(0);
        }
        printf("Waiting to enter in critical section...\n");
        ExclAccessOpen();
        printf("In critical section...\n");
        c = getch();
        ExclAccessClose();
    }while (1);
}

void GetUsers (char users[MAX_CONJ][USER_LEN], int *hosts)
{
    int i;
    printf("\nNumber of processes: \n");
    scanf("%d",hosts);
    for(i=0;i<*hosts;i++){
        printf("User name for stations:");
        scanf("%s", users[i]);
    }
}

void GetObjects (int *noOb)
{
    printf("Number of objects = \n");
    scanf("%d",noOb);
}

void ExitOnError(int val){
    printf(" Eroare %d %d\n",val, code);
    ExclAccessEnd();
    exit(1);
}
```

B2. Implementarea algoritmului Lamport pentru obținerea excluderii reciproce - varianta Unix

```
.....
Server care implementeaza algoritmul Lamport
-Suportul teoretic este prezentat in capitolul 4, la 4.3.2.1.
..... */

#include <stdio.h>
#include <string.h>
#include <signal.h>
#include <unistd.h>
#include <sys/time.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <netinet.h>
#include <arpa/inet.h>
#include <netdb.h>
#include <sys/ipc.h>
#include <sys/sem.h>
#include <fcntl.h>
#include <sys/types.h>
#include <sys/stat.h>
#include <stdarg.h>
#include "Lampuer.h"

#define BROADCAST_ADDR "193.226.37.255"
#define REMOTE_PORT 6560
#define LOCAL_PORT 6561
#define BROADCAST_PORT 6562
#define FIFO_NAME "enterChan"
#define PERMS 0666
#define MAX_STRING 40
#define MAX_HOST_NO 20
#define MAX_HOST_NAME 256
#define MAX_TYPE_MES_LEN 20
#define LEN 400
#define TRUE 1
#define FALSE 0
#define DEBUG TRUE

/* .....
* Functii si variabile interne serverului
* ..... */

typedef struct{ /* structura elementelor cozii de cereri */
    int ts; /* marca de timp */
    int conn; /* identificatorul host-ului */
}queue;
void CreateBroadcastChan(void); /* creare canal broadcast */
int CreateRemoteChan(void); /* creare canal remoteChan*/
int CreateLocalChan(void); /* creare canal localChan */
int CreateClientChan(void);
void CreateEnterChan(void) /* creare canal enterChan*/;
void CloseEnterChan(void);
void FindAddress(void); /* determina adresele host-urilor*/
void InOrder(queue q[MAX_HOST_NO]); /* ordoneaza coada*/
int max(int a,int b);
void Debug(const char *c, ...);
int enterChan ; /* canal implementat prin FIFO */
int bsd; /* soclu pt. transmisii broadcast */
struct sockaddr_in sad, cad,cla,sva,I,Cad; /* adrese Inet */
int myId; /* identificatorul host-ului local */char
connTab[MAX_HOST_NO][MAX_HOST_NO];
int processNo; /* numarul de noduri participante la alg.*/
queue reqQueue[MAX_HOST_NO]; /* coada de cereri */ int
queueLength; /* lungimea cozii de cereri */

void main (int argc, char *argv[] )
{
    struct sockaddr_in sdr; /* adresa Inet */
    fd_set rfd; /* read file descriptor set*/
    fd_set afd; /* active file descriptor set*/
    int rChan,IChan,cli; /* identificatori canalelor */
    int sval,I,Cad; /* folosite la recvfrom */
}
```

Anexa B: Sincronizare in sisteme distribuite

```

int sender, ts; /* folosite pentru receptia mesajelor */
int LocalTS=0, HighTS = 0; /* retin marcele de timp */
int inAccess = FALSE; /* flag pentru intrarea in SC */
char sBf[LEN], rBf[LEN]; /* buffere */
char mes[MAX_TYPE_MES_LEN]; /* buffer */
int replyCounter = 0; /* contorizeaza mesajele REPLY */
int nfds, err, i, j; /* de lucru */

/* .....
faza de initializare: se prileia din linia de comanda identificatorul
host-ului local si numarul de procese.
Se creaza canalele.
Fiecare server trimite broadcast identificatorul sau si adresa sa
Se pun in asteptare de mesaje canalele remoteChan, localChan
..... */
if(argc==3) {
    printf ("Usage : [name] [identifer] [hosts] \n");
    exit(0);
}
myId = atoi(argv[1]);
processNo = atoi(argv[2]);
CreateBroadcastChan();
CreateEnterChan();
getchar();
FindAddress();
for(i = 0; i < processNo; i++)
    Debug("Host= %d of address= %s \n", i, connTab[i]);
rChan = CreateRemoteChan();
lChan = CreateLocalChan();
cli = CreateClientChan();
nfds = getdtablesize();
FD_ZERO(&afds);
FD_SET(rChan, &afds);
FD_SET(lChan, &afds);

/* .....
faza de tratare a mesajelor primite de la clientul local
sau din retea, de la serverele de pe celelalte host-uri
..... */
while(TRUE){
    Debug("Waiting messages..\n");
    bcopy(char *) &afds, (char *) &rfd, sizeof(rfd));
    if (select(nfds, &rfd, (fd_set *) 0, (fd_set *) 0, (struct timeval *) 0) < 0){
        Debug("Select error \n");
        exit(1);
    }
    if(FD_ISSET(rChan, &rfd)){ /* s-a primit un mesaj din retea */
        sval = sizeof(sva);
        err = recvfrom(rChan, rBf, LEN, 0, (struct sockaddr*)&sva, &sval);
        if(err < 0)
            Panic(ERROR14);
        scanf(rBf, "%2d%3d%s", &sender, &ts, mes);
        Debug("Received from %d %s ts = %d..\n", sender, mes, ts);
        if(strcmp(mes, "REQUEST") == 0){ /* tratare REQUEST */
            reqQueue[queueLength].ts = ts;
            reqQueue[queueLength].conn = sender;
            queueLength++;
            InOrder(reqQueue);
            Debug("After reorder..\n");
            for(i=0; i < queueLength; i++)
                Debug("Q %dCon= %d ts = %d\n", i, reqQueue[i].conn, reqQueue[i].ts);
            HighTS = max(HighTS, ts); /* actualizare ts */
            if(sender != myId){ /* transmite REPLY */
                bzero((char *) &sdr, sizeof(sdr));
                sdr.sin_family = AF_INET;
                sdr.sin_addr.s_addr = inet_addr(connTab[sender]);
                sdr.sin_port = htons(REMOTE_PORT);
                printf(sBf, "%2d %3d %s", myId, 1, "REPLY");
                Debug("Sending: %s\n", sBf);
                sendto(cli, sBf, LEN, 0, (struct sockaddr*)&sdr, sizeof(sdr));
            }
        }
        else
            if(strcmp(mes, "REPLY") == 0){ /* tratare REPLY retea */
                replyCounter--;
                Debug("replyCounter = %d\n", replyCounter);
                if((replyCounter == 0) && (reqQueue[0].conn == myId)){
                    inAccess = TRUE; /* s-a obtinut accesul in CS */
                    printf(sBf, "%s", "enter");
                    write(enterChan, sBf, 6); /* deblocare client */
                }
            }
        else
            if(strcmp(mes, "RELEASE") == 0){ /* tratare RELEASE retea */
                for(i=0; i < processNo; i++)
                    if(reqQueue[i].conn == sender)
                        for(j = i; j < queueLength-1; j++)
                            reqQueue[j] = reqQueue[j+1];
                reqQueue[queueLength-1].ts = 0;
                reqQueue[queueLength-1].conn = 0;
                queueLength--;
                Debug(" RELEASE on remote , c nt= %d\n", replyCounter);
                if((replyCounter == 0) && (reqQueue[0].conn == myId)){
                    inAccess = TRUE; /* s-a obtinut accesul in CS */
                    printf(sBf, "%s", "enter");
                    write(enterChan, sBf, 6); /* deblocare client */
                }
            }
        }
    }
    if (FD_ISSET(lChan, &rfd)){ /* mesaj de la clientul local */
        lCadi = sizeof(lCad);
        err = recvfrom(lChan, rBf, LEN, 0, (struct sockaddr*)&lCad, &lCad);
        if(err < 0)
            Panic(ERROR15);
        scanf(rBf, "%s", mes);
        Debug("Received from my client mes= %s\n", mes);
        if(strcmp(mes, "REQUEST") == 0){ /* tratare REQUEST local */
            LocalTS = HighTS + 1; /* pregatire ts */
            replyCounter = processNo-1;
            Debug(" ReplyCounter = %d\n", replyCounter);
            for(i=0; i < processNo; i++){ /* transmite cererea in retea */
                bzero((char *) &sdr, sizeof(sdr));
                sdr.sin_family = AF_INET;
                sdr.sin_addr.s_addr = inet_addr(connTab[i]);
                sdr.sin_port = htons(REMOTE_PORT);
                printf(sBf, "%2d %3d %s", myId, LocalTS, "REQUEST");
                Debug("Sending: %s la %d\n", sBf, i);
                sendto(cli, sBf, LEN, 0, (struct sockaddr*)&sdr, sizeof(sdr));
            }
        }
        else
            if(strcmp(mes, "RELEASE") == 0){ /* tratare RELEASE local */
                inAccess = FALSE;
                for(i=0; i < processNo; i++)
                    if(reqQueue[i].conn == myId)
                        for(j = i; j < queueLength-1; j++)
                            reqQueue[j] = reqQueue[j+1];
                reqQueue[queueLength-1].ts = 0;
                reqQueue[queueLength-1].conn = 0;
                queueLength--;
                for(i=0; i < processNo; i++) /* transmite mesajul in retea */
                    if(i != myId){
                        bzero((char *) &sdr, sizeof(sdr));
                        sdr.sin_family = AF_INET;
                        sdr.sin_addr.s_addr = inet_addr(connTab[i]);
                        sdr.sin_port = htons(REMOTE_PORT);
                        printf(sBf, "%2d %3d %s", myId, 1, "RELEASE");
                        Debug("Sending: %s la %d\n", sBf, i);
                        sendto(cli, sBf, LEN, 0, (struct sockaddr*)&sdr, sizeof(sdr));
                    }
            }
        }
    }
}
/* .....
Functii folosite de server
..... */
void InOrder(queue ql[MAX_HOST_NO])
/* ordoneaza coada de cereri dupa marcele de timp */
{
    queue mx;
    int sort, i;
}

```

```

do{
    sort = TRUE;
    for(i=0;i<queueLength-2;i++){
        if(qL[i].ts > qL[i+1].ts){
            mx = qL[i];
            qL[i] = qL[i+1];
            qL[i+1] = mx;
            sort = FALSE;
        }
    }
}while(!sort);
}

void CreateBroadcastChan(void)
/* creaza canalul pentru transmisii broadcast */
{
    int i;
    bzero((char *)&sad, sizeof(sad));
    sad.sin_family = AF_INET;
    sad.sin_addr.s_addr = inet_addr(BROADCAST_ADDR);
    sad.sin_port = htons(BROADCAST_PORT);
    if((bsd = socket(AF_INET, SOCK_DGRAM, 0)) < 0)
        Panic(ERROR1);
    i = 1;
    if(setsockopt(bsd, SOL_SOCKET, SO_BROADCAST, (char *)
        &i, sizeof(i)) {
        Panic(ERROR2);
    }
    bzero((char *)&cad, sizeof(cad));
    cad.sin_family = AF_INET;
    cad.sin_addr.s_addr = htonl(INADDR_ANY);
    cad.sin_port = htons(BROADCAST_PORT);
    if(bind(bsd, (struct sockaddr *)&cad, sizeof(cad)) < 0)
        Panic(ERROR3);
}

void FindAddress(void)
/* memoreaza in tabela de adrese adresele si identificatorii */
/* tuturor host-urilor */
{
    struct in_addr adr;
    struct sockaddr cadr;
    struct hostent *host;
    int cadr;
    char HostName[MAX_HOST_NAME];
    char sBf[LEN];
    char code[MAX_STRING]; char string[MAX_STRING];
    int idRec,i;

    if (gethostname(HostName,256) != 0)
        Panic(ERROR4);
    if((host = gethostbyname(HostName)) == NULL)
        Panic(ERROR5);
    if(host != NULL)
        adr.s_addr = ((struct in_addr*)(host->h_addr_list[0]))->s_addr;
    if(strcmp(code,inet_ntoa(adr)) == NULL)
        Panic(ERROR6);
    strcpy(conTab[myId], code);
    sprintf(sBf,"%2d %s",myId,code);
    sendto(bsd,sBf,LEN,0,(struct sockaddr *)&sad,sizeof(sad));
    for(i=0;i<processNo;i++){
        cadr = sizeof(cad);
        recvfrom(bsd, sBf, LEN, 0, (struct sockaddr *)&cadr, &cadr);
        sscanf(sBf,"%2d%s",&idRec,string);
        if(idRec != myId)
            strcpy(connTab[idRec],string);
    }
}

int CreateRemoteChan(void)
/* creaza canalul pentru transmisii in retea */
{
    int s;
    if((s = socket(AF_INET, SOCK_DGRAM, 0)) < 0)
        Panic(ERROR7);
    bzero((char *)&sad, sizeof(sad));
    sva.sin_family = AF_INET;
    sva.sin_addr.s_addr = htonl(INADDR_ANY);
    sva.sin_port = htons(REMOTE_PORT);

    if(bind(s, (struct sockaddr *)&sva, sizeof(sva)) < 0)
        Panic(ERROR8);
    return s;
}

int CreateLocalChan(void)
/* creaza canalul pentru receptia mesajelor de la clientul local */
{
    int s;
    if((s = socket(AF_INET, SOCK_DGRAM, 0)) < 0)
        Panic(ERROR9);
    bzero((char *)&ICad, sizeof(ICad));
    ICad.sin_family = AF_INET;
    ICad.sin_addr.s_addr = htonl(INADDR_ANY);
    ICad.sin_port = htons(LOCAL_PORT);
    if(bind(s, (struct sockaddr *)&ICad, sizeof(ICad)) < 0)
        Panic(ERROR10);
    return s;
}

int CreateClientChan()
/* creaza canalul pentru emisia mesajelor in retea */
{
    int s;
    if((s = socket(AF_INET, SOCK_DGRAM, 0)) < 0)
        Panic(ERROR11);
    bzero((char *)&cla, sizeof(cla));
    cla.sin_family = AF_INET;
    cla.sin_addr.s_addr = htonl(INADDR_ANY);
    cla.sin_port = htons(REMOTE_PORT);
    return s;
}

int max(int a,int b)
{
    if(a > b) return a;
    else return b;
}

void CreateEnterChan(void)
{
    unlink(FIFO_NAME);
    if( mknod(FIFO_NAME, S_IFIFO | PERMS, 0)<0)
        Panic(ERROR12);
    if((enterChan = open(FIFO_NAME, O_RDWR))<0)
        Panic(ERROR13);
}

void CloseEnterChan(void)
{
    unlink(FIFO_NAME);
}

void Panic( int err){
    printf("Exit: %d error\n", err);
    exit(0);
}

void Debug(const char *c, ...){
    va_list args;
    if (DEBUG){
        va_start(args,c);
        vprintf(c, args);
        va_end(args);
    }
}

```


B3. Implementarea algoritmului Ricard si Agrawala pentru obtinerea excluderii reciproce varianta Novell NetWare

Program care implementeaza algoritmul Ricard si Agrawala
-Suportul teoretic este prezentat in capitolul 4, la 4.3.2.2

```
#include <nxt.h>
#include <nit.h>
#include <stdio.h>
#include <stdlib.h>
#include <dos.h>
#include <bios.h>
#include <mem.h>
#include <string.h>
#include "exclacc.h"
```

```
#define TRUE 1
#define FALSE 0
#define MAX_NODES 30
#define DATA_LENGTH 80
#define MAX_CON 30
```

```
/* .....
* ..... Functii si variabile folosite
* ..... */
```

```
int Initialize(void); /* asigura initializare IPX */
int InitReceive(void); /* asigura initializarea unui ECB pt. emisie */
int InitSend(void); /* asigura initializarea unui ECB pt. receptie */
int GetConnectionNumbers(char sers[MAX_CON][USER_LEN]);
void far _loadds ReceiveESR(void); /* rutina ESR */
void SendRequest(int serverConnNo, char *message);
void SendReply(IPXHeader *ppacketHeaderRec, char *message);
void RicAgrawRequest(void);
void RicAgrawEnter(void);
void RicAgrawRelease(void);
void ReceiveRequest(char *message);
void ReceiveReply(void);
void RicAgrawClose(void);
```

```
int noNodes; /* numarul de statii */
int localTS; /* valori pt. marcele de timp */
int highTS=0;
int replyCounter; /* contorizeaza mesajele de reply */
int request;
int laterTab[MAX_NODES];
int inSection; /* indica intrarea intr-o sectiune critica */
int later;
int noOb;
```

```
WORD connectionList[MAX_CON];
IPXHeader packetHeaderRec;
char packetDataSnd[DATA_LENGTH]; /* buffere */
char packetDataRec[DATA_LENGTH];
char packetDataReply[DATA_LENGTH];
char bufferSnd[DATA_LENGTH];
ECB ecbEmission;
ECB ecbReception;
WORD socketDest=0x5000;
WORD socketSnd=0x6000;
WORD socketRec = 0x5000;
int timeStamp, sender;
int code;
```

```
/* .....
* ..... Functia de conectare la sistem
* ..... */
int ExclAccessInitialize(void)
{
    int val;

    if ((val = Initialize()) != 0) return val;
    if ((val = InitSend()) != 0) return val;
    if ((val = InitReceive()) != 0) return val;
    return 0;
}

int SetUsers(char users[MAX_CON][USER_LEN], int hosts, int noObjects)
{
    int value;
    noNodes = hosts;
    value = GetConnectionNumbers(users);
    if (value != 0) return value;
    noOb = noObjects;
    return 0;
}

int GetConnectionNumbers(char users[MAX_CON][USER_LEN])
{
    int i=0;
    WORD numberOfConnections = 0;
    WORD maxConnections = 20;
    for(i=0; i<noNodes; i++){
        GetObjectConnectionNumbers(users[i], OT_USER,
            &numberOfConnections, &connectionList[i], maxConnections);
        if(numberOfConnections == 0) return i;
    }
    return 0;
}

/* .....
* ..... Functia pentru obtinerea accesului in sectiunea critica
* ..... */
void ExclAccessOpen()
{
    RicAgrawRequest();
    RicAgrawEnter();
}

/* .....
* ..... Functia pentru deconectarea de la sistem
* ..... */
void ExclAccessEnd(void)
{
    IPXCancelEvent(&ecbReception);
    IPXCloseSocket(socketRec);
    IPXCloseSocket(socketSnd);
}

/* .....
* ..... Functia pentru iesirea din sectiunea critica
* ..... */
void ExclAccessClose(void)
{
    int j;
    disable();
    request = FALSE;
    for(j=1; j<noNodes; j++){
        if(laterTab[connectionList[j]-1])
            SendRequest(connectionList[j], "Reply");
        laterTab[connectionList[j]-1] = FALSE;
    }
    inSection = FALSE;
    enable();
}

int Initialize(void)
{
    if ((code = IPXInitialize()) != 0) return 1;
    return 0;
}
```

```

}

int InitSend(void)
{
    if ((code = IPXOpenSocket((BYTE*)&socketSnd,0x00)) != 0)
        return 2;
    return 0;
}

int InitReceive(void)
{
    int code;

    /* initializarea structurilor de date */
    ecbReception.ESRAddress = ReceiveESR;
    ecbReception.fragmentCount = 2;
    ecbReception.socketNumber = socketRec;
    ecbReception.fragmentDescriptor[0].address =
    &packetHeaderRec;
    ecbReception.fragmentDescriptor[0].size = sizeof(IPXHeader);
    ecbReception.fragmentDescriptor[1].address = (BYTE *)
    packetDataRec;
    ecbReception.fragmentDescriptor[1].size = DATA_LENGTH;

    /* deschiderea soclului si intrarea in "ascultare" de pachete */
    if ((code = IPXOpenSocket((BYTE*)&socketRec,0x00)) != 0)
    return 3;
    IPXListenForPacket(&ecbReception);
    IPXRelinquishControl();
    return 0;
}

void far _loadds ReceiveESR(void)
{
    disable();
    if (ecbReception.completionCode == 0){
        if(stromp(packetDataRec, "Reply") == 0)
            ReceiveReply();
        else
            ReceiveRequest(packetDataRec);
    }
    enable();
    IPXListenForPacket(&ecbReception);
}

void SendRequest(int serverConnNo, char * message)
{
    BYTE serverNode[6];
    BYTE serverNetwork[4];
    IPXHeader packetHeader;
    WORD socketDest=0x5000,socket;

    GetInternetAddress(serverConnNo,(char *)
        serverNetwork,(char *)serverNode,&socket);
    movmem(serverNetwork, packetHeader.destination.network,4);
    movmem(serverNode, packetHeader.destination.node,6);
    movmem(&socketDest, packetHeader.destination.socket,2);
    movmem(serverNode, ecbEmission.immediateAddress,6);
    movmem(message, packetDataSnd, DATA_LENGTH);
    ecbEmission.ESRAddress = NULL;
    ecbEmission.fragmentCount = 2;
    ecbEmission.socketNumber = socketSnd;
    ecbEmission.fragmentDescriptor[0].address = &packetHeader;
    ecbEmission.fragmentDescriptor[0].size = sizeof(IPXHeader);
    ecbEmission.fragmentDescriptor[1].address = (BYTE *)
        packetDataSnd;
    ecbEmission.fragmentDescriptor[1].size = DATA_LENGTH;
    while (ecbEmission.inUseFlag != 0) IPXRelinquishControl();
    IPXSendPacket(&ecbEmission);
}

void SendReply(IPXHeader *ppacketHeaderRec, char *message)
{
    IPXHeader packetHeader;
    movmem(ppacketHeaderRec->source.network,
        packetHeader.destination.network,4);
    movmem(ppacketHeaderRec->source.node,
        packetHeader.destination.node,6);
    movmem(&socketDest, packetHeader.destination.socket,2);
    movmem(message, packetDataReply, DATA_LENGTH);
    ecbEmission.ESRAddress = NULL;
    ecbEmission.fragmentCount = 2;
    ecbEmission.socketNumber = socketSnd;
    ecbEmission.fragmentDescriptor[0].address = &packetHeader;
    ecbEmission.fragmentDescriptor[0].size = sizeof(IPXHeader);
    ecbEmission.fragmentDescriptor[1].address =
        (BYTE *) packetDataReply;
    ecbEmission.fragmentDescriptor[1].size = DATA_LENGTH;
    while (ecbEmission.inUseFlag != 0) IPXRelinquishControl();
    IPXSendPacket(&ecbEmission);
}

void RicAgrawRequest(void)
{
    int j;
    disable();
    localTS = highTS + 1;
    request = TRUE;
    replyCounter = noNodes - 1;
    for(j=0;j<noNodes;j++){
        if(connectionList[j] != connectionList[0]){
            memcpy(&bufferSnd[0],&localTS,2);
            memcpy(&bufferSnd[2],&connectionList[0],2);
            SendRequest(connectionList[j],bufferSnd);
        }
    }
    enable();
}

void RicAgrawEnter(void)
{
    int ind1;

    do{
        disable();
        ind1 = replyCounter;
        enable();
    }while (!bioskey(1) && (ind1 != 0));
    if (ind1) {
        printf("Cancelled by user intervention!\n");
        exit(1);
    }
    inSection = TRUE;
}

void ReceiveRequest(char *message)
{
    memcpy(&timeStamp,&message[0],2);
    memcpy(&sender,&message[2],2);
    highTS = max(highTS,timeStamp);
    later = (request && ((timeStamp > localTS) ||
        ((timeStamp == localTS) && (sender->connectionList[0])))
        || inSection;
    if (later)
        laterTab[sender-1] = TRUE;
    else
        SendReply(&packetHeaderRec,"Reply");
}

void ReceiveReply(void)
{
    replyCounter--;
}

```

B4. Implementarea algoritmului pentru obținerea excluderii reciproce în accesul la k resurse -varianta Windows

```

Server care implementeaza algoritmul excluderii
reciproce in accesul la K resurse
-Suportul teoretic este prezentat in capitolul 4, la 4.3.2.3.
..... */
#include <windows.h>
#include <winsock.h>
#include <stdio.h>
#include <conio.h>
#include <string.h>
#include <process.h>
#include <stdlib.h>
#include <sys/types.h>
#include <fcntl.h>
#include <sys/stat.h>
#include "raywerr.h"

#define BROADCAST_ADDR "193.226.37.255"
#define BROADCAST_PORT 6562
#define REMOTE_PORT 6560
#define LOCAL_PORT 6561
#define MAX_HOST_NO 20
#define MAX_STRING 20
#define MAX_BUFFER 256
#define LEN 400
#define TRUE 1
#define FALSE 0
#define DEBUG TRUE
..... */
* Functii si variabile interne serverului *
..... */

void CreateBroadcastChan(void);
SOCKET CreateRemoteChan(void);
SOCKET CreateLocalChan(void);
SOCKET CreateClientChan(void);
void FindAddress(void);
int NotInCS(void);
void Debug(const char *c,...);
struct sockaddr_in sad, cad, cla, sva, lCad;
char connTab[MAX_HOST_NO][MAX_HOST_NO];
char HostName[MAX_BUFFER];
int processNo, myID;
int replyCounter[MAX_HOST_NO];
int bsd;

void main(int argc, char *argv[])
{
    struct sockaddr_in sdr; /* adresa Inet */
    fd_set rfd; /* read file descriptor set */
    fd_set afds; /* active file descriptor set */
    SOCKET rChan, lChan, cli; /* identificatorii canalelor */
    int sval, lCad; /* folosite la recvrom */
    int sender, ts; /* folosite pentru receptia mesajelor */
    int LocalITS=0, HighTS=0; /* retin marcile de timp */
    BOOL inAccess=FALSE; /* flag pentru intrarea in SC */
    char sBf[LEN], rBf[LEN]; /* buffere */
    HANDLE hStdin, hStdout, enterChan; /* enterChan este semafor */
    WSADATA WSAData;
    DWORD result;
    BOOL later, reqAccess=FALSE; /* flag-uri */
    char temp[LEN], mes[MAX_STRING]; /* buffere */
    int laterTab[MAX_HOST_NO]; /* no. mes. Reply de asteptat */
    int objNo; /* no. de obiecte */
    int nfds, err, i, status; /* var. de lucru */
    ..... */
faza de initializare: se preria din linia de comanda identificatorul
host-ului local si numarul de procese si numarul de obiecte
Se creeaza canalele.
Fiecare server trimite broadcast identificatorul sau si adresa sa
Se pun in asteptare de mesaje canalele remoteChan, localChan
..... */
FreeConsole();
AllocConsole();
SetConsoleTitle(" Server Exclusive Access");
hStdin=GetStdHandle(STD_INPUT_HANDLE);
hStdout=GetStdHandle(STD_OUTPUT_HANDLE);
if ((status = WSASStartup(MAKEWORD(1,1), &WSAData) == 0)
    Debug("WSA ok \n");
else
    Panic(ERROR12);
if(argc!=4) {
    printf("Usage : [name] [identifier] [hosts] [objects]\n");
    exit(0);
}
for(i = 0; i<MAX_HOST_NO; i++)
    laterTab[i]=0;
enterChan = CreateSemaphore(NULL, 0L, 1L, "LocalSem");
myID = atoi(argv[1]);
processNo = atoi(argv[2]);
objNo = atoi(argv[3]);
CreateBroadcastChan();
ReadConsole(hStdin, temp, 10, &result, NULL);
FindAddress();
for(i = 0; i<processNo; i++)
    Debug(" Host %d of address %s \n", i, connTab[i]);
getchar();
closesocket(bsd);
rChan = CreateRemoteChan();
lChan = CreateLocalChan();
cli = CreateClientChan();
FD_ZERO(&afds);
FD_SET(rChan, &afds);
FD_SET(lChan, &afds);
..... */
faza de tratare a mesajelor primite de la clientul local
sau din retea, de la serverele de pe celelalte host-uri
..... */
while(TRUE){
    Debug("Waiting messages... \n");
    CopyMemory(char *)&rfd, (char *)&afds, sizeof(rfd);
    if(select(nfds, &rfd, (fd_set *)0, (fd_set *)0, (struct timeval *)0)<0){
        err = WSAGetLastError();
        Debug("Select error: %d\n", err);
    }
    if(FD_ISSET(rChan, &rfd)){ /* s-a primit un mesaj din retea */
        sval = sizeof(sva);
        err = recvfrom(rChan, rBf, LEN, 0, (struct sockaddr*)&sva, &sval);
        if(err<0)
            Debug("Recvrom -remote- error\n");
            sscanf(rBf, "%2d%3d%s", &sender, &ts, mes);
            Debug("%s from %d ts = %d.. \n", mes, sender, ts);
            if(strcmp(mes, "REQUEST") == 0){ /* tratare REQUEST */
                HighTS = max(HighTS, ts);
                later = (inAccess) && ((ts > LocalITS) ||
                    ((ts == LocalITS) && (sender > myID)));
                if (later) /* se amina transmiterea mes. REPLY */
                    laterTab[sender]++;
                Debug("laterTab = ");
                for(i=0; i<10; i++)
                    Debug(" %d", laterTab[i]);
                Debug("\n");
            }
            else /* se poate transmite REPLY */
                memset((void FAR *)&sdr, 0, sizeof(sdr));
                sdr.sin_family = AF_INET;
                sdr.sin_addr.s_addr = inet_addr(connTab[sender]);
                sdr.sin_port = htons(REMOTE_PORT);
                sprintf(sBf, "%2d %3d %s", myID, 1, "REPLY");
                Debug("Sending: %s\n", sBf);
                sendto(cli, sBf, LEN, 0, (struct sockaddr*)&sdr, sizeof(sdr));
            }
        }
    }
    else
        if(strcmp(mes, "REPLY") == 0) /* tratare REPLY */
            replyCounter[sender] = replyCounter[sender] - ts;
            if(reqAccess && (NotInCS()) >= processNo - objNo){

```

