

UNIVERSITATEA "POLITEHNICA"
TIMIȘOARA

BIBLIOTECA CENTRALĂ

Nr. inv. 555.04x

Dulap 369 Lit. C

UNIVERSITATEA „POLITEHNICA” TIMIȘOARA
DE AUTOMATICĂ ȘI CALCULATOARE

Viorel COIFAN

CONTRIBUȚII LA OPTIMIZAREA PERFORMANȚEI RAPORTATĂ LA DEPENDABILITATE A SISTEMELOR MULTIPROCESOR CU MEMORIE PARTAJATĂ

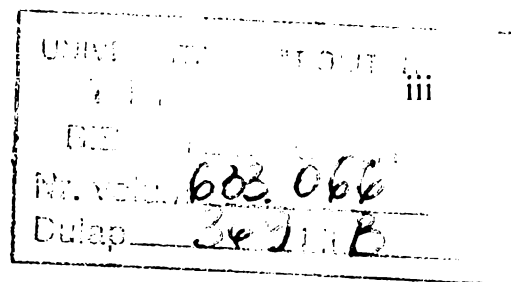
Dizertație pentru obținerea titlului de doctor inginer

BIBLIOTECA CENTRALĂ
UNIVERSITATEA "POLITEHNICA"
TIMIȘOARA

Coordonator:
Prof. dr. ing. Crișan STRUGARU

Timișoara, 2001

PREFAȚĂ



În această dizertație sunt tratate subiecte referitoare la ameliorarea *performanței* sistemelor de calcul multiprocesor în funcție de *dependabilitate*.

Dependabilitatea este un termen introdus de Laprie în 1985 [Lapr85] ce înglobează conceptele de *fiabilitate*, *disponibilitate*, *siguranță în funcționare*, *mentenabilitate*, *performabilitate* și *testabilitate*. Aceste concepte sunt exemple de măsuri utilizate pentru cuantificarea dependabilității unui sistem de calcul.

Capitolul 1, intitulat Introducere, este rezervat prezentării unor probleme fundamentale în abordarea domeniului proiectării sistemelor de calcul multiprocesor. Sunt trecute în revistă subiecte legate de tendințele arhitecturale în dezvoltarea acestor sisteme, reliefându-se strânsa legătură dintre dezvoltarea tehnologică și cea arhitecturală. *Convergența arhitecturilor paralele* exprimă o realitate ce se referă la structura de comunicație, spațiul de adrese partajate, transferul de mesaj, etc. Vechile taxonomii, ca de exemplu cea a lui Flynn, încep să aibă o arie restrânsă de aplicabilitate în comparație cu noile abordări din domeniu.

Coerența și consistența memoriilor cache reprezintă probleme fundamentale în definirea și proiectarea sistemelor multiprocesor și, prin urmare, tratarea lor reprezintă subiecte privilegiate.

Rețelele de interconectare asigură într-un grad ridicat nivelul de performanță al acestor sisteme. *Reducerea latenței și mărirea benzii de trecere* reprezintă obiective ce sunt permanent în atenția proiectanților acestor rețele.

Toleranța la defect este prezentată atât ca și problematică majoră în proiectarea sistemelor multiprocesor, cât și ca subiect privilegiat în analiza de performanță.

În Capitolul 2, intitulat Analiza arhitecturilor sistemelor de calcul MIMD, sunt dezvoltate subiecte care contribuie la o corectă analiză a arhitecturilor sistemelor multiprocesor din clasa MIMD. Ele se referă la granulația de paralelism, latența memoriei, rețele din interconectare, mecanisme de sincronizare a proceselor.

Capitolul 3, intitulat Analiza de performanță a sistemelor de calcul din clasa MIMD, este dedicat analizei de performanță a sistemelor de calcul din clasa MIMD.

După prezentarea tehnologiei de bază pentru evaluarea performanțelor rețelelor de interconectare (RIN) se analizează performanțele principalelor tipuri de rețele: cu *magistrală unică*, cu *magistrală multiplă*, *crossbar* și *multinivel*. Sunt prezentate expresiile analitice ale *lățimii de bandă*, *capacității de trecere*, *probabilității de acceptanță*. În continuare se face o analiză comparativă a facilităților hardware oferite de aceste tipuri de rețele.

Analiza de performanță a arhitecturilor sistemelor multiprocesor (SMP) cu *memorie partajată* și a celor cu *memorie partajată distribuită* urmează în mod firesc. Se prezintă *in extenso* mecanisme și protocoale de menținere a coerenței memoriilor cache și modele de consistență a memoriilor.

Concluziile acestui capitol se referă în special la relația dintre caracteristicile arhitecturale și mecanismele de menținere a coerenței.

În Capitolul 4 care este intitulat Analiza arhitecturilor sistemelor de calcul din clasa MIMD tolerante la defect se prezintă o analiză a arhitecturilor sistemelor de calcul *MIMD tolerante la defect*. Terminologia prezentată ca și enumerarea obiectivelor toleranței la defecte au drept scop reliefarea conceptului de dependabilitate și a tehnicilor de evaluare a ei.

Detectarea defectelor prin duplicare și comparare, prin diagnoză cât și controlul erorii în memoriile de mare viteză sunt enumerate în continuare.

Strategiile de recuperare din defect, crearea și exploatarea punctelor de verificare, recuperarea prin *rulare înapoi* și recuperarea prin *avansare* sunt analizate distinct, atât pentru sistemele multiprocesor cu memorie partajată cât și pentru cele cu transfer de mesaje.

Reconfigurarea sistemelor multiprocesor constituie un paragraf distinct. Sunt prezentate și analizate tehnici de reconfigurare pentru diverse tipuri de RIN, insistându-se în mod particular pe reconfigurarea RIN *multinivel*. În continuare sunt concluzionate câteva aspecte referitoare la proiectarea RIN-MN tolerante la defect.

Analiza și modelarea dependabilității pentru sistemele multiprocesor cu memorie partajată sunt prezentate în Capitolul 5.

Sunt descrise modele de fiabilitate și tehnici de evaluare a fiabilității pentru diferite sisteme multiprocesor bazate pe *magistrală*, *crossbar*, RIN-MN, *hipercub*, etc. Sunt deduse expresiile ce definesc fiabilitatea de terminal pentru diferite arhitecturi.

Modelele de disponibilitate sunt prezentate, utilizând lanțuri Markov, evidențiindu-se contribuția ratelor de defectare respectiv a ratelor de reparare ale diferitelor componente.

În continuare sunt discutate câteva modele de dependabilitate raportate la performanță, introducându-se parametri ca și *disponibilitatea lățimii de bandă*, *performanța acumulată la un moment dat* și *performabilitatea*.

Modelarea disponibilității pentru un multiprocesor cu memorie partajată orientat în jurul unei RIN-MN reprezintă un subiect care a fost tratat cu atenție. Sunt prezentate și analizate subsistemele procesor-memorie, RIN-MN în diferite configurații cu 16, 64 de noduri și în final un model generalizat.

Capitolul 6, intitulat Modelul experimental, este destinat prezentării unui model experimental constituit dintr-un sistem multiprocesor cu memorie partajată distribuită. El este bazat pe o rețea de interconectare, de tip magistrală multinivel. Se prezintă structura sistemului și se scot în evidență avantajele utilizării acestui tip de rețea în comparație cu o rețea de interconectare clasică. Se dezvoltă modelul analitic al acestui sistem și apoi se calculează parametrii de performanță ca și: *factorul de utilizare a sistemului*, *timpul de răspuns*, *fiabilitatea*, etc. Simularea funcționării a fost efectuată sub PROTEUS, un simulator dezvoltat la MIT.

Capitolul 7, intitulat Concluzii generale, prezintă concluziile lucrării de dizertație, contribuțiile autorului și tendințe viitoare de cercetare.

Prezenta lucrare de dizertație a fost elaborată sub îndrumarea competentă și exigentă a conducătorului științific, dl. prof. dr. ing. Crișan STRUGARU. Autorul îi mulțumește în mod deosebit pentru sprijinul permanent acordat, pentru generozitatea și răbdarea pe care a manifestat-o în toată această perioadă.

Conceptul de dependabilitate, ca și un criteriu de analiză a performanței sistemelor multiprocesor, i-a fost sugerat autorului de prof. dr. ing. Mircea VLĂDUȚIU cu care acesta a purtat numeroase discuții interesante referitoare la proiectarea și analiza sistemelor tolerante la defect. Colaborarea cu domnia sa a constituit un privilegiu pentru care autorul îi mulțumește și pe această cale.

Lucrarea elaborată beneficiază de referințe bibliografice la zi. Selectarea bibliografiei, obținerea în timp util a informației necesare, contactele cu colegii din universități de prestigiu ca și: Texas A&M University, University of Minnesota, Princeton University și Stanford University ar fi fost imposibile fără sprijinul colegial și cordial acordat de prof. dr. ing. Petru ELES de la Universitatea din Linköpings (Suedia), fost student și coleg al autorului. Îi mulțumesc încă odată pentru generozitatea și prietenia sa care mă onorează.

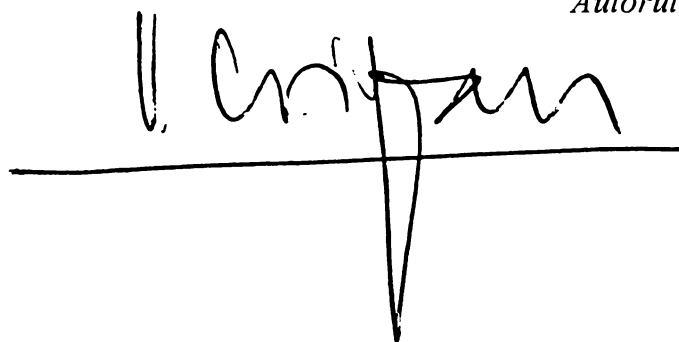
Dactilografierea, prezentarea grafică și în general toate facilitățile oferite de o tehnoredactare modernă și profesionistă au fost asigurate de dl. ing. Norbert HANIGOVSZKI căruia autorul îi mulțumește pentru răbdarea dar și pentru ingeniozitatea dovedită în mod constant.

Pe tot parcursul elaborării acestei lucrări, autorul a beneficiat de prietenia și de sprijinul moral constant al colegilor prof. dr. ing. Octavian PROȘTEAN și prof. dr. ing. Ștefan HOLBAN de la Facultatea de Automatică și Calculatoare.

În fine, dar nu în ultimul rând, autorul își manifestă întreaga sa grațitudine față de familia sa – mamă, soție și copii – care l-a înțeles și l-a sprijinit în toată această perioadă, renunțând de multe ori la micile bucurii cotidiene pe care le-ar fi binemeritat din plin.

Mulțumesc tuturor, care prin cuvinte, fapte și atitudini au contribuit la finalizarea acestei dizertații.

Autorul



Familiei mele

CUPRINS

	Pag.
1. INTRODUCERE	1
1.1 Tendințe arhitecturale	1
1.2 Convergența arhitecturilor paralele	1
1.2.1 Arhitectura de comunicație	3
1.2.2 Spațiul de adresă partajat	3
1.2.3 Transferul de mesaj	9
1.2.4 Convergența	11
1.3 Coerența și consistența memoriilor cache	12
1.3.1 Modele de consistență	14
1.3.2 Coerența memoriilor cache.....	15
1.3.3 Relația între modelele de consistență și coerența memoriilor cache	17
1.3.4 Scheme de asigurare a coerenței memoriilor cache	17
1.3.5 Factori ce afectează mecanismele de menținere a coerenței memoriilor cache	20
1.4 Rețele de interconectare	20
1.4.1 Caracteristici	20
1.4.2 Performanță de comunicare de bază în RJN	22
1.5 Toleranța la latență	27
1.5.1 Cerințe fundamentale, beneficii și limitări ale toleranței la latență	29
1.6 Toleranța la defect	32
1.6.1 Implicații ale toleranței la defect în arhitecturile MIMD ...	34
2. ANALIZA ARHITECTURILOR SISTEMELOR DE CALCUL DIN CLASA MIMD	35
2.1 Aspecte generale	35
2.2 Granulația de paralelism	40
2.3 Latența memoriei și rezolvarea ei	42

		x
2.3.1	Comutarea de proces	43
2.3.2	Utilizarea arhitecturilor LOAD/STORE	43
2.3.3	Coerența memoriilor cache	44
2.4	Rețele de interconectare	58
2.4.1	Terminologie	58
2.4.2	Sistemul multiprocesor cu memorie partajată orientate pe magistrală	64
2.4.3	Rețele de interconectare <i>statice</i>	71
2.4.4	Rețele de interconectare <i>dinamice</i>	75
2.5	Mecanisme de sincronizare a proceselor la arhitecturi din clasa MIMD	94
2.5.1	Moduri de sincronizare	94
2.5.2	Sincronizarea bazată pe variabile partajate	96
2.5.3	Sincronizarea bazată pe transfer de mesaje	102
2.6	Concluzii	106
3.	ANALIZA DE PERFORMANȚĂ A SISTEMELOR DE CALCUL DIN CLASA MIMD	109
3.1	Analiza de performanță la RIN	109
3.1.1	Terminologie de bază utilizată în evaluarea performanțelor RIN	109
3.1.2	Analiza performanțelor magistralei unice	111
3.1.3	Analiza performanțelor magistralei multiple	115
3.1.4	Analiza performanțelor RIN <i>crossbar</i>	125
3.1.5	Analiza performanțelor RIN <i>multinivel</i>	128
3.1.6	Analiza comparativă a performanțelor RIN de tip <i>crossbar, magistrală unică și multinivel</i>	130
3.2	Analiza de performanță a sistemelor de memorie în arhitecturile SMP din clasa MIMD	135
3.2.1	Arhitecturi cu memorie partajată centralizată	138
3.2.2	Arhitecturi cu memorie partajată distribuită	152
3.2.3	Modele de consistență a memoriei	164
3.3	Concluzii	173
4.	ANALIZA ARHITECTURILOR SISTEMELOR DE CALCUL DIN CLASA MIMD TOLERANTE LA DEFECT	177
4.1	Terminologie	177

4.1.1	Obiectivele toleranței la defect la sistemele multiprocesor	180
4.1.2	Parametri de evaluare a dependabilității	181
4.2	Toleranța la defect prin utilizarea redundanței	185
4.2.1	Toleranța la defect prin utilizarea <i>redundanței statice</i>	185
4.2.2	Toleranța la defect prin utilizarea <i>redundanței dinamice</i> ..	192
4.3	Detectarea defectelor în sistemele multiprocesor	193
4.3.1	Detectarea defectelor prin <i>duplicare și comparare</i>	193
4.3.2	Detectarea defectului utilizând <i>diagnosticarea și tehnici de codificare</i>	195
4.3.3	Controlul erorii în memoriile de mare viteză	196
4.4	Strategii de recuperare pentru sistemele multiprocesor	199
4.4.1	Recuperarea din defect prin <i>rulare înapoi</i>	208
4.4.2	Recuperarea din defect prin <i>rulare înainte</i>	216
4.5	Reconfigurarea în sistemele multiprocesor	226
4.5.1	Reconfigurarea în SMP-MP bazate pe <i>magistrală</i>	227
4.5.2	Reconfigurarea rețelelor de tip <i>crossbar</i>	228
4.5.3	Reconfigurarea rețelelor <i>hipercub</i>	229
4.5.4	Reconfigurarea în rețelele <i>plasă</i>	236
4.5.5	Reconfigurarea în rețelele <i>arbore</i>	238
4.5.6	Reconfigurarea în rețelele <i>multinivel</i>	241
4.6	Concluzii	260
5.	MODELAREA DEPENDABILITĂȚII PENTRU SISTEMELE MULTIPROCESOR CU MEMORIE PARTAJATĂ	263
5.1	Etapele preliminare în descrierea modelului	263
5.2	Modele de fiabilitate.....	265
5.2.1	Fiabilitatea de tip <i>terminal</i>	265
5.2.2	Fiabilitatea de tip <i>multiterminal</i>	267
5.2.3	Fiabilitatea <i>bazată pe sarcină</i>	268
5.2.4	Fiabilitatea de tip <i>rețea</i>	272
5.3	Modele de disponibilitate	272
5.3.1	Sisteme orientate pe <i>crossbar</i>	273
5.3.2	Sisteme orientate pe <i>magistrală</i>	273
5.3.3	Sisteme orientate pe <i>RIN-MN</i>	273

5.4	Modele de dependabilitate raportate la performanță.....	275
5.5	Modelarea disponibilității utilizând descompunerea de sistem pentru SMP-MP, orientate pe RIN-MN.....	279
5.5.1	Descrierea sistemului.....	279
5.5.2	Tehnica de modelare.....	280
5.5.3	Un model generalizat.....	290
5.6	Concluzii	294
6.	MODELUL EXPERIMENTAL	297
6.1	Arhitectura SMP-MP utilizat ca și model	297
6.2	Structura rețelei de interconectare de tip RMM	301
6.2.1	Algoritm de rutare pentru RMM	315
6.2.2	Analiza de performanță a RMM	321
6.2.3	Rezultate numerice și discuții	332
6.2.4	Toleranța la defect și fiabilitatea RMM	340
6.3	Protocoale de menținere a coerenței	344
6.4	Simularea pilotată de execuție	352
6.4.1	Construcția simulării	352
6.4.2	Interfața cu rețeaua	353
6.4.3	Parametrii de simulare	354
6.4.4	Protocolul de menținere a coerenței și sincronizarea	355
6.4.5	Programe de evaluare a performanțelor	355
6.4.6	Rezultate ale simulării	356
6.5	Efectul politicilor de management al memoriei	358
6.5.1	Algoritm pentru determinarea amplasării paginilor bazată pe acces	360
6.6	Concluzii	361
7.	CONCLUZII GENERALE	363
	ANEXA 1	373
	ANEXA 2	379
	ANEXA 3	383
	BIBLIOGRAFIE	389

LISTĂ DE TABELE

TABEL		Pagina
I	Rezumatul caracteristicilor hardware a trei RIN	111
II	Rezultatul analizelor RIN <i>crossbar</i>	127
III	Sumarul facilităților hardware a trei RIN	131
IV	Scalarea componentelor de calcul, de comunicare și a raportului	138
V	Mesajele transmise în noduri pentru menținerea coerenței	156
VI	Ordonările impuse de diferite modele de consistență	169
VII	Condiții de defectare	204
VIII	Patru scheme diferite cu <i>rulare înainte</i>	219
IX	Tabelul de căutare a conexiunilor	285
X	Tabelul nodurilor de comutare	291
XI	Lungimile de cale pentru diferite direcționări	320
XII	Numărul de destinații pentru diferite dimensiuni ale RIN	328
XIII	Dimensiunea optimă a unui comutator pentru eficacitatea costului	339
XIV	Utilizarea procesului și timpul de răspuns la RMM și RMB	340
XV	Parametrii de simulare	354
XVI	Caracteristici ale aplicațiilor utilizate în evaluare	356
XVII	Latențele medii ale mesajelor utilizând RMM și RMB	357

LISTĂ DE FIGURI

Figura		Pagina
1.1	Numărul de procesoare în SMP-MP comerciale	2
1.2	Lățimea de bandă a SMP-MP disponibile pe piață	2
1.3	Model tipic de memorie pentru programe paralele cu memorie partajată	4
1.4	Extensia unui sistem prin adăugarea de module procesor	5
1.5	Scheme de interconectare tipice în SMP-MP	6
1.6	Organizarea unui SMP-MP scalabil organizat NUMA	8
1.7	Un model de abstractizare a transferului de mesaj la nivel de utilizator	10
1.8	Un exemplu de problemă de coerență a memoriei cache	16
1.9	Comportamentul la saturație pentru o rețea de interconectare tipică	26
1.10	Limitări și beneficii ale toleranței la latență	30
1.11	Oferta de SMP în funcție de dependabilitate / performanță	33
2.1	Relația dintre caracteristicile algoritmilor și arhitecturile SMP	36
2.2	SMP-MP din subtaxonul MIMD-T	37
2.3	SMP-MD din subtaxonul MIMD-L	37
2.4	SMP-MP cu arhitectură tip ”sală de dans”	37
2.5	O clasificare a SMP paralel/distribuite	39
2.6a	SMP-MP cu memorii cache private	46
2.6b	SMP-MP cu memorii cache partajate	47
2.6c	SMP-MP cu memorii cache private și partajate	47
2.7	SMP-MP cu magistrală unică utilizând tehnica <i>snooping</i>	47
2.8	Diagrama de stări pentru un protocol <i>invalidare de scriere</i>	49
2.9	Taxonomia schemelor de asigurare a coerenței la memoriile cache	51
2.10	Condiții pentru un acces la un bloc cache perimat	57
2.11	Sumarul mecanismelor de asigurare a coerenței memoriilor cache mazate pe <i>autoinvalidare</i>	58

2.12	Clasificare a rețelelor de interconectare – RIN	60
2.13	Structura unui nod de interconectare la RIN	61
2.14	Exemple de RIN <i>statice</i>	62
2.15	Exemple de interconectări	63
2.16	SMP-MP orientat pe magistrală	65
2.17	Structuri de memorii ierarhizate utilizate pentru optimizarea traficului	68
2.18	SMP-MP cu <i>magistrală multiplă</i>	71
2.19	RIN <i>arbore</i>	72
2.20	Hipercub cu 16 procesoare	74
2.21	Comutatoare	76
2.22	Stările unui comutator 2 X 2	77
2.23	Comanda unei intrări la un comutator 2 X 2	78
2.24	Conexiuni în cubul de ordin 3	79
2.25	Interconectare <i>amestec perfect</i>	80
2.26	Interconectare <i>butterfly</i>	81
2.27	Interconectare <i>bit oglindit</i>	82
2.28	Rețea Illiac	83
2.29	Interconectare bazată pe <i>deplasare</i>	83
2.30	Schema bloc a unei RIN <i>dinamică</i>	86
2.31	RIN ce realizează funcția <i>amestec perfect</i>	86
2.32	RIN <i>crossbar</i>	88
2.33	Memorii multiport	88
2.34	Schema bloc a unei RIN <i>dinamice</i>	89
2.35	RIN <i>omega</i> ce conectează o intrare cu o ieșire	90
2.36	Situație conflictuală într-o RIN <i>omega</i>	91
2.37	RIN <i>delta</i> pentru $n=2$	92
2.38	Clasificarea topologiilor RIN	93
2.39	Procese partajând o structură de date D	95
2.40	SMP-TM cu 8 noduri PM organizate în cub boolean	103
2.41	Structura proceselor Occam și canalul lor de comunicare	104

3.1	Structura timpului de procesare	112
3.2	Saturația magistralei	114
3.3	Raportul între numărul de procesoare în așteptare și cele active	115
3.4	Structura unui SMP-MP cu magistrală multiplă	116
3.5	Structura unei magistrale multiple	117
3.6	Comportarea asimptotică a parametrului BW	122
3.7	Schema bloc a unui arbitru pentru magistrala multiplă	123
3.8	Structura registrului de stare	124
3.9	Un model cu șiruri de așteptare pentru SMP-MP <i>crossbar</i>	128
3.10	P_A funcție de talia sistemului	133
3.11	BW funcție de p pentru un sistem sincron	134
3.12	P_u funcție de p	134
3.13	Structura de bază a unui SMP-MP cu ierarhie de memorie	136
3.14	Arhitectura de bază pentru un SMP-TM cu memorie partajată distribuită	136
3.15	Problema coerenței pentru o locație de memorie cache	139
3.16	Protocol de invalidare în cazul unei magistrale cu monitorizare	142
3.17	Protocol de <i>actualizare la scriere</i>	142
3.18	Mecanism de asigurare a coerenței	144
3.19	Protocol cu <i>invalidare de scriere</i> pentru MCH cu rescriere	145
3.20	Diagrama de stare a coerenței MCH	146
3.21	Variația ratelor de eșec în funcție de numărul de procesoare	148
3.22	Variația ratelor de eșec în funcție de talia memoriei cache	148
3.23	Rata de eșec în funcție de talia blocului	149
3.24	Traficul pe magistrală în funcție de talia blocului	149
3.25	Componentele ratei de eșec în funcție de talia memoriei cache	150
3.26	Componentele ratei de eșec în funcție de talia blocului	150
3.27	Rata de eșec în funcție de talia memoriei cache	151
3.28	Evoluția traficului în porțiunea de sistem de operare	151
3.29	Arhitectura SMP-TM cu directori atașați nodului	154
3.30	Diagrama tranzițiilor de stare pentru un bloc cache	158

3.31	Diagrama tranzițiilor de stare pentru directori la o memorie cache individuală	159
3.32	Rata de eșec în funcție de numărul de procesoare	162
3.33	Scăderea ratelor de eșec pe măsura creșterii taliei memoriei cache	162
3.34	Rata de eșec în funcție de talia blocului într-un SMP cu 64 de procesoare și 128 Kbyte memorie cache	163
3.35	Creșterea numărului de bytes per referință de dată în funcție de creșterea taliei blocului	164
3.36	Un exemplu de consistență secvențială	165
3.37	Reducerea numărului de ordine impuse pe măsura relaxării modelelor	170
3.38	Performanța modelelor de consistență relaxate	172
4.1	Arborele dependabilității	179
4.2	Un sistem SMR	186
4.3	Agrementul Bizantin	189
4.4	Agrementul Bizantin după efectuarea algoritmului	190
4.5	Un Agreement Bizantin imposibil	190
4.6	Exemplu de mapare a grafurilor duplicate	195
4.7	Crearea de puncte de verificare bazate pe procesor	201
4.8	Tehnici tolerante la defect pentru copierea în cascadă	203
4.9	Conceptul de bază pentru crearea unui PV virtual	205
4.10	Maparea unei singure pagini	206
4.11	Prima referire după un PV	207
4.12	Pagină referită anterior	207
4.13	O diagramă de timp pentru un calcul distribuit	211
4.14	Tăieturi consistente și inconsistente	212
4.15	Un exemplu de mesaj <i>orfan</i>	213
4.16	Un exemplu de mesaj <i>pierdut</i>	214
4.17	Situația blocajului activ	214
4.18	Relansarea concurentă în schemele RFCS	218
4.19	Tehnica de relansare concurentă în schemele RFCS	219

4.20	Schema optimistă cu sau fără modul de rezervă	220
4.21	Scheme pesimiste	221
4.22	Diferite scheme ce utilizează <i>rularea înainte</i>	222
4.23	Comparația în performanță dintre schemele optimiste și pesimiste	223
4.24	Comparații ale siguranței în funcționare între schemele optimiste și pesimiste	223
4.25	Întârzierea permanentă la ieșirile schemelor cu <i>rulare înapoi</i>	224
4.26	Întârzierea temporară la ieșirile schemelor cu <i>rulare înainte</i>	225
4.27	RIN <i>crossbar</i> tolerantă la defect	229
4.28	Reconfigurarea unui <i>hipercub</i>	232
4.29	Reconfigurarea unui <i>hipercub</i> utilizând noduri de extensie	234
4.30	Direcționarea în <i>hipercuburi</i>	235
4.31	Reconfigurarea unei RIN <i>plasă</i>	237
4.32	Reconfigurarea unei RIN <i>plasă</i> prin mapare logică	238
4.33	Reconfigurarea în RIN <i>arbore</i>	240
4.34	O diagramă generică pentru RIN <i>multinivel</i>	241
4.35	RIN <i>omega</i> 16 X 16	247
4.36	RIN-MN cu identificarea subseturilor conjugate de comutatoare	248
4.37	Exemplu de izolare a unui comutator defect	248
4.38	Scheme de realizare a unor conexiuni multiple	251
4.39	Un exemplu de rețea ASEN-Max	252
4.40	O rețea ASEN-2 pentru $N=16$	252
4.41	Grafurile de redundanță pentru rețelele ASEN 2 și ASEN 4	257
4.42	Probabilitatea de terminal efectivă pentru ASEN	258
4.43	Parametrul <i>MTTF</i> pentru diferite RIN-MN	258
4.44	Probabilitatea de acceptanță pentru diferite RIN-MN	259
5.1	O diagramă bloc a fiabilității unei RIN	267
5.2	Un SMP-MP cu RIN <i>butterfly</i> de dimensiuni 16 X 16	270
5.3	Un lanț Markov pentru un subsistem procesor	273
5.4	Un lanț Markov pentru o RIN-MN	275
5.5	Fiabilitatea unui SMP-MP de talie 16 X 16	277

5.6	Dependabilitatea raportată la performanță	277
5.7	Distribuția mediată în timp a BW	277
5.8	Decompoziția unui sistem	280
5.9	Lanț Markov pentru subsistemul SPM	281
5.10	Lanț Markov pentru RIN-MN	286
6.1	Arhitectura SMP experimental	299
6.2	Arhitectura SMP cu ierarhie de memorie	300
6.3	Un SMP de talie 16 X 16 bazat pe RMM	302
6.4	O rețea RMM 16 X 16 ce utilizează comutatoare 4 X 4	302
6.5	O rețea RMM 16 X 16 ce utilizează comutatoare 2 X 2	304
6.6	Diagrama bloc a unui comutator 2 X 2 din rețeaua RMM	305
6.7	Diagrama bloc a unui controller de trafic	306
6.8	Înlănțuirea pachetelor într-un comutator RMM	308
6.9	Rutarea RD în RMM	310
6.10	Rutarea RI în RMM	311
6.11	Rutarea BD în RMM	312
6.12	Rutarea BI în RMM	313
6.13	Cozi de așteptare în comutatoare RMM și RMB	329
6.14	Comparație a rezultatelor analizei și simulării pentru P_u la RMM	334
6.15	Comparație a rezultatelor analizei și simulării pentru T_r la RMM	334
6.16	Comparație între utilizările procesorului, variind m	335
6.17	Comparație între timpii de răspuns, variind m	335
6.18	Utilizarea comutatorului în RMM și în RMC variind dimensiunea lui	336
6.19	Eficacitatea costului în RMM și în RMC variind dimensiunea	336
6.20	Utilizarea procesorului: scalabilitatea RMM în raport cu RMB	338
6.21	Timp de răspuns: scalabilitatea RMM în raport cu RMB	338
6.22	Graf de redundanță în RMM	342
6.23	Fiabilitatea diferitelor RIN pentru comutatoare 2 X 2	343
6.24	Fiabilitatea diferitelor RIN pentru comutatoare 4 X 4	343
6.25	Un arbore ce conectează modulul M_4 cu toate procesoarele	345
6.26	Diagrame ale tranzițiilor de stare	346

6.27	Configurarea adresei fizice	350
6.28	Organizarea directorului într-un comutator RMM	351
6.29	Simulator pilotat prin execuție	352
6.30	Configurația unui nod și interfața cu rețeaua	353
6.31	Rezultate ale simulării bazate pe execuție	357
6.32	Organizare de tip IGR a unui SMP-MPD	359

LISTA PRINCIPALELOR ABREVIERI

ASEN	Rețea de interconectare augmentată bazată pe amestec perfect
BD	Buclare directă
BI	Buclare indirectă
CCH	Coerența memoriei cache
LAN	Rețea cu arie locală
M	Memorie
MC	Memorie comună
MCH	Memorie cache
MIMD	SMP din taxonul „instrucții multiple, date multiple”
MIMD-L	SMP din taxonul MIMD, slab cuplat
MIMD-T	SMP din taxonul MIMD, strâns cuplat
MK	Lanț Markov
ML	Memorie locală
MM	Magistrală multiplă
MP	Memorie partajată
MPS	Multiprocesor simetric
NUMA	Acess non-uniform la memorie
P	Procesor
<i>P</i>	Proces
PE	Element de procesare
PM	Cuplu procesor-memorie
PV	Punct de verificare
RD	Rutare directă
RI	Rutare indirectă
RIN	Rețea de interconectare
RIN-CB	Rețea de interconectare <i>crossbar</i>
RIN-MM	Rețea de interconectare cu magistrale multiple
RIN-MN	Rețea de interconectare multinivel
RIN-MN(B)	Rețea de interconectare multinivel bidirecțională

RMB	Rețea de interconectare multinivel bidirecțională
RMM	Rețea de interconectare bazată pe magistrală multinivel
SAC-MCH	Schemă de asigurare a coerenței memoriilor cache
SE	Element de comutare
SMP	Sistem multiprocesor
SMP-MP	Sistem multiprocesor cu memorie partajată
SMP-MPD	Sistem multiprocesor cu memorie partajată distribuită
SMP-TM	Sistem multiprocesor cu transfer de mesaje
SW	Comutator
UMA	Acces uniform la memorie

1. INTRODUCERE

Calculul paralel pe scară largă devine un instrument indispensabil în diferite domenii ale științei și tehnologiei cu un grad ridicat de diversitate. Proiectarea și evaluarea sistemelor multiprocesor reprezintă abordări complexe și necesită explorarea unui spațiu de proiectare foarte larg.

1.1 TENDINȚE ARHITECTURALE

Noile tendințe în arhitecturile SMP sunt în strânsă legătură cu dezvoltarea tehnologică în sensul convertirii acesteia în performanță și capabilitate. Fundamental, cele două căi prin care un volum mare de resurse (de exemplu mai mulți tranzistori) ameliorează performanța, sunt *paralelismul* și *localitatea*. Cele două abordări concurează în ultimă instanță pentru aceleași resurse. În general, cea mai bună performanță este obținută printr-o strategie combinată ce alocă resurse pentru a exploata un anumit grad de paralelism și un anumit grad de localitate. Paralelismul și localitatea interacționează pe diferite niveluri, de la circuit integrat până la structură.

În fig. 1.1 se prezintă evoluția în timp a sistemelor multiprocesor, evidențiindu-se creșterea numărului de procesoare.

În fig. 1.2 se evidențiază lățimea de bandă în diferite multiprocesoare comerciale. Se remarcă că după o creștere lentă în decurs de câțiva ani, a început o nouă eră în 1991, cea a sistemelor bazate pe magistrală ce suportă creșterea substanțială a numărului de procesoare foarte rapide.

1.2 CONVERGENȚA ARHITECTURILOR PARALELE

Din punct de vedere istoric, mașinile paralele s-au dezvoltat în jurul câtorva concepte arhitectonice și cele mai multe discuții se poartă și astăzi în jurul taxonomiei acestor proiecte. Proiectanții sunt influențați de aceleași forțe tehnologice și de necesități similare generate de aplicații. Din această cauză a apărut un pronunțat fenomen de convergență în domeniu. El se manifestă mai ales în arhitecturile aparținând ultimului deceniu.

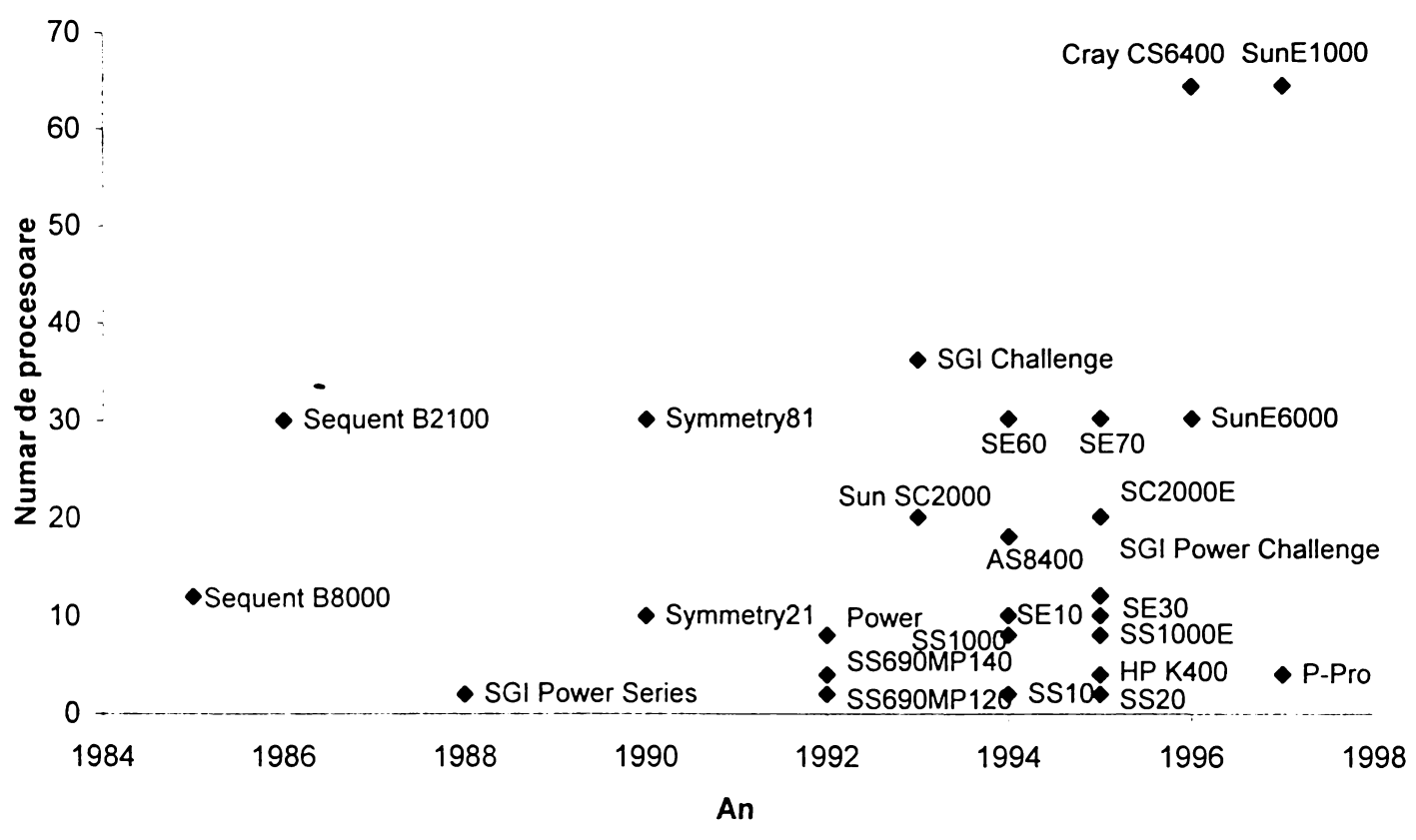


Fig. 1.1 Numărul de procesoare în SMP-MP comerciale

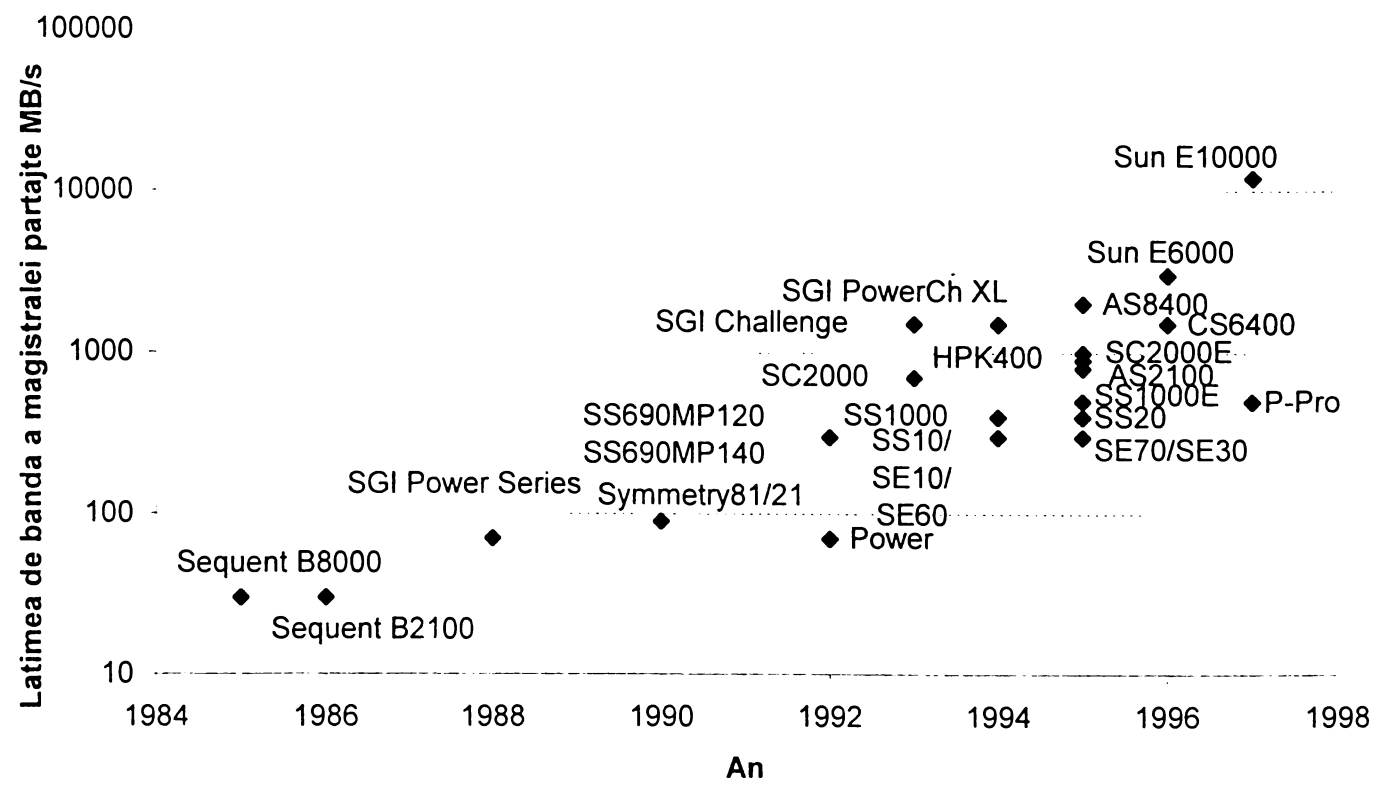


Fig. 1.2 Lățimea de bandă a SMP-MP disponibile pe piață

1.2.1 Arhitectura de comunicație

Plecând de la plastica definiție a lui Almasi și Gottlieb în care un calculator paralel este văzut ca și ”o colecție de elemente de procesare ce comunică și cooperează pentru a rezolva rapid mari probleme” [AG89], este rațional să se privească o arhitectură paralelă ca fiind extensia unei arhitecturi convenționale. În arhitectura paralelă se utilizează comunicarea și cooperarea între procesoare. În esență, arhitectura paralelă extinde conceptele uzuale ale arhitecturii calculatorului cu cele ale unei arhitecturi de sisteme de comunicație. Arhitectura calculatorului are două fațete distincte. Una se referă la frontiera hard/soft și la frontiera utilizator/sistem. Cea de-a doua se referă la structura organizatorică ce realizează aceste cerințe pentru a obține o performanță ridicată la un cost convenabil.

O *arhitectură de comunicație* definește modelul de comunicare de bază și sincronizarea structurând modulele ce realizează aceste operații. Modelele de programare în arhitecturile paralele sunt: *adresa partajată, transferul de mesaje, prelucrarea în paralel a datelor, multiprogramarea.*

Din perspectiva istorică a unui model de programare se poate gândi arhitectura ca fiind îmbinarea armonioasă a *modelului de programare*, a *modelului de comunicare* și a *organizării sistemului*. Rezultă, în această manieră, arhitecturi cu memorie partajată, cu transfer de mesaje, etc.

1.2.2 Spațiul de adresă partajat

Una dintre cele mai importante clase de sisteme paralele o reprezintă *multiprocesoarele cu memorie partajată* (SMP-MP). O proprietate fundamentală a acestei clase este apariția comunicării în mod implicit ca și un rezultat al instrucțiilor de acces la memoria convențională. Un model primar de programare pentru aceste sisteme este în esență unul de partajare a timpului pe un singur procesor, cu excepția faptului că paralelismul real înlocuiește interpătrunderea în timp. În mod formal, un *proces* este un spațiu de adresă virtuală și unul sau mai multe fluxuri de comandă. Procesele pot fi configurate de o astfel de manieră încât porțiuni ale spațiului de adresă virtuală sunt partajate, ceea ce înseamnă că sunt mapate la o locație fizică comună, așa cum se sugerează în fig. 1.3.

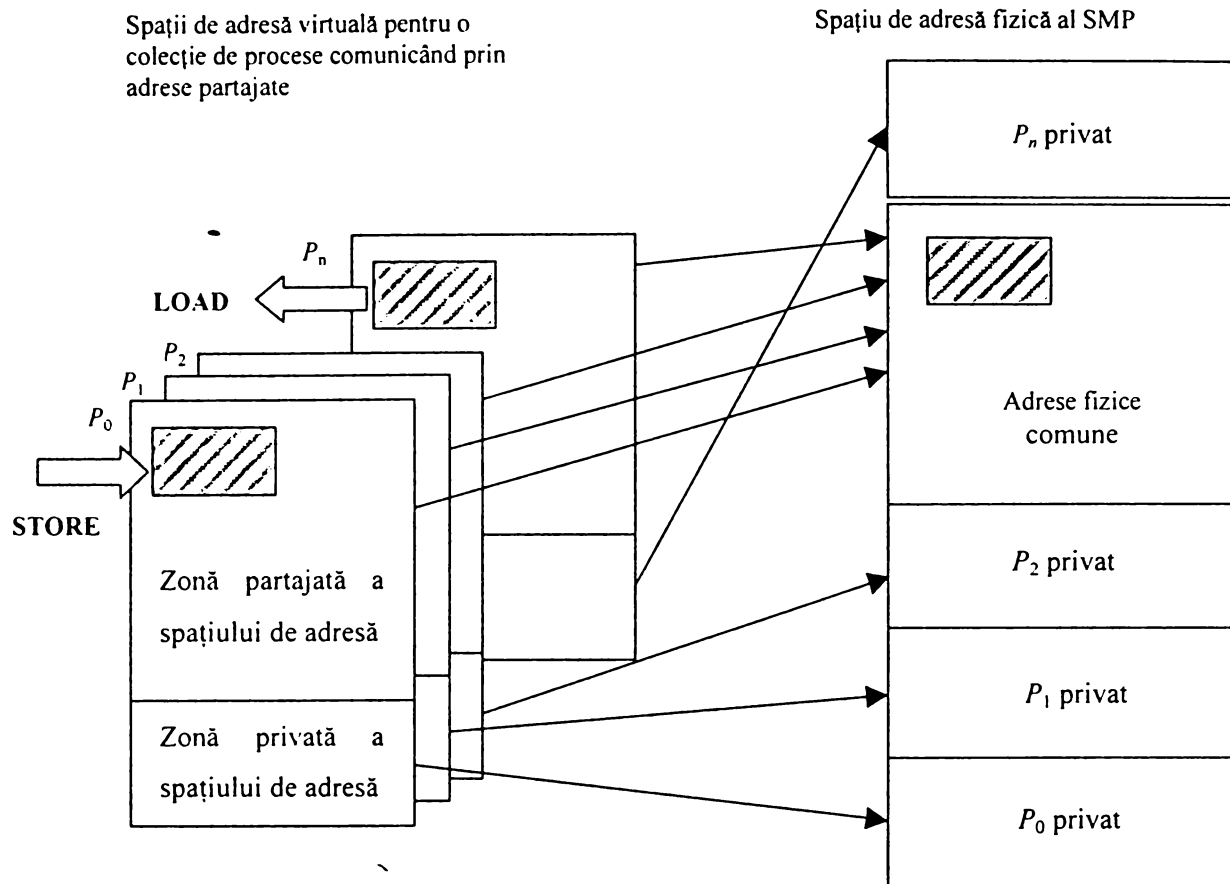


Fig. 1.3 Model tipic de memorie pentru programe paralele cu memorie partajată

Cooperarea și coordonarea pentru fluxuri sunt realizate prin citire și scriere de variabile partajate ca și de pointeri ce se referă la adresele partajate. *Scrierile la o adresă logic partajată de către un flux sunt vizibile pentru citirile altor fluxuri.*

Arhitectura de comunicare utilizează operațiuni convenționale referitoare la memorie pentru a asigura comunicarea printre adresele partajate în aceeași măsură ca și operațiuni speciale la nivel atomic pentru sincronizare.

Hardware-ul de comunicație pentru SMP-MP reprezintă o extensie naturală a sistemului de memorie ce se regăsește în cele mai multe dintre calculatoare. În esență, toate sistemele de calcul paralele permit unui procesor și unui set de controlere I/E să acceseze o colecție de module de memorie printr-un sistem de interconectare. Capacitatea memoriei este mărită prin adăugarea de module suplimentare. Această capacitate suplimentară poate sau nu să mărească lățimea de bandă a memoriei disponibile, depinzând de organizarea specifică a sistemului. Capacitatea I/E se mărește prin

adăugarea de dispozitive I/E la controlerele I/E sau prin inserarea de noi controlere. Există două căi posibile de creștere a capacității de procesare: să se aștepte apariția unor procesoare mai rapide sau să se adauge mai multe procesoare. Configurația cadru a unui astfel de SMP-MP este prezentată în fig. 1.4.

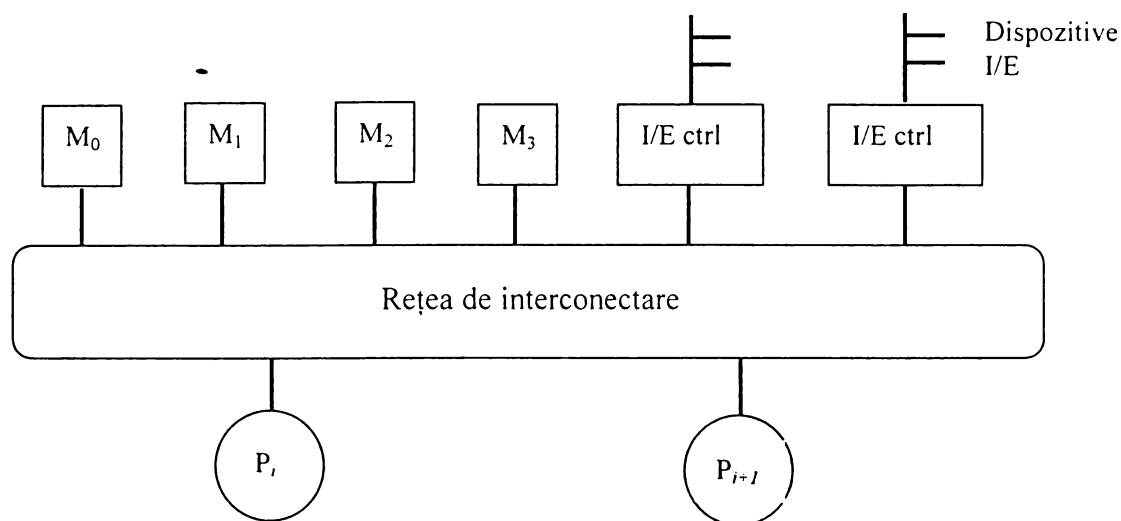
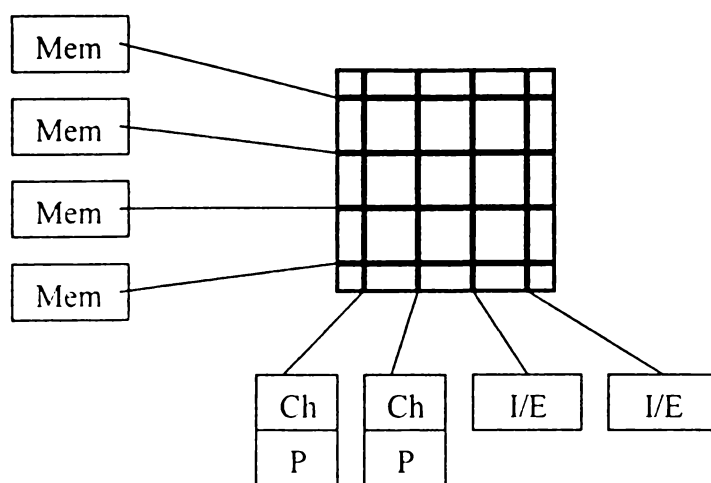
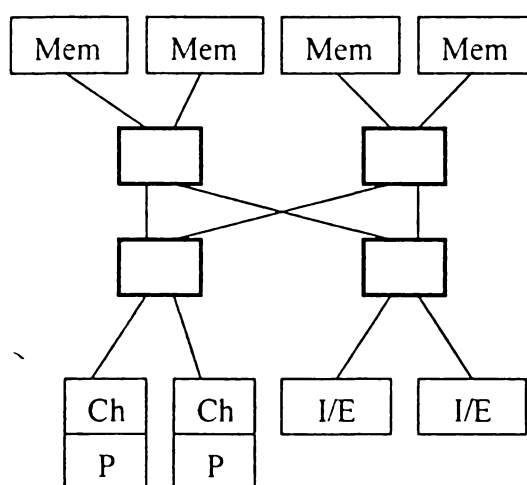


Fig. 1.4 Extensia unui sistem într-un SMP-MP prin adăugarea de module procesor

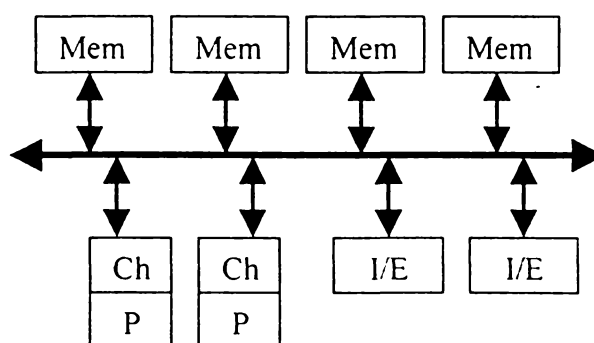
Din această configurație au derivat și altele, în funcție de avansul tehnologic.

În fig. 1.5 a), b), c) se prezintă cele trei configurații de bază ale SMP-MP.

Inițial sistemele au fost organizate ca și în fig. 1.5 a) unde o conexiune *crossbar* conectează CPU și câteva canale I/E la câteva module de memorie. Adăugarea de noi procesoare a reprezentat în primă instanță o problemă de expandare a dimensiunii comutatorului; structura hardware de acces la memorie a rămas nemodificată. Dimensiunea și costul procesorului au limitat aceste structuri timpurii la un număr mic de procesoare, dar pe măsură ce densitatea de încapsulare hardware s-a ameliorat și costurile au scăzut, au putut fi proiectate sisteme mai mari. Costul scalării unei rețele *crossbar* a devenit un factor limitativ și în cele mai multe dintre cazuri rețeaua a fost înlocuită cu *rețeaua de interconectare multinivel* (RIN-MN) ca și în fig. 1.5 b). În acest caz costul a evoluat mult mai încet în funcție de numărul de porturi. Aceste economii s-au făcut însă pe cheltuiala unei latențe mai mari și a unei lățimi de bandă mai mici. Abilitatea de a accesa toată memoria direct de la fiecare procesor are câteva avantaje: fiecare procesor poate rula orice proces sau poate manipula orice eveniment I/E iar structurile de date pot fi partajate în cadrul sistemului de operare (SO).

a) Comutator *crossbar*

b) Rețea de interconectare multinivel



c) Interconectare de tip magistrală

Fig. 1.5 Scheme de interconectare tipice în SMP-MP

Utilizarea masivă a SMP-MP a început odată cu revoluția microprocesorului pe 32 bit la mijlocul anilor '80, atunci când procesorul, cache-ul, unitatea de virgulă flotantă și unitatea de management a memoriei au putut fi încorporate pe o singură placă. Multe sisteme de rang mediu, incluzând minicalculatoare, servere, stații de lucru și PC-uri sunt organizate în jurul unor magistrale ca și în fig. 1.5 c). Mecanismul de acces standard al magistralei permite oricărui procesor să acceseze orice adresă fizică din sistem. Ca și în sistemele bazate pe comutatoare, toate locațiile de memorie sunt echidistante față de toate procesoarele. Astfel toate procesoarele experimentează același timp de acces sau aceeași latență la o referire de memorie. Această configurație este adeseori denumită *multiprocesor simetric (MPS)*.

Factorii ce limitează numărul de procesoare ce pot fi suportate într-o organizare bazată pe magistrală sunt diferiți de aceia dintr-o structură bazată pe comutator. Adăugarea de procesoare la comutator este o operație costisitoare; oricum, lățimea de bandă crește cu numărul de porturi. Costul adăugării unui procesor la magistrală este redus, dar lățimea de bandă este fixată. Divizând această lățime de bandă fixată pentru mai multe procesoare se limitează scalabilitatea abordării. Problema este parțial rezolvată prin introducerea memoriilor cache ce reduc solicitările pentru lățime de bandă dar introduc alte limitări. Datorită datelor replicate în cache-urile locale, apare însă problema consistenței și coerenței, ce vor fi tratate în continuare.

Scalabilitatea alături de latență și de lățime de bandă, reprezintă un alt parametru important în evaluarea performanței unui SMP. SMP-MP orientate pe magistrală nu sunt scalabile deoarece au lățime de bandă fixă. *Crossbar*-ul nu scalează convenabil deoarece costul evoluează cu patratul numărului de porturi. Există multe alternative de RIN scalabile în care lățimea de bandă crește pe măsură ce se adaugă procesoare iar costul nu crește excesiv.

O abordare naturală în construcția unor SMP-MP scalabile este menținerea *accesului uniform la memorie (UMA)* ca și în arhitectura tip "sală de dans" prezentată în fig. 1.4 și furnizarea unei interconectări scalabile între procesoare și memorii. Fiecare acces la memorie este convertit într-o tranzacție-mesaj prin RIN, mult mai profund decât s-ar fi putut face în cazul MPS. Dezavantajul major al acestor abordări este acela că latența RIN este experimentată la fiecare acces la memorie și trebuie să fie furnizată o mare lățime de bandă fiecărui procesor.

O abordare alternativă o reprezintă interconectarea completă a procesoarelor, fiecare cu memoria sa locală, ca și în fig. 1.6. În această abordare cu *acces neuniform la memorie* (NUMA), controlerul memoriei locale determină dacă este cazul să se facă un acces la memoria locală sau o tranzacție-mesaj cu un controler al unei memorii situate la distanță. Accesarea memoriei locale este mai rapidă decât accesarea memoriei la distanță. Sistemul I/E poate fi fie parte din fiecare nod, sau poate fi consolidat în noduri I/E speciale. Accesesele la datele private, ca și în cazul codurilor sau al stivei, pot fi frecvent operate local. La fel se prezintă situația în cazul datelor partajate ce sunt înmagazinate în nodul local.

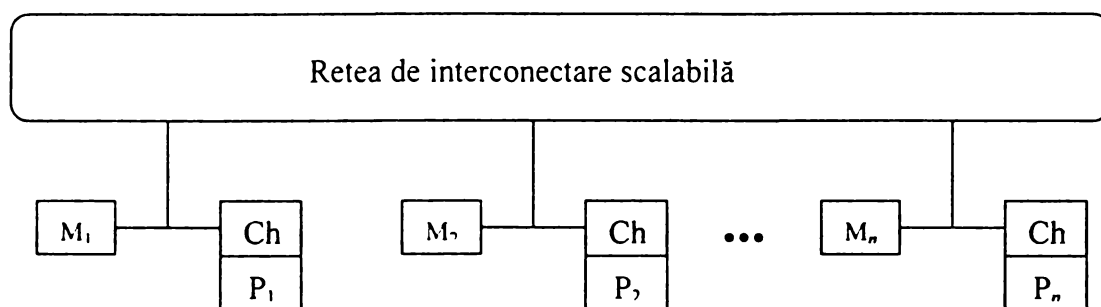


Fig. 1.6 Organizarea unui SMP-MP scalabil cu acces neuniform la memorie (NUMA)

Abilitatea de a accesa rapid memoria locală nu crește semnificativ timpul de acces la datele situate la distanță. Se reduce însă timpul mediu de acces, mai ales în situația în care o mare parte din accese sunt pentru datele locale. De asemenea se reduce cererea de lățime de bandă, plasată pe RIN. Deși a avea toate datele partajate situate echidistant față de orice procesor reprezintă o anumită simplificare în proiectare, abordarea NUMA a devenit prevalentă în majoritatea sistemelor moderne. Cităm sistemele CRAY T3E și SGI Origin.

În concluzie, comunicarea și cooperarea în modelul de programare cu spațiu de adresă partajată constau în citiri și scrieri de variabile partajate; aceste operații sunt mapate direct la un model de comunicare constând în instrucții LOAD/STORE ce accesează un spațiu global de adrese partajate. Fiecare procesor poate *nominaliza* fiecare adresă fizică din mașină. Un proces poate nominaliza toate datele pe care le partajează cu altele în cadrul spațiului de adrese virtuale.

Datele sunt transferate fie ca și primitive în setul de instrucții (byte, cuvinte, etc.) fie în blocuri cache. Fiecare proces efectuează operații referitoare la memorie pe adrese

aparținând spațiului de adrese virtuale. Procesul de translatăre a adresei identifică o adresă fizică, ce poate fi locală sau la distanță față de procesor și poate fi partajată cu alte procese. Translatărea de adresă realizează protecția în cadrul spațiului de adresă partajată, ca și în cazul uniprocsoarelor, deoarece un proces poate doar accesa datele în spațiul de adresă virtuală.

Efăcitatea SMP-MP depinde de latența înglobată în accesele la memorie ca și de lățimea de bandă a transferului de date ce poate fi suportat de sistem.

1.2.3 Transferul de mesaj

O a doua clasă importantă de mașini paralele, denumită *arhitecturi cu transfer de mesaje*, utilizează calculatoare complet echipate ca și blocuri constitutive și furnizează comunicare între procesoare în termeni de operații explicite de I/E. Diagrama bloc de nivel înalt pentru un *sistem multiprocesor cu transfer de mesaje* (SMP-TM) este în esență similară cu abordarea NUMA în cazul SMP-MP, prezentată în fig. 1.6. Prima diferență constă în modul de integrare al comunicării la nivel de I/E și nu de memorie. Acest stil de proiectare are multe în comun cu rețelele de stații de lucru, sau *grupuri (clusters)*, exceptând faptul că gradul de compactare a nodurilor este mult mai masiv, neexistând monitor sau tastatură la nivelul nodului, iar RIN are o capabilitate mult mai mare decât LAN-urile standard. Integrarea dintre procesor și rețea este mult mai puternică decât în structurile tradiționale I/E, deoarece transferul de mesaje este în mod fundamental o comunicare procesor la procesor.

În transferul de mesaje, există o distanță substanțială între modelul de programare și primitivele hardware, existând o comunicare la nivel de utilizator efectuată prin SO sau apeluri din bibliotecă, ceea ce provoacă multe acțiuni la nivel inferior.

Cele mai comune acțiuni de comunicare la nivelul utilizatorului sunt variante de SEND/RECEIVE. În forma sa cea mai simplă, SEND specifică un tampon de date local al cărui conținut trebuie transmis și un proces receptor, de obicei rezident într-un procesor situat la distanță. RECEIVE specifică un proces transmițător și un tampon local în care data transmisă va fi depozitată. Interacțiunea dintre un SEND și un RECEIVE cauzează un transfer de date de la un proces la altul, situație prezentată în fig. 1.7.

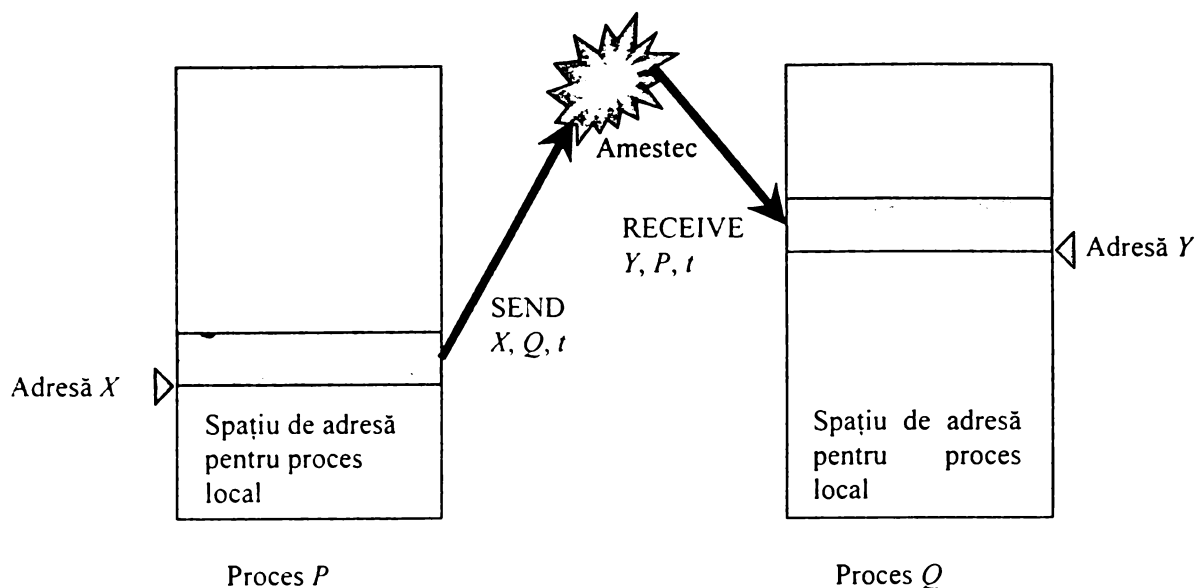


Fig. 1.7 Un model de abstractizare a transferului de mesaj de tip SEND/RECEIVE la nivel de utilizator

În cele mai multe dintre SMP-TM, operația SEND permite de asemenea atașarea la mesaj a unui indicator sau *etichetă de rutare* iar la recepție se specifică regula de dezmembrare a mesajului. Astfel, programul utilizatorului nominalizează adresele locale și intrările într-un spațiu abstract proces-etichetă. Există câteva posibile variante de sincronizare, depinzând de faptul că SEND se mai execută în timp ce RECEIVE a fost executat, când tamponul de transmisie este disponibil pentru reutilizare sau când cererea a fost acceptată.

SMP-TM din primele generații au furnizat primitive hardware ce erau foarte apropiate de nivelul de comunicare SEND/RECEIVE la nivel de utilizator, având însă câteva restricții adiționale. Multe SMP-TM în fazele incipiente au fost *hipercuburi* unde fiecare nod este conectat la n alte noduri dintr-un total de 2^n noduri. Altele au fost *plase*, unde nodurile sunt conectate la vecini pe două sau trei dimensiuni. Topologia RIN a fost importantă în SMP-TM incipiente deoarece doar procesoarele învecinate ar fi putut fi nominalizate în cadrul operațiilor SEND/RECEIVE.

Topologia fizică a RIN a impus la început denumirea algoritmilor paraleli. Ei au fost denumiți: *inel*, *hipercub*, etc. Accentul pe topologia RIN a fost semnificativ redus odată cu introducerea RIN cu un grad ridicat de generalitate, care asigurau transferul în loturi a mesajului prin fiecare dintre mecanismele de rutare ce formau RIN. Câteva exemple sunt: IBM-SP2, RS6000, Intel Paragon, etc.

1.2.4 Convergența

Evoluția SMP în termeni de hardware și software a estompat frontiera, odată clar delimitată, dintre SMP-MP și SMP-TM. Se constată o anumită convergență între cele două arhitecturi ce apar nu numai în concepte ci și în realizările structurale [CSG99].

În primul rând, la nivelul *operațiunilor de comunicare* disponibile la procesele utilizatorilor există:

- Transferuri de mesaje tradiționale, SEND/RECEIVE, ce sunt suportate de cele mai multe din SMP-MP prin tampoane de stocare partajate. SEND presupune scrierea de date sau o indicare a lor în tampon; RECEIVE implică citirea datelor din tampoanele partajate. Fanioanele și blocajele sunt utilizate pentru a comanda accesul spre tampon și pentru a indica evenimente ca și sosirea mesajelor.
- La un SMP-TM, un proces utilizator poate construi un spațiu global de adresă prin transportul pointerilor specificând procesul și o adresă virtuală locală în acest proces. Accesul la o astfel de adresă globală poate fi obținut software printr-o tranzacție de mesaj explicită. O citire logică se realizează prin transmiterea unei cereri spre procesul conținând obiectul și recepționând un răspuns.
- Un spațiu de adresă virtuală poate fi stabilit la un SMP-TM la nivel de pagină. O colecție de procese are o regiune de adrese partajate dar, pentru fiecare proces, doar paginile ce sunt locale sunt accesibile. După un acces la o pagină ratată, de obicei situată la distanță, apare un defect de pagină și SO angajează nodul de la distanță într-o tranzacție de mesaj cu scopul de a transfera pagina și de a o proiecta în spațiul de adresă al utilizatorului.

La nivelul *organizării sistemului*, o convergență substanțială a apărut în ultimii ani. Arhitecturile SMP-TM apar, la nivel de schemă bloc, identice cu arhitecturile scalabile NUMA. În cazul memoriilor partajate, RIN a fost integrată cu controlerul cache sau cu cel al memoriei în scopul de a observa eșecurile de acces la cache și pentru a conduce la generarea unui mesaj pentru accesul la memoria dispusă într-un nod situat la distanță. În SMP-TM, interfața cu RIN este de obicei asigurată de un dispozitiv I/E. Tendința a fost de a integra acest dispozitiv mult mai profund în sistemul de memorie și de a transfera datele direct de la și spre spațiul de adrese al utilizatorului. Unele proiecte prevăd

transferul DMA prin RIN de la memoria unui sistem la memoria altui sistem, astfel ca interfața RIN să fie puternic integrată cu sistemul de memorie.

În plus față de convergența SMP-TM scalabile și a SMP-MP, au apărut LAN-uri bazate pe comutator ce includ Ethernet, ATM, și alte medii de comunicare, oferind interconectări scalabile ce se apropie de ceea ce mașinile paralele tradiționale pot oferi. Aceste noi rețele sunt utilizate pentru a conecta colecții de sisteme în *grupuri* care pot opera ca și o mașină paralelă asupra unor probleme mari sau ca și multe mașini individuale în multiprogramare.

În concluzie, *transferul de mesaj* și *spațiul de adresă partajată* reprezintă două modele distincte de programare, fiecare furnizând o bine definită paradigmă pentru *partajare*, *comunicare* și *sincronizare*. Prin urmare, este normal să se considere ambele aspecte într-un cadru comun. Integrarea mult mai puternică a comunicării în sistemul de memorie tinde să reducă latența tranzacțiilor prin RIN și să amelioreze lățimea de bandă ce poate fi furnizată sau acceptată de RIN.

1.3 COERENȚA ȘI CONSISTENȚA MEMORIILOR CACHE

Reducerea *latenței* comunicației reprezintă una dintre finalitățile pe care un proiectant de arhitecturi de sisteme de calcul paralele trebuie să o aibă în vedere. Această reducere se efectuează atât la nivel de ierarhie de memorie cât și la cel al rețelei de interconectare. În ceea ce privește ierarhia de memorie, un rol major îl au memoriile cache private situate pe nivelul superior, adică cel mai apropiat de procesor, și care reduc efectiv întârzierea medie a accesului la memorie. În sistemele multiprocesor, datorită dispersiei datelor printre procesoare și datorită problemei coerenței memoriilor cache, efectul benefic al introducerii memoriilor cache nu este atât de evident ca și în cazul sistemelor uniprocesor. În literatură [Lil93], [CF78], [CFKA90], [Choi96] sunt descrise o varietate largă de mecanisme de menținere a coerenței cache în SMP-MP de capacitate mare. Pentru a se putea face o comparare și în același timp o evaluare a performanței acestor mecanisme, arhitectul de sistem de calcul trebuie să ia în considerare următoarele:

- *strategia de detectare a coerenței* prin care posibilele accese incoerente la memorie sunt detectate fie statistic în timpul compilării, fie dinamic în timpul rulării programului

- *strategia de întărire a coerenței*, manifestată prin protocoale de actualizare sau de invalidare și care este utilizată pentru a se asigura că intrările în memoriile cache neactualizate nu vor fi niciodată referite de către un procesor
- *precizia informației referitoare la partajarea blocurilor* care dă o măsură a costului de implementare și a performanței unui mecanism de menținere a coerenței
- *dimensiunea blocului cache* ce afectează performanța sistemului de memorie

Secvența adreselor de memorie generate de către un program afișează proprietățile localității *temporale* și *spațiale* [Smit82]. *Localitatea temporală*, sau localitatea în timp, înseamnă că memoriile ce au fost recent adresate de către un program, vor fi din nou referite în viitorul foarte apropiat. *Localitatea spațială* înseamnă că adresele referite de către un program într-o perioadă scurtă de timp vor cuprinde o porțiune relativ mică a întregului spațiu de adresă. Memoriile cache private care sunt mici, rapide și situate lângă un procesor, exploatează aceste proprietăți de referire la memorie pentru a reduce timpul mediu necesar accesării memoriei principale care uzual are o capacitate mult mai mare. Prin înmagazinarea temporară în cache a unei copii a unei valori provenind din memoria principală, care a fost activ referită de către program, memoriile cache amortizează timpul cerut pentru copierea conținutului unei locații de memorie de la o memorie principală mai lentă într-o memorie cache mai rapidă. Datorită faptului că multiplele copii ale unei locații de memorie partajată pot fi rezidente simultan în câteva cache-uri diferite, memoriile cache private introduc problema *coerenței* în care există posibilitatea pentru ca diferite copii cache-ate să aibă valori diferite în același moment de timp. Este responsabilitatea mecanismului de menținere a coerenței cache – SAC-MCH – de a asigura că ori de câte ori un procesor citește o locație de memorie el recepționează valoarea corectă.

Există două aspecte importante referitoare la problema coerenței cache. Prima se referă la *modelul sistemului de memorie* prezentat programatorului. Al doilea aspect important îl reprezintă *meccanismul* utilizat de către sistem pentru menținerea coerenței între memoriile cache și memoria principală.

1.3.1 Modele de consistență

Una dintre definițiile unui SMP ce utilizează memorii cache coerente este aceea a unui sistem care garantează că ”... valoarea returnată a unei instrucțiuni Load este întotdeauna valoarea dată de cea mai din urmă instrucțiune Store ce a operat asupra aceleiași adrese” [CF78]. Dificultatea cu această definiție este că semnificația lui ”cel mai din urmă”-nu este precis definită atunci când instrucțiunile de încărcare și de memorare apar în diferite procesoare ce rulează asincron, unul în raport cu celălalt. Datorită întârzierilor și temporizărilor în diferite porțiuni ale RIN ce interconectează binomul procesor – memorie și în cadrul însăși a procesoarelor și memoriilor, fiecare procesor și fiecare modul de memorie poate observa o ordonare diferită a evenimentelor. *Modelul de consistență* al unui SMP definește punctul de vedere al programatorului referitor la ordonarea în timp a evenimentelor ce apar în diferite procesoare. Aceste evenimente includ citiri și scrieri în și din memorie, operații de sincronizare.

Din punctul de vedere al programatorului referitor la sistemul de memorie, modelul de *consistență secvențială* definește o ordonare strictă a secvenței de execuție a operațiilor de memorie permise în cadrul unui procesor și între procesoare. Se spune că un SMP este consistent secvențial dacă ”... rezultatul oricărei execuții a unui program este același ca și în cazul în care operațiile tuturor procesoarelor au fost executate într-o anumită ordine secvențială și operațiile fiecărui procesor individual apar în această secvență în ordinea specificată de către programul său” [Lamp79]. Această strictă ordonare a acceselor la memorie impune o penalizare severă în performanță prin limitarea drastică a suprapunerilor permise între operațiile de memorie livrate de către un procesor individual și de către alte procesoare.

Modelul de consistență cu *ordonare slabă* [DSB88] relaxează ordinea garantată a evenimentelor modelului de consistență secvențială pentru a permite o mai mare suprapunere a scrierilor și citirilor din memorie. Cu acest model doar accesele la memorie către variabilele de sincronizare definite de programator sunt garantate că vor apare într-o ordine secvențială consistentă. În rest, toate referirile la memorie produse de procesoare diferite pot apare în oricare ordine arbitrară.

Modelul de *consistență prin eliberare* [GGH91] relaxează și mai mult constrângerile de ordonare asupra variabilelor de sincronizare, divizând operația de sincronizare în operații separate de *achiziționare* și *eliberare*. Operația de *achiziționare* este furnizată de către un

procesor când el dorește să obțină acces exclusiv la un anumit obiect din memoria partajată. Pentru a preveni interferența cu alt procesor ce poate în mod curent să aibă un acces exclusiv la obiectul partajat, procesorul trebuie să aștepte ca operațiunea de achiziționare să se termine înainte de inițierea oricărei referiri la memoria partajată. Operația de *eliberare*, pe de altă parte, este utilizată pentru furnizarea unui acces exclusiv la obiectul partajat din memorie. Această fracționare a operației de sincronizare în două faze separate, permite acestui model de consistență să obțină o mai mare suprapunere a operațiilor de memorie efectuate de toate procesoarele decât în cazul modelelor de consistență secvențiale, respectiv de ordonare slabă [TH90], [GGH91], [GGHM91], [ZB92], [Marq89].

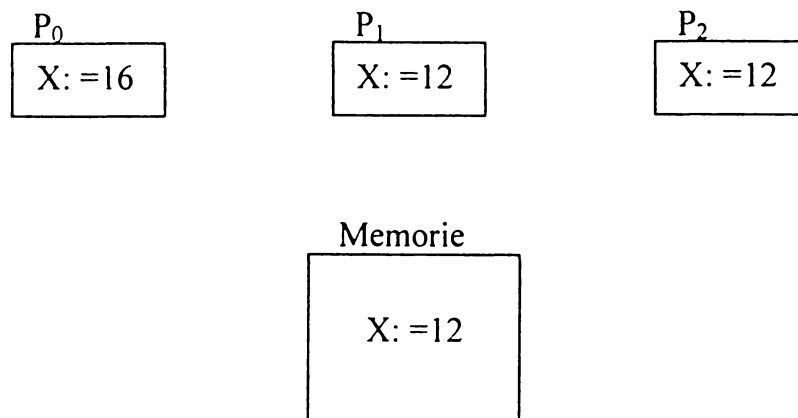
1.3.2 Coerența memoriilor cache

Referitor la coerența memoriilor cache, trebuie subliniat faptul că atunci când o locație de memorie este scrisă de către oricare procesor, faptul că valoarea din acea locație a fost schimbată trebuie să fie comunicat tuturor celorlalte procesoare care dețin o copie cache a acestei locații pentru a se asigura că nici unul dintre ele nu utilizează o versiune perimată. Coerența este menținută de către scheme de asigurare a coerenței memoriilor cache (SAC-MCH). Se consideră un sistem cu trei procesoare, fiecare cu o memorie cache privată în care se execută secvența de citiri și de scrieri prezentată în fig. 1.8 a). După ce primele două citiri au fost efectuate la t_2 , cache-urile ambelor procesoare P_0 și P_1 vor conține valoarea "12" pentru variabila înmagazinată în locația X a memoriei ca și în fig. 1.8 b). La momentul t_3 procesorul P_0 scrie în această locație schimbând valoarea în "16". Într-un sistem fără SAC-MCH, această valoare va fi actualizată doar în cache-ul lui P_0 astfel că atunci când P_1 recitește X la momentul t_4 , el va citi vechea valoare de "12" din cache-ul propriu, situație prezentată în fig. 1.8 c).

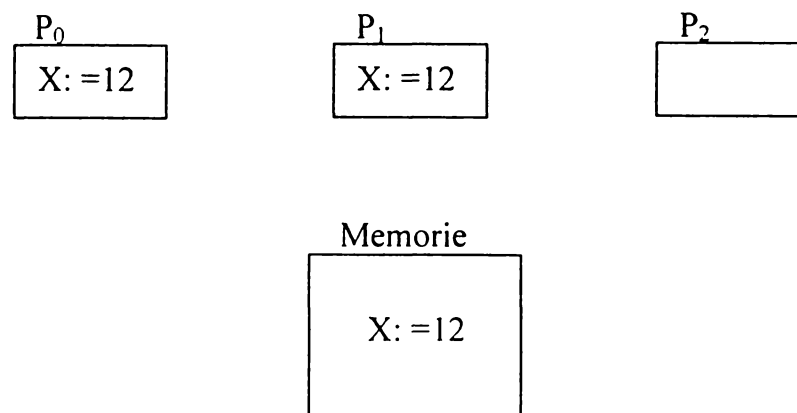
Similar, P_2 de asemenea va citi vechea valoare deoarece memoria principală nu a fost actualizată cu ultima valoare scrisă de P_0 . SAC-MCH este necesar pentru a asigura ca valorile perimate ale lui X în cache-urile altor procesoare și din memoria principală, de exemplu valoarea "12", nu vor fi propagate în viitoarele operații de citire.

Timp	P ₀	P ₁	P ₂
t_1	citește X		
t_2		citește X	
t_3	scrie 16,X		
t_4	citește X	citește X	citește X

a) Secvențe de citiri și scrieri



c) Conținuturile memoriilor cache după citiri la timpul t_4



b) Conținuturile memoriilor cache după citire la timpul t_2

Fig. 1.8 Un exemplu de problemă de coerență a memoriei cache

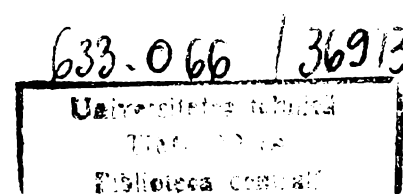
Se crează o relație între modelul de consistență de memorie și SAC-MCH. Această relație se referă în primul rând la faptul că SAC-MCH asigură că toate memoriile cache privesc toate operațiunile de scriere la un bloc specific în aceeași ordine logică. Modelul de consistență definește programatorului ordinea de scrieri în diferite blocuri, ca și când ar fi percepute de către fiecare din procesoare. Aceasta înseamnă că dacă programatorul urmează regulile modelului de consistență pentru sistemul ce urmează să fie utilizat, SAC-MCH forțează ca valoarea returnată de către oricare operațiune LOAD să fie valoarea garantată de către modelul de consistență.

1.3.3 Relația între modelele de consistență și coerența memoriilor cache

Într-un sistem ce garantează consistența secvențială, SAC-MCH asigură ca efectele fiecărei operațiuni de scriere la o locație din memoria partajată să fie propagate către toate cache-urile înainte de următoarea scriere la aceeași locație de către orice procesor care ar putea urma. Aceste accese sunt *puternic ordonate* [DSB88]. Pe de altă parte, într-un sistem bazat pe un model de consistență slab ordonat, doar accesele la variabilele de sincronizare predefinite sunt puternic ordonate. Ordonarea acceselor la locațiile memoriei partajate care sunt efectuate de diferite procesoare, pot apare în oricare ordine. Astfel, modelul de consistență cu ordonare slabă asigură că valorile de date în memoriile cache sunt coerente doar în punctele de sincronizare.

1.3.4 Scheme de asigurare a coerenței memoriilor cache

Pentru a rezolva problema coerenței memoriilor cache au fost propuse o varietate de scheme de menținere a coerenței. Soluția optimă pentru un sistem multiprocesor dat depinde de câțiva factori ca și mărimea sistemului exprimată în număr de procesoare, utilizarea anticipată a sistemului și costul dorit al acestuia. În principiu, în literatură se descriu trei tipuri principale de scheme de asigurare a coerenței memoriilor cache (SAC-MCH). SAC-MCH orientat pe *monitorizare (snooping)* se bazează pe un mediu de interconectare partajat între procesoare și modulele de memorie ca de exemplu magistrala comună, care permite fiecărui procesor să monitorizeze toate tranzacțiile la memoria partajată. Un procesor monitorizează toate referirile la memoria celorlalte procesoare și poate detecta când un bloc care este cache-at a fost schimbat de către alt procesor. Apoi



el invalidează [GVW89], [Good83] propria sa copie cache-ată de o astfel de manieră încât următoarea referire la acest bloc va forța un eșec de acces și valoarea curentă va fi obținută din memoria principală sau din altă memorie cache. În mod alternativ, procesorul poate să actualizeze direct copia sa cache-ată cu o nouă valoare disponibilă pe magistrală. Deoarece magistrala partajată difuzează efectele fiecărei operații de scriere într-o locație de memorie partajată tuturor memoriilor cache, în același ciclu în care se produce scrierea propriu-zisă, acest SAC-MCH implementează un model de consistență puternic ordonată.

În timp ce aceste mecanisme de monitorizare sunt relativ simplu de implementat, magistrala partajată poate să devină o limitare severă a performanței. Pentru a reduce situația conflictuală de pe o singură magistrală, sistemul Wisconsin Multicube [Good88] a propus o rețea n dimensională de magistrale cu procesoarele localizate în punctele de intersecție ale magistralelor și cu modulele de memorie la punctele terminale. Magistralele adiționale furnizează o lățime de bandă a memoriei mai mare în detrimentul unui protocol de menținere a coerenței mult mai complicat. O altă abordare [Arch88] grupează procesoarele în jurul unor magistrale mai mici separate și menține coerența printre procesoarele fiecărui grup. O ierarhie adițională de magistrale este introdusă pentru a menține coerența între grupuri. O altă abordare adaugă o magistrală specială de menținere a coerenței [BLA89], [Marq89] pentru a deplasa traficul de actualizare a coerenței de la citirea datelor din magistrala memoriei primare.

Deoarece toate aceste scheme utilizează magistrale adiționale pentru a crește lățimea de bandă între procesoare și memoria partajată, performanța lor – în ultimă instanță – va fi limitată de către posibilitatea apariției de conflicte de acces pe magistrală când există prea multe procesoare și de asemenea limitarea apare și datorită dificultății tehnologice de a construi aceste magistrale de mare viteză. Prin urmare, SAC-MCH bazate pe monitorizare sunt limitate la sistemele multiprocesor de talie redusă.

O altă metodă de a preveni problema saturării magistralei constă în a înlocui magistrala cu o rețea de interconectare ca de exemplu RIN *omega*, RIN *plasă* sau un *hipercub*. Aceste rețele furnizează o lățime de bandă mai mare între procesoare și memorii decât aceea a unei magistrale partajate dar în același timp măresc timpul de întârziere în accesul la memorie. Această întârziere mai mare intensifică dorința de a utiliza memorii private prin eliminarea mecanismului prin care procesoarele monitorizează tranzacțiile la

memoria partajată. Atât mecanismele hardware bazate pe *director* cât și abordările *orientate pe compiler* constituie alternative pentru menținerea coerenței în aceste sisteme.

În cadrul unui SAC-MCH bazat pe *director* [Lil93] un procesor trebuie să comunice cu un director comun ori de câte ori acțiunea procesorului poate cauza o inconsistență între cache-ul său și alte cache-uri și memorie. Directorul menține informația referitoare la procesoarele ce au copia blocurilor deoarece diferite procesoare pot avea o copie a aceluiași bloc cache-at în același moment de timp. Înainte ca un procesor să poată să scrie un bloc, el trebuie să solicite un acces exclusiv la director pentru acest bloc. Înainte ca directorul să asigure acest acces exclusiv, el trimite un mesaj către toate procesoarele însoțit de o copie cache-ată a blocului, forțând fiecare procesor să-și invalideze propria copie. După recepționarea validării transferului de către toate aceste procesoare, directorul asigură un acces exclusiv procesorului aflat în ciclul de scriere. Când un procesor încearcă să citească un bloc care este exclusiv într-un procesor diferit, el va trimite o cerere de deservire a eșecului în acces către acest director. Directorul va trimite apoi un mesaj către procesorul cu copia exclusivă, forțându-l să rescrie noua valoare în memorie. După recepționarea acestei noi valori, directorul transmite o copie a blocului către procesorul solicitant. Schemele bazate pe director diferă în cantitatea de informație ce pot să o mențină referitoare la blocurile partajate și în ce măsură sunt utilizate operații de invalidare sau de actualizare pentru a asigura coerența.

Aceste scheme bazate pe director implementează un model de consistență puternic ordonată. Un model de consistență slab ordonată poate fi implementat cu un director prin deținerea unor informații referitoare la întârzierea procesorului, aflat într-o operațiune de scriere, atunci când el accesează o variabilă de sincronizare. Deoarece modelul de consistență slabă întârzie operațiunile de scriere ale procesorului doar în cazul variabilelor de sincronizare, el prezintă o performanță mai ridicată decât a modelului puternic ordonat.

SAC-MCH orientat pe *compiler* determină în timpul compilării care blocuri cache pot să devină perimate. Instrucțiuni speciale sunt apoi inserate în codul generat în scopul de a fi executate de fiecare dintre procesoare pentru a preveni ca să utilizeze posibile date perimate. Unul dintre cele mai simple dintre aceste mecanisme [Weid86] utilizează

invalidarea nediscriminată a datelor cache-ate în scopul întăririi coerenței prin intermediul unui model de consistență slab ordonat.

1.3.5 Factori ce afectează mecanismele de menținere a coerenței memoriilor cache

Cel mai important element în alegerea unui SAC-MCH îl constituie în mod uzual performanța sa exprimată în reducerea întârzierii medii atunci când se extrag date din memorie. Un alt factor important îl reprezintă costul implementării, exprimat în cât de multă memorie este necesară pentru a înmagazina blocul cache ce deține informația de partajare și de către complexitatea părții de comandă. Diferitele SAC-MCH au diferite abordări în cost și performanță, ceea ce face dificilă evaluarea lor. Parametrii ce afectează SAC-MCH pot fi rezumați astfel:

- *strategia de detectare a coerenței* ce reprezintă strategia prin care SAC-MCH detectează un posibil acces incoerent la memorie. Acest acces poate fi detectat *dinamic* în timpul rulării sau *statistic* în timpul compilării.
- *strategia de întărire a coerenței* obținută prin actualizare sau invalidare și care este utilizată pentru a se asigura că intrările unei memorii cache ce conține informație perimată nu vor fi referite niciodată de către un procesor.
- *precizia informației de partajare a blocului* ce poate fi schimbată în raport cu costul de implementare și de performanță a SAC-MCH.
- *talia blocului cache* ce afectează într-o oarecare măsură performanța sistemului de memorie.

1.4 REȚELE DE INTERCONECTARE

Rețeaua de interconectare (RIN) reprezintă o componentă esențială a SMP. Deoarece performanța unei aplicații ce rulează pe un astfel de sistem depinde într-o mare măsură de performanța accesului la memoria situată la distanță, *latența* RIN joacă un rol extrem de important. Prin urmare, proiectarea componentelor RIN trebuie să fie studiată și evaluată de o astfel de manieră încât să furnizeze o *latență minimă* și o *mare lățime de bandă*.

1.4.1 Caracteristici

O rețea de interconectare este caracterizată de: *topologie*, *algoritm de rutare*, *strategie de comutare* și *mecanism de comandă a fluxului de date* [CSG99].

- *Topologia* reprezintă structura de interconectare a grafului rețelei; ea poate fi *regulată* sau *neregulată*. Cele mai multe SMP utilizează structuri regulate. Adeseori se face o distincție între rețelele directe și cele indirecte; RIN *directe* au câte un nod conectat la fiecare comutator în timp ce RIN *indirecte* au nodurile conectate doar la un subset specific de comutatoare care formează muchiile rețelei. Cele mai multe dintre SMP utilizează strategii mixte, așa că distincția cea mai critică este între cele două tipuri de noduri: nodurile procesoare ce generează și deplasează traficul, în timp ce comutatoarele permit doar traversarea pentru acest trafic. Exemple de rețele directe sunt multe dar cele mai cunoscute sunt *hipercubul* utilizat în SGI Origin 2000, rețeaua *plasă* 2D utilizată la sistemul Intel Paragon și rețeaua *tor* 3D utilizată la CRAY T3E. RIN indirecte sunt construite utilizând nivele de comutatoare de comunicație între nodurile sursă și cele destinație. Exemple de rețele indirecte includ numeroase variante constructive ale rețelelor de interconectare ca și Omega, Delta, Baseline, Benes, etc.
- *Algoritmul de rutare* determină care mesaje pot parcurge graful rețelei. Algoritmul de rutare restricționează setul căilor posibile la un set mai redus de căi permise.
- *Strategia de comutare* determină modul în care datele dintr-un mesaj traversează calea de parcurs. Există două strategii de comutare de bază: în *comutarea de circuit* se stabilește o cale de la sursă la destinație și ea este menținută până când mesajul este transferat prin circuit. În *comutarea de pachet* mesajul ce trebuie transmis este fragmentat în pachete mai mici ce se înlănțuie secvențial. Pachetele sunt rutate individual de la sursă la destinație, fiind rutate independent datorită informațiilor de rutare care sunt conținute. Comutarea de pachet permite o mai bună utilizare a resurselor RIN deoarece conexiunile și tamponul sunt ocupate doar pe timpul traversării de către pachet. Această strategie *store and forward* se aplică la nivele intermediare de elemente de comutare.
- În tehnica comutării prin *canale virtuale*, se înmagazinează un pachet la un nivel intermediar de comutatoare doar dacă calea corespunzătoare de ieșire este ocupată. Se utilizează tehnica de rutare *wormhole* [NK93], [CSG99], [Da192] în care mesajele sunt divizate într-un bloc ce conține informații de rutare și câteva

blocuri de date. Aceste blocuri sunt transmise prin RIN într-o manieră *pipeline*, în care blocurile de date urmează calea creată de blocul conducător.

- *Mecanismul de comandă al fluxului* determină când un mesaj, sau fragmente din el, traversează ruta stabilită. În particular, comanda fluxului este necesară ori de câte ori două sau mai multe mesaje tind să utilizeze aceeași resursă a RIN, de exemplu un canal, în același timp. Unul din fluxurile din trafic poate fi stopat, întârziat în tampoane, deturnat pe o cale alternativă sau pur și simplu îndepărtat. Fiecare dintre aceste opțiuni reclamă cerințe specifice în proiectarea comutatorului și influențează aspecte ale subsistemului de comunicare. Unitatea minimă de informație ce poate fi transferată de-a lungul unei conexiuni este denumită *flit* (flow control unit).

1.4.2 Performanță de comunicare de bază în RIN

Este foarte important să se înțeleagă relația care există între performanță și funcționalitate la nivelul subsistemului global de comunicare. În cazul RIN această discuție trebuie purtată din perspectiva latenței și a lățimii de bandă.

Latența

Timpul pentru a transfera n bytes de informație de la o sursă la o destinație, are patru componente după cum urmează:

$$\text{Timp}(n)_{S-D} = \text{Timp inițializare transfer} + \text{Întârziere de rutare} + \text{Grad de ocupare a canalului} + \text{Întârziere de conflict} \quad (1.1)$$

Timpul de inițializare a transferului reprezintă timpul consumat de procesor pentru inițializarea unui transfer prin RIN. El poate fi un cost fixat dacă procesorul pur și simplu dă comanda de inițializare a comunicării sau poate fi liniar cu n dacă procesorul trebuie să copieze datele înainte de a inițializa transferul.

Ocuparea canalului reprezintă timpul consumat pentru date în parcurgerea prin cea mai lentă componentă a căii de comunicare. De exemplu fiecare conexiune ce este traversată în RIN va fi ocupată pentru un timp n/BW unde BW este lățimea de bandă a conexiunii.

Întârzierea în rutare reprezintă timpul consumat în RIN în mod distinct pentru identificarea următorului comutator ce trebuie traversat.

Întârzierea generată de conflict reprezintă timpul consumat în rezolvarea situației conflictuale generată de solicitarea simultană a aceleiași căi de comunicare în RIN. Întârzierea în rutare și ocuparea canalului formează așa-numita *latență neîncărcată* a rețelei pentru dimensiuni de mesaje tipice [CSG99].

Latența neîncărcată pentru un pachet de n byte, incluzând antetul, într-o rutare *store and forward* este:

$$T_{sf}(n, h) = h \left(\frac{n}{b} + \Delta \right) \quad (1.2)$$

unde:

h – distanță de rutare ce reprezintă numărul de canale din calea de parcurgere

b – lățimea de bandă a liniei dintr-un canal

Δ – întârziere adițională în rutare

Această ecuație sugerează că topologia RIN este un element extrem în determinarea latenței RIN deoarece topologia determină în mod fundamental distanța de rutare h . În realitate situația este mult mai complicată depinzând de mai mulți factori printre care strategia de comutare este determinantă. De exemplu, în literatură [CS99] se citează latența neîncărcată în RIN în unități de timp de rețea, τ , pentru un mesaj de n byte ce parcurge distanța h într-o RIN comutată în circuit ca fiind:

$$T_{CS}(n, h) = \frac{n}{b} + h\Delta \quad (1.3)$$

În această ecuație se observă situația în care la creșterea lungimii mesajului, distanța de rutare și prin urmare topologia, devine o fracțiune nesemnificativă a latenței de comunicare.

În ceea ce privește latența în RIN cu comutare de pachet, se observă că un mesaj lung poate fi fragmentat în câteva pachete mici ce parcurg rețeaua într-o manieră *pipeline*. În acest caz, latența neîncărcată este:

$$T_{sf}(n, h, n_p) = \frac{n - n_p}{b} + h \left(\frac{n_p}{b} + \Delta \right) \quad (1.4)$$

unde n_p reprezintă dimensiunea fragmentelor de mesaj.

Această abordare a fost adoptată software în rețelele de comunicare tradiționale ca de exemplu Internet.

Discuția precedentă referitoare la latența de comunicare se referă la situația când un mesaj curge de la sursă la destinație fără a întâlni alt trafic în cale. În această situație RIN poate fi văzută simplu ca și o structură pipeline. Diferitele strategii de comutare și de rutare schimbă structura efectivă a pipeline-ului în timp ce topologia, lățimea de bandă a conexiunii și fragmentarea determină adâncimea și timpul la nivel de etaj al RIN.

Desigur, rațiunea pentru care RIN sunt atât de interesante din punct de vedere arhitectural este aceea că ele nu sunt simple structuri pipeline ci mai curând rețele întretesute de porțiuni de mai multe pipeline-uri. Motivația principală pentru utilizarea unei RIN în locul unei magistrale este aceea de a permite multiple transferuri de date de pot apărea simultan. Aceasta înseamnă că un mesaj care parcurge rețeaua poate să intre în coliziune cu altele și să apară un conflict pentru resurse. Comportamentul în cazul situației conflictuale depinde de câteva fațete ale proiectului rețelei – topologia, strategia de comutare și algoritmul de rutare – dar elementul esențial îl reprezintă faptul că la orice moment dat un canal poate fi ocupat doar de un mesaj. Dacă două mesaje tind să utilizeze același canal simultan, unul trebuie să fie îndepărtat. În mod tipic, fiecare comutator furnizează câteva mijloace de arbitrarie pentru canalele de ieșire. Prin urmare un comutator va selecta unul dintre pachetele ce sosesc la intrare și care concurează pentru fiecare ieșire iar celelalte vor fi înlăturate pentru o tratare ulterioară în aceeași manieră.

În rețelele de comunicare din calculatoarele paralele, un pachet este mai repede blocat în tampon decât îndepărtat; această procedură presupune existența unei tranzacții între portul de ieșire și portul de intrare de-a lungul unei conexiuni, ceea ce reprezintă în ultimă instanță un control al fluxului la nivel de conexiune.

Există o numeroasă literatură care tratează atât de importantul, însă în același timp atât de delicatul subiect al latenței, dar credem că în lucrarea de referință [CSG99] sunt sintetizate principalele aspecte legate de latență și de lățimea de bandă într-o corectă și performantă proiectare a RIN.

Lățimea de bandă

Lățimea de bandă este un parametru critic în evaluarea programelor paralele, pe de o parte datorită faptului că o lățime de bandă mai mare diminuează gradul de ocupare, iar pe de altă parte datorită faptului că o lățime de bandă mai mare reduce probabilitatea de conflict și, în ultimă instanță, datorită faptului că porțiuni din program pot să prelucreze

un volum mare de date fără să aștepte ca transmisia datelor individuale să se termine. Deoarece RIN se comportă ca și structurile pipeline, este posibil să se furnizeze lățime de bandă mai mare chiar și când latența este mare.

Este util să se privească lățimea de bandă din două puncte de vedere: lățimea de bandă "globală" disponibilă la toate nodurile din cadrul RIN și lățimea de bandă individuală "locală" disponibilă la un nod. Dacă volumul de comunicare total al unui program este de M bytes iar lățimea de bandă a comunicării din RIN este B bytes/secundă atunci este evident că timpul de comunicare este cel puțin M/B secunde. Pe de altă parte, dacă toată comunicarea se produce spre sau dinspre un singur nod, această estimare este prea optimistă, în sensul că timpul de comunicare va fi determinat de lățimea de bandă a unui singur nod.

Lățimea de bandă totală a tuturor canalelor în rețea este Cb bytes/secundă sau Cw biți/ciclu sau C unități fizice de comunicare/ciclu. Dacă fiecare din cele N noduri livrează un pachet la fiecare M cicluri cu o distanță medie de rutare h , atunci fiecare pachet ocupă în medie h canale pentru $l=n/w$ cicluri iar sarcina totală a RIN este: Nhl/M unități fizice/ciclu. Utilizarea medie a canalului este:

$$\rho = M \frac{C}{Nhl} < 1 \quad (1.5)$$

O modalitate de a privi această formulă este una în termeni de număr de conexiuni/nod, C/N , ce reflectă lățimea de comunicație disponibilă în medie pe fiecare nod. Lățimea de comunicație este consumată într-o proporție directă cu distanța de rutare și cu talia mesajului. Numărul de legături/nod este o proprietate statică a topologiei. Distanța de rutare medie este determinată de topologie, de algoritmul de rutare, de configurația de comunicație a programului și de maparea programului în sistem. O bună localitate a comunicației poate utiliza un h mic în timp ce o comunicație aleatoare va traversa distanța medie care se poate întinde pe întreg diametrul rețelei. Dimensiunea mesajului este determinată de comportamentul programului și de modelul comunicației.

În practică, câțiva factori limitează utilizarea canalului ρ ce devine mult subunitar. Sarcina poate să nu fie perfect echilibrată pe toate căile de comunicare ale RIN. Chiar dacă sarcina este echilibrată, algoritmul de rutare poate să reducă utilizarea conexiunilor pentru o configurație particulară de comunicare utilizată de program. Și chiar în situația când toate conexiunile sunt utilizabile și sarcina este echilibrată, pot apare variații

stochastice în ceea ce privește sarcina și conflictul pentru resurse de nivel inferior. Toți acești factori afectează *punctul de saturare* ce reprezintă lățimea totală a canalului ce poate fi furnizată în mod util. Așa cum se ilustrează în fig.1.9, dacă solicitarea de lățime de bandă plasată pe rețea de către procesoare – denumită *lățime de bandă oferită* – este moderată, latența rămâne coborâtă și lățimea de bandă furnizată crește cu lățimea de bandă oferită.