```

inAccess = TRUE;
reqAccess = FALSE;
ReleaseSemaphore(enterChan,1,NULL);
}
}
}
if(FD_ISSET(lChan,&rfd)){ /* s-a primit un mesaj de la client */
lCadl = sizeof(lCad);
err = recvfrom(lChan,rBf,LEN,0,(struct sockaddr*)&lCad,&lCadl);
if(err<0)
    Debug("Recvfrom -local- error\n");
    sscanf(rBf,"%s",mes);
    Debug("Received from my cli %s.\n",mes);
if(strncmp(mes,"REQUEST") == 0){ /* tratare REQUEST */
    reqAccess = TRUE;
    LocalTS = HighTS + 1;
    for(i=0;j<processNo;i++){
        if(i!= myID){
            memset((void FAR *)&sdr, 0, sizeof(sdr));
            sdr.sin_family = AF_INET ;
            sdr.sin_addr.s_addr = inet_addr(connTab[i]);
            sdr.sin_port = htons(REMOTE_PORT);
            sprintf(sBf,"%2d %3d %s",myID,LocalTS,"REQUEST");
            Debug("Sending:%s\n",sBf);
            sendto(cli,sBf,LEN,0,(struct sockaddr*)&sdr,sizeof(sdr));
            replyCounter[i]++;
        }
    }
}
else
if(strncmp(mes,"RELEASE") == 0){ /* tratare RELEASE */
inAccess = FALSE;
for(i=0;j<processNo;i++){
if(laterTab[i] != 0){
    memset((void FAR *)&sdr, 0, sizeof(sdr));
    sdr.sin_family = AF_INET ;
    sdr.sin_addr.s_addr = inet_addr(connTab[i]);
    sdr.sin_port = htons(REMOTE_PORT);
    sprintf(sBf,"%2d %3d %s",myID,laterTab[i],"REPLY");
    Debug("Sending:%s\n",sBf);
    sendto(cli,sBf,LEN,0,(struct sockaddr*)&sdr,sizeof(sdr));
    laterTab[i] = 0;
}
}
}
}
}
/* .....
*          Functii folosite de server
* ..... */
int NotInCS(void) /* contorizeaza nr. de procese care nu */
{ /* sunt in sectiune critica */
    int cnt,j;
    cnt = 0;
    for(j=0;j<processNo;j++){
        if((j != myID) && (replyCounter[j] == 0))
            cnt++;
    }
    return cnt;
}
void CreateBroadcastChan(void)
{
    int i;
    memset((void FAR *)&sad, 0, sizeof(sad));
    sad.sin_family = AF_INET ;
    sad.sin_addr.s_addr = inet_addr(BROADCAST_ADDR);
    sad.sin_port = htons(BROADCAST_PORT);
    if((bsd = socket(AF_INET, SOCK_DGRAM, 0)) < 0)
        Panic(ERROR1);
    i = 1;
    if(setsockopt(bsd, SOL_SOCKET, SO_BROADCAST, (char *)
        &i, sizeof(i)))
        Panic(ERROR2);
    memset((void FAR *)&cad, 0, sizeof(cad));
    cad.sin_family = AF_INET ;
    cad.sin_addr.s_addr = htonl(INADDR_ANY);
    cad.sin_port = htons(BROADCAST_PORT);
    if(bind(bsd, (struct sockaddr *)&cad, sizeof(cad)) < 0)
        Panic(ERROR3);
}
void FindAddress(void) /* determina adresele host-urilor */
{
    struct in_addr adr;
    struct sockaddr cadr;
    struct hostent *host;
    int cadr;
    char code[MAX_STRING], string[MAX_STRING];
    int idRec,i;
    char sBf[LEN];
    if(gethostname(HostName,256) != 0)
        Panic(ERROR4);
    if((host = gethostbyname(HostName)) == NULL)
        Panic(ERROR5);
    if(host != NULL)
        adr.s_addr = ((struct in_addr *) (host->h_addr_list[0]))->s_addr;
    if(strcpy(code,inet_ntoa(adr)) == NULL)
        Panic(ERROR6);
    strcpy(connTab[myID], code);
    sprintf(sBf,"%2d %s",myID,code);
    sendto(bsd,sBf,LEN,0,(struct sockaddr*)&sad,sizeof(sad));
    for(i=0;i<processNo;i++){
        cadr = sizeof(cad);
        recvfrom(bsd, sBf, LEN, 0, (struct sockaddr*)&cadr, &cadr);
        sscanf(sBf,"%2d%s",&idRec,string);
        if(idRec != myID)
            strcpy(connTab[idRec],string);
    }
}
SOCKET CreateRemoteChan(void) {
    SOCKET s;
    if((s = socket(AF_INET, SOCK_DGRAM, 0)) < 0)
        Panic(ERROR7);
    memset((void FAR *)&sva, 0, sizeof(sva));
    sva.sin_family = AF_INET ;
    sva.sin_addr.s_addr = htonl(INADDR_ANY);
    sva.sin_port = htons(REMOTE_PORT);
    if(bind(s, (struct sockaddr *)&sva, sizeof(sva)) < 0)
        Panic(ERROR8);
    return s;
}
SOCKET CreateLocalChan(void) {
    SOCKET s;
    if((s = socket(AF_INET, SOCK_DGRAM, 0)) < 0)
        Panic(ERROR9);
    memset((void FAR *)&lCad, 0, sizeof(lCad));
    lCad.sin_family = AF_INET ;
    lCad.sin_addr.s_addr = htonl(INADDR_ANY);
    lCad.sin_port = htons(LOCAL_PORT);
    if(bind(s, (struct sockaddr *)&lCad, sizeof(lCad)) < 0)
        Panic(ERROR10);
    return s;
}
SOCKET CreateClientChan() {
    SOCKET s;
    if((s = socket(AF_INET, SOCK_DGRAM, 0)) < 0)
        Panic(ERROR11);
    memset((void FAR *)&cla, 0, sizeof(cla));
    cla.sin_family = AF_INET ;
    cla.sin_addr.s_addr = htonl(INADDR_ANY);
    cla.sin_port = htons(REMOTE_PORT);
    return s;
}
void Debug(const char *c,...){
    va_list args;
    if (DEBUG) {
        va_start(args,c);
        vprintf(c,args);
        va_end(args);
    }
}
void Panic(int err){
    printf("Exit: %d error\n", err);
    exit(0);
}
}

```

B5. Implementarea algoritmului bazat pe un token circulant pentru obținerea excluderii reciproce variante Novell NetWare

```

.....
Server care implementeaza algoritmul bazat pe
un token circulant
-Suportul teoretic este prezentat in capitolui 4, la 4.3.3.1.
.....
#include <nxt.h>
#include <nit.h>
#include <stdio.h>
#include <dos.h>
#include <mem.h>
#include <string.h>
#include "exclacc.h"
* .....
* Functii si variabile interne serverului
* .....
#define DATA_LENGTH 80
int Initialize(void); /* asigura initializare IPX */
int InitSend(void); /* asigura initializarea unui ECB pt. emisie */
int InitReceive(void); /* asigura initializarea unui ECB pt. receptie */
int GetConnectionNumbers
(char users[MAX_CON][USER_LEN]);
void far _loadds ReceiveESR(void); /* rutina ESR */
void SendToken(char *message); /* trimite token-ul la vecin */
void TokenRequest(void); /* cere token-ul */
void TokenEnter(void); /* blocare pina la primirea token-ului */
void ReceiveToken(void); /* apelata la primirea token-ului */
int noNodes; /* numarul de statii */
int request; /* cerere token */
int enter; /* indica accesul in SC */
int myPosition, nextPosition;
int noOb;
WORD connectionList[30]; /* nr de ale conexiune statiilor */
IPXHeader nextAddress; /* zone pt. header */
IPXHeader packetHeaderSnd;
IPXHeader packetHeaderRec;
char packetDataRec[DATA_LENGTH]; /* buffere */
char packetDataToken[DATA_LENGTH];
char bufferSnd[DATA_LENGTH];
ECB ecbEmission;
ECB ecbReception;
WORD socketRec = 0x5000;
WORD socketDest=0x5000;
WORD socketSnd=0x6000, socket;
int code;

int ExclAccessInitialize(void) /* conectare la sistem */
{
    int val;
    if ((val = Initialize()) != 0) return val;
    if ((val = InitSend()) != 0) return val;
    if ((val = InitReceive()) != 0) return val;
    if (connectionList[myPosition] == connectionList[0])
        SendToken("OK");
    return 0;
}

int SetUsers(char users[MAX_CON][USER_LEN],
int hosts, int noObjects)
/* cunoscind numele user-ilor de pe fiecare statie, determina */
/* numerele de conexiune pentru toate statiile */
{
    int value;
    noNodes = hosts;
    value = GetConnectionNumbers(users);
    if (value != 0) return value;
    noOb = noObjects;
    return 0;
}

int Initialize(void) /* initializare IPX */
{
    if ((code = IPXInitialize()) != 0)
        return 1;
    else
        return 0;
}

int InitSend(void)
{
    if ((code = IPXOpenSocket((BYTE *)&socketSnd,0x00)) != 0)
        return 2;
    else
        return 0;
}

int InitReceive(void)
{
    int code;
    /* initializarea structurilor de date */
    ecbReception.ESRAddress = ReceiveESR;
    ecbReception.fragmentCount = 2;
    ecbReception.socketNumber = socketRec;
    ecbReception.fragmentDescriptor[0].address =
        &packetHeaderRec;
    ecbReception.fragmentDescriptor[0].size = sizeof(IPXHeader);
    ecbReception.fragmentDescriptor[1].address = (BYTE *)
        packetDataRec;
    ecbReception.fragmentDescriptor[1].size = DATA_LENGTH;
    /* deschiderea sociului si intrarea in "ascultare" de pachete */
    if ((code = IPXOpenSocket((BYTE *)&socketRec,0x00)) != 0)
        return 3;
    IPXListenForPacket(&ecbReception);
    IPXRelinquishControl();
    return 0;
}

int GetConnectionNumbers(char users[MAX_CON][USER_LEN])
{
    int i=0; int sort;
    WORD numberOfConnections = 0, saved;
    WORD maxConnections = MAX_CON;
    WORD myConnectionNo;
    BYTE serverNode[6];
    BYTE serverNetwork[4];
    int j;
    for(i=0;i<noNodes;i++)
    {
        do{
            GetObjectConnectionNumbers(users[i], OT_USER,
            &numberOfConnections,&connectionList[i], maxConnections);
            if(numberOfConnections == 0)
                return i;
        }while(numberOfConnections == 0);
    }
    myConnectionNo = connectionList[0];
    do{ /* ordoneaza tabela de conexiuni */
        sort=TRUE;
        for(i=0;i<noNodes-1;i++)
            if(connectionList[i] > connectionList[i+1]){
                saved = connectionList[i];
                connectionList[i]=connectionList[i+1];
                connectionList[i+1]=saved;
                sort=FALSE;
            }
    }while(!sort);
    for(j=0;j<noNodes;j++){
        if(connectionList[j] == myConnectionNo)
            myPosition = j;
    }
    if(myPosition == noNodes -1)
        nextPosition = 0;
    else
        nextPosition = myPosition+1;
    GetInternetAddress(connectionList[nextPosition],(char *)
        serverNetwork,(char *)serverNode,&socket);
    movmem(serverNetwork,
        packetHeaderSnd.destination.network,4);
    movmem(serverNode, packetHeaderSnd.destination.node,6);
    movmem(serverNode, ecbEmission.immediateAddress,6);
    movmem(&socketDest, packetHeaderSnd.destination.socket,2);
    return 0;
}

void ExclAccessEnd(void)

```

B6. Implementarea algoritmului Chandy pentru obținerea excluderii reciproce -varianta Unix

```

{
    IPXCancelEvent(&ecbReception);
    IPXCloseSocket(socketRec);
    IPXCloseSocket(socketSnd);
}
void ExclAccessOpen(void) /* pentru obtinerea accesului in SC */
{
    TokenRequest();
    TokenEnter();
}
void ExclAccessClose(void) /* pentru iesirea din SC */
{
    disable();
    request = FALSE;
    SendToken("OK");
    enter = FALSE;
    enable();
}
void far _loads ReceiveESR(void)
/* apelata la sosirea unui pachet */
{
    disable();
    if (ecbReception.completionCode == 0)
        if (strcmp(packetDataRec, "OK") == 0)
            ReceiveToken();
    }
    enable();
    IPXListenForPacket(&ecbReception);
}
void SendToken(char *message)
{
    movmem(message, packetDataToken, DATA_LENGTH);
    ecbEmission.ESRAddress = NULL;
    ecbEmission.fragmentCount = 2;
    ecbEmission.socketNumber = socketSnd;
    ecbEmission.fragmentDescriptor[0].address =
        &packetHeaderSnd;
    ecbEmission.fragmentDescriptor[0].size = sizeof(IPXHeader);
    ecbEmission.fragmentDescriptor[1].address = (BYTE *)
        packetDataToken;
    ecbEmission.fragmentDescriptor[1].size = DATA_LENGTH;
    while (ecbEmission.inUseFlag != 0) IPXRelinquishControl();
    IPXSendPacket(&ecbEmission);
}

void TokenRequest(void)
{
    request = TRUE;
}

void TokenEnter(void) /* asteptare pina la intrarea in SC */
{
    int enter1;
    do{
        disable();
        enter1=enter;
        enable();
    }while(!enter1);
}

void ReceiveToken(void) /* apelata la primirea token-ului */
{
    int localRequest;
    localRequest = request;
    if(!localRequest)
        SendToken("OK");
    else
        enter = TRUE;
}
}
}

#include <stdio.h>
#include <string.h>
#include <signal.h>
#include <unistd.h>
#include <sys/time.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <arpa/inet.h>
#include <netdb.h>
#include <sys/ipc.h>
#include <sys/sem.h>
#include <fcntl.h>
#include <sys/types.h>
#include <sys/stat.h>
#include <errno.h>
#include <stdarg.h>
#include "Chanduer.h"

#define BROADCAST_ADDR "193.226.37.255"
#define REMOTE_PORT 6560
#define LOCAL_PORT 6561
#define BROADCAST_PORT 6562
#define FIFO_NAME "enterChan"
#define PERMS 0666
#define MAX_HOST_NO 20
#define MAX_HOST_NAME 256
#define MAX_TYPE_MES_LEN 10
#define LEN 400
#define MAX_NO_PACKETS 100
#define MAX_STRING 20
#define REQUEST 1
#define TOKEN 0
#define TRUE 1
#define FALSE 0
#define DEBUG TRUE

/* Functii si variabile interne serverului */
.....
void CreateBroadcastChan(void); /* creare canal broadcast */
int CreateRemoteChan(void); /* creare canal remoteChan */
int CreateLocalChan(void); /* creare canal localChan */
int CreateClientChan(void);
void CreateEnterChan(void); /* creare canal enterChan */;
void CloseEnterChan(void);
void FindAdresse(void); /* determina adresele host-urilor */
void Debug(const char *c, ...);

int enterChan ; /* canal implementat prin FIFO */
int bsd; /* soclu pt. transmisii broadcast */
struct sockaddr_in sad, cad,clia,sva,ICad; /* adrese Inet */
char connTab[MAX_HOST_NO][MAX_HOST_NO];
int processNo; /* numarul de noduri participante la alg.*/
int myID; /* identificatorul host-ului local */

void main(int argc, char *argv[] )
{
    struct sockaddr_in sdr; /* adresa Inet */
    int rChan,lChan,cli; /* identicatorii canalelor */
    fd_set rfd; /* read file descriptor set*/
    fd_set afds; /* active file descriptor set*/
    int sval,ICadl; /* folosite la rcvfrom */
    int inAccess = FALSE; /* indica accesul in SC */
    int tokenOwner = FALSE; /* indica posesia token-ului */
    int knid; /* tipul mesajelor */
    char sBf[LEN],sBf2[LEN]; /* buffere */
}

```

Anexa B: Sincronizare in sisteme distribuite

```

char rBf[LEN],sBf3[LEN],mes[MAX_TYPE_MES_LEN];
int sender; /* folosit pentru receptia mesajelor */
int reqQ[MAX_NO_PACKETS]; /* coada de cereri */
int p = 0; /* nr. de cereri ale host-ului */
int T[MAX_HOST_NO];
int nfd, err, poz , i, j; /* var de lucru */
char *pdest;

/* .....
faza de initializare: se preria din linia de comanda identificatorul
host-ului local si numarul de procese.
Se creeza canalele.
Fiecare server trimite broadcast identificatorul sau si adresa sa
Se pun in asteptare de mesaje canalele remoteChan, localChan
Hostul cu identificatorul cel mai mic primeste token-ul
..... */
if(argc!=3) {
    printf("Usage : [name] [identifer] [hosts] \n");
    exit(0);
}
for(i=0;i<MAX_NO_PACKETS;i++)
    reqQ[i] = 0;
for(i=0;i<MAX_NO_PACKETS;i++)
    T[i] = 0;
myID = atoi(argv[1]);
processNo = atoi(argv[2]);
CreateBroadcastChan();
CreateEnterChan();
getchar();
FindAdresse();
for(i = 0;i<processNo;i++){
    Debug("Host %d of adrese %s \n",i,connTab[i]);
    rChan = CreateRemoteChan();
    lChan = CreateLocalChan();
    cli = CreateClientChan();
    nfd = getdtablesize();
    FD_ZERO(&afds);
    FD_SET(rChan,&afds);
    FD_SET(lChan,&afds);
    if(myID == 0){
        tokenOwner = TRUE;
    }
}
/* .....
faza de tratare a mesajelor primite de la clientul local
sau din retea, de la serverele de pe celelalte host-uri
..... */
while(1){
    Debug("Waiting messages...n");
    bcopy((char *)&afds, (char *)&rfd, sizeof(rfd));
    if(select(nfd,&rfd,(fd_set *)0,(fd_set *)0,(struct timeval *)0)<0)
        Debug("Select error %d \n",errno);
    if(FD_ISSET(rChan,&rfd)){ /* mesaj din retea */
        sval = sizeof(sva);
        err=recvfrom(rChan,rBf,LEN,0,(struct sockaddr*)&sva, &sval);
        if(err<0)
            Debug("eroare la recvfrom \n");
        Debug("Received : %s \n",rBf);
        sscanf(rBf,"%2d-%2d",&knd,&sender);
        Debug("Message type : %d \n",knd);
        Debug("sender : %d \n",sender);
        switch(knd){
            case 0:
                strcpy(mes,"TOKEN");
                break;
            case 1:
                strcpy(mes,"REQUEST");
                break;
            default:
                Debug(" Unknown type\n");
                break;
        }
        Debug("Received from %d mes:= %s.. \n",sender,mes);
        if(strcmp(mes,"REQUEST") == 0){ /* tratare REQUEST rem */
            j = 0;
            sscanf(rBf,"%2d-%2d-%2d",&knd,&sender,&reqQ[sender]);
            Debug("In reqQ[%d] = %d \n",sender,reqQ[sender]);
            if(tokenOwner && (!inAccess)){
                if(reqQ[myID] != 0){
                    inAccess = TRUE;
                    T[myID]=reqQ[myID];
                    Debug("Push in T[%d]=%d",myID,reqQ[myID]);
                    sprintf(sBf,"%s","enter");
                    write(enterChan,sBf,6); /* deblocheaza clientul local */
                }
            }
            else
            do{
                if(reqQ[j]!=0){ /* paseaza token-ul */
                    bzero((char *)&sdr, sizeof(sdr));
                    sdr.sin_family = AF_INET;
                    sdr.sin_addr.s_addr = inet_addr(connTab[j]);
                    sdr.sin_port = htons(REMOTE_PORT);
                    poz= sprintf(sBf,"%2d-%2d-", TOKEN,myID);
                    for(i=0;i<processNo;i++){
                        poz+=sprintf(sBf+poz,"%2d-",T[i]);
                        Debug("Sending token with sBf =:%s\n",sBf);
                        sendto(cli,sBf,LEN,0,(struct sockaddr*)&sdr,sizeof(sdr));
                        tokenOwner = FALSE;
                    }
                    j++;
                }while((j<processNo) && tokenOwner);
            }
            else
            if(strcmp(mes,"TOKEN") == 0){ /* tratare TOKEN remote */
                Debug("ReceivedToken :%s ",rBf);
                pdest = strchr(rBf,' ');
                for(i=0;i<processNo;i++){ /* actualizeaza T */
                    pdest = strchr(pdest+1,' ');
                    sscanf(pdest+1,"%2d",&T[i]);
                }
                for(i=0;i<10;i++){
                    Debug(" T[%d] = %d",i,T[i]);
                    Debug("\n");
                    T[myID] = reqQ[myID];
                    inAccess = TRUE;
                    tokenOwner = TRUE;
                    sprintf(sBf,"%s","enter");
                    write(enterChan,sBf,6); /* deblocheaza clientul local */
                }
            }
        }
        if(FD_ISSET(lChan,&rfd)){ /* mesaj de la clientul local */
            lCadr = sizeof(lCadr);
            err=recvfrom(lChan,rBf,LEN,0,(struct sockaddr*)&lCadr,&lCadr);
            if(err<0)
                Debug("Recvfrom -local- error \n");
            sscanf(rBf,"%s",mes);
            Debug("Received from a local client %s\n",mes);
            if(strcmp(mes,"REQUEST")==0){
                p++; /* contorizeaza o alta cerere */
                for(i=0;i<processNo;i++){ /* trimite cererea in retea */
                    bzero((char *)&sdr, sizeof(sdr));
                    sdr.sin_family = AF_INET;
                    sdr.sin_addr.s_addr = inet_addr(connTab[i]);
                    sdr.sin_port = htons(REMOTE_PORT);
                    sprintf(sBf3,"%2d-%2d-%2d",REQUEST,myID,p);
                    Debug("Sending REQ:%s la %d\n",sBf3,i);
                    sendto(cli,sBf3,LEN,0,(struct sockaddr*)&sdr,sizeof(sdr));
                }
            }
            else
            if(strcmp(mes,"RELEASE") == 0){ /* tratare RELEASE local */
                inAccess = FALSE;
                reqQ[myID] = 0;
                j = 0;
                for(i=0;i<10;i++){
                    Debug(" reqQ[%d] = %d",i,reqQ[i]);
                }
                Debug("\n");
                for(i=0;i<10;i++){
                    Debug(" T[%d] = %d",i,T[i]);
                }
                Debug("\n");
                do{
                    if(reqQ[j]!=0)
                        if(reqQ[j] > T[j]){
                            bzero((char *)&sdr, sizeof(sdr));
                        }
                }
            }
        }
    }
}

```


B7. Implementarea semafoarelor distribuite -varianta Windows

```

Server pentru semafoare distribuite
-Serverul este utilizat pentru controlul accesului la o baza de
date intr-o retea locala folosind semafoare distribuite
(anexa C4)
-Suportul teoretic este prezentat in capitolul 4, la 4.3.4.1
..... */
#include <windows.h>
#include <winsock.h>
#include <stdio.h>
#include <string.h>
#include <process.h>
#include <stdlib.h>
#include <sys/types.h>
#include <fcntl.h>
#include <sys/stat.h>
#include "semwerr.h"

#define BROADCAST_ADDR "193.226.37.255"
#define BROADCAST_PORT 6560
#define LOCAL_PORT 6561
#define MAX_SEM 16 /* nr. maxim de semafoare */
#define MAX_MES_QUEUE 100 /* lung. max. a cozii de mes. */
#define MAX_ACK_QLEN 1000 /* lung. max. a cozii de ack. */
#define LEN 100 /* lungimea max.ima a unui mesaj */
#define TRUE 1
#define FALSE 0
#define DEBUG 1

/* ..... **
 * Functii si variabile interne serverului
 * ..... */
typedef enum { UP, DOWN, ACK } kind; /* actiuni asupra sem. */
typedef struct /* structura unui element din coada de op. necon. */
{
    kind k; /* tipul mesajului */
    int sender; /* identificatorul emitorului */
    int ts; /* marca de timp */
    int sem; /* semaforul */
} QUEUE;
typedef struct /* structura unui element din coada de Ack */
{
    int sender; /* identificatorul emitorului */
    int ts; /* marca de timp */
    int sem; /* semaforul */
} ACKNOWLEDGE;

SOCKET CreateRemoteChan(void);
SOCKET CreateLocalChan(void);
void CreateAnswerChan(void);
void CloseAnswerChan(void);
void InsertQueue(int sender, kind k, int ts, int sem);
void RemoveQueue(int poz);
void InsertAck(int sender, int ts, int sem);
int FindFullyAck(kind k, int sem);
void Debug(const char *c, ...);

HANDLE answerChan; /* canalul answerChan */
QUEUE mesQueue[MAX_MES_QUEUE]; /* coada de op */
ACKNOWLEDGE ackList[MAX_ACK_QLEN];
int nUp[MAX_SEM], nDown[MAX_SEM];
int maxHosts; /* numarul participantilor la algoritmul */
int maxQueue, maxAck; /* dim. curenta a cozilor */
LPSTR lpAnswer;
int ifa;
struct sockaddr_in sar,sal; /* adrese pentru lucrul cu socket */

void main(int argc, char* argv[])
{
    SOCKET rChan,lChan; /* canalele pentru comunicatii */
    struct sockaddr_in sdr;
    int sarl,sall;
    HANDLE hStdout;
    HANDLE hAnswerChan; /* answerChan se implementeaza
    /* cu un semafor si o zona de memorie partajata */
    char sBf[LEN],rBf[LEN]; /* buffere */
    fd_set rfd; /* read file descriptor set */
    fd_set afds; /* active file descriptor set */
    int faUp=0, faDown=0; /* no. mesajelor fully ack. */
    int myid; /* identificatorul host-ului */
    int localTS=0; /* marca de timp locala */
    WSADATA WSAData;
    int status;
    int Wait = FALSE; /* folosit la deblocare clientului */
    int SemWait = -1; /* folosit la deblocare clientului */
    int nfds, err; /* variabile de lucru */
    int nr, sem;
    int ts, sender, char ck; /* componente in mesaje */
    kind k;
    int i;

    /* ..... **
    faza de initializare: se preria din linia de comanda identificatorul
    host-ului local si numarul de procese.
    Se creeaza canalele.
    Se pun in asteptare de mesaje canalele remoteChan, localChan
    ..... */
    FreeConsole();
    AllocConsole();
    SetConsoleTitle("Win32 - Server");
    hStdout=GetStdHandle(STD_OUTPUT_HANDLE);
    if ((status = WSASStartup(MAKEWORD(1,1), &WSAData) == 0)
        Panic(ERROR4);
    else
        Debug("WSA not OK\n");
    if(argc==3) {
        printf("Usage : [name] [identifier] [hosts]\n");
        exit(0);
    }
    myid=atoi(argv[1]);
    maxHosts=atoi(argv[2]);
    hAnswerChan = CreateSemaphore(NULL, 0L, 1L, "semhelp");
    localTS = 0;
    for(i=0;i<MAX_SEM; i++)
        nDown[i]= 0; nUp[i]= 0;
    CreateAnswerChan();
    rChan = CreateRemoteChan();
    lChan = CreateLocalChan();
    FD_ZERO(&afds);
    FD_SET(rChan,&afds);
    FD_SET(lChan,&afds);

    /* ..... **
    faza de tratare a mesajelor primite de la clientul local
    sau din retea, de la serverele de pe celelalte host-uri
    ..... */
    while(TRUE){
        Debug("Waiting messages...\n");
        CopyMemory((char *)&rfd, (char *)&afds, sizeof(rfd));
        if(select(nfds,&rfd,(fd_set *)&0,(fd_set *)&0,(struct timeval *)&0)<0){
            err = WSAGetLastError();
            Debug("Select error : %d",err);
        }
        if(FD_ISSET(rChan,&rfd)){ /* s-a primit un mesaj remote */
            sarl = sizeof(sar);
            err=recvfrom(rChan, rBf,LEN,0,(struct sockaddr *)&sar, &sarl);
            if(err<0)
                Debug("Recvfrom -remote- error\n");
            sscanf(rBf,"%2d-%c-%2d-%10d",&sender,&ck,&sem,&ts);
            switch(ck) {
                case 'V': k = UP; break;
                case 'P': k = DOWN; break;
                case 'A': k = ACK; break;
            }
            Debug("%c for %d from %d host ts=%d...\n",ck,sem,sender,ts);
            localTS = max(localTS, ts+1); /* actualizare TS */
            localTS++;
            if((sender == myid)&&(k == UP)){ /* se memoreaza deblocare */

```

```

SemWait = sem;
Wait = TRUE;
}
if((k == DOWN) || (k == UP)) { /* s-a primit DOWN sau UP */
    InsertQueue(sender, k, ts, sem); /* inserare mesaj in coada */
    memset((void FAR *)&sdr, 0, sizeof(sdr));
    sdr.sin_family = AF_INET; /* si trimite confirmare */
    sdr.sin_addr.s_addr = inet_addr(BROADCAST_ADDR);
    sdr.sin_port = htons(BROADCAST_PORT);
    sprintf(sBf, "%2d-%c-%2d-%10d", myid, 'A', sem, localTS);
    Debug("Sending ACK for %d sem %d host ... \n", sem, myid);
    Debug("Sent message %s \n", sBf);
    sendto(rChan, sBf, LEN, 0, (struct sockaddr*)&sdr, sizeof(sdr));
    localTS++;
}
else
if(k == ACK) { /* s-a primit confirmare de operatie */
    InsertAck(sender, ts, sem); /* inerea za confirmarea */
    ifa=0;
    do { /* cauta daca operatia UP devine complet confirmata */
        faUp = FindFullyAck( UP, sem );
        if(faUp != -1) {
            RemoveQueue(faUp);
            nUp[sem]++;
            if(Wait && (SemWait == sem)){
                Wait = FALSE;
                nr = nUp[sem]-nDown[sem];
                sprintf(lpAnswer, "%d", nr);
                ReleaseSemaphore(hAnswerChan, 1, NULL);
            }
        }
    } while(faUp != -1);
    Debug("Value nUp=%d for %d sem. \n", nUp[sem], sem);
    ifa = 0;
    do { /* cauta daca operatia DOWN devine complet confirmata */
        faDown = FindFullyAck( DOWN, sem );
        Debug(" nUp=%d ; nDown=%d \n", nUp[sem], nDown[sem]);
        if(( !faDown != -1) && (nUp[sem] > nDown[sem])) {
            ifa--;
            sender=mesQueue[faDown].sender;
            RemoveQueue( faDown );
            nDown[sem]++;
            if(sender == myid) {
                nr = nUp[sem]-nDown[sem];
                sprintf(lpAnswer, "%d", nr); /* Deblochez la semul
local sem */
                Debug("Unblocking my client who's waiting at DOWN \n");
                ReleaseSemaphore(hAnswerChan, 1, NULL);
            }
        }
    } while(faDown != -1);
    Debug("Value nDown=%d for %d sem. \n ", nDown[sem], sem);
}
}

if(FD_ISSET(lChan, &rfd)) /* s-a primit un mesaj local */
    sall = sizeof(sal);
err = recvfrom(lChan, rBf, LEN, 0, (struct sockaddr*)&sal, &sall);
if(err<0){
    Debug("Recvfrom -local- error\n");
}
sscanf(rBf, "%2d-%c", &sem, &ck);
switch(ck) {
    case 'V': k = UP; break;
    case 'P': k = DOWN; break;
}
Debug("Received from client %c type for %d sem \n", ck, sem);
if((k== DOWN) || (k==UP)){ /* trimite mes. in retea */
    memset((void FAR *)&sdr, 0, sizeof(sdr));
    sdr.sin_family = AF_INET;
    sdr.sin_addr.s_addr = inet_addr(BROADCAST_ADDR);
    sdr.sin_port = htons(BROADCAST_PORT);
    if(k==DOWN)
        sprintf(sBf, "%2d-%c-%2d-%10d", myid, 'P', sem, localTS);
    else
        sprintf(sBf, "%2d-%c-%2d-%10d", myid, 'V', sem, localTS);
    Debug("Sent message: %s \n", sBf);
    sendto(rChan, sBf, LEN, 0, (struct sockaddr*)&sdr, sizeof(sdr));
    localTS++;
}
}
CloseAnswerChan();
}
/* .....
* ..... Functii folosite de server .....
* ..... */
void InsertQueue(int sender, kind k, int ts, int sem)
/* InsertQueue - insereaza un element in coada de op. */
/* astfel incat coada sa ramana ordonata */
{
    int i;
    if((k==UP) || (k==DOWN)) {
        if(maxQueue == 0) {
            mesQueue[0].sender = sender;
            mesQueue[0].k = k;
            mesQueue[0].ts = ts;
            mesQueue[0].sem = sem;
            maxQueue = 1;
        }
        else {
            i = 0;
            while((i < maxQueue) && (ts >= mesQueue[i].ts))
                i++;
            for(j = maxQueue-1; j >= i; j--)
                mesQueue[j+1] = mesQueue[j];
            maxQueue++;
            mesQueue[i].sender = sender;
            mesQueue[i].k = k;
            mesQueue[i].ts = ts;
            mesQueue[i].sem = sem;
        }
    }
}
void InsertAck(int sender, int ts, int sem)
/* insereaza un element in coada de confirmari */
{
    ackList[maxAck].sender = sender;
    ackList[maxAck].ts = ts;
    ackList[maxAck].sem = sem;
    maxAck++;
}
int FindFullyAck(kind k, int sem)
/* FindFullyAck - returneaza operatiile de tipul complet */
/* confirmate din coada de operatii */
{
    int i, j, t;
    for(i=ifa; i<maxQueue; i++) {
        if((mesQueue[i].sem==sem) && (mesQueue[i].k==k)) {
            for(j=0; j<maxHosts; j++)
                for(t=0; t<maxAck; t++)
                    if(( ackList[t].sender == j) && (ackList[t].ts>mesQueue[i].ts) &&
(ackList[t].sem == sem ))
                        break;
                    if(t==maxAck)
                        break;
                }
            if(j==maxHosts) {
                if(k==UP) ifa=0;
                else ifa=i+1;
                return i;
            }
        }
    }
    return -1;
}
void RemoveQueue(int poz)
/* RemoveQueue - sterge elementul din pozitia poz din coada */

```

```

{
int i;
for(i = poz; i < maxQueue; i++)
    mesQueue[i] = mesQueue[i+1];
maxQueue--;
}
SOCKET CreateRemoteChan(void)
/* creaza canalul pentru transmisii broadcast */
{
SOCKET s;
if((s = socket(AF_INET, SOCK_DGRAM, 0)) < 0)
    Panic(ERROR1);
memset((void FAR *)&sar, 0, sizeof(sar));
sar.sin_family = AF_INET;
sar.sin_addr.s_addr = htonl(INADDR_ANY);
sar.sin_port = htons(BROADCAST_PORT);
if(bind(s, (struct sockaddr *)&sar, sizeof(sar)) < 0)
    Panic(ERROR2);
return s;
}

SOCKET CreateLocalChan(void)
/* creaza canalul pentru receptia mesajelor de la clientul local */
{
SOCKET s;
if((s = socket(AF_INET, SOCK_DGRAM, 0)) < 0)
    Panic(ERROR3);
memset((void FAR *)&sal, 0, sizeof(sal));
sal.sin_family = AF_INET;
sal.sin_addr.s_addr = htonl(INADDR_ANY);
sal.sin_port = htons(LOCAL_PORT);
if(bind(s, (struct sockaddr *)&sal, sizeof(sal)) < 0)
    Panic(ERROR4);
return s;
}

void CreateAnswerChan(void)
/* creaza zona de memorie partajata prin care se */
/* comunica cu clientul */
{
answerChan = CreateFileMapping((HANDLE)0xFFFFFFFF,
    NULL, PAGE_READWRITE, 0, 0x1000, "SEMDEMO");
if(answerChan == NULL)
    Panic(ERROR5);
lpAnswer = (LPSTR)MapViewOfFile(answerChan,
    FILE_MAP_WRITE, 0, 0, 0);
if(lpAnswer == NULL)
    Panic(ERROR6);
}

void CloseAnswerChan(void)
{
UnmapViewOfFile(lpAnswer);
}

void Debug(const char *c, ...){
va_list args;
va_start(args,c);
vprintf(c,args);
va_end(args);
}

void Panic( int err){
printf("Exit: %d error \n", err);
exit(0);
}
}

/* .....
Funcțiile puse la dispozitia unui client
care foloseste semafoarele distribuite - varianta Windows