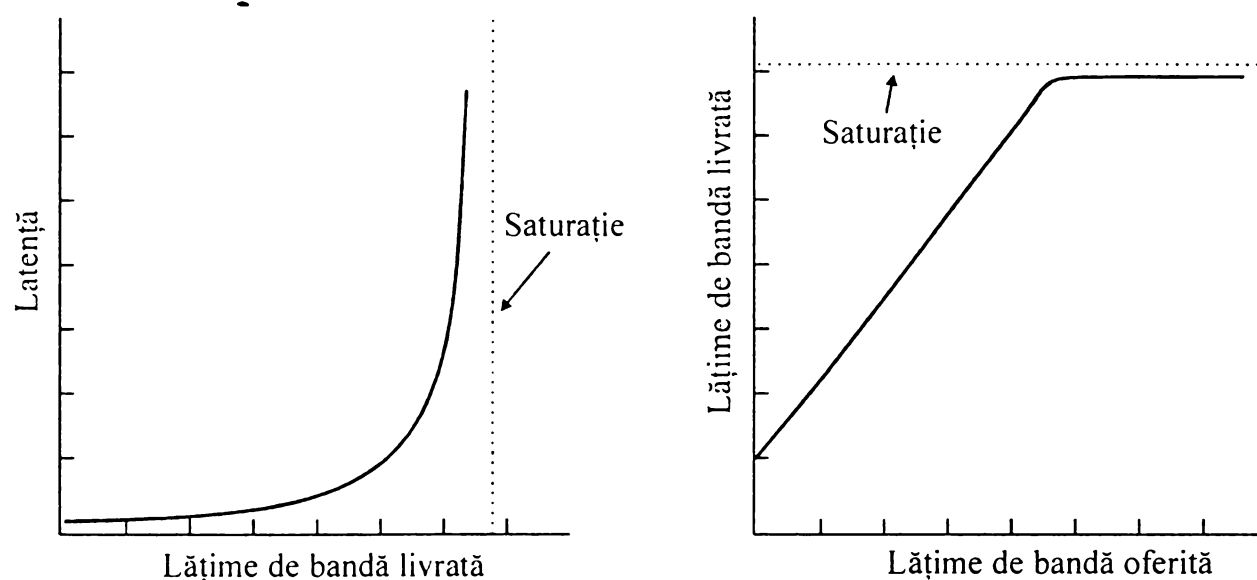


Fig. 1.9 Comportamentul la saturație pentru o rețea de interconectare tipică

La un anumit punct, solicitarea pentru o lățime de bandă superioară nu face decât să crească conflictul pentru resurse iar latența crește dramatic. Creșterea lățimii de bandă oferită nu crește lățimea de bandă furnizată. Pentru a preveni această situație de saturare a rețelei se limitează fie solicitările plasate de către procesoare, fie se asigură o lățime de bandă de comunicare amplă.

Modelarea performanței la rețelele mașinilor paralele demonstrează că latența și lățimea de bandă a rețelelor reale depind sub toate aspectele de schema rețelei. Modelarea performanței rețelelor este ea însăși o arie bogată cu o literatură voluminoasă [BYA89], [BINK97], [DB94], [KS83], [YLL90], [BRA89]. Aceste studii consideră rețeaua ca fiind izolată de context și sunt bazate pe modele analitice ce presupun modele simple ale traficului de date prin rețea. Este dificil să se formalizeze toate detaliile unei rețele într-un simplu model analitic sau într-unul de simulare. Prin urmare, aceste modele nu iau în considerare diferitele tipuri de mesaje, traficul de invalidare asociat cu coerența memoriilor cache și cu sincronizarea, acestea fiind specifice unui mediu multiprocesor. În plus, modelele prezentate evaluează rețeaua în termeni de *capacitate de trecere* și de

întârziere a mesajului, în loc de *timp de răspuns* al aplicației. În fine, rețelele de interconectare au evoluat semnificativ odată cu introducerea unor tehnici revoluționare cum ar fi *canalele virtuale* [Dal92].

Este important de adăugat că performanța nu este singurul criteriu dominant în proiectele de rețele. Alte două criterii critice sunt *costul* și *toleranța la defect* asupra cărora se va insista în capitolele următoare.

1.5 TOLERANȚA LA LATENȚĂ

Există diferențe de ritm în ceea ce privește creșterea vitezei microprocesoarelor în comparație cu timpul de acces al memoriilor utilizate. În timp ce viteza a crescut mai mult de 10% pe decadă, timpul de acces a crescut doar pe jumătate. Prin urmare, latența accesului la memorie exprimată în termeni de cicli de procesor crește cu mai mult de 20% într-o perioadă de 10 ani. În sistemele bazate pe magistrală, latența crește datorită procedurii de monitorizare a magistralei. În sistemele cu memorie distribuită, latența RIN, interfața rețelei se adaugă la timpul de accesare a memoriei locale din nod. Memoriile cache ajută la reducerea frecvenței acceselor cu latență ridicată dar ele nu reprezintă o soluție universală: ele nu reduc comunicația inerentă și programele au o rată de eșecuri semnificativă.

În proiectarea SMP, unul dintre scopurile protocoalelor de menținere a coerenței este acela de a reduce frecvența evenimentelor cu latență mare și cu solicitări de lățime de bandă impuse în mediul de comunicare în timp ce se asigură un model de programare convenabil. Scopul proiectării suportului hardware în cadrul SMP este acela de a reduce latența accesului de date în timp ce se menține o lățime de bandă scalabilă mare. În mod uzual se poate ameliora lățimea de bandă prin tehnici hardware, de exemplu utilizând conexiuni mai largi sau topologii mai generoase. Dar latența rămâne limitarea fundamentală.

În proiectarea arhitecturii unui SMP există trei căi pentru a reduce latența efectivă a accesului de date – primele două căzând în responsabilitatea sistemului iar a treia în responsabilitatea aplicației.

1. *Reducerea timpului de acces la fiecare nivel al unei ierarhii extinse de memorie.* Această tehnică pretinde acordarea unei atenții deosebite în ceea ce privește eficientizarea fiecărui pas din calea de acces. Interfața procesor-cache

poate fi proiectată să funcționeze foarte strâns. Controllerul cache trebuie să acționeze rapid în cazul unui eșec de acces pentru a reduce penalizarea în timp atunci când accesează nivelul următor. Interfața cu RIN poate fi strâns cuplată cu nodul și proiectată pentru a forma, livra și manipula rapid tranzacțiile prin rețea. Rețeaua propriu-zisă poate fi proiectată pentru a reduce întârzierile de rutare, timpul de transfer și întârzierile datorită congestiei de trafic. Printr-o proiectare atentă trebuie să se înțeleagă însă că nu pot fi depășite latențele inerente tehnologiei. Evident, tolerarea latenței prin această tehnică adaugă costuri suplimentare.

2. *Structurarea sistemului pentru a reduce frecvența acceselor cu latență ridicată.* Aceasta este sarcina de bază a replicării automate ca și aceea ce se produce în memoriile cache unde se iau în considerare avantajele localității spațiale și temporale în configurația de acces a programului, în sensul că cele mai importante date trebuie să fie cât mai aproape de procesor. Se poate face replicarea mult mai efektivă prin adaptarea structurii sistemului multiprocesor; de exemplu prin furnizarea unei cantități substanțiale de capacitate de înmagazinare fiecărui nod al sistemului.
3. *Structurarea aplicației pentru a reduce frecvența acceselor cu latență ridicată.* Aceasta presupune descompunerea și atribuirea sarcinii de calcul procesoarelor în scopul reducerii comunicării inerente și structurării configurațiilor de acces pentru a crește localitatea spațială și temporară.

În plus față de comunicare și accesul de date, există alte evenimente cu un potențial de latență ridicată ca de exemplu sincronizarea pentru care trebuie făcute eforturi similare. Aceste eforturi în ameliorarea latenței sistemului și a aplicației, deseori nu sunt suficiente. În acest sens o altă abordare [CSG99] este aceea de a *tolera latența remanentă*. Aceasta înseamnă estomparea latenței din drumul critic al procesorului suprapunând-o cu calculul sau cu alte evenimente cu latență ridicată. Procesorului i se permite să efectueze alte sarcini utile sau chiar acces de date și comunicare în timp ce evenimentul cu latență ridicată este în evoluție. Cheia tolerării latenței este în ultimă instanță *paralelismul* deoarece activitățile suprapuse trebuie să fie independente una de cealaltă.

1.5.1 Cerințe fundamentale, beneficii și limitări ale toleranței la latență

Cerințe

Tolerarea latenței necesită un *extraparalelism*, o *lățime de bandă mărită* și, în cele mai multe dintre cazuri, un *hardware mult mai sofisticat și protocoale elaborate*.

- *Extraparalelism*. Deoarece activitățile suprapuse trebuie să fie independente una față de cealaltă, paralelismul în aplicație trebuie să fie mai mare decât numărul de procesoare utilizate. Un paralelism suplimentar poate fi explicat sub forma unor fluxuri adiționale, ca și în sistemele multiflux, sau poate să existe în cadrul unui flux. Chiar comunicarea a două cuvinte din același proces în paralel implică o lipsă a unei serializări totale dintre ele, și prin urmare un extraparalelism.
- *Lățime de bandă mărită*. În timp ce toleranța la latență poate să reducă timpul de execuție, nu se reduce însă cantitatea de comunicație. Aceeași comunicare efectuată într-un timp mai mic, înseamnă o rată de comunicare mai mare și, prin urmare, se impun lățimi de bandă mai mari în arhitectura de comunicație.
- *Un hardware mai sofisticat și protocoale mai elaborate*. Cu excepția cazului producerii unor mesaje mai mari, procesorul trebuie să i se permită să efectueze o operație cu latență mai mare înainte ca operația să se termine.

Aceste cerințe implică ca toate tehnicile de tolerare a latenței să aibă un cost relativ semnificativ. Prin urmare, proiectantul trebuie să utilizeze tehnici algoritmice pentru a reduce frecvența evenimentelor cu latență ridicată înainte de a se baza pe tehnici de tolerare a latenței. Cu cât mai puține sunt evenimentele cu latență ridicată, cu atât mai puțin agresivă va fi tendința de a estompa această latență.

Beneficii potențiale

O simplă analiză poate oferi limitele în beneficiile de performanță ce pot fi așteptate de la tolerarea latenței. Preuspunându-se că nu se utilizează tehnici de tolerare a latenței, timpul de execuție, așa cum este el văzut de un procesor, are următorul profil: T_C cicli sunt consumați în calcul local, T_{OV} cicli sunt consumați în procesarea mesajului, T_{OCC} cicli în asistarea comunicării și T_l cicli în așteptarea transmisiei mesajului în RIN. Dacă se presupune că celelalte resurse pot fi perfect suprapuse cu activitatea din procesorul principal atunci creșterea potențială de viteză poate fi determinată prin aplicarea simplă a legii lui Amdahl. Procesorul poate fi ocupat pentru $T_C + T_{OV}$ cicli; latența maximă ce

poate fi estompată de către procesor este $T_I + T_{OCC}$ prin urmare creșterea maximă de viteză datorită estompării latenței este:

$$\left(1 + \frac{T_{OCC} + T_I}{T_C + T_{Ov}}\right) \quad (1.6)$$

Această limită reprezintă un prag superior deoarece se presupune o perfectă suprapunere a resurselor și nu este impus nici un cost suplimentar datorită toleranței la latență.

Cât de multă latență poate fi efectiv estompată depinde de mulți factori ce implică aplicația și arhitectura. Caracteristici de aplicație relevante includ structura comunicației și cât de multe alte activități pot fi suprapuse cu ea.

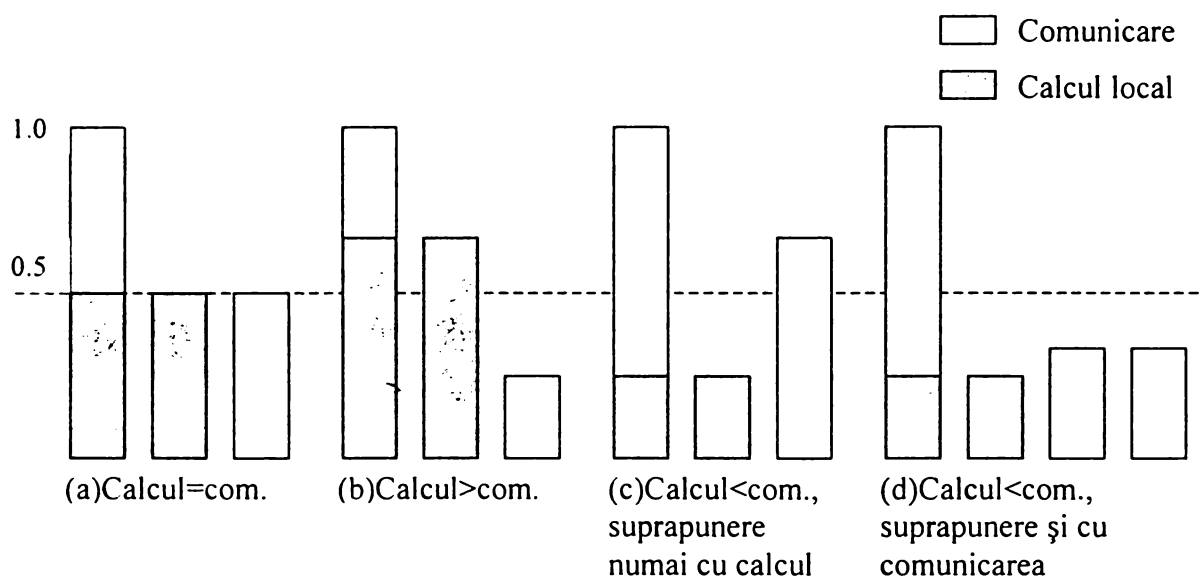


Fig.1.10 Limitări și beneficii ale toleranței la latență

Limitări

Limitările majore pot fi divizate în trei clase: *limitări de aplicație*, *limitări ale arhitecturii de comunicație* și *limitări de procesor*.

Limitări de aplicație. Cantitatea de timp de calcul independent ce este disponibilă pentru suprapunere cu latența poate fi limitată, astfel că nu toată latența poate fi estompată iar procesorul va trebui să stagneze pentru o anumită perioadă de timp. Chiar dacă programul are suficientă sarcină utilă și extraparalelism, structura programului poate să facă dificil pentru sistem sau pentru programator identificarea operațiunilor concurente și să orchestreze suprapunerea.

Limitări ale arhitecturii de comunicație. Arhitectura de comunicație poate să restricționeze numărul de mesaje sau numărul de cuvinte ce pot fi extrase dintr-un nod la

un moment dat și, prin urmare, parametrii de performanță ai arhitecturii de comunicație pot limita latența ce trebuie să fie estompată.

În literatură [CSG99] se tratează pe larg subiecte legate de limitarea latenței ce poate fi mascată datorită specificităților arhitecturale.

Limitări ale procesorului. În SMP-MP ce utilizează mecanisme de asigurare a coerenței memoriilor cache, localitatea spațială asigurată prin blocuri cache lungi permite extragerea la un moment dat a mai multor cuvinte de comunicație. Pentru a estompa latența, un procesor și subsistemul său cache trebuie să permită apariția mai multor eșecuri în acces la memoria cache în mod simultan. Acest lucru este însă costisitor așa că sistemele procesor-cache au limitări relativ reduse în ceea ce privește numărul de eșecuri extrase.

Este clar din cele prezentate că transformarea unei arhitecturi de comunicație într-una eficientă reprezintă o operațiune foarte importantă pentru tolerarea efectivă a latenței.

În continuare menționăm că o tratare sistematică a relației complexe dintre arhitectură, toleranță la latență și performanță excede scopul acestei dizertații, ea putând însă genera subiecte majore de reflexie pentru proiectanții de sisteme de calcul paralele.

În concluzie, trebuie să remarcăm că în situația creșterii discrepanței dintre vitezele procesoarelor și timpii de acces la memorie și timpii de comunicație, toleranța latenței va deveni un subiect critic pentru viitoarele multiprocesoare. Au fost dezvoltate multe tehnici de tolerare a latenței și fiecare își are avantaje și dezavantaje relative. Toate aceste tehnici se bazează pe concurența în exces în programe de aplicație în cadrul numărului de procesoare utilizate și toate tind să crească cererile de lățime de bandă plasată pe arhitectura de comunicație.

Pentru multiprocesoarele coerente din punct de vedere cache, tehnicile de tolerare a latenței sunt suportate hardware atât de procesor cât și de sistemul memoriei cache și conduc la un spațiu generos de alternative de proiectare. Multe dintre aceste tehnici de toleranță a latenței suportate de către hardware sunt de asemenea aplicabile la uniprocesoare. Tehnici ca: *planificarea dinamică*, *modele de consistență a memoriei relaxate*, *preaducerea datelor din memorie*, etc. sunt frecvent întâlnite în arhitecturile microprocesoarelor contemporane. Cea mai generală tehnică de mascare a latenței, cea care presupune crearea de fluxuri multiple, nu este încă populară din punct de vedere comercial deoarece ea nu a fost încă testată pe uniprocesoare.

În pofida acestui spațiu generos de alternative de proiectare multe dintre problemele cotidiene ale toleranței latenței sunt de asemenea probleme susținute prin tehnici software. Automatizarea și simplificarea suportului software cerut pentru asigurarea toleranței la latență este o sarcină departe de a fi fost rezolvată definitiv. În ultimă instanță, raportul dintre tehnicile de tolerare a latenței suportate prin resurse hardware și cele suportate software reprezintă un subiect deschis în arhitecturile paralele.

1.6 TOLERANȚA LA DEFECT

Indiferent de tendințele arhitecturale abordate în proiectarea și construcția SMP, *toleranța la defect* va reprezenta unul dintre principalele obiective urmărite de proiectanți. Noile tehnologii, indiferent cât de performante sunt, nu asigură o siguranță în funcționare absolută. Prin urmare, luarea în considerare a toleranței la defect în stabilirea indicelui de performanță reprezintă o necesitate stringentă ce nu poate fi ignorată în nici un fel.

Fiabilitatea și disponibilitatea sistemelor de calcul constituie aspecte majore în proiectarea lor. Relativ recenta descoperire a unei erori de proiectare în microprocesorul Pentium a revelat complexitatea măsurilor de asigurare a fiabilității în utilizarea microprocesoarelor avansate.

Oportunitatea pentru o recuperare rapidă în cazul apariției unei defecțiuni este consistent ajustată de apariția microprocesoarelor de mare viteză, dar apar noi provocări în ceea ce privește asigurarea unei sincronizări fabile. O altă provocare pentru proiectarea SMP tolerante la defect o reprezintă creșterea utilizării sistemelor masiv paralele. În dorința de a obține cea mai mare performanță posibilă cu ajutorul arhitecturilor inovative, aceste sisteme sunt adeseori bazate pe tehnologii neprobate. Nu este surprinzător faptul că diferențele în dependabilitate, prezentate în fig. 1.11 sunt atât de mari. Un cost de necontestat al lipsei de toleranță la defect este pierderea puterii de calcul în calculul de mare performanță. O mică fracțiune din timpul operațional pierdut datorită defectelor ar putea genera rapid o pierdere în performanță care ar transforma SMP într-un singur calculator.

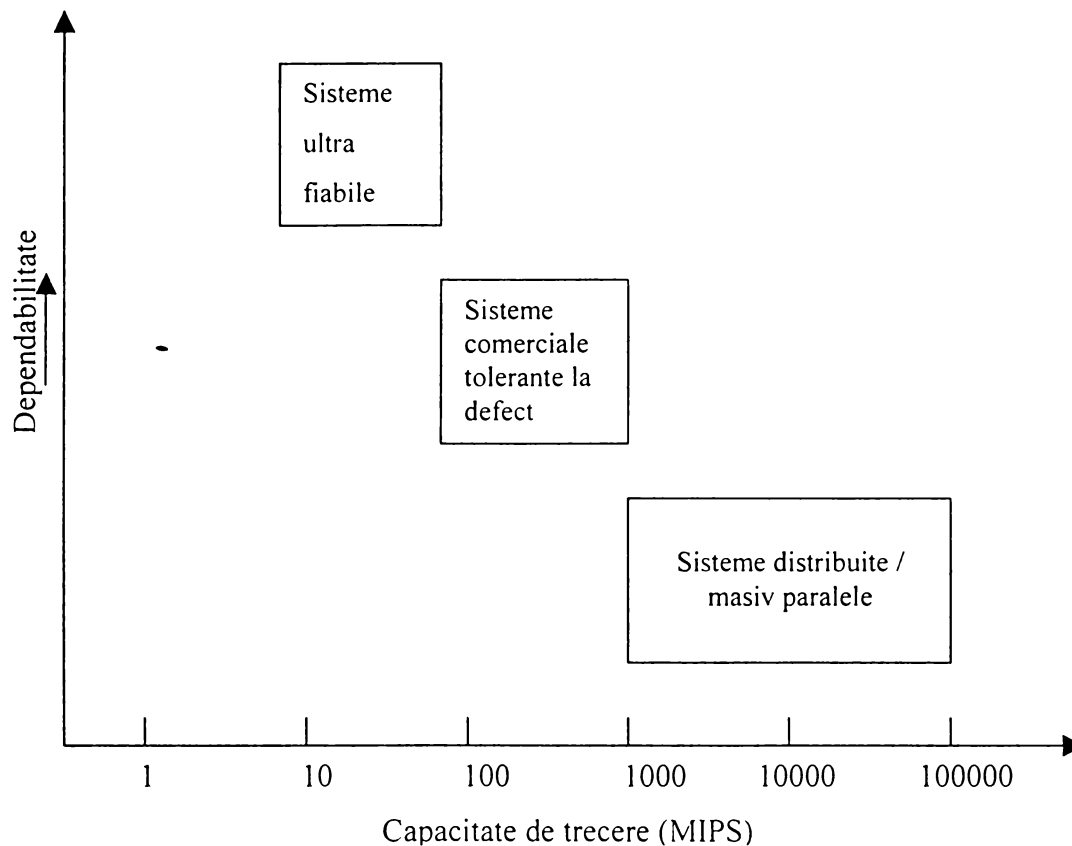


Fig. 1.11 Oferta de SMP în funcție de dependabilitate/performanță

Prin urmare, trebuie asigurată toleranța la defect fără o pierdere semnificativă în performanță.

Toleranța la defect în SMP complexe trebuie să încorporeze atât tehnici hardware cât și software. Cerințele de dependabilitate sunt în principal ierarhizate astfel, în funcție de nivelul căruia i se adresează:

- La nivel de *circuit* și de *tehnologie*, ținând cont de tehnicile de proiectare a circuitelor pentru a ameliora testabilitatea și fiabilitatea modulelor individuale.
- La nivel de *arhitectură a nodului*, ținând cont de proiectarea capsulelor VLSI ce implementează un nod de procesare.
- La nivel de *arhitectură între noduri*, specificându-se modul în care aceste noduri sunt conectate împreună și maniera în care sistemul va fi reconfigurat în prezența procesoarelor defecte. Se ia în considerare topologia de interconectare.
- La nivel de *sistem de operare*, ținând cont de crearea și exploatarea *punctelor de verificare – PV*, de *rulare înapoi* după ce o porțiune din sistem a fost identificată ca fiind defectă. Sarcina procesorului defect se distribuie altor procesoare.
- La nivel de *aplicație* unde utilizatorul verifică rezultatele calculului.

În mod tipic, toleranța la defect este căutată la fiecare nivel al proiectului. Defectele nedetectate la nivelele inferioare trebuie să fie manipulate la nivele superioare.

1.6.1 Implicații ale toleranței la defect în arhitecturile MIMD

În SMP-MP procesoarele sunt *anonime* în sensul că, *in extremis*, toate datele și instrucțiunile staționează într-o memorie comună disponibilă fiecărui procesor. SMP-MP dispune de memorii cache pentru a micșora latența și a crește lățimea de bandă. Cu cât talia memoriilor cache crește rolul memoriei centrale se reduce iar SMP-MP tinde să se comporte ca și un SMP-TM. Memoria cache reduce astfel anonimatul procesorului. Consistența memoriilor cache ridică probleme speciale în ceea ce privește toleranța la defect deoarece în timpul defectărilor, este afectată coerența cache între diferite memorii cache dispuse pe procesoare diferite. În Capitolul 4, în cadrul secțiunii ce tratează recuperarea din defect la SMP-MP, va fi discutat acest subiect în detaliu. O altă problemă importantă ce apare la SMP-MP, o reprezintă *propagarea defectului*. Un procesor defect poate potențial să scrie în orice locație a memoriei centrale și prin urmare să afecteze conținutul memoriei, mai puțin în situația în care se efectuează o verificare înainte de fiecare scriere în memoria principală.

Datorită faptului că SMP-TM comunică prin mesaje, sunt necesare în primul rând verificări ale corectitudinii transmiterii mesajelor. Acestea sunt mai puțin frecvente decât cele efectuate asupra execuției instrucțiilor. Astfel, din punct de vedere al protecției la erori, SMP-TM se comportă mai bine.

Abordările de bază specifice în ceea ce privește toleranța la defect în SMP sunt două. În *redundanța statică* în care N copii ale fiecărui procesor sunt utilizate iar gradul minim de replicare este triplicarea.

În *redundanța dinamică*, prima dată este detectat procesorul defect care ulterior este înlocuit cu unul de rezervă.

2. ANALIZA ARHITECTURILOR SISTEMELOR DE CALCUL DIN CLASA MIMD

2.1 ASPECTE GENERALE

Paralelismul este un termen general utilizat pentru caracterizarea unei varietăți de simultaneități ce apar în calculatoarele moderne. Avantajul major oferit de paralelism îl reprezintă viteza crescută a sistemelor de calcul bazate pe el.

Procesarea paralelă include atât *algoritmi paraleli* cât și *arhitecturi paralele*.

Un algoritm paralel poate fi văzut ca și o colecție de module independente dintre care unele pot fi executate simultan. Modulele comunică între ele în timpul executării unui algoritm. Deoarece algoritmiile evoluează într-un mediu paralel din punct de vedere hardware, care constă din module de procesare ce schimbă între ele date, este o cale naturală abordarea relației dintre *spațiul algoritmului* și *spațiul hardware*. Cu cât această relație este mai bună și mai nuanțată, cu atât mai bine se pot proiecta sisteme de calcul paralel.

Criteriile principale după care se pot caracteriza algoritmiile sunt :

- *Granularitatea modulului (module granularity);*
- *Controlul concurenței (concurrency control);*
- *Mecanismul de utilizare a datelor (data mechanism);*
- *Geometria de comunicare (communication geometry);*
- *Dimensiunea algoritmului (algorithm size).*

Arhitecturile paralele pot fi caracterizate utilizând următoarele criterii :

- *Complexitatea procesorului;*
- *Modul de operare;*
- *Structura memoriei;*
- *Rețeaua de interconectare;*
- *Numărul de procesoare și talia memoriei.*

În literatură [Mold 93] se explicitează relația dintre caracteristicile din cele două spații.

În tabelul din fig. 2.1 se prezintă succint aceste relații :

<i>Algoritmi paraleli</i>	<i>Arhitecturi paralele</i>
Granularitatea modulului	Complexitatea procesorului
Controlul concurenței	Modul de operare
Mecanismul de date	Structura memoriei
Geometria de comunicare	Rețeaua de interconectare
Talia algoritmului	Numărul de procesoare și talia memoriei

Fig. 2.1 Relația dintre caracteristicile algoritmilor paraleli și arhitecturile paralele

În literatură [Dasg 89], [Hwan 93], [Prad 96] etc se utilizează mai multe taxonomii pentru clasificarea sistemelor de calcul.

Ne interesează evoluția taxonului MIMD – Multiple Instruction Multiple Data – definit prima dată de Flynn (1966) și apoi nuanțat de Hwang-Briggs (1984). Aceștia modifică clasificarea lui Flynn introducând subtaxonii:

- MIMD–T (Multiple Instruction Multiple Data – Tightly) – Sistem strâns cuplat;
- MIMD–L (Multiple Instruction Multiple Data – Loosely) – Sistem slab cuplat.

Această partajare s-a făcut în funcție de modul de utilizare a memoriei de către procesoare, respectiv de gradul ei de utilizare în procesul de comunicare.

Organizarea sistemelor de calcul din subclasa MIMD – T este prezentată în fig. 2.2, iar a celor din subclasa MIMD – L în fig. 2.3.

Se observă că în primul caz mediul de conectare – *magistrala* – este între procesoare și memorie, în timp ce în al doilea caz memoria este atașată fiecărui procesor și mediul de conectare este *rețeaua de interconectare* ce face conectarea modulelor combinate.

Există serioase dezbateri asupra modului de organizare a memoriilor în sistemele paralele de talie mare. Dezbateră se centrează, și nu întotdeauna justificat, pe o falsă dihotomie: *memoria partajată versus memoria distribuită*. Memoria partajată înseamnă un singur spațiu de adresă, presupunând implicit o comunicare prin instrucții LOAD/STORE. În opoziție, în raport cu o singură adresă se află spații multiple de adrese private, ceea ce implică comunicare prin SEND/RECEIVE. Memoria distribuită se referă la localizarea fizică a memoriei. Dacă memoria fizică este divizată în module, cu câteva plasate lângă fiecare procesor, atunci memoria fizică este distribuită.

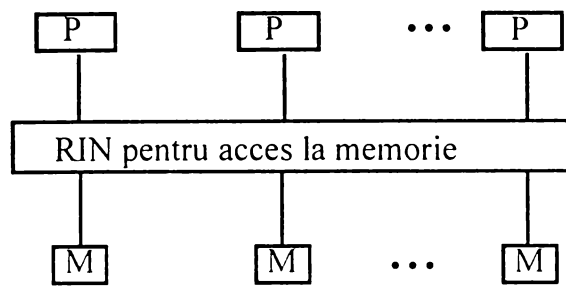


Fig 2.2 SMP cu memorie partajată, din subtaxonul MIMD-T

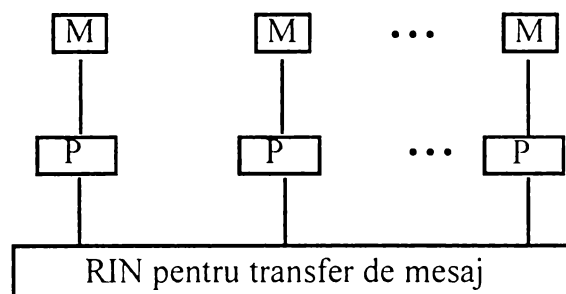
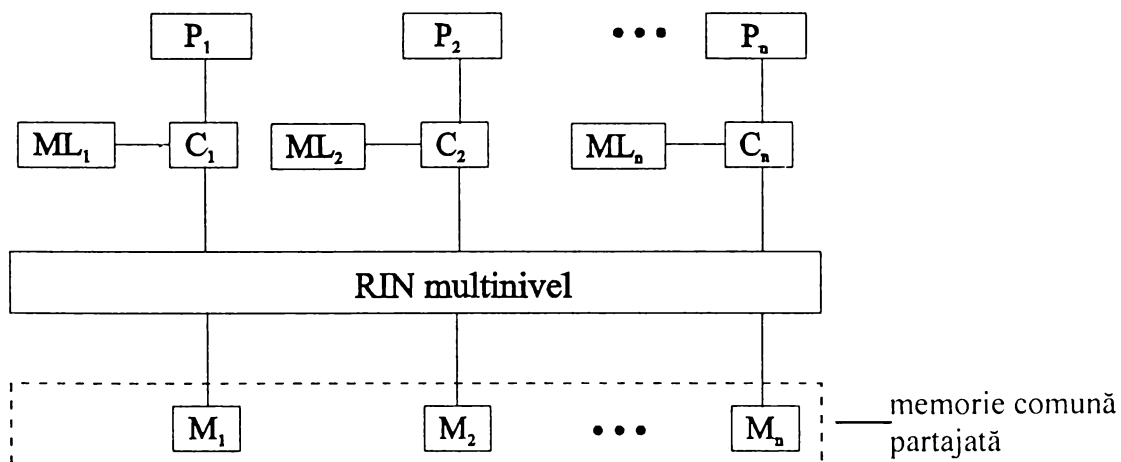


Fig. 2.3 SMP cu memorie distribuită, din subtaxonul MIMD-L



P - procesor
 M - modul de memorie
 C - memorie cache privată
 ML - memorie locală

Fig. 2.4 SMP-MP cu arhitectură tip "sală de dans"

Opusul memoriei distribuite îl constituie *memoria centralizată*, unde timpul de acces la o locație de memorie fizică este același pentru toate procesoarele deoarece fiecare acces se face prin rețeaua de conectare ca și în fig. 2.4. Această arhitectură este –după unii autori– denumită *sală de dans (dance hall)* cu procesoarele dispuse pe o latură iar memoriile pe alta. Analogia este mai mult decât evidentă.

Spațiul singular de adresă versus spații multiple de adresă și memorie distribuită versus memorie centralizată sunt subiecte ortogonale: *Mașinile MIMD pot avea un singur spațiu de adresă și o memorie distribuită sau spații cu adresă privată multiple și o memorie fizică centralizată.*

Pentru simplificare însă, vom considera uzual sistemele din subclasa MIMD – T ca având o memorie globală partajată prin intermediul căreia procesoarele își asigură *comunicarea și coordonarea.*

Le vom denumi *sisteme multiprocesor cu memorie centrală partajată (SMP-MCP)* sau pe scurt *multiprocesoare*. Un multiprocesor se consideră a fi *fiabil pentru multiprocesare (multiprocessing reliable)* dacă și numai dacă poate fi utilizat atât pentru operații MIMD cât și pentru M/SISD (Multiple/Single Instruction Single Data) conform [Hwan 93].

Reamintim că atât sistemele M/SISD cât și cele SISD sunt cazuri degenerate ale celor MIMD.

Exemple de calculatoare din subclasa MIMD-T sunt: C.mmp construit la Carnegie Mellon University, Multimax de la Encore Computer, FX de la Alliant, etc.

În sistemele din subclasa MIMD-L, memoria principală este partiționată și atașată procesoarelor. Acestea împart același spațiu adresă de memorie. Vor fi denumite *sisteme multiprocesor cu memorie centrală distribuită (SMP-MCD).*

Exemple de arhitecturi MIMD-L sunt Cm* de la CMU, Butterfly de la BBN Laboratories, RP3 de la IBM ș.a. Atât calculatoarele MIMD-T cât și cele MIMD-L vor fi denumite *multiprocesoare, sisteme multiprocesor cu memorie partajată* și abreviate SMP-MP.

Multicalculatoarele se caracterizează prin faptul că memoria nu este partajată. Interacțiunea dintre procesoare se produce prin mesaje ce sunt transferate între procesoarele sursă și destinație (noduri). Aceste mesaje traversează legătura ce conectează direct două noduri și pot să traverseze prin câteva noduri într-o manieră

înmagazinează și înaintează (*store and forward*) înainte de a-și atinge destinația. Multicalculatoare – *sisteme multiprocesor ce transferă mesaje (SMP-TM)* se bazează uzual pe topologii ca și *inel, arbore, stea, hipercub*, etc.

Exemple de multicalculatoare sunt: Caltech Cosmic Cube, Intel IPCS1, AMETEK System 1, Intel Tuchstone ș.a.

Prin urmare o clășificare a sistemelor paralel distribuite este [BYA 89] prezentată în fig. 2.5.

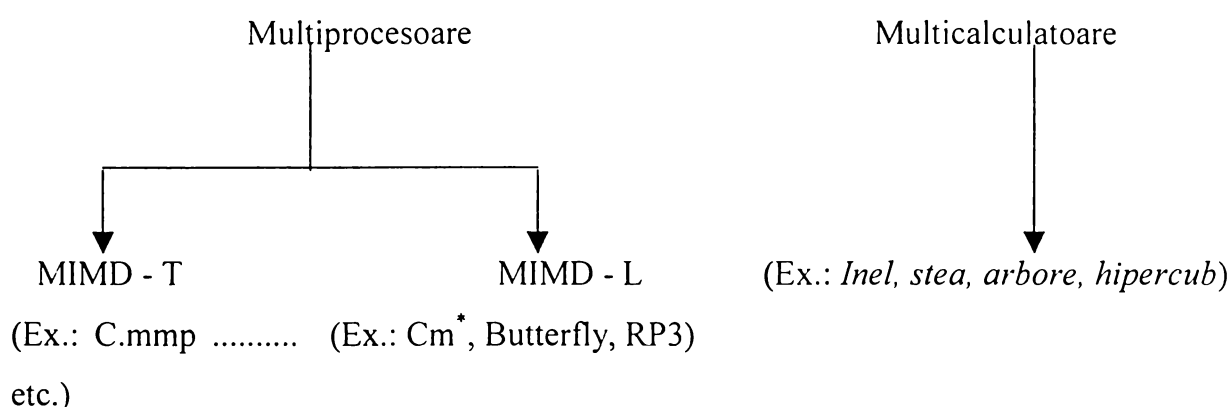


Fig. 2.5 O clășificare a sistemelor paralel/distribuite.

Clasificarea nu include *ariile de calculatoare (array computers)* și cele cu *prelucrare pe loturi (pipeline computers)* cât și *rețelele locale de calculatoare (Local Area Networks)*. Ariile de calculatoare și cele cu prelucrare pe loturi sunt cu prelucrare paralelă dar nu distribuită în timp ce LAN fac parte din prelucrarea distribuită dar nu și din cea paralelă. Proiectarea oricărei arhitecturi multiprocesor implică rezolvarea unui set fundamental de probleme, unele dintre ele fiind noi pentru contextul multiprocesării în timp de altele sunt moștenite din domeniul uniprocessor.

Cele mai importante dintre aceste probleme sunt:

1. *Granularitatea paralelismului* ce trebuie suportat de multiprocesor.
2. *Mecanismele de sincronizare* pentru accesul la memoria partajată și comunicarea între procese.
3. *Latența memoriei* și rezolvarea sa.
4. Problema *rețelei de interconectare*.
5. *Strategiile de planificare* a procesorului/procesoarelor.

6. *Limbajele de programare concurente și tehnicile de programare.*

Problemele de la 1 la 4 sunt specific arhitecturale. Arhitecturile sistemelor paralel / distribuite pot fi analizate și sub alte aspecte cum ar fi:

- *expandabilitatea* – măsură a abilității sistemului de a prelua sarcini mărite;
- *capacitatea* – proprietatea rețelei de a putea mări cantitatea de date transferate atunci când se adaugă noi procesoare la sistem;
- *flexibilitatea* – posibilitatea de adaptare a sistemului la o modificare a fluxului de date.

Opinăm însă că o abordare modernă a sistemelor paralel/distribuite trebuie să se facă prin prisma *dependabilității*, definită simplu în [Prad96] ca fiind calitatea serviciilor oferite de un sistem de calcul.

Sistemele paralel/distribuite sunt adecvate pentru implementarea *toleranței la defecte (fault-tolerance)* datorită resurselor hardware multiple din componență. Comportarea *rezistentă la defecte (fault resilient)* decurge din conectivitatea nodurilor și a muchiilor grafului asociat sistemului. În multe situații toleranța este posibilă doar cu reducerea performanțelor până la o limită admisibilă. Are loc o așa numită *degradare grațioasă (graceful degradation)*.

Intr-un sistem distribuit, diagnoza permite localizarea defectelor astfel încât sistemul să poată fi reconfigurat în jurul unităților operaționale.

Degradarea grațioasă se obține prin *configurarea automată și capacitatea de recuperare* furnizate cu ajutorul unei mentenanțe on-line.

2.2 GRANULAȚIA DE PARALELISM

Paralelismul în calcule poate fi examinat la câteva nivele depinzând de complexitatea dorită. Acestea sunt: *program (job)*, *sarcină (task)*, *proces (process)*, *variabilă (variable)* și *bit*. Criteriul cel mai uzual de clasificare a sistemelor de calcul paralel este *granularitatea (granularity)* interacțiunilor dintre activitățile care sunt executate în paralel.

Considerând un sistem cu N procesoare în care o parte a timpului de lucru al fiecărui procesor este afectată schimburilor de date și a informațiilor de control cu alte procesoare se poate scrie:

$$T^c = \sum_{i=1}^N t_i^c \quad \text{și} \quad T^w = \sum_{i=1}^N t_i^w$$

unde:

t_i^c - timpul consumat de P_i pentru comunicare;

t_i^w - timpul consumat de P_i pentru calcul;

T^c - timpul consumat de sistem pentru comunicare;

T^w - timpul consumat de sistem pentru calcul.

Considerând că nu există o lege de repartizare a timpului de prelucrare pentru fiecare procesor, se poate scrie *timpul total de prelucrare*:

$$T = T^c + T^w. \quad (2.1)$$

Granulația denumită și *granulație de prelucrare (grain of computation)* este:

$$\gamma = \frac{T^w}{T^c}. \quad (2.2)$$

Granularitatea exprimă caracterul reunit al unui program paralel, a arhitecturii pe care se execută și a dimensiunii problemei de rezolvat. Datorită acestui caracter reunit, dacă se modifică un parametru – de exemplu problema de rezolvat – se modifică și granularitatea pentru o anumită perioadă de timp.

Din această cauză, pe lângă *granularitatea potențială (potential granularity)* se introduce și *granularitatea actuală (actual granularity)* ca și o medie a timpului afectat pentru două comunicări succesive:

$$\gamma' = \frac{T^w}{S} \quad (2.3)$$

unde:

$$S = \sum_{i=1}^N S_i$$

reprezintă numărul total de comutări de la calcul la comunicare pentru toate procesoarele P_i .

În general procesarea pe un sistem de calcul paralel se situează între două extreme:

- durate scurte de calcul + durate scurte de comunicare;
- durate lungi de calcul + durate lungi de comunicare.

În primul caz se definește o *granularitate fină (fine granularity)* iar în al doilea caz o *granularitate grosolană (coarse granularity)*. Deși granularitatea este un parametru complex, unele clase de sisteme favorizează un anumit timp de granulație. Cu cât mai fin este gradul granulației cu atât crește costul sistemului de comunicare și deci costurile de comunicație.

Arhitectura din subclasa MIMD – T (*tightly coupled*) se caracterizează prin granulație fină ($\gamma < 10$) iar cele din clasa MIMD – L (*loosely coupled*) prin granulație grosieră ($\gamma > 10$). Evident granularitatea poate varia de la sistem la sistem.

Prin urmare multiprocesoarele din subclasa MIMD – T dispun de un *control centralizat* în timp ce multicalculatoarele din subclasa MIMD – L funcționează pe principiul *autonomiei cooperative*.

2.3 LATENȚA MEMORIEI ȘI REZOLVAREA EI.

Prin *latență (latency)* se înțelege timpul consumat între o solicitare făcută de un procesor la memorie și recepționarea de către procesor a datelor asociate solicitării.

Latența memoriei poate reprezenta un potențial de strangulare a traficului și trebuie tratată în cel mai serios mod.

În cazul SMP-MP problema latenței devine acută datorită:

- prezenței mai multor procesoare ce pot solicita accesul la același modul de memorie, ceea ce poate conduce la *conflict de memorie (memory contention)*;
- prezenței RIN care introduce o întârziere inerentă între un procesor ce solicită un acces la memorie și primirea datei;
- potențialul conflictual în orice componentă particulară a RIN propriu-zise.

Cu cât mai mare va fi talia SMP-MP – în termeni de număr de procesoare – cu atât mai mare va fi latența medie.

În acest paragraf se vor prezenta câteva tehnici de *evitare* sau de *tolerare* a efectelor latenței memoriei.

2.3.1 Comutarea de proces

O metodă posibilă de mascare a efectului latenței este de a efectua un proces sau un comutator “*context*” (*context switch*) ori de câte ori un proces P_i solicită un acces la memorie: P_i este suspendat și procesorul care-l execută pe P_i este reafectat la unul din procesele P_j aflat în starea de *așteptare pentru reluare* (*waiting for run*).

Abordarea acestei metode într-un context multiprocesor ridică următoarea problemă. Starea procesului suspendat P_i , conținută în registrele procesorului, va trebui ea însăși să fie salvată și starea procesului P_j va trebui să fie încărcată în registrele acestui procesor în scopul efectuării unui comutator de context.

Dacă starea procesului este salvată în memorie, atunci însăși această procedură va implica acces la memorie ceea ce va agrava problema latenței. Se poate evita această situație dacă se prevăd seturi multiple de registre de procesor care să mențină stări multiple de proces. Evident crește complexitatea sistemului.

Pentru ca această metodă să fie rentabilă este necesar ca tot timpul să existe un număr suficient de procese concurente disponibile: *Ori de câte ori un proces efectuează o solicitare de acces la memorie trebuie să existe un set nevid de procese ce așteaptă să fie executate.*

2.3.2 Utilizarea arhitecturilor de tip LOAD/STORE

În literatură [Dasg 89] se citează cu titlu de exemplu o metodă de reducere a latenței utilizată la calculatorul CRAY-1, constând din utilizarea instrucțiilor de tip LOAD/STORE: *Singurele instrucții ce fac referire la memorie sunt instrucții de tip LOAD/STORE; toate celelalte instrucții ce se referă la date utilizează registrele procesorului.*

Această observație se aplică și mașinilor RISC.

De fapt, utilizarea unui set larg de registre programabile (sau a unui sistem local de înmagazinare a informației) sugerează noțiunea de memorie cache, singura excepție fiind aceea că registrele sunt programabile. Rezumând, problema se pune ca fiecare procesor să aibă memoria sa locală și orice modificare de date, copie a unei locații din memoria principală, să fie actualizată în orice altă memorie locală unde ea se află.

2.3.3 Coerența memoriilor cache (CCH)

În SMP-MP o metodă de a reduce *latența* memoriei partajate și - *în extenso* - a întregului sistem constă în utilizarea memoriilor cache - C - asociate fiecărui procesor - amplasate într-un sistem ierarhizat.

Copii multiple ale aceluiași bloc de memorie din MP pot fi rezidente în mai multe memorii cache private. Modificarea de către un procesor a oricărei copii a acestui bloc partajat în memoria sa cache privată, va schimba aceste date partajate într-o valoare absolută în memoria principală sau în oricare alta C.

Un sistem de memorie este *coerent* dacă valoarea returnată în urma execuției unei instrucții LOAD este întotdeauna identică cu valoarea scrisă în aceeași adresă de către ultima instrucție STORE. [CF78].

Sursele inconsistențelor datelor pot fi mai multe.

În fig. 2.6 se prezintă schematic structura unui sistem multicache, având n memorii cache C_i pentru $i=1, \dots, n$, și o memorie principală partajată de toate procesoarele P_i .

Fie:

X - adresă fizică a unui bloc din memoria principală

y_i - adresa acestui bloc în C_i

Problema coerenței apare atunci când:

1. Un bloc X este memorat în C_i la adresa y_i și în C_j la y_j . O inconsistență între y_i și y_j poate apare deoarece:
 - Un proces P ce migrează de la procesorul P_i spre P_j , modifică blocul y_i într-un fel și blocul y_j în alt fel;
 - Un proces ce rulează în P_i modifică y_i , și alt proces din P_j modifică diferit y_j .
2. Un bloc X nu este apt să se sincronizeze în dinamica modificărilor cu memoriile C. De exemplu, o copie a lui X este în y_i , dar când y_i este modificat, X nu este actualizat. Dacă blocul X este solicitat în P_j , atunci ca și un rezultat al erorii, X este copiat în C_j , dar acesta nu este ultimul X .

Diferite variante de configurare a SMP-MP cu memorii cache sunt prezentate în fig. 2.6.a, 2.6.b și 2.6.c.

Cel mai popular protocol pentru a menține CCH, se numește *snooping*. În fig. 2.7 se prezintă un SMP-MP ce utilizează o singură magistrală și în care se evidențiază modul de acces al memoriilor cache la memoria principală.

Toate dispozitivele de comandă ale memoriilor cache (DCC) monitorizează magistralele pentru a determina dacă în C există copii ale blocului protejat.

Acestă monitorizare – *snooping* - este utilizată pentru SMP-MP cu o singură magistrală, arhitecturi tipice pentru anii '80. Mașinile din acea perioadă au evoluat de la unipresoare prin adăugarea de procesoare multiple la o singură magistrală. Apoi s-au adăugat memoriile cache pentru ameliorarea performanțelor și apoi s-au proiectat DCC ce să monitorizeze informația pe magistrala partajată.

Menținerea CCH are două componente: *citiri* și *scrieri*. Copiile multiple nu constituie problemă pentru citire, dar un procesor trebuie să aibă un acces exclusiv atunci când scrie un cuvânt. Procesoarele trebuie să conțină cea mai recentă copie la citirea unui obiect, astfel că toate procesoarele trebuie să obțină noile valori după o scriere. Prin urmare protocoalele de monitorizare trebuie să localizeze toate C care protejează un obiect ce trebuie scris. Consecința unei scrieri a unei date protejate este fie *invalidarea* tuturor celorlalte copii, fie *actualizarea* copiilor partajate cu valoarea ce a fost scrisă.

Biții de stare dintr-un bloc cache sunt utilizați pentru *protocoale de monitorizare - snooping-* și informația este utilizată pentru monitorizarea activităților magistralei. La un eșec în citire toate C verifică să vadă dacă ele au o copie a blocului solicitat și apoi efectuează acțiuni specifice, ca de exemplu furnizarea datelor către C care a eșuat la citire. Similar, la o scriere, toate C verifică existența copiei la fiecare din ele și apoi acționează, fie invalidând, fie actualizând copia lor cu noua valoare.

Monitorizarea interferează rar cu accesele procesorului datorită existenței unui port suplimentar atașat fiecărui P - prin care se investighează porțiunea de etichetă referitoare la stare.

Protocoalele de monitorizare sunt:

- *invalidare de scriere / (write-invalidare)*: Procesorul care efectuează o operație de scriere provoacă invalidarea tuturor copiilor în alte C înainte de schimbarea copiei sale locale; apoi se actualizează data locală până când alt procesor o solicită. Procesorul care scrie trimite un semnal de invalidare pe magistrală, și toate C verifică

existența propriilor copii; dacă este așa, ele trebuie să invalideze blocul conținând cuvântul. Astfel, această schemă permite existența mai multor procesoare ce citesc, dar doar a unuia ce scrie;

- *actualizare de scriere / (write-update)* : Procesorul care scrie transmite noua dată prin magistrală; Toate copiile sunt apoi actualizate cu noua valoare. Această schemă cunoscută sub numele de *difuzare de scriere*, transmite continuu comenzi de scriere spre data partajată în timp ce procedura *invalidare de scriere*, șterge toate celelalte copii, astfel că există doar o copie locală pentru scrieri ulterioare.

Protocolul *actualizare de scriere* este similar cu protocolul *write – through* deoarece toate scrierile parcurg magistrala pentru a actualiza copiile datelor partajate. Protocolul *invalidare de scriere* utilizează magistrala doar la prima scriere pentru a invalida celelalte copii, și prin urmare scrierile următoare nu vor afecta traficul magistralei. În consecință protocolul *invalidare de scriere* are avantajul de a reduce accesul la magistrală, deci o reducere a traficului, în timp ce *actualizare de scriere (scrie – actualizează)* are avantajul de a face ca noile valori să apară în memoriile C, prin urmare se poate reduce latența.

SMP–MP comerciale bazate pe magistrală unică partajată utilizează protocolul *invalidare de scriere* ca și protocol standard de asigurare a CCH.

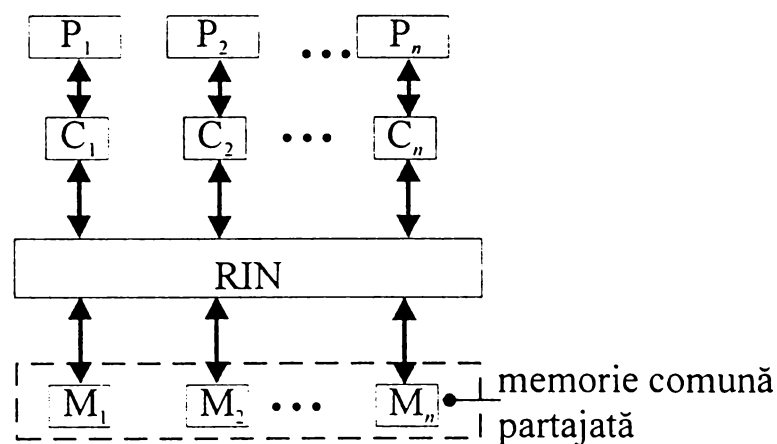


Fig. 2.6.a SMP - MP cu MCH private

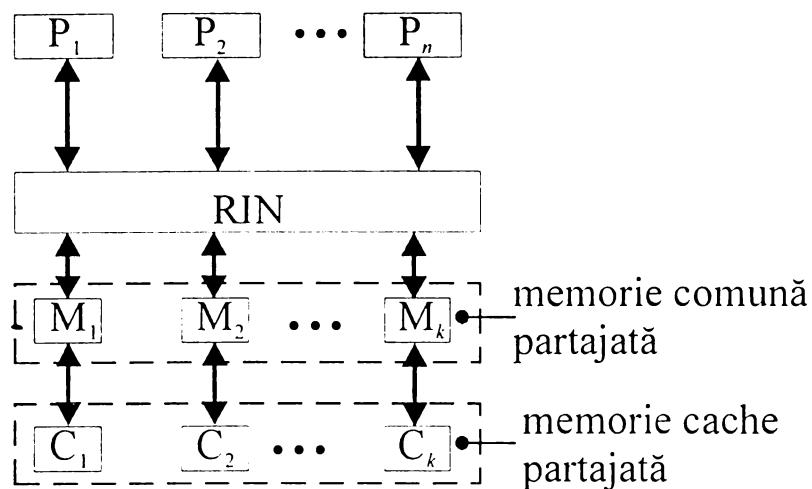


Fig. 2.6.b SMP - MP cu memorii cache partajate

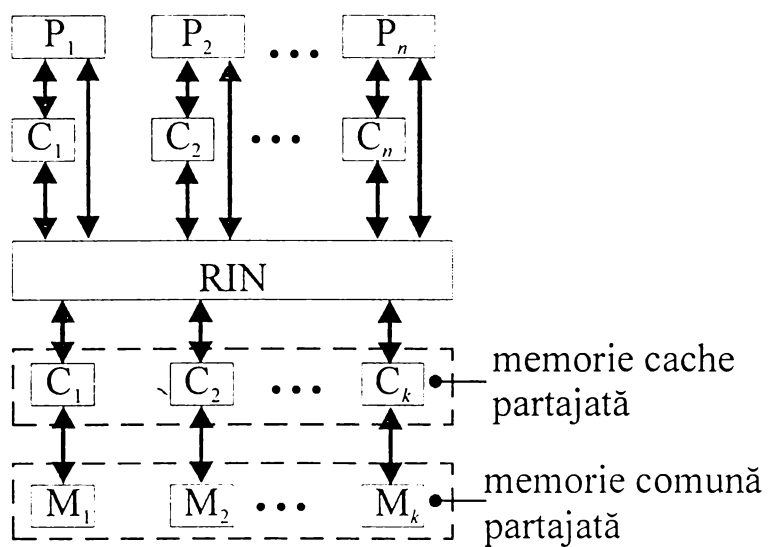
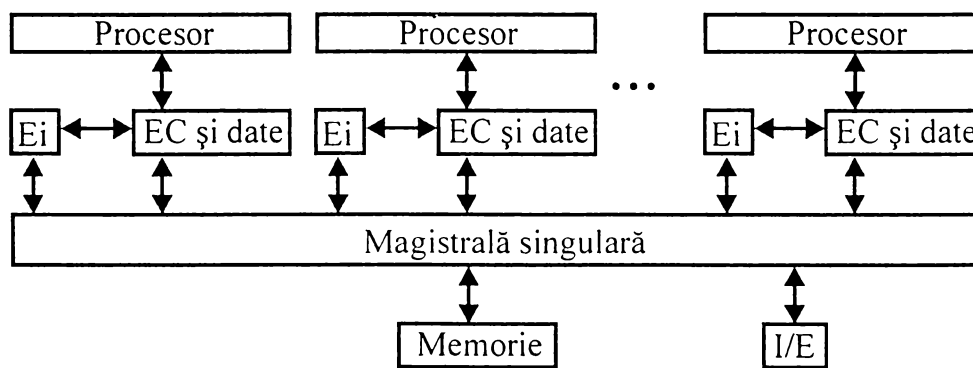


Fig. 2.6.c SMP - MP având memorii cache private și partajate



Ei - etichetă pentru înregistrare

EC - etichetă pentru cache-are

Fig. 2.7 Structura unui SMP-MP cu magistrală unică ce utilizează tehnica *snooping* pentru asigurarea coerenței

Un exemplu de protocol *invalidare de scriere*

În fig. 2.8 se prezintă o diagramă de stări finite pentru un protocol *invalidare de scriere* bazat pe o politică de “*rescriere*”. Fiecare bloc cache este în unul din cele trei stări:

1. *citește numai*: Acest bloc cache este “curat” (nemodificat) / (clean) și poate fi partajat de alte procesoare;
2. *citește / scrie*: Acest bloc este “viciat” (modificat) / (dirty) și nu este partajat;
3. *invalid*: Acest bloc nu are date valide.

Cele trei stări sunt duplicate în figură, pentru a arăta tranziții bazate pe acțiuni ale procesorului în opoziție cu tranziții bazate pe operațiuni de magistrală.

În realitate există doar un singur automat cu stări finite pentru fiecare C_i , cu stimuli provenind fie de la P_i , fie de pe magistrală.

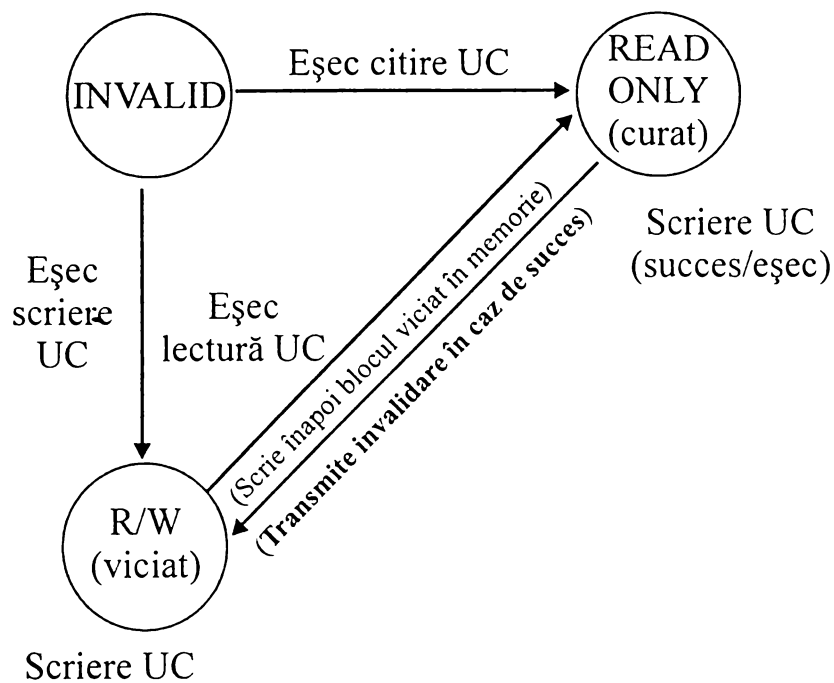
Tranziții de stare ale unui C_i apar la eșecul în citire, la eșecul în scriere sau în scrieri reușite. Succesele la citire nu schimbă starea lui C_i .

Să analizăm ce se întâmplă în cazul unui eșec la citire.

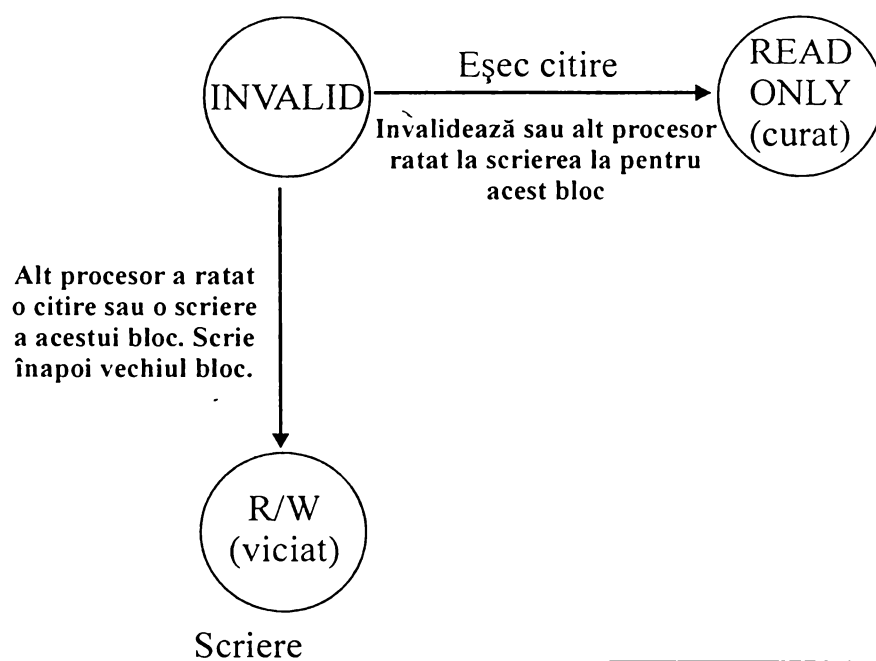
Cînd P_i are un eșec la citire, el va schimba starea blocului din C_i în Citește Numai și va rescrie vechiul bloc dacă acesta a fost în starea Citește/Scrie (viciat). Toate celelalte memorii cache din restul procesoarelor monitorizează eșecul la citire pentru a observa dacă acest bloc este în memoriile C private. Dacă unul are o copie și este în starea Citește/Scrie, atunci blocul este scris în M și, în acest protocol, blocul este poziționat în starea Invalid. Eșecul în citire este apoi satisfăcut citind din memorie.

Pentru a scrie un bloc, procesorul câștigă accesul la magistrală, trimite un semnal de invalidare, scrie în acest bloc și îl plasează în starea Citește / Scrie. Datorită faptului că alte memorii cache monitorizează magistrala, toate C verifică să vadă dacă au copii ale acestui bloc; dacă da, atunci îl invalidează.

Evident, după cum se observă, există multe variante de scheme de asigurare a CCH. Ele vor fi discutate ulterior.



Tranziții ale stării lui C_i utilizând semnale de la P_i



Tranziții ale stării lui C_i utilizând semnale provenind din magistrală

Observații:

1. Săgețile simple și textele nesubliniate specifică acțiuni ce se produc în memorii C fără coerență.
2. Săgețile și textele îngroșate reprezintă acțiuni adăugate pentru asigurarea CCH.

Fig. 2.8. Diagrama de stări pentru un protocol *invalidare de scriere*

Scheme de asigurare a coerenței la memoriile cache (SAC- MCH)

Pentru a păstra coerența memoriei cache este necesar să ne asigurăm că toate copiile păstrate în C, ale unei locații din M sunt identice.

Protocoale realizate hardware sunt utilizate în arhitecturile cu magistrale partajate. În general acestea utilizează un mecanism de transmitere asigurat de magistrala partajată pentru a menține coerența memoriilor cache multiple (MCHM).

În aceste protocoale, fiecare operație de scriere într-o locație din M - partajată între procesoare - este transmisă prin magistrală pentru celelalte procesoare ca să *invalideze* sau să *actualizeze* propriile lor copii ale locației, dacă ele sunt prezente.

Dacă într-o memorie ierarhizată dintr-un SMP - MP coexistă copii multiple ale aceleiași date, accesul la aceste copii trebuie să fie manipulat și/sau coordonat atent pentru a preveni posibilitatea accesării unei copii pierdute.

Reamintim că un sistem este *coerent*, dacă valoarea returnată la o instrucție de tip LOAD este întotdeauna identică cu valoarea scrisă la aceeași adresă de către ultima instrucție STORE.

În fig. 2.9 se prezintă o taxonomie a diferitelor scheme de asigurare a coerenței (SAC-MCH).

În literatură [SLM89], [CBZ95], [Wils87] se inventariază câteva din scheme propunându-se soluții de extindere la SMP - MP de talie mare.

Aceasta taxonomie se bazează pe ideea că dacă întărirea coerenței nu este importantă, atunci este indiferent dacă invalidările sunt efectuate de un procesor local (de ex. auto-invalidare) sau de către alte procesoare (invalidări induse).

Schemele sunt în continuare clasificate în funcție de solicitarea de-a avea sau nu resurse partajate și, dacă este așa, ce tip de resurse trebuie alocate pentru mărirea coerenței.

A. SAC- MCH bazate pe ne-invalidare

În literatură [SLM 89] este indicată metoda *difuzare de scriere / write-broadcast* în care MCH private sunt menținute coerente prin forțarea dispozitivelor de comandă ale MCH să informeze alte MCH din sistem asupra fiecărei scrieri de-o astfel de manieră încât fiecare MCH să poată actualiza propria sa copie, dacă o astfel de copie există.

Cea mai simplă schemă dintre cele bazate pe metoda *difuzare de scriere*, este *difuzare de scriere cu scriere prin tot spațiul*. / *write-broadcast with write-through*.

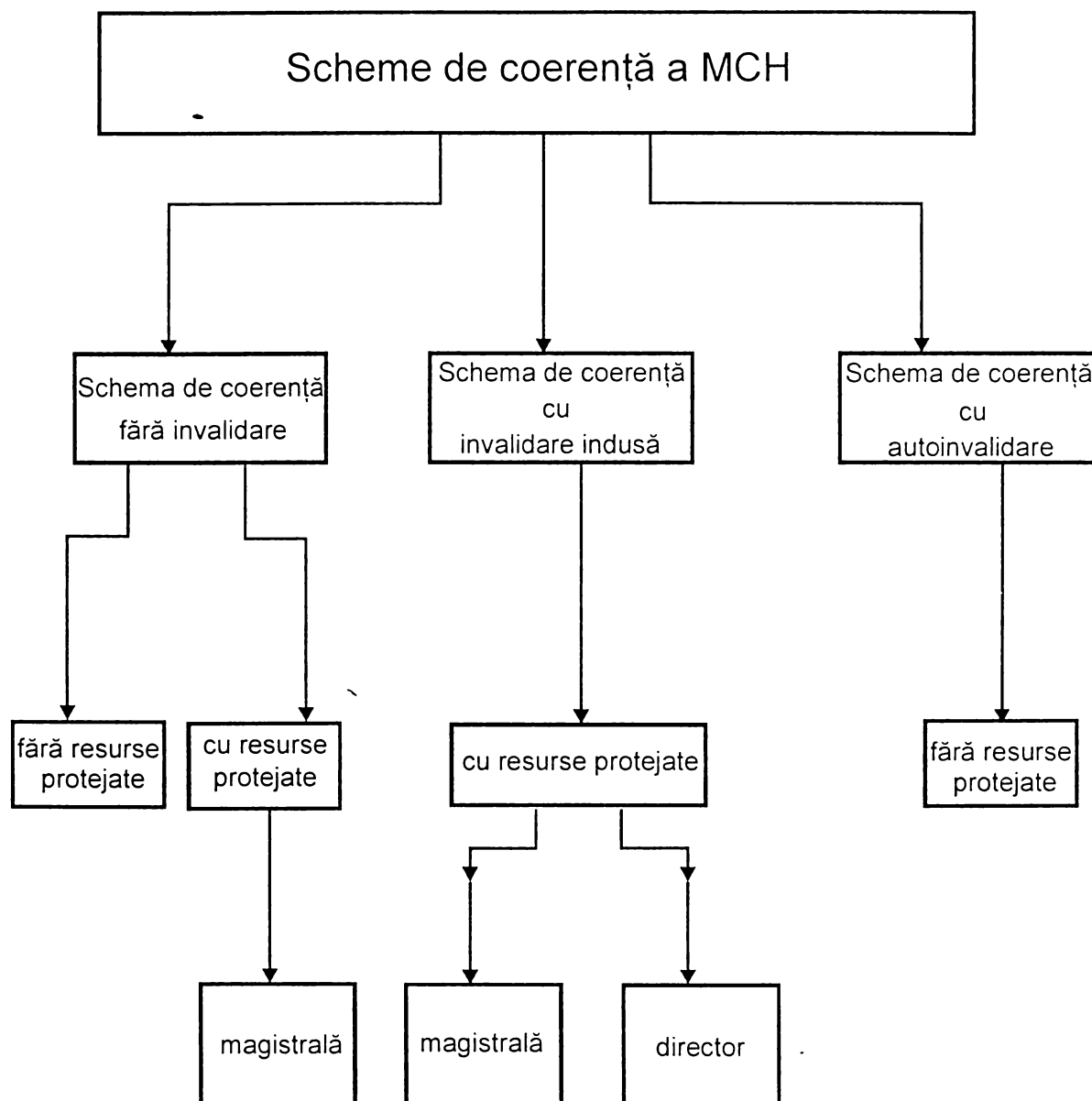


Fig. 2.9 Taxonomia schemelor de asigurare a coerenței memorilor cache

În această schemă, toate scrierile sunt imediat reflectate în MP - memoria globală partajată. Fiecare MCH din sistem urmărește aceste tranzacții *scrie-în-memorie* și își actualizează propria sa copie. Această procedură de urmărire implică existența unui mediu cu magistrală partajată.

Este de notat faptul că traficul de scrieri va satura curând magistrala de memorie chiar și în condițiile existenței unui număr redus de procesoare, ceea ce limitează această abordare.

Pentru a reduce utilizarea magistralei de memorie, multe SAC-MCH bazate pe tehnica *difuzare de scriere* sunt proiectate pentru a diminua numărul de cicluri difuzare de scriere. De exemplu, operații *difuzare de scriere* sunt efectuate doar atunci când blocul ce este scris se regăsește și în alte MCH.

Chiar dacă schema *difuzare de scriere* a fost utilizată doar în SMP-MP bazate pe magistrală, ea poate fi utilizată la orice SMP-MP. Protocolul va fi *bazat pe director / directory-based* cu dispozitivele de comandă ale memoriei având responsabilitatea aducerii la zi a copiilor MCH în cadrul procesului de scriere.

După recepționarea unei cereri de scriere de la un procesor, controler-ul de memorie va trimite cereri de actualizare a intrării cache corespunzătoare către MCH ce au o astfel de intrare.

Acest lucru se întâmplă după recepționarea semnalelor de confirmare a transferurilor de la toate MCH spre care s-au transmis semnale de actualizare pe care controler-ul de memorie le trimite ca și confirmare pentru cererea de scriere a procesorului solicitat. Un astfel de protocol este inefficient și nu va mai fi luat în considerare.

B. SAC-MCH bazate pe invalidare indusă.

În schemele de asigurare a coerenței bazate pe *invalidare - indusă/induce invalidation* intrările în MCH private pot fi invalidate prin mesaje generate de alte procesoare. În această clasă de scheme de asigurare a coerenței, resursele partajate sunt uzual solicitate pentru a decide dacă este nevoie de invalidare, pentru a decide pe cine să invalideze și/sau să actualizeze invalidările. În cele ce urmează se vor analiza schemele de asigurare a coerenței MCH din această clasă în concordanță cu tipul de resursă partajată pe care acestea le solicită.

B.1 Scheme bazate pe magistrală

Într-un SMP cu magistrală partajată cu o schemă de asigurare a coerenței bazată pe *invalidare indusă*, tranzacții de scriere sunt transmise pe magistrală. Fiecare controler al

MCH scanează magistrala și invalidează o intrare de mixare în MCH, dacă o astfel de intrare există. Magistrala este de asemenea utilizată să alimenteze cea mai actuală dată fie din memorie fie din MCH privată la un eșec în citire.

În această clasă de scheme, o stare este uzual asociată cu fiecare bloc cache în fiecare MCH.

În mod uzual starea indică când un bloc este valid, partajat sau exclusiv, și/sau dacă memoria este de actualizat în raport cu blocul cache corespunzător ("nemodificat") sau nu ("modificat").

De exemplu, starea locală în schemele de asigurare a coerenței de tip *write-once* poate fi una din următoarele : (1) *invalidă*, (2) *nemodificat-partajată*, (3) *nemodificat -exclusivă*, (4) *nemodificat-exclusivă*.

Toate controler-ele MCH în sistem lucrează împreună, fiecare cu fiecare și cu memoria pentru a asigura un sistem de memorie coerent.

Cele mai multe dintre aceste scheme permit ca multiple copii ale unui bloc " nemodificat " să fie conținute în MCH private, dar cel mult un MCH privat poate conține o copie "modificată ".

B.2 Scheme bazate pe director

Într-o schemă de asigurare a coerenței MCH cu *invalidare indusă bazată pe director*, (*directory-based*) câte o intrare *director (index)*, ce este memorată în controlerul memoriei, este asociată fiecărui bloc de memorie. Această intrare codifică starea blocului. Starea este utilizată pentru a decide dacă există o solicitare pentru invalidări în cadrul unei tranzacții date de scriere în bloc și, dacă este așa, să localizeze MCH privată care are o copie a blocului ce trebuie invalidată.

Este de asemenea utilizată pentru a indica dacă blocul de memorie corespunzător este perimat (stale) sau nu, și dacă este așa, să localizeze MCH privată care este garantată să aibă cea mai recentă copie a blocului. În plus față de stare din directorul global, în mod uzual se asociază o stare locală cu fiecare MCH privat.

Această stare locală este utilizată să permită unei MCH private să proceseze cele mai multe cereri de la procesorul asociat fără implicarea altor acțiuni globale.

În [CF 78], Censier și Feautrier au propus o SAC- MCH cu invalidare indusă, bazată pe director. În schema lor, fiecare intrare director constă dintr-*un bit de semnalizare a stării* "modificat " (dirty) și o *etichetă de n-bit*, unde *n* este numărul MCH private din sistem. Al *i*-lea bit din eticheta unui bloc de memorie este poziționat dacă MCH privată corespunzătoare are o copie validă a blocului. Bitul de semnalizare indică dacă blocul de memorie este actualizat sau nu.

Starea locală poate fi una din următoarele trei stări :

- (1) *invalidă (invalid)*;
- (2) *nemodificat-partajat (clean-shared)*;
- (3) *modificat -exclusiv (dirty-exclusive)*.

La o citire reușită, data solicitată este imediat furnizată procesorului solicitant și nici o acțiune adițională nu este solicitată.

- Dacă există o scriere reușită și MCH are o copie *modificat-exclusiv* a blocului, scrierea poate să continue local fără să implice nici o acțiune globală. Dar dacă există o scriere reușită și MCH are o copie *nemodificat-partajat* a blocului, MCH va transmite o cerere către controlerul MCH pentru a invalida copii ale blocului prezente în alte MCH private.

La un eșec la citire, bitul de semnalizare și eticheta de prezență, asociată cu blocul de memorie solicitat sunt examinate pentru a vedea care MCH privată are copii valide ale blocului. Dacă o oarecare MCH privată are o copie *viciat-exclusiv* a blocului solicitat, blocul este rescris în memoria globală și utilizat să furnizeze date spre MCH solicitant. Biții "viciați" ai blocului sunt puși pe 0 în directorul general și în starea locală a MCH ce are o copie "viciată". Stările locale în aceea MCH și în cea solicitantă devin *nemodificat-partajate*. În alte cazuri (de ex., când nici una sau mai multe MCH au copii *nemodificat-partajate* ale blocului), blocul solicitat este furnizat de memoria globală și este încărcat în MCH ca și *nemodificat-partajat*.

La un eșec la scriere se testează din nou bitul de semnalizare și eticheta de prezență pentru a se vedea dacă una din MCH private are o copie validă a blocului solicitat.

Există trei cazuri ce trebuie analizate.

Dacă nici una din ele nu are o copie validă, blocul este furnizat de memoria globală. Dacă cel puțin o MCH privată are o copie *nemodificat-partajată* a blocului, toate copiile blocului sunt invalide și blocul solicitat este furnizat de memoria globală.

Dacă o MCH privată oarecare o copie *modificat-exclusivă* a blocului, copia modificată este rescrisă în memoria globală înainte ca ea să fie invalidată. Noul bloc scris este utilizat pentru a alimenta date spre MCH solicitant.

În toate cele trei cazuri menționate anterior, bitul de semnalizare global al stării modificate aparținând blocului solicitat este poziționat și blocul solicitat este încărcat în MCH solicitant ca și *modificat-exclusiv*.

Pentru a limita cantitatea informației globale de stare conținută în director, Archibald și Baer [SLM 89] propun o SAC- MCH bazată pe director denumită *schema cu director format din doi biți*.

În această schemă, informația directorului global constă din doi biți ce codifică una din cele patru stări posibile :

- (1) Blocul nu este prezent în nici o MCH;
- (2) Blocul este într-un MCH precis și nu este modificat;
- (3) Blocul este prezent în nici - unul sau în mai multe MCH și nu este modificat;
- (4) Blocul este prezent exact într-un MCH și este modificat ("viciat").

Deoarece directorul nu mai conține identități ale MCH conținând un bloc dat în schema propusă de autori, sunt necesare transmisiile pentru a efectua invalidări și pentru a solicita o rescriere a blocului modificat.

C. SAC-MCH bazate pe auto-invalidare

Chiar dacă SAC- MCH bazate pe invalidări induse pot fi eficiente în asigurarea unei rate înalte de succes datorită abilității lor de-a urmări dinamic traseul stării fiecărui bloc, traficul rețelei generate pentru solicitări de invalidare și pentru manipularea informației de stare locală/globală poate fi substanțial.

În SAC- MCH cu *auto-invalidare*, intrările în MCH pot fi invalidate doar de procesoare locale și nici o informație manipulată global nu este asociată cu blocuri MCH sau cu blocuri din memoria globală. Aceasta elimină traficul suplimentar în detrimentul unei gestionări mai puțin eficiente a MCH. În a decide ce și când să invalidezi pentru a

menține concret coerența, SAC- MCH cu *auto-invalidare* utilizează câteva informații din cadrul structurii programului paralel obținute în timpul compilării. Aceste scheme sunt uneori menționate ca fiind scheme asistate software datorită dependenței lor de compilator.

În literatură [SLM 89] se propune o *schemă cu invalidare selectivă rapidă*.

Fiecare referire la o operație de citire este marcată în cadrul compilării ca fiind fie *citire din memorie* (MCH poate să conțină o copie perimată a datei de citit) sau citire din MCH (la o apelare reușită, se garantează că intrarea MCH este actualizată).

În fig. 2.10 se prezintă două condiții posibile pentru un acces perimat.

În prima condiție, o intrare cache este încărcată printr-o scriere de către P_i și devine perimată când alt P_j transmite o scriere către aceeași locație de memorie. Dacă lui P_i îi este permis să acceseze copia perimată în propria MCH, el ar viola criteriul de corectitudine a coerenței MCH. În mod similar, un acces perimat este posibil când intrarea cache este încărcată printr-o operațiune eșuată de citire și devine perimată printr-o scriere provenind de la alt procesor. De notat că aceste condiții pretind ca fiecare acces la o variabilă partajată să fie marcată ca *citire din memorie* odată ce variabila a fost scrisă de două ori de către diferite procesoare.

Suplimentar cu schema de date pentru accese la date perimate, *un bit de schimb* se asociază cu fiecare cuvânt. Scopul lui principal este de-a permite reținerea referințelor *citește din memorie* ce apar mai mult de o singură dată la un moment dat într-un ciclu. Biții de schimb a MCH sunt poziționați la începutul fiecărui ciclu.

O *citire din cache (cache-read)* este executată ca de obicei, indiferent de starea bitului de schimb. La o *citire din memorie (memory-read)* dacă există un eșec sau un acces reușit cu starea bitului de schimb modificată, atunci cuvântul este încărcat în cache și se inițializează bitul de schimb. Această situație implică ca fiecare prin acces la o variabilă partajată după începutul unui ciclu să fie servit de memoria globală odată ce variabila a fost scrisă de două ori de către procesoare diferite. Aceasta limitează scalabilitatea schemei. La o *citire din memorie*, când a avut loc un acces reușit și bitul de schimb asociat a fost inițializat, solicitarea de acces este satisfăcută de către cache. În plus, fiecare scriere reușită inițializează biții de schimb asociați și actualizează memoriile principale și cache. În [SLM 89] se prezintă o schemă de asigurare a coerenței cu auto-

invalidare care permite operarea memoriilor cache în interiorul secțiunilor critice. În această schemă, operarea memoriilor cache în interiorul secțiunilor critice a fost posibilă prin utilizarea identificatorilor de tipul *one time identifier* (OTI) ce sunt garantați să fie unici în timpul existenței lor. Identificatorul OTI este utilizat pentru a se distinge intrările de actualizare ale memoriilor cache (încărcate în cache în secțiunea critică curentă) de acelea perimate (încărcate în cache în timpul secțiunilor critice trecute).

În acest scop, fiecare cache are un câmp OTI și un OTI distinct este folosit pentru fiecare execuție de secțiuni critice. La fiecare acces cache, OTI curent este comparat cu câmpul OTI în blocul cache corespunzător. Dacă câmpurile se potrivesc, solicitarea este satisfăcută de către memoria cache. Astfel, blocul este încărcat din memoria globală și câmpul OTI este poziționat ca fiind OTI curent. În tabelul din fig.2.11, citat în [SLM89], se prezintă un sumar a SAC - MCH cu *auto-invalidare* discutate în acest paragraf.

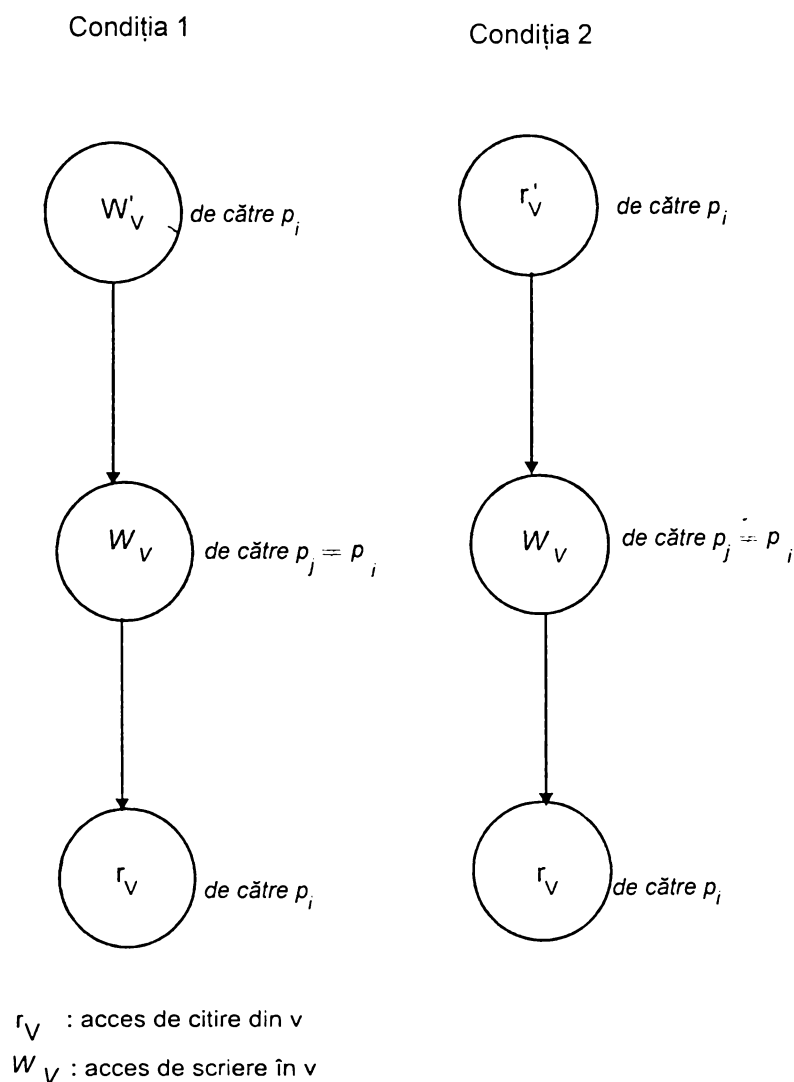


Fig. 2.10 Condiții pentru un acces la un bloc cache perimat

	Veidenbaum	McAuliffe	Lee	Veidenbaum Cheong	Cytron	Smith
unitate de întărire a coerenței	toate variabilele	fiecare variabilă	fiecare variabilă	fiecare referință	fiecare referință	fiecare pagină
regiune de întărire a coerenței	buclă	unitate de calcul (UC)	secțiune	program	program	secțiune critică
Detecția incoerenței	structură cu buclă încuibată	utilizare de variabile în UC	utilizare de variabilă într-o secțiune	analiză de flux	analiză de dependență de date	OTI
loc de invalidare	frontiere ale buclelor	frontiere ale UC-urilor	frontiere ale secțiunilor	frontieră de bucle	fiecare referință	sfârșit al secțiunii critice
politică de actualizare a memoriei globale	scrie prin	rescrie	rescrie	scrie prin	schemă hibridă	scrie prin

Fig. 2.11 Sumar al SAC-MCH cu *autoinvalidare*

2.4 REȚELE DE INTERCONECTARE

2.4.1 Terminologie

O rețea de interconectare RIN reprezintă o conexiune complexă de comutatoare și legături ce permite comunicarea de date între procesoare și memorii în SMP-MP - multiprocesoare - și între elemente de procesare în SMP -TM - multicalculatoare.

Reamintim că în subclasa SMP-MP toate procesoarele P partajează o memorie comună în timp ce la SMP-TM fiecare procesor are memoria sa proprie cât și alte resurse de I/E. De aceea vom utiliza denumirea de *element de procesare* -PE.

Există o relație directă între arhitectura unui sistem MIMD și topologia RIN utilizată.

SMP-MP utilizează RIN de tip *crossbar*, *magistrală multiplă*, *multinivel* etc. în timp ce SMP-TM utilizează topologii ca și : *stea*, *inel*, *arbore*, *hipercub*.

Compararea, clasificarea și ordonarea RIN se fac în scopul evaluării performanțelor și a stabilirii tipurilor de aplicații. Se utilizează mai multe criterii.

În literatură [BYA 89], [Prad 96], [Popa 96] etc. se admite că trebuiesc luate în considerare :

- filozofia temporizării;
- metodologia de comutare;

- strategia de comandă.

De aici rezultă principalele criterii de analiză :

- *modul de operare;*
- *tehnica de comutare;*
- *strategia de comandă;*
- *topologia;*
- *gradul de interconectare;*
- *complexitatea;*
- *gradul de blocare;*
- *capacitatea de reconfigurare.*

Modul de operare se referă la natura comunicării între moduri : *sincronă* și *asincronă*.

Modul de operare sincron presupune existența unui ceas central care fixează pașii în care se produce transferul de date.

Varianta asincronă asigură comunicările utilizând protocoale hand-shaking cu interblocare. Se asigură o bună expandabilitate și modularitate dar cresc dificultățile de proiectare.

Tehnica de comutare stabilește modul de direcționare a informației. Există două tehnici utilizate: *comutare de circuite* și *comutare de pachete*.

Comutarea de circuite presupune ca înainte de începerea transferului comutatoarele să fie astfel poziționate încât să stabilească o cale de comunicare între transmițător și receptor. Aceasta rămâne nemodificată pe toată durata efectuării transferului. Datorită existenței unui timp de stabilire a conexiunii cu implicare directă asupra latenței RIN, metoda este potrivită pentru blocuri lungi. Ele pot fi de dimensiuni variabile.

Comutarea de mesaje, pachete de date, presupune secționarea blocului ce trebuie transmis în entități de lungime fixă - pachete - care se transferă prin RIN după procedura "*memorează și transmite (store - and forward)*".

Un comutator care recepționează pachetul îl memorează și caută să stabilească o conexiune spre nivelul următor. Dacă nu găsește o cale liberă și un comutator liber, reține pachetul și așteaptă. În caz contrar, existând o cale liberă, pachetul este transmis la nivelul următor. Tehnica este utilizată pentru mesaje scurte.

Strategia de comandă definește modul prin care semnalele de comandă direcționează fluxul de date generat în RIN și stabilește modul de comandă a comutatoarelor.

Informația de comandă, care poziționează comutatoarele, poate fi generată de un controler central sau poate fi conținută în mesajul ce este transferat prin RIN.

În consecință controlul poate fi *centralizat* sau *distribuit*. Controlul centralizat, datorită existenței controlerului central, introduce întârzieri în stabilirea căii de comunicație pentru că elementul de control trebuie să cunoască toate cererile procesoarelor și să rezolve un eventual conflict. De asemenea introduce o strangulare a sistemului (system bottleneck) cu efecte negative directe asupra dependabilității sistemului.

Controlul distribuit descentralizează decizia de comutare la nivelul comutatoarelor în sensul că se atribuie un controler simplificat fiecărei componente a sistemului.

În aplicații ale multiprocesoarelor, comanda unui RIN *crossbar* este uzual centralizată iar comanda unui RIN-MN (Multiple Level Interconnection Network) este descentralizată.

Bazat pe caracteristicile operaționale prezentate mai sus, RIN pot fi clasificate în opt categorii diferite pentru o topologie dată.

Schema de clasificare este prezentată în fig. 2.12

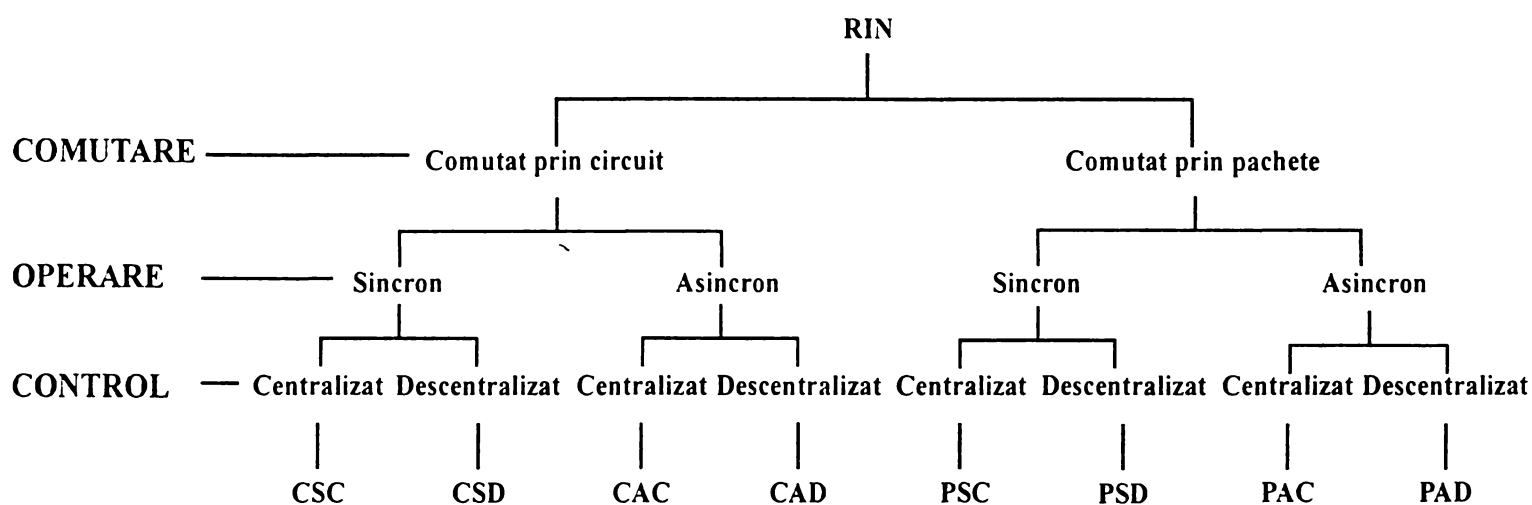


Fig. 2.12 Clasificare a RIN

Topologia definește forma unui RIN. Aceasta poate fi descrisă printr-un graf în care *nodurile* sunt comutatoarele iar *arcele* reprezintă legăturile dintre ele. De fapt trebuie notat că un nod reprezintă un cuplu procesor-memorie-comutator. Comutatoarele apar însă ca noduri ale grafului.

Topologiile pot fi *regulate* și *neregulate*. Din punct de vedere a facilităților de reconfigurare ele pot fi *statice* sau *dinamice*.

RIN statice (static interconnection network) sunt acelea care permit numai legături fixe (uni sau bidirecționale) între două noduri. Aceste RIN sunt tipice pentru multicalculatoare adică SMP – TM.

RIN statice se pot clasifica în funcție de criteriul spațial în topologii:

- pe o dimensiune;
- pe două dimensiuni;
- pe trei dimensiuni;
- pe mai mult de trei dimensiuni.

Exemplu de topologii RIN statice:

monodimensională	-	<i>magistrala comună</i> ;
bidimensională	-	<i>inel, stea, arbore, plasă, sistolă</i>
tridimensională	-	<i>cub</i>
multidimensională	-	<i>hipercub</i>

În fig. 2.13 se prezintă structura unui nod ce este interconectat în RIN.

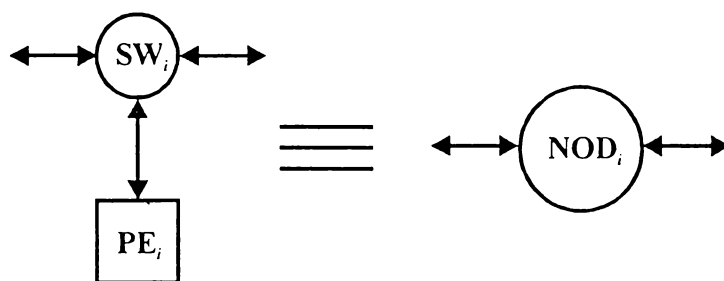


Fig 2.13 Structura unui nod de interconectare la RIN

În fig. 2.14 se prezintă câteva exemple de RIN statice. Acestea sunt *magistrala comună*, *inel*, *stea*, *arbore*, *plasă*, *sistolică*, *cub 3D*, *hipercub*. Se remarcă diferențele în topologie între aceste tipuri de interconectări. La o extremitate se găsește RIN – *magistrală comună* în care toate nodurile comunică între ele prin intermediul unei magistrale partajate în timp, iar la cealaltă se găsesc RIN de tip *cub*, respectiv *hipercub* în care fiecare nod este interconectat cu celelalte.

În cadrul sistemelor multiprocesor cu memorie partajată, cea mai simplă conexiune o reprezintă *magistrala comună* a cărei analiză se va face ulterior. Sistemele de calcul distribuite sunt caracterizate prin RIN *plasă* și RIN *cub*, respectiv *hipercub*.

Rețelele cu arii locale – LAN – sunt conectate în general prin intermediul RIN-*stea* respectiv RIN-*inel* în timp ce RIN-*sistolică* face obiectul unei clase separate de mașini paralele.

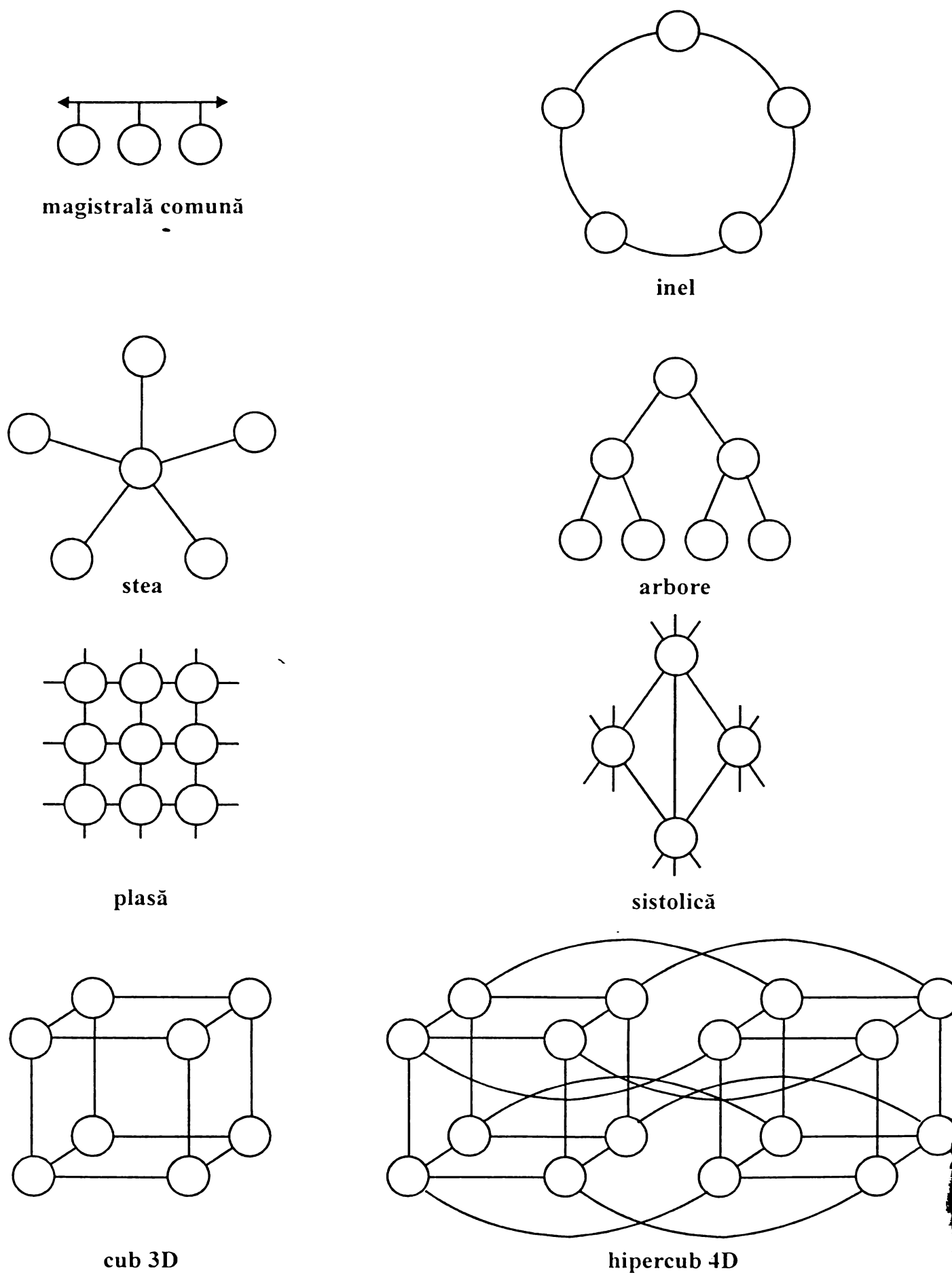


Fig. 2.14 RIN statice. Exemple

Topologiile dinamice se caracterizează prin existența comutatoarelor active ce stabilesc calea de comunicație între două elemente de procesare.

RIN dinamice pot fi configurate în topologii pe:

- un nivel
- mai multe nivele
- tip *crossbar*.

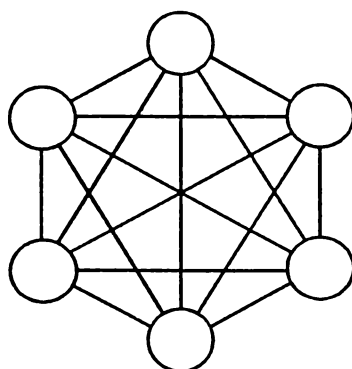
Topologiile pe un nivel având dispuse comutatoarele pe acesta, sunt cu recirculare în sensul că un mesaj ca să ajungă la destinație poate să parcurgă de mai multe ori RIN.

Topologiile organizate pe mai multe niveluri sunt formate din cupluri de *conexiuni de permutare* și *etaje de comutare* dispuse pe niveluri succesive. Exemple sunt RIN de tip *banyan*, *delta*, *omega*.

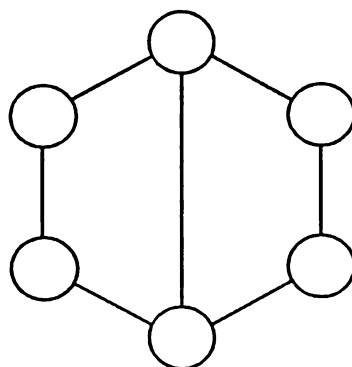
Topologia *crossbar* asigură căi de comunicație separate pentru fiecare pereche de parteneri. Este cea mai complexă RIN dar și cea mai scumpă.

Gradul de interconectare împarte topologiile RIN în două: cu *interconectare totală* și cu *interconectare parțială*.

În fig. 2.15 se prezintă exemple de topologii cu cele două tipuri de interconectare.



**RIN cu
interconectare totală**



**RIN cu
interconectare parțială**

Fig. 2.15 Exemple de interconectări

Complexitatea este un criteriu de analiză în strânsă corelare cu costul RIN.

Numărul de comutatoare poate fi determinat pe baza următoarelor tipuri de funcții:

- $f(N)$ - depinde liniar de numărul N al partenerilor;
- $f(N \times \log N)$ - depinde logaritmic de numărul partenerilor;
- $f(N^2)$ - dependență patrată.

RIN de tip *stea*, *inel*, *plasă*, *cub* și cele *dinamice* pe un singur nivel au o complexitate dată de $f(n)$.

În a doua categorie intră RIN dinamice pe mai multe nivele. Baza logaritmului este dată de numărul intrărilor și ieșirilor comutatorului. Pentru un comutator cu M intrări și M ieșiri funcția va fi de tipul $f(N \times \log_M N)$.

În literatura de specialitate [Hwan 93] se consideră că aceste tipuri de RIN optimizează raportul performanță/cost.

RIN a căror complexitate și cost a topologiei sunt determinate de $f(N^2)$ au performanțe maxime oferind conectivitate totală în timp minim. Datorită costului mare, devin nerentabile pentru un N mare.

2.4.2 Sisteme multiprocesor cu memorie partajată orientate pe magistrală

Cea mai simplă conexiune între procesoare și alte unități funcționale în cadrul unui SMP – MP o reprezintă *magistrala partajată (shared-bus)*.

Magistrala reprezintă o cale de comunicare unică la care se conectează unități funcționale care pot fi:

- procesoare;
- memorii;
- porturi de I/E.

Magistrala comună este o unitate pasivă neexistând comutatoare. Controlul transferului de date/instrucții se face de către interfețele existente la fiecare unitate conectată pe magistrală.

În fig. 2.16 se prezintă configurația unui SMP – MP orientat pe magistrală.

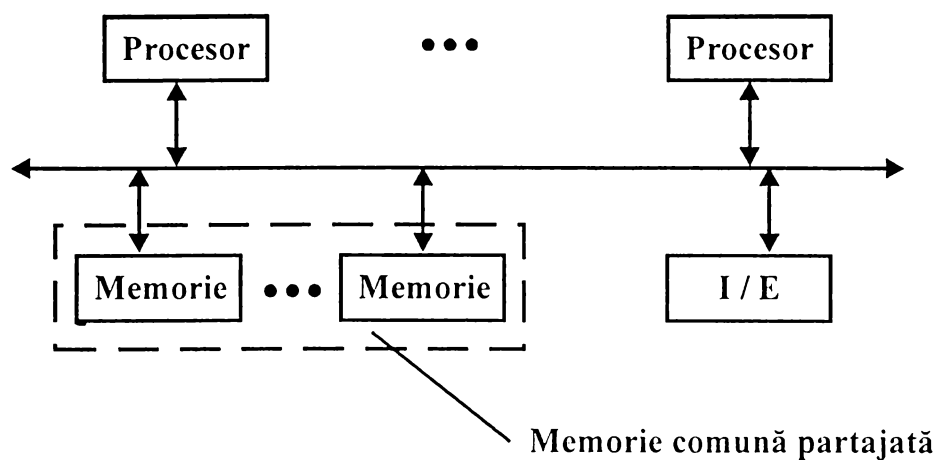


Fig. 2.16 SMP – MP orientat pe magistrală

Magistrala este *comună* pentru că la ea se conectează toate resursele hard ale sistemului și este *partajată în timp* pentru că la un moment dat doar două unități funcționale sunt logic conectate între ele efectuându-se un transfer de date.

Structura magistralei este construită în jurul conceptului *master-slave*, unde o unitate funcțională de tip *master* preia controlul magistralei și a unității funcționale de tip *slave*, după ce i-a decodificat adresa și efectuează transferul. Această procedură - bazată pe un protocol de tip *hand-shake* – permite conectarea modulelor de diferite viteze prin intermediul magistralei. Ea este prin urmare *asincronă*.