-Suportul teoretic este prezentat in capitolul 4, la 4.3.4.1
-Functiile care pot fi folosite de un client sunt:
int SemaphoreInitialize(void)
int SemaphoreUP(int sid)
int SemaphoreDown(int sid)
void SemaphoreEnd(void)

- Comunicarea cu serverul se face printr-o zona de memorie
partajata iar sincronizarea server - client se realizeaza cu
ajutorul unui semafor

.....*/
#include <winsock.h>
#include <windows.h>
#include "resource.h"
#include "semaforw.h"

int OpenLocalChan(void);
int CreateLocalChan(void);
int GetServerAddress(void);

LPSTR lpAnswer;
int lChan;
struct sockaddr_in sad;

int SemaphoreDown(int sid)
/* Semaphore_Down- operatia DOWN asupra sem.sid */
{
char info[LEN];
int nr;
if ( sid < 0 || sid >MAX_SEM-1)
    return -1;
sprintf(info,"%2d-%c",sid, 'P');
sendto(lChan,info, LEN, 0,
    (struct sockaddr *)&sad, sizeof(sad));

/* blocare pina la primirea raspunsului de la server */
WaitForSingleObject(hAnswerChan,INFINITE);

/* urmeaza citirea din zona de memorie partajata */
/* a valorii semaforului */
sscanf(lpAnswer,"%d",&nr);
return nr;
}

int SemaphoreUP(int sid)
/* SemaphoreUp- operatiaUP asupra sem. sid */
{
char info[LEN];
int nr;
if ( sid < 0 || sid >MAX_SEM-1)
    return -1;
sprintf(info,"%2d-%c",sid, 'V');
sendto(lChan,info, LEN, 0,
    (struct sockaddr *)&sad,sizeof(sad));

/* blocare pina la primirea raspunsului de la server */
WaitForSingleObject(hAnswerChan,INFINITE);
sscanf(lpAnswer,"%d",&nr);
return nr;
}

int SemaphoreInitialize(void)
/* pentru conectare la server */
{
int errCode;
if ((errCode = CreateLocalChan()) != 0)
    return errCode;
if ((errCode = OpenLocalChan()) != 0)
    return errCode;
return 0;
}

char * GetMyAddress(void)

```

B8. Implementarea contorilor de eveniment distribuiți -varianta UNIX

```

{
    int errCode;
    if ((errCode = GetServerAddress()) != 0)
        return NULL;
    else
        return MyServerAddress;
}

int GetServerAddress(void)
/*determina adresa serverului propriu*/
{
    struct in_addr adr;
    char HostName[MAX_HOST_NAME];
    struct hostent *host;
    if (gethostname(HostName,MAX_HOST_NAME) != 0)
        return ERROR2;
    if((host = gethostbyname(HostName)) == NULL)
        return ERROR3;
    if(host != NULL)
        adr.s_addr = ((struct in_addr *) (host->h_addr_list[0]))->s_addr;
    if(strcmp(MyServerAddress,inet_ntoa(adr)) == NULL)
        return ERROR4;
    return 0;
}

int OpenLocalChan(void)
/* deschide semaforul si socul pentru comunicare
   cu serverul ; daca nu reuseste, intoarce un cod de eroare
*/
{
    int errCode; /*cod eroare */
    hAnswerChan = OpenSemaphore(SYNCHRONIZE,TRUE,
    "semhelp");
    if (hAnswerChan == NULL)
        return ERROR1;
    if ((errCode=GetServerAddress()) != 0)
        return errCode;
    memset((void FAR *)&sad, 0, sizeof(sad)) ;
    sad.sin_family = AF_INET ;
    sad.sin_addr.s_addr = inet_addr(MyServerAddress) ;
    sad.sin_port = htons(LOCAL_PORT) ;
    if((lChan = socket(AF_INET, SOCK_DGRAM, 0)) < 0)
        return ERROR5;
    return 0;
}

int CreateLocalChan(void)
/* obtine accesul la zona de memorie partajata */
{
    HANDLE hmmf; /* handler pentru mem. partajata */
    hmmf = CreateFileMapping((HANDLE)0xFFFFFFFF,
    NULL, PAGE_READWRITE, 0, 0x2000, "SEMDEMO");
    if(hmmf == NULL)
        return ERROR6;
    lpAnswer = (LPSTR)MapViewOfFile(hmmf, FILE_MAP_WRITE,
    0, 0, 0);
    if(lpAnswer == NULL)
        return ERROR7;
    return 0;
}

void SemaphoreEnd(void)
{
    closesocket(lChan);
    UnmapViewOfFile(lpAnswer);
}

}

/* *****
   Server pentru contori de eveniment distribuiti
   -Serverul este folosit in anexa B10 pentru implementarea
   excluderii reciproce cu ajutorul contorilor de evenimente distribuiti
   si a secventiatorilor distribuiti
   - Suportul teoretic este prezentat in capitolul 4, la 4.3.4.2-
   ***** */
#include <stdio.h>
#include <string.h>
#include <signal.h>
#include <unistd.h>
#include <sys/time.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <arpa/inet.h>
#include <netdb.h>
#include <sys/pc.h>
#include <sys/sem.h>
#include <fcntl.h>
#include <sys/stat.h>
#include <errno.h>
#include <stdarg.h>
#include "countuer.h"

#define BROADCAST_ADDR "193.226.37.255"
#define BROADCAST_PORT 6560
#define LOCAL_PORT 6561
#define FIFO_NAME "answerChan"
#define PERMS 0666
#define MAX_COUNTER 10 /* nr. max de contori */
#define MAX_MES_QUEUE 100 /* dim max a cozii de mesaje */
#define MAX_ACK_QLEN 1000 /* dim max a listei de confir.*/
#define ACKN 4
#define ADVANCE 1
#define LEN 400
#define TRUE 1
#define FALSE 0
#define DEBUG TRUE

/* *****
   * Functii si variabile interne serverului
   ***** */

typedef enum { ADV, READ, AW, ACK } kind ;
typedef struct /* structura unui element din coada de op. necon.*/
    int sender; /* identificatorul emitatorului */
    int ts; /* marca de timp */
    int counter; /* semaforul */
}QUEUE;
typedef struct { /* structura unui element din coada de Ack*/
    int sender; /* identificatorul emitatorului */
    int ts; /* marca de timp */
    int counter; /* contorul */
}ACKNOWLEDGE;

int CreateRemoteChan(void);
int CreateLocalChan(void);
void CreateAnswerChan(void);
void CloseAnswerChan(void);
void InsertQueue(int sender, int ts, int counter);
void RemoveQueue(int poz);
void InsertAck(int sender, int ts, int counter);
int FindFullyAck(int sender, int counter);
void Debug(const char *c, ...);

int answerChan ; /* handler pipe */
QUEUE mesQueue[MAX_MES_QUEUE];
ACKNOWLEDGE ackList[MAX_ACK_QLEN];
int maxQueue, maxAck;
int maxHosts;
struct sockaddr_in sar,sal;

```

Anexa B: Sincronizare in sisteme distribuite

```

void main( int argc, char *argv[])
{
    struct sockaddr_in sdr;
    int sari,sall;
    int rChan,iChan;
    fd_set rfds;          /* read file descriptor set*/
    fd_set afds;         /* active file descriptor set*/
    int nE[MAX_COUNTER], nEVal[MAX_COUNTER];
    char sBuf[LEN],rBuf[LEN]; /* buffere */
    int localTs=0;       /* marca de timp locala */
    int myID;           /* identificatorul host-ului */
    int sender, ts;     /* componente in mesaje */
    int poz;           /* variabile de lucru */
    int cr, counter;
    int val=0;
    int i, nfds, err;
    kind k;
    int kd;

    /* .....
faza de initializare: se preria din linia de comanda identificatorul
host-ului local si numarul de procese.
Se creeaza canalele.
Se pun in asteptare de mesaje canalele remoteChan, localChan
..... */
    if(argc!=3) {
        Debug("Usage : [nume] [nr.identif.] [nr.total de procese] \n");
        exit(0);
    }
    for(i=0;i<MAX_COUNTER;i++)
        nEVal[i]=-1;
    CreateAnswerChan();
    myID = atoi(argv[1]);
    maxHosts = atoi(argv[2]);
    rChan = CreateRemoteChan();
    lChan = CreateLocalChan();
    nfds = getdtablesize();
    FD_ZERO(&afds);
    FD_SET(rChan,&afds);
    FD_SET(lChan,&afds);

    /* .....
faza de tratare a mesajelor primite de la clientul local
sau din retea, de la serverele de pe celelalte host-uri
..... */
    while(TRUE){
        Debug("Waiting messages...\n");
        bcopy((char *)&afds, (char *)&rfds, sizeof(rfds));
        if(select(nfds,&rfds,(fd_set *)0,(fd_set *)0,(struct timeval *)0)<0)
            Debug("Select error %d \n",errno);
        if(FD_ISSET(rChan,&rfds)){ /* mesaj din retea */
            sari = sizeof(sar);
            err=recvfrom(rChan,rBuf,LEN,0,(struct sockaddr*)&sar,&sari);
            if(err<=0){
                Debug("Recvfrom -remote- error%d\n",err);
            }
            sscanf(rBuf,"%2d-%2d-%6d-%2d",&kd,&sender,&ts,&cr);
            Debug("From %d %d ts = %d counter %d...\n",sender,kd,ts,cr);
            switch(kd){
                case 1: k = ADV;
                    break;
                case 4: k = ACK;
                    break;
            }
            localTs=max(localTs,ts+1); /* actualizare marca de timp */
            localTs++;
            if(k==ADV){ /* tratare advance */
                InsertQueue(sender,ts,cr); /* insereaza mesajul in coada */
                bzero((char *)&sdr, sizeof(sdr));
                sdr.sin_family = AF_INET;
                sdr.sin_addr.s_addr = inet_addr(BROADCAST_ADDR);
                sdr.sin_port = htons(BROADCAST_PORT);
                sprintf(sBuf,"%2d-%2d-%6d-%2d",ACKN,myID,localTs,cr);
                Debug("Sending:%s\n",sBuf); /* trimite confirmarea */
                sendto(rChan,sBuf,LEN,0,(struct sockaddr*)&sdr,sizeof(sdr));
            }
            localTs++;
        }
        else
            if(k==ACK){ /* tratare ACK */
                InsertAck(sender,ts,cr);
                for(i=0;i<maxAck;i++)
                    Debug("ackList = %d ",ackList[i].ts);
                Debug("\n");
                for(i=0;i<maxQueue;i++)
                    Debug("mesQueue = %d ",mesQueue[i].ts);
                Debug("\n");
                poz = FindFullyAck(sender, cr);
                if(poz>=0){
                    counter = mesQueue[poz].counter;
                    nE[counter]++;
                    RemoveQueue(poz);/*poz = pozitia celui scos din coada*/
                    Debug("nE[%d]= %d \n",counter,nE[counter]);
                    if((nE[counter]>= nEVal[counter]) && (nEVal[counter] != -1)){
                        Debug("nEVal = %d\n",nEVal[counter]);
                        nEVal[counter] = -1;
                        sprintf(sBuf,"%s","gata");
                        write(answerChan,sBuf,6);
                    }
                }
            }
        }
        if(FD_ISSET(lChan,&rfds)){ /* mesaj de la clientul local */
            sall = sizeof(sal);
            err=recvfrom(lChan,rBuf,LEN, 0, (struct sockaddr*)&sal, &sall);
            if(err<=0)
                Debug("Recvfrom - local -error\n");
            sscanf(rBuf,"%2d",&kd);
            switch(kd){
                case 2: k=AW;
                    break;
                case 3: k = READ;
                    break;
                case 1:k=ADV;
                    break;
            }
            if(k==ADV){ /* tratare advance */
                sscanf(rBuf,"%2d-%2d",&kd,&cr);
                bzero((char *)&sdr, sizeof(sdr));
                sdr.sin_family = AF_INET;
                sdr.sin_addr.s_addr = inet_addr(BROADCAST_ADDR);
                sdr.sin_port = htons(BROADCAST_PORT);
                sprintf(sBuf,"%2d-%2d-%6d-%2d",ADVANCE,myID,localTs,cr);
                Debug("Sending:%s\n",sBuf);
                sendto(rChan,sBuf,LEN,0,(struct sockaddr*)&sdr,sizeof(sdr));
                localTs++;
            }
            if(k==READ){ /* tratare read */
                sscanf(rBuf,"%2d-%2d",&kd,&cr);
                Debug("Received %d and counter = %d...\n",kd,cr);
                Debug("nE[%d] pe READ local este : %d \n",cr,nE[cr]);
                sprintf(sBuf,"%d",nE[cr]);
                write(answerChan,sBuf,6);
            }
        }
        else
            if(k==AW){ /* tratare await */
                sscanf(rBuf,"%2d-%2d-%2d",&kd,&cr,&val);
                Debug(" %d with counter = %d and val=%d...\n",kd,cr,val);
                nEVal[cr] = val;
                Debug("nEVal=%d\n",nEVal[cr]);
                if(nE[cr] >= nEVal[cr]){
                    nEVal[cr] = -1;
                    sprintf(sBuf,"%s","gata");
                    write(answerChan,sBuf,6);
                }
            }
        }
    }
    /* .....
Functii folosite de server
..... */
}

```

```

void InsertQueue(int sender, int ts, int counter)
/* insereaza un element in coada */
/* astfel incat coada sa ramana ordonata */
{
    int i,j;
    if(maxQueue == 0) {
        mesQueue[0].sender = sender;
        mesQueue[0].ts = ts;
        mesQueue[0].counter = counter;
        maxQueue = 1;
    }
    else {
        i = 0;
        while((i < maxQueue) && (ts >= mesQueue[i].ts))
            i++;
        for(j = maxQueue-1; j >= i; j--)
            mesQueue[j+1] = mesQueue[j];
        maxQueue++;
        mesQueue[j].sender = sender;
        mesQueue[j].ts = ts;
        mesQueue[j].counter = counter;
    }
}

void InsertAck(int sender, int ts, int counter)
/* insereaza in coada de Ack o confirmare */
{
    ackList[maxAck].sender = sender;
    ackList[maxAck].ts = ts;
    ackList[maxAck].counter = counter;
    maxAck++;
}

int FindFullyAck(int sender, int counter)
/* FindFullyAck - returneaza operatiile de tipul complet */
/* confirmate din coada de operatii */
{
    int i,j,t;
    for(i=0;i<maxQueue;i++){
        if(mesQueue[i].counter == counter)
            for(j=0;j<maxHosts;j++){
                for(t=0;t<maxAck;t++)
                    if((ackList[t].sender == j) && (ackList[t].ts>mesQueue[i].ts)
                        && (ackList[t].counter == counter))
                        break;
                if(t==maxAck)
                    break;
            }
            if(j==maxHosts)
                return i;
    }
    return -1;
}

void RemoveQueue(int poz)
/* sterge elementul din pozitia poz din coada */
{
    int i;
    for(i = poz; i < maxQueue; i++)
        mesQueue[i] = mesQueue[i+1];
    maxQueue--;
}

int max(int a,int b)
{
    if(a > b)
        return a;
    else
        return b;
}

int CreateRemoteChan(void)
{
    int s, i;
    if(s = socket(AF_INET, SOCK_DGRAM, 0) < 0)
        Panic(ERROR1);
    i = 1;
}

if(setsockopt(s, SOL_SOCKET, SO_BROADCAST, (char *) &i,
sizeof(i)))
    Panic(ERROR7);
bzero((char *)&sar, sizeof(sar));
sar.sin_family = AF_INET;
sar.sin_addr.s_addr = htonl(INADDR_ANY);
sar.sin_port = htons(BROADCAST_PORT);
if(bind(s, (struct sockaddr *)&sar, sizeof(sar)) < 0)
    Panic(ERROR2);
return s;
}

int CreateLocalChan(void)
{
    int s;
    if(s = socket(AF_INET, SOCK_DGRAM, 0) < 0)
        Panic(ERROR3);
    bzero((char *)&sal, sizeof(sal));
    sal.sin_family = AF_INET;
    sal.sin_addr.s_addr = htonl(INADDR_ANY);
    sal.sin_port = htons(LOCAL_PORT);
    if(bind(s, (struct sockaddr *)&sal, sizeof(sal)) < 0)
        Panic(ERROR4);
    return s;
}

void CreateAnswerChan(void)
/* creeaza si deschide answerChan */
{
    unlink(FIFO_NAME);
    if( mknod(FIFO_NAME, S_IFIFO | PERMS, 0) < 0)
        Panic(ERROR5);
    if((answerChan = open(FIFO_NAME, O_RDWR)) < 0)
        Panic(ERROR6);
}

void CloseAnswerChan(void)
/* sterge pipe */
{
    unlink(FIFO_NAME);
}

void Debug(const char *c, ...)
{
    va_list args;
    va_start(args,c);
    vprintf(c, args);
    va_end(args);
}

void Panic( int err){
    printf("Exit: %d error \n", err);
    exit(0);
}
/* .....
    Functiile puse la dispozitia unui client
    care foloseste contori de eveniment distribuiti - varianta Unix
    -Suportul teoretic este prezentat in capitolul 4, la 4.3.4.2
    -Functiile care pot fi folosite sunt de un client :

int EventCountersAdvance(int contor)
int EventCountersAwait(int contor, int val)
int EventCountersRead(int contor)
int EventCountersInitialize(void)
void EventCountersEnd(void)
- Comunicarea cu serverul se face printr-un pipe FIFO
    .....

#include "counter.h"
int IChan;
struct sockaddr_in sad;
int answerChan; /* handler answerChan */
int GetServerAddress(void);
int OpenAnswerChan(void);
int EventCountersAdvance(int counter) /* Op. ADVANCED */
{
    char info[LEN];
}

```



```

int err;
if (counter > MAX_CONTOR-1 || counter < 0)
    return -1;
sprintf(info,"%2d-%2d",ADVANCE,counter);
sendto(lChan,info,LEN,0,(struct sockaddr *)&sad,sizeof(sad));
return 0;
}
int EventCountersAwait(int counter, int val) /*Op. AWAIT*/
{
char info[LEN];
if (counter > MAX_CONTOR-1 || counter < 0)
    return -1;
sprintf(info,"%2d-%2d-%2d",AWAIT,counter,val);
sendto(lChan,info,LEN,0,(struct sockaddr *)&sad, sizeof(sad));
/* Asteptare deblocare */
read(answerChan,info,6);
return 0;
}
int EventCountersRead(int counter) /* OP. READ */
{
char info[200];
int val;
if (counter > MAX_CONTOR-1 || counter < 0)
    return -1;
sprintf(info,"%2d-%2d",READ,counter);
sendto(lChan,info, LEN, 0,(struct sockaddr *) &sad,sizeof(sad));
read(answerChan,info,6);
sscanf(info,"%d",&val);
return val;
}
int EventCountersInitialize(void)/*realizeaza conectarea la server*/
{
int errCode;
if ((errCode=OpenAnswerChan()) != 0)
    return errCode;
if ((errCode= GetServerAddress()) != 0)
    return errCode;
memset(&sad, 0, sizeof(sad));
sad.sin_family = AF_INET;
sad.sin_addr.s_addr = inet_addr(MyServerAddress);
sad.sin_port = htons(LOCAL_PORT);
if(!lChan= socket(AF_INET, SOCK_DGRAM, 0) < 0)
    return ERRORS;
return 0;
}
void EventCountersEnd(void){ /* deconectare de la server */
close(lChan);
unlink(FIFO_NAME);
}
int GetServerAddress(void) /*determina adr,serverului propriu*/
{
struct in_addr adr;
char HostName[MAX_HOST_NAME];
struct hostent *host;
if (gethostname(HostName,MAX_HOST_NAME) != 0)
    return ERROR2;
if((host = gethostbyname(HostName)) == NULL)
    return ERROR3;
if(host != NULL)
    adr.s_addr = ((struct in_addr *) (host->h_addr_list[0]))->s_addr;
if(strcmp(MyServerAddress,inet_ntoa(adr)) == NULL)
    return ERROR4;
return 0;
}
char * GetMyAddress(void){
int errCode;
if ((errCode = GetServerAddress()) != 0)
    return NULL;
else
    return MyServerAddress;
}
int OpenAnswerChan(void) /* deschide answerChan */
{
if((answerChan= open(FIFO_NAME, O_RDWR))<0)
    return ERROR1;
return 0;
}

```

B9. Implementarea secventiatorilor distribuiji -varianta UNIX

```

/* .....
Server pentru secventiatori distribuiti
-Serverul este folosit in anexa B10 pentru implementarea
excluderii reciproce cu ajutorul contorilor de evenimente distribuiti
si a secventiatorilor distribuiti
- Suportul teoretic este prezentat in capitolul 4, la 4.3.4.3-
.....
*/
#include <stdio.h>
#include <string.h>
#include <signal.h>
#include <unistd.h>
#include <sys/time.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <arpa/inet.h>
#include <netdb.h>
#include <sys/ipc.h>
#include <sys/sem.h>
#include <fcntl.h>
#include <sys/stat.h>
#include <errno.h>
#include <stdarg.h>
#include sequer.h

#define BROADCAST_ADDR "193.226.37.255"
#define BROADCAST_PORT 6570
#define LOCAL_PORT 6571
#define FIFO_NAME "answerChan"
#define PERMS 0666
#define MAX_MES_QUEUE 100 /* dim max a cozii de mesaje */
#define MAX_ACK_QLEN 1000 /* dim max a listei de confir.*/
#define ACKN 4
#define TICKET 5
#define LEN 100
#define TRUE 1
#define FALSE 0
#define DEBUG TRUE

/* .....
* Functii si variabile interne serverului
.....
*/
typedef enum { TICK, ACK } kind;
typedef struct /* struct. unui element din coada de op. neconf.*/
{
int sender; /* identificatorul emitorului */
int ts; /* marca de timp */
} QUEUE;
typedef struct{ /* structura unui element din coada de Ack*/
int sender; /* identificatorul emitorului */
int ts; /* marca de timp */
} ACKNOWLEDGE;int CreateRemoteChan(void);

int CreateLocalChan(void);
void CreateAnswerChan(void);
void CloseAnswerChan(void);
void InsertQueue(int sender, int ts);
void RemoveQueue(int poz);
void InsertAck(int sender, int ts);
int FindFullyAck(int sender);
void Debug(const char *c, ...);
int max(int a, int b);

int answerChan;
QUEUE mesQueue[MAX_MES_QUEUE];
ACKNOWLEDGE ACKList[MAX_ACK_QLEN];
struct sockaddr_in sar,sal;
int sari, salli;
int maxQueue, maxAck;
int maxHosts;
int myID;
int maxACK;

```

```

int poz;
void main(argc,argv)
int argc;
char *argv[];
{
    int rChan,lChan;
    struct sockaddr_in sdr;
    fd_set rfds; /* read file descriptor set*/
    fd_set afds; /* active file descriptor set*/
    char sBuf[LEN],rBuf[LEN]; /* buffere */
    int localTS=0; /* marca de timp locala */
    int nfds, err; /* variabile de lucru */
    int nT=0;
    int typeMes, ts, sender, recSender; /* componente in mesaje */
    kind k;

    /* .....
faza de initializare: se preria din linia de comanda identificatorul
host-ului local si numarul de procese.
Se creaza canalele.
Se pun in asteptare de mesaje canalele remoteChan, localChan
..... */

if(argc!=3) {
    Debug("Usage : [name] [identifier] [hosts] \n");
    exit(0);
}
CreateAnswerChan();
myID = atoi(argv[1]);
maxHosts = atoi(argv[2]);
rChan = CreateRemoteChan();
lChan = CreateLocalChan();
nfds = getdtablesize();
FD_ZERO(&afds);
FD_SET(rChan,&afds);
FD_SET(lChan,&afds);

/* .....
faza de tratare a mesajelor primite de la clientul local
sau din retea, de la serverele de pe celelalte host-uri
..... */

while(TRUE){
    Debug("Waiting messages...\n");
    bcopy((char *)&afds, (char *)&rfds, sizeof(rfds));
    if(select(nfds,&rfds,(fd_set *)0,(fd_set *)0,(struct timeval *)0)<0)
        Debug("Select error \n");
    if(FD_ISSET(rChan,&rfds)){ /* mesaj din retea */
        sarl = sizeof(sar);
        err=recvfrom(rChan,rBuf,LEN,0,(struct sockaddr*)&sar,&sarl);
        if(err<0)
            Debug("Recvfrom - remote- error\n");
        scanf(rBuf,"%2d-%2d-%6d",&typeMes,&recSender,&ts);
        switch(typeMes){
            case 4: k=ACK;
                    break;
            case 5: k=TICK;
                    break;
        }
        localTS=max(localTS,ts+1); /* actualizare marca de timp */
        localTS++;
    }
    if(k==TICK){ /* tratare TICK */
        InsertQueue(recSender,ts);
        memset(&sdr, 0, sizeof(sdr));
        sdr.sin_family = AF_INET;
        sdr.sin_addr.s_addr = inet_addr(BROADCAST_ADDR);
        sdr.sin_port = htons(BROADCAST_PORT);
        sprintf(sBuf,"%2d-%2d-%6d",ACKN,myID,localTS);
        Debug("Sending:%s\n",sBuf); /* trimite confirmarea */
        sendto(rChan,sBuf,LEN,0,(struct sockaddr*)&sdr,sizeof(sdr));
        localTS++;
    }
    else
    if(k==ACK){ /* tratare ACK */
        InsertAck(recSender,ts);
        poz = FindFullyAck(recSender);
        if(poz>=0){
            nT++;
            sender = mesQueue[poz].sender;
            RemoveQueue(poz);
            if(sender == myID) {
                sprintf(sBuf,"%d",nT);
                write(answerChan,sBuf,6);
            }
        }
    }
    if(FD_ISSET(lChan,&rfds)){ /* mesaj de la clientul local */
        sall = sizeof(sal);
        err =recvfrom(lChan,rBuf,LEN,0,(struct sockaddr*)&sall,&sall);
        if(err<=0)
            Debug("Recvfrom - local -error\n");
        sscanf(rBuf,"%2d",&typeMes);
        switch(typeMes){
            case 5: k=TICK;
                    break;
        }
    }
    if(k==TICK){ /* tratare TICK */
        memset(&sdr, 0, sizeof(sdr));
        sdr.sin_family = AF_INET;
        sdr.sin_addr.s_addr = inet_addr(BROADCAST_ADDR);
        sdr.sin_port = htons(BROADCAST_PORT);
        sprintf(sBuf,"%2d-%2d-%6d",TICKET,myID,localTS);
        Debug("Sending:%s\n",sBuf);
        sendto(rChan,sBuf,LEN,0,(struct sockaddr*)&sdr,sizeof(sdr));
        localTS++;
    }
}

void InsertQueue(int sender, int ts)
/* insereaza un element in coada astfel incat coada */
/* sa ramana ordonata */
{
    int i,j;
    if(maxQueue == 0) {
        mesQueue[0].sender = sender;
        mesQueue[0].ts = ts;
        maxQueue = 1;
    }
    else {
        i = 0;
        while((i < maxQueue) && (ts >= mesQueue[i].ts))
            i++;
        for(j = maxQueue-1; j >= i; j--)
            mesQueue[j+1] = mesQueue[j];
        maxQueue++;
        mesQueue[i].sender = sender;
        mesQueue[i].ts = ts;
    }
}

void InsertAck(int sender, int ts)
{
    ACKList[maxACK].sender = sender;
    ACKList[maxACK].ts = ts;
    maxACK++;
}

int FindFullyAck(int sender)
/* returneaza inregistrările de tipul */
/* fully acknowledgement din coada */
{
    int i,j,t;
    for(i=0;i<maxQueue;i++){
        for(j=0;j<maxHosts;j++){
            for(t=0;t<maxACK;t++){
                if((ACKList[t].sender==i)&&(ACKList[t].ts > mesQueue[i].ts))
                    break;
                if(t==maxACK)
                    break;
            }
        }
        if(j==maxHosts)
            return i;
    }
    return -1;
}

```

```

void RemoveQueue(int poz)
/* sterge elementul din pozitia poz din coada */
{
    int i;
    for(i = poz; i < maxQueue; i++)
        mesQueue[i] = mesQueue[i+1];
    maxQueue--;
}

int max(int a, int b)
{
    if(a>=b)
        return a;
    else
        return b;
}

int CreateRemoteChan(void)
{
    int s, i;
    if((s = socket(AF_INET, SOCK_DGRAM, 0)) < 0)
        Panic(ERROR1);
    i = 1;
    if(setsockopt(s, SOL_SOCKET, SO_BROADCAST, (char *)&i,
        sizeof(i)))
        Panic(ERROR7);
    bzero((char *)&sar, sizeof(sar));
    sar.sin_family = AF_INET;
    sar.sin_addr.s_addr = htonl(INADDR_ANY);
    sar.sin_port = htons(BROADCAST_PORT);
    if(bind(s, (struct sockaddr *)&sar, sizeof(sar)) < 0)
        Panic(ERROR2);
    return s;
}

int CreateLocalChan(void)
{
    int s;
    if((s = socket(AF_INET, SOCK_DGRAM, 0)) < 0)
        Panic(ERROR3);
    bzero((char *)&sal, sizeof(sal));
    sal.sin_family = AF_INET;
    sal.sin_addr.s_addr = htonl(INADDR_ANY);
    sal.sin_port = htons(LOCAL_PORT);
    if(bind(s, (struct sockaddr *)&sal, sizeof(sal)) < 0)
        Panic(ERROR4);
    return s;
}

void CreateAnswerChan(void)
/* creeaza si deschide answerChan */
{
    unlink(FIFO_NAME);
    if(mknod(FIFO_NAME, S_IFIFO | PERMS, 0) < 0)
        Panic(ERROR5);
    if((answerChan = open(FIFO_NAME, O_RDWR)) < 0)
        Panic(ERROR6);
}

void CloseAnswerChan(void)
/* sterge pipe */
{
    unlink(FIFO_NAME);
}

void Debug(const char *c, ...)
{
    va_list args;
    va_start(args, c);
    vprintf(c, args);
    va_end(args);
}

void Panic(int err){
    printf("Exit: %d error\n", err);
    exit(0);
}
}

/* .....
    Functiile puse la dispozitia unui client
    care foloseste secventiatori distribuiti - varianta Unix
    -Suportul teoretic este prezentat in capitolul 4, la 4.3.4.3
    -Functiile care pot fi folosite sunt de un client :
    int SequencerInitialize(void)
    int SequencerTick(void)

    - Comunicarea cu serverul se face printr-un pipe FIFO
    .....*/

#include <stdio.h>
#include <string.h>
#include <signal.h>
#include <unistd.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <arpa/inet.h>
#include <netdb.h>
#include <sys/ipc.h>
#include <sys/sem.h>
#include <fcntl.h>
#include <sys/stat.h>
#include "tickeru.h"