De asemenea, o caracteristică importantă a magistralei este abilitatea sa de a permite conectarea mai multor unități funcționale de tip *master* în vederea multiprocesării.

Un procesor se poate găsi în două ipostaze: *master* sau *slave*. O memorie nu poate fi decât *slave* iar o unitate funcțională de I/E, de obicei în postură de *slave*, poate fi *master* doar dacă are facilități DMA.

Dacă mai multe unități *master* solicită accesul la magistrală este necesară *arbitrarea* solicitărilor pentru a se evita situația de *conflict*.

Conflictul poate apare dacă:

- două sau mai multe unități *master* doresc să comunice cu aceste unități *slave*;
- două sau mai multe unități *master* necesită aceeași cale de comunicare.

Dat fiind faptul că toate transferurile se efectuează prin magistrala comună, aceasta limitează drastic performanțele sistemului.

Debitul global al sistemului (overall throughput) este limitat de apariția *strangulării la nivel de magistrală (bus-bottleneck)*, produs de solicitari simultane multiple ale masterilor.

În paragraful 3.2 vom detalia analiza de performanțe ale acestei conexiuni.

Este evident că reducerea traficului pe magistrală presupune o mai mare utilizare a resurselor proprii atașate fiecărui procesor. Prin urmare, crearea unui sistem ierarhizat de memorii între procesoare și magistrala comună devine o necesitate.

În [Coif 94] sunt prezentate mai multe soluții de ierarhizare a memoriilor.

Din punct de vedere fizic pot exista, în raport cu un procesor, *memorii locale* și *memorii externe*.

Un procesor poate adresa memoria sa locală prin intermediul unui singur nivel de conectare- *magistrala locală*- și o memorie externă prin intermediul mai multor nivele de conectare – *magistrale locale* și *magistrala globală*. Pot exista memorii care să nu poată fi adresate de unele procesoare.

Din punct de vedere logic, memoriile se împart în *private* și *comune*. Cele private pot fi adresate doar de procesoarele la care sunt conectate.

Procesoarele execută programe rezidente în memoriile locale. Ele cooperează prin transmiterea de mesaje prin intermediul memoriei comune partajate. Memoria comună (MC) poate fi implementată fie utilizând una sau mai multe module de memorie externă tuturor procesoarelor, fie distribuind-o în partea neprivată a memoriilor locale (ML).

Un procesor poate fi în una din următoarele stări:

- *activă* – execută un program din memoria sa privată;
- *de adresare* – cooperează cu alte procesoare utilizând MC;
- *de așteptare* – așteaptă să adreseze o MC;
- *blocată* – un procesor este blocat de un alt procesor care adresează segmentul de MC din memoria sa locală – ML.

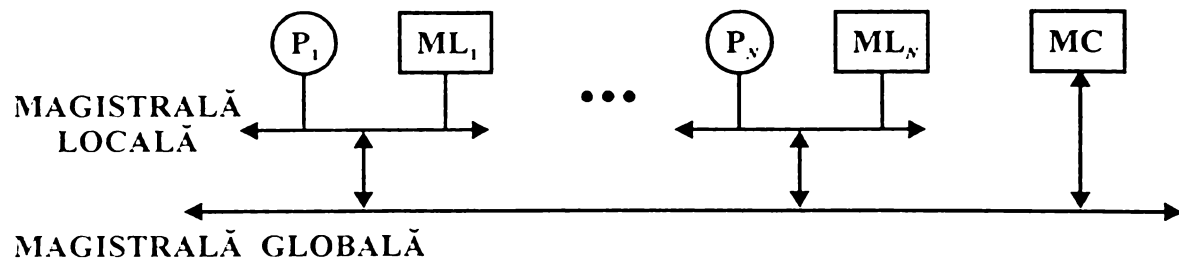
Structura 1 este prezentată în fig. 2.17 a). Se caracterizează prin existența unei memorii comune MC, externe oricărui procesor și adresabilă doar prin magistrala globală. Fiecare procesor P dispune de o memorie locală ML, neadresabilă decât de procesorul la care este atașată. Conflictul apare doar la adresarea memoriei MC.

Structura 2 este prezentată în fig. 2.17 b). Se observă distribuția MC. Fiecare ML este configurată din două zone: una privată MP și una comună MC. Fiecare P este conectat la ML proprie prin intermediul unei magistrale locale. Pentru ca un P_i să adreseze o zonă comună din ML_j atașat unui alt procesor P_j se vor utiliza: magistrala locală a lui P_i , magistrala globală și magistrala locală a lui P_j . Un procesor care primește accesul la magistrala globală deține prioritate pentru accesul la orice resursă. El devine prioritar în raport cu alte procesoare.

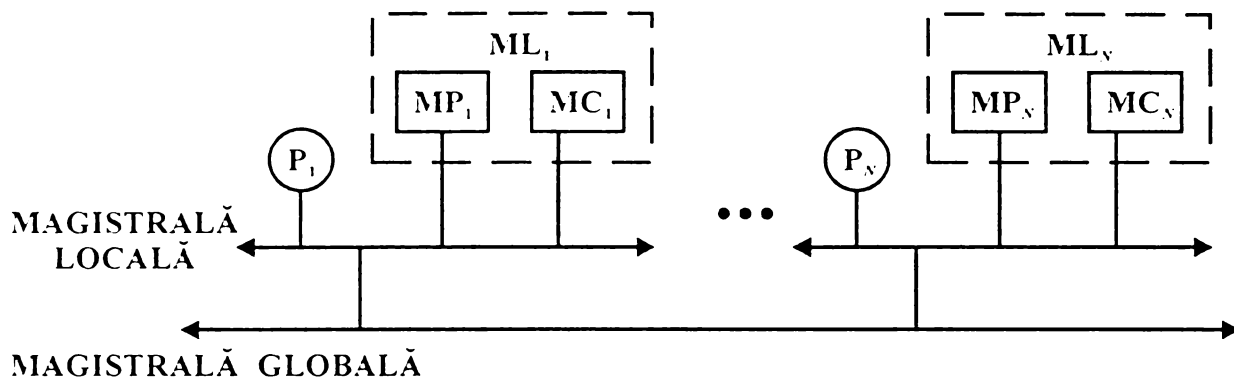
Dacă aceste procesoare erau în stare *activă* ele trec în stare *blocată*. Dacă erau în stare de *așteptare* rămân în continuare în această stare dar eliberează magistrala locală. Acest protocol previne blocarea și ameliorează performanțele.

Structura 3 este prezentată în fig. 2.17 c). Este o ameliorare a structurii 2 prin utilizarea unei memorii *dual-port* pentru a implementa zona comună a memoriilor locale ale fiecărui procesor. Zonele comune devin direct adresabile prin magistrala globală. Se scurtează timpul de transfer. Nu apare conflict nici la nivelul magistralelor locale și nici la cel al memoriei dual-port care suportă două accese simultane. Conflictul apare doar la nivelul magistralei globale. În această structură un procesor nu poate adresa zona comună a propriei ML decât prin intermediul magistralei globale.

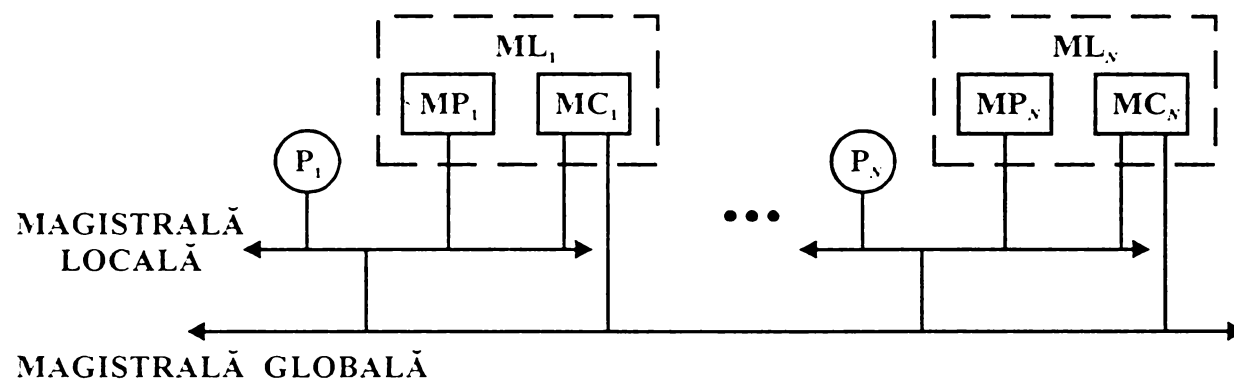
Structura 4 este prezentată în fig. 2.17 d). Reprezintă o variantă a structurii 3 în cazul în care nu există disponibilă o memorie dual-port. Zonele de memorie comună sunt externe tuturor procesoarelor. Ca și la structura 3 procesoarele nu pot adresa zonele comune ale propriilor memorii prin intermediul magistralei globale ci dispun de o cale dedicată acestui scop. Doar un singur P_i poate adresa la un moment dat o memorie comună. Conflictul poate apare la nivelul memoriilor comune și a magistralei globale. Prioritatea este dată accesului ce vine de pe magistrala globală.



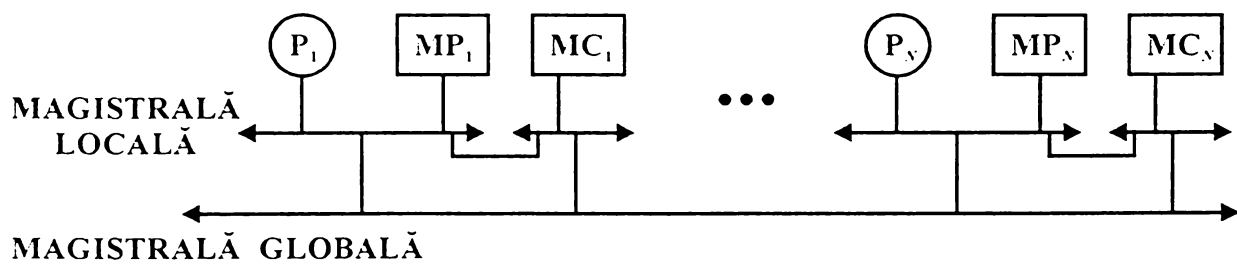
a) Structura 1



b) Structura 2



c) Structura 3



d) Structura 4

Fig. 2.17 Memorii ierarhizate pentru optimizarea traficului pe magistrală

Structura generală a unei magistrale cuprinde fire pe care circulă următoarele tipuri de semnale:

- adrese;
- date;
- control;
- întreruperi;
- acces la magistrală;
- tensiuni electrice de referință, etc.

În [Coif 94] sunt prezentate mai multe tipuri de magistrale dar ne vom referi doar la MULTIBUS, familia de magistrale dezvoltată de firma INTEL și care face obiectul standardului IEEE 796.

La MULTIBUS spațiul de memorie maxim admisibil este 1MO iar cel de I/E de 64 K port-uri.

Modulele ce se pot conecta pe magistrală sunt de tip *master* și *slave*. Un modul *master* poate comanda liniile de adrese și de control și deci poate prelua controlul magistralei. Un modul *slave* nu are această posibilitate. În cazul existenței mai multor masteri este necesară arbitrarea. Sunt permise două tipuri de arbitrări: serială și paralelă. Tactul magistralei secvențiază operația de arbitrare.

Se pot efectua transferuri de date pe 8 sau 16 bit, existând două linii *A0* și */BHEN* care indică lungimea cuvântului transferat. Există patru posibilități:

- transfer pe 16 bit;
- transfer de octet par, pe jumătatea inferioară a liniilor de date;
- transfer de octet impar, pe jumătatea superioară a liniilor de date;
- transfer de octet impar, pe jumătatea inferioară a liniilor de date.

O caracteristică a MULTIBUS-ului este că poate insera o memorie ROM suprapusă peste una RAM sau peste alta ROM. Facilitatea este utilă în cazul existenței unei rutine *bootstrap*, în cazul plasării porturilor în spațiul de memorie sau pentru testare.

MULTIBUS suportă două tipuri de întreruperi: *vectorizate* și *nevectorizate*. Există 8 linii pentru cereri de întreruperi. Cele vectorizate cer transferul unui vector de la un *master* către un *slave* utilizând secvența de acceptare indicată de activarea liniei */INTA*. Cele

nevectorizate presupun că localizarea sursei de întrerupere se face în cadrul masterului fără a mai fi necesar transferul unui vector.

MULTIBUS este compusă din 73 de linii după cum urmează:

- 20 de adrese;
- 16 de date;
- 1 pentru inițializare;
- 2 pentru inhibare;
- 1 pentru stabilirea, împreună cu Ao, a lungimii cuvântului transferat;
- 7 pentru arbitrare și control;
- 5 pentru comanda transferurilor de date;
- 9 pentru sistemul de întreruperi;
- 4 pentru sesizarea stării tensiunii de alimentare;
- 4 auxiliare;
- 4 pentru alimentare: +5V, -5V, +12V, 0V.

Există terminale pentru dezvoltări ulterioare.

Caracteristicile RIN de timp magistrală comună sunt:

- avantaje:
 - simplitate logică și constructivă;
 - posibilitate de extensie cu module definite și proiectate în jurul ei (around it). Familia MULTIBUS este un bun exemplu.
- dezavantaje:
 - saturarea rapidă a magistralei cu repercursiuni negative asupra latenței RIN;
 - toleranța minimă la defect, orice defecțiune conducând la căderea catastrofală a întregului sistem.

RIN de tip *magistrală comună multiplă* păstrează principiul magistralei comune singulare dar îi și ameliorează performanțele din punct de vedere al toleranței la defecte.

În fig. 2.18 se prezintă schematic structura unui SMP – MP cu magistrală multiplă.

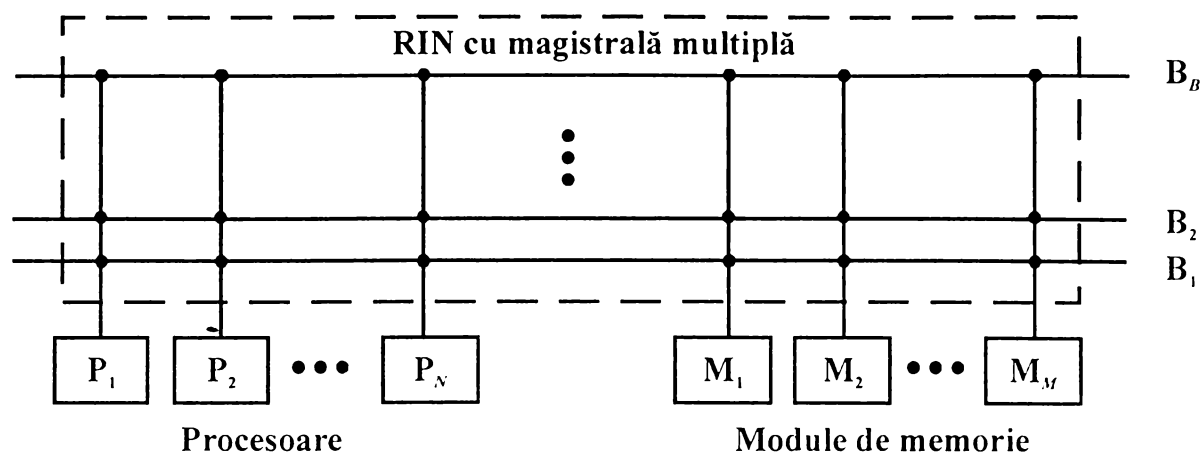


Fig. 2.18 SMP- MP cu magistrală multiplă

Se produce o decongestionare a traficului deci o scădere a latenței RIN și în același timp o ameliorare a toleranței la defect. Crește complexitatea și costul.

În cap. 3.3 se studiază amănunțit acest RIN din punct de vedere al performanțelor și se descrie structura unui SMP – MP construit în jurul unei astfel de magistrale.

2.4.3 Rețele de interconectare statice

Rețelele de interconectare statice (static interconnection network) sunt acelea care permit numai legături fixe (uni sau bidirecționale) între două noduri. Acestea sunt potrivite pentru conexiuni procesor – procesor, adică în SMP – TM.

În literatură [Prad 96], [Hwan 93] sunt descrise pe larg topologii de RIN-stactice. În continuare ne vom referi la topologiile *arbore (tree)*, *plasă (mesh)* și *hipercub (hypercube)* deoarece aici se pot face intervenții semnificative pentru exploatarea potențialului de toleranță la defecte prin reconfigurare.

RIN arbore

Sunt într-o oarecare măsură similare cu cele de tip *stea* în sensul că nodul rădăcină (apex) se folosește pentru intermedierea transferului de la o ramură a arborelui la cealaltă. În general de la nodul sursă se parcurge arborele în sus până când se întâlnește un nod comun cu calea până la nodul destinație, după care se parcurge această cale descendent. Numărul de legături al RIN este $N-1$, adică costurile cresc liniar cu dimensiunea. Defectarea unui nod critic – apex-ul sau unul apropiat lui – are urmări catastrofale asupra întregii RIN.

Într-o rețea regulată cu b ramuri sub fiecare nod și cu j niveluri, numărul total N de noduri este:

$$N = F(j) = 1 + b + b^2 + \dots + b^{j-1} = \frac{b^j - 1}{b - 1} \quad (2.4)$$

În literatură [Witt 81] se calculează intensitate traficului de mesaje pe o legătură de deasupra unui nod din nivelul k . Se consideră un timp de $N-1$ unități în care fiecare nod va trimite în medie câte un mesaj la celălalt. În acest timp fiecare nod din subarborele de sub nodul considerat - care are $j - (k-1)$ niveluri - va schimba o pereche de mesaje între ele și restul de noduri ale RIN. Prin urmare densitatea $T(k)$ a traficului pe linia de deasupra nodului de nivel k va fi :

$$T(k) = \frac{2F(j)(j-k+1)[N - F(j)(j-k+1)]}{N-1} \quad (2.5)$$

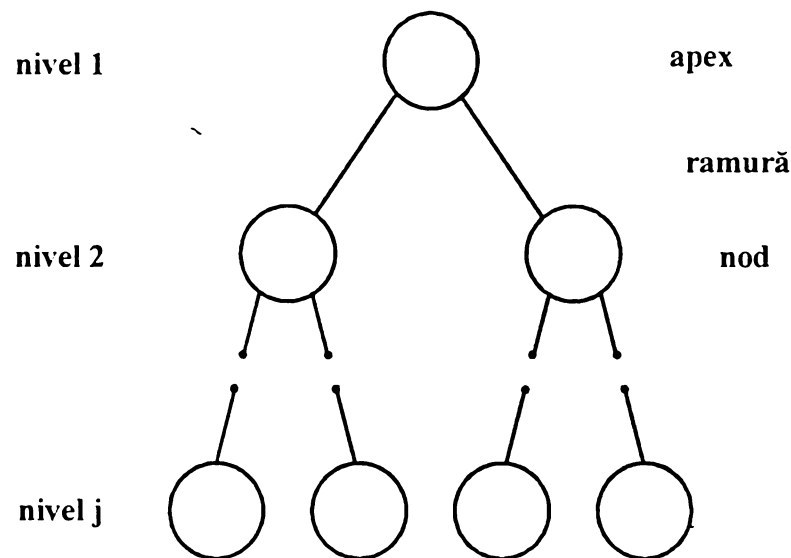


Fig. 2.19 RIN arbore

RIN plasă

O arie de procesoare bidimensională de tip plasă care interconectează vecinii cei mai apropiați (nearest neighbour mesh) se poate obține dacă fiecare nod se conectează prin câte o legătură la cei patru vecini, iar capetele libere pot fi recirculate pe partea opusă. Dacă pe o latură sunt W procesoare atunci numărul total de procesoare este $N = W^2$ iar în

cazul general al unei *plase regulate* într-un spațiu D - dimensional acest număr este $N = W^D$.

Distanța maximă de transmitere a unui mesaj în orice direcție este de $W/2$, iar lungimea medie de transmitere a mesajelor este $DW/4$. [Witt 81]. Deoarece fiecare nod se conectează la fiecare din cele D dimensiuni la câte doi vecini, există $2D$ legături la un nod.

Pentru RIN vor exista în total $2DN/2 = DN$ legături.

Densitatea medie a transferului pentru fiecare legătură este :

$$T = \frac{ND(W/4)}{ND} = \frac{W}{4} \quad (2.6)$$

Un nod al rețelei trebuie să dispună de $2D$ porturi de conectare la vecini. Dacă W crește cresc întârzierile în transmiterea mesajelor.

RIN hiper cub

Structura de tip *hipercub* (*hypercube*) se caracterizează prin prezența a $N = 2^D$ procesoare, interconectate sub forma unui cub binar D dimensional. Fiecare nod are legături directe și separate la alte D noduri vecine. Adresele a două noduri diferă printr-un bit. Există în total 2^D adrese binare asociate nodurilor.

Hipercubul generalizat (*generalized hypercube*) are lățimea W (noduri pe latură) dar păstrează adiacența adreselor. Numărul total de noduri este $N = W^D$ ca și la RIN *plasă*. Conectarea legăturilor necesită doar D porturi/nod în timp ce la *plasă* sunt necesare $2D$ /nod. Există în total ND/W laturi. Numărul mediu de pași pentru transmiterea unui mesaj este de $D(W-1)/W$ [Witt 81].

Densitatea traficului pe fiecare latură este :

$$T = \frac{ND(W-1)/W}{ND/W} = W - 1 \quad (2.7)$$

Extensia dimensiunii unui *hipercub* este mai dificilă decât a unei rețele *plasă*, deoarece creșterea numărului de noduri care folosesc o latură W este limitată de T admisibil, iar creșterea lui D presupune mărirea numărului de porturi la fiecare nod.

Hipercubul are proprietatea că el poate fi definit inductiv.

Un *hipercub* de ordinul 0 este un nod, în timp ce *hipercubul* de ordin $n+1$ este construit din două *hipercuburi* de ordin n și conectând nodurile lor corespunzătoare.

În fig. 2.20 se prezintă un *hipercub* de ordinul 4.

De notat că nodul unui *hipercub* constă dintr-un procesor, memoria sa locală, porturi de comunicare pentru alte noduri și porturi I/E.

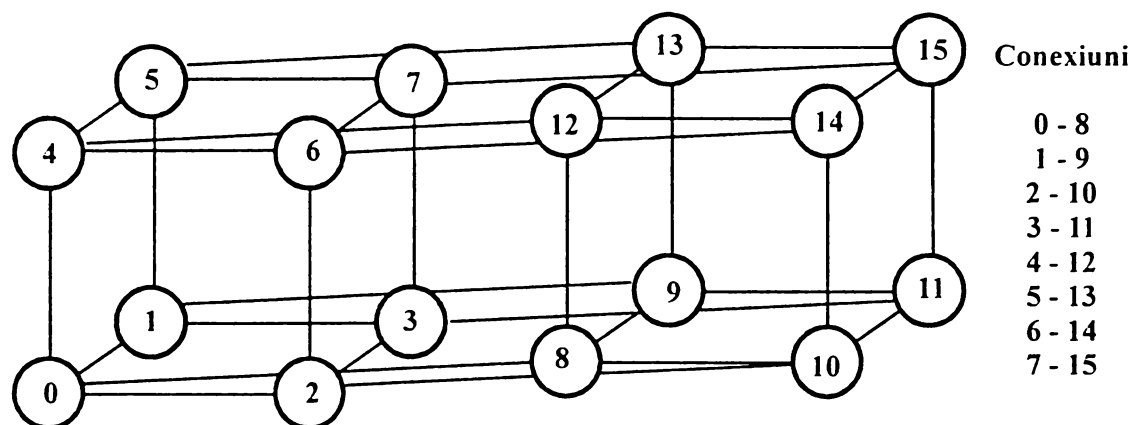


Fig. 2.20 Un hiper-cub cu 16 procesoare

RIN *hipercub* are câteva proprietăți interesante :

1. Cu cât numărul de procesoare crește, crește logaritmic și numărul de conexiuni și hardul asociat și prin urmare sisteme cu număr mare de procesoare sunt realizabile.
2. Un *hipercub* este un superset a altor RIN statice ca și *inel*, *arbore* etc. deoarece toate acestea pot fi transformate în *hipercub* ignorând unele conexiuni.
3. *Hipercuburile* sunt scalabile, o proprietate ce rezultă direct din proprietatea de recursivitate a lui.
4. Schemele de direcționare sunt simple. O politică de rutare a mesajului poate fi aceea prin care se trimite un mesaj la vecinul a cărui etichetă binară se potrivește cu eticheta destinației finale în următorul rang binar.

Lungimea căii de comunicație între două noduri este exact numărul de biți prin care etichetele diferă. Maximum este n iar media $n/2$.

Există numeroase căi posibile între două noduri; această redundanță se folosește pentru creșterea lui T și a toleranței la defecte.

RIN *hipercub* a fost utilizată la SMP- TM elaborate de Caltech, AMETER, Intel etc.

În literatură [Mold 93] se face o prezentare sistematică a evoluției acestei topologii până la Intel Touchstone Delta - o familie de supercalculatoare ce și-a făcut debutul în 1990.

Protocolul de rutare a mesajelor utilizează pentru prima dată tehnica *wormhole* propusă de Dally și Seity în 1987 și implementată la SMP-TM AMETEK 2010.

Rutarea *wormhole* diferă de comutarea de pachete ce utilizează tehnica *store-and-forward*. Pachetele sunt manipulate prin comutatoare speciale și nu software ca la majoritatea sistemelor bazate pe hipercub. Când pachetul din fruntea mesajului găsește un canal ocupat, pachetele din secvență sunt blocate. Pachetul din fruntea mesajului, uzual constând din câțiva bit, poziționează comutatoarele astfel încât restul pachetelor poate să urmeze fără a se consuma timp. Această tehnică de rutare reduce memoria asociată nodului și în același timp software-ul necesar pentru manipularea mesajelor.

2.4.4 Rețele de interconectare dinamice

Într-o RIN dinamică conexiunea între două noduri (procesor – memorie la SMP–MP și procesor – procesor la multicalculatoare) se stabilește prin intermediul unor comutatoare electronice denumite *bloc de comutare (switching box)* sau mai simplu, *comutatoare*. Schimbarea stării acestora poate face posibilă realizarea comunicației între diferite perechi de noduri. RIN care admit conexiuni simultane între toate combinațiile de intrări și ieșiri se numesc *rețele fără blocare*. Ele devin *strict fără blocare* dacă indiferent de conexiunile existente și fără a le perturba, întotdeauna se poate stabili o cale între oricare intrare nefolosită și oricare ieșire nefolosită. Unele rețele fără blocare permit efectuarea unei conexiuni I/E fără afectarea celor existente numai dacă acestea din urmă au fost înfăptuite pe baza unui *algoritm de rutare (routing algorithm)* prescris. Aceste RIN sunt considerate a fi *fără blocare într-un sens larg (wide sense nonblocking)*.

Pentru a reduce numărul de comutatoare majoritatea rețelelor sunt *rețele cu blocare*, adică nu fac posibile toate combinațiile de conexiuni I/E simultan. Dacă căile blocate pot fi totuși restabilite și redirecționate prin modificarea stării comutatoarelor, atunci rețeaua se numește *rearanjabilă*.

La o RIN dinamică multinivel (RIN-MN), cea mai complexă RIN, se disting următoarele elemente componente:

- *comutatorul*;

- *conexiunea* ce asigură realizarea funcției de interconectare (permutare);
- *nivelul de comutare* sau *etajul de comutare*.

Asociat acestora, se studiază *mecanismul de rutare a informației*.

Comutatorul

Comutatorul stabilește calea de conectare între o intrare și o ieșire. Complexitatea comutatorului depinde de tipul de RIN în care este utilizat: mică pentru RIN – *crossbar* și mare pentru RIN *multinivel*.

Literatura de specialitate [Dasg 89], [Hwan 93], [Popa 96], [BWK 98] ș.a. prezintă *in extenso* mai multe tipuri. Reținem trei de tipul:

patratic, arbitru, distribuitor, prezentate în fig. 2.21

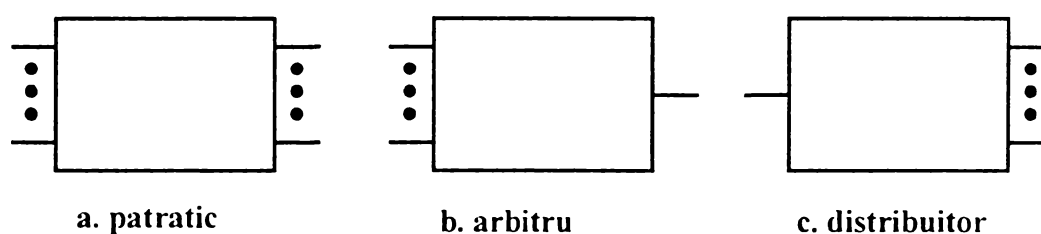


Fig 2.21 Comutatoare

Comutatorul de tip *patratic* are același număr de intrări și de ieșiri.

Comutatorul de tip *arbitru* (multiplexor) asigură selectarea uneia dintre intrări și transferul ei la unica ieșire.

Comutatorul de tip *distribuitor* alocă unei intrări o anumită ieșire.

Cel mai utilizat este comutatorul patratic cu două intrări și două ieșiri, 2 x 2. ce poate fi de două tipuri: cu două sau cu patru stări.

În fig. 2.22 se prezintă un astfel de comutator și cele patru stări posibile ale sale cu transfer:

- direct;
- în cruce;
- superior;
- inferior.

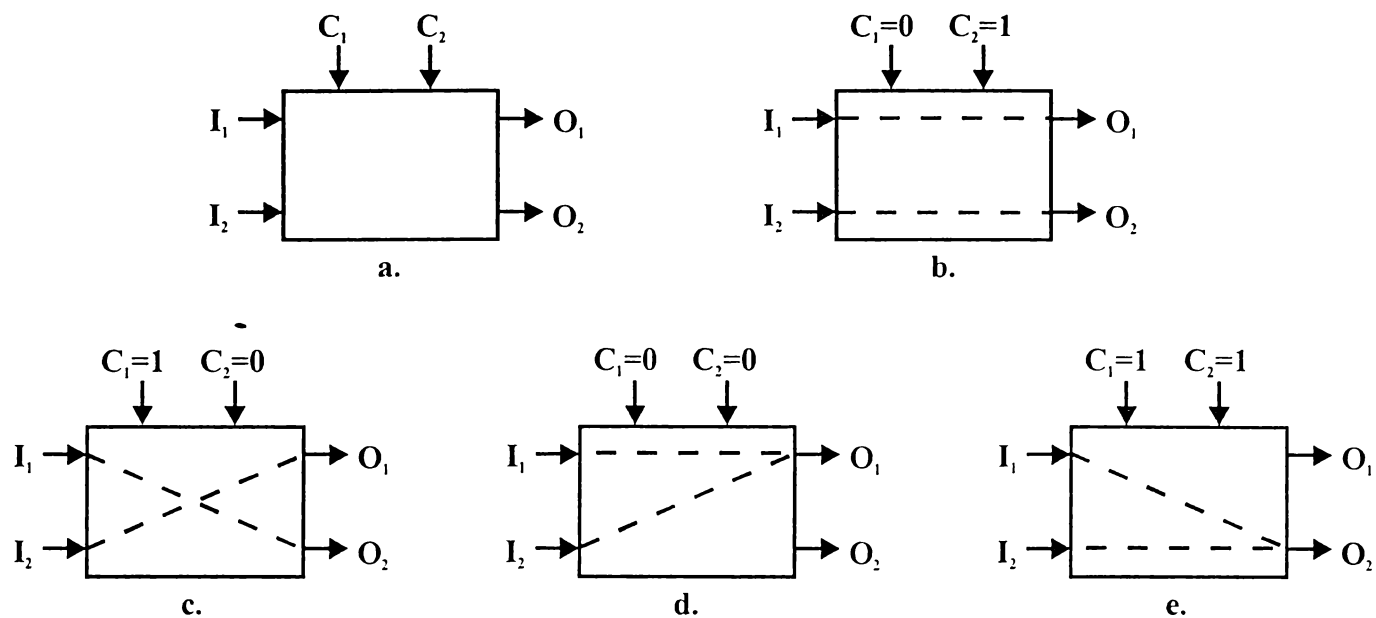


Fig. 2.22 Comutator 2 x 2

- reprezentarea comutatorului;
- transfer direct;
- transfer în cruce;
- transfer superior;
- transfer inferior.

Stările descrise la b) și c) nu sunt conflictuale, asigurând comutatorului atributul de *fără blocare* în timp ce stările prezentate la d) și e) sunt conflictuale deoarece ambele intrări concură pentru O_1 respectiv O_2 . În acest caz comutatorul este *cu blocare*.

Uzual, în cazul comutatoarelor fără blocare, se utilizează o singură intrare de comandă c care respectă specificațiile:

```

if  $c = 0$  then /* conectare directă */
    new  $O_1 = I_1$  &
    new  $O_2 = I_2$ 
else /* conectare încrucișată */
    new  $O_1 = I_2$  &
    new  $O_2 = I_1$  .
  
```

În acest caz comutatorul va avea doar două stări ca și în fig. 2.23.

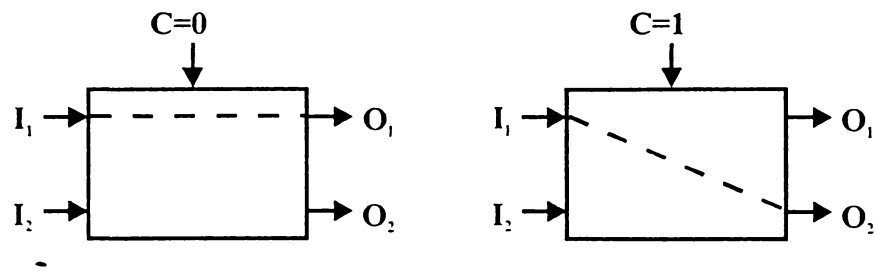


Fig. 2.23 Comanda unei intrări

Pentru evitarea situațiilor conflictuale prezentate mai sus se utilizează comutatoare ce au incorporate tamponare pentru memorarea pachetului ce trebuie transferat. Aceste tamponare sunt memorii FIFO cu o capacitate depinzând de intensitatea traficului și de lungimea mesajelor.

Funcții de interconectare

Funcțiile de interconectare realizează permutări ale ordinii procesoarelor.

Permutările sunt funcții bijective simple în care fiecare intrare este mapată în una și numai una dintre ieșiri.

Ele au un rol important în proiectarea RIN. Prin combinarea mai multor bijecții, se pot obține conexiuni complexe.

Când se execută o *funcție de rutare* f într-o RIN, datele de la nodul PE_i sunt transmise la nodul $PE_{j(i)}$. Această rutare a datelor apare simultan pentru toate nodurile active PE_i din set. Un nod inactiv PE poate să recepționeze date de la alt PE când se execută o funcție de rutare, dar el nu poate să transmită date. Este necesară asigurarea unei adrese la fiecare PE din matrice. Se pot utiliza mai multe scheme de adresare. Procesoarele pot fi adresate utilizând câmpuri de adresă de m bit, în acest caz $N = 2^m$. În alte situații când nodurile PE sunt dispuse într-un câmp bidimensional, se utilizează coordonate carteziene.

Fie n o adresă a unui nod PE exprimat printr-un câmp de adresă de m bit:

$$n = (b_m, b_{m-1}, \dots, b_1).$$

Interconectările din rețea pot fi descrise ca și permutări asupra reprezentărilor binare ale intrărilor.

În literatură [Mold 93] se descriu mai multe funcții de interconectare.

Schimb

Permutarea *schimb* (*exchange*) este definită ca:

$$E_k = (b_m, \dots, \bar{b}_k, \dots, b_1) \quad 1 \leq k \leq m \quad (2.8)$$

Structura RIN ce se obține aplicând $E_k(x)$ pentru $m = 3$ se numește *cub de ordinul 3*.

$$E_1(x) = \bar{b}_3 b_2 \bar{b}_1 \quad \text{corespunde legăturilor horizontale}$$

$$E_2(x) = b_3 \bar{b}_2 b_1 \quad \text{corespunde legăturilor oblice}$$

$$E_3(x) = \bar{b}_3 b_2 b_1 \quad \text{corespunde legăturilor verticale.}$$

În fig. 2.24 se prezintă conexiunile din cubul de ordin 3.

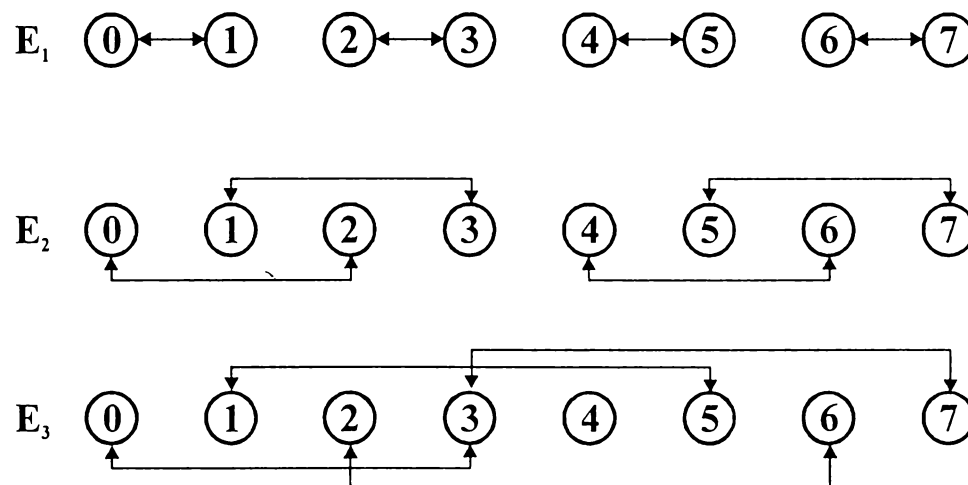


Fig. 2.24 Conexiuni în cubul de ordinul 3

În general într-o RIN bazată pe această funcție de permutare, cele N noduri sunt conectate prin $N \times 2^{N-1}$ legături. Ordinul k al cubului este $k = \log_2 N$.

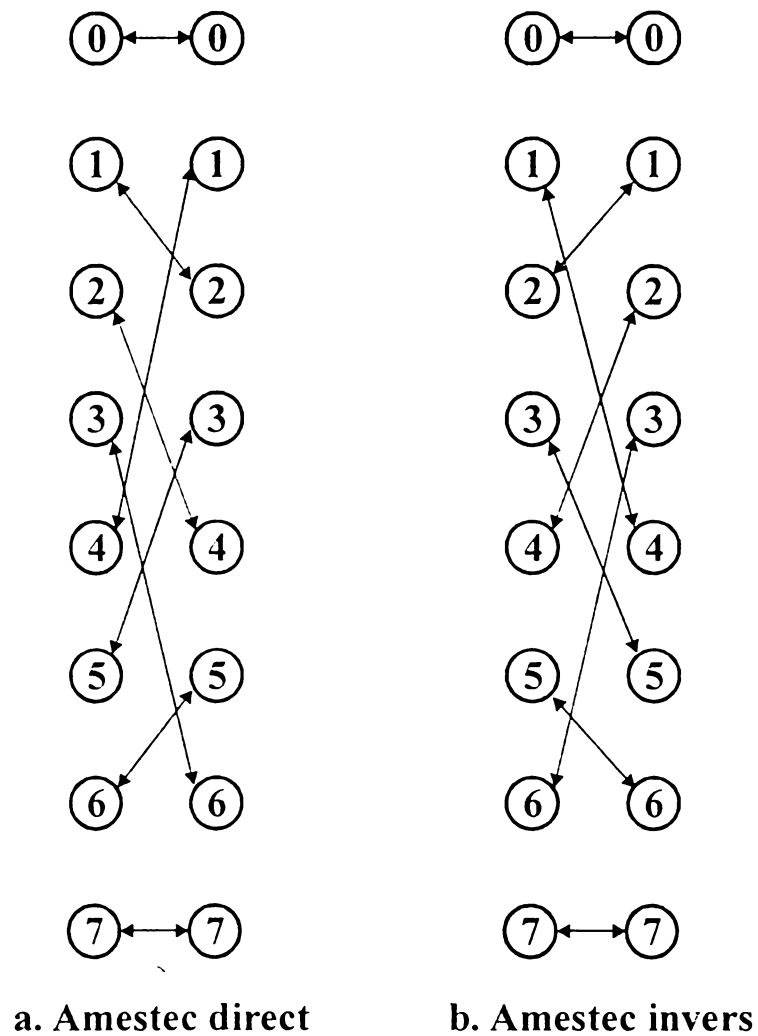
Amestec perfect

Permutarea *amestec perfect* (*perfect shuffle*) este definită ca:

$$S(x) = (b_{m-1}, b_{m-2}, \dots, b_1, b_m) \quad \text{- direct} \quad (2.9)$$

$$S'(x) = (b_1, b_m, \dots, b_3, b_2) \quad \text{- invers} \quad (2.10)$$

RIN bazată pe $S(x)$ în cazul $m = 3$ este prezentată în fig. 2.25 a iar pe $S'(x)$ în fig. 2.25 b.

Fig. 2.25 Interconectare *amestec perfect*

Se pot defini subseturi respectiv superseturi ale lui $S(x)$. A k -a subfuncție S_k este:

$$S_k(x) = (b_m, \dots, b_{k+1}, b_{k-1}, \dots, b_1, b_k) \quad (2.11)$$

iar a k -a suprafuncție $S^k(x)$ este:

$$S^k(x) = (b_{m-1}, \dots, b_{m-k+1}, b_m, b_{m-k}, \dots, b_1). \quad (2.12)$$

Fluture

Permutarea *fluture* (*butterfly*) este definită ca și:

$$B(x) = (b_1, b_{m-1}, \dots, b_2, b_m). \quad (2.13)$$

Interconectarea bazată pe această funcție, pentru cazul $m = 3$ este prezentată în fig. 2.26.

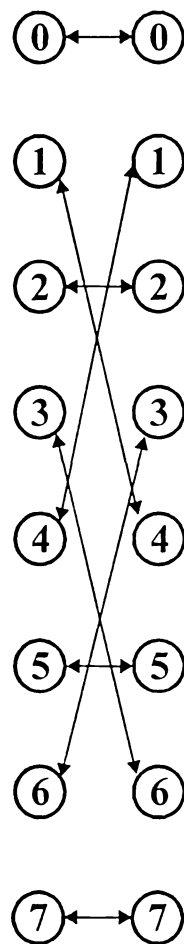


Fig. 2.26 Interconectarea fluture

Se pot defini în manieră similară:

$$B_k(x) = (b_m, \dots, b_{k+1}, b_1, b_{k-1}, \dots, b_k) \quad (2.14)$$

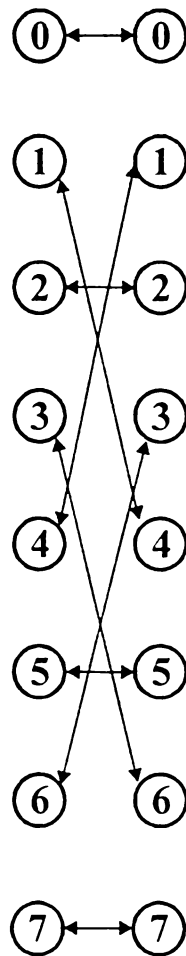
$$B^k(x) = (b_{m-k+1}, \dots, b_{m-k+2}, b_m, b_{m-k}, \dots, b_1) \quad (2.15)$$

Bit oglindit

Permutarea *bit oglindit* (*bit reversal*) se definește:

$$R(x) = (b_1, b_2, \dots, b_m). \quad (2.16)$$

Pentru cazul $m = 3$ conexiunea este prezentată în fig. 2.27.

Fig. 2.27 Interconectare *bit-oglindit*

Rețeaua Illiac

Aceasta este utilizată în calculatorul Illiac IV și se definește astfel:

$$\begin{aligned}
 R_{+1}(x) &= (x + 1) \bmod N \\
 R_{-1}(x) &= (x - 1) \bmod N \\
 R_{+r}(x) &= (x + r) \bmod N \\
 R_{-r}(x) &= (x - r) \bmod N.
 \end{aligned}
 \tag{2.17}$$

Pentru cazul în care $N = 16$ și $r = 4$, această rețea este prezentată în fig. 2.28.

Ea este o RIN parțial conectată. Când $N = 16$, patru noduri *PE* pot fi atinse într-un pas, șapte în doi pași și unsprezece *PE* în trei pași.

În general sunt necesari $I \leq \sqrt{N} - 1$ pași pentru a conecta fiecare nod *PE*.

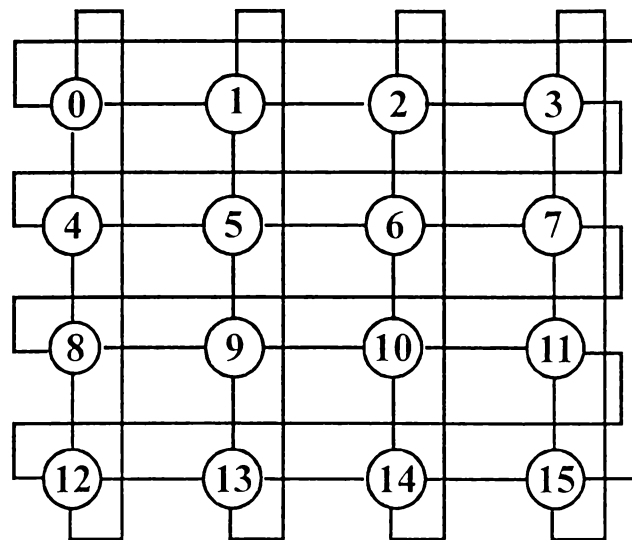


Fig. 2.28 Reprezentarea rețelei Illiac

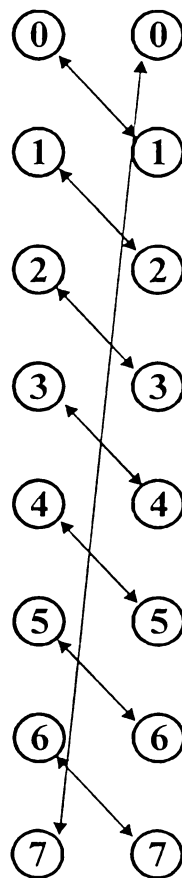
Deplasare

Permutarea *deplasare* (*shift*) este definită ca și:

$$SH(x) = |x + 1|_N \quad (2.18)$$

unde $| \cdot |_N$ înseamnă *modulo* N .

Interconectarea bazată pe deplasare, în cazul $m = 3$, este prezentată în fig. 2.29.

Fig. 2.29 Interconectarea bazată pe *deplasare*

Se pot determina realții diferite între permutări.

De exemplu pentru $m = 3$,

$$R_3 = B_3.$$

Pentru anumite interconectări, ca de exemplu la rețeaua Illiac, este preferabil să se utilizeze direct indexul fiecărui PE_x , pentru $0 \leq x \leq N - 1$. Conexiunile de la PE_x sunt definite de către elementul de procesare destinație.

Mecanisme de rutare a informației

Prin *mecanism de rutare a informației* se înțelege totalitatea procedurilor / algoritmilor utilizate pentru controlul transferului informației de la un nod la următorul.

Un astfel de mecanism trebuie să asigure:

- integritatea informației
- calea de date optimă între transmițător și receptor
- evitarea saturării căilor de date
- găsirea căii alternative în cazul defectării unei căi de date.

Alegerea mecanismului de rutare este impusă de caracteristicile generale ale RIN.

În continuare se descriu câteva mecanisme de direcționare:

Rutarea prin saturare. Dacă un nod primește un mesaj și dacă nu-i este destinat îl va transmite simultan către toți vecinii, cu excepția transmițătorului. Avantajul principal îl reprezintă simplitatea iar dezavantajul major creșterea traficului din care o mare parte se datorează metodei.

Rutarea aleatoare constă în modul aleator în care se aleg căile de date. Nodul care primește un mesaj ce nu îi este destinat îl va transmite în continuare după o anumită lege de distribuție. Avantajele și dezavantajele sunt ca și la metoda anterioară.

Rutarea fixă presupune generarea unei *matrici de rutare* pentru o topologie dată.

Matricea se construiește din atâtea linii și coloane câte noduri sunt. la intersecția unei linii cu o coloană găsiindu-se acel nod care este intermediarul între nodul care corespunde liniei și cel care corespunde coloanei.

Pentru RIN dinamice, în literatură [Hwan 93], [Popa 96], [Mold 93] se prezintă mai multe mecanisme:

- *Utilizarea adresei destinației pentru stabilirea căii de comunicație.*

Procesorul transmițător generează o *etichetă de rutare (routing tag)* care are un număr de biți egal cu numărul de niveluri. Fiecare rang controlează câte un nivel. La intrarea fiecărui nivel j , comutatorul verifică al j -lea bit al etichetei, începând cu primul nivel, analizând c.m.s. bit. Dacă bitul este 0, comutatorul selectează portul superior; dacă bitul este 1 se selectează portul inferior.

- *Calcularea funcției SAU-EXCLUSIV între adresele sursei și destinației.*

Combi-nația rezultată va controla poziționarea comutatoarelor. Dacă adresele sursei și destinației sunt:

$$a_n a_{n-1} \cdots a_i \cdots a_1 \text{ respectiv } b_n b_{n-1} \cdots b_i \cdots b_1$$

atunci rezultatul este:

$$c_n c_{n-1} \cdots c_i \cdots c_1$$

unde $c_i = a_i \oplus b_i$, controlează nivelul i .

- *Utilizarea unei funcții de control.*

Această funcție are atâtea ranguri câte niveluri de comutatoare există. Valoarea 0 a unui rang va stabili conexiunea directă iar valoarea 1 conexiunea în cruce. Generarea funcției de control este sarcina modulului ce asigură controlul întregii rețele.

A. RIN dinamice pe un singur nivel.

O RIN dinamică pe un singur nivel este formată din:

- *conexiunile de realizare a funcției de interconectare;*
- *nivelul de comutare conținând comutatoare.*

Schema bloc este prezentată în fig. 2.30.

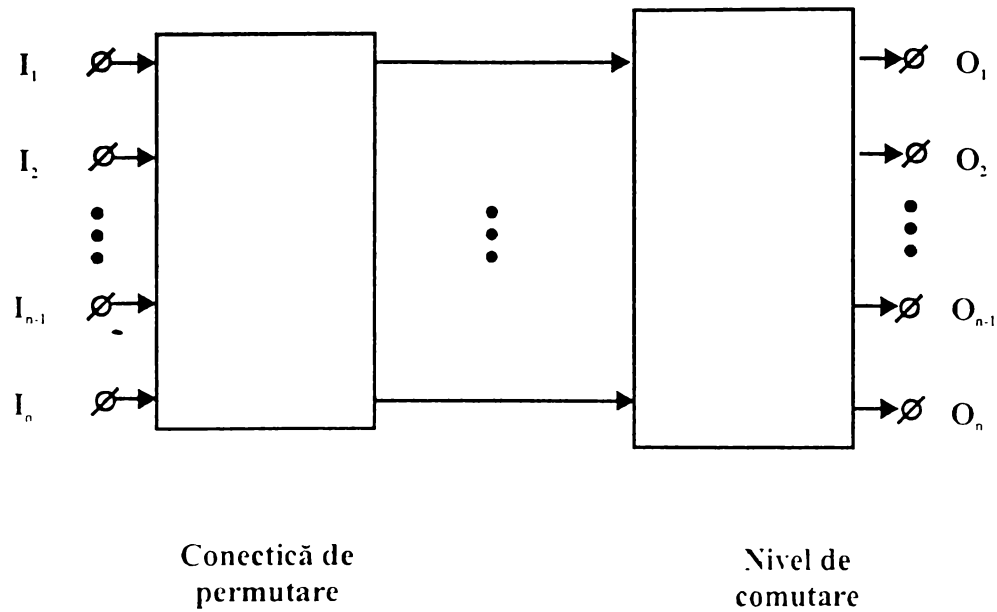


Fig. 2.30 Schema bloc a unui RIN dinamic *cu un singur nivel*.

Un exemplu clasic este aceea a RIN bazată pe funcția de interconectare *amestec perfect* și varianta *amestec perfect invers*.

Schemele celor două variante de RIN sunt prezentate în fig. 2.31 a și b.

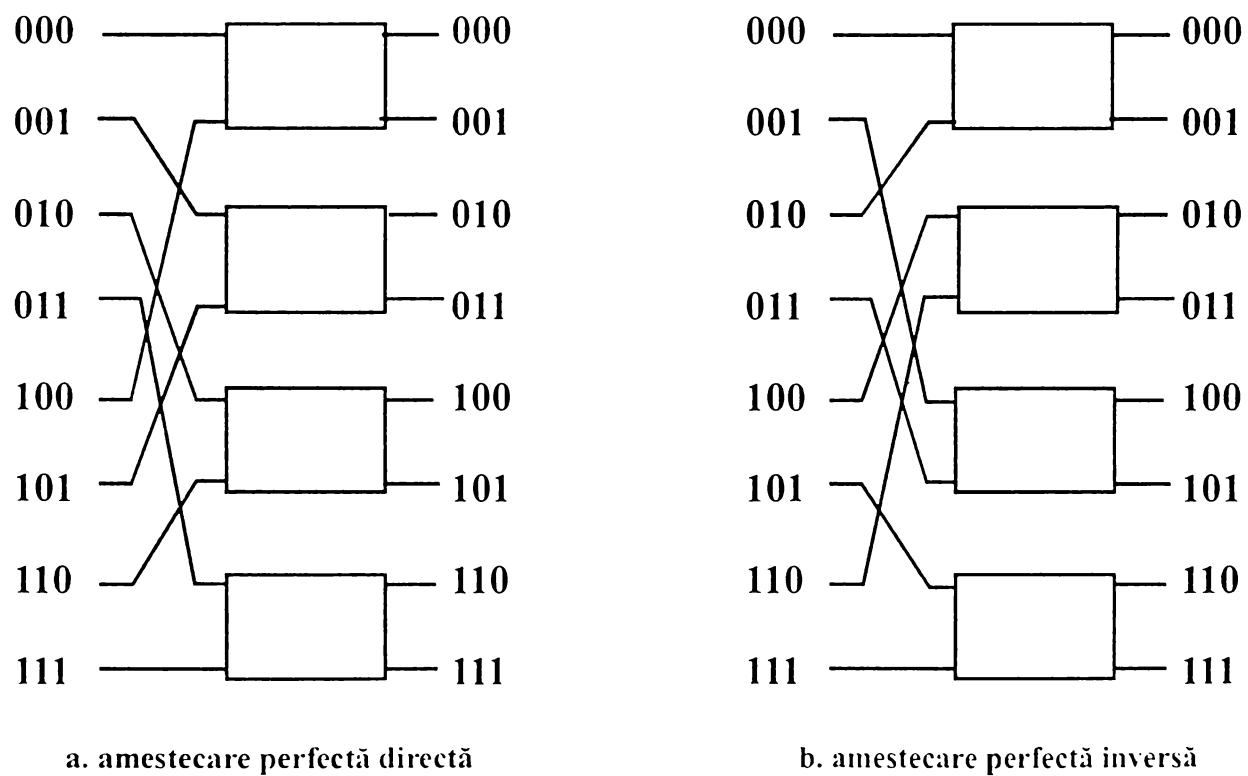


Fig. 2.31 RIN ce realizează funcția *amestec perfect*

RIN dinamice pe un *singur nivel* au o complexitate medie dar oferă posibilitați limitate de conectare, numărul permutărilor posibil la o trecere fiind mic. Conectivitatea totală se poate obține recirculând informația prin rețea, de unde și denumirea de RIN cu *recirculare*.

Datorită faptului că la fiecare permutare noua adresă se obține din precedenta printr-o rotire rezultă că numărul maxim de recirculări este egal cu $\log_2 m$ unde m reprezintă numărul de ranguri al adresei.

RIN *crossbar*

Este o RIN dinamică pe un singur nivel, deplin conectată (fully connected). Este o rețea fără blocare permițând conectarea simultană a N elemente de procesare PE la M memorii. Accesele simultane la toate modulele de memorie conferă RIN un mare debit de prelucrare a informației. Situația conflictuală poate apare atunci când mai mult elemente de procesare solicită același modul de memorie.

Dacă $N = M$ atunci se utilizează N^2 comutatoare, ceea ce face costul $O(N^2)$ prohibitiv față de $\frac{N}{2} \log N$ a unui RIN *omega*.

Un PE este format dintr-un procesor P cu resurse locale de memorie ML conectate în jurul unei magistrale locale. RIN *crossbar* se utilizează în SMP–MP cu număr redus de procesoare. Un exemplu îl constituie sistemul C.mmp [Dasg 89] care utilizează 16 module procesoare PDP 11 și 16 module de memorie partajate având o capacitate totală de 32 Mbyte.

Structura unei RIN *crossbar* de dimensiuni $N \times M$ este prezentată în fig. 2.32.

Intr-o astfel de rețea funcția de comutare și comanda sa sunt distribuite de-a lungul celor $N \times M$ comutatoare. Este posibilă concentrarea funcției de comutare la fiecare dintre cele M module de memorie, rezultând un *sistem cu memorie multiport*.

Structura generală a unui astfel de sistem este prezentată în fig. 2.33.

Fiecare modul de memorie are n porturi, un port / procesor. Comutatorul asociat fiecărui modul de memorie conține și logica de arbitrare.

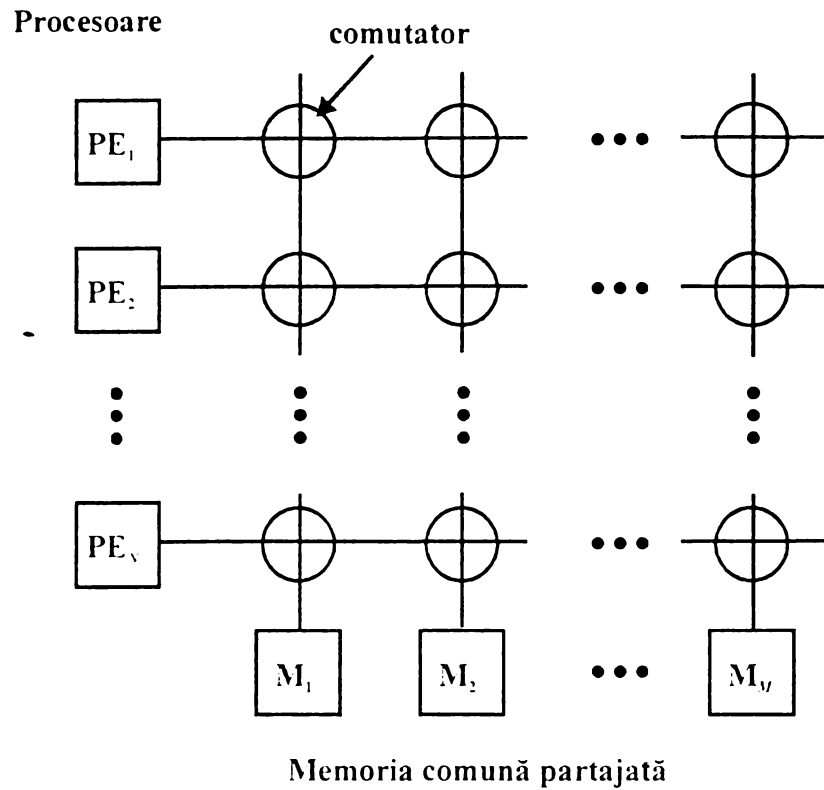
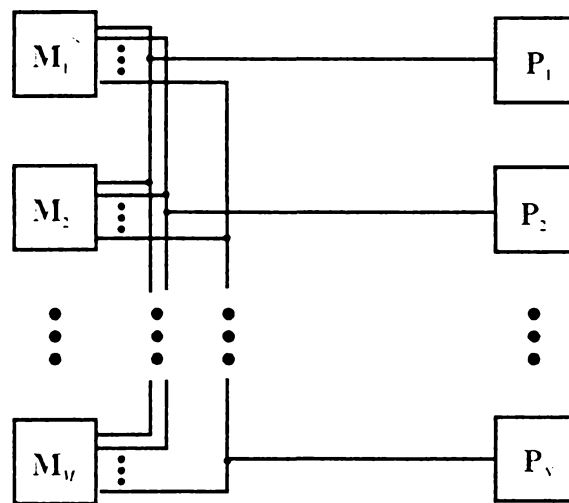
Fig. 2.32 RIN *crossbar*

Fig. 2.33 Memorii multiport

B. RIN dinamice multinivel

Asigură un bun compromis performanță / cost permițând stabilirea unei căi de comunicație între intrare și ieșire într-un timp mai scurt decât cel obținut într-o RIN dinamică mononivel cu N intrări și N ieșiri, există $\log_2 N$ niveluri de comutare 2×2 .

Pe un nivel numărul de comutatoare este $N/2$.

Shcema bloc a unui RIN dinamic *multinivel* este prezentată în fig. 2.34.

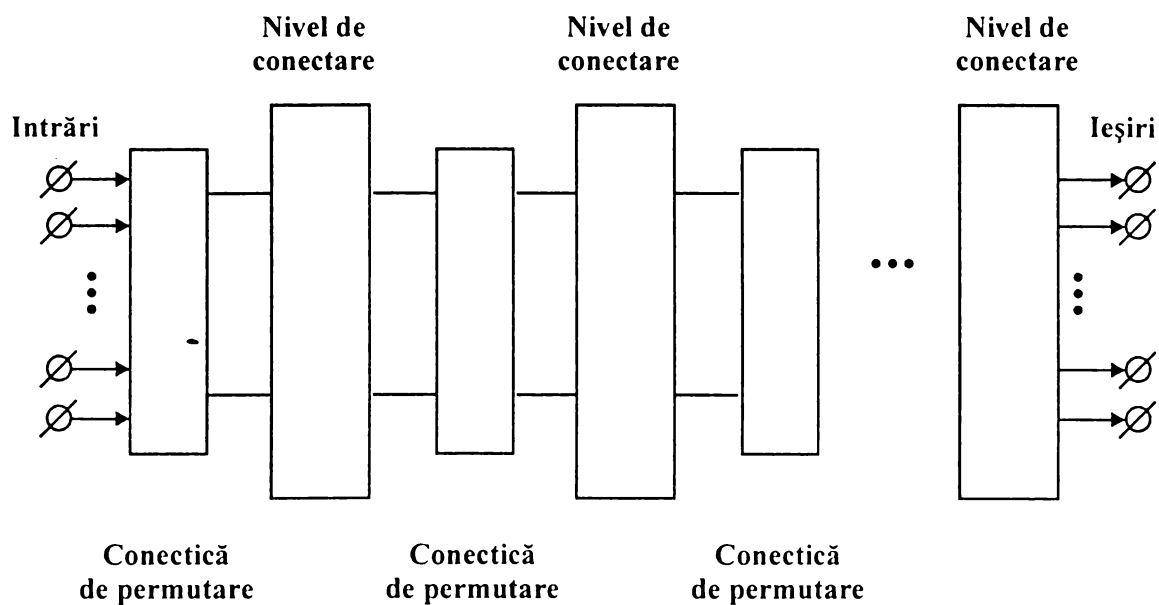


Fig. 2.34 Schema bloc a unui RIN dinamic *multinivel*

Diferențe importante apar între RIN dinamice *pe un singur nivel* și RIN dinamice *multinivel*:

- O RIN dinamică *multinivel* poate fi utilizată atât la transferul prin comutare de circuite cât și la cel prin comutare de pachete. O RIN *pe un singur nivel* nu poate transfera decât pachete.
- O RIN dinamică *multinivel* are viteza de transfer mai mică deoarece informația parcurge mai multe niveluri înainte de a ajunge la destinație, pe când la o RIN *pe un singur nivel*, informația sau părți din informație ajung la destinație după fiecare iterație.

RIN *omega*

Această RIN a fost studiată de D. Lawrie și a fost utilizată în multiprocesorul Cedar și o variantă ameliorată, la mașina NYU Ultracomputer. Este descrisă pe larg în [Hwa 93], [Mold 93], [MHB 86] ș.a.

În varianta sa inițială conectează N procesoare la N memorii.

Se caracterizează prin:

- lățime de bandă proporțională cu N ;
- pentru $N = 2^k$ există $k = \log_2 N$ niveluri identice;
- între două niveluri adiacente conexiunea este de tip *amestec perfect*;
- $N/2$ comutatoare pe fiecare nivel;
- $N \log_2 N$ comutatoare per ansamblu;
- comandă independentă a comutatoarelor pe fiecare nivel;

- latența, logaritmică în N .

Schema unei RIN *omega* 8 x 8 este prezentată în fig. 2.35.

Ea constă din trei niveluri fiecare conținând patru comutatoare 2 x 2.

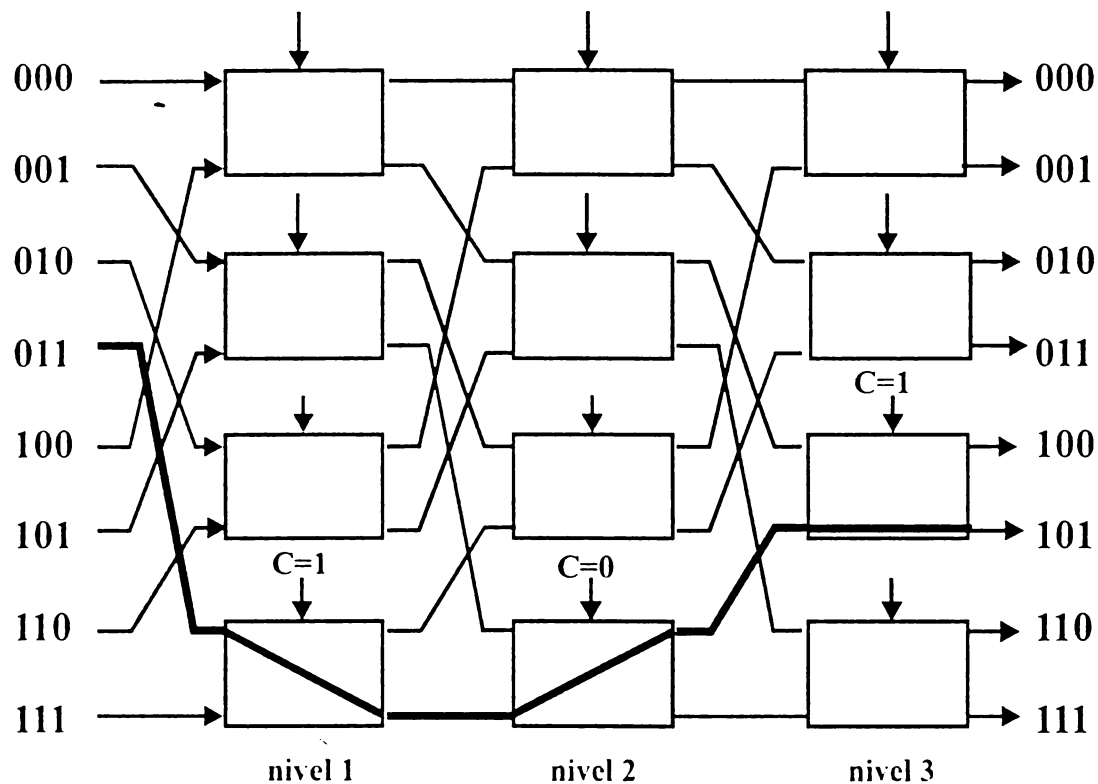


Fig. 2.35 RIN *omega* conectând intrarea 3 cu ieșirea 5

Procesorul conectat la intrarea 3 (011) generează o etichetă care este reprezentarea binară a destinației, în cazul nostru 5 (101). Conectarea unui comutator dispus pe nivelul i este acționat de al i -lea bit al etichetei controlizat de la C.m.s.b. În figură se vede traseul parcurs de un mesaj de la intrarea 3 la ieșirea 5.

În general considerându-se că se transmit date de la dispozitivul de intrare având numărul I la dispozitivul de ieșire având numărul O și că reprezentarea binară este:

$$I = i_1 i_2 \dots i_k \quad \text{și} \quad O = o_1 o_2 \dots o_k.$$

Se poate formaliza astfel:

```

For  $j = 1 \dots k$  do
    If  $O_j = o$  then pune pe 0,  $c$  de pe nivelul  $j$  conectat la  $I$ 
    else pune pe 1,  $c$  de pe nivelul  $j$  conectat la  $I$ 
endfor

```

Reamintim că la $c = 0$ se conectează intrarea cu ieșirea superioară iar la $c = 1$ cu ieșirea inferioară în cadrul comutatorului.

Inițial, dispozitivul I este conectat la un comutator dispus pe nivelul 1 a RIN. Apoi algoritmul de direcționare permite comutatoarelor relevante de pe nivelurile 1, 2, ..., k să fie poziționate succesiv astfel ca transmisia să fie efectuată de la I la O.

În RIN *omega* există doar o singură cale între orice intrare și orice ieșire. Prin urmare pentru o mapare dată a unui set de intrări la un set de ieșiri, algoritmul va produce un set unic de căi de comunicație pentru a realiza această mapare. Există un potențial conflictual deoarece o anumită pereche de astfel de căi poate să intre în conflict cu alta, deoarece se solicită simultan două stări diferite ale aceluiași comutator.

De exemplu maparea intrărilor (0, 4) la ieșirile (4, 7) va produce un astfel de conflict.

Formalizând, o RIN *omega* nu poate produce toate mapările posibile de la intrare la ieșire. În literatură [Hwan 93] se citează câteva clase de mapări ce pot fi realizate cu RIN *omega*.

În fig. 2.36 se prezintă situația conflictuală menționată anterior.

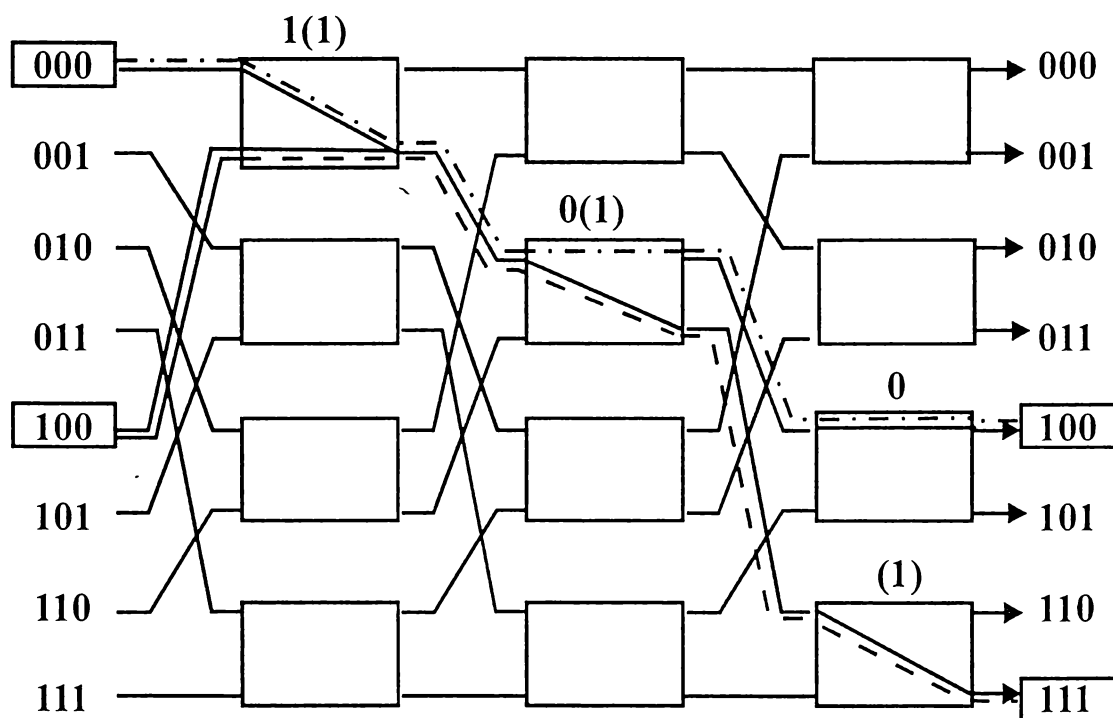


Fig. 2.36 Situația conflictuală la RIN *omega* pentru conexiunile $0 \rightarrow 4$ și $4 \rightarrow 7$.

RIN *delta*

O RIN *delta* poate conecta $N=a^n$ intrări la $M=b^n$ ieșiri prin n niveluri de comutatoare *crossbar* $a \times b$. *Rețeaua cu amestec generalizat (generalized shuffle network) (GSN)* este capabilă să conecteze orice $M = m_1 * m_2 * \dots * m_r$ intrări la

$N = n_1 * n_2 * \dots * n_r$ ieșiri prin intermediul a r niveluri de comutatoare. Al i -lea nivel utilizează $m_i \times n_i$ comutatoare *crossbar* și este precedat de o GSN care în esență este un supraset a interconectărilor *omega* și *delta*.

Aceasta este cea mai generalizată versiune a unei RIN–MN ce permite talii diferite în ceea ce privește numărul de intrări și ieșiri. Toate celelalte RIN dinamice pot fi obținute din această configurație alegând în mod corespunzător m_i și n_i . De exemplu când $m_i = a$ și $n_i = b$ pentru toți i , rezultă o RIN *delta*. Când $m_i = n_i = 2$ pentru toți i , rezultă o RIN *omega*. Dacă $r = 1$ se obține RIN *crossbar*.

Când $M = M * 1$ și $N = 1 * N$ se obține o *magistrală partajată*.

În fig. 2.37 se prezintă structura unei RIN *delta* particularizată pentru $n = 2$.

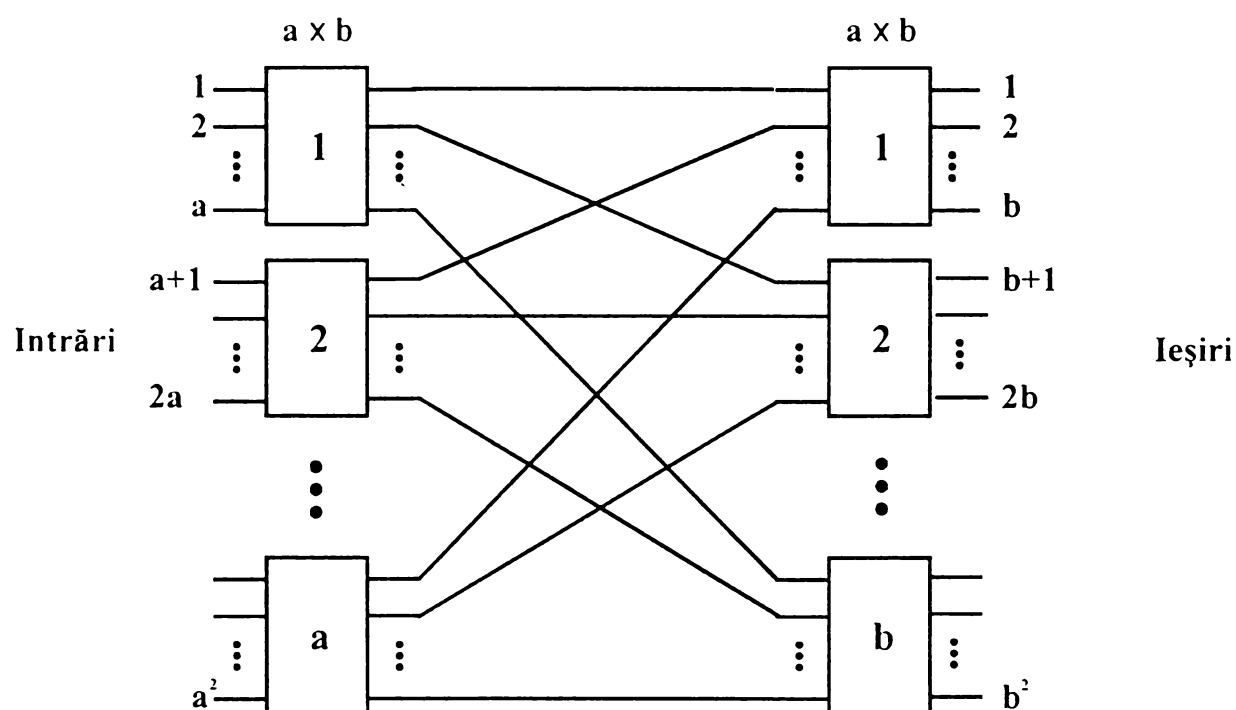


Fig. 2.37 RIN *delta* pentru $n = 2$

Un dezavantaj serios al RIN – MN îl constituie existența unei singure căi de date între o intrare și o ieșire. A fost necesar să se incorporeze resurse ce să facă RIN–MN tolerantă la defecte. Ele vor fi analizate ulterior.

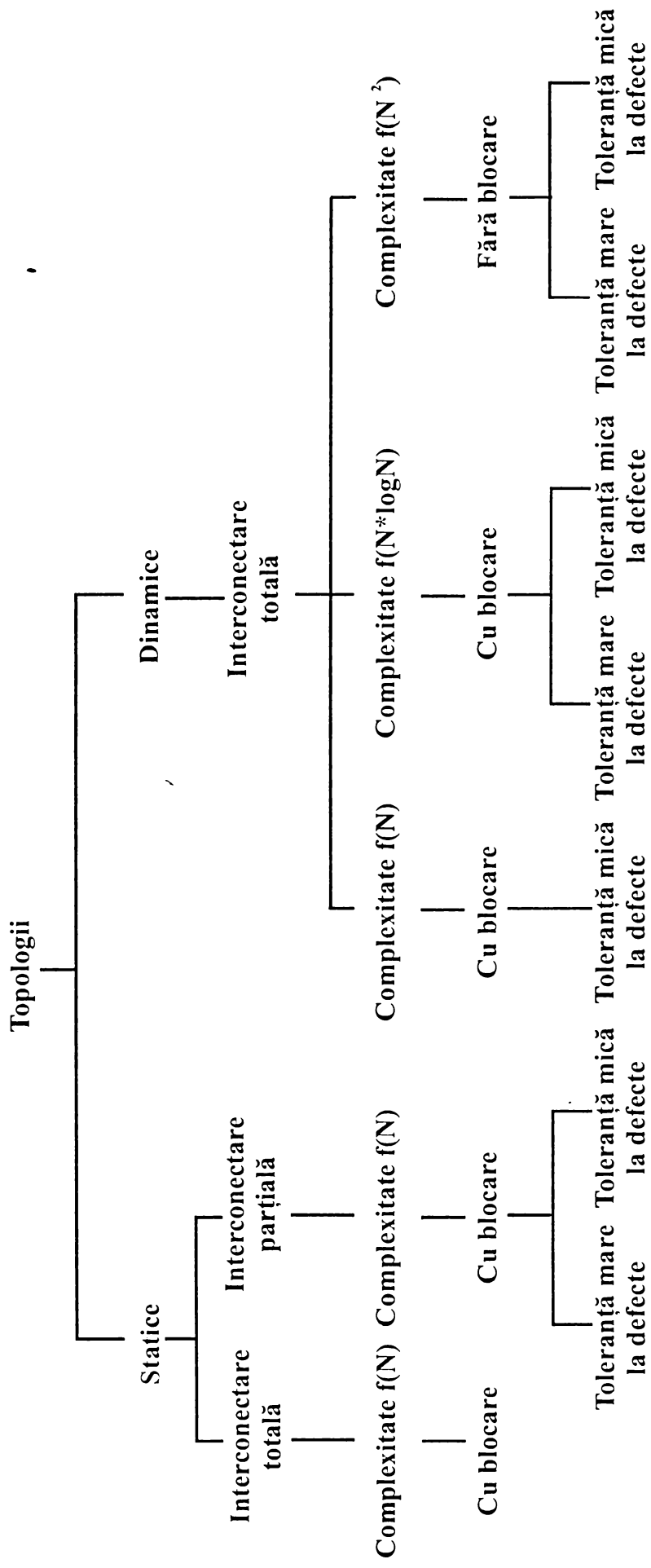


Fig. 2.38 Clasificarea topologiilor RIN

2.5 MECANISME DE SINCRONIZARE A PROCESELOR LA ARHITECTURI DIN CLASA MIMD

Arhitecturile MIMD trebuie să conțină capabilități pentru facilitarea comunicării între procese. O axiomă fundamentală ce guvernează problematica comunicării între procese este următoarea:

- *Nu se poate face nici o presupunere referitoare la vitezele relative ale proceselor*

Prin urmare mecanismele de sincronizare sunt programate în interiorul proceselor și nu sunt impuse de procesoare.

Problematica comunicării între procese se reduce astfel la problematica sincronizării proceselor.

Interacțiunea dintre procese paralele se manifestă în sincronizarea lor în scopul de-a asigura:

1. *cooperarea* proceselor pentru a facilita comanda secvenței ce asigură rezultate corecte
2. *competiția* proceselor pentru accesul la resursele partajate prin furnizarea comenzii accesului pentru acestea

În literatură [Dasg 89], [Mold 93] se analizează diferite mecanisme de sincronizare pentru arhitecturi MIMD-T și MIMD-L, *in extenso* pentru SMP - MP și respectiv SMP - TM - multicalculatoare.

Tradițional aceste mecanisme se bazează pe utilizarea:

- *variabilelor partajate* - utilizate la SMP - MP;
- *transferurilor de mesaje* - utilizate la SMP - TM.

2.5.1 Moduri de sincronizare

În fig.2.39 se prezintă situația în care două procesoare P_1 și P_2 execută procesele P_1 și P_2 . P_1 și P_2 sunt două procese independente ce sunt procesate în paralel. În anumite puncte ale execuției lor, atât P_1 cât și P_2 solicită modificarea structurii comune de date D rezidentă în memoria comună M .

Pentru a preveni apariția *nedeterminării funcționale* - situație în care D este modificată de o manieră incontrolabilă - este necesară tratarea segmentelor de cod CS_1 și CS_2 ca fiind *regiuni critice*.

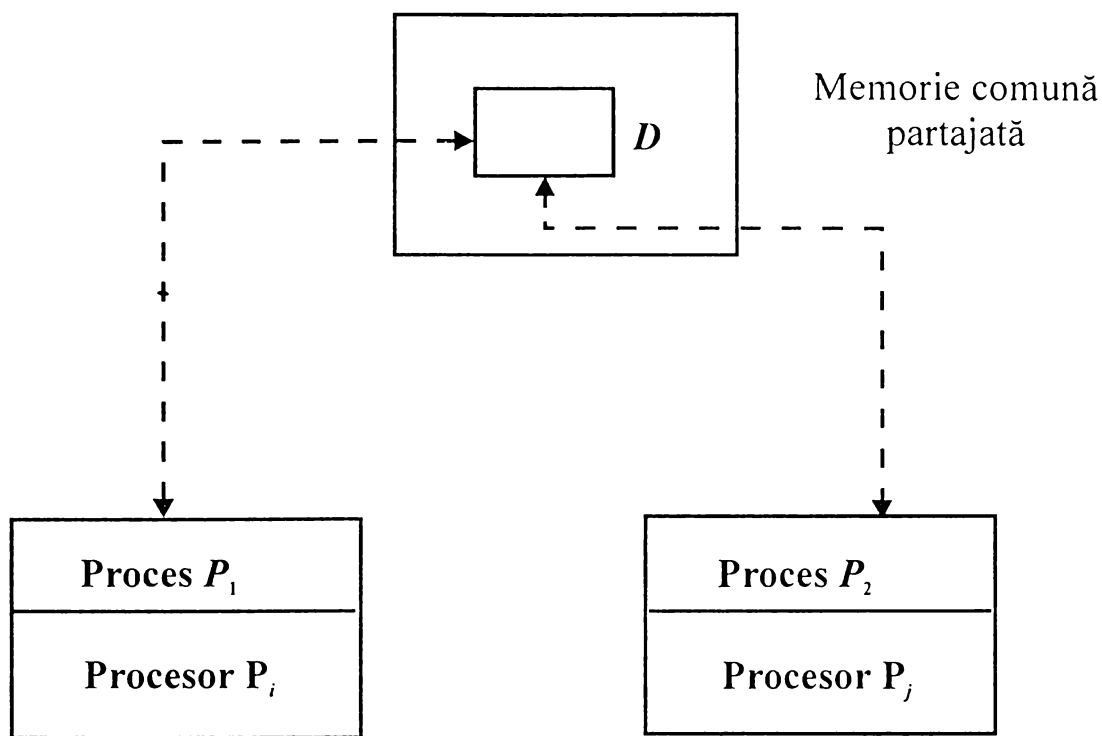


Fig. 2.39 Procese partajând o structură de date D

O *regiune critică* este o secțiune de cod în cadrul unui sistem concurrent ce este astfel controlată încât la un moment dat doar un proces poate intra în ea.

Regiunile critice trebuie tratate ca fiind unități *indivizibile* iar modificarea structurilor de date se face într-o manieră de *excludere mutuală*.

În 1968 Dijkstra a tratat sistematic problemele sincronizării în regiunile critice devenind astfel un clasic.

Actualmente se preferă utilizarea termenului de *secțiune critică* în locul celui de *regiune critică*.

Implementarea celor două cerințe simultan obligatorii - *indivizibilitatea* și *excluderea mutuală* - pentru o secțiune critică se face prin intermediul unui mecanism de *blocare* format din operații indivizibile (atomice) și o "cheie" k ce poate fi *fanion*, *semafor* etc.

În cazul prezentat ca și exemplu, structurile de date D i se asociază "cheia" k și operațiile atomice LOCK (k) și UNLOCK (k) ce poziționează k pe 1 respectiv pe 0.

Examinarea și poziționarea lui k reprezintă o operație atomică în sensul că ea se efectuează pe durata unui singur impuls de orologiu și nu poate fi întreruptă.

Structura unei astfel de operații atomice este:

Loop: [if $k = 0$ then $k = 1$ else] go to Loop;

unde [] indică indivizibilitatea.

Prin urmare un proces poate să execute LOCK (k) atunci când este în starea deblocată, adică când $k = 0$. În final, după efectuarea lui LOCK (k), $k = 1$ și toate celelalte procese sunt excluse de la efectuarea acestei operații. Ele trebuie să aștepte într-o stare de *busy waiting*. Operația UNLOCK (k) efectuează $k = 0$.

Secțiunile critice în P_1 și P_2 pot fi implementate astfel:

P_1 : LOCK (k) secțiune critică CS_1 UNLOCK (k)	P_2 : LOCK (k) secțiune critică CS_2 UNLOCK (k)
---	---

P_1 și P_2 au fost sincronizate în CS_1 și CS_2 într-o manieră de excludere mutuală.

În situația *sincronizării condiției* se utilizează aceleași mecanisme de blocare ca și la excluderea mutuală. Se asociază două chei, k_1 și k_2 , structurii de date D . P_1 deblochează k_2 pentru a semnală lui P_2 că a depus o nouă valoare în D în timp ce P_2 deblochează k_1 pentru a informa P_1 că a consumat ultima valoare plasată în D . Presupunând că k_1 și k_2 sunt inițial poziționați pe 0 și 1, segmentele relevante de cod în P_1 și P_2 sunt:

P_1 : LOCK (k_1) produce valoare și pune în D UNLOCK (k_2)	P_2 : LOCK (k_2) consumă valoare din D UNLOCK (k_1)
--	---

Aceste mecanisme de blocare se realizează prin intermediul *primitivelor de sincronizare* ce utilizează variabile partajate.

În general se face distincție între două mecanisme diferite de sincronizare:

- *sincronizare la nivel de comandă* - se utilizează întreruperile procesorului și comenzi START, RESET;
- *sincronizări la nivel de date* - utilizată ori de câte ori una sau mai multe variabile dintr-un task sunt partajate;

2.5.2 Sincronizarea bazată pe variabile partajate

Implementarea operațiilor LOCK/UNLOCK trebuie să îndeplinească două condiții:

1. Operațiile să fie suficient de primitive ca ele însuși să fie atomice;
2. Să se cunoască condiția de deblocare a procesului blocat.

Primitivele de sincronizare sunt:

- TEST - AND - SET (TAS)
- FETCH - AND - ADD (FAA)
- COMPARE - AND - SWAP (CAS)

Primitiva TEST - AND - SET (TAS)

Este utilizată pentru coordonarea sau sincronizarea unor procese cât și pentru crearea altora noi. TAS asigură excluderea mutuală prin executarea unei instrucții ca și o operație neîntreruptibilă, adică pe durata unui singur impuls de ceas.

Sintaxa este:

$$C := \text{TEST - AND - SET } (k)$$

unde:

k - variabilă de memorie

C - variabilă locală de procesor

Semantica este:

$$\text{TEST - AND - SET } (\tilde{k})$$

$$\{ temp \leftarrow k; k \leftarrow 1; \text{RETURN } k; \}$$

Acces la variabila partajată C

$$\text{RESET } (k)$$

$$\{ k \leftarrow 1; \}$$

Operația LOCK poate fi acum implementată ca și:

$$\text{LOCK } (k) : \text{repeat } C := \text{TEST - AND - SET } (k) \text{ until } C = 0$$

iar UNLOCK are forma:

$$\text{UNLOCK } (k) : k := 0$$

La IBM 360/370 această operație se regăsește sub forma instrucției TS iar la i386 ca și TAS.

Inconveniente majore ale primitivei TAS sunt:

- procesorul P_i în care P_1 așteaptă un LOCK se găsește în starea *busy-waiting* în care nu efectuează nimic;

- fiecare procesor aflat în starea *busy-waiting* consumă cicluri de memorie utilizabile mai eficient în alt mod.

Primitiva de tipul FETCH - AND - ADD (FAA)

Scopul operațiilor de sincronizare este de a *serializa* două sau mai multe segmente de cod paralele implicând o structură de date partajate D .

În timp ce FAA se bazează pe ideea utilizată în TAS, modul ei de operare este diferit. Ca răspuns la n accese simultane, FAA acționează dând fiecărui proces un număr unic și apoi le extrage serial.

Se obține astfel un anumit grad de simultaneitate a acceselor la o variabilă partajată în timp ce se păstrează cerința de serializare.

Principiul *serializării* este: dacă x reprezintă o variabilă partajată și mai multe operații FAA adresează simultan pe x , efectul acestor operații este exact acela ce s-ar fi obținut dacă aceste operații ar fi apărut într-o ordine serială nespacificată.

Forma operației FAA este:

FETCH - AND - ADD (x, e)

unde:

x - variabila partajată

e - expresie întreagă

Semantica este:

FETCH - AND - ADD (x, e)

$$\{ temp \leftarrow x, x \leftarrow temp + e; \\ \text{RETURN } temp \}$$

Execuția FAA provoacă ca valoarea lui x să fie returnată iar x să fie înlocuit cu $x + e$ ca și o singură operație indivizibilă.

Cunoscând că x este variabila partajată și că P_i și P_j execută simultan P_1 și P_2 , atunci:

$$P_i : temp_i \leftarrow \text{FETCH - AND - ADD } (x, e_i)$$

$$P_j : temp_j \leftarrow \text{FETCH - AND - ADD } (x, e_j)$$

apoi fie:

$$temp_i \leftarrow x \text{ sau } temp_i \leftarrow x + e_j$$

$$temp_j \leftarrow x + e_j \quad temp_j \leftarrow x$$

În orice caz valoarea lui x la sfârșit devine $x + e_i + e_j$.

FETCH - AND - ADD ($x, true$) este echivalent cu TEST - AND - SET (x)

$$\{ temp \leftarrow x, \\ x \leftarrow true \\ \text{RETURN } temp \}$$

În literatură [Dasg89], se prezintă o generalizare a primitivei FAA. Ea este:

FETCH - AND - OP (x, e)

unde:

x - variabila partajată

e - expresie compatibilă ca li tip cu x

Această operație returnează valoarea lui x și în același timp o înlocuiește cu valoarea OP (x, e).

Acum primitiva:

C: TEST - AND - SET (k)

poate fi realizată de instrucția:

C: FETCH - AND - OR ($k, TRUE$)

unde (k) este o variabilă booleană.

Pentru a reliefa avantajul primitivei FAO față de TAS considerăm două procese concurente P_1 și P_2 ce operează excluderea mutuală asupra unei variabile utilizând secțiuni critice.

Utilizând TAS, avem:

versiunea 1 (inițial, $k = 0$)

P_1 :

repeat C_1 : TEST - AND - SET (k) **until** $C_1 = 0$ ("LOCK")

 secțiunea critică CS_1

$k := 0$ ("UNLOCK")

P_2 :

repeat C_2 : TEST - AND - SET (k) **until** $C_2 = 0$ ("LOCK")

 secțiunea critică CS_2

$k := 0$ ("UNLOCK")

Utilizând FAO, avem:

versiunea 2 (inițial, $k = FALSE$)

$P_1: \dots$

repeat C_1 : FETCH – AND – OR (k , TRUE) **until** $C_1 = \text{FALSE}$

secțiunea critică CS_1

$k := \text{FALSE}$

$P_2: \dots$

repeat C_2 : FETCH – AND – OR (k , TRUE) **until** $C_2 = \text{FALSE}$

secțiunea critică CS_2

$k := \text{FALSE}$

În versiunea 2, dacă se încearcă execuții simultane a două FAO, proprietatea de serializare va fi garantată fie pentru:

$C_1 = \text{FALSE}$, $C_2 = \text{TRUE}$

în care caz P_1 va intra în CS_1 , sau:

$C_1 = \text{TRUE}$, $C_2 = \text{FALSE}$

în care caz P_2 va intra în CS_2

În oricare din cazuri $k = \text{TRUE}$.

Ceea ce este important, cele două actualizări ale lui k vor fi combinate într-un singur acces la memorie.

În contrast modificarea lui k în varianta 1 se face de către două operații TAS în mod serial.

Primitiva COMPARE – AND – SWAP (CAS)

CAS este o primitivă de sincronizare ce blochează o variabilă partajată, o actualizează și apoi o deblochează.

Semantica sa este:

COMPARE – AND- SWAP (r -old, r -new, addr)

{temp ← m (addr): if temp = r – old

then { m (addr) ← r – new; $z = \leftarrow 1$; }

else { r – old ← m (addr); $z = 0$; }

}

Variabila partajată m (*addr*) reprezintă conținutul unei locații de memorie denumite *addr*, și $r - \text{old}$ și $r - \text{new}$ sunt două registre utilizate pentru a înmagazina valoarea curentă și noua valoare a variabilei partajate. Succesul operației CAS este indicat de fanionul Z. CAS este mult mai puternică decât TAS deoarece după ce noua valoare a variabilei partajate este disponibilă, într-o singură operație neîntreruptibilă această primitivă de sincronizare readuce variabila partajată, verifică dacă valoarea ei a fost schimbată, și dacă nu, efectuează actualizarea.

Metoda numărătorului

Sincronizarea prin *metoda numărătorului* (*contorului*) este utilizată pentru un număr mai mare procese. De exemplu patru procese P_1, P_2, P_3 și P_4 accesează o variabilă partajată. Un contor este asociat variabilei partajate. Fiecare proces P_n poate să acceseze locația variabilei partajate doar când contorul conține același număr ca și cheia asigurată aceluși proces.

Semafoare

Recunoscând dezavantajul că un procesor ce execută LOCK se găsește în starea *busy-waiting*, Dijkstra (1968) a propus două noi operații indivizibile notate cu P și V, definite asupra unei variabile speciale denumită *semafor*.

Un semafor este o variabilă întregă ne-negativă. Dându-se un semafor S , operațiile $P(S)$ și $V(S)$ sunt definite astfel:

```

P(S) :  if S > 0
        then S := S - 1
        else blochează procesorul ce execută P(S) și plasează-l într-un “șir de
         așteptare” asociat cu S;
        Alocă procesorul la alt proces din “șirul operațional”

V(S) :  S := S + 1
        if există oricare alte procese în “șirul de așteptare” asociat cu S
        then selectează un astfel de proces din “șirul de așteptare” și plasează-l în
         “șirul operațional”

```

Cu alte cuvinte, $P(S)$ va bloca procesul apelant atâta timp cât S are o valoare pozitivă.

$V(S)$ va face ca un proces ce a fost anterior blocat pe acest semafor să fie plasat în starea *ready* astfel încât să fie eventual alocat unui procesor disponibil. Astfel condiția *busy-wait* este evitată.

Utilizând un semafor, problema excluderii mutuale poate fi rezolvată astfel. Presupunem că $S = 1$.

Atunci:

$P_1: \dots\dots\dots$ $P(S)$ secțiunea critică CS_1 $V(S)$ $\dots\dots\dots$	$P_2: \dots\dots\dots$ $P(S)$ secțiunea critică CS_1 $V(S)$ $\dots\dots\dots$
---	---

Similar, problema sincronizării condiției descrisă anterior poate fi soluționată utilizând două semafoare S_1 și S_2 . Presupunem că $S_1 = 1$ și $S_2 = 0$. Atunci:

$P_1: \dots\dots\dots$ $P(S_1)$ produce valoarea și depozitează în D $V(S_2)$ $\dots\dots\dots$	$P_2: \dots\dots\dots$ $P(S_2)$ consumă valoare în D $V(S_1)$ $\dots\dots\dots$
--	---

Un semafor el însuși este o variabilă partajată, și P și V ca operații asupra lui S trebuie să satisfacă proprietatea de excludere mutuală. Prin urmare P și V trebuie să fie programate ca și secțiuni critice utilizând primitiva TAS. Deoarece acestea vor fi secțiuni critice foarte mici, timpul de busy – wait va fi mic.

2.5.3 Sincronizarea bazată pe transfer de mesaje

Sistemele multiprocesor cu transfer de mesaje - SMP-TM - sunt arhitecturi aparținând taxonului MIMD-L.

Reamintim că din punct de vedere structural un SMP -TM este format din perechi procesor/memorie - PM - interconectate prin RIN. Comunicarea datelor se face prin *transfer de mesaje* și nu prin variabile partajate. Lungimea mesajelor diferă, dar în mod uzual fiecare mesaj constă dintr-un număr de pachete de dimensiune fixă. Comunicarea între calculatoare urmează un protocol de comunicare predeterminat. Conexiunile între

calculatoare pot fi seriale de mare viteză sau paralele. Prin urmare nodul activ într-un SMP - TM este un calculator.

În fig. 2.40 se prezintă o structură simplă formate din 8 PM conectate sub forma unui *cub boolean*, o structură particular a *hipercubului*.

Procesele sunt distribuite printre PM -uri. Fiecare proces are acces direct la memoria sa locală.

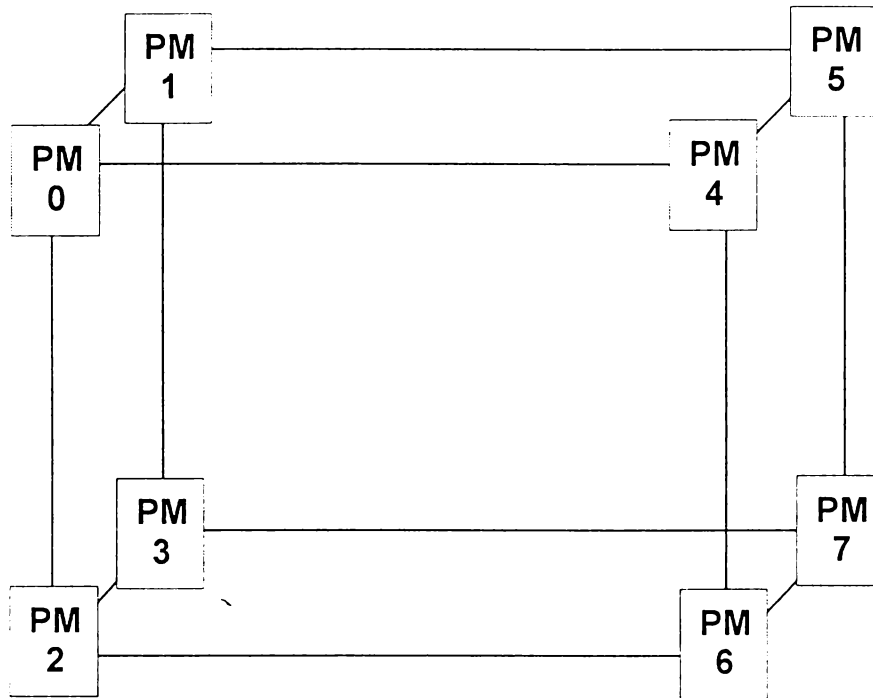


Fig. 2.40 SMP- TM cu 8 PM organizate în cub boolean

Sincronizarea proceselor se bazează pe principiul *transferului de mesaje*.

Operațiile primitive pentru transferul mesajelor sunt în mod uzual de tip SEND/RECEIVE.

Un proces trimite un mesaj executând o operație de forma:

SEND expresie TO procesul de destinație

ce face ca mesajul - valoarea expresiei - să fie transmis spre procesul de destinație.

Un mesaj este recepționat de un proces când se execută:

RECEIVE variabilă FROM procesul sursă

ce face ca mesajul din procesul sursă să fie recepționat în variabilă.

Când un proces *PS* emite și altul *PR* recepționează se stabilește un *canal de comunicație*.

Cele mai restrictive canale de comunicație sunt acelea în care procesele au o nominalizare explicită de destinație / sursă. Astfel de protocoale se numesc cu *nominalizare directă*.

În SMP -TM bazate pe protocoale de tip *nominalizare directă*, mijloacele pentru sincronizarea proceselor emitente și a celor de receptare pot fi afectate în mai multe moduri.

În sistemele *sincrone* nu există buffere asociate cu canalul de comunicație. Prin urmare executarea unui SEND de către procesorul *PS* este întotdeauna întârziată până când *PR* execută RECEIVE. *PS* este sursa iar *PR* destinația mesajului.

În sistemele *asincrone* se asociază buffere canalelor de comunicație. Operația SEND este întârziată doar când bufferul este plin și operația RECEIVE va fi întârziată doar când bufferul este gol.

Un model de transfer de mesaj este cel utilizat în *transputer* utilizând limbajul de programare Occam. Acest limbaj permite programarea concurentă.

În exemplul anterior prezentat în fig. 2.39 problema producător - consumator consta în faptul că P_1 produce un flux de valori ce sunt consumate de P_2 în ordinea în care valorile sunt produse.

În modelul cu variabilă partajată, valorile produse au fost depozitate în structura partajată D care de asemenea servea ca sursă de valori pentru P_2 . Accesul corect la D presupune utilizarea *sincronizării de condiție*.

Utilizând Occam, aceeași problemă poate fi rezolvată astfel. Două *procese* evoluează în paralel și comunică prin *canalul X*. (fig. 2.41).

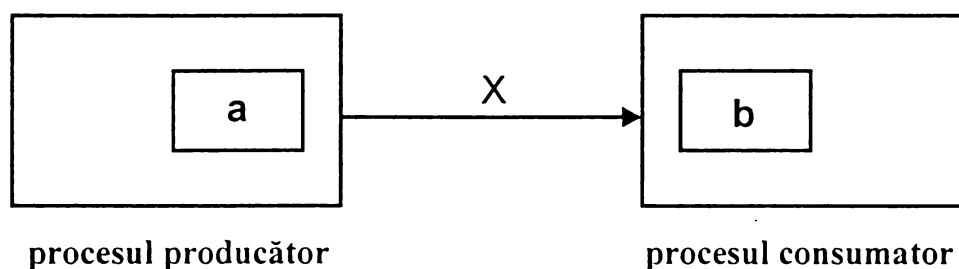


Fig. 2.41 Structura proceselor Occam și canalul lor de comunicare

Compoziția întregului proces paralel este indicată sub forma:

```
par
    [componentă
    [componentă
```

unde componentele sunt indicate în raport cu cuvântul cheie *par*. Fiecare din cele două componente ale procesului paralel este un proces secvențial iterativ, indicat de forma textuală:


```

chan x
par
  while true
    var a
    seq
    _ calculează valoarea și plasează în a
    x ! a
  while true
    var b
    seq
    x ? b
    consumă valoarea în b

```

Procesul producător execută o instrucție de ieșire:

$x ! a$

care transmite variabilă locală în x_j în timp ce procesul consumator execută:

$x ? b$

ce cauzează ca valoarea din calcul x să fie citită în b , variabila locală.