#define BROADCAST_ADDR "193.226.37.255"
#define LOCAL_PORT 6571
#define FIFO_NAME "answerChan"
#define PERMS 0666
#define LEN 100
#define MAX_HOST_NAME 256
#define TICKET 5
#define MAX_FIFO_LEN 6

int GetServerAddress(void);
int SequencersInitialize(void);
int SequencerTicket(void);
int SequencerEnd(void);
int OpenAnswerChan(void);

int IChan;
struct sockaddr_in sad;
int answerChan;

int SequencerTicket(void) /*determinare valoare secventiator */
{
    char info[LEN];
    int val;
    sprintf(info, "%2d", TICKET);
    sendto(IChan, info, LEN, 0, (struct sockaddr *)&sad, sizeof(sad));
    read(answerChan, info, MAX_FIFO_LEN);
    sscanf(info, "%d", &val);
    return val;
}

int SequencersInitialize(void) /*realizeaza conectarea la server*/
{
    int errCode;
    if ((errCode=OpenAnswerChan()) != 0)
        return errCode;
    if ((errCode=GetServerAddress()) != 0)
        return errCode;
    memset(&sad, 0, sizeof(sad));
    sad.sin_family = AF_INET;
    sad.sin_addr.s_addr = inet_addr(MyServerAddress);
    sad.sin_port = htons(LOCAL_PORT);
    if((IChan= socket(AF_INET, SOCK_DGRAM, 0)) < 0)
        return ERROR5;
    return 0;
}

int SequencerEnd(void) /* deconectare de la server */
{
    close(IChan);
    unlink(FIFO_NAME);
}

char * GetMyAddress(void){
    int errCode;
}

```

```

if ((errCode = GetServerAddress()) != 0)
    return NULL;
else
    return MyServerAddress;
}

int GetServerAddress(void) /*determina adr,serverului propriu*/
{
    struct in_addr adr;
    char HostName[MAX_HOST_NAME];
    struct hostent *host;
    if (gethostname(HostName,MAX_HOST_NAME) != 0)
        return ERROR2;
    if (host = gethostbyname(HostName)) == NULL)
        return ERROR3;
    if (host != NULL)
        adr.s_addr = ((struct in_addr *) (host->h_addr_list[0]))->s_addr;
    if (strcpy(MyServerAddress,inet_ntoa(adr)) == NULL)
        return ERROR4;
    return 0;
}

int OpenAnswerChan(void) /* deschide answerChan */
{
    if ((answerChan = open(FIFO_NAME, O_RDWR)) < 0)
        return ERROR1;
    return 0;
}

char MyServerAddress[MAX_SERVER_ADDR];
int SequencerTicket(void);
int EventCounterAdvance(int);
int EventCounterAwait(int, int);
int EventCounterRead(int);
int GetServerAddress(void);
int ExclAccessInitialize(void);
void ExclAccessEnd(void);

int IChanE,IChanT ;
int answerChanE,answerChanT ;
struct sockaddr_in_saa,sat;

int SequencerTicket(void)
/*Operatie de TICKET*/
{
    char info[LEN];
    int val;
    sprintf(info,"%2d",TICKET) ;
    sendto(IChanT,info, LEN, 0,(struct sockaddr *) &sat,sizeof(sat)) ;
    read(answerChanT,info,MAX_FIFO_LEN);
    sscanf(info,"%d",&val); .
    return val;
}

int EventCounterAdvanced(int contor)
/* Operatie de ADVANCED - */
{
    char bf[LEN];
    if (contor > MAX_CONTOR-1 || contor < 0)
        return -1;
    sprintf(bf,"%2d-%2d",ADVANCE,contor);
    sendto(IChanE,bf,LEN,0,(struct sockaddr *)&saa,sizeof(saa));
    return 0;
}

int EventCounterAwait(int contor, int val)
/*Operatie de AWAIT*/
{
    char bf[LEN];
    if (contor > MAX_CONTOR-1 || contor < 0)
        return -1;
    sprintf(bf,"%2d-%2d-%2d",AWAIT,contor,val);
    sendto(IChanE,bf, LEN, 0,(struct sockaddr *)&saa, sizeof(saa));
    read(answerChanE,bf,MAX_FIFO_LEN);
    return 0;
}

int EventCounterRead(int contor)
/*Operatie de READ*/
{
    char bf[LEN];
    int val;
    if (contor > MAX_CONTOR-1 || contor < 0)
        return -1;
    sprintf(bf,"%2d-%2d",READ,contor) ;
    sendto(IChanE,bf, LEN, 0,(struct sockaddr *) &saa,sizeof(saa)) ;
    read(answerChanE,bf,MAX_FIFO_LEN);
    sscanf(bf,"%d",&val) ;
    return val;
}

char * GetMyAddress(void){
    int errCode;
    if ((errCode = GetServerAddress()) != 0)
        return NULL;
    else
        return MyServerAddress;
}

int GetServerAddress(void) /*determina adr,serverului propriu*/
{
    struct in_addr adr;
    char HostName[MAX_HOST_NAME];
    struct hostent *host;
    if (gethostname(HostName,MAX_HOST_NAME) != 0)

```

B10. Implementarea excluderii reciproce folosind secvențiatori și contori de evenimente distribuiți varianta Unix

/*

Se folosesc cele doua servere prezentate anterior:
cel pentru secvențiatori distribuiti si cel pentru contori de
evenimente distribuiti
Funcțiile puse la dispozitia programelor client sunt:

```

int EventCounterAdvanced(int contor)
int EventCounterAwait(int contor, int val)
int EventCounterRead(int contor)
int SequencerTicket(void)
char * GetMyAddress(void)
int ExclAccessInitialize(void);
void ExclAccessEnd(void);
.....*/

#include <stdio.h>
#include <string.h>
#include <signal.h>
#include <unistd.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <arpa/inet.h>
#include <netdb.h>
#include <sys/ipc.h>
#include <sys/sem.h>
#include <fcntl.h>
#include <sys/stat.h>
#include "evsequer.h"

#define LOCAL_PORT_EVENT 6561
#define LOCAL_PORT_SEQUENCER 6571
#define FIFO_NAMEE "eventChan"
#define FIFO_NAME_T "sequencerChan"
# define MAX_CONTOR 10
#define MAX_SERVER_ADDR 20
#define LEN 400
#define MAX_FIFO_LEN 6
#define MAX_HOST_NAME 256
#define ADVANCE 1
#define AWAIT 2
#define READ 3
#define TICKET 5

```

```

return ERROR3;
if((host = gethostbyname(HostName)) == NULL)
return ERROR4;
if(host != NULL)
adr.s_addr = ((struct in_addr *) (host->h_addr_list[0]))->s_addr;
if(strcmp(MyServerAddress,inet_ntoa(adr)) == NULL)
return ERROR5;
return 0;
}

int ExclAccessInitialize(void)
{
int errCode;
if ((errCode =OpenAnswerChannels()) != 0)
return errCode;
if ((errCode=GetServerAddress()) !=0 )
return errCode;
memset(&saa, 0, sizeof(saa) );
saa.sin_family = AF_INET ;
saa.sin_addr.s_addr = inet_addr(MyServerAddress) ;
saa.sin_port = htons(LOCAL_PORT_EVENT) ;
if((IChanE = socket(AF_INET, SOCK_DGRAM, 0) < 0)
return ERROR6;
memset(&sat, 0, sizeof(sat) );
sat.sin_family = AF_INET ;
sat.sin_addr.s_addr = inet_addr(MyServerAddress) ;
sat.sin_port = htons(LOCAL_PORT_SEQUENCER) ;
if((IChanT = socket(AF_INET, SOCK_DGRAM, 0) < 0)
return ERROR7;
return 0;
}

int OpenAnswerChannels()
/* creaza si deschide canalele answer*/
{
if((answerChanE = open(FIFO_NAMEE, O_RDWR)<0)
return ERROR1;
if((answerChanT = open(FIFO_NAME_T, O_RDWR)<0)
return ERROR2;
return 0;
}

void ExclAccessEnd(void)
/* sterge canalele answer */
{
unlink(FIFO_NAMEE);
unlink(FIFO_NAME_T);
close(IChanE);
close(IChanT);
}
}
.....
Programul client
.....
}
.....
#include <stdio.h>
#include <X11/Intrinsic.h>
#include <X11/StringDefs.h>
#include <X11/Shell.h>
#include <X11/Xaw/List.h>
#include <X11/Xaw/Command.h>
#include <X11/Xaw/Paned.h>
#include <X11/Xaw/Box.h>
#include <X11/Xaw/Dialog.h>
#include <X11/Xaw/AsciiText.h>
#include <string.h>
#include <sys/types.h>
#include <sys/socket.h>
#include <netinet/in.h>
#include <arpa/inet.h>
#include <netdb.h>
#include <signal.h>
#include "adtic.h"
#define LEN 100

Widget
widQuit,widInfo,widView, widPaned,widBoxMenu,
widBoxView, widP, widV, widTopLevel,
widRead,widRequest,widRelease,widInit,label;

int eventNo=1,eventVal;

void CallbackInit
(Widget w,XtPointer client_data,XtPointer call_data);
void CallbackInfo
(Widget w,XtPointer client_data,XtPointer call_data);
void CallbackRead
(Widget w,XtPointer client_data,XtPointer call_data);
void CallbackRequest
(Widget w,XtPointer client_data,XtPointer call_data);
void CallbackRelease
(Widget w,XtPointer client_data,XtPointer call_data);
void CallbackQuit
(Widget w,XtPointer client_data,XtPointer call_data);

main(argc, argv)
int argc;
char **argv;
{
XtAppContext app_context;
if (ExclAccessInitialize())
exit(1);
widTopLevel = XtVaAppInitialize(
&app_context,
"Main",
NULL,0,
&argc, argv,
NULL,
NULL);
widPaned = XtVaCreateManagedWidget(
"mainPaned",
panedWidgetClass,
widTopLevel,
NULL);
widBoxMenu = XtVaCreateManagedWidget(
"boxMenu",
boxWidgetClass,
widPaned,
XtOrientation, XtorientHorizontal,
NULL);
widInfo = XtVaCreateManagedWidget(
"Info",
commandWidgetClass,
widBoxMenu,
NULL);
widInit = XtVaCreateManagedWidget(
"Init",
commandWidgetClass,
widBoxMenu,
NULL);
widRead = XtVaCreateManagedWidget(
"Read",
commandWidgetClass,
widBoxMenu,
NULL);
widRequest = XtVaCreateManagedWidget(
"Request",
commandWidgetClass,
widBoxMenu,
NULL);
widRelease = XtVaCreateManagedWidget(
"Release",
commandWidgetClass,
widBoxMenu,
NULL);
widQuit = XtVaCreateManagedWidget(
"Quit",
commandWidgetClass,
widBoxMenu,
NULL);
label = XtVaCreateManagedWidget(
"label",
}

```

```
        labelWidgetClass,
        widPaned,
        XtNlabel,"n          \n",
        XtNwidth,400,
        NULL);
XtAddCallback(widInfo, XtNcallback, CallbackInfo, 0);
XtAddCallback(widQuit, XtNcallback, CallbackQuit, 0);
XtAddCallback(widInit, XtNcallback, CallbackInit, 0);
XtAddCallback(widRead, XtNcallback, CallbackRead, 0);
XtAddCallback(widRequest, XtNcallback, CallbackRequest, 0);
XtAddCallback(widRelease, XtNcallback, CallbackRelease, 0);
XtRealizeWidget(widTopLevel);
XtAppMainLoop(app_context);
}

void CallbackQuit(w, client_data, call_data)
Widget w;
XtPointer client_data, call_data;
{
    ExclAccessEnd();
    exit(0);
}

void CallbackInfo( w, client_data, call_data)
Widget w;
XtPointer client_data, call_data;
{
    char msg[LEN],*adr;
    adr=GetMyAddress();
    sprintf(msg,"Server address %s",adr);
    XtVaSetValues(label,XtNlabel,msg,NULL);
}

void CallbackInit( w, client_data, call_data)
Widget w;
XtPointer client_data, call_data;
{
    char msg[LEN];
    int i;
    EventCounterAdvanced(eventNo);
    sprintf(msg," Counter = %d",eventNo);
    XtVaSetValues(label,XtNlabel,msg,NULL);
}

void CallbackRead( w, client_data, call_data)
Widget w;
XtPointer client_data, call_data;
{
    char msg[LEN];
    int i;
    i = EventCounterRead(eventNo);
    sprintf(msg,"%d",i);
    XtVaSetValues(label,XtNlabel,msg,NULL);
}

void CallbackRequest( w, client_data, call_data)
Widget w;
XtPointer client_data, call_data;
{
    char msg[LEN];
    eventVal = Ticket();
    EventCounterAwait(eventNo, eventVal);
    XtVaSetValues(label,XtNlabel,"I am in critical section...",NULL);
}

void CallbackRelease( w, client_data, call_data)
Widget w;
XtPointer client_data, call_data;
{
    char msg[LEN];
    EventCounterAdvanced(eventNo);
    XtVaSetValues(label,XtNlabel,"Done release...",NULL);
}
```

ANEXA C

Implementarea unui toolkit pentru controlul accesului concurent la o bază de date în rețele Novell, Unix și Windows

C1. Implementarea primei soluții bazate pe semafoare în rețele Novell

Supportul teoretic este prezentat în capitolul 5, la 5.1.1.

În prima soluție se folosesc liste de identificatori ai proceselor clienți care sunt gestionate FIFO

Procesele care cer acces exclusiv sunt mai prioritare

Funcțiilor care pot fi apelate de programele utilizator sunt:

```
int ExclOrSharedAcclInitialize(void); /* realizeaza conectarea la sistem */
void ExclOrSharedAccessEnd(void); /* realizeaza deconectarea de la sistem */
int OpenSharedAccess(void); /* realizeaza obtinerea accesului pentru operatii partajabile asupra a bazei da date */
int CloseSharedAccess(void); /* realizeaza terminarea accesului i partajabil asupra a bazei da date */
int OpenExclAccess(void); /* realizeaza obtinerea accesului pentru operatii exclusive asupra a bazei da date */
int CloseExclAccess(void); /* realizeaza terminarea accesului exclusiv asupra a bazei da date */

#include <nif.h>
#include <stdlib.h>
#include <string.h>

#define SHARED_ACCESS 1 /* codficarea celor doua moduri de lucru: shereable, exclusive */
#define EXCL_ACCESS 2
#define TIMEOUT 360 /* timpul de asteptare maxim la un semafor */
#define MAX_CON 50 /* numarul maxim de noduri care pot participa la algoritim */
#define MAX_LENGTH 30 /* lungimea unui mesaj */
#define TRUE 1
#define FALSE 0
int ExclOrSharedAccessInitialize(void); /* realizeaza conectarea la sistem */
void ExclOrSharedAccessEnd(void); /* realizeaza deconectarea de la sistem */
int OpenSharedAccess(void); /* realizeaza obtinerea accesului pentru operatii partajabile asupra a bazei da date */
int CloseSharedAccess(void); /* realizeaza terminarea accesului i partajabil asupra a bazei da date */
int OpenExclAccess(void); /* realizeaza obtinerea accesului pentru operatii exclusive asupra a bazei da date */
int CloseExclAccess(void); /* realizeaza terminarea accesului exclusiv asupra a bazei da date */
int GetPriorityListProcess(int); /* functie apelabila local; nu este pusa la dispozitia clientilor */

long hNoSharedAccess, hNoExclAccess, hInExclAccess; /* handlers pentru semafoare */
long hSharedDataProtect, hWaitingSharedAccess[MAX_CON+1], hWaitingExclAccess[MAX_CON+1], hFlag[MAX_CON+1];
WORD pNoSharedAccess, pNoExclAccess, pInExclAccess;
WORD pSharedDataProtect, pWaitingSharedAccess[MAX_CON+1], pWaitingExclAccess[MAX_CON+1], pFlag[MAX_CON+1];
int vNoSharedAccess, vNoExclAccess, vInExclAccess, vFlag[MAX_CON+1];
int ConnectionNumber;
int code;
char name1[MAX_LENGTH] = "WaitingSharedAccess"; /* nume pentru semafoarele utilizate */
char name2[MAX_LENGTH] = "WaitingExclAccess";
char name3[MAX_LENGTH] = "Flag";
char temp[2] = "A";

int ExclOrSharedAccessInitialize(void) /* realizeaza conectarea la sistem */
{
    int i;
    if ((code = OpenSemaphore("NoSharedAccess", 0, &hNoSharedAccess, &pNoSharedAccess)) != 0) return 1;
    if ((code = OpenSemaphore("NoExclAccess", 0, &hNoExclAccess, &pNoExclAccess)) != 0) return 2;
    if ((code = OpenSemaphore("InExclAccess", 1, &hInExclAccess, &pInExclAccess)) != 0) return 3;
    if ((code = OpenSemaphore("SharedDataProtect", 1, &hSharedDataProtect, &pSharedDataProtect)) != 0) return 4;
    for (i=1; i<=MAX_CON; i++){
        temp[0] = 'A'+i;
        if ((code = OpenSemaphore((char *)strcat(name1, temp), 0, &hWaitingSharedAccess[i], &pWaitingSharedAccess[i])) != 0) return 5;
        strcpy(name1, "WaitingSharedAccess");
        if ((code = OpenSemaphore((char *)strcat(name2, temp), 0, &hWaitingExclAccess[i], &pWaitingExclAccess[i])) != 0) return 6;
        strcpy(name2, "WaitingExclAccess");
        if ((code = OpenSemaphore((char *)strcat(name3, temp), 0, &hFlag[i], &pFlag[i])) != 0) return 7;
        strcpy(name3, "Flag");
    }
    ConnectionNumber = (int)GetConnectionNumber();
    return 0;
}
```


Anexa C: Implementarea unui toolkit pentru controlul accesului concurrent la o baza de date

```
void ExclOrSharedAccessEnd(void) /* realizeaza deconectarea de la sistem */
{
    int i;
    CloseSemaphore(hNoSharedAccess); /* inchide semafoarele utilizate */
    CloseSemaphore(hNoExclAccess);
    CloseSemaphore(hSharedDataProtect);
    CloseSemaphore(hInExclAccess);
    for (i=1; i<=MAX_CON; i++){
        CloseSemaphore(hWaitingSharedAccess[i]);
        CloseSemaphore(hWaitingExclAccess[i]);
        CloseSemaphore(hFlag[i]);
    }
}

int OpenSharedAccess(void) /* realizeaza obtinerea accesului pentru operatii partajabile asupra a bazei da date */
{
    if ((code = WaitOnSemaphore(hSharedDataProtect, TIMEOUT)) != 0) return 8;
    if ((code = ExamineSemaphore(hNoExclAccess, &vNoExclAccess, &pNoExclAccess)) != 0) return 9;
    if (vNoExclAccess == 0) {
        if ((code = SignalSemaphore(hNoSharedAccess)) != 0) return 10;
        if ((code = SignalSemaphore(hWaitingSharedAccess[ConnectionNumber])) != 0) return 11;
    }
    else if ((code = SignalSemaphore(hFlag[ConnectionNumber])) != 0) return 12;
    if ((code = SignalSemaphore(hSharedDataProtect)) != 0) return 13;
    if ((code = WaitOnSemaphore(hWaitingSharedAccess[ConnectionNumber], TIMEOUT)) != 0) return 14;
    return 0;
}

int CloseSharedAccess(void) /* realizeaza terminarea accesului i partajabil asupra a bazei da date */
{
    int k;
    if ((code = WaitOnSemaphore(hSharedDataProtect, TIMEOUT)) != 0) return 15;
    if ((code = WaitOnSemaphore(hNoSharedAccess, TIMEOUT)) != 0) return 16;
    if ((code = ExamineSemaphore(hNoSharedAccess, &vNoSharedAccess, &pNoSharedAccess)) != 0) return 17;
    if ((code = ExamineSemaphore(hNoExclAccess, &vNoExclAccess, &pNoExclAccess)) != 0) return 18;
    if ((vNoSharedAccess == 0) && (vNoExclAccess > 0)) {
        k = GetPriorityListProcess(EXCL_ACCESS);
        if (k < 0) return abs(k);
        if ((code = WaitOnSemaphore(hInExclAccess, TIMEOUT)) != 0) return 19;
        if ((code = SignalSemaphore(hWaitingExclAccess[k])) != 0) return 20;
    }
    if ((code = SignalSemaphore(hSharedDataProtect)) != 0) return 21;
    return 0;
}

int OpenExclAccess(void) /* realizeaza obtinerea accesului pentru operatii exclusive asupra a bazei da date */
{
    if ((code = WaitOnSemaphore(hSharedDataProtect, TIMEOUT)) != 0) return 22;
    if ((code = SignalSemaphore(hNoExclAccess)) != 0) return 23;
    if ((code = ExamineSemaphore(hNoSharedAccess, &vNoSharedAccess, &pNoSharedAccess)) != 0) return 24;
    if ((code = ExamineSemaphore(hInExclAccess, &vInExclAccess, &pInExclAccess)) != 0) return 25;
    if ((vNoSharedAccess == 0) && (vInExclAccess == 1)) {
        if ((code = WaitOnSemaphore(hInExclAccess, TIMEOUT)) != 0) return 26;
        if ((code = SignalSemaphore(hWaitingExclAccess[ConnectionNumber])) != 0) return 27;
    }
    else {
        if ((code = SignalSemaphore(hFlag[ConnectionNumber])) != 0) return 28;
        if ((code = SignalSemaphore(hFlag[ConnectionNumber])) != 0) return 29;
    }
    if ((code = SignalSemaphore(hSharedDataProtect)) != 0) return 30;
    if ((code = WaitOnSemaphore(hWaitingExclAccess[ConnectionNumber], TIMEOUT)) != 0) return 31;
    return 0;
}

int CloseExclAccess(void) /* realizeaza terminarea accesului exclusiv asupra a bazei da date */
{
    int k;
    if ((code = WaitOnSemaphore(hSharedDataProtect, TIMEOUT)) != 0) return 32;
    if ((code = WaitOnSemaphore(hNoExclAccess, TIMEOUT)) != 0) return 33;
    if ((code = SignalSemaphore(hInExclAccess)) != 0) return 34;
    if ((code = ExamineSemaphore(hNoExclAccess, &vNoExclAccess, &pNoExclAccess)) != 0) return 35;
    if (vNoExclAccess > 0) {
        k = GetPriorityListProcess(EXCL_ACCESS);
        if (k < 0) return abs(k);
        if ((code = WaitOnSemaphore(hInExclAccess, TIMEOUT)) != 0) return 36;
        if ((code = SignalSemaphore(hWaitingExclAccess[k])) != 0) return 37;
    }
}
```

```
else
while((k = GetPriorityListProcess(SHARED_ACCESS)) > 0){
if (k < 0) return abs(k);
if ((code = SignalSemaphore(hNoSharedAccess)) != 0) return 38;
if ((SignalSemaphore(hWaitingSharedAccess[k])) != 0) return 39;
}
if ((code = SignalSemaphore(hSharedDataProtect)) != 0) return 40;
return 0;
}

int GetPriorityListProcess(int type) /* extrage din lista de identificatori ai clientilor urmatorul client pentru a fi deblocat */
{
int i;
if(type == SHARED_ACCESS) /* extrage un client pentru acces partajabil */
for(i=1;i<=MAX_CON; i++){
if ((code = ExamineSemaphore(hFlag[i],&vFlag[i],&pFlag[i])) != 0) return -41;
if(vFlag[i]==1){
if ((code = WaitOnSemaphore(hFlag[i],TIMEOUT)) != 0) return -42;
return i;
}
}
}
else /*extrage un client pentru acces exclusiv */
for(i=1;i<=MAX_CON; i++){
if ((code = ExamineSemaphore(hFlag[i],&vFlag[i],&pFlag[i])) != 0) return -43;
if(vFlag[i]==2){
if ((code = WaitOnSemaphore(hFlag[i], TIMEOUT)) != 0) return -44;
if ((code = WaitOnSemaphore(hFlag[i], TIMEOUT)) != 0) return -44;
return i;
}
}
}
return 0;
}
```

C2. Implementarea celei de-a doua soluții bazate pe semafoare în rețele Novell

.....
Suportul teoretic este prezentat in capitolul 5, la 5.1.1.

Procesele care cer acces exclusiv sunt mai prioritare.

Functiilor care pot fi apelate de programele utilizator sunt:

```
int ExclOrSharedAccessInitialize(void); /* realizeaza conectarea la sistem */
void ExclOrSharedAccessEnd(void); /* realizeaza deconectarea de la sistem */
int OpenSharedAccess(void ); /* realizeaza obtinerea accesului pentru operatii partajabile asupra a bazei da date */
int CloseSharedAccess(void ); /* realizeaza terminarea accesului i partajabil asupra a bazei da date */
int OpenExclAccess(void ); /* realizeaza obtinerea accesului pentru operatii exclusive asupra a bazei da date */
int CloseExclAccess(void ); /* realizeaza terminarea accesului exclusiv asupra a bazei da date */
.....

#include <nit.h>
#include <stdlib.h>

#define TIMEOUT 360 /* timpul de asteptare maxim la un semafor */
#define FALSE 0
#define TRUE 1

int ExclOrSharedAccessInitialize(void); /* realizeaza conectarea la sistem */
void ExclOrSharedAccessEnd(void); /* realizeaza deconectarea de la sistem */
int OpenSharedAccess(void ); /* realizeaza obtinerea accesului pentru operatii partajabile asupra a bazei da date */
int CloseSharedAccess(void ); /* realizeaza terminarea accesului i partajabil asupra a bazei da date */
int OpenExclAccess(void ); /* realizeaza obtinerea accesului pentru operatii exclusive asupra a bazei da date */
int CloseExclAccess(void ); /* realizeaza terminarea accesului exclusiv asupra a bazei da date */

int WaitWithCondition(long hExclusive, long hWaiting, long hCounter); /* functii apelabile local; nu sunt pusa la dispozitia clientilor */
int SignalWithCondition(long hExclusive, long hWaiting, long hCounter);

long hNoSharedAccess, hNoExclAccess, hStopExclAccess, hStopSharedAccess; /* handlere pentru semafoare */
long hExclusive, hExclDataProtect, hSharedDataProtect; /* handlere pentru semafoare */
WORD pNoSharedAccess, pNoExclAccess, pStopExclAccess, pStopSharedAccess;
WORD pExclusive, pExclDataProtect, pSharedDataProtect;
int vNoSharedAccess, vNoExclAccess;
int code;
```

Anexa C: Implementarea unui toolkit pentru controlul accesului concurrent la o baza de date

```
int ExclOrSharedAccessInitialize(void) /* realizeaza conectarea la sistem */
{
    if ((code = OpenSemaphore("Exclusive", 1, &hExclusive, &pExclusive)) != 0) return 1;
    if ((code = OpenSemaphore("NoSharedAccess", 0, &hNoSharedAccess, &pNoSharedAccess)) != 0) return 2;
    if ((code = OpenSemaphore("NoExclAccess", 0, &hNoExclAccess, &pNoExclAccess)) != 0) return 3;
    if ((code = OpenSemaphore("SharedDataProtect", 1, &hSharedDataProtect, &pSharedDataProtect)) != 0) return 4;
    if ((code = OpenSemaphore("ExclDataProtect", 1, &hExclDataProtect, &pExclDataProtect)) != 0) return 5;
    if ((code = OpenSemaphore("StopExclAccess", 1, &hStopExclAccess, &pStopExclAccess)) != 0) return 6;
    if ((code = OpenSemaphore("StopSharedAccess", 1, &hStopSharedAccess, &pStopSharedAccess)) != 0) return 7;
    return 0;
}

void ExclOrSharedAccessEnd(void) /* realizeaza deconectarea de la sistem */
{
    CloseSemaphore(hExclusive);
    CloseSemaphore(hNoSharedAccess);
    CloseSemaphore(hNoExclAccess);
    CloseSemaphore(hSharedDataProtect);
    CloseSemaphore(hExclDataProtect);
    CloseSemaphore(hStopSharedAccess);
    CloseSemaphore(hStopExclAccess);
}

int OpenSharedAccess(void) /* realizeaza obtinerea accesului pentru operatii partajabile asupra a bazei da date */
{
    int val;
    if ((code = WaitOnSemaphore(hExclusive, TIMEOUT)) != 0) return 8;
    if ((code = WaitOnSemaphore(hStopSharedAccess, TIMEOUT)) != 0) return 9;
    if ((val = WaitWithCondition(hSharedDataProtect, hStopExclAccess, hNoSharedAccess)) != 0)
        return val;
    if ((code = SignalSemaphore(hStopSharedAccess)) != 0) return 10;
    if ((code = SignalSemaphore(hExclusive)) != 0) return 11;
    return 0;
}

int CloseSharedAccess(void) /* realizeaza terminarea accesului i partajabil asupra a bazei da date */
{
    return SignalWithCondition(hSharedDataProtect, hStopExclAccess, hNoSharedAccess);
}

int OpenExclAccess(void) /* realizeaza obtinerea accesului pentru operatii exclusive asupra a bazei da date */
{
    int val;
    if ((val = WaitWithCondition(hExclDataProtect, hStopSharedAccess, hNoExclAccess)) != 0)
        return val;
    if ((code = WaitOnSemaphore(hStopExclAccess, TIMEOUT)) != 0) return 12;
    return 0;
}

int CloseExclAccess(void) /* realizeaza terminarea accesului exclusiv asupra a bazei da date */
{
    if ((code = SignalSemaphore(hStopExclAccess)) != 0) return 13;
    return SignalWithCondition(hExclDataProtect, hStopSharedAccess, hNoExclAccess);
}

int WaitWithCondition(long hExclusive, long hWaiting, long hCounter)
{
    WORD pCounter;
    int vCounter;

    if ((code = WaitOnSemaphore(hExclusive, TIMEOUT)) != 0) return 14;
    if ((code = ExamineSemaphore(hCounter, &vCounter, &pCounter)) != 0) return 15;
    if (vCounter == 0) if ((code = WaitOnSemaphore(hWaiting, TIMEOUT)) != 0) return 16;
    if ((code = SignalSemaphore(hCounter)) != 0) return 17;
    if ((code = SignalSemaphore(hExclusive)) != 0) return 18;
    return 0;
}

int SignalWithCondition(long hExclusive, long hWaiting, long hCounter)
{
    WORD pCounter;
    int vCounter;

    if ((code = WaitOnSemaphore(hExclusive, TIMEOUT)) != 0) return 19;
    if ((code = WaitOnSemaphore(hCounter, TIMEOUT)) != 0) return 20;
    if ((code = ExamineSemaphore(hCounter, &vCounter, &pCounter)) != 0) return 21;
    if (vCounter == 0) if ((code = SignalSemaphore(hWaiting)) != 0) return 22;
    if ((code = SignalSemaphore(hExclusive)) != 0) return 23;
    return 0;
}
```

C3. Implementarea celei de-a treia soluții bazate pe semafoare în rețele Novell

```
.....
Supportul teoretic este prezentat in capitolul 5, la 5.1.1.
Procesele care cer acces partajat sunt mai prioritare.
Funcțiilor care pot fi apelate de programele utilizator sunt:
int ExclOrSharedAccessInitialize(void); /* realizeaza conectarea la sistem */
void ExclOrSharedAccessEnd(void); /* realizeaza deconectarea de la sistem */
int OpenSharedAccess(void); /* realizeaza obtinerea accesului pentru operatii partajabile asupra a bazei da date */
int CloseSharedAccess(void); /* realizeaza terminarea accesului i partajabil asupra a bazei da date */
int OpenExclAccess(void); /* realizeaza obtinerea accesului pentru operatii exclusive asupra a bazei da date */
int CloseExclAccess(void); /* realizeaza terminarea accesului exclusiv asupra a bazei da date */
.....