Înstrucțiile de intrare/ieșire sunt două din cele trei primitive din Occam.

Sincronizarea este obținută prin faptul că o instrucție de intrare în care x este implicat, nu se termină decât atunci când o instrucție de ieșire, ce implică același canal, nu este executată.

2.6 CONCLUZII

Din analiza arhitecturilor sistemelor multiprocesor se pot desprinde câteva concluzii referitoare atât la ariile de aplicabilitate cât și la tendințele de evoluție ulterioare.

SMP-MP *simetrice* reprezintă o extensie naturală a stațiilor de lucru și a calculatoarelor personale. O aplicație secvențială poate fi rulată în totalitate nemodificată și un beneficiu în performanță se obține prin utilizarea avantajului de a avea o mare cantitate de memorie principală partajată și capacități I/E tipic disponibile în astfel de mașini. Aplicațiile paralele sunt de asemenea relativ ușor de implementat deoarece toate datele partajate sunt direct accesibile de la toate procesoarele utilizând instrucțiuni obișnuite de tip LOAD și

STORE. O paralelizare graduală este posibilă prin paralelizarea selectivă în mod intensiv a unor aplicații secvențiale. Pentru aplicații multiprogramate un avantaj substanțial îl reprezintă granularitatea fină în care resursele pot fi partajate printre aplicațiile de proces și de către sistemul de operare.

Deoarece procesorul, sistemul de memorie, circuitele integrate și tehnologiile de asamblare continuă să se amelioreze rapid, sarcina proiectantului se focalizează pe găsirea unor structuri performante. SMP cu număr redus de procesoare vor continua să fie importante din cel puțin trei aspecte. Primul se referă la faptul că ele oferă o combinație atractivă cost-performanță. Al doilea se referă la faptul că astăzi microprocesoarele sunt proiectate să fie pregătite pentru includerea în structuri multiprocesor, dispunând de tot setul de semnale necesare pentru interfațare. A treia rațiune se referă la faptul că tehnologia software esențială pentru mașinile paralele – compilatoare, sisteme de operare, limbaje de programare – se dezvoltă rapid pentru astfel de SMP-MP de scară redusă.

Aspectele referitoare la proiectare sunt importante la nivelele moderate de paralelism. Desigur, alternativele de proiectare optimă se pot schimba. De exemplu, deși nu sunt foarte populare, este posibil ca memoriile cache partajate la un anumit nivel de ierarhie să devină foarte atractive atunci când tehnologia de încapsulare a modulelor multichip va deveni ieftină sau când se vor putea încapsula procesoare multiple într-o singură capsulă.

Deși este un mecanism puternic, o interconectare de tip magistrală partajată are în mod clar limitări în ceea ce privește lățimea de bandă pe măsură ce numărul de procesoare sau viteza procesorului cresc. Oricum, soluția generală o reprezintă construcția mașinilor coerente din punct de vedere cache, scalabile, în care memoria să fie în mod fizic distribuită printre noduri și în care să se utilizeze o interconectare scalabilă.

Aceste abordări trebuie corelate cu protocoalele de menținere a coerenței ce nu se bazează pe monitorizare. În proiectarea arhitecturilor cu memorie partajată coerente din punct de vedere cache, indiferent de interconectarea utilizată, plasarea unor capacități de memorare în astfel de RIN reprezintă o opțiune majoră ce ameliorează performanța. Indiferent însă de opțiunea de proiectare, aceste proiecte bazate pe magistrală oferă blocul constructiv de bază pentru arhitecturi paralele pe scară largă.

O evoluție firească o reprezintă arhitecturile paralele pe scară largă ce sunt construite cu noduri de utilizare generală ce au în componență o ierarhie de memorie locală completă

augmentată de o interfață de comunicare pretabilă la rețele scalabile. Oricum, o gamă largă de opțiuni de proiectare este disponibilă pentru interfața de asistare a comunicării. Proiectarea este puternic influențată de locul în care interfețele de asigurare a comunicării sunt plasate în arhitectura nodului: la nivel de procesor, la nivel de controller al memoriei cache, la nivel de magistrală a memoriei sau al magistralei de I/E. De asemenea o puternică influență o reprezintă *arhitectura de comunicare și modelul de programare*. Modelele de programare sunt implementate în mașinile paralele de talie mare utilizând protocoale construite pe baza tranzacțiilor primitive de la nivel de rețea. Riscurile în implementarea unui astfel de protocol se referă la numărul mare de tranzacții ce trebuie să fie gestionate simultan și la arbitrarea globală care în general nu este disponibilă.

Ca și un rezultat se constată evoluția a cel puțin două direcții în care arhitecturile paralele pe scală largă evoluează. Mașini orientate în general pe conceptele de transfer de mesaj și apariția unor grupuri de procesoare sau de memorii reprezintă o direcție. Cealaltă direcție este reprezentată de mașini care integrează puternic rețeaua de comunicare în sistemul de memorie pentru a furniza accese coerente din punct de vedere cache la un spațiu fizic de adresă partajată. Desigur, riscurile în promovarea unor mecanisme sofisticate de menținere a coerenței se referă la limitările în ceea ce privește lățimea de bandă scalabilă și latența redusă.

Sistemele scalabile ce suportă un spațiu de adresă partajată coerentă reprezintă o parte importantă a mediului de multiprocesare deoarece ele combină ușurința programării unui model de programare cu spațiu de adresă partajată cu avantajele în ceea ce privește scalabilitatea unei memorii distribuite și a interconectării aferente. Suportul hardware pentru asigurarea coerenței cache asigură evoluția unor protocoale bazate pe director atât în situația când sistemele se bazează pe memorie cât și pe memoria cache.

În timp ce suportul hardware pentru asigurarea coerenței cache reprezintă o opțiune de proiectare costisitoare, apariția standardelor și a faptului că microprocesoarele însăși oferă suport pentru coerența cache, ameliorează această opțiune. Odată ce protocolul de menținere a coerenței este disponibil la nivel de microprocesor, proiectanții pot dezvolta protocoale la nivel de multiprocesor și arhitecturi de comunicație încă de la început.

3. ANALIZA DE PERFORMANȚĂ A SISTEMELOR DE CALCUL DIN CLASA MIMD

3.1 ANALIZA DE PERFORMANȚĂ LA RIN

3.1.1 Terminologie de bază utilizată în evaluarea performanțelor RIN

În evaluarea performanțelor RIN se utilizează mai mulți parametri derivați din metrica general aplicabilă în determinarea dependabilității sistemelor de calcul.

În analiza sistemelor *sincrone* se utilizează următorii parametri:

- *Lățimea de bandă a memoriei (BW)* - numărul mediu de module active de memorie într-un ciclu de transfer al unei RIN. Termenul "activ" înseamnă că un procesor a efectuat cu succes o operație referitoare la memorie (citire sau scriere) în modulul definit ca fiind activ. *BW* se utilizează în analizarea RIN *sincrone*.

În analiza sistemelor *asincrone*, se utilizează următorul parametru:

- *Capacitatea de trecere (Thr)* - numărul mediu de pachete furnizate de către RIN în unitatea de timp.

În RIN utilizate în SMP - MP, *Thr* reprezintă numărul mediu de accese la memorie efectuate în unitatea de timp.

În continuare definim:

- *Probabilitatea de acceptanță (P_A)* - raportul dintre lățimea de bandă prevăzută și numărul așteptat de cereri de acces generate per ciclu.

$$P_A = \frac{BW}{Np}$$

- *Utilizarea procesorului (P_u)* - valoarea prevăzută a procentajului din timpul în care un procesor este activ. Un procesor este "activ" atunci când își efectuează calculul intern fără accesarea memoriei principale.
- *Puterea de procesare (P_p)* - extensie a P_u , reprezentând suma utilizărilor procesoarelor împărțită de numărul de procesoare din sistem.

Rezumând:

$$P_u = \frac{BW}{N\lambda T} \quad - \text{ sisteme sincrone}$$

$$P_u = \frac{Thr}{\lambda} \quad - \text{ sisteme asincrone}$$

unde:

N - nr. de procesoare

T - timpul afectat unei operații de citire/scriere

λ - rata de cerere de acces la memorie

Studiul performanței unui sistem de calcul se bazează pe diverse tehnici.

Modelul analitic reprezintă o tehnică des utilizată sub raportul costului.

Pentru a face un astfel de model practicabil se fac următoarele ipoteze de aproximare.

- *modelul de referință uniformă - MRU*. Se presupune că toate procesoarele sunt identice și atunci când un procesor efectuează o solicitare de acces la memoria principală, cererea va fi direcționată la oricare din cele M module de memorie cu aceeași probabilitate $1/M$. Aceasta înseamnă că adresa de destinație a unei cereri de acces la memorie este uniform distribuită printre cele M module de memorie.

Dacă sistemul de memorie este întrețesut această presupunere este aproximativ corectă.

Dacă sistemul de memorie nu este întrețesut, atunci există o localizare de referință și o abordare în care se definește o memorie ca fiind "favorită", este mai corectă.

- *Rata de cerere de acces la memorie λ* a unui procesor identifică cât de des un procesor accesează memoria principală. În mod indirect reflectă timpul mediu de execuție a unei instrucții.

În sistemele sincrone, acest parametru poate fi specificat ca și o probabilitate în care un procesor generează o cerere de acces la memorie la începutul unui ciclu.

În sistemele asincrone, o cerere de acces la memorie ar putea fi generată în orice moment deoarece nu există un ceas central. Se poate imagina o evoluție exponențială a timpului în sensul că durata dintre execuția unei cereri și următoarea cerere de acces la memoria principală reprezintă o variabilă aleatoare distribuită exponențial.

- *Ipoteza independenței cererii de acces* - denumită aproximarea lui Strecker - este utilizată în analiza sistemelor sincrone.

Se presupune că o cerere de acces la memorie generată într-un ciclu este independentă de cererile din ciclurile anterioare.

În realitate însă o cerere de acces rejectată în ciclul anterior va fi relansată în ciclul curent.

TABEL I. Rezumatul caracteristicilor hardware ale celor trei RIN

	<i>crossbar</i>	RIN - MN	Magistrală multiplă
Nr. de comutatoare	$N \times M$	$N \log N$	$B \times (N+M)$
Încărcarea magistralelor	N	1	B
Nr. de fire	M	N	B
Nr. de arbitrii	$M-1$ din N	$N \log N-1$ din 2	1 " B din M " $M-1$ din N
Toleranța la defecte	Acceptabilă	Slabă	Bună

3.1.2 Analiza performanțelor magistralei unice

SMP - MP cu *magistrală unică* reprezintă cele mai simple arhitecturi de sisteme multiprocesor.

La aceste sisteme, RIN este constituită dintr-o magistrală unică, partajată în timp.

Procesoarele - module *master* - solicită accesul la memoria comună - resursă *slave* - și apoi efectuează transferul de date.

Există două situații distincte în care se satisface o cerere de acces la memorie.

În primul caz, un singur procesor, la un moment dat, solicită și obține accesul la un modul de memorie, după care se efectuează operația de citire sau scriere din /în memorie. Este cazul banal.

În al doilea caz, două sau mai multe procesoare solicită simultan acces la memorie. Este o situație conflictuală ce se rezolvă prin arbitrar, după un anumit algoritm de alocare a resursei solicitate procesorului declarat câștigător.

Reamintim că la SMP - MP cu *magistrală unică*, la un moment dat, un singur procesor operează pe magistrală cu modulul de memorie solicitat.

Evident că magistrala unică reprezintă resursa ce limitează drastic performanțele acestor sisteme.

Un procesor se poate afla în una din următoarele stări :

- *activ (active)* - procesorul lucrează cu memoria sa locală - ML;
- *în accesare (accessing)* - procesorul schimbă date cu alte procesoare prin operații de scriere în sau de citire din zona de memorie comună;
- *în coadă de așteptare (queued)* - procesorul așteaptă obținerea unui acces la MC.

Pentru a determina *încărcarea sistemului (workload)* în literatură [HP96], [Kak96] etc. se prezintă două abordări posibile : *deterministă* și *probabilistică*.

În abordarea *deterministă*, încărcarea fiecărui procesor se consideră cunoscută și pe această bază se poate dezvolta o secvență precisă a operațiilor efectuate de fiecare procesor pentru a caracteriza comportarea întregului sistem. Parametrii de performanță se pot evalua prin simulare, ceea ce limitează această tehnică la cazul concret aflat sub observare.

În abordarea *probabilistică* a sistemului, se dezvoltă un model statistic care încearcă să capteze esențialul în funcționarea SMP, fără a fi cunoscute detaliile individuale ale încărcării lui. Parametrii de încărcare se definesc ca și variabile aleatoare cu o anumită distribuție.

Rezultatele furnizate au avantajul de-a reprezenta relații funcționale între parametrii de sistem și performanță. De asemenea, un astfel de model stohastic permite studiul comparativ al diferitelor soluții arhitecturale cu referire la obiectivele de performanță impuse.

În cadrul modelelor de performanță se pot admite două moduri de funcționare ale procesoarelor : *asincronă* sau *sincronă*.

Uzual se utilizează abordarea de tip *sincron* care consideră sistemul la nivelul instrucțiilor, reprezentate ca operații sincrone pentru toate procesoarele. Axa timpului este divizată în intervale elementare - *sloturi*. (fig. 3.1).

La începutul unui slot procesoarele pot emite cereri de acces la memoria comună. Mecanismul de arbitraj selectează cererile ce pot fi onorate și procesorul câștigător se conectează la modulele de memorie pe toată durata slotului.

Timpul de procesare al fiecărui procesor P este împărțit în :

- *timp activ (t_a)* - procesorul lucrează singur sau cu memoria sa locală;
- *timp de lucru cu memoria (t_m)* - procesorul, în urma unui acces reușit, lucrează cu MC.

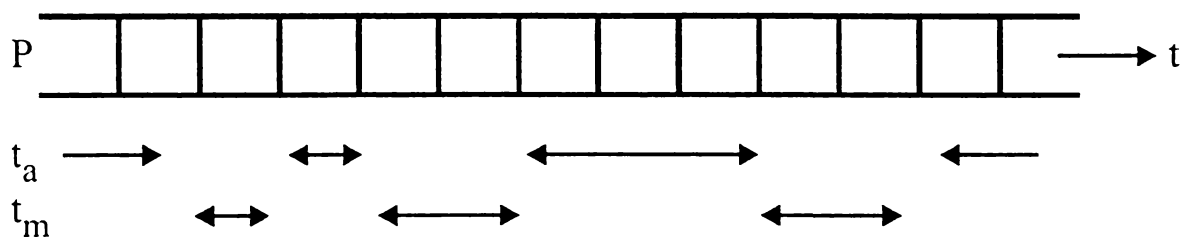


Fig. 3.1 Structura timpului de procesare

Ciclul de execuție al procesorului se studiază sub aspectul raportului dintre *timpul activ al procesorului și timpul de citire/scriere în MC.*

Fie un SMP - MP *cu magistrală unică* ce conține N procesoare, P_1, P_2, \dots, P_N , fiecare având propria memorie cache - MCH - și toate conectate la o memorie comună MC formată din M blocuri întrețesute, M_1, M_2, \dots, M_M .

Fie r probabilitatea ca un procesor oarecare P_i să solicite acces la MC la începutul unui slot.

Atunci Nr reprezintă numărul așteptat de solicitări de acces la MC.

Dacă $r = 1$ atunci $t_a = 1$ iar pentru $r < 1$ acest timp va avea o distribuție geometrică.

BW reprezintă lățimea de bandă a subsistemului compus din o magistrală unică ce conectează M blocuri de memorie, adică numărul așteptat de blocuri de memorie ocupate. Acesta reprezintă gradul de ocupare al memoriei.

Datorită conflictelor de acces există $BW < Nr$.

Probabilitatea ca să existe o solicitare de acces de la P_i către un anumit bloc M_j este r/M .

Probabilitatea ca să nu existe solicitări de la P_i către M_j este $1-r/M$ și prin urmare probabilitatea ca nici unul dintre procesoare să solicite M_j la începutul unui slot este $(1 - r/M)^N$.

Fie E_j evenimentul ce constă în faptul că există cel puțin o solicitare pentru M_j .

Atunci :

$$\Pr[E_j] = 1 - \left(1 - \frac{r}{M}\right)^N = q \quad (3.1)$$

pentru orice i .

Lățimea de bandă va fi :

$$BW = 1 - \left(1 - \frac{r}{M}\right)^N = q \quad (3.2)$$

O rată ridicată a cererilor de acces duce la saturarea rapidă a magistralei.

Capacitatea de trecere Thr se poate exprima ca și o funcție de numărul de instrucții executate pe secundă de către SMP.

Fie :

- T_S - debitul SMP
- T_p - debitul unui procesor P_i
- b - factorul de utilizare al magistralei ce exprimă numărul de cicluri de magistrală solicitate de fiecare procesor

În fig. 3.2 se prezintă variația funcției $T_S = f(N)$ pentru diferite valori ale lui b .

Se observă o ameliorare semnificativă a performanțelor atunci când b scade prin utilizarea memoriilor locale.

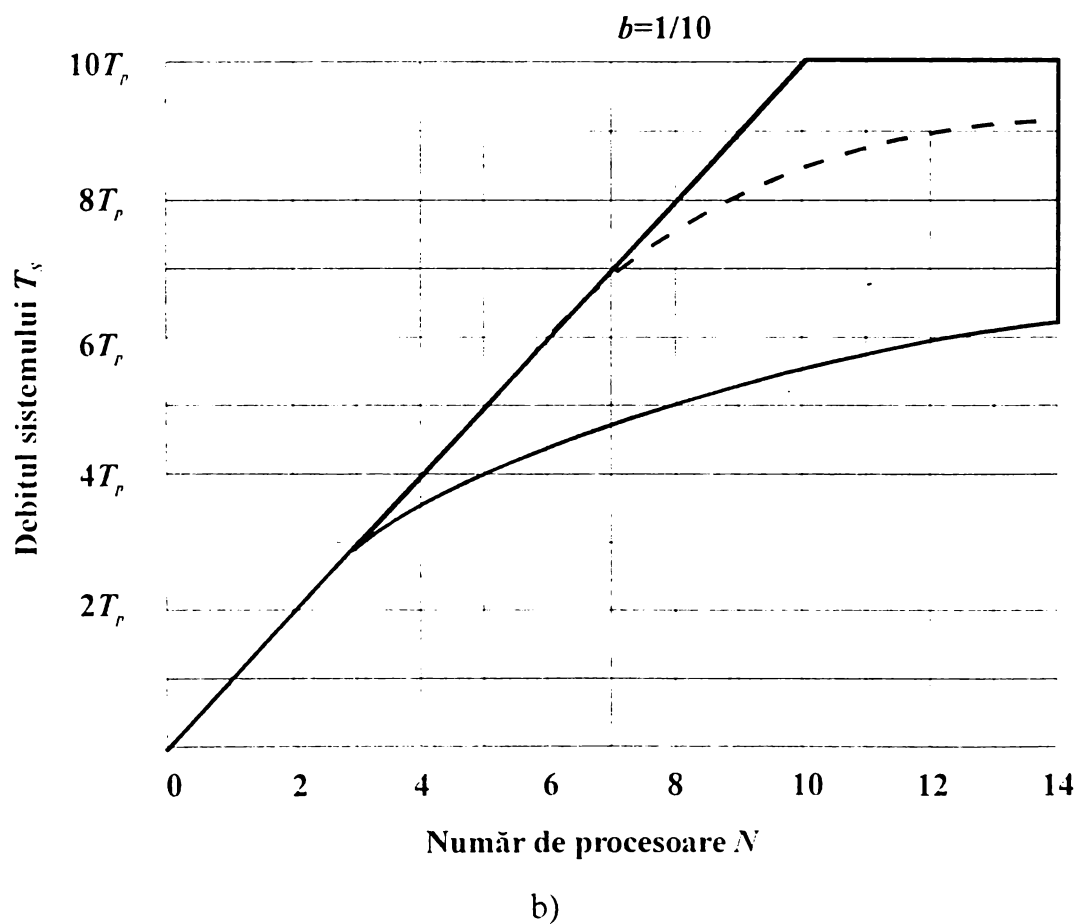
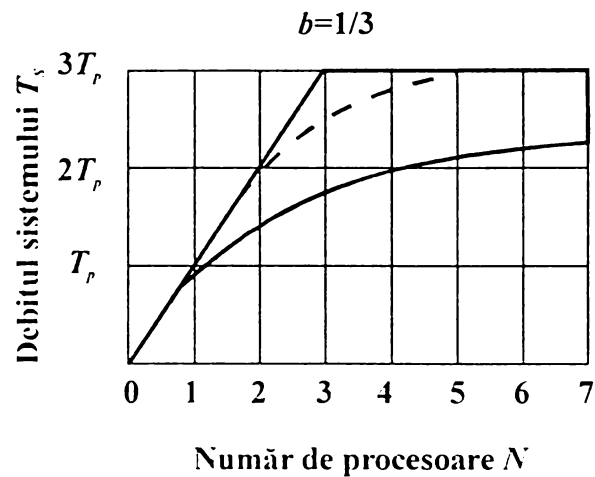


Fig. 3.2 Saturația magistralei

Dacă se consideră că W este numărul de procesoare aflate în așteptare atunci se poate exprima :

$$W = \begin{cases} 0 & \text{dacă } N \leq \frac{1}{b} \\ N - \frac{1}{b} & \text{dacă } N > \frac{1}{b} \end{cases} \quad (3.3)$$

Atunci când la un SMP apare saturarea magistralei, traficul mare de pe magistrală va conduce la creșterea numărului de procesoare aflate în așteptare și prin urmare se pierde din eficiența dobândită prin distribuirea sarcinilor la mai multe procesoare.

Curba caracteristică $W = f(N, b)$ este prezentată în fig. 3.3.

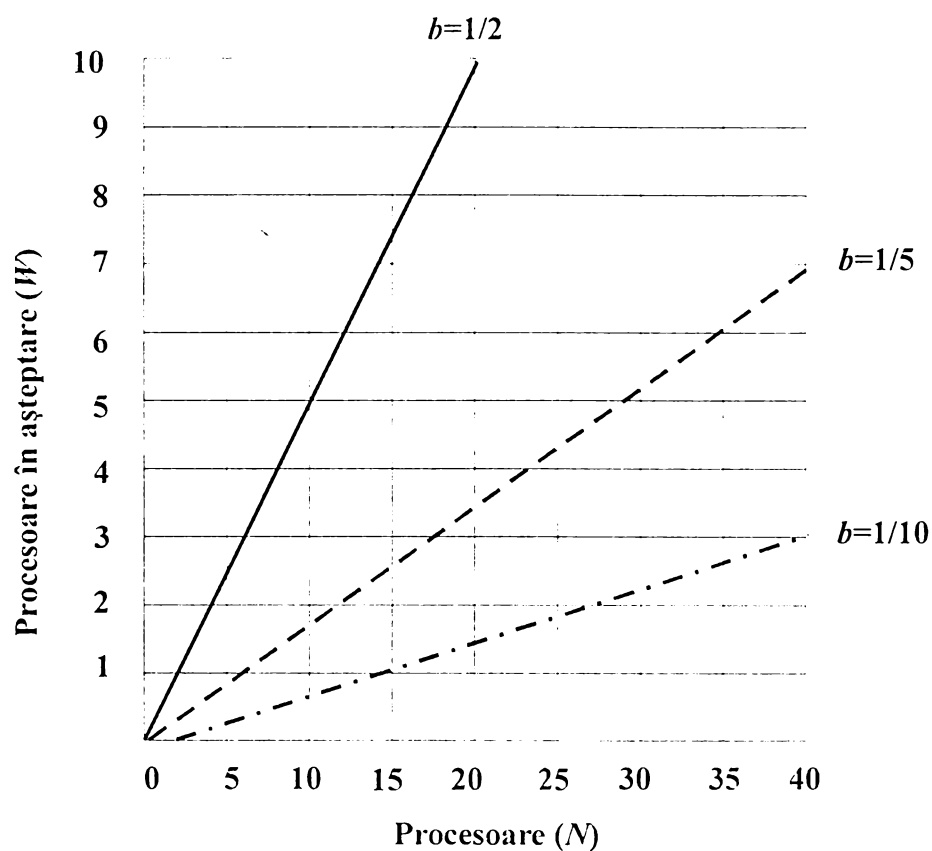


Fig. 3.3. Raportul între procesoarele în așteptare și cele active

3.1.3 Analiza performanțelor magistralei multiple.

În cazul SMP-MP o sursă majoră de limitare a performanțelor și implicit a fiabilității și a toleranței de defecte o reprezintă magistrala unică.

Una dintre metodele de ameliorare a performanțelor o reprezintă creșterea capacității de transfer a magistralei.

O soluție luată în considerare [MHW87], [YB91] o reprezintă *multiplicarea numărului de magistrale* în scopul asigurării unei fiabilități superioare și a unei *BW* mărite.

Arhitectura propusă reprezintă cazul conectării *complete* și este prezentată în fig. 3.4.

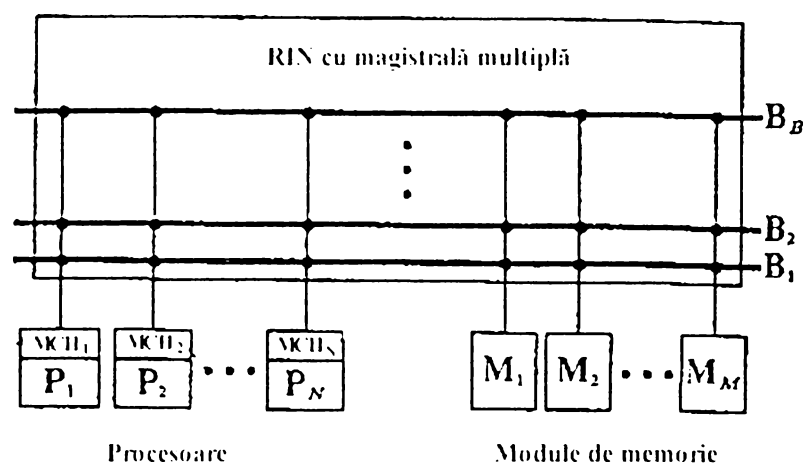


Fig. 3.4. Structura unui SMP-MP cu *magistrală multiplă*

Ea conține N procesoare, P_1, P_2, \dots, P_N , fiecare având propria memorie cache - MCH - și toate conectate la o memorie partajată - MP - prin intermediul a B magistrale, B_1, B_2, \dots, B_B . MP constă din M blocuri întreșesute, M_1, M_2, \dots, M_M pentru a permite solicitări concurente de acces la MP, memoria partajată. Această tehnică evită pierderea de performanță ce apare dacă accesesele sunt serializate - cazul unui singur bloc de memorie.

Fiecare procesor este conectat la fiecare magistrală și la fiecare bloc de memorie. Când un procesor dorește să aibă acces la un anumit modul, el are B magistrale din care poate alege una pentru a realiza conectarea. Astfel fiecare cuplu procesor - memorie este conectat prin câteva căi redundante ceea ce implică un comportament mai tolerant la defectarea uneia sau a mai multor magistrale, evident cu prețul degradării performanței sistemului în termeni de latență și cost.

În fig.3.5.b se prezintă cazul unei conectări *paralele* printr-o magistrală multiplă. Fiecare memorie este conectată doar la un subset din magistrala multiplă. În particular, o configurație completă cu magistrală multiplă având $B \approx N/2$ are aproximativ aceeași lățime de bandă ca și un *crossbar* $N \times M$.

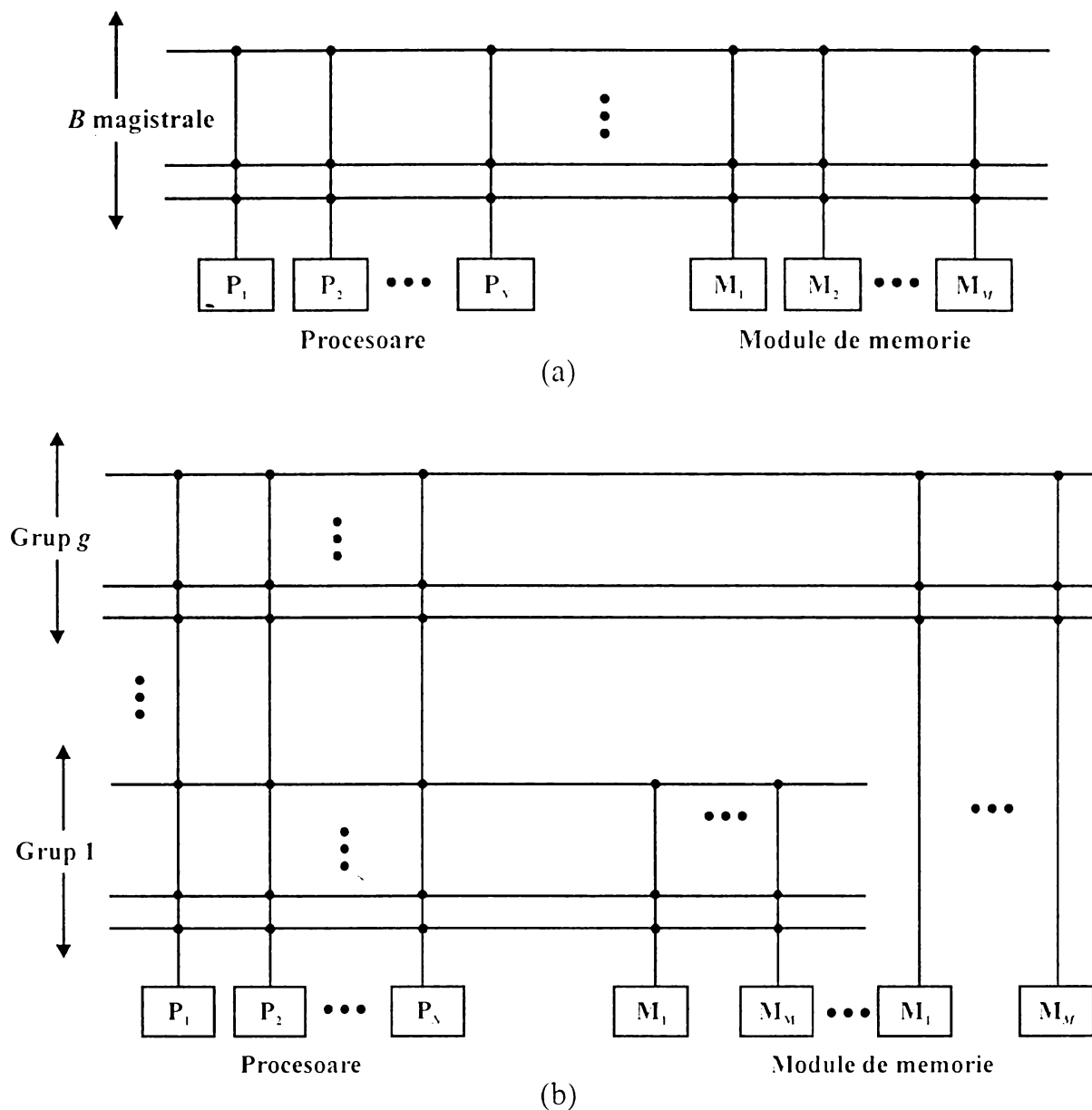


Fig. 3.5 Structura unei magistrale multiple

În SMP cu *magistrală multiplă* trebuie adoptată o politică de alocare a magistralelor disponibile la procesoarele ce solicită acces la memorie. În particular această politică trebuie aplicată în cazul când $N > B$.

Procedura de alocare se implementează prin arbitri hardware ce conduce însă la creșterea complexității rețelei de interconectare.

Tehnicile de interconectare a procesoarelor multiple și în particular, conectarea lor la o memorie partajată, au fost studiate în literatură [MHW87], [Pat81]. Cel mai studiat tip, probabil, este *rețeaua multinivel*. Acest tip de conectare are mai puține componente decât o rețea *crossbar* dar reține mai multe dintre proprietățile de conectivitate specifice ei.

În contrast, SMP-MP utilizează arhitecturi de interconectare bazate pe magistrale convenționale.

Interconectarea prin magistrale multiple reține facilitățile unei singure magistrale dar permite construcția unor sisteme mai mari ce pot concura cu sistemele bazate pe rețele multinivel în ceea ce privește puterea de calcul.

A. Operarea sistemului cu magistrală multiplă

Importanța memoriilor cache rezidă din faptul că utilizarea lor reduce traficul pe magistrală.

Se pot construi astfel sisteme mai performante fără să crească numărul de magistrale.

Creșterea numărului și a timpului de acces la memoria partajată MP depinde de algoritmul de utilizare a memoriei cache. Un aspect cheie al acestui algoritm este politica de actualizare a MP când o solicitare de acces este adresată memoriei cache - MCH.

Cele mai utilizate tehnici sunt *recopierea* și *rescrierea*, unde un bloc solicitat din MCH este copiat în MP numai când blocul este dealocat din MCH.

Două surse de conflict legate de solicitările la memorie sunt prezente în sistemul din fig. 3.4.

Prima, mai multe solicitări pot fi adresate aceluiași modul de memorie. A doua sursă, capacitatea disponibilă a magistralei poate fi insuficientă pentru a se acomoda la toate solicitările.

În mod corespunzător, alocarea unei magistrale la un procesor necesită un proces în două faze după cum urmează :

- (1). Conflictele de acces la memorie sunt rezolvate mai întâi de M -arbitrii "1 din N ", câte unul pentru fiecare modul de memorie. Fiecare arbitru "1 din N " selectează o solicitare din până la maximum N posibile pentru a permite accesul la blocul de memorie;
- (2). Solicitățile de acces la memorie selectate de arbitrii de memorie sunt apoi alocate la o magistrală de câte un arbitru " B din M ". Acesta poate selecta până la B solicitări din M arbitri de memorie.

Ipoteza operării asincrone a liniilor de adrese și de date permite ca procesul de arbitrare să se suprapună peste transferul datelor.

Aparent s-ar părea că ar trebui să existe o magistrală pentru fiecare bloc de memorie.

Aceasta ar conduce la necesitatea unui al doilea nivel de arbitrare. Configurația optimă poate necesita mai puține magistrale decât numărul de blocuri de memorie.

B. Modelarea performanțelor

Cererile de acces ale unui procesor la memorie pot fi modelate ca și o secvență de încercări Bernoulli independente. Se pot dezvolta expresii pentru un număr mediu de solicitări de acces produse de N procesoare la începutul fiecărui ciclu de memorie.

Fie p probabilitatea ca un procesor oarecare P_i să solicite acces la MP la începutul unui ciclu de memorie (aceasta este o încercare de tip Bernoulli), atunci Nr reprezintă numărul așteptat de solicitări la MP.

Unele cereri vor fi întotdeauna blocate datorită celor două tipuri de conflicte menționate anterior, indiferent de dimensiunile lui B și M .

BW reprezintă lățimea de bandă a subsistemului compus din B magistrale și din M blocuri de memorie, adică numărul așteptat de blocuri de memorie ocupate. De asemenea BW poate fi considerat ca reprezentând numărul așteptat de accese reușite la memorie.

Datorită conflictelor de acces $BW < Nr$.

Cînd două sau mai multe procesoare intenționează să acceseze memoria, se utilizează un arbitru pentru selectarea doar a unei cereri.

Dacă presupunem că blocurile de memorie sunt întrepătrunse pe spațiul biților mai puțin semnificativi ai adresei blocului cache și că ciclurile de extragere a instrucțiilor și accesul la date se întrepătrund, atunci empiric se poate sugera că solicitările la toate blocurile de memorie sunt egal probabile.

Probabilitatea ca o solicitare de acces de la P_i către un bloc anume M_j este r/M independent de i sau j .

În literatură [MHB86] dezvoltarea unei expresii pentru BW se face în doi pași, corespunzător celor două niveluri de arbitrare.

(1) Arbitrarea memoriei

Probabilitatea să nu existe solicitări de la P_i către M_j este $1 - p/M$, și prin urmare probabilitatea ca nici unul dintre procesoare să solicite M_j la începutul unui ciclu de memorie este $(1 - p/M)^N$.

Fie E_j evenimentul constând în faptul că există cel puțin o solicitare pentru M_j .

Cu alte cuvinte că arbitrul "1 din N ", pentru M_j , extrage o solicitare de acces.

Atunci :

$$\Pr[E_j] = 1 - \left(1 - \frac{p}{M}\right)^N = q \quad (3.4)$$

pentru orice i .

Dacă evenimentele E_j , $j = 1, \dots, M$, se presupun ca fiind independente și dacă există întotdeauna un număr suficient de magistrale, adică $B \geq M$, atunci numărul așteptat de blocuri de memorie ocupate este :

$$BW_S = \sum_{j=1}^M \Pr[E_j] = M \left[1 - \left(1 - \frac{p}{M}\right)^N \right] \quad (3.5)$$

unde S înseamnă "suficiente magistrale".

În cazul unui N mai mare, expresia are limita inferioară aproximativă

$$BW_S \approx M \left(1 - e^{-\frac{Nr}{M}} \right) \quad (3.6)$$

Ecuția (3.6) poate fi evaluată prin estimarea lui p , a capacității MC, a raportului α dintre timpul ciclului de memorie și timpul ciclului de procesor.

În cel mai general caz, când $B < M$, se ajunge la arbitrarea magistralilor.

(2) Arbitrarea magistralilor

Presupunerea că evenimentele E_j sunt independente nu permite să exprimăm $f(i)$ ca exact i dintre arbitrii de memorie să producă o solicitare la începutul unui ciclu de memorie astfel:

$$f(i) = \binom{M}{i} \Pr[E_i]^i (1 - \Pr[E_i])^{M-i} = \binom{M}{i} q^i (1 - q)^{M-i} \quad (3.7)$$

În cazul unde $i \leq B$, există suficiente magistrale care să manipuleze solicitările de acces și arbitrul " B din M " nu blochează nici o solicitare. În cazul unde $i > B$, toate magistralele B sunt utilizate și arbitrul " B din M " blochează $(i - B)$ din solicitări.

Tinând cont de aceste două cazuri se poate scrie expresia pentru numărul de blocuri de memorie utilizate astfel:

$$BW = \sum_{i=1}^B if(i) + \sum_{i=B+1}^M Bf(i) \quad (3.8)$$

Când $B = M$, $BW = BW_S$.

În general însă, $Nr > BW_S > BW$ pentru $r > 0$.

Aceste inegalități corespund conflictelor ce fac ca solicitările de acces la memorie să fie blocate în timpul arbitrării memoriei ($Nr > BW_S$) și în timpul arbitrării magistralei ($BW_S > BW$).

Mudge [MHB86] a observat că ecuația (3.6) se bazează pe două presupuneri: *independența temporală* și *independența spațială*.

Independența temporală pretinde ca solicitările succesive de acces la memorie să fie independente.

Independența spațială, corespunde evenimentelor E_i , independente, o ipoteză cu validitate limitată.

Dacă notăm

$$\mathbf{F}(B) = \sum_{i=B}^M f(i) \quad (3.9)$$

atunci

$$BW = B\mathbf{F}(B) + \sum_{i=1}^{B-1} if(i) \quad (3.10)$$

Se poate generaliza (3.10) pentru cazul magistralelor partajate (fig. 3.5.b). Analiza interferenței de memorie este aceeași ca și în cazul anterior deoarece este independentă de configurația magistralei.

Analiza interferenței magistralelor necesită modificări. Fie cele B magistrale divizate în g grupe egale, presupunând că g este un factor depinzând de B și M . Cu $m = M/g$ și $b = B/g$ ecuațiile (3.9) și (3.10) devin

$$fg(i) = \binom{m}{i} q^i (1-q)^{m-i} \quad (3.11)$$

$$\mathbf{F}_g(B) = \sum_{i=b}^m fg(i) \quad (3.12)$$

În mod corespunzător

$$BW'_g = g \left[b\mathbf{F}_g(B) + \sum_{i=1}^{b-1} if_g(i) \right] \quad (3.13)$$

care reprezintă de g ori lățimea de bandă a oricărui din cele g subsisteme formate din N procesoare, b magistrale și m memorii.

În fig. 3.6 se prezintă comportamentul BW în funcție de numărul de magistrale.

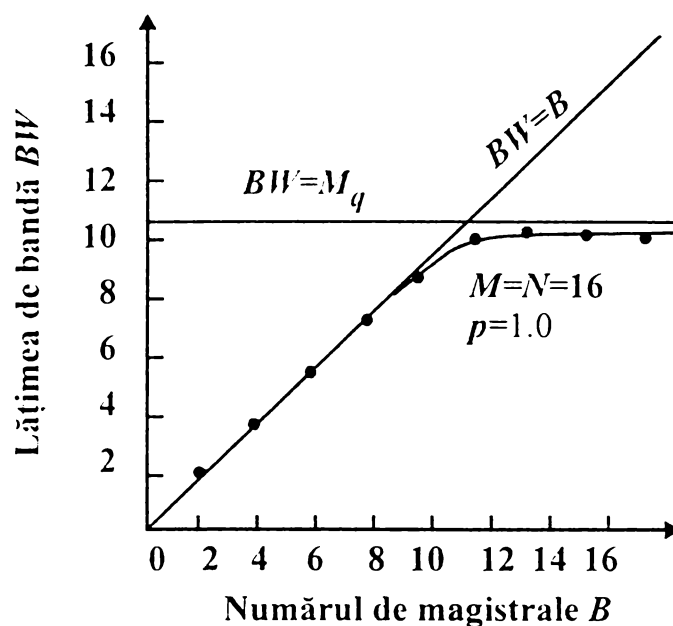


Fig. 3.6 Comportarea asimptotică a lui BW

C. Proiectarea circuitelor de arbitrare

Se utilizează două tipuri de arbitri:

“1 din N ” – selectează procesorul din cele ce solicită accesul

“ B din M ” – alocă magistralele la acele procesoare care au reușit să obțină accesul la memorie.

Arbitrii “1 din N ”

Fiecare procesor P_i are o linie de solicitare R_i și o linie de alocare G_i . Procesorul P_i solicită un acces la memorie prin activarea liniei R_i și arbitrul indică alocarea blocului de memorie activând G_i .

Schemele de arbitrare se pot împărți în trei categorii: *cu prioritate fixă*, *în inel* și *structură arborescentă*.

Arbitrii cu *prioritate fixă* sunt simpli și rapizi, dar ei au dezavantajul că nu asigură șanse egale în procesul de arbitrare. Procesoarele cu prioritate scăzută pot fi forțate să aștepte nedefinit dacă procesoarele cu prioritate ridicată mențin memoria ocupată.

Arbitrii cu *structură inelară* oferă prioritate procesoarelor pe baza de rotație, cu prioritatea cea mai coborâtă data procesorului ce a utilizat cel mai recent modulul de memorie solicitat.

Se garantează astfel tuturor procesoarelor că într-un interval finit de timp vor accesa memoria, dar timpul de arbitrare crește linear cu numărul de procesoare.

Un arbitru "1 din N" cu structura arborescentă este în general un arbore binar de adâncime $\log_2 N$ construit din module de arbitrare "1 din 2". Schema unui astfel de arbitru este prezentată în fig. 3.7.

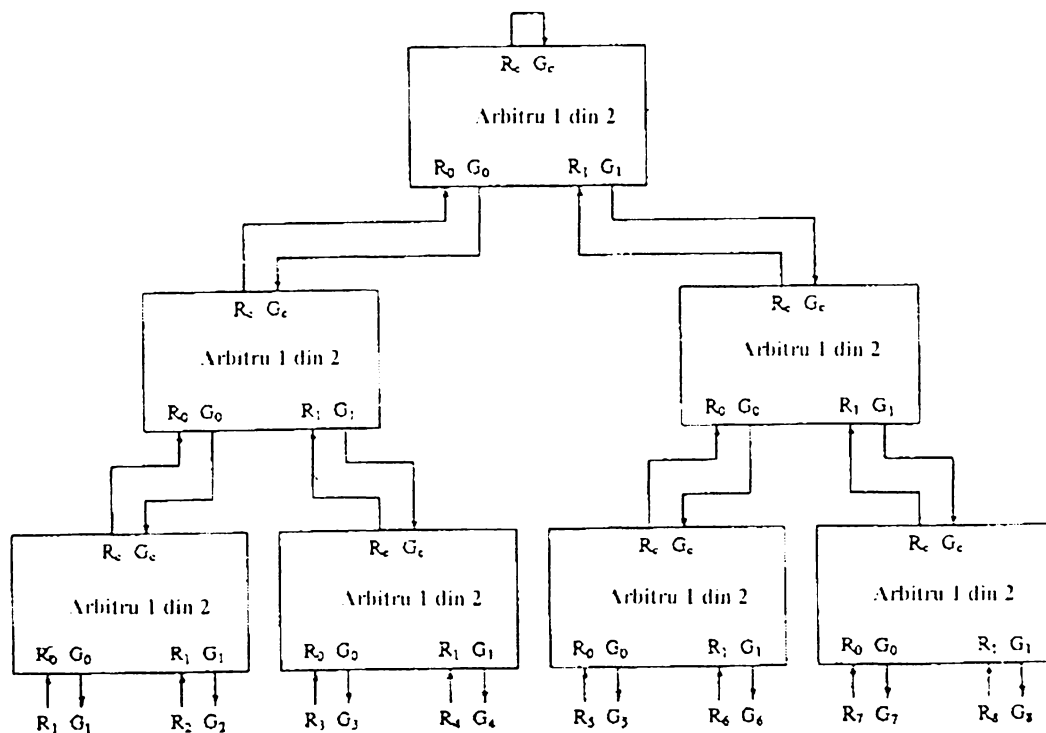


Fig. 3.7. Schema bloc a arbitrului pentru magistrala multiplă

Acest tip de arbitru este mai rapid decât arbitrii în inel deoarece timpul de arbitrare crește ca și $O(\log_2 N)$ în loc de $O(N)$.

Memorarea stării este necesară pentru o abordare corectă ținând cont de asignările anterioare ale magistralei. După fiecare ciclu de arbitrare, prioritatea cea mai mare este dată modulului ce urmează imediat celui ce a fost servit - de fapt se aplică o politică standard de tip *round - robin*.

Un ciclu de arbitrare începe având toate magistralele disponibile. Registrul de stare, prezentat în fig.3.8, identifică modulul de arbitrare A_j cu cea mai mare prioritate prin poziționarea e_j la intrarea aceluia modul.

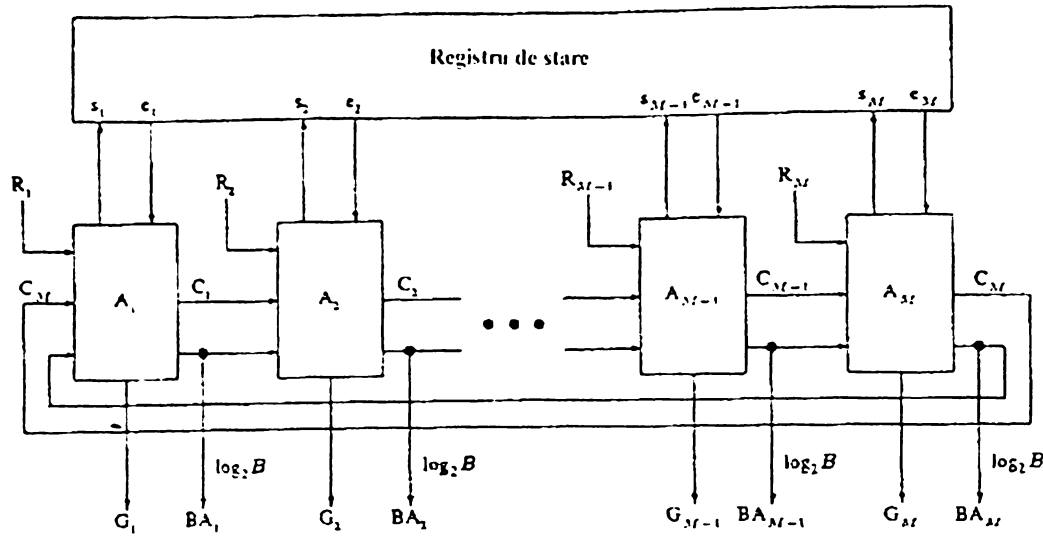


Fig. 3.8. Structura registrului de stare

Arbitrarea începe cu acest modul și parcurge inelul de la stânga spre dreapta, la fiecare modul de arbitrare, intrarea R_i este examinată pentru a se observa dacă modulul de memorie M_i solicită magistrala. Dacă o solicitare este prezentă și magistrala este disponibilă, atunci adresa primei magistrale disponibile este plasată pe ieșirea BA_i și G_i este poziționat.

BA_i este de asemenea transmis următorului modul pentru a indica magistralei cu indicele cel mai mare, că a fost asignată. Dacă unui modul nu i-a fost alocată o magistrală, atunci ieșirea corespunzătoare $BA_i = BA_{i-1}$, adică cu intrarea sa.

Dacă modulului i -a fost alocată o magistrală atunci ieșirea $BA_i = BA_{i-1}$.

Când $BA_i = B$ toate magistralele au fost utilizate și procesul de asignare încetează.

Modulul cu cea mai mare prioritate A_i , indicat de semnalul e_i , ignoră intrarea sa BA_{i-1} și începe asignarea magistralei cu prima magistrală poziționând $BA_i = 1$.

Intrarea C_{i-1} a fiecărui modul este un semnal provenind de la modulul anterior ce indică că acesta și-a completat asignarea magistralei.

Arbitrarea continuă secvențial prin toate modulele până când toate magistralele au fost asignate, sau până când toate solicitările au fost satisfăcute.

Ultimul modul își poziționează propriul semnal s_i . Acesta este înregistrat în registrul de stare, care îl utilizează pentru a selecta următoarea ieșire e_i , astfel încât următorul ciclu de arbitrare va începe imediat după acela în care s-a asignat ultima magistrală.

Referitor la performanțele arbitrilor "B din M", în literatura dedicată acestui subiect [LV 82] se evaluează întârzierea în arbitrare în raport cu M , numărul modulelor de arbitrare.

Prin combinarea a g din aceste module într-un singur modul, întârzierea este redusă cu un factor ce depinde de g .

Utilizând circuite PLA în implementarea schemelor din această categorie, circuite având o întârziere de 3Δ , întârzierea schemei de arbitrare devine $(3M/g)\Delta$.

În schema din fig. 3.7, unde $M = 16$ și $g = 4$, întârzierea devine aproximativ 12Δ .

Întârzierea pe ansamblul de arbitrare format din arbitrul "1 din N " și din "B din M " devine $18\Delta + 12\Delta = 30\Delta$.

Datorită caracterului asincron arbitrarea se poate suprapune peste accesul la magistrală ceea ce conduce la concluzia că durata ciclului memorie - magistrală poate fi în final de 30Δ .