#include <nit.h>
#include <stdlib.h>
#define TIMEOUT 360 /* timpul de asteptare maxim la un semafor */
#define FALSE 0
#define TRUE 1

int ExclOrSharedAccessInitialize(void); /* realizeaza conectarea la sistem */
void ExclOrSharedAccessEnd(void); /* realizeaza deconectarea de la sistem */
int OpenSharedAccess(void); /* realizeaza obtinerea accesului pentru operatii partajabile asupra a bazei da date */
int CloseSharedAccess(void); /* realizeaza terminarea accesului i partajabil asupra a bazei da date */
int OpenExclAccess(void); /* realizeaza obtinerea accesului pentru operatii exclusive asupra a bazei da date */
int CloseExclAccess(void); /* realizeaza terminarea accesului exclusiv asupra a bazei da date */

long hNoSharedAccess,hPrioritySharedAccess,hWaitingExclAccess,hExclusive;
WORD pNoSharedAccess,pExclusive,pWaitingExclAccess,pPrioritySharedAccess;
int vNoSharedAccess;
int code;

int ExclOrSharedAccessInitialize(void) /* realizeaza conectarea la sistem */
{
if ((code = OpenSemaphore("NoSharedAccess",0,&hNoSharedAccess,&pNoSharedAccess)) != 0) return 1;
if ((code = OpenSemaphore("Exclusive",1,&hExclusive,&pExclusive)) != 0) return 2;
if ((code = OpenSemaphore("WaitingExclAccess",1,&hWaitingExclAccess,&pWaitingExclAccess)) != 0) return 3;
if ((code = OpenSemaphore("PrioritySharedAccess",1,&hPrioritySharedAccess,&pPrioritySharedAccess)) != 0) return 4;
return 0;
}
void ExclOrSharedAccessEnd(void) /* realizeaza deconectarea de la sistem */
{
CloseSemaphore(hNoSharedAccess);
CloseSemaphore(hExclusive);
CloseSemaphore(hWaitingExclAccess);
CloseSemaphore(hPrioritySharedAccess);
}
int OpenSharedAccess(void) /* realizeaza obtinerea accesului pentru operatii partajabile asupra a bazei da date */
{
if ((code = WaitOnSemaphore(hExclusive,TIMEOUT)) != 0) return 5;
if ((code = SignalSemaphore(hNoSharedAccess)) != 0) return 6;
if ((code = ExamineSemaphore(hNoSharedAccess, &vNoSharedAccess, &pNoSharedAccess)) != 0) return 7;
if (vNoSharedAccess == 1) if ((code = WaitOnSemaphore(hPrioritySharedAccess,TIMEOUT)) != 0) return 8;
if ((SignalSemaphore(hExclusive)) != 0) return 9;
return 0;
}
int CloseSharedAccess(void) /* realizeaza terminarea accesului i partajabil asupra a bazei da date */
{
if ((code = WaitOnSemaphore(hExclusive,TIMEOUT)) != 0) return 10;
if ((code = WaitOnSemaphore(hNoSharedAccess,TIMEOUT)) != 0) return 11;
if ((code = ExamineSemaphore(hNoSharedAccess, &vNoSharedAccess, &pNoSharedAccess)) != 0) return 12;
if (vNoSharedAccess == 0) if ((code = SignalSemaphore(hPrioritySharedAccess)) != 0) return 13;
if ((code = SignalSemaphore(hExclusive)) != 0) return 14;
return 0;
}
int OpenExclAccess(void) /* realizeaza obtinerea accesului pentru operatii exclusive asupra a bazei da date */
{
if ((code = WaitOnSemaphore(hWaitingExclAccess,TIMEOUT)) != 0) return 15;
if ((code = WaitOnSemaphore(hPrioritySharedAccess,TIMEOUT)) != 0) return 16;
return 0;
}
int CloseExclAccess(void) /* realizeaza terminarea accesului exclusiv asupra a bazei da date */
{
if ((code = SignalSemaphore(hPrioritySharedAccess)) != 0) return 17;
if ((code = SignalSemaphore(hWaitingExclAccess)) != 0) return 18;
return 0;
}
}
.....
```

C4. Implementarea soluției bazate pe un proces coordonator în rețele Novell

```

...../
Suportul teoretic este prezentat in capitolul 5, la 5.1.2.
Procesele care cer acces exclusiv sunt mai prioritare.
*Funcțiilor care pot fi apelate de programele utilizator sunt:
int ExclOrSharedAccessInitialize(void); /*conectare la sistem */
void ExclOrSharedAccessEnd(void); /* deconectare */
int OpenSharedAccess(void); /* obtinere acces partajat */
int CloseSharedAccess(void); /* obtinere acces partajat */
int OpenExclAccess(void); /* obtinere acces exclusiv */
int CloseExclAccess(void); /* obtinere acces exclusiv */
...../
}

Program client.
Accesul este obtinut de la un server a carui localizare se
determina cunoscind numele de user sub care s-a deschis
asiunea la statia serverului. Receptia mesajelor se face in mod
asincron, intr-o functie EventServiceRoutine.
...../
#include <nxt.h>
#include <nit.h>
#include <string.h>
#include <mem.h>
#include <stdio.h>
#include <dos.h>

#define WAIT_FOR_MESSAGE 0
#define MAX_NO_PACKETS 10
#define DATA_LENGTH 15
#define USER_NAME_LENGTH 40
#define TRUE 1
#define FALSE 0
#define RECEIVE_ERROR 99
#define SEND_ERROR 98

int ExclOrSharedAccessInitialize(void);
void ExclOrSharedAccessEnd(void);
int OpenSharedAccess(void);
int CloseSharedAccess(void);
void far _loads ReceiveESR(void);
int InitReceive(void);
void WaitResponse(void);
void ReceiveResponse(char *msg);
int SendRequest(char request, int cliConnectionNo, int
servConnectionNo);
void FindDestinationAddress( int *servConnectionNo );

WORD clientSocketRec=0x5000;
IPXHeader packetHeader;
ECB ecbR;
char packetData[DATA_LENGTH];
char queue[MAX_NO_PACKETS];
int cliConnNo, servConnNo;
int ind, code;

int ExclOrSharedAccessInitialize(void)
{
int val;
if ((code = IPXInitialize()) != 0) return 1;
cliConnNo = (int)GetConnectionNumber();
FindDestinationAddress(&servConnNo);
if ((val = InitReceive()) != 0) return val;
return 0;
}

void ExclOrSharedAccessEnd(void)
{
IPXCancelEvent(&ecbR);
IPXCloseSocket(clientSocketRec);
}

int OpenSharedAccess(void)
{
char buffer;
int val;
if ((val = SendRequest('r',cliConnNo,servConnNo)) != 0)
return val;
WaitResponse();
ReceiveResponse(&buffer);
if(buffer != '@'){
code = 0;
return RECEIVE_ERROR;
}
return 0;
}

int CloseSharedAccess(void)
{
char buffer;
int val;
if ((val = SendRequest('R',cliConnNo,servConnNo)) != 0)
return val;
WaitResponse();
ReceiveResponse(&buffer);
if(buffer != '@'){
code = 0;
return RECEIVE_ERROR;
}
return 0;
}

int OpenExclAccess(void)
{
char buffer;
int val;
if ((val = SendRequest('w',cliConnNo,servConnNo)) != 0)
return val;
WaitResponse();
ReceiveResponse(&buffer);
if (buffer != '@'){
code = 0;
return RECEIVE_ERROR;
}
return 0;
}

int CloseExclAccess(void)
{
char buffer;
int val ;
if ((val = SendRequest('W',cliConnNo,servConnNo)) != 0)
return val;
WaitResponse();
ReceiveResponse(&buffer);
if (buffer != '@'){
code = 0;
return RECEIVE_ERROR;
}
return 0;
}

int InitReceive(void)
{
if ((code = IPXOpenSocket((BYTE*)&clientSocketRec,0x00))!=0)
return 2;
ind = 0;
/* initializarea structurilor de date */
ecbR.ESRAddress = ReceiveESR;
ecbR.fragmentCount=2;
ecbR.socketNumber=clientSocketRec;
ecbR.fragmentDescriptor[0].address=&packetHeader;
ecbR.fragmentDescriptor[0].size=sizeof(IPXHeader);
ecbR.fragmentDescriptor[1].address=(BYTE *) packetData;
ecbR.fragmentDescriptor[1].size=DATA_LENGTH;
}

```

Anexa C: Implementarea unui toolkit pentru controlul accesului concurrent la o baza de date

```

/*deschiderea soclului si intrarea in "ascultare" de pachete*/
IPXListenForPacket(&ecbR);
IPXRelinquishControl();
return 0;
}

void far _loadds ReceiveESR(void)
{
    if(ind <MAX_NO_PACKETS && ecbR.completionCode ==0){
        queue[ind] = packetData[0];
        ind++;
    }
    IPXListenForPacket(&ecbR);
}

void WaitResponse(void)
{
    int ind1;

    do{
        disable();
        ind1=ind;
        enable();
    }while (ind1 == WAIT_FOR_MESSAGE);
}

void ReceiveResponse(char *msg)
{
    int i;
    disable();
    if (ind>0) {
        *msg=queue[0];
        for (i=0; i<ind; i++) queue[i] = queue[i+1];
        ind--;
    }
    enable();
}

int SendRequest(char request, int cliConnNo, int servConnNo)
{
    BYTE serverNode[6];
    BYTE serverNetwork[4];
    WORD serverSocket = 0x5000,Socket;
    WORD clientSocketSnd = 0x6000;
    IPXHeader packetHeader;
    ECB ecb;
    char packetData[DATA_LENGTH];
    char buff[20];

    GetInternetAddress(servConnNo,(char *) serverNetwork,
                      (char *)serverNode,&Socket);
    movmem(serverNetwork,packetHeader.destination.network,4);
    movmem(serverNode,packetHeader.destination.node,6);
    movmem(&serverSocket,packetHeader.destination.socket,2);
    movmem(serverNode,ecb.immediateAddress,6);
    if((code =IPXOpenSocket((BYTE*)&clientSocketSnd,0x00))!=0)
    {
        IPXCancelEvent(&ecbR);
        return 3;
    }
    buff[0] = request;
    itoa(cliConnNo,&buff[1],10);
    buff[strlen(buff)+1] = '\0';
    strcpy(packetData,buff);
    ecb.ESRAddress = NULL;
    ecb.fragmentCount = 2;
    ecb.socketNumber = clientSocketSnd;
    ecb.fragmentDescriptor[0].address = &packetHeader;
    ecb.fragmentDescriptor[0].size = sizeof(IPXHeader);
    ecb.fragmentDescriptor[1].address = (BYTE *) packetData;
    ecb.fragmentDescriptor[1].size = strlen(packetData)+1;
    IPXSendPacket(&ecb);

    while(ecb.inUseFlag != 0) IPXRelinquishControl();
    if(ecb.completionCode != 0){
        IPXCancelEvent(&ecbR);
        IPXCloseSocket(clientSocketRec);
    }
}

code = 0;
return SEND_ERROR;
}
IPXCloseSocket(clientSocketSnd);
return 0;
}

void FindDestinationAddress( int *servConnNo )
{
    char name[USER_NAME_LENGTH] ;
    WORD numberOfConnections ;
    WORD connectionList[1] ;
    WORD maxConnections = 1 ;

    do{
        printf("\nLogin (user name for server): ");
        scanf("%s",name) ;

        GetObjectConnectionNumbers(name,OT_USER,&numberOfCon
        nections, connectionList, maxConnections);
        if(numberOfConnections == 0) printf("User name not
        available. Try again!\n");
    }while (numberOfConnections == 0);
    *servConnNo = *connectionList ;
}

/*****
Program server.
Localizarea sa se determina cunoscind numele de user sub care
s-a deschis sesiunea pe stia serverului
Receptia mesajelor se face in mod asincron, intr-o functie
EventServiceRoutine.
La receptia pachetelor, acestea sunt depuse intr-un buffer de
unde se preiau apoi pentru tratare.
*****/

#include <nxt.h>
#include <nit.h>
#include <stdio.h>
#include <bios.h>
#include <string.h>
#include <mem.h>
#include <conio.h>
#include <dos.h>
#include <stdlib.h>

#define WAIT_FOR_MESSAGE 0
#define DATA_LENGTH 10
#define MAX_NO_PACKETS 10
#define MAX_NO_CONN 30
#define TRUE 1
#define FALSE 0

void Initialize(void);
void far _loadds ReceiveESR(void);
void InitReceive(void);
void WaitRequest(void);
void queueExtract(char *msg);
void ReceiveRequest(char *prequest);
void SendResponse(char ok[], int nrConn);
int GetPriProcess(int I[]);
int EmptyList(int I[]);
void OpenSharedAccess(void);
void CloseSharedAccess(void);
void OpenExclAccess(void);
void CloseExclAccess(void);

WORD serverSocketRec = 0x5000;
IPXHeader packetHeader;
ECB ecbR;
char packetData[DATA_LENGTH];
char *queue[MAX_NO_PACKETS];
int ind;
int SharedAccessNo,ExclAccessNo;
int inExclAccess;
int ShAcclList[MAX_NO_CONN],ExclAcclList[MAX_NO_CONN];
int cliConnNo;

```

Anexa C: Implementarea unui toolkit pentru controlul accesului concurrent la o baza de date

```
void main(void)
{
    int i;
    char request;
    Initialize(); /* initializarea comunicatiei */
    clrscr();
    SharedAccessNo = 0; /* no de procese shared */
    ExclAccessNo = 0; /* no de procese exclusiv */
    inExclAccess = FALSE; /* acces liber */
    for(i=0;i<MAX_NO_CONN;i++){ /* liste vide */
        ShAccList[i] = 0;
        ExclAccList[i] = 0;
    }
    InitReceive(); /* initializare receptie */
    do{ /* bucla de primire mesaje */
        WaitRequest(); /* asteptare cerere */
        ReceiveRequest(&request); /* receptie cerere */
        switch(request){ /* test tip cerere */
            case 'r':
                OpenSharedAccess();
                break;
            case 'R':
                CloseSharedAccess();
                break;
            case 'w':
                OpenExclAccess();
                break;
            case 'W':
                CloseExclAccess();
                break;
        }
    }while(1);
}

void Initialize(void) /* initializare sistem comunicatie */
{
    if(IPXInitialize()!=0) {
        printf("\nIPX not installed. Abort!\n");
        exit(1);
    }
}

void InitReceive(void) /* initializare receptie */
{
    int code,i;

    if((code=IPXOpenSocket((BYTE*)&serverSocketRec,0x00))!=0) {
        printf("\nOpen socket error: %d. Abort!",code);
        exit(1);
    }
    ind=0;
    /* rezervarea cozii de pachete - pachetele se buffereaza */
    for (i=0; i<MAX_NO_PACKETS; i++)
        if((queue[i] = (char*)malloc(DATA_LENGTH))!=NULL){
            printf("Memory allocation error. Abort!\n");
            exit(1);
        }
    /* initializarea structurilor de date */
    ecbR.ESRAddress=ReceiveESR;
    ecbR.fragmentCount=2;
    ecbR.socketNumber=serverSocketRec;
    ecbR.fragmentDescriptor[0].address=&packetHeader;
    ecbR.fragmentDescriptor[0].size=sizeof(IPXHeader);
    ecbR.fragmentDescriptor[1].address=(BYTE *) packetData;
    ecbR.fragmentDescriptor[1].size=DATA_LENGTH;

    /*deschiderea soclului si intrarea in "ascultare" de pachete*/
    IPXListenForPacket(&ecbR);
    IPXRelinquishControl();
}

void far _loadds ReceiveESR(void)
/* rutina la nsata in cadrul intreruperii la primirea unui mesaj */
{
    if (ind < MAX_NO_PACKETS && ecbR.completionCode == 0){
        strcpy(queue[ind],(char*)packetData);
        ind++;
    }
    IPXListenForPacket(&ecbR);
}

void WaitRequest(void)
/* functie pentru asteptarea sosirii unui mesaj */
{
    int ind1;

    do{
        disable();
        ind1=ind;
        enable();
    }while ((ind1 == WAIT_FOR_MESSAGE) && !bioskey(1));
    if (bioskey(1)){
        printf("\nServer terminated by user intervention!\n");
        IPXCancelEvent(&ecbR);
        IPXCloseSocket(serverSocketRec);
        exit(1);
    }
}

void queueExtract(char *msg)
/*extrage din coada de pachete urmatoarea cerere netratata*/
{
    int i;
    disable();
    if(ind>0){
        strcpy(msg,queue[0]);
        for(i=0;i<ind;i++)
            strcpy(queue[i],queue[i+1]);
        ind--;
    }
    enable();
}

void ReceiveRequest(char *prequest)
/* preia tipul urmatorului mesaj */
{
    char tmp[30];

    queueExtract(tmp);
    cliConnNo = atoi(&tmp[1]);
    printf("%d ",cliConnNo);
    *prequest = tmp[0];
}

void SendResponse(char ok[], int cliConnNo)
/* trimite un pachet la apelantul unei cereri */
{
    WORD serverSocketSnd = 0x6000;
    WORD clientSocket = 0x5000,Socket;
    BYTE clientNode[6];
    BYTE clientNetwork[4];
    IPXHeader packetHeader;
    ECB ecb;
    char packetData[DATA_LENGTH];
    int code;
    strcpy(packetData, ok);
    GetInternetAddress(cliConnNo,(char *) clientNetwork,
        (char *)clientNode,&Socket);
    movmem(clientNetwork,packetHeader.destination.network,4);
    movmem(clientNode,packetHeader.destination.node,6);
    movmem(&clientSocket,packetHeader.destination.socket,2);
    movmem(clientNode,ecb.immediateAddress,6);

    ecb.ESRAddress = NULL;
    ecb.fragmentCount = 2;
    ecb.socketNumber = serverSocketSnd;
    ecb.fragmentDescriptor[0].address = &packetHeader;
    ecb.fragmentDescriptor[0].size = sizeof(IPXHeader);
    ecb.fragmentDescriptor[1].address = (BYTE *) packetData;
    ecb.fragmentDescriptor[1].size = strlen(packetData)+1;
    if((code=IPXOpenSocket((BYTE*)&serverSocketSnd,0x00))!=0)
    {
        IPXSendPacket(&ecb);
    }
}
```


Anexa C: Implementarea unui toolkit pentru controlul accesului concurrent la o baza de date

```
while(ecb.inUseFlag!=0) IPXRelinquishControl();
}
else {
    printf("Open socket error: %d. Abort!\n",code);
    IPXCancelEvent(&ecbR);
    IPXCloseSocket(serverSocketRec);
    exit (1);
}
IPXCloseSocket(serverSocketSnd);
}

int EmptyList(int list[])
/* testeaza daca o lista de identificatori este vida */
{
    int i;
    for (i=0; i<MAX_NO_CONN; i++)
        if (list[i] != 0)
            return 0;
    return 1;
}

int GetPriProcess(int list[])
/* preia identificatorul urmatorelui apelant
din lista de identificatori */
{
    int i;
    for (i=1; i<=MAX_NO_CONN; i++)
        if (list[i] != 0){
            list[i] = 0;
            return i;
        }
    return 0;
}

void OpenSharedAccess(void)
/* realizeaza obtinerea accesului pentru operatii partajabile asupra
a bazei da date */
{
    if (ExclAccessNo == 0){
        SharedAccessNo++;
        SendResponse("@",cliConnNo);
    }
    else
        ShAccList[cliConnNo] = cliConnNo;
}

void CloseSharedAccess(void)
/* realizeaza terminarea accesului i partajabil asupra a bazei da
date */
{
    SharedAccessNo--;
    SendResponse("@",cliConnNo);
    if((SharedAccessNo == 0) && (!EmptyList(ExclAccList))){
        cliConnNo = GetPriProcess(ExclAccList);
        inExclAccess = TRUE;
        SendResponse("@",cliConnNo);
    }
}

void OpenExclAccess(void)
/* realizeaza obtinerea accesului exclusiv asupra a bazei da date
*/
{
    ExclAccessNo++;
    if((SharedAccessNo == 0) && (inExclAccess == FALSE)){
        inExclAccess = TRUE;
        SendResponse("@",cliConnNo);
    }
    else
        ExclAccList[cliConnNo] = cliConnNo;
}

void CloseExclAccess(void)
/* realizeaza terminarea accesului exclusiv asupra a bazei da
date */
{
    ExclAccessNo--;
    inExclAccess = FALSE;
    SendResponse("@",cliConnNo);
    if(!EmptyList(ExclAccList)){
        cliConnNo = GetPriProcess(ExclAccList);
        inExclAccess = TRUE;
        SendResponse("@",cliConnNo);
    }
    else
        while(!EmptyList(ShAccList)){
            cliConnNo = GetPriProcess(ShAccList);
            SharedAccessNo++;
            SendResponse("@",cliConnNo);
        }
}
}
```

C5. Implementarea celei de-a treia soluții bazate pe semafoare în rețele Unix și Windows

.....
 Varianta Unix

Suportul teoretic este prezentat în capitolul 5, la 5.1.1.
 Implementarea este descrisă la 5.1.3.2.
 Procesele care cer acces partajat sunt mai prioritare.
 Se folosesc semafoare distribuite a caror implementare a fost prezentată în anexa B7
 Funcțiile care pot fi apelate de programele utilizator sunt:

```
int ExclOrSharedAccessInitialize(void); /*conectare la sistem */
void ExclOrSharedAccessEnd(void); /* deconectare */
void ExclOrSharedAccessStart(void) /* initializare semaf. */
void OpenSharedAccess(void); /* obtinere acces partajat */
void CloseSharedAccess(void); /* terminare acces partajat */
void OpenExclAccess(void); /* obtinere acces exclusiv */
void CloseExclAccess(void); /* terminare acces exclusiv */
.....
```

.....
 Descrierea funcțiilor puse la dispoziția programelor client

.....
 /* definirea semafoarelor distribuite folosite */

```
int hPrioritySharedAcc= 3;
int hNoSharedAccess = 4;
int hWaitingExclAccess = 1;
int hExclusive = 2;
int NoSharedAccess;
```

```
void OpenSharedAccess(void){
SemaphoreDown( hExclusive );
NoSharedAccess = SemaphoreUp( hNoSharedAccess );
if (NoSharedAccess==1)SemaphoreDown( hPrioritySharedAcc);
SemaphoreUp( hExclusive);
}
```

```
void CloseSharedAccess(void){
SemaphoreDown( hExclusive );
NoSharedAccess = SemaphoreDown( hNoSharedAccess );
if (NoSharedAccess==0) SemaphoreUp( hPrioritySharedAcc);
SemaphoreUp( hExclusive );
}
```

```
void OpenExclAccess(void){
SemaphoreDown( hWaitingExclAccess );
SemaphoreDown( hPrioritySharedAcc);
}
```

```
void CloseExclAccess(void){
SemaphoreUp( hPrioritySharedAcc);
SemaphoreUp( hWaitingExclAccess );
}
```

```
int ExclOrSharedAccessInitialize(void){
int errCode;
if (errCode = SemaphoreInitialize()) != 0)
return errCode;
return 0;
}
```

```
void ExclOrSharedAccessStart(void){
SemaphoreUp( hPrioritySharedAcc);
SemaphoreUp( hWaitingExclAccess );
SemaphoreUp( hExclusive );
}
```

```
void ExclOrSharedAccessEnd(void){
SemaphoreEnd();
}
```

.....
 Programul client

.....

```
#include <stdio.h>
#include <X11/Intrinsic.h>
#include <X11/StringDefs.h>
#include <X11/Shell.h>
#include <X11/Xaw/List.h>
#include <X11/Xaw/Command.h>
#include <X11/Xaw/Paned.h>
#include <X11/Xaw/Box.h>
#include <X11/Xaw/Dialog.h>
#include <X11/Xaw/AsciiText.h>
#include "shexcl.h"
```

```
# define MAX_MES 50
```

```
typedef struct _ViewText{
XawTextPosition      nrChText;
unsigned int         nrRwText;
XawTextPosition      chPos;
unsigned int         rwPos;
unsigned int         rowsNr;
char                 *text;
}ViewText;
```

```
Widget widQuit,widInfo,widView, widPaned,
widBoxMenu, widBoxView, widP, widV, widTopLevel,
open_shared,open_excl,close_shared,close_excl,init;
void WriteMessage(Widget widDest,char *wString);
void AdjustDisplayPos(ViewText *display, int rowsNr);
```

```
void CallbackOpenShared
(Widget w,XtPointer client_data,XtPointer call_data);
void CallbackCloseShared
(Widget w,XtPointer client_data,XtPointer call_data);
void CallbackOpenWrite
(Widget w,XtPointer client_data,XtPointer call_data);
void CallbackCloseWrite
(Widget w,XtPointer client_data,XtPointer call_data);
void CallbackInit
(Widget w,XtPointer client_data,XtPointer call_data);
void CallbackInfo
(Widget w,XtPointer client_data,XtPointer call_data);
void CallbackQuit
(Widget w,XtPointer client_data,XtPointer call_data);
```

```
String SemList[MAX_SEM+1] = {NULL};
```

```
int heightFont;
```

```
main(argc, argv)
int argc;
char **argv;
{
Dimension heightApp = 300,widthApp = 600,
heightView, heightMenu;
XtAppContext app_context;
if ( ExclOrSharedAccessInitialize())
exit(1);
widTopLevel = XtVaAppInitialize(
&app_context,
"Main",
NULL,0,
&argc, argv,
NULL,
NULL);
XtVaSetValues(widTopLevel,
XtNx, 300,
XtNy, 400,
XtNwidth, widthApp,
XtNheight, heightApp,
NULL);
widPaned = XtVaCreateManagedWidget(
"mainPaned",
panedWidgetClass,
```

Anexa C: Implementarea unui toolkit pentru controlul accesului concurrent la o baza de date

```
    widTopLevel,
    NULL);
widBoxMenu = XtVaCreateManagedWidget(
    "boxMenu",
    boxWidgetClass,
    widPaned,
    XINorientation, XtorientHorizontal,
    NULL);
widInfo = XtVaCreateManagedWidget(
    "Info",
    commandWidgetClass,
    widBoxMenu,
    NULL);
open_shared= XtVaCreateManagedWidget(
    "OpenShared",
    commandWidgetClass,
    widBoxMenu,
    NULL);
close_shared= XtVaCreateManagedWidget(
    "CloseShared",
    commandWidgetClass,
    widBoxMenu,
    NULL);
open_excl = XtVaCreateManagedWidget(
    "OpenExcl",
    commandWidgetClass,
    widBoxMenu,
    NULL);
close_excl= XtVaCreateManagedWidget(
    "CloseExcl",
    commandWidgetClass,
    widBoxMenu,
    NULL);
init = XtVaCreateManagedWidget(
    "init",
    commandWidgetClass,
    widBoxMenu,
    NULL);
widQuit = XtVaCreateManagedWidget(
    "Quit",
    commandWidgetClass,
    widBoxMenu,
    NULL);
XtVaGetValues( widQuit, XInheight, &heightMenu, NULL);
XawPanedSetMinMax( widBoxMenu, heightMenu+10,
                    heightMenu+10 );
widBoxView = XtVaCreateManagedWidget(
    "boxView",
    formWidgetClass,
    widPaned,
    NULL);
widP = XtVaCreateManagedWidget(
    "widP",
    listWidgetClass,
    widBoxView,
    XtNlist, SemList,
    XtNwidth, 4,
    XtNheight, 0,
    XINdefaultColumns,1,
    XINforceColumns,True,
    NULL);
widView = XtVaCreateManagedWidget(
    "Message",
    asciiTextWidgetClass,
    widBoxView,
    NULL);
XtVaGetValues(widView,
    XtNheight,&heightFont,
    NULL);
XtVaGetValues(widP,
    XtNheight,&heightView,
    NULL);
XtVaSetValues(widView,
    XINeditType, XawtextEdit,
    XINscrollVertical,XawtextScrollAlways,
    XINfromHoriz, widP,
    XtNheight, heightView,
    NULL);
widV = XtVaCreateManagedWidget(
    "widP",
    listWidgetClass,
    widBoxView,
    XtNlist, SemList,
    XtNwidth, 4,
    XtNheight, 0,
    XINdefaultColumns,1,
    XINforceColumns,True,
    NULL);
XtAddCallback(open_shared,XtNcallback,
    CallbackOpenShared,0);
XtAddCallback(close_shared,XtNcallback,
    CallbackCloseShared,0);
XtAddCallback(open_excl,XtNcallback,CallbackOpenWrite,0);
XtAddCallback(close_excl,XtNcallback,CallbackCloseWrite,0);
XtAddCallback(init,XtNcallback,CallbackInit,0);
XtAddCallback(widInfo, XtNcallback, CallbackInfo, 0);
XtAddCallback(widQuit, XtNcallback, CallbackQuit, 0);

XtOverrideTranslations( widView,
    XtParseTranslationTable("<Key>:");
XtRealizeWidget(widTopLevel);
XtAppMainLoop(app_context);
}

void CallbackInit(w,client_data,call_data)
Widget w;
XtPointer client_data,call_data;
{
    ExclOrSharedAccessStart();
}

void CallbackQuit(w, client_data, call_data)
Widget w;
XtPointer client_data, call_data;
{
    ExclOrSharedAccessEnd();
    exit(0);
}

void CallbackInfo( w, client_data, call_data)
Widget w;
XtPointer client_data, call_data;
{
    char msg[MAX_MES], *adr;
    adr=GetMyAddress();
    sprintf(msg,"Server address %s",adr);
    WriteMessage( widView, msg);
}

void CallbackOpenShared(w,client_data,call_data)
Widget w;
XtPointer client_data,call_data;
{
    OpenSharedAccess();
    WriteMessage(widView, "In shared access...");
}

void CallbackCloseShared(w,client_data,call_data)
Widget w;
XtPointer client_data,call_data;
{
    CloseSharedAccess();
    WriteMessage(widView, "Done shared access...");
}

void CallbackOpenWrite(w,client_data,call_data)
Widget w;
XtPointer client_data,call_data;
{
    OpenExclAccess();
    WriteMessage(widView, "In exclusive access...");
}
void CallbackCloseWrite(w,client_data,call_data
```

Anexa C: Implementarea unui toolkit pentru controlul accesului concurent la o baza de date

```

Widget w;
XtPointer client_data,call_data;
{
  CloseExclAccess();
  WriteMessage(widView, " Done exclusive access..");
}
void AdjustDisplayPos( display, rowsNr)
ViewText *display;
int rowsNr;
{
  XawTextPosition iPos;
  if( rowsNr != display->rowsNr ){
    display->rowsNr = rowsNr;
    for( iPos = display->nrChText; iPos && rowsNr; iPos-- )
      if( display->text[iPos] == '\n' )
        rowsNr--;
    display->chPos = iPos ? iPos + 2 : 0L ;
    display->rwPos=display->nrRwText-( display->rowsNr -rowsNr );
  }
  else
    if( display->rwPos + display->rowsNr <= display->nrRwText ){
      display->rwPos ++ ; /* Scroll o linie */
      for( ; display->text[display->chPos++] != '\n' ; ) ;
    }
}

void WriteMessage(widDest,wString)
Widget widDest;
char *wString;
{
  static ViewText display = { 0, 0, 0, 0, 0, NULL};
  char newLine[ 100];
  XawTextBlock newText;
  Boolean newDispPos;
  Dimension heightWid;
  XawTextPosition displayPos;
  strcpy( newLine, wString);
  strcat( newLine, "\n");
  XtVaGetValues( widDest,
                XtNdisplayPosition, &displayPos,
                XtNheight, &heightWid,
                XtNstring, &display.text,
                NULL);
  newText.firstPos = 0;
  newText.format = FMT8BIT;
  newText.ptr = newLine;
  newText.length = strlen(newLine);
  newDispPos = ( displayPos >= display.chPos );
  XawTextReplace( widDest, display.nrChText, display.nrChText,
                  &newText );
  display.nrChText += newText.length;
  display.nrRwText ++;
  AdjustDisplayPos( &display, heightWid/heightFont );
  if(newDispPos)
    XtVaSetValues( widDest, XtNdisplayPosition,
                  display.chPos, NULL);
}
}
/*****
Descrierea functiilor puse la dispozitia programelor client
*****/
#include <stdio.h>
#include <stdlib.h>
#include "semaforw.h"
#include "generic.h"

int hWaitingExclAccess = 1; /* handlerse semafoare */
int hExclusive = 2;
int hPrioritySharedAcc= 3;
int hNoSharedAccess = 4;
int NoSharedAccess; /* variabila partajata */

void OpenSharedAccess(void){
  SemaphoreDown( hExclusive );
  NoSharedAccess = SemaphoreUp( hNoSharedAccess );
  if (NoSharedAccess==1)SemaphoreDown( hPrioritySharedAcc);
  SemaphoreUp( hExclusive);
}

void CloseSharedAccess(void){
  SemaphoreDown( hExclusive );
  NoSharedAccess = SemaphoreDown( hNoSharedAccess );
  if (NoSharedAccess==0) SemaphoreUp( hPrioritySharedAcc);
  SemaphoreUp( hExclusive );
}

void OpenExclAccess(void){
  SemaphoreDown( hWaitingExclAccess );
  SemaphoreDown( hPrioritySharedAcc);
}

void CloseExclAccess(void){
  SemaphoreUp( hPrioritySharedAcc);
  SemaphoreUp( hWaitingExclAccess );
}

int ExclOrSharedAccessInitialize(void){
  int errCode;
  if ((errCode = SemaphoreInitialize()) != 0)
    return errCode;
  return 0;
}

void ExclOrSharedAccessStart(void){
  SemaphoreUp( hPrioritySharedAcc);
  SemaphoreUp( hWaitingExclAccess );
  SemaphoreUp( hExclusive );
}

void ExclOrSharedAccessEnd(void){
  SemaphoreEnd();
}
}
/*****
Programul client
*****/
Varianta Windows
*****

Suportul teoretic este prezentat in capitolul 5, la 5.1.1.
Implementarea este descrisa la 5.1.3.2.
Procesele care cer acces partajat sunt mai prioritare.
Se folosesc semafoare distribuite a caror implementare a fost
prezentata in anexa B7
Funcitiile care pot fi apelate de programele utilizator sunt:

int ExclOrSharedAccessInitialize(void); /*conectare la sistem */
void ExclOrSharedAccessEnd(void); /* deconectare */
void ExclOrSharedAccessStart(void) /* initializare semaf.*/
void OpenSharedAccess(void); /* obtinere acces partajat */
void CloseSharedAccess(void); /* terminare acces partajat */
void OpenExclAccess(void); /* obtinere acces exclusiv */
void CloseExclAccess(void); /* terminare acces exclusiv */
*****

LRESULT CALLBACK WndProc( HWND hWnd, UINT uMsg,
                          WPARAM wParam, LPARAM lParam );

HINSTANCE hInst;
LPCTSTR lpszAppName = "MyApp";
LPCTSTR lpszTitle = "Shared-Exclusive";
int APIENTRY WinMain( HINSTANCE hInstance, HINSTANCE hPrevInstance,
                     LPTSTR lpCmdLine, int nCmdShow)
{
  MSG msg;
  HWND hWnd;
  WNDCLASS wc;
  WSADATA wSData;
  UNREFERENCED_PARAMETER( lpCmdLine );
  UNREFERENCED_PARAMETER( hPrevInstance );
}

```

```

if( GetVersion() & 0x80000000 && (GetVersion() & 0xFF) ==3){
    MessageBox( NULL, "This application cannot run on Windows
    3.1.n" "This application will now terminate.", "WinApp",
    MB_OK | MB_ICONSTOP | MB_SETFOREGROUND );
    return( 1 );
}
hInst = hInstance;
wc.style = CS_HREDRAW | CS_VREDRAW;
wc.lpszWndProc = (WNDPROC)WndProc;
wc.cbClsExtra = 0;
wc.cbWndExtra = 0;
wc.hInstance = hInstance;
wc.hIcon = LoadIcon( hInstance, lpszAppName);
wc.hCursor = LoadCursor(NULL, IDC_ARROW);
wc.hbrBackground = (HBRUSH)(COLOR_WINDOW+1);
wc.lpszMenuName = MAKEINTRESOURCE(IDR_MENU1);
wc.lpszClassName = lpszAppName;
RegisterClass( &wc );
hWnd = CreateWindow( lpszAppName,
    lpszTitle,
    WS_OVERLAPPEDWINDOW,
    CW_USEDEFAULT, 0,
    CW_USEDEFAULT, 0,
    NULL,
    NULL,
    hInstance,
    NULL
);
if ( !hWnd )
    return( FALSE );
ShowWindow( hWnd, nCmdShow );
WSAStartup(MAKEWORD(1,1), &WSAData );
if (ExclOrSharedAccessInitialize()){
    MessageBox( NULL, "This application will now terminate.",
    "WinApp", MB_OKIMB_ICONSTOPIMB_SETFOREGROUND);
    exit(1);
}
while( GetMessage( &msg, NULL, 0, 0 )
    DispatchMessage( &msg );
return( msg.wParam );
}

LRESULT CALLBACK WndProc(HWND hWnd, UINT uMsg,
    WPARAM wParam, LPARAM lParam)
{
    HDC hdc;
    char str[10];
    int nr;
    switch (uMsg)
    {
    case WM_CREATE:
        return DefWindowProc( hWnd, uMsg, wParam, lParam );
        break;
    case WM_DESTROY:
        SemaphoreEnd();
        PostQuitMessage( 0 );
        break;
    case WM_COMMAND:
        switch ( LOWORD( wParam ) )
        {
        case ID_INIT:
            ExclOrSharedAccessStart();
            break;
        case ID_SHARED:
            SetWindowText( hWnd, "Waiting for shared access" );
            OpenSharedAccess();
            hdc = GetDC(hWnd);
            TextOut(hdc,5,5,"",10);
            sprintf(str,"%d",nr);
            TextOut(hdc,5,5,str,strlen(str));
            ReleaseDC(hWnd,hdc);
            SetWindowText( hWnd, "In shared access..." );
            Sleep(10000);
            MessageBox( NULL, "Continue?", "WinApp",
            MB_OKIMB_ICONQUESTIONIMB_SETFOREGROUND );
            hdc = GetDC(hWnd);
            TextOut(hdc,5,5,"",10);
            sprintf(str,"%s", "End");
            TextOut(hdc,5,5,str,strlen(str));
            ReleaseDC(hWnd,hdc);
            CloseSharedAccess();
            hdc = GetDC(hWnd);
            TextOut(hdc,5,5,"",10);
            sprintf(str,"%d",nr);
            TextOut(hdc,5,5,str,strlen(str));
            ReleaseDC(hWnd,hdc);
            SetWindowText( hWnd, "Done shared access" );
            break;
        case ID_EXCLUSIVE:
            SetWindowText( hWnd, "Waiting for exclusive access..." );
            OpenExclAccess();
            SetWindowText( hWnd, "In exclusive access..." );
            Sleep(10000);
            MessageBox( NULL, "Continue?", "WinApp",
            MB_OKIMB_ICONQUESTIONIMB_SETFOREGROUND);
            hdc = GetDC(hWnd);
            TextOut(hdc,5,5,"",10);
            sprintf(str,"%s", "End");
            TextOut(hdc,5,5,str,strlen(str));
            ReleaseDC(hWnd,hdc);
            CloseExclAccess();
            SetWindowText( hWnd, "Done exclusive access" );
            break;
        case ID_EXIT:
            ExclOrSharedAccessEnd();
            DestroyWindow( hWnd );
            break;
        }
    }
    break;
    default:
        return DefWindowProc( hWnd, uMsg, wParam, lParam );
    }
    return 0;
}

```

ANEXA D

Implementarea sistemului DOSTP în rețele UNIX și WINDOWS

D1. Definirea clasei CObject

```

/* .....
   Definitia clasei CObject in Unix
   ..... */
class CObject :public CDebug, public CError {
public:
CObjectManager(int debug=1, FILE *f=stdout); #constructor
virtual int TreatRequest(PACKET *)=0;
virtual int PrepareToCommit(void)=0; // tratare faza prepare
virtual int Commit(void)=0; // tratare faza commit
virtual int Abort(void)=0; // tratare faza abort
virtual int Interrupt(void)=0; // nod. per. nu rasp.
virtual int Continue(int, int)=0; // recuperare A si C
virtual int ObjectInit(void)=0; // recuperare initiala
virtual int Finish(void)=0; // apel destructori
virtual void SetTransNo(int, int)=0; // setare transld si pri
virtual GetInitData(char *)=0; // obt. val obj pt depan.
virtual void Startup(void)=0; // initializari
virtual void Cleanup(void)=0; // dealocari finale
};
class CError { // pentru tratare erori
public:
int m_errCode=0; // codul erorii
virtual void Error(int)=0;
};
class CDebug( // pentru informatii la depanarea progr.
public:
int m_debug; // flag, cu depanare sau fara
FILE* m_debugFd; // fisierul in care se depune informatia
int m_pid; // pid-ul
void Debug(char*,...); // functia de scriere in fisier
CDebug(int debug=1,FILE* debugFd=stdout) // constructor
{m_pid=getpid();m_debug=debug;m_debugFd=debugFd;};
void SetPid() {m_pid=getpid();} // fixeaza pid
void Init (int debug=1,FILE* debugFd=stdout) // initializari
{m_debug=debug;m_debugFd=debugFd;};
void Flush(){fflush(m_debugFd);}; // golire buffer
};
/* .....
   Definitia clasei CObject in Windows
   ..... */
class CObject :public CDebug, public CError {
public:
CObjectManager(int debug=1, FILE *f=stdout); #constructor
virtual int TreatRequest(PACKET *)=0;
virtual int PrepareToCommit(void)=0; // tratare faza prepare
virtual int Commit(void)=0; // tratare faza commit
virtual int Abort(void)=0; // tratare faza abort
virtual int Interrupt(void)=0; // nod. per. nu rasp.
virtual int Continue(int, int)=0; // recuperare A si C
virtual int ObjectInit(void)=0; // recuperare initiala
virtual int Finish(void)=0; // apel destructori
virtual void SetTransNo(int, int)=0; // setare transld si pri
virtual GetInitData(char *)=0; // obt. val obj pt depan.
virtual void Startup(void)=0; // initializari
virtual void Cleanup(void)=0; // dealocari finale
virtual CObject * NewCopy(void)=0; // singura fct.noua
};
class CError { // pentru tratare erori
public:
int m_errCode=0; // codul erorii
virtual void Error(int)=0;
};
class CDebug(
public:
int m_debug;

```

```

FILE* m_debugFd;
int m_pid;
void Debug(char*,...);
CDebug(int debug=1,FILE* debugFd=stdout) // constructor
{m_pid=GetCurrentThreadId();
_debug=debug; m_debugFd=debugFd;};
void SetPid(){m_pid=GetCurrentThreadId();}
void Init(int debug=1,FILE* debugFd=stdout)
{m_debug=debug;m_debugFd=debugFd;};
void Flush(){fflush(m_debugFd);};
};

```

D2. Definirea clasei CVector

```

/* .....
   Definitia clasei CVector in Unix
   ..... */
class Cvector : public CObject {
// membri privati:
char m_serverName[MAX_OBJ_LEN];
char m_dataFile[MAX_PATH], m_statFile[MAX_PATH];
char m_logFilePC[MAX_PATH], m_logFileAC[MAX_PATH];
char m_tsFile[MAX_PATH]; // numele fisierelor
FILE * m_lfPC, *m_lfAC, *m_tf; //LogFile, TsFile
int m_df, m_sf; // descriptori pt dataFile
long m_logPos; // poz articolelor in fisierul PCLOG
int *m_pData; // pointer la zona de mem. cu fis de date
STATINF *m_pTrStatus; // pointer la mem. cu inf de stare
int m_semObl; // semafor pt. acc. excl. la inf. de stare
int m_semTs; // semafoare pt. acc. excl. la fis. log
int m_transid; // identificatorul tranzactiei
int m_pri; // prioritatea tranzactiei
int m_transld; // indexul in vectorul de tranz al acestei tranz
WR m_vwPos[MAX_OBJ]; // perechi index,value pt. write
int m_wld; // urm intrare ibera din m_vwPos
int m_vrPos[MAX_OBJ]; // pozitiile citite
int m_rld; // urm intrare libera din m_vrPos
int m_status; //starea tranzactiei :
// OK_PC, FAIL_PC, WORKING_PC

// metode private
void ClearLogPC(void); // rescrie PCLOG
void InitialRec(void); // recupereaza tranz. din ACLOG
int Read(VEC_REQUEST *p); // citire o pozitie
int Write(VEC_REQUEST *p); // scriere o pozitie
int AllocInVectTz(void); // aloca o intr. dn vect. cu tr. active
int ReleaseInVectTz(int); // elibereaza o intr. din vect. de tr.
void RecordOldTz(void); // act. tz. old citind PCLOG
void OldTzStart(int index); // act. starea unei tz. old la WORK
void OldTzDone(int index); // act. starea unei tz. old la DONE
short OldTzStat(int idTranz); // return. starea unei tz old
void InOldTz(void); // obt. acc. excl. cu m_semTs
void OutOldTz(void); // elib. acc. excl. cu m_semTs
long LogPcAddr(int idTz, int *index, int *prior); // ret. adr. din PC
int inTF(void); // ret. daca tz. a fost terminata
int inVectTz(void); // ret. daca tz este in lucru
void ExclAcces(void); // obt. acc. excl. cu m_semObl
void ReleaseAcces(void); // elib. acc. excl. cu m_semObl
void SemWait(int); // astept. la sem. param.
void SemSignal(int); // Up la semaf. param.
off_t DimFile(int ); // determ. dim. fis. de date
void DCVector(void); // rol de destructor
void PrintStatus(void); // inf. pt. depanare

public:
INCLUDE_SYSTEM (CVector)
CVector (char *name, char *dataFile, char *logFilePC, char *
logFileAC, char * tzFile, char * statFile );
CVector();

```

```

virtual int TreatRequest(PACKET *); // tratarea unui pckt SHM
virtual int PrepareToCommit(void); // faza prepare
virtual int Commit(void); // terminare anormal tr.
virtual int Abort(void); // terminare anormal tr.
virtual int Interrupt(void); // pentru cind TM nu mai rasp.
int Continue(int tz,int op); // recuperarea tz. respective
int ObjectInit(); // recuperarea si rescrie PCLOG
virtual void SetTransNo(int transNo,int pri)
    (m_transid=transNo;m_pri=pri); // setare transld, pri
virtual void Startup(void); // initializari
virtual void Cleanup(void); // dealocari
virtual void Finish(void); // dealocari finale
void GetBaseData(char *); // afiseaza pt. depanare val. ob.
void Error(int); // afisarea erorilor
};

/* .....
    Definitia clasei CVector in Windows
    ..... */
class Cvector : public CObject {
// membri privati
char m_serverName[MAX_OBJ_LEN], m_statFile[MAX_PATH];
char m_logFilePC[MAX_PATH], m_logFileAC[MAX_PATH];
char m_tsFile[MAX_PATH], // numele fisierelor
FILE * m_ifPC,*m_ifAC,*m_tf; // LogFile, TsFile
long m_logPos; // poz articolelor in fisierul PCLOG
HANDLE m_df; //handle pt dataFile
HANDLE m_hData; // handle al zonei de mem partajata pt ob
int *m_pData; // pointer la zona de mem. cu fis de date
STATINF *m_pTrStatus; // pointer la zona de mem. cu inf de stare
HANDLE m_semObld; // pt acc. excl. inf de stare
HANDLE m_semTz; // pt acc. excl. la inf. de stare pt. tz old
int m_transid; // identificatorul tranzactiei
int m_pri; // prioritatea tranzactiei
int m_transidx; // indexul in vectorul de tranz al acestei tranz
WR m_vWPos[MAX_OBJ]; // perchi index,value pt. write
int m_wldx; // urm intrare ibera din m_vWPos
int m_vRPos[MAX_OBJ]; // pozitiile citite
int m_rldx; // urm intrare libera din m_vRPos
int m_status; //starea tranzactiei :
// OK_PC, FAIL_PC, WORKING_PC

// metode private
void ClearLogPC(void); // rescrie PCLOG
void InitialRec(void); // recupereaza tranz. din ACLOG
int Read(VEC_REQUEST *p); // citire o pozitie
int Write(VEC_REQUEST *p); // scriere o pozitie
int AllocInVectTz(void); // aloca o intr. dn vect. cu tr. active
int ReleaseInVectTz(int); // elibereaza o intr. din vect. de tr.
void RecordOldTz(void); // act. tz. old citind PCLOG
void OldTzStart(int index); // act. starea unei tz. old la WORK
void OldTzDone(int index); // act. starea unei tz. old la DONE
short OldTzStat(int idTranz); // return. starea unei tz old
void InOldTz(void); // obt. acc. excl. cu m_semTs
void OutOldTz(void); // elib. acc. excl. cu m_semTs
long LogPcAddr(int idTz, int *index,int *prior); // ret. adr. din PC
int InTF(void); // ret. daca tz. a fost terminata
int InVectTz(void); // ret. daca tz este in lucru
void ExclAcces(void); // obt. acc. excl. cu m_semObld
void ReleaseAcces(void); // elib. acc. excl. cu m_semObld
void SemWait(HANDLE); // astept. la sem. param.
void SemSignal(HANDLE); // Up la semaf. param.
off_t DimFile(int); // determ. dim. fis. de date
void DCVector(void); // rol de destructor la Quit
void PrintStatus(void); // inf. pt. depanare

public:
INCLUDE_SYSTEM (CVector)
Cvector (char *name, char *dataFile, char * logFilePC,char *
logFileAC,char * tzFile, char * statFile); // constructor
Cvector();
CBase * NewCopy(void); // creeaza dinamic o alta instanta a ob.
virtual int TreatRequest(PACKET *); // tratarea unui pckt SHM
virtual int PrepareToCommit(void); // faza prepare
virtual int Commit(void); // terminare anormal tr.
virtual int Abort(void); // terminare anormal tr.
virtual int Interrupt(void); // pentru cind TM nu mai rasp.
int Continue(int tz,int op); // recuperarea tz. respective
int ObjectInit(void); // recuperarea si rescrie PCLOG
virtual void SetTransNo(int transNo,int pri)
    {m_transid=transNo;m_pri=pri}; // setare transld, pri
virtual void Startup(void); // initializari inainte de rec. pckt SHM
virtual void Cleanup(void); // dealocari dupa abort, commit SHM
virtual void Finish(void); // dealocari finale
void GetBaseData(char * str); // ret. valorile ob. pt. depanare
void Error(int); // afisarea erorilor
};

D3. Definirea clasei CObjectManager

/* .....
    Definitia clasei CObjectManager in Unix
    ..... */
class CObjectManager: public CError,public CDebug{
// membri privati
Widget m_label,m_quit,m_box;
XtInpuitld m_sockld;
int m_aut; // flag care indica daca ob. s-a conectat sau nu
int m_finished; // folosit pt. a evita dubla dealocare:quit si ies.
int m_transNo; // numarul asociat tranzactiei
CClientList m_clientList; // lista tranzactiilor ce acceseaza ob.
char pathName[MAX_PATH]; // nume pt. dir. curent
struct sockaddr_un m_serv_addr,m_cli_addr; // adrese
int m_clilen; // dim. adresei
int m_chan_U_OBJ; // canalul static (beg. abort, com)
int m_chan_f_CS; // pentru comunic. cu CS
BT_OBJECT m_mes; // cererea curenta SNextRequest*/

public:
CObject * m_pObject; // pointer la ob. gestionat
char m_objName[MAX_OBJ_LEN];
void Xinit(XtAppContext context,Widget pane);
int Init (XtAppContext context, Widget pane, char* name, int
debug=1,FILE * debugFd=stdout);
int ConnectObjectToSystem(void); // fctia pt. conectarea ob.
int CreateSock(void); // crearea canalului propriu
int RecordObj(void); // inregistreaza ob. la CM
void DisconnectObjectFromSystem(void); // deconectare sist.
void TransactionList(void); // afisarea tranzactiilor active
int AddClient(char * name,int transld,int pri); // adauga o tr.
int RemoveClient(char * name); // elimina un client name
int RemoveClient(int transld); // elimina un client transld
int SNextRequest(void); // trateaza urm.cerere din socket
void Error(int); // traterea erorilor
int GetTzNo(){return m_tzNo;}; // ret. nr tranzactie
void DCObjectManager(void); // simuleaza destructor
SigFunc * signal(int signo,SigFunc * func);
friend void sig_chld(int signo);
friend void Quit(Widget w,XtPointer client_data, XtPointer
call_data);

friend int SNextRequest(XtPointer,int * ,XtInpuit *); //
trateaza urmatoarea cerere care vine in canalul static
};

/* .....
    Definitia clasei CObjectManager in Windows
    ..... */
class CObjectManager: public CError,public CDebug{
// membri privati
int m_aut; // flag care indica daca ob. s-a conectat sau nu
int m_finished; // folosit pt. a evita dubla alocare: quit si ies.
CClientList m_clientList; // lista tranzactiilor ce acceseaza ob.
HANDLE m_listMutex; // semaf. pentru acc. excl.
SOCKADDR_IN m_serv_addr,m_cli_addr; // adrese ln et
int m_clilen; // dim. adresei
SOCKET m_chan_K_OBJ; // canalul static (beg. abort, com)
SOCKET m_chan_K_CS; // pentru comunic. cu CS
SOCKET m_newsocKfd; // canal dinamic
u_short m_myPort; // portul asociat managerului de ob.
HANDLE m_copyMutex,m_gCopyMutex; // pt. copierea inf.
BT_OBJECT m_mes; // cererea curenta SNextRequest*/

```



```
public:
CObject * m_pObject; // pointer la ob. gestionat
char m_objName[MAX_OBJ_LEN];
int Init (CObject *pServer,char* name,int port,HANDLE
gCopyMutex, int debug, FILE* debugFd); // initializari
int ConnectObjectToSystem(void); // fctia pt. conectarea ob.
int CreateSock(void); // crearea canalului propriu
int RecordObj(void); // inregistreaza ob. la CM
void DisconnectObjectFromSystem(void); // deconectare sist.
void TransactionList(void); // afisarea tranzactiilor active
int AddClient(char * name,int transld,int pri); // adauga o tr.
int RemoveClient(char* name); // elimina un client name
int RemoveClient(int transld); // elimina un client transld
int SHMServer(void); // creaza un fir pt. comun. SHM
int RecoveryServer(void); // creaza un fir pt. recuperare
int SNextRequest(void) // trateaza urm.cerere din socket
BT_OBJECT GetMessage(void); // ret. mes. primi t in socket
HANDLE GetCopyMutex(void); // ret. handle-ul pt CopyM
HANDLE GetGCopyMutex(void); // ret. handle-ul pt gCopyM
SOCKET GetNewSockFd(void); // ret. m_newsocfd
CObject * GetPObject(void); // ret. m_pObject
void Error(int); // tratarea erorilor
void DObjectManager(void); // simuleaza destructor
};
```

D4. Definirea clasei CDispatcher

```
/* .....
Definitia clasei CDispatcher in Unix
..... */
class CDispatcher: public CError , public CDebug{
// membri privati
int m_transld; // identificatorul tranzactiei
int m_pri; // prioritatea tranzactiei
char m_path[MAX_PATH]; // cale socket unix
struct sockaddr_un m_serv_addr; // adresa
CObject *m_pObject; // pointer la obiect ;
BT_OBJECT m_mes; // cererea curenta SNextRequest
PACKET m_pcket; // pckt primit in SHM
int m_semld; // id. pt. setul de semaf. pt. SHM
int m_semNr; // nr. semaf. pt. obj. din set
int m_shmld; // id. pt. mem. partajata
SHM * m_pShm; // pointer la mem. partajata
int m_queueNo; // nr. cozii rezervate in tranz. pt. ob.
int m_queueIdx ; // indicele pkt-ului urm. din coada
int m_waiting; // flag pentru bucla de asteptare

public:
int Init(int debug=1,FILE *debugFd=stdout) // initializari
{m_queueIdx =0;m_waiting=1;CDebug::Init(debug,debugFd);}
int WrongTzNo(void); // verifica nr. tz. din pckt cu m_transld
void SetPacketErr(int code){m_pckt.text.errCode=code ;}
void MainJob(void); // bucla principala a serverului
int Attach(void); // atasare la SHM
void Detatch(void){ shmctl((char *)m_pShm);}
int ResolveTrans(BT_OBJECT *char *) ; // gestioneaza tr.
int WaitForAMessage(void); // asteptarea unui pckt in SHM
void TreatMessage(void); // tratare mesaje
void FirstRequest(void); // tratare prima cerere la ob
void Request(void); // tratare cerere
void PrepareToCommitTr(void); // tratare prepare
void CommitTr(void); // tratare commit
int AbortTr(void); // tratare abort
int InterruptTr(void); // tratare interrupt
int SendPckt(void); // emisie pckt la CS prin SHM
int Up_S(int); // Up la semaf.param
int Down_S(int); // Down la semaf.param
void FreeSlot(int); // Elibereaza un slot din SHM
int GetSlot(void); // Aloca un slot in SHM
void SetClasDeriv (CObject * athis){m_pServer=athis;}
void SetTransNo(int); // stabileste nr. tranzactiei
void FinishTr(void){m_pObject->Finish();}
void GetBaseData(char* p){m_pObject->GetBaseData(p);}
void RecoveryTr(BT_OBJECT* int ,char *) ; // recupere tr.
void Error(int code); // tratarea erorilor
};
```

```
/* .....
Definitia clasei CDispatcher in Windows
..... */
class CDispatcher: public CError , public CDebug{
// membri privati
int m_transld; // identificatorul tranzactiei
int m_pri; // prioritatea tranzactiei
CObject *m_pObject; // pointer la obiect ;
BT_OBJECT m_mes; // cererea curenta SNextRequest
PACKET m_pcket; // pckt primit in SHM
SHM * m_pShm; // pointer la mem. partajata
HANDLE m_sShm, m_sCs, m_sOb, m_sSlot, m_hShm
int m_queueNo; // nr. cozii rezervate in tranz. pt. ob.
int m_queueIdx ; // indicele pkt-ului urm. din coada
int m_waiting; // flag pentru bucla de asteptare
CObjectManager * m_pManager; // pointer la manager

public:
int Init(int debug=1,FILE *debugFd=stdout) // initializari
{m_queueIdx =0;m_waiting=1;CDebug::Init(debug,debugFd);}
int WrongTzNo(void); // verifica nr. tz. din pckt cu m_transld
void SetPacketErr(int code){m_pckt.text.errCode=code ;}
void MainJob(void); // bucla principala a serverului
int Attach(void); // atasare la SHM
void Detatch(void){
UnmapViewOfFile( LPSTR) m_pShm );
CloseHandle(m_sShm);
CloseHandle(m_sMC);
CloseHandle(m_sOB);
CloseHandle(m_sSlot);
CloseHandle(m_hShm);
}
int ResolveTrans(BT_OBJECT *char *) ; // gestioneaza tr.
int WaitForAMessage(void); // asteptarea unui pckt in SHM
void TreatMessage(void); // tratare mesaje
void FirstRequest(void); // tratare prima cerere la ob
void Request(void); // tratare cerere
void PrepareToCommitTr(void); // tratare prepare
void CommitTr(void); // tratare commit
int AbortTr(void); // tratare abort
int InterruptTr(void); // tratare interrupt
int SendPckt(void); // emisie pckt la CS prin SHM
int Up_S(int); // Up la semaf.param
int Down_S(int); // Down la semaf.param
void FreeSlot(int); // Elibereaza un slot din SHM
int GetSlot(void); // Aloca un slot in SHM
void SetClasDeriv (CObject * athis){m_pServer=athis;}
void SetTransNo(int); // stabileste nr. tranzactiei
void FinishTr(void){m_pObject->Finish();}
void GetBaseData(char* p){m_pObject->GetBaseData(p);}
void RecoveryTr(BT_OBJECT* int ,char *) ; // recupere tr.
void Error(int code); // tratarea erorilor
void SendMessageToParent(int);
void SetPMng(CObjectManager* pMng){m_pManger=pMng;}
};
```

D5. Definierea macroinstrucțiunilor pentru conectarea automată a obiectelor în sistemul DOSTP

```
/*.....
      Definierea macroinstrucțiunilor în Unix
      *
.....*/
#define INCLUDE_SYSTEM(className)\
  ObjectManager * m_master;\
  void friend Init(className* p, XtAppContext, Widget );

#define DEFINE_SYSTEM(className)\
  void Init(className* s, XtAppContext context, Widget pane)\
  { ObjectManager* master=new(ObjectManager);\
    s->m_master=master;\
    master->m_pObject=s;\
    master->Init(context, pane, s->m_serverName, s->m_debug,\
                s->m_debugFd);\
    s->Debug(" Master executed Init ");\
  }

#define START_SYSTEM(server, context, pane) \
  Init(server, context, pane)

#define STOP_SYSTEM(obj) \
  obj->m_master->DisconnectObjectFromSystem(void)

/*.....
      Definierea macroinstrucțiunilor în Windows
      *
.....*/

#define INCLUDE_SYSTEM(className)\
  u_short m_port;\
  HANDLE m_copyMutex;\
  void friend Init(className* pObject, u_short port);\
  void friend ThreadMaster(className *);

#define DEFINE_SYSTEM(className)\
  void Init(className* pObject, u_short port)\
  { pObject->m_port=port;\
    char name[MAX_OBJ_LEN];\
    strcpy(name, pObject->m_serverName);\
    strcat(name, "copyMutex");\
    pObject->m_copyMutex=CreateSemaphore(NULL, 0, 1, name);\
    if (!pObject->m_copyMutex)\
      return ;\
  HANDLE h;\
  DWORD idT;\
  h=CreateThread((NULL, 0, (LPTHREAD_START_ROUTINE) \
                ThreadMaster, pObject, 0, &idT);\
  if (!h) {MsgBox("Create Thread error \n");}\
  }\
  void ThreadMaster(className *pObject)\
  {\
    ObjectManager master;\
    master.Init(pObject, pObject->m_serverName,\
              pObject->m_port, pObject->m_copyMutex,\
              pObject->m_debug, pObject->m_debugFd);\
    HANDLE h; \
    DWORD idT; \
    h=CreateThread (NULL, 0, (LPTHREAD_START_ROUTINE)\
                  ThreadDisplay, &master, 0, &idT);\
    if (!h) {MsgBox(" Create Thread error n");}\
    while(TRUE)\
      master.SNextRequest();\
  }

#define START_SYSTEM(pObject, port) \
  Init(pObject, port)

#define OUT_SYSTEM(obj) \
  obj->m_master.DisconnectObjectFromSystem(void)
```

ANEXA E

Comparație între performanțele tranzacțiilor distribuite și cele ale tranzacțiilor locale în sistemul DOSTP

Această anexă prezintă comparativ performanțele tranzacțiilor distribuite față de cele locale. Pentru testare s-au folosit 3 manageri de obiecte, într-o prima variantă fiecare pe câte o stație și într-o altă variantă toți local, pe aceeași stație. S-a folosit același program client însă utilizând 1,2,3,4,5 tranzacții care implică cele 3 obiecte. S-au afișat următoarele valori (vezi tabelul următor): timpul total cât a durat execuția programului client, timpul utilizator în care s-a executat programul client și timpul sistem în care s-a executat programul client. Timpul a fost reprezentat în tick-uri de 1/100 sec. Rezultatele sunt descrise în tabelul următor.

1) Tranzacții asincrone locale

	(1)	(2)	(3)	(4)	(5)
Timp total	2,192	3,372	4,347	5,247	6,582
Timp utilizator	1	2	2	4	4
Timp sistem	4	6	10	14	16

2) Tranzacții sincrone locale

	(1)	(2)	(3)	(4)	(5)
Timp total	2,258	3,439	4,713	5,450	6,626
Timp utilizator	1	2	2	3	3
Timp sistem	4	6	12	15	18

3) Tranzacții asincrone distribuite

	(1)	(2)	(3)	(4)	(5)
Timp total	952	1,369	1,896	2,333	2,797
Timp utilizator	1	2	2	2	2
Timp sistem	4	8	11	13	16

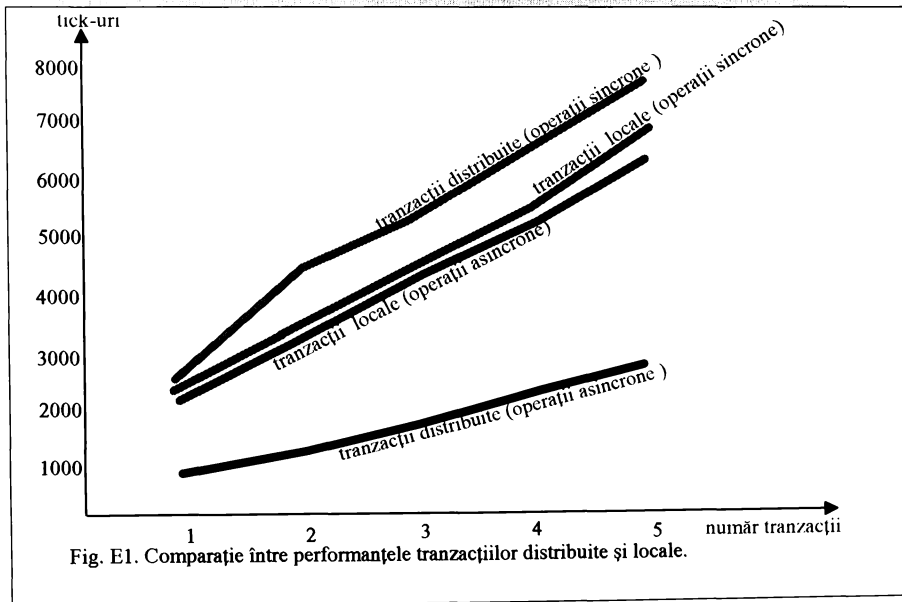
4) Tranzacții sincrone distribuite

	(1)	(2)	(3)	(4)	(5)
Timp total	2,557	4,574	4,980	6,788	7,243
Timp utilizator	2	2	2	3	3
Timp sistem	5	6	13	17	19

Anexa E: Comparație între performanțele tranzacțiilor distribuite și cele ale tranzacțiilor locale în DOSTP

Se observă că:

- (1) Timpul necesar realizării tranzacțiilor distribuite sincrone este cel mai mare;
 - (2) Timpul necesar realizării tranzacțiilor distribuite asincrone este cel mai mic;
 - (3) Timpul necesar realizării tranzacțiilor care se execută local este aproximativ de 2.5 (numărul de obiecte implicate în tranzacție este 3 dar intervine și timpul necesar protocolului two phase commit) ori mai mare decât timpul necesar realizării aceluiași tranzacții pe alte stații.
- Aceste concluzii sunt reprezentate în figura următoare.



ANEXA F

Comparație între performanțele variantelor distribuite și semidistribuite ale sistemului TGC

Accastă anexă prezintă rezultatele obținute în urma testării sistemului TGC. Testele s-au realizat prin execuția unor programe client care trimit în buclă continuă mesaje către un program client receptor. Pentru comparație, programele emițătoare s-au executat în mai multe variante: pe aceeași stație cu programul emițător (1, 2, 3 și 4 programe emițătoare) și, la distanță, pe 1, 2, 3 sau 4 stații; tot pentru comparație, s-au înregistrat rezultatele pentru trimiterea a 50, 100, 200, 300, 400, 500 și 1000 mesaje. S-au afișat următoarele valori (vezi tabelul următor):

- pentru programele client emițătoare: (1) numărul de mesaje trimise; (2) numărul de mesaje retrimise în urma unei erori de tipul " server de comunicație ocupat"; (3) timpul total cât a durat execuția; (4) timpul utilizator în care s-a executat programul; (5) timpul sistem în care s-a executat programul;

- pentru programul receptor: (6) numărul de mesaje recepționate; (7) numărul de mesaje pierdute; (8) timpul total cât a durat execuția; (9) timpul utilizator în care s-a executat programul; (10) timpul sistem în care s-a executat programul. Timpul a fost reprezentat în tick-uri de 1/100 sec.

Rezultatele sunt descrise în tabelele următoare.

1. Varianta distribuită cu transmisie broadcast

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)
Un client la distanță	50	0	31	4	3	50	0	31	0	1
	100	0	62	1	2	100	0	61	2	3
	200	0	128	1	7	200	0	127	5	3
	300	0	200	8	9	300	0	198	3	7
	400	0	301	12	10	400	0	308	8	9
	500	0	391	8	9	500	0	389	11	9
2 clienți la distanță	1,000	0	949	22	29	1,000	0	948	15	9
	50	0	36	2	2	100	0	176	3	0
	100	0	81	7	2	200	0	304	5	7
	200	0	209	5	5	400	0	518	7	14
	300	0	483	6	5	600	0	862	10	9
	400	0	698	8	10	800	0	1,245	16	23
3 clienți la distanță	500	0	1,119	14	11	1,000	0	1,610	27	27
	1,000	0	2,712	31	28	2,000	0	2,900	46	60
	50	0	38	1	1					
	100	0	89	1	1					
	200	0	247	5	8					
	300	0	597	6	8					
3 clienți la distanță	400	0	848	7	9					
	500	0	1,309	7	10					
	1,000	0	2,823	27	27					
	50	0	59	1	1	150	0	118	3	4
	100	0	225	1	3	300	0	276	4	3
	200	0	749	1	4	600	0	796	13	5
3 clienți la distanță	300	0	769	10	8	900	0	1,640	18	21
	400	0	844	12	17	1,200	0	1,822	26	36
	500	0	1,449	16	17	1,500	0	2,460	23	22
	1,000	0	3,705	28	27	3,000	0	6,236	47	70
	50	0	62	0	4					
	100	0	228	1	2					
3 clienți la distanță	200	0	738	1	3					
	300	1,467	1,325	31	41					
	400	977	1,599	29	26					
	500	1,287	2,223	21	34					
	1,000	2,749	5,526	80	93					

Anexa F: Comparație între performanțele variantelor distribuite și semidistribuite ale sistemului TGC

	50	0	64	1	1					
	100	0	166	4	4					
	200	0	682	4	12					
	300	1,528	1,327	36	50					
	400	0	1,062	30	30					
	500	1,315	3,261	37	37					
	1,000	2,758	5,673	69	68					
4 clienti la distanță	50	0	25	2	0	200	0	224	5	4
	100	0	52	2	2	400	0	328	6	11
	200	58	258	8	4	800	0	745	13	23
	300	578	838	19	23	1,200	0	1,763	28	18
	400	1,517	2,901	37	35	1,600	0	4,692	44	40
	500	5,681	4,568	112	168	2,000	0	5,601	43	52
	1,000	18,243	14,356	349	510	4,000	0	14,976	105	192
	50	0	25	4	0					
	100	0	68	1	2					
	200	77	348	8	7					
300	583	1,417	21	26						
400	1,578	4,360	44	45						
500	0	3,659	22	25						
1,000	13,002	10,854	262	362						
50	0	38	1	0						
100	0	90	3	2						
200	0	282	3	1						
300	877	1,367	32	34						
400	0	3,204	7	11						
500	5,660	5,244	121	110						
1,000	16,044	14,672	235	498						
50	0	33	0	1						
100	0	87	2	1						
200	0	360	4	4						
300	807	1,439	18	38						
400	0	3,466	10	9						
500	6,566	5,246	131	190						
1,000	19,595	14,672	325	565						
Un client local	50	0	22	1	1	50	0	21	1	1
	100	0	50	3	4	100	0	50	1	2
	200	0	106	6	0	200	0	105	1	5
	300	0	189	5	6	300	0	189	5	8
	400	0	250	13	11	400	0	251	8	8
	500	0	336	9	15	500	0	336	9	4
	1,000	0	717	12	17	1,000	0	717	21	26
	50	0	22	1	1	100	0	140	1	3
	100	0	53	1	2	200	0	157	2	6
	200	0	115	2	0	400	0	362	10	8
300	0	206	5	7	600	0	429	7	14	
400	0	253	4	7	800	0	584	18	19	
500	212	512	17	15	1,000	0	786	19	40	
1,000	1,034	1,705	48	50	2,000	0	2,119	34	25	
50	0	30	0	2						
100	0	59	3	2						
200	0	118	3	5						
300	0	206	5	10						
400	0	292	6	8						
500	0	555	15	16						
1,000	0	2,823	44	60						
3 clienti locali	50	0	21	2	1	150	0	251	3	3
	100	0	50	3	1	300	0	310	6	5
	200	0	115	3	2	600	0	625	11	13
	300	0	203	5	9	900	0	792	13	18
	400	61	333	10	12	1,200	0	951	23	38
	500	279	673	13	22	1,500	0	1,250	40	35
	1,000	1,056	2,510	21	39	3,000	0	3,864	54	75

anexa F-2

Anexa F: Comparație între performanțele variantelor distribuite și semidistribuite ale sistemului TGC

	50	0	26	3	1					
	100	0	50	3	3					
	200	0	133	8	6					
	300	0	208	7	9					
	400	256	535	16	10					
	500	249	620	19	17					
	1,000	1,866	3,510	39	41					
	50	0	21	0	1					
	100	0	47	1	2					
	200	0	127	8	3					
	300	0	212	6	8					
	400	191	463	13	3					
	500	0	232	4	9					
	1,000	1,778	3,414	37	42					
4 clienti locali	50	0	21	1	1	200	0	346	2	7
	100	0	55	0	3	400	0	450	7	10
	200	0	125	4	4	800	0	732	11	13
	300	0	251	7	5	1,200	0	1,171	24	30
	400	78	343	10	5	1,600	0	1,506	24	32
	500	93	439	13	8	2,000	0	2,254	25	41
	1,000	509	2,181	34	28	4,000	0	5,164	90	98
	50	0	29	0	2					
	100	0	55	2	3					
	200	0	126	5	4					
	300	0	254	6	7					
	400	402	666	10	22					
	500	646	1,188	21	24					
	1,000	2,215	4,642	50	60					
	50	0	22	0	1					
	100	0	59	3	4					
	200	0	121	3	4					
	300	0	221	5	13					
	400	538	823	19	16					
	500	959	1,574	20	27					
	1,000	2,120	4,631	52	53					
	50	0	28	0	2					
	100	0	54	1	1					
	200	0	132	1	3					
	300	0	232	4	9					
	400	220	581	19	24					
	500	767	1,440	14	23					
	1,000	1,381	4,060	25	34					

1. Varianta distribuită cu transmisie în inel

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)
Un client la distanță	50	0	147	0	1	50	0	145	1	0
	100	0	290	4	2	100	0	288	1	0
	200	0	580	11	9	200	0	578	1	2
	300	0	862	3	10	300	0	859	3	2
	400	0	1,144	5	4	400	0	1,142	4	3
	500	0	1,440	7	1	500	0	1,440	7	1
	1,000	0	2,868	23	10	1,000	0	2,865	5	11
2 clienti la distanță	50	0	212	2	0	100	0	347	5	2
	100	0	442	0	4	200	0	534	8	6
	200	0	886	9	4	400	0	970	10	14
	300	0	1,381	13	13	600	0	1,592	12	15
	400	0	1,854	7	7	800	0	2,017	13	30
	500	0	2,284	9	11	1,000	0	2,451	22	25
	1,000	0	4,290	15	14	2,000	0	4,489	53	66
	50	0	206	1	1					
	100	0	435	0	4					
	200	0	880	0	5					

Anexa F: Comparație între performanțele variantelor distribuite și semidistribuite ale sistemului TGC

	300	0	1,375	8	15					
	400	0	1,854	9	9					
	500	0	2,280	7	6					
	1,000	0	4,270	6	8					
3 clienti la distanță	50	0	292	2	0	150	0	456	3	4
	100	0	573	3	6	300	0	928	6	8
	200	0	1,176	7	4	600	0	1,485	5	9
	300	0	1,775	7	4	900	0	2,121	13	15
	400	0	2,355	11	13	1,200	0	2,675	14	16
	500	0	2,977	12	13	1,500	0	3,433	23	21
	1,000	0	5,951	18	13	3,000	0	7,702	44	58
	50	0	293	1	0					
	100	0	582	1	3					
	200	0	1,182	2	3					
	300	0	1,784	8	10					
	400	0	2,357	9	10					
	500	0	2,926	10	11					
	1,000	0	6,060	34	20					
	50	0	288	1	1					
	100	0	578	4	1					
	200	0	1,165	2	3					
	300	0	1,763	4	4					
	400	0	2,347	2	3					
	500	0	2,981	3	5					
	1,000	0	6,058	15	23					
4 clienti la distanță	50	0	275	1	3	200	0	767	5	7
	100	0	570	3	1	400	0	1,119	5	14
	200	0	1,161	2	3	800	0	1,684	13	19
	300	0	1,746	3	8	1,200	0	2,493	27	41
	400	0	2,245	6	4	1,600	0	3,146	31	47
	500	0	3,019	8	4	2,000	0	3,628	46	41
	1,000	0	6,263	15	10	4,000	0	7,131	60	73
	50	0	286	1	1					
	100	0	582	0	1					
	200	0	1,191	12	2					
	300	0	1,774	13	9					
	400	0	2,453	13	10					
	500	0	3,043	12	13					
	1,000	0	6,289	30	28					
	50	0	296	1	2					
	100	0	577	6	6					
	200	0	1,193	5	2					
	300	0	1,769	4	8					
	400	0	2,455	9	5					
	500	0	3,048	7	4					
	1,000	0	6,263	14	19					
	50	0	287	3	5					
	100	0	567	5	9					
	200	0	1,175	2	6					
	300	0	1,758	8	2					
	400	0	2,422	6	9					
	500	0	3,020	4	7					
	1,000	0	6,189	17	16					
Un client local	50	0	25	0	2	50	0	20	0	1
	100	0	35	2	3	100	0	33	1	3
	200	0	85	4	5	200	0	84	4	6
	300	0	161	12	10	300	0	161	3	9
	400	0	196	13	11	400	0	197	5	15
	500	0	264	14	11	500	0	265	8	13
	1,000	0	554	28	28	1,000	0	555	16	18
2 clienti locali	50	0	28	1	3	100	0	39	2	5
	100	0	52	4	0	200	0	83	1	3
	200	0	89	5	3	400	0	235	4	9

anexa F-4

Anexa F: Comparatie între performantele variantelor distribuite și semidistribuite ale sistemului TGC

	300	0	132	6	6	600	0	307	16	21
	400	0	183	5	13	800	0	416	13	28
	500	0	366	13	19	1,000	0	576	24	32
	1,000	0	975	17	25	2,000	0	1,248	33	39
	50	0	28	0	3					
	100	0	62	3	4					
	200	0	98	7	6					
	300	0	137	8	8					
	400	0	210	10	18					
	500	0	403	8	13					
	1,000	0	979	23	18					
3 clienti locali	50	0	16	1	1	150	0	248	3	5
	100	0	50	1	0	300	0	292	12	16
	200	0	106	3	3	600	0	530	19	19
	300	0	185	9	4	900	0	718	30	14
	400	0	256	9	10	1,200	0	863	19	37
	500	0	387	15	9	1,500	0	927	28	33
	1,000	71	997	19	27	3,000	0	1,783	55	64
	50	0	16	1	1					
	100	0	36	1	2					
	200	0	95	1	3					
	300	0	199	6	6					
	400	0	265	7	12					
	500	0	457	8	14					
	1,000	71	1,221	20	74					
	50	0	32	3	1					
	100	0	41	2	1					
	200	0	104	3	6					
	300	0	168	6	8					
	400	0	221	10	7					
	500	0	328	6	12					
	1,000	73	1,002	27	21					
4 clienti locali	50	0	22	3	1	200	0	339	5	5
	100	0	38	1	1	400	0	444	8	5
	200	0	103	5	4	800	0	623	25	24
	300	0	195	7	10	1,200	0	781	29	21
	400	0	281	9	11	1,600	0	1,089	31	33
	500	0	313	12	16	2,000	0	1,369	32	44
	1,000	234	1,357	18	32	4,000	0	3,069	67	84
	50	0	22	2	1					
	100	0	64	1	3					
	200	0	118	4	5					
	300	0	153	6	7					
	400	0	297	5	7					
	500	0	335	9	12					
	1,000	677	2,293	36	54					
	50	0	24	1	2					
	100	0	86	4	1					
	200	0	114	4	4					
	300	0	286	4	11					
	400	0	244	6	7					
	500	0	398	8	11					
	1,000	678	2,280	41	37					
	50	0	25	1	1					
	100	0	79	1	1					
	200	0	97	5	3					
	300	0	286	7	4					
	400	0	241	5	8					
	500	0	412	8	21					
	1,000	485	1,829	42	43					

Anexa F: Comparatie între performanțele variantelor distribuite și semidistribuite ale sistemului TGC

1. Varianta semidistribuită

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)
Un client la	50	0	82	0	0	50	0	82	0	0
distanță	100	0	190	2	1	100	0	190	0	2
	200	0	454	1	2	200	0	454	1	2
	300	0	729	1	2	300	0	728	1	2
	400	0	956	1	2	400	0	954	0	2
	500	0	1,129	5	3	500	0	1,440	2	5
	1,000	0	2,321	9	14	1,000	0	2,865	4	11
2 clienti la	50	0	71	4	1	100	0	220	5	2
distanță	100	0	235	0	2	200	0	338	8	16
	200	0	554	0	3	400	0	970	10	14
	300	0	918	6	8	600	0	978	12	15
	400	0	1,353	3	2	800	0	1,372	19	20
	500	0	1,606	4	5	1,000	0	1,702	26	20
	1,000	0	3,331	5	5	2,000	0	3,354	65	55
	50	0	73	0	1					
	100	0	276	0	1					
	200	0	484	2	2					
	300	0	904	0	2					
	400	0	1,319	9	17					
	500	0	1,528	2	7					
	1,000	0	3,232	1	5					
3 clienti la	50	0	123	0	0	150	0	230	6	2
distanță	100	0	316	0	1	300	0	743	6	8
	200	0	767	1	1	600	0	859	13	8
	300	0	1,251	1	0	900	0	1,289	14	20
	400	0	1,594	0	2	1,200	0	1,616	20	33
	500	0	1,936	4	2	1,500	0	2,309	23	36
	1,000	0	4,303	9	13	3,000	0	4,477	44	58
	50	0	124	0	1					
	100	0	316	0	3					
	200	0	721	1	4					
	300	0	1,261	5	2					
	400	0	1,619	11	12					
	500	0	2,057	8	15					
	1,000	0	4,262	9	10					
	50	0	92	0	0					
	100	0	301	0	1					
	200	0	730	1	1					
	300	0	1,141	2	3					
	400	0	1,455	1	3					
	500	0	1,939	2	3					
	1,000	0	4,197	4	7					
Un client	50	0	100	0	0	50	0	98	0	0
local	100	0	201	0	0	100	0	198	0	0
	200	0	400	0	0	200	0	397	0	0
	300	0	600	0	0	300	0	598	0	0
	400	0	802	0	0	400	0	799	0	0
	500	0	1,000	0	0	500	0	998	0	0
	1,000	0	2,002	0	0	1,000	0	2,000	0	0
2 clienti	50	0	139	0	0	100	0	197	0	0
locali	100	0	356	0	0	200	0	396	0	0
	200	0	765	0	0	400	0	798	0	0
	300	0	1,134	0	0	600	0	1,196	0	0
	400	0	1,572	0	0	800	0	1,594	0	0
	500	0	1,935	0	0	1,000	0	1,996	0	0
	1,000	0	3,947	0	0	2,000	0	3,998	0	0
	50	0	140	0	0					
	100	0	358	0	0					
	200	0	764	0	0					
	300	0	1,136	0	0					

Anexa F: Comparație între performanțele variantelor distribuite și semidistribuite ale sistemului TGC

	400	0	1,576	0	0					
	500	0	1,937	0	0					
	1,000	0	3,947	0	0					
3 clienti	50	0	167	0	0	150	0	265	0	0
locali	100	0	369	0	0	300	0	457	0	0
	200	0	776	0	0	600	0	855	0	0
	300	0	1,151	0	0	900	0	1,259	0	0
	400	0	1,527	0	0	1,200	0	1,693	0	0
	500	0	1,935	0	0	1,500	0	2,075	0	0
	1,000	0	3,947	0	0	3,000	0	4,078	0	1
	50	0	192	0	0					
	100	0	394	0	0					
	200	0	796	0	0					
	300	0	1,202	0	0					
	400	0	1,598	0	0					
	500	0	1,994	0	0					
	1,000	0	4,004	0	0					
	50	0	165	0	0					
	100	0	375	0	0					
	200	0	778	0	0					
	300	0	1,185	0	0					
	400	0	1,576	0	0					
	500	0	1,981	0	0					
	1,000	0	3,981	0	0					
4 clienti	50	0	169	0	0	200	0	397	0	0
locali	100	0	357	0	0	400	0	554	0	0
	200	0	747	0	0	800	0	925	0	0
	300	0	1,156	0	0	1,200	0	1,372	0	0
	400	0	1,553	0	0	1,600	0	1,744	0	0
	500	0	1,910	0	0	2,000	0	2,166	0	1
	1,000	0	3,951	0	0	4,000	0	4,151	1	1
	50	0	198	0	0					
	100	0	394	0	0					
	200	0	803	0	0					
	300	0	1,194	0	0					
	400	0	1,594	0	0					
	500	0	1,969	0	0					
	1,000	0	4,001	0	1					
	50	0	195	0	0					
	100	0	362	0	0					
	200	0	801	0	0					
	300	0	1,197	0	0					
	400	0	1,601	0	0					
	500	0	2,000	0	0					
	1,000	0	4,004	0	0					
	50	0	181	0	0					
	100	0	397	0	0					
	200	0	785	0	0					
	300	0	1,153	0	0					
	400	0	1,559	0	0					
	500	0	1,995	0	1					
	1,000	0	3,947	0	1					

Se observă că:

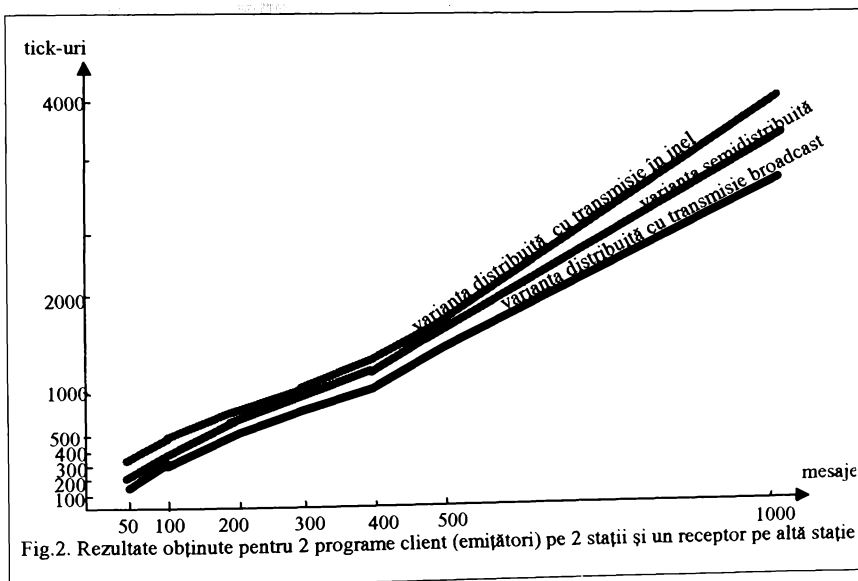
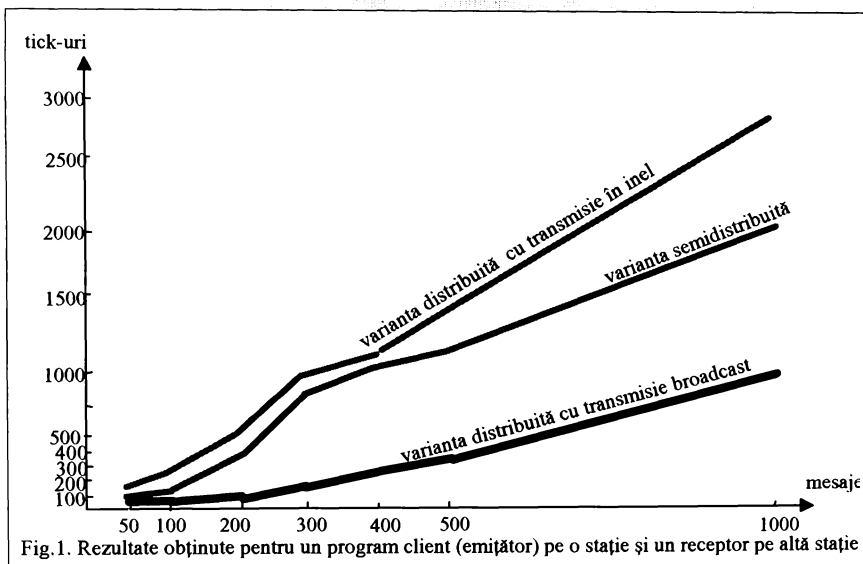
- (1) pentru varianta distribuită broadcast se obțin cei mai buni la transmisia mesajelor, atîta timp cît serverul nu este sufocat și nu se fac retransmisiile (retransmisiile se fac în situația în care în care serverul returnează o eroare de tip EBUSY);
- (2) deoarece s-a ales pentru dimensiunea bufferului de mesaje din token o dimensiune de cel puțin de două ori mai mare decît numărul maxim de stații care se pot conecta la sistem, pentru varianta distribuită cu transmisie în inel nu se poate obține situația în care nu mai este loc pentru o cerere în token.
- (3) timpul obținut pentru varianta în inel este cel mai slab; rezultatul este însă compensat de faptul că în această variantă nu se poate obține situația de "server de comunicație sufocat";
- (4) pentru varianta semidistribuită se obțin timpi intermediari între cele două variante distribuite.

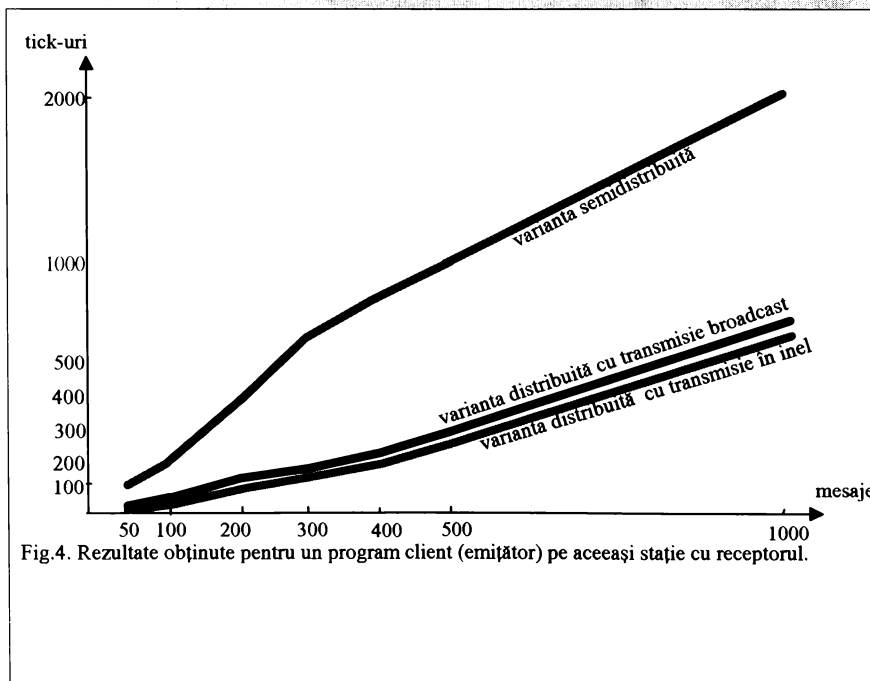
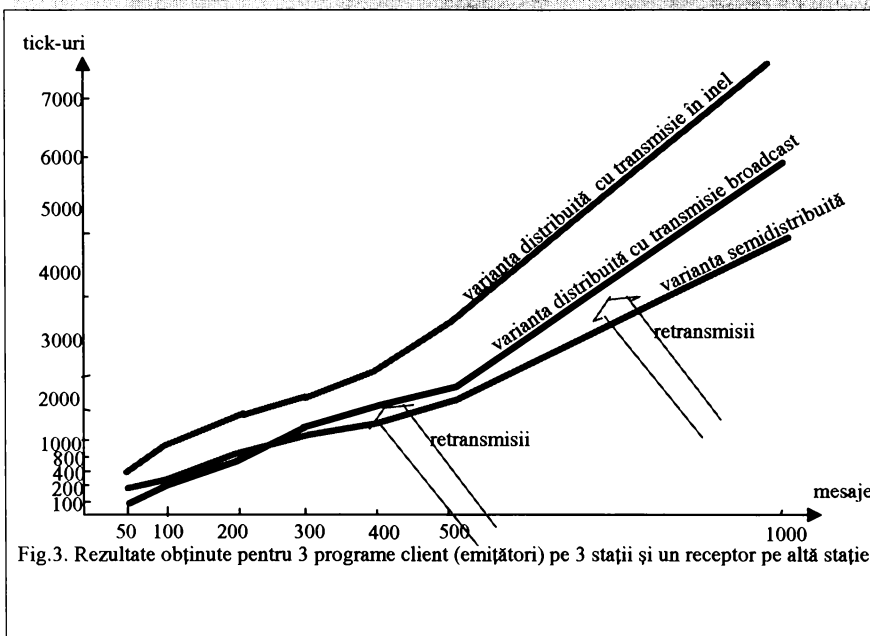
Anexa F: Comparație între performanțele variantelor distribuite și semidistribuite ale sistemului TGC

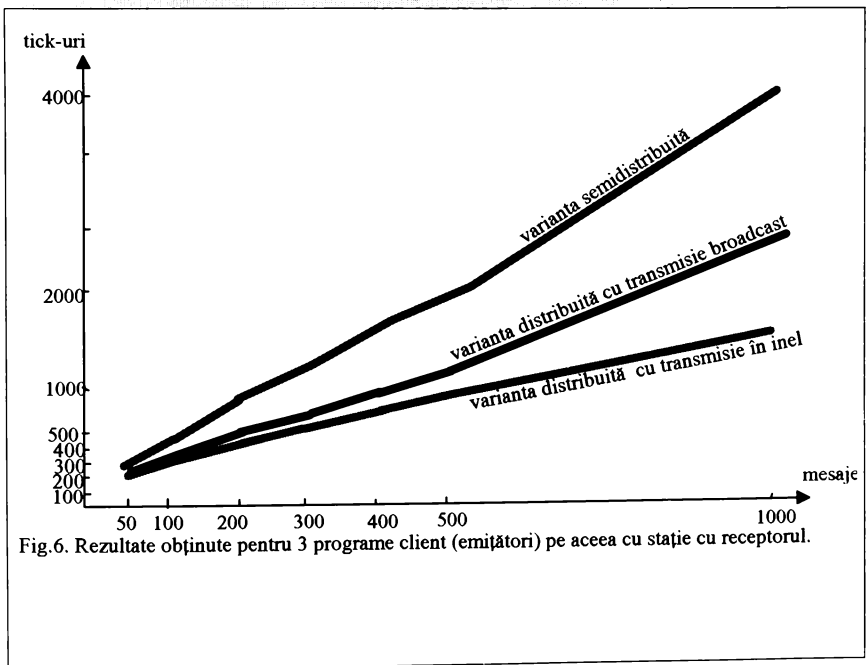
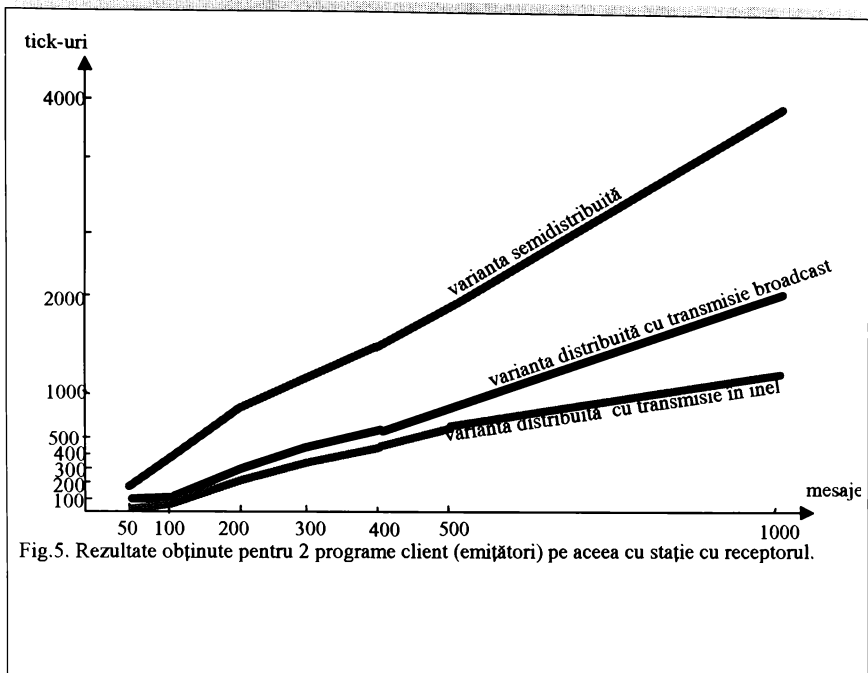
(5) în cazul cînd în sistem există conectată o singură stație, cea mai rapidă variantă este varianta distribuită cu transmisie în inel, deoarece serverul de comunicație nu se mai execută (returnează o confirmare pentru orice cerere de la managerul de grupuri). În acest caz, se observă că cea mai proastă comportare o are varianta semidistribuită, pentru care trebuie să lucreze două servere.

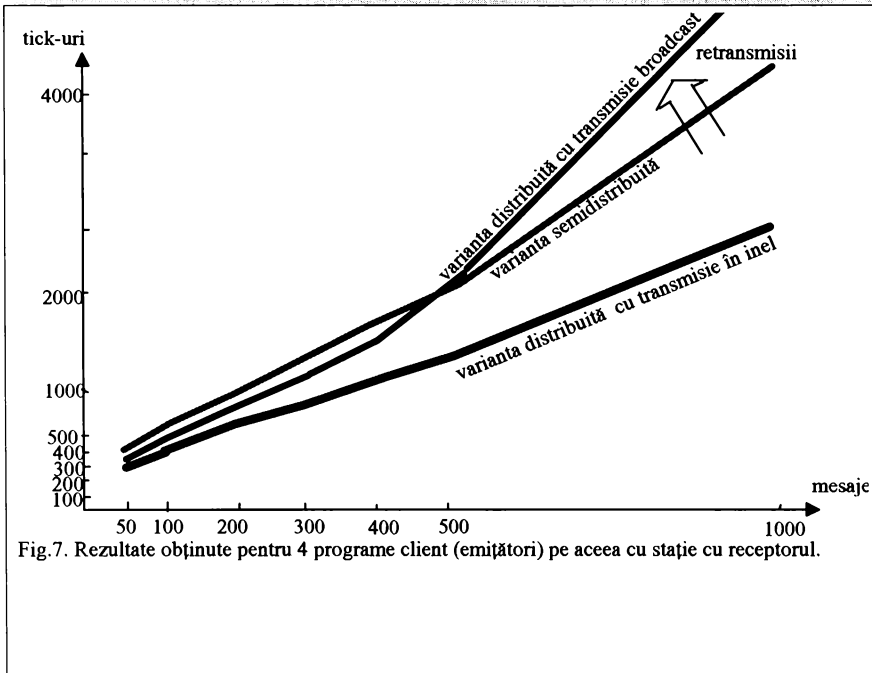
(6) în cazul cînd în sistem există conectată o singură stație, în varianta cu transmisie în inel, începînd de la 4 procese locale au loc retransmisii locale către managerul de grupuri, datorită depășirii capacității bufferului de mesaje.

În figurile următoare se prezintă centralizat aceste concluzii.









Bibliografie.

- [1] [ABB93] Assenmacher H., Breitbach T., Buhler P., Hubsch V., Schwarz R.: "PANDA-Supporting Distributed Programming in C++", in *Object Oriented Programming*, pp361-383, Springer Verlag, 1993
- [2] [ABN85] Almes G.T., Black. A.P., Lazawska E.D. and Noe J.D.: "The Eden System: A Tehcnical Review", *IEEE Transaction on Software Engineering* SE-11 (1), pp. 43-59,1985
- [3] [AC93] Atkins N.S., Coady M.S.: "Adaptable Concurrency Control for Atomic Data Types", *ACM Transaction on Computer Systems*, Vol.10, No.3, pp 190-225, 1993
- [4] [ACF87] Artsy Y. Chang H.Y. and Finkel R.: "Interprocess Communication in Charlotte", *IEEE Software*, pp 22-28, 1987
- [5] [ADB87] Ahamad M., Dasgupta P., Le Blanc R.J., and Wilkis, C.T. : "Fault taulerant computing in object-based distributed systems", *IEEE 6th Simposium of Reliability in Distributed Software and Database Systems*, 1987
- [6] [ALP90] Ancilotti, P., Lazzerini, B., Prete, C., A, and Sacchi, M. : "A Distributed Commit Protocol for a Multicomputer System", *IEEE Transactions on Software Engineering*, vol.16, No.5, pp 718-725, 1990
- [7] [And90] Anderson A.: "Multiple Processing. A System Overview.", McGrawHill 90
- [8] [And91] Andrews G.R.: "Paradigms for Process Interaction in Distributed Programs", *ACM Computing Surveys*, vol 23, 1, 49-90, 1991
- [9] [AS83] Andrews G.R., and Schneider F.B.: " Concepts and Notations for Concurrent Programming", *ACM Computing Surveys*, 15(1), pp 3-49, 1983
- [10] [AW94] Attiya, H., Walch, J.: "Sequential Consistency versus Linearizability", *ACM Transactions on Computer Systems*, Vol.12, No.2, pp 91-122, 1994
- [11] [Bac89] Bach M.: "The Design of the Unix Operating System", Prentice Hall, 1989
- [12] [BBK91] Balter R., Banatre J.B., Krakowiask : "Construction des systemes repartis", Ed. Collection Didactic, Inria 91
- [13] [C87] Bernshad B.N., Ching D.T., Lazawska E. D. , Sanislo E.D. and Schwarz M.; "A Remote Procedure Call Protocol Facility for Interconnecting Heterogeneous Computer Systems", *IEEE Transactions on Software Engineering* SE-13 (8), pp 880-894, 1987
- [14] [BDW89] Barbacci M., Doubleday D., Weinstock Ch.: "The Durra Language and Runtime Environment: Tools for PMS-level Programming", *Working Conference on Decentralized Systems*.Proceedings, pp1-16, Lyon, 1989
- [15] [Bel 89] Belzile, C., Coulas M., MacEwen G.H., Marquis G.: "RNET: A Hard Real Time Distributed Programming System", *Proceeding of the 1989 Real-Time Systems Symposium*, pag 2-13., 1989
- [16] [BHG87] Bernstein Ph., Hadzilacos V., Goodman N.,; "Concurency Control and Recovery in Database Systems", Addison Wesley 87
- [17] [BF80] Bullis K., Franta W.: " Implementation of Evencounts in a Broadcast Network", *Computer Networks*, 4, 57-69, 1980
- [18] [BF85] LeBlanc T., Friedberg S.: "HPC: A Model of Structure and Change in Distributed Systems", *IEEE Transaction on Computers*, vol. C-34, pp 1114-1130, dec. 1985
- [19] [BG81] Bernstein P.A. and Goodman N. : "Concurency control in distributed database systems', *ACM Computing Surveys.*, 13, 2 (Jun), 185-221, 1981
- [20] [BGH93] Bever M., Geihs K., Heuser L., Muhlhauser M., Schill A.: " Distributed Systems, OSF DCE and Beyond", in *Distributed Systems*, pp1-20, Springer Verlag, 93
- [21] [Bir91] Birmann K.: "Fault-Tolerance in Sixth Generation Operation Systems", in *Operating Systems of the 1990s and Beyond*, pp154-161, Springer-Verlag, 1991
- [22] [BJ87] Birman K.P., Joseph T.A.: "Reliable Communications in The Presence of Failures", *ACM Transaction on Computer Systems*, vol.5, Nr.1, 46-47, Feb.87
- [23] [BJ90] Birman, K.P., Joseph, T.A.: "Communication support for reliable distributed computing", *Fault Taulerant Distributed Computing*, Springer-Verlag, 1990
- [24] [BK85] Bancilhon F.K., and Korth H.: "A Model of CAD Transactions", *Proceeding of the 11th International Conference on Very Large Databases*, pp.25-33, 1985
- [25] [BKT85] Brown M.R., Kolling K.N., and Taft, E.A. : "The Alpine file system", *ACMTOCS*, vol.3, no.4, pp.261-193, Nov. 1985

- [26] [BK91] Barghouti N.S., and Kaiser G.E. : "Concurrency Control In Advanced Database Applications", *ACM Computing Surveys*, vol.23, No.4, pp 269-294, 1991
- [27] [BN84] Birell A. D. and Nelson B. J.: "Implementing Remote Procedure Calls", *ACM Transaction on Computers Systems* 2 (1) pp 39-59, Feb 1984
- [28] [BT84] Brancha G., Toneg S.: "A distributed algorithm for generalized deadlock detection", *Proceedings of the Symposium on Principles of Distributed Computing*, ACM, New York, pp 285-301, 1984
- [29] [CC91] Chin R., Chonson S.: "Distributed Object-Based Programming Systems", *ACM Computing Surveys*, Vol.23, No.1, pp 91-124, March 1991
- [30] [CDI90] Ciciani B., Dias D. Iyer, B., Yu, P. : "A Hibrid Distributed Centralized System Structure for Transaction Processing", *IEEE Transactions on Software Engineering*, vol.16, No.8, pp 791-805, 1990
- [31] [CG87] Cellary, Gelenbe.: Concurrency Control in Distributed Database Systems", Addison Wesley 87
- [32] [Cha82] Chandy K.M.: "A Mutual Exclusion Algorithm for Distributed Systems", *Techn. Rep.* University of Texas
- [33] [Che84] Cheriton D.R.: " The V Kernel : A Software Base for Distributed System", *IEEE Software*, pp 19-42, 1984
- [34] [CL85] Chandy K.M., and Lamport, L.: "Distributed snapshots: Determining global states of distributed systems", *ACM Transactions on Programming Language Systems*, 3, 1, 63-75, Feb. 1985
- [35] [CL87] Chin L., Liu M.T.: "An optimistic concurrency control mechanism without freezing for distributed database systems", in *Proceeding Conf. on Data Engineering*, pp.322-329, LA, 1987
- [36] [CO82] Ceri S., Owicki S.: "On the Use of Optimistic Methods for Concurrency Control in Distributed Databases", *Proc. 6th Berkeley WorkShop on Distributed Data Management and Computer Networks*, pp.117-130, Berkeley, 1982
- [37] [CS91] Comer D., Stevens D., "Internetworking with TCP/IP", Prentice Hall, 1991
- [38] [CM82] Chandy K.M., and Misra, J. : "A Distributed algorithm for detecting resource deadlocks in distributed systems", *Proceedings of the Symposium on Principles of Distributed Computing*, ACM, New York, pp 157-164, 1982
- [39] [CM84] Chang J., Maxemohuk N.F.: "Reliable Broadcast Protocols", *ACM Transaction on Computer Systems*, vol.2, Nr.3, 251-273, Aug. 84
- [40] [CM86] Chandy K.M., and Misra J.: "An exemple of stepwise refinement of distributed programs: Quiescence detection", *ACM Transactions on Programming Language Systems*, 8, 3, 326-343, July 1986
- [41] [CMH83] Chandy K.M., Misra J., and Haas L.M., : "Distributed deadlock detection", *ACM Transactions on Computer Systems*, 1, 2, pp 147-156, May 1983
- [42] [CPS93] Cabrera L.F., McPherson J. Schwarz, P.M., Wyllie, J.: "Implementing Atomicity in Two Systems: Techniques, Tradeoffs and Experience", *IEEE Trans. on Software Engineering*, Vol.19, No.10, pp 950-601, Oct 1993
- [43] [CP84] Ceri S., Pellagatti, G.: "Distributed Databases: Principles and Systems", McGraw-Hill, 1984
- [44] [CZ83] Cheriton D.R. and Zwaenepoel W.: "The Distributed V Kernel and its Performance for Diskless Workstations", *Proceedings of the Ninth ACM Symposium on Operating Systems Principles*, Bretton Woods, New Hampshire, pp. 129-140, 1983
- [45] [Cri89] Crichlow A.: " Introduction to Distributed and Parallel Computing", McGraw-Hill. 1989
- [46] [CTM92] Cristea V., Tapus N., Moisa T., Damian V.: " Retele de calculatoare", Ed Teora 1992
- [47] [Cri96] Cristian F.: "Synchronous and Asynchronous Group Communication", *Communication on the ACM*, vol 39, no 4, april 1996
- [48] [Dav81] Davies D. W.: *Distributed Systems-Architectures and Implementation*", Lectures Notes in Computer Science, Springer Verlag, Berlin 1981
- [49] [Dav94] Davis R.: " Windows NT Network Programming", Addison-Wesley 1994
- [50] [Dei90] Deitel H.: "An Introduction to Operating System", Addison Wesley 1990
- [51] [DG89] Delattre E., Geib J.M.: "OMPHALE, an active objects-based distributed operating system", *Working Conference on Decentralized Systems*, Proceedings, Lyon 89, pag 327-340
- [52] [Dil93] Dille J.: "Object-Oriented Distributed Computing With C++ and OSF DCE", in *DCE-The OSF Distributed Computing Environment Client Server Mode and Beyond*, pp256-266, Springer-Verlag, 1993
- [53] [DM96] Dolev D., Lalki D.: "The Transis Approach to High Availability Cluster Communication". *Communication on the ACM*, vol 39, no 4, pp54-63, april 96

- [54] [DVP88] Dumitrescu V., Teodor V., Paiu O., Stefanescu V.: " Inițiere în informatica distribuită", Ed. tehnica, București, 1988
- [55] [ED88] Elmagarmid A.,K., and Datta A.,K.: "Two-Phase Deadlock Detection Algorithm", *IEEE Transactions on Computers*, Vol.37, No.11, pp 1454-1458, 1988
- [56] [EGT76] Eswaran K.P., Gray J.N., Lorie R.A., Traiger, I.L.: "The Notion of Consistency and Predicate locks in a Database System", *Communication on the ACM*, 19(11), 624-633, 1976
- [57] [FHR93] Fazzolare M., Humn B., Ranson D.: "Object Oriented Extendibility in Hermes/ST, a Transactional Distributed Programming Environment", in *Object Based Distributed Programming*, pp240-262, Springer Verlag, 1993
- [58] [For88] Fortier P.: " Design of Distributed Operating Systems", McGraw Hill, 1988
- [59] [FSK89] Finkel R., Scott M.L., Kalsow W.K., Artsy Y. and Chang H.Y.: " Experience with Charlotte : Simplicity and Function in a Distributed Operating System", *IEEE Transactions on Software Engineering*, SE - 15, 6, pp 676-685, 1989
- [60] [GDT92] Grassi V.,Donatiello L.,Tucce S.: "On the Optimal Checkpointing of Critical Tasks and Transactions", *IEEE Transactions on Software Engineering*, vol 18, no 1, 1992
- [61] [Gue94] Guerraio R.: "Atomic Object Composition", in *Object Oriented Programming*, pp118-139, Springer Verlag, 1994
- [62] [GM83] Garcia-Molina H.:"Using semantic Knowledge for transaction processing in a distributed database", *ACM Transactions on Database Systems*, 2 (June), pp.186-213, 1983
- [63] [GMo82] Garcia Molina H.: "Election in a Distributed Computing System", *IEEE Transactions on Computers*, C-31(1), 48-59, 1982
- [64] [GS91] Garcia Malina H., Spanster A.: "Ordered and Reliable Multicast Communication", *ACM Transaction on Computer Systems*, vol. 9 nr.3, 242-271, Aug.1991
- [65] [Gif79] Gifford D.K., "Weighted Voting for Replicated Data", *Proceedings of the Seventh Symposium on Operating Systems Principles*, ACM, 150-162, 1979
- [66] [Gie91] Gien M.: "Next Generation Operating Systems Architecture", in *Operating Systems of the 1990s and Beyond*, pp227-232, Springer-Verlag, 1991
- [67] [GS80] Gligor V., Shattick S.: "On deadlock detection in distributed databases", *IEEE Transactions on Software Engineering*, SE-6, 5 Sept.1980
- [68] [Gos91] Goscinski A.: "Distributed Operating Systems- The logical design", Addison Wesley, 1991
- [69] [Hag87] Hagmann R. : "Reimplementing the Cedar file system using logging and group commit", *Proceedings of the 11th SOSP*, ACM, pp.155-162, Nov. 1987
- [70] [Hal88] Halshall F.: "Data Communications Computer Network and OSI", second edition, Addison Wesley, 1988
- [71] [Hin96] Hinden R.M. : " IP Next Generation Overview", *Communication on the ACM*, vol. 39, No.6, Jul.96
- [72] [HJP87] Helary J., Jard C., Planzean, N., and Rayna, M., 1987 : "Detection of stable properties in distributed application", *Proceedings of the Symposium on Principles of Distributed Computing*, ACM, New York, pp 125-136, 1987
- [73] [HT83] Haerder T., Reuter A.:"Principles of Transactions-Oriented Database Recovery", *ACM Computing Surveys*, vol.15, No.4, pp 287-289, Dec. 1983
- [74] [HSE91] Hoffman F., Schlenk P., Eirich T., "Encapsulation and Interaction in Future Operating Systems", in *Operating Systems of the 1990s and Beyond*, pp241-245, Springer-Verlag, 1991
- [75] [HW88] Herlity M.P., and Wehl W.E.: "Hybrid Concurrency for Abstract Data Types", *Proceedings of the 7th ACM Symposium on Principles of Database Systems*, ACM Press, pp.201-210, 1988
- [76] [JB89] Joseph T.A. and Birman K.P.: "Reliable Broadcast Protocols" in *Distributed Systems*, ed. Mullender S.J., ACM Press, pp293-316, 1989
- [77] [JLH88] Jul E., Levy H., Hutchinson N. and Black A.: "Fine-grained Mobility in the Emerald System", *ACM Transactions on Computer Systems* 6(1), pp 109-133, 1988
- [78] [Jon79] Jones A.: "The Object Model: A Conceptual Tool for Structuring Software", in *Operating Systems-An Advanced Course*, Springer-Verlag,1979
- [79] [JRT85] Jones M.B., Rashid R.F. and Thompson M.R.: "Matchmaker: An Interface Specification Language for Distributed Processing", *Proceedings of the 12th ACM Symposium on Principles on Programming Language*, pp. 225-235, 1985
- [80] [JW93] Johnson B. David, Zwaenepoel W.: "The Peregrine High-Performance RPC Systems", *Software Practice and Experience*, vol. 23(2)[119]
- [81] [Jur84] Jurca I.: "Sisteme de operare", curs, Reprografie, Universitatea Politehnica Timisoara, 1984
- [82] [Jur92] Jurca I.: "Programarea orientata pe obiecte în limbajul C++", ed. Eurobit, Timisoara,1992

- [83] [Kap87] Kapp E. : "Deadlock Detection in Distributed Databases", *ACM Computing Surveys*, Vol.19, Nr.4, Dec.1987
- [84] [Koh81] Kohler W.H.: "A Survey of Techniques for Synchronization and Recovery in Distributed Computer Systems", *ACM Computing Surveys*, 13, 2, 149-183, 1981
- [85] [KR81] Kung H.T., and Robison J.T.: "Optimistic Methods for Concurrency Control", *ACM Transactions on Database Systems*, 6(2), pp.213-226, 1981
- [86] [KS89] Kuman A., Stonebraker M.: "Performance Consideration of an Operating System Transaction Manager", *IEEE Transactions on Software Engineering*, vol.15, No.6, June 1989, pp 705-715
- [87] [Lag79] Lagally, K: "Synchronization in a Layered System", in *Operating Systems-An Advanced Course*, Ed. Flynn, M. J., Gray, J. N., Jones, A. et al. , Springer-Verlag 1979
- [88] [Lam78] Lamport L.: Time, Clocks and the Ordering of Events in a Distributed System", *Communication on the ACM*, 21, 7 July 78
- [89] [Lam81] Lampon E.W. : "Atomic Transactions", *Distributed Systems Architecture and Implementation*, Springer-Verlag, Berlin 246-264, 1981
- [90] [LG85] Lee I., Goldwasser, S.M.: "A Distribute Testbed for Active Sensory Processing", *International Conference on Robotics and Automation*, pp925-930, 1985.
- [91] [LLA81] Lazawska E.D., Levy H.M., Almes G.T., Fisher M.I., Fowler R.I., Vestal S. C.: "The Architecture of the Eden Systems", *Proceedings of the Eight Symposium on Operating Systems Principles*, Pacific Grove, California, pp. 148-159, 1981
- [92] [LRM93] Lashkari Y., Ramachandar V., Malpani S., Mehndiratta S.L.: "Vartalaap: a Distributed Multicast Communication System", *Software Practice and Experience*, vol.23(7), 799-811, Jul. 1993
- [93] [Lan83] LeLann G.: "Synchronization", in *Local Area Networks: An Advanced Course*, Lectures Notes in Computer Science, Springer-Verlag, pp361-395, 1983
- [94] [LCJ87] Liskov, B., Curtis, D., Johnson, P. and Scheifler, R.: "Implementations of Argus", *ACM Proceedings 12th Symposium on Operating System Principles*, pp 111-122, 1987
- [95] [Lis88] Liskov B.: "Distributed Programming in Argus", *Communication on the ACM* 31. 3, 300-312, March1988
- [96] [Lis90] Liskov B.: "Argus", *Fault Tolerant Distributed Computing*, Ed. Simons, B., Spector, A., Springer-Verlag, pp108-115, 1990
- [97] [LS83] Liskov B. and Scheifler R.: "Guardians and Actions: Linguistic Support for Robust. Distributed Programs", *ACM Transactions on Programming Languages and Systems* 5 (3), pp 381-404, July 1983
- [98] [LS90] Levy E., Silberschatz A.: "Distributed File System: Concepts and Examples", *ACM Computing Surveys*, vol 22, No 4, pp 321-372, 1990
- [99] [LM92] Lungu M., Musatescu C.: "Sistem in timp real pentru implementarea unei metode de calcul ierarhizat", *Simpozionul National de Teoria Sistemelor, Robotica, Calculatoare si Informatica de proces (SINTES 6)*, ed. a 6-a Craiova, 28-29 martie 1992, pp82-88
- [100] [Lyn83] Lynch N.A. : "Multilevel atomicity: A new corectness criterion for database concurrency control", *ACM Transactions on Database Systems.*, 8, 4, pp.484-502, 1983
- [101] [Mar87] Martin B. : "Modelling concurrent activities with nested objects.", *Proceedings of the 7th International Conference on Distributed Systems*, IEEE Computer Society Press, pp.432-439, 1987
- [102] [MK83] Magee J.,Kramer, J.: "Dynamic Configuration for Distribute Real-Time Systems", *Proceedings of the 1983 Real-Time Systems Symposium*, pag 277-288, 1983
- [103] [Moc93] Mock M.: "DCE++. Distributed C++ Objects Using OSF DCE", in *Distributed Computing Environment*", pp242-256, Springer Verlag, 1993
- [104] [Mos85] Mos J.E.B. : "Nested Transaction: An approach to reliable distributed computing", The MIT Press, Cambridge, March, 1985
- [105] [ML83] Mohan C., Lindsay B.: "Efficient commit protocols for the tree of Processes model of distributed transactions", *ACM Proceedings of the 2nd SIGALT/ SIGOPS Symposium of Principles of Distributed Computing*, pp 76-88, Aug. 1983
- [106] [MOO87] Maekawa M., Oldehoeft A.E., Oldehoeft R.R.: "Operating System Advanced Concepts", The Benjamin Cummings Publishing Co., 1987
- [107] [MM79] Menasce D., and Muntz R.: "Locking and Deadlock Detection in Distributed Databases", *IEEE Transactions on Software Engineering* SE- 5, 6, 3 May 1979
- [108] [MPP87] Miller B.P., Presotto D.L., Powel M.L., DEMOS/MP: "The Development of a Distributed Operating System", *Software-Practice and Experience*, 17 (4), pp. 277-290. 1987

- [109] [MMA96] Moser L.E., Melliar P.M., Agarwal A.D., Budhia R.K., Lingley-Papadopoulos C.A.: "Totem: A Fault-Tolerant Multicast Group Communication System", *Communication on the ACM*, vol 39, no 4, pp54-63, april 96
- [110] [Mus92] **Musatescu C.**, Marian Gh., Padeanu L.: "Sistemul de operare SINIX-Modalitate de implementare a sistemului de operare UNIX.", *Simpozionul National de Teoria Sistemelor, Robotica, Calculatoare si Informatica de proces (SINTES 6)*, ed. a 6-a Craiova, 28-29 martie 1992, pp178-184
- [111] [Mus93a] **Musatescu C.**, Mocanu M., Badica C.: "Symmetric and Unsymmetric Cryptosystems in Distributed UNIX-like Systems", *The International Conference on Applied and Thoretical Electrotechnics*, (ICATE 93)18-20 nov. 93, Craiova, pp 315-320
- [112] [Mus93b] **Musatescu C.**, Mocanu M., Badica C.: "Authentication and Encryption in UNIX-like Systems", *Romanian Open Systems Event ROSE-93*, Cluj 29-30 sept,1993
- [113] [Mus93d] **Musatescu C.**: "Stadiul actual al dezvoltarii sistemelor distribuite", referat doctorat. Timisoara 1993
- [114] [Mus94b] **Musatescu C.**: "Concepte si metode folosite in proiectarea aplicatiilor distribuite", *Simpozionul National de Teoria Sistemelor, Robotica, Calculatoare si Informatica de proces (SINTES 7)*, ed. a 7-a Craiova, 20-23 mai 1994, pp 30-36
- [115] [Mus94c] **Musatescu C.**: "Controlul concurentei in aplicatii complexe", *Simpozionul National de Teoria Sistemelor, Robotica, Calculatoare si Informatica de proces (SINTES 7)*, ed. a 7-a Craiova, 20-23 mai 1994, pp 36-44
- [116] [Mus94d] **Musatescu C.**: "Aplicatii in retele locale de calculatoare", Referat doctorat, Timisoara, 1994
- [117] [MMG94] Marian Gh., **Musatescu C.**, Grosu M., Lungu M.: "Sistem distribuit pentru dirijarea optima a circulatiei in intersectii, sincronizarea semafoarelor de pe fluxurile de circulatie si supravegherea functionarii acestora", grant no 3006/14c/B9/1994, beneficiar Ministerul Invatamintului, 1994
- [118] [Mus95a] **Musatescu C.**: "Support for Distributed Transaction in Distrib", *10-th International Conference on Control Systems and Computer Science*, Bucharest, 24-26 May 1995, pp287-291
- [119] [Mus95b] **Musatescu C.**: "Distrib-A Distributed Object-Oriented System Based on Transaction Processing", *5-th Symposium Automatic Control and Computer Science*, Iasi 26-27 oct 1995, pp157-163
- [120] [Mus95c] **Musatescu C.**: "Sistemul de operare Novell NetWare 3.11. Apeluri sistem", Reprografia Universitatii din Craiova, 1995
- [121] [Mus95d] **Musatescu C.**: "Gestionarea fisierelor in retele de calculatoare", *Else*, no 5, pp2-19, dec 1994
- [122] [Mus96a] **Musatescu C.**, Burdescu D.: "DOSTP-A Distributed Object Oriented System", The International Conference of tehcnical Informatics, Timisoara Conti-94, 14-16 nov 1996, pp 25-33
- [123] [Mus96b] **Musatescu C.**, Burdescu D.: "The Client Server Model in a Distributed Object Oriented System Based on Transaction Processing", *International Symposium On Systems Theory*, SINTES 8, Craiova 6-7 iunie 1996, pp 21-28
- [124] [MMB96] Marian, **Musatescu C.**, Budica I.: "A Strategy for Traffic Control System", *International Symposium On Systems Theory*, SINTES 8, Craiova 6-7 iunie 1996, pp143-152
- [125] [Mus96d] **Musatescu C.**: "A Novell Implementation of a Distributed Object Oriented System Based on Transaction Processing", *International Symposium On Systems Theory*, SINTES 8, Craiova 6-7 iunie 1996, pp 180-188
- [126] [Mus96e] **Musatescu C.**, Mocanu M., Olteanu C.: "TDL: A Run-Time Environment for Connecting Tasks in a Unix Local NetWork", *International Symposium On Systems Theory*, SINTES 8, Craiova 6-7 iunie 1996, pp309-316
- [127] [Mus96f] **Musatescu C.**: "Aspecte ale proiectarii sistemelor de operare", *Tech. Rep. Tempus*, Reprografia Univ. Craiova, 1996
- [128] [Mus96g] **Musatescu C.**, Burdescu D.: " The Design and the Implementation of a DOSTP- A Distributed Object Oriented System Based On Transactions Processing", *Facta Universitatis, series: Electronics and Energetics*, vol 9, no 2, 1996, pp 149-172, Nis, Yugoslavia
- [129] [Mus96i] **Musatescu C.**: "Designing a Distributed Object Oriented System", *Research Rep.*, Technological Educational Institute of Piraeus, Faculty of engineering, Department of Computer Systems, Athens, Greece, may 1996
- [130] [Mus96j] **Musatescu C.**: "Aspecte ale modelului client server in programarea distribuita", *Else*, no 9, pp33-52, 1996
- [131] [Mus97] **Musatescu C.**, Burdescu D.,Hamburg I.: "Object Oriented System for Distributed Transactions during the Development of a Product within a group of Companies". *1-st bibliografie-5*

- International Conference on Engineering Design and Automation*, Bangkok, Thailand, March 18-21, 1997
- [132] [Mull89] Mullender S.J., *Distributed Systems*, ACM Press, 1989
- [133] [Mull89b] Mullender S.J., "Interprocess Communications", in *Distributed Systems*, ed. Mullender S.J., pp37-63, ACM Press, 1989
- [134] [Mul90] Mullender, S. J.: "Amoeba-A Distributed Operating System for the 1990s", *Computer Magazine*, May 1990
- [135] [Neh91] Nehmer J.: "The Immortality of Operating Systems or Is Research in Operating Systems Still Justified?", in *Operating Systems of the 1990s and Beyond*, pp77-82, Springer-Verlag, 1991
- [136] [Nut92] Nutt G.: "Centralized and Distributed Operating Systems", Prentice Hall, 1992
- [137] [OCD88] Ousterhout J.K., Cherenen A.R., Douglish F., Nelson M.N. and Welch B.B.: "The Sprite Network Operating system", *Computer*, 21(2), pp 23-26,1988
- [138] [Pow83] Powell M.L. and Miller B.P.: "Process Migration in DEMOS/MP", in *Proceedings of the 9th ACM Symposium on Operating System Principles*, Bretton Woods, New Hampshire, pp 110-119,1983
- [139] [PSU86] Pradel U., Schlageter G., and Unland R. : "Redesign of optimistic methods: Improving performance and availability", *Proceedings of the 2nd International Conference on Data Engineering IEEE Computer Society Press*, pp.466-473, 1986
- [140] [PG93] Paunescu F., Golesteanu D. P., "Sisteme cu prelucrare distribuita si aplicatiile lor", Ed tehnica 93
- [141] [RA81] Ricard, G., Agrawala R. K.: "An optimal Algorithm for Mutual Exclusion in Computer Network", *Communication on the ACM*,22(2), 24(1), 9-17, 1981
- [142] [Ras81] Rashid R.F.: "An Inter-Process Communication Facility for UNIX", *Local networks for Computer Communications*, Springer-Verlag, pp 319 -354, 1981
- [143] [Ras88] Rashid R.F.: "Machine - Independent Virtual Memory Management for Paged Uniprocessor and Multiprocessor Architectures", *IEEE Transactions on Computers*, C - 37(8), pp 896 - 908, 1988
- [144] [Ray87] Raymond K.: "Multiple Entries with Ricard and Agrawala Distributed Mutual Exclusion Algorithm", Techn. Report, Department of Computer Science, The University of Queensland, 1987
- [145] [RB88] Ramamritham K., Badrinat B.R. : "Synchronizing Transactions on Objects:, *IEEE Transactions on Computer*, Vol.37, No.5, 541-549, May 1988
- [146] [RBM96] Robert van Renesse, Birman K.P., Maffei S.: "Horus: A Flexible Group Communication System", *Communication on the ACM*, vol 39, no 4, pp54-63, April 96
- [147] [RK79] Reed D.P., Kanodia R. K.: "Synchronization with Eventcounts and Sequencers", *Communication on the ACM*, 22(2), 115-123, 1979
- [148] [RSB90] Ramanathan P., Shin K.G. and Butler R.W.: "Fault Taulerant Clock Synchronization in Distributed Systems", *IEEE Computer*, vol 23, pp 33-42, oct 1990
- [149] [SBD86] Spector A.Z., Bloch, J.J., Daniels, D., Dravers D.S., Duchamp, D., Eppinger J.L., Menes S.G., and Thomson D.S. : "The Camelot Project", *Tech. Rep. CMU-CS-86*, 166, Departament of Computer Science, Carnegie Mellon University, Pittsburgh, Penn., 1986
- [150] [SC91] Shrivastava S., McCue D.: "Operating System Support for Object Oriented Distributed Systems", in *Operating Systems of the 1990s and Beyond*, pp256-257, Springer-Verlag, 1991
- [151] [SDE85] Spector A.Z., Daniels D.S, Duchamp D., Eppinger J.L., and Pausch R.: "Distributed Transactions for Reliable Systems", *Tech. Rep. CMU-CS-85*, 117, Departament of Computer Science, Carnegie Mellon University, Pittsburgh, Penn., 1985
- [152] [SGA87] Salem K., Garcia-Molina H., and Alonso R.: "Altruistic locking: A strategy for coping with long-lived transactions", *Proceedings of the 2nd International WorkShop on High Performance Transaction Systems*, pp.19.1-19.24, 1987
- [153] [SH82] Shoch J.F. and Hupp J.A. : "The Worm' Programs - Early Experience With a Distributed Computation", *Communications of the ACM*, 25,3, pp 173-80, 1982
- [154] [SK89] Stonebraker M., Kumar A.: "Performance Consideration for an Operating System Transaction Manager", *IEEE Transactions on Software Engineering*, vol 15. no. 6, pp705-713, jun 1989
- [155] [SR95] Shirley J., Rosenberry W.: "Microsoft RPC Programming Guide", O'Reilly&Associates.Inc. 1995
- [156] [SP88] Silberschatz A., Peterson J., "Operating Systems Concepts", Addison Wesley. 1988
- [157] [Sch91] Schmutz H.: "Autonomous Heterogeneous Computing. Some Open Problems". in *Operating Systems for 90s and Beyond*", pp63-67, Springer Verlag, 1991

- [158] [Spe89a] Spector A. Z." Achieving Applications Requirements", in *Distributed Systems*, ed. Mullender S.J., ACM Press, pp19-35, 1989
- [159] [Spe89b] Spector A. Z."Distributed Transaction Processing Facilities", in *Distributed Systems*, ed. Mullender S.J., ACM Press, pp191-215, 1989
- [160] [Spe90] Spector A.Z.: "TABS", *Fault Tolerant Distributed Computing*, Ed. Simons, B., Spector, A., Springer-Verlag, pp115-124, 1990
- [161] [Spe91] Spector A.Z.: "Some Thoughts on Systems Challenges for the 1990", in *Operating Systems of the 1990s and Beyond*. pp223-226. Springer-Verlag, 1990
- [162] [SP84] Schwarz P.M., Spector A.Z.: "Synchronizing shared abstract abstract types", *ACM Transactions on Computer Systems*. vol 3, pp 223-250, 1984
- [163] [SR96] Schiper A., Raynal M.: " From Group Communicatin to Trabsactions in Distributed Systems", *ACM Transactions on Computer Systems*, 2 (4), pp. 277-88, 1984
- [164] [SRC84] Saltz I. H. , Reed D. P. , Clark D. D.: End-To-End Argumenys in System Design , *ACM Transactions on Computer Systems*. 2 (4), pp. 277-88, 1984
- [165] [SS85] Spector A., Schwarz P.: "Transactions: A construct for Reliable Distributed Computing", Tech. Rep. CMU-CS-85-143, 1985
- [166] [Ste90] Stevens W.R.: "Unix Networking programming", Prentice Hall, 1990
- [167] [Sun85] Sun Microsystems: "Remote Procedure Call Protocol Specification", Sun Microsystem Inc; 1985
- [168] [SK85] Suzuki I., Kasan T.: " A Distributed Mutual Exclusion Algorithm", *ACM Transactions on Computer Systems*, 3(4), 344-349, 1985
- [169] [Sta89] Stalligs W.: " Handbook of computer communication standards", Howard W. Sams & Company, 1989
- [170] [Svo84] Svobodova L.: "File Servers for Network - Based Distributed Systems", *Computing Surverys*,16(4), pp 353-98, 1984
- [171] [SZ89] Skarra A.H., and Zdenik S.B.: "Concurrency control and object-oriented databases", *Object-Oriented Concepts, Databases and Applications*, ACM Press, pp.395-421, NY, 1989
- [172] [THG89] Tag M.,Hebrad P., Geib J.M.: "Experience with CSA-An Object-Oriented Distributed System", *Working Conference on Decentralized Systems*, Proceedings, pp309-326, Lyon,89
- [173] [TCM94] Taylor P., Cabrill V., Mock M.: "Combining Object-oriented systems and open transactions processing", *The Computer Journal*, Vol.37, No.6, pp 488-498, 1994
- [174] [Tan88] Tanenbaum R.S.: "Computer Networks", Second Edition,1988
- [175] [Tan93] Tanenbaum R.S.: "Modern Operating Systems", Prentice Hall 1993
- [176] [Tan96] Tanenbaum R.S.: "Distributed Operating Systems", Prentice Hall 1996
- [177] [Tok91] Tokoro M.: "Toward Computing Systems fot the 2000s", in *Operating Systems of the 1990s and Beyond*, pp233-239, Springer-Verlag, 1991
- [178] [TR85] Tanenbaum R.S. and van Renesse R.: "Distributed Operating Systems", *ACM Computing Surveys*, 17 (4), pp. 419-470, 1985
- [179] [UB87] Uyless B.: " Computer Networks: Protocol, standards and interfaces", Prentice Hall, 87
- [180] [WPE83] Walker B., Popek G., English R., Kline C. and Thiel G., " The Locus Distributed Operating System", *Proceedingd of thr 9th ACM Symposium on Operating Systems Principles*, Bretton Woods, N.N., pp 49 -70, 1983
- [181] [Wei89a] Weihl, W.E." Remote Procedure Call", in *Distributed Systems*, ed. Mullender S.J., ACM Press, pp65-87, 1989
- [182] [Wei89b] Weihl, W.E."Using Transactions in Distributed Systems", in *Distributed Systems*, ed. Mullender S.J.,ACM Press, pp218-237, 1989
- [183] [Wei89c] Weihl, W.E."Theory of Nested Transactions", in *Distributed Systems*, ed. Mullender S.J., pp237-263, ACM Press, 1989
- [184] [Wet91] Wettstein H.: "The Explainable Operating System", in *Operating Systems of the 1990s and Beyond*, Springer-Verlag, pp247-248, 1991
- [185] [Wit80] Wittie L.D. and von Tilborg A.M.: "MICROS.A Distributed Object Operating System for MICRONET. A Reconfigurable Network Computer", *IEEE Transactions on Computer*, 29(12), pp 1133-44, 1980
- [186] [YEK87] Yeh S., Elles, C., Ege, A. and Korth, H.: "Performance analysis of two concurrency control schemes for design environements", *Tech. Rep. STP-036-87*, MCC, Austin, Texas, 1987
- [187] [Zah91] Zahorjan J., " *Operating Systems on the 1990s and Beyond*, in *Operating Systems of the 1990s and Beyond*, pp53-55, Springer-Verlag, 1991
- [188] [Zay87] Zayas E.R.: "Attacking the Process Migration Bottleneck", in *Proceedings of the 11 th ACM Symposium on Operating System Principles*, Austin,Texas, pp 13-24, 1987