3.1.4 Analiza performanțelor rețelelor de interconectare de tip *crossbar*

O rețea de interconectare de tip *crossbar* (RIN-CB) este constituită dintr-o matrice de perechi de contacte operate individual în care există câte o pereche pentru fiecare combinație de I/E. Schema de principiu este prezentată în fig. 2.32

O RIN-CB cu N intrări și M ieșiri este denumită $N \times M$. Atât timp cât nu există interferență de memorie printr-un set de solicitări de acces la memorie generate de către procesoare, adică nu există două sau mai multe procesoare ce solicită accese la același modul de memorie, se pot stabili toate conexiunile în mod simultan.

Această disponibilitate este însă penalizată cu un cost ridicat al conectării care este ($O(NM)$). Deși RIN-CB poate asigura simultan toate conexiunile, BW este cu mult mai scăzută decât capacitatea teoretic posibilă. Această scădere se datorează interferenței de memorie cauzată de natura aleatorie a cererilor de acces într-un mediu multiprocesor. Prin urmare analiza de performanță a unui RIN-CB se reduce la analiza interferenței de memorie.

În literatură [BYA 89] se prezintă mai multe modele pentru interferența de memorie. În general, însă operațiile sistemului sunt approximate de procese stochastice după cum urmează: La începutul unui ciclu de sistem, un P_i selectează un modul M_j în mod aleator și produce o solicitare de acces la acest modul cu o probabilitate p . Dacă se produc mai mult de o solicitare la același modul M_j , controlerul memoriei va selecta aleator una și procesoarele respinse vor reîncerca în ciclul următor. Comportamentul procesoarelor este considerat independent și statistic identic, așa cum este și comportamentul modulelor de memorie.

Pentru studiul interferenței de memorie se aplică modele de lanțuri Markov discrete. Se admite că un modul M este caracterizat de *durata ciclului* t_c , ce constă din *timpul de acces* t_a , urmat de un *timp de rescriere* t_w .

Adică:

$$t_c = t_a + t_w$$

Comportamentul procesorului P este modelat ca și o secvență ordonată, constând dintr-o solicitare de acces la memorie urmată de un *timp de execuție* t_p .

Parametrul t_p este măsurat din momentul în care data solicitată a fost obținută în urma cererii anterioare până la momentul în care următoarea cerere a fost transmisă memoriei.

În sistemele reale, $t_w = t_p$, ceea ce este echivalent cu situația în care P_i generează o cerere de acces la M_j la începutul fiecărui ciclu de memorie.

Complexitatea modelului interferenței memoriei este ameliorată dacă se presupune că un procesor încetează cererea și generează o nouă cerere independentă la începutul ciclului următor (ipoteza independenței cererilor de acces).

Pentru un sistem cu N procesoare și M module de memorie, dacă un procesor generează o cerere de acces cu probabilitatea p într-un ciclu direcționat spre fiecare modul de memorie cu o probabilitate egală (modelul de referință uniformă), atunci lățimea de bandă a memoriei este:

$$BW = M \left(1 - \left(1 - \frac{p}{M} \right)^N \right) \quad (3.14)$$

Explicația este următoarea: Dacă p/M este probabilitatea ca un procesor să solicite un acces la un anumit modul de memorie, $[1 - (p/M)]^N$ reprezintă probabilitatea ca niciunul din cele N procesoare să solicite acel modul de memorie într-un anumit ciclu.

Paranteza $\left(1 - \left(1 - \frac{p}{M} \right)^N \right)$ indică probabilitatea ca cel puțin o solicitare de acces să fie transmisă.

În fine, factorul de multiplicare M dă numărul așteptat de module distincte de memorie ce au fost solicitate într-un ciclu, și prin urmare lățimea de bandă BW .

Ecuția (3.14) reprezintă un indicator simplu de previzionare a performanței unui RIN-CB. O ipoteză mai rafinată prezentată în [BYA 89] introduce problema *referinței neuniforme* (*non uniform reference model*). În acest model –MRNU– se introduce

conceptul de *memorie favorită a unui procesor*. Modulul de memorie solicitat cel mai des de un procesor este denumit *memoria favorită a procesorului*.

Fie m probabilitatea cu care un procesor își adresează memoria sa favorită, presupunând că procesorul generează o solicitare de acces într-un ciclu.

Atunci lățimea de bandă pentru un SMP-MP cu RIN-CB de dimensiuni $N \times M$ este:

$$BW = N \left[1 - (1 - pm) \left(1 - p \frac{1 - m}{N - 1} \right)^{N-1} \right] \quad (3.15)$$

Se furnizează de asemenea soluții pentru cazurile $M \leq N$ și $M \geq N$.

Dacă se substituie $m=1/M$, analiza se reduce la modelul MRU.

Ecuțiile sunt valabile doar în cazul sistemelor *sincrone* și cu comutare prin circuite.

În cazul sistemelor *asincrone* comutate prin circuite se utilizează modele în care procesoarele sunt modelate ca și un set de servere și memoriile ca și un set de șiruri de așteptare cu disciplină FIFO.

Modelul este prezentat în fig. 3.9. În Tabel II se prezintă sintetic diverse tehnici analitice de evaluare a RIN-CB, acuratețea modelării, metrica de evaluare a performanței și costurile de calcul implicate.

TABEL II. Rezumat al analizelor RIN - *crossbar*

	RIN- <i>crossbar</i> sincronă		RIN- <i>crossbar</i> asincronă
Tehnica de analiză	Șir Markov discret	Probabilistică cu ipoteză de independență	Rețea de șiruri de așteptare
Reprezentarea parametrilor	Rată de cerere acces	Probabilitatea de cerere acces	Timp de așteptare
Parametrii de performanță	BW sau P_A	BW sau P_A	P_u sau P_W
Acuratețea	Exactă	Bună	Exactă
Costul de calcul	Foarte ridicat	Coborât	Moderat

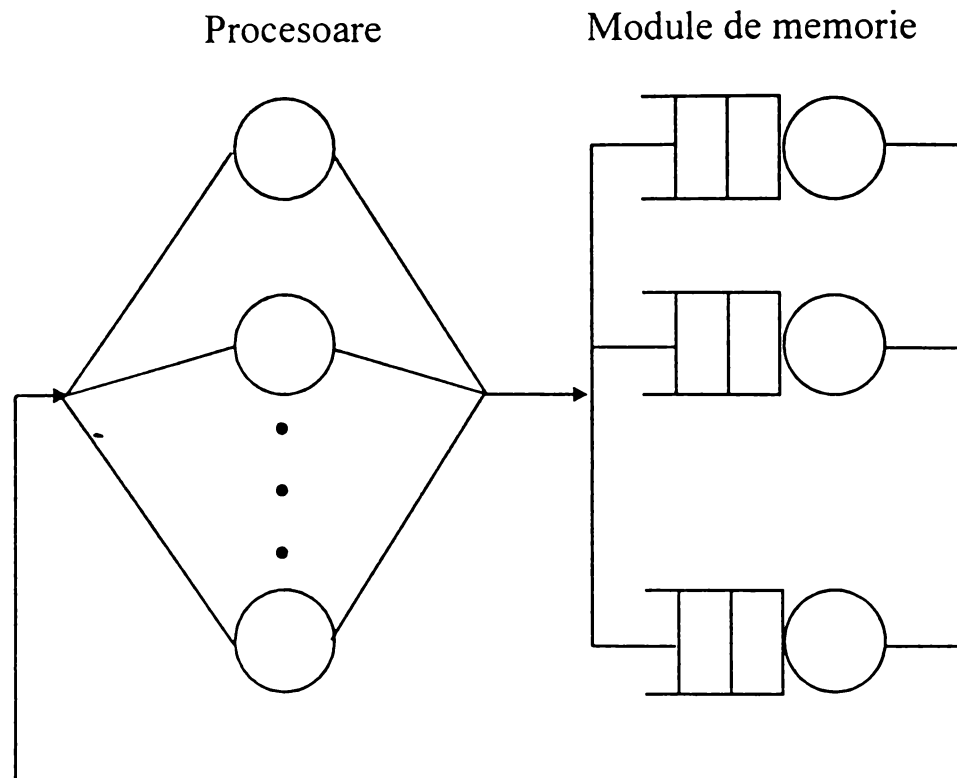


Fig. 3.9 Un model cu cozi de așteptare pentru SMP - MP bazate pe RIN - CB

3.1.5 Analiza performanțelor rețelelor de interconectare *multinivel*

Rețelele de interconectare *multinivel* (RIN-MN) au apărut relativ recent și au avantajul principal în raportul cost / performanțe acceptabil. Ele permit o multitudine de conexiuni unu la unu ale procesoarelor cu modulele de memorie, în timp ce costul hardware se reduce la $O(N \log N)$ în comparație cu $O(N^2)$ pentru o RIN - CB.

O RIN-MN de talie $N \times N$ conectează N procesoare la N memorii. Dacă N este o putere a lui 2, atunci se utilizează $\log_2 N$ etaje de comutatoare 2×2 cu $N/2$ comutatoare per nivel. Fiecare comutator are două intrări și două ieșiri. Conexiunea dintre o intrare și o ieșire este stabilită depinzând de un bit de control furnizat de intrare. Când $C=0$, intrarea este conectată la ieșirea superioară, iar când $C=1$, ea este conectată la ieșirea inferioară ca și în fig. 2.22.

O RIN-MN *omega*, prezentată în fig. 2.35, este caracterizată de o interconectare de tip "*amestec perfect*" ce precede fiecare nivel de comutatoare.

Procesorul solicitat generează o etichetă ce este o reprezentare binară a destinației. Conectarea unui comutator de pe nivelul i este apoi însoțită de către al i -lea bit a acestei etichete binare contorizat de la bitul cel mai semnificativ.

Conexiunea dintre intrarea 3 și ieșirea 5 (101_2) este prezentată prin linie întărită în fig. 2.35. Proprietatea de autorutare a RIN-MN evită necesitatea existenței unui dispozitiv central de comandă, ceea ce reprezintă un avantaj pentru SMP-MP. Prin urmare

discuțiile asupra performanțelor prezentate în această secțiune se vor concentra în mod unic asupra schemelor descentralizate de comandă.

Dimensiunea comutatorului în RIN-MN nu trebuie să fie restricționată la 2×2 .

În literatură [BYA89], [MD96], se citează cazul mașinii Butterfly ce conectează N intrări la N ieșiri utilizând comutatoare *crossbar* 4×4 și un număr de $\log_4 N$ de niveluri cu $N/4$ comutatoare per nivel.

O RIN-MN *delta* poate conecta $N = a^n$ intrări la $M = b^n$ ieșiri.

În general se pot conecta orice $M = m_1 * m_2 * \dots * m_r$ intrări la $N = n_1 * n_2 * \dots * n_r$ ieșiri prin intermediul a r niveluri de comutatoare.

Toate proiectele de sisteme bazate pe RIN-MN suferă însă de un dezavantaj major. Există doar o singură cale de la o intrare la o ieșire. Din acest motiv este necesar să se încorporeze câteva facilități de toleranță la defecte în aceste rețele astfel că cel puțin o defecțiune să fie tolerată într-un comutator sau într-un traseu de conectare.

În literatură [BYA 89] se sugerează o abordare probabilistică în analiza RIN-MN de tip "*delta*". Metoda presupune un *model de referință uniform* - MRU - și ipoteza *independenței cererilor de acces la memorie*.

Se consideră o RIN-MN *delta* de dimensiuni $a^n \times b^n$ constituită din n module *crossbar* $a \times b$. Fiecare nivel a RIN-MN *delta* este comandat de un bit cu destinație distinctă (codificat în baza b) pentru poziționarea comutatoarelor $a \times b$ individuale. Deoarece destinațiile sunt independente și uniform distribuite, cererile de acces la oricare modul $a \times b$ sunt independente și uniform distribuite pentru cele b destinații diferite.

Ecuția (3.14) poate fi aplicată la orice element în RIN-MN *delta*.

Numărul așteptat de solicitări ce traversează spre cele b ieșiri, este obținut punând $N = a$ și $M = b$ în ecuația (3.14).

Împărțind acest număr prin b se obține probabilitatea de solicitare a accesului a oricărei dintre cele b linii de ieșire ale unui comutator $a \times b$, ca și o funcție a probabilității de intrare. Deoarece ieșirea unui nivel este intrarea următorului nivel, se poate evalua de o manieră recursivă probabilitatea de ieșire a oricărui nivel pornind cu nivelul 1. Dacă p_i este probabilitatea ca să existe o solicitare la ieșirea unui comutator de pe nivelul i , atunci:

$$p_i = 1 - \left(1 - \frac{p_i}{b}\right)^a \quad (3.16)$$

pentru $1 \leq i \leq n$.

În particular, probabilitatea de ieșire a nivelului final determină lățimea de bandă a RIN-MN *delta*, ce este $BW = p_n b^n$.

Această tehnică analitică a fost utilizată pentru evaluarea diferitelor tipuri de RIN-MN.

În literatură [BYA 89] se prezintă și extinderea analizei pentru cazul *memoriei favorite*.

Pentru RIN-MN de dimensiuni $N \times N$, cu $N = a^n$, procesoarele sunt definite ca fiind conectate la memoriile lor favorite când toate comutatoarele sunt conectate direct, adică intrarea i a unui comutator este conectată la ieșirea i a circuitului.

Într-o RIN-MN *omega*, MM_j devine memoria favorită a procesorului P_j . Fie q_{i-1} probabilitatea că există o cerere favorită de acces la intrarea nivelului i . Într-o RIN-MN *delta* de dimensiuni $N \times N$, unde $N = a^n$, putem exprima:

$$p_i = 1 - (1 - p_{i-1} q_{i-1}) \left(1 - p_{i-1} \frac{1 - q_{i-1}}{a - 1} \right) \quad (3.17)$$

Mai sunt necesare aproximativ șase alte ecuații pentru a determina q_{i-1} la nivelul i^3 și, în final, $BW = N p_n$.

Analiza de mai sus este valabilă pentru RIN-MN sincrone cu comutare de pachete în următoarele condiții:

- pachetele sunt generate doar la începutul ciclului RIN-MN;
- comutatoarele nu au circuite tampon, astfel încât pachetele sunt aleatoriu alese în caz de conflicte și pachete ce nu au fost comutate sunt pierdute.

RIN-MN *delta* au mai fost studiate și cu ajutorul rețelelor Petri ce oferă facilități grafice pentru o descriere precisă a operațiilor de sistem. De asemenea performanțele pentru sisteme mici pot fi mai ușor evaluate.

Aceste modele bazate pe grafuri permit descrierea clară a concurenței, conflictelor și sincronizării task-urilor.

Oricum complexitatea rețelelor Petri crește exponențial cu creșterea în dimensiuni a sistemului.

3.1.6 Analiza comparativă a performanțelor RIN de tip *crossbar*, *magistrală multiplă și multinivel*

În finalul acestui paragraf se impune analiza comparativ a performanțelor pentru trei tipuri de RIN specifice SMP-MP.

Nu se ia în considerare *magistrala unică* deoarece performanțele ei sunt mult mai slabe comparativ cu acelea ale RIN indicate mai sus.

Analiza se bazează pe compararea costului circuitelor și a performanțelor sistemului pus în discuție.

Tabelul prezintă facilitățile hardware în mod selectiv pentru cele trei tipuri de RIN.

TABEL III. Sumarul facilităților hardware ale celor trei RIN

	<i>crossbar</i>	MN	MM
Nr. de comutatoare sau de conexiuni	$N \times M$	$N \log N$	$B \times (N+M)$
Încărcarea magistralelor	N	1	B
Nr. de fire	M	N	B
Arbitru	M "1 din N "	$N \log N$ "1 din 2"	1 " B din M " și M "1 din N "
Toleranța la defect și expandabilitate	Acceptabilă	Slabă	Bună

Numărul de elemente de comutare utilizate în RIN *crossbar* de dimensiuni $N \times M$ este $N \times M$ în comparație cu $N \log N$ ce sunt prezente în RIN – MN.

Numărul de conexiuni necesare într-o RIN – MM de dimensiuni $N \times M \times B$ este proporțional cu $B(N+M)$. Deoarece fiecare magistrală într-un SMP cu RIN – MM trebuie să conecteze $N+M$ module, încărcarea magistralei este proporțională cu $N+M$, în timp ce sarcina magistralei la o RIN – MN este 1 datorită conectării *unu la unu*.

Toate cele trei tipuri de RIN necesită arbitrii de complexități diferite pentru rezolvarea conflictelor de acces.

Expandabilitatea și siguranța în funcționare sunt două facilități hardware foarte importante. În acest context o RIN – MN este foarte convenabilă, în raport cu celelalte două, datorită proprietății de *reconfigurabilitate* și a *existenței căilor multiple de date* între fiecare procesor și memorie. Se mai poate opera chiar și în modul degradat după apariția unei defecțiuni a unui subset de magistrale.

În literatură [BYA 89] se prezintă rezultatele mai multor cercetări în domeniul RIN – MN tolerante la defect și reconfigurabile. Toleranța la defect a RIN – MN poate fi realizată prin adăugarea de circuite suplimentare ca de exemplu niveluri suplimentare

sau căi dublate de date. Este mai simplu să se reconfigureze un *crossbar* decât o RIN – MN.

În situația apariției unui defect, o anumită coloană sau o anumită linie pot fi îndepărtate și RIN poate opera în modul degradat.

În ceea ce privește situația comparativă a performanțelor celor trei tipuri de RIN, figurile 3.10, 3.11 și 3.12 sunt edificatoare.

În fig. 3.10, cele două curbe sunt trasate conform ecuațiilor (3.14) și (3.15). Distanța dintre cele două curbe crește pe măsură ce talia sistemului crește. P_A rămâne constant la *crossbar* atunci când talia SMP – MP devine foarte mare. P_A descrește la RIN – MN când talia sistemului crește.

În fig. 3.11 se evidențiază saturarea rapidă a BW la magistrala unică și creșterea lui în raport direct cu creșterea numărului de magistrale.

În fig. 3.12 prezintă comparativ comportarea RIN sincrone cu pachete comutate. Timpii ce caracterizează ciclul de procesor, întârzierea transferului de magistrală ca și timpul total de transfer al unui pachet ideal între o intrare și o ieșire al unui RIN – MN se presupun fixați în raport cu timpul ciclului de sistem. Timpul de acces la memorie se presupune egal cu patru cicluri de sistem. În timpul unui acces la memorie, procesorul ce trimite o solicitare de acces la memorie rămâne în așteptare până când accesul la memorie este finalizat.

Rata de solicitare a memoriei ce caracterizează un procesor, analizată prin prisma unui RIN, se prezintă în fig. 3.12. De remarcat că un SMP – MP bazat pe RIN – MM cu doar patru magistrale poate să realizeze aproximativ aceleași performanțe ca și un SMP – MP bazat pe RIN *crossbar* în timp ce costurile hardware sunt semnificativ diminuate.

Toate analizele anterior prezentate sunt pertinente în mediile sincrone cu comutare de circuite sau cu comutare de pachete. Pentru RIN - MN de talie mare, controlul centralizat, global al operațiilor rețelei este dificil. Prin urmare trebuie să se ia în considerare proiecte asincrone pentru SMP de talie mare.

Al doilea dezavantaj al acestei analize este ignorarea timpului de așteptare al unui procesor pentru epuizarea accesului la memorie, dar se presupune sosiri continue la intrare după o distribuție Poisson.

Al treilea dezavantaj este presupunerea unui acces la o memorie favorită. Configurațiile de referire la memorie sunt puternic dependente de program și pot fi arbitrare.

În literatură [BYA89], [YLL90], [YB91], se descrie metode bazate pe *cozi de așteptare* și pe *analiza a valorii medii* a RIN - MN operate asincron.

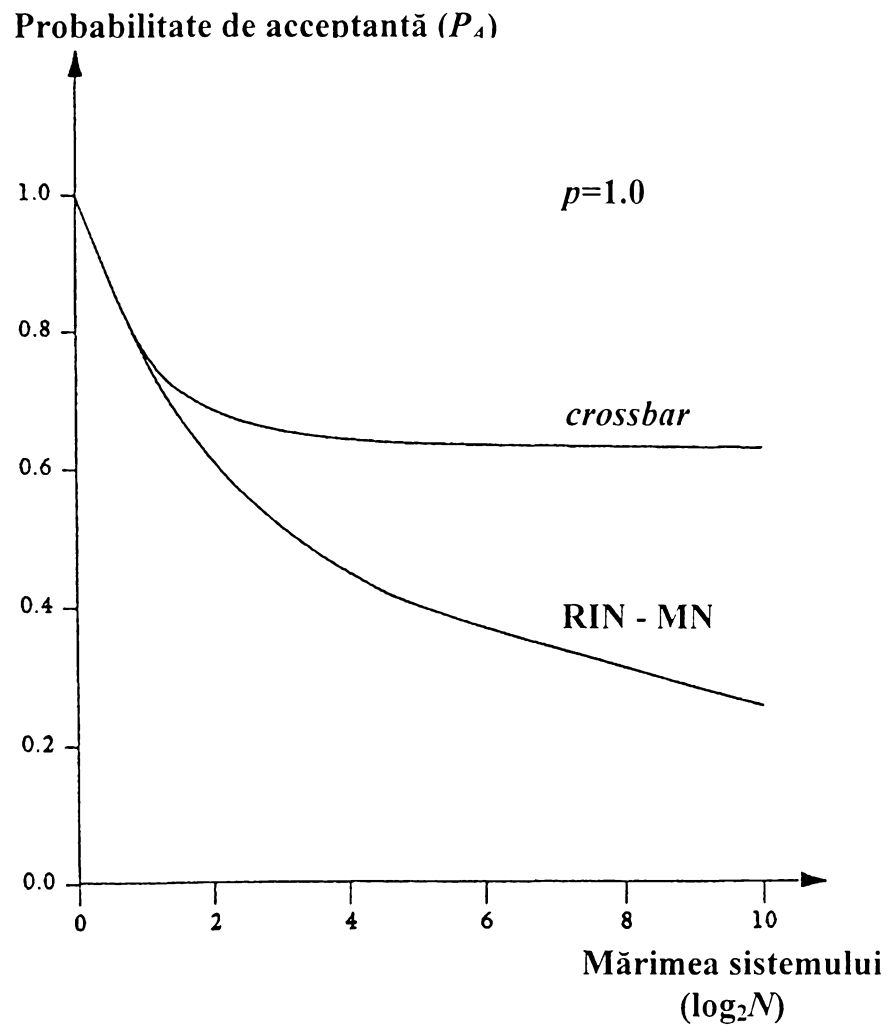


Fig. 3.10 Probabilitatea de acceptanță ca și funcție de mărimea sistemului pentru sisteme sincrone comutate prin circuit

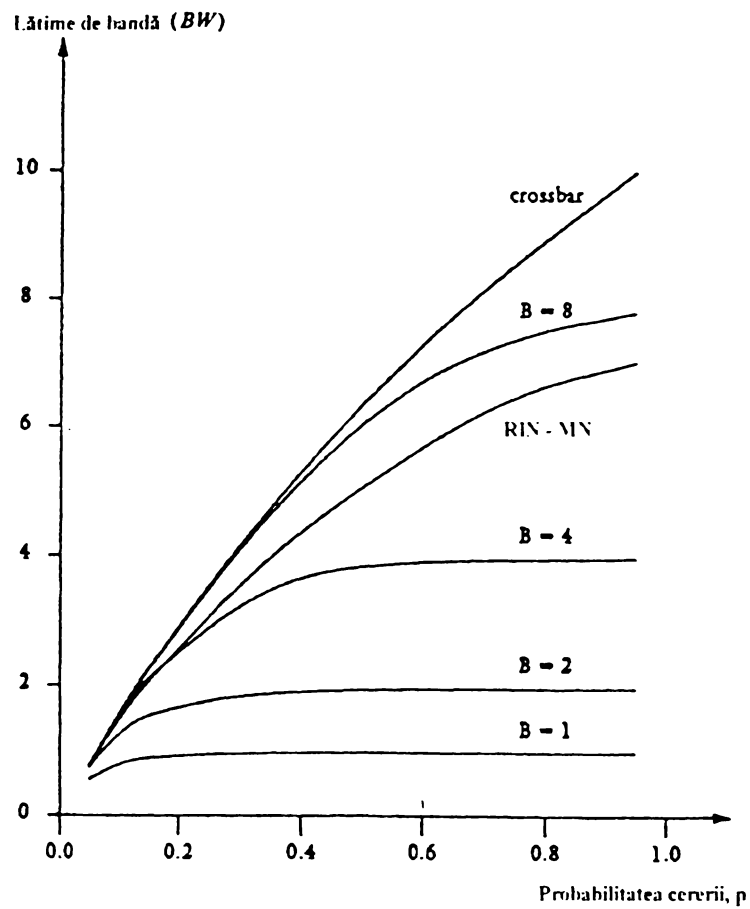


Fig. 3.11. Lățimea de bandă a memoriei ca și o funcție de probabilitate de solicitare de acces pentru un sistem sincron 16 x 16 comutat prin circuit

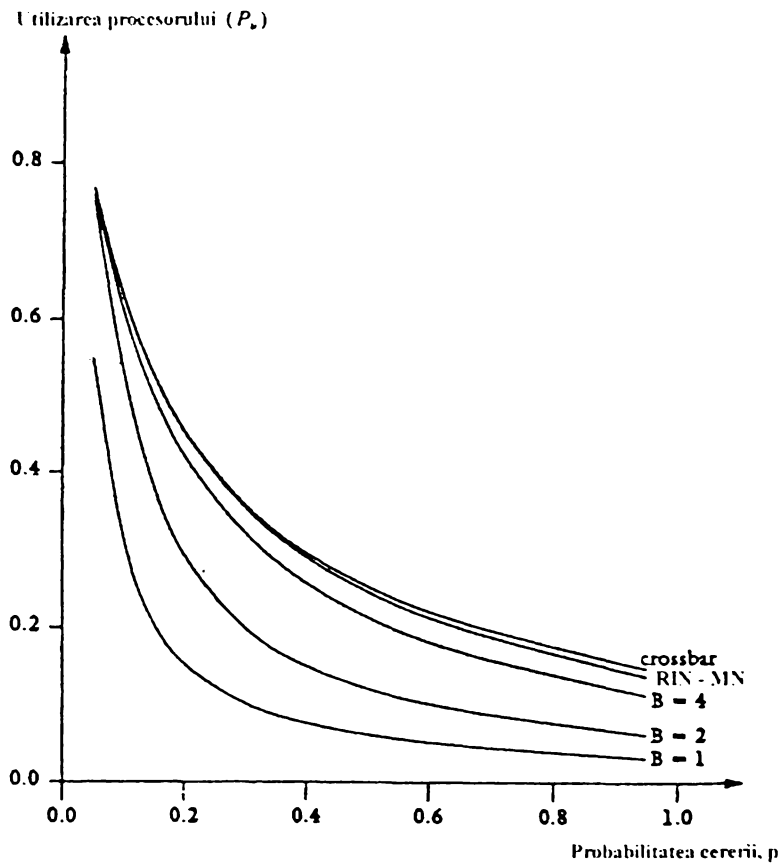


Fig 3.12 Utilizarea procesorului ca și o funcție a probabilităților de solicitare de acces pentru un sistem sincron 16 x 16 comutat prin pachet

3.2 ANALIZA DE PERFORMANȚĂ A SISTEMELOR DE MEMORII ÎN ARHITECTURILE SMP DIN CLASA MIMD

În capitolul 2 s-au analizat caracteristicile structurale ale arhitecturilor SMP din clasa MIMD grupate în subtaxonii MIMD-T și MIMD-L.

Reamintim că, în funcție de organizarea memoriei, sistemele din subtaxonul MIMD-T sunt cu *memorie partajată centralizată (centralized shared-memory)* în timp ce acelea din MIMD-L sunt cu *memorie partajată distribuită (distributed shared-memory)*.

SMP-MP se caracterizează prin faptul că au o singură memorie principală ce are un timp de acces uniform pentru fiecare procesor. Din acest punct de vedere ele sunt mașini cu *acces la memorie uniform (uniform memory access - UMA)*.

Structura de bază a unui SMP-MP este prezentată în fig. 3.13 Se remarcă faptul că mai multe subsisteme procesor - cache împart aceeași memorie fizică, uzual conectate printr-o magistrală.

SMP-MP sunt mașini cu memorie fizic distribuită. Distribuirea memoriei este necesară pentru a se putea conecta un număr mare de procesoare ce reclamă o lățime de bandă mare.

Distribuirea memoriei printre noduri oferă beneficii majore. Primul, este o cale cu costuri reduse pentru a scala memoria, dacă cele mai multe dintre accese sunt la memoria locală din nod. Al doilea, se reduce latența pentru accese la memoria locală.

Dezavantajul major îl constituie faptul că procedura de comunicare între procesoare devine mai complexă.

Arhitectura de bază a unui sistem SMP-MP este prezentată în fig. 3.13

În aceste arhitecturi nodurile de comunicație pot conține unul sau mai multe procesoare împreună cu memoria și circuitele de interfață cu RIN. Se pot crea *grupuri (cluster)* de procesoare conectate pe mai multe niveluri de RIN.

Diferențele arhitecturale majore ce se disting printre mașinile cu memorie distribuită constau în modul în care se produce comunicarea și arhitectura logică a memoriei distribuite.

Există două abordări arhitecturale alternative în ceea ce privește comunicarea datelor între procesoare.

Prima constă în faptul că memoriile fizic separate pot fi adresate ca și un singur spațiu de adresă logic partajat, ceea ce înseamnă că orice procesor poate accesa orice locație de memorie, dacă condițiile de acces sunt îndeplinite.

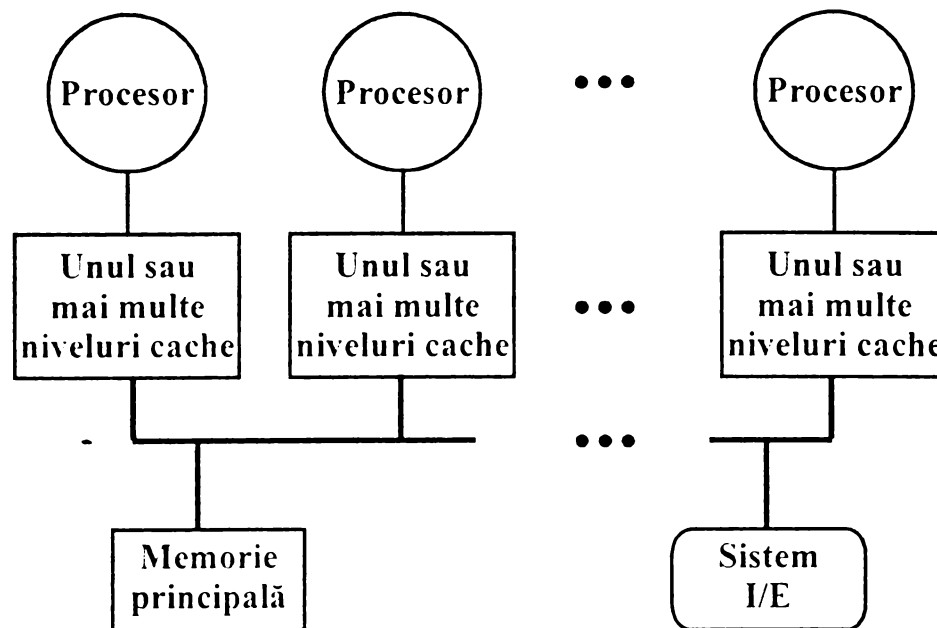


Fig. 3.13 Structura de bază a unui SMP-MP

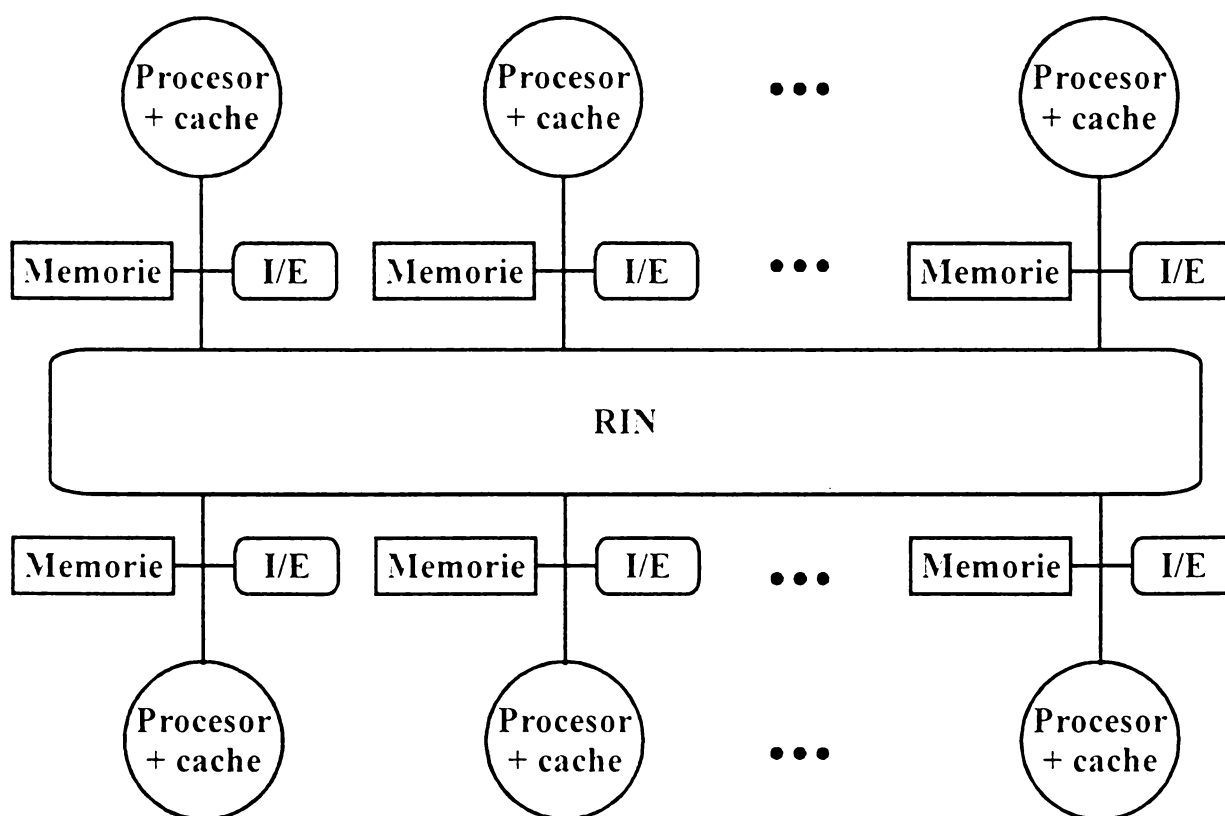


Fig. 3.14 Arhitectura de bază pentru un SMP-TM, cu memorie partajată distribuită

Denumirea de *memorie partajată* se referă la faptul că *spațiul de adresă* este partajat: adică aceeași adresă fizică solicitată de două procesoare se referă la aceeași locație în memorie. Memoria partajată nu înseamnă că există o singură memorie centrală. Prin

urmare sistemele cu memorie partajată pot fi cu *acces uniform la memorie (uniform memory access - UMA)* sau cu *acces neuniform la memorie (non - uniform memory access - NUMA)*.

Soluția alternativă pentru comunicarea datelor o oferă spațiul de adresă ce constă din spații de adresă private multiple ce sunt logic disjuncte și nu pot fi adresate de către un procesor de la distanță. La aceste mașini, aceeași adresă fizică în două procesoare diferite se referă la două locații diferite în două memorii diferite. Fiecare modul procesor - memorie este un calculator separat și aceste sisteme se mai numesc *multicalculatoare*. Sunt arhitecturile clasice pentru SMP-TM !

În paragraful 3.1 s-au studiat performanțele RIN.

În acest paragraf ne vom focaliza eforturile pentru evidențierea performanțelor oferite de memoriile cache în diferite configurații.

În literatură [HP 96] se indică mai multe procedee de-a evalua performanțele diferitelor arhitecturi de SMP. De remarcat că există multe aplicații ale calculului paralel dintre care amintim : transformarea Fourier rapidă (FFT), descompunerea LU, aplicații ale algoritmului Barnes-Hut pentru rezolvarea problemelor legate de evoluția galaxiei, tehnici Gauss-Seidel pentru rezolvarea ecuațiilor diferențiale eliptice etc.

O caracteristică cheie în determinarea performanțelor programelor paralele îl constituie raportul dintre calcul și comunicare.

Dacă raportul este mare înseamnă că aplicația are o mulțime de calcule pentru fiecare dată comunicată.

Este situația ideală. Într-un mediu de procesare paralelă se urmărește tot timpul cum se modifică raportul calcul/comunicare atunci când se adaugă noi procesoare, se mărește dimensiunea problemei sau se manifestă ambele tendințe.

În TABEL IV se prezintă modul de evoluție al raportului calcul/comunicare în funcție de diferite tipuri de aplicații.

Se notează :

p - numărul de procesoare;

n - mărimea setului de date.

TABEL IV. Scalarea componenteii de calcul, a celei de comunicare și a raportului calcul/comunicare

Aplicație	Scalarea calculului	Scalarea comunicării	Scalarea calcul/comunicare
FFT	$\frac{n \log n}{p}$	$\frac{n}{p}$	$\log n$
LU	$\frac{n}{p}$	$\frac{\sqrt{n}}{\sqrt{p}}$	$\frac{\sqrt{n}}{\sqrt{p}}$
Barnes	$\frac{n \log n}{p}$	$\frac{\sqrt{n(\log n)}}{\sqrt{p}}$	$\frac{\sqrt{n}}{\sqrt{p}}$
Ocean	$\frac{n}{p}$	$\frac{\sqrt{n}}{\sqrt{p}}$	$\frac{\sqrt{n}}{\sqrt{p}}$

3.2.1 Arhitecturi cu memorie partajată centralizată

Aceste arhitecturi au dominat perioada anilor '80 datorită evoluției spectaculoase a circuitelor VLSI în general și a microprocesoarelor în special.

Caracteristica principală - din punct de vedere al organizării ierarhiei memoriilor - constă în plasarea nivelului superior format din cache-uri pe aceleași module cu procesoarele și a nivelului inferior - memorie RAM principală - modul (e) separat (e).

Arhitectura SMP-MP suportă plasarea în memorii cache atât a *datelor private* cât și a celor *partajate*. Datele private sunt utilizate de un singur procesor în timp ce datele partajate sunt utilizate de mai multe procesoare, asigurându-se comunicarea între procesoare prin citiri și scrieri ale datelor partajate. Când o dată privată este plasată în memoria cache, locația sa *migrează* în cache, reducându-se timpul de acces mediu atât cât o solicită lățimea de bandă a memoriei.

Deoarece nici un alt procesor nu utilizează aceste date private, comportamentul programului este identic cu acela dintr-un uniprocessor.

Când datele partajate sunt plasate în cache, valoarea partajată trebuie să fie *replicată* în mai multe memorii cache. În plus față de reducerea latenței accesului și a lățimii de bandă a memoriei, această replicare de asemenea furnizează o reducere a conflictelor ce pot apărea în procesul de citire a datelor partajate de către mai multe procesoare simultan.

Menționăm că plasarea în memorii cache a datelor partajate introduce problema *coerenței memoriei cache*.

Performanțelor arhitecturilor cu memorie partajată sunt în relație directă cu tehnicile de menținere a coerenței utilizate în acestea.

Introducerea memoriilor cache într-o ierarhie de memorii cauzează apariția problemei coerenței pentru operații de I/E deoarece viziunea asupra memorie prin cache poate fi diferită de viziunea asupra memoriei obținută prin subsisteme de I/E.

Figura 3.15 ilustrează un exemplu simplu de apariție a incoerenței.

Timp	Eveniment	Conținutul cache pentru CPU A	Conținutul cache pentru CPU B	Conținutul memoriei la adresa X
0				1
1	CPU A citește X	1		1
2	CPU B citește X	1	1	1
3	CPU A scrie 0 în X	0	1	0

Fig. 3.15 Problema coerenței cache pentru o locație de memorie, citită și scrisă de două procesoare

Reamintim că toate operațiile de citire/scriere în memoria partajată situată pe nivelul inferior în cadrul unei ierarhii de memorii se fac prin intermediul cache-urilor plasate pe nivelul superior, logic și fizic amplasat lângă procesoare.

Un sistem de memorie este coerent dacă :

1. O operație de citire efectuată de un procesor P_i , la o locație X ce urmează după o scriere de către P_i în X, fără să existe alte scrieri ale lui X de către alt procesor ce să apară între scrierea și citirea efectuate de P_i , întotdeauna returnează valoarea scrisă de P_i .
2. O citire efectuată de către un procesor P_i la o locație X ce urmează o scriere efectuată de alt procesor P_j în X, returnează valoarea scrisă dacă citirea și scrierea sunt separate suficient și nu apar alte scrieri în X între cele două accese.
3. Scrierile în aceeași locație de memorie sunt *serializate* : aceasta înseamnă că două scrieri în aceeași locație efectuate de către oricare alte două procesoare sunt privite în aceeași ordine de către toate procesoarele.

Modelul de consistență a memoriei oferă informații despre momentul și condițiile în care o valoare scrisă trebuie să fie văzută de către utilizator.

Coerența și consistența sunt complementare :

Coerența definește comportamentul scrierilor și citirilor în aceeași locație de memorie, în timp ce *consistența* definește comportamentul citirilor și scrierilor în raport cu accesul la alte locații de memorie.

În capitolul 2, paragraful 2.3, sunt descrise tehnici și scheme de bază pentru întărirea coerenței.

În acest capitol vom dezvolta tehnici și scheme de asigurare a coerenței la memoriile cache (SAC-MCH) specifice arhitecturilor cu memorie partajată.

A. Scheme de asigurare a coerenței memoriilor cache (SAC-MCH)

Într-un SMP-MP cu coerență asigurată, memoriile cache asigură atât *migrarea* cât și *replicarea* datelor partajate. Memoriile cache coerente asigură migrarea, deoarece o dată poate fi deplasată într-un cache local și utilizată apoi într-o manieră transparentă; această procedură reduce latența la accesul unei date partajate ce este alocată unui procesor fizic situat la distanță. De asemenea MCH asigură replicarea pentru datele partajate ce sunt simultan citite, deoarece MCH produce o copie a datei partajate în memoria cache locală (MCH – L). Replicarea reduce atât latența accesului cât și conflictul la o dată partajată ce trebuie citită.

Migrarea și replicarea sunt subiecte critice pentru stabilirea criteriilor de performanță în accesarea datelor partajate.

Prin urmare în SMP-MP de talie redusă se adoptă soluții hardware prin introducerea de protocoale de menținere a coerenței.

În literatură [HP 96] sunt descrise protocoale specifice pentru asigurarea coerenței memoriilor cache la SMP-MP.

Protocoalele de menținere a coerenței la memoriile cache (cache - coherence protocols) utilizează tehnica trasării stării pentru orice partajare a unui bloc de date.

În literatură [HP 96], [HP 94], [SLM 89] etc, se inventariază mai multe protocoale dar ele pot fi grupate în două clase ce utilizează diferite tehnici de trasare a stării de partajare.

Aceste sunt :

- *Bazate pe director (directory - based)* - starea de partajare a unui bloc de memorie este menținută într-o locație denumită *director*; este tipică pentru arhitecturile cu memorie partajată scalabile, adică SMP-MP.
- *Monitorizare (snooping)* - Fiecare cache ce are o copie a datelor provenind de la un bloc din memoria comună partajată are de asemenea o copie a stării de partajare a blocului. Nici o stare centralizată nu este păstrată. Toate controlerile atașate memoriilor cache *monitorizează (snoop)* magistrala partajată pentru a determina dacă au sau nu o copie a blocului ce este solicitat pe magistrală. Aceste protocoale bazate pe monitorizare sunt populare pentru SMP-MP deoarece aceste arhitecturi utilizează microprocesoare și cache-uri atașate unei singure memorii partajate și pot utiliza conexiuni fizice pre-existente de la magistrală la memorie.

Există două metode de a menține coerența.

Prima metodă constă în asigurarea că un procesor are un acces exclusiv la o dată înainte ca ea să fie scrisă.

Protocolul de invalidare a scrierii (write invalidate protocol) se bazează pe invalidarea altor copii la o operație de scriere. Este cel mai utilizat protocol atât pentru scheme bazate pe director cât și pentru cele bazate pe monitorizare.

Accesul exclusiv asigură că nu există alte copii de citit sau de scris ale unei date atunci când apare scrierea : toate celelalte copii din cache-uri ale datei implicate sunt invalidate.

Acest protocol asigură coerența astfel. Se consideră o scriere urmată de-o citire de către alt procesor : deoarece scrierea solicită un acces exclusiv, orice copie deținută de procesorul efectuând citirea trebuie să fie *invalidată*. Astfel când apare citirea, copia lipsește din cache și se forțează căutarea unei noi copii a datei.

Pentru o scriere, este necesar ca procesorul ce scrie să aibă un acces exclusiv, prevenind orice alt procesor de-a ajunge în starea de-a scrie simultan. Dacă două procesoare intenționează să scrie aceeași dată simultan, unul dintre ele va câștiga cursa provocând invalidarea copiei celuilalt procesor. Pentru ca celălalt procesor să-și completeze scrierea, este necesar să obțină o copie nouă a datei, care trebuie să conțină acum data actualizată.

Prin urmare acest protocol se bazează pe *serializarea scrierii*.

În fig. 3.16 se prezintă un exemplu a unui protocol de invalidare pentru o magistrală cu linie de monitorizare utilizând rescrierea în memoriile cache.

Activitate a procesorului	Activitate a magistralei	Conținut cache CPU A	Conținut cache CPU B	Conținut locație X din memorie
				0
CPU A citește X	Lipsă cache pentru X	0		0
CPU B citește X	Lipsă cache pentru X	0	0	0
CPU A scrie 1 în X	Invalidare pentru X	1		0
CPU B citește X	Lipsă cache pentru X	1	1	1

Fig. 3.16 Protocol de *invalidare* utilizat pe o magistrală cu monitorizare pentru un singur bloc cache utilizând *rescrierea*

Protocolul de actualizare la scriere (*write update protocol, write broadcast protocol*) reprezintă o alternativă la protocolul de invalidare a scrierii și constă în actualizarea tuturor copiilor deținute în memoriile cache ale datei ce este scrisă.

Pentru a se păstra sub control necesitățile dictate de lățimea de bandă, este util să se traseze dacă un cuvânt din cache este sau nu partajat, adică dacă este conținut în alte memorii cache. Dacă nu este conținut în alte MCH, atunci nu este necesar să se actualizeze oricare alte MCH.

În fig. 3.17 se prezintă un exemplu de *protocol de actualizare la scriere*. Acest protocol este mult mai puțin utilizat față de cel de invalidare a scrierii datorită diferențelor semnificative în performanțe.

Activitate a procesorului	Activitate a magistralei	Conținut cache CPU A	Conținut cache CPU B	Conținut locație X din memorie
				0
CPU A citește X	Lipsă cache pentru X	0		0
CPU B citește X	Lipsă cache pentru X	0	0	0
CPU A scrie 1 în X	Actualizare X	1	1	1
CPU B citește X		1	1	1

Fig. 3.17 Protocol de *actualizare la scriere* pe o magistrală cu monitorizare pentru un singur bloc cache cu *rescriere*

Diferențele în performanțe dintre cele două protocoale apar datorită a trei caracteristici :

1. Scrieri multiple la același cuvânt, fără citiri ce să intervină, pretind transmisii multiple de scrieri într-un *protocol de actualizare*, dar doar o singură invalidare inițială într-un *protocol de invalidare de scriere*.
2. În cazul existenței de blocuri cache având cuvinte multiple, fiecare cuvânt scris într-un bloc cache pretinde o transmisie a semnalului de scriere într-un *protocol de actualizare*, în timp ce doar prima scriere la orice cuvânt din bloc necesită să se genereze o invalidare într-un *protocol de invalidare*. Un *protocol de invalidare* acționează asupra blocurilor cache, în timp ce un *protocol de actualizare* trebuie să acționeze asupra cuvintelor individuale.
3. Întârzierea dintre scrierea unui cuvânt într-un procesor și citirea valorii scrise în alt procesor este uzual mai mică într-o schemă de actualizare la scriere, deoarece data scrisă este imediat actualizată în memoria cache a procesorului ce citește.

Cele mai utilizate protocoale sunt cele bazate pe *invalidare*.

B. Tehnici de implementare de bază

Cheia implementării unui protocol de invalidare într-un SMP-MP de talie redusă o constituie utilizarea magistralei pentru efectuarea operațiilor de invalidare.

Pentru a efectua o invalidare procesorul obține accesul la magistrală și transmite adresa pentru a fi invalidată pe magistrală. Toate procesoarele monitorizează continuu magistrala urmărind adresele. Procesoarele verifică dacă adresa de pe magistrală este în cache-urile proprii. În caz afirmativ, data corespunzătoare din cache este invalidată. Semnalizarea accesului – asigurată de către magistrală – de asemenea forțează serializarea scrierilor, deoarece când două procesoare concurează pentru a scrie în aceeași locație, unul trebuie să obțină accesul la magistrală înaintea celuilalt. Primul procesor ce a câștigat competiția va provoca invalidarea copiei celuilalt procesor, obligând serializarea scrierilor. O consecință a acestei scheme este că o scriere la o dată partajată nu poate fi efectuată până când nu se obține accesul la magistrală.

Pentru un cache cu rescriere problema găsirii celei mai recente date este mai dificilă, deoarece cea mai recentă dată poate fi mai probabil în cache decât în memorie. Fiecare procesor monitorizează fiecare adresă plasată pe magistrală. Dacă un procesor

constată că are o copie viciată a blocului din cache solicitat, el furnizează acest bloc ca și răspuns la o solicitare de citire și cauzează abandonarea accesului la memorie.

Deoarece MCH cu rescriere nu au nevoie de lățime de bandă semnificativă pentru memorie sunt preferate în SMP.

În literatură [HP 96] sunt descrise tehnici de implementare a *MCH multinivel cu rescriere (write - back multilevel caches)*.

Un protocol de menținere a coerenței bazat pe magistrală se implementează uzual incorporând un controler cu stări finite în fiecare nod. Acest controler răspunde la solicitările procesoarelor și ale magistralei, schimbând starea blocului cache selectat, utilizând magistrala pentru a avea acces la magistrală sau pentru a o invalida.

În fig. 3.18 se prezintă cereri generate de către un modul procesor-cache într-un nod cât și acelea generate de magistrală. În acest protocol nu se face distincție între un succes la scriere (write hit) și un eșec la scriere (write miss) la un bloc cache partajat. În ambele cazuri un astfel de acces este tratat ca și un eșec la scriere. Când o astfel de scriere este plasată pe magistrală, oricare dintre procesoare cu copii ale blocului cache le invalidează.

Solicitare	Sursă	Funcție
Succes la citire	Procesor	Citește data în cache
Succes la scriere	Procesor	Scrie data în cache
Eșec la citire	Magistrală	Solicită data din cache sau din memorie
Eșec la scriere	Magistrală	Solicită data din cache sau memorie. Invalidează

Fig. 3.18 Mecanism de asigurare a coerenței cu solicitări generate de modulul procesor-cache

În fig. 3.19 se prezintă o diagramă de tranziție a unui automat cu stări finite pentru un singur bloc cache ce utilizează un protocol de invalidare la scriere și un cache cu rescriere. Din rațiuni de simplitate cele trei stări ale protocolului sunt dublate pentru a reprezenta tranziții bazate pe solicitări din CPU (în partea stângă) în opoziție cu tranziții bazate pe solicitări din magistrală (partea dreaptă). Starea în fiecare nod reprezintă starea blocului selectat specificat de către procesor sau de către cererea provenind din magistrală. Există doar un automat cu stări finite per cache, cu stimuli provenind fie din CPU atașat, fie din magistrală.

În fig.3.20 se prezintă combinația dintre diagramele prezentate anterior, pentru a forma o singură diagramă de stare pentru fiecare cache.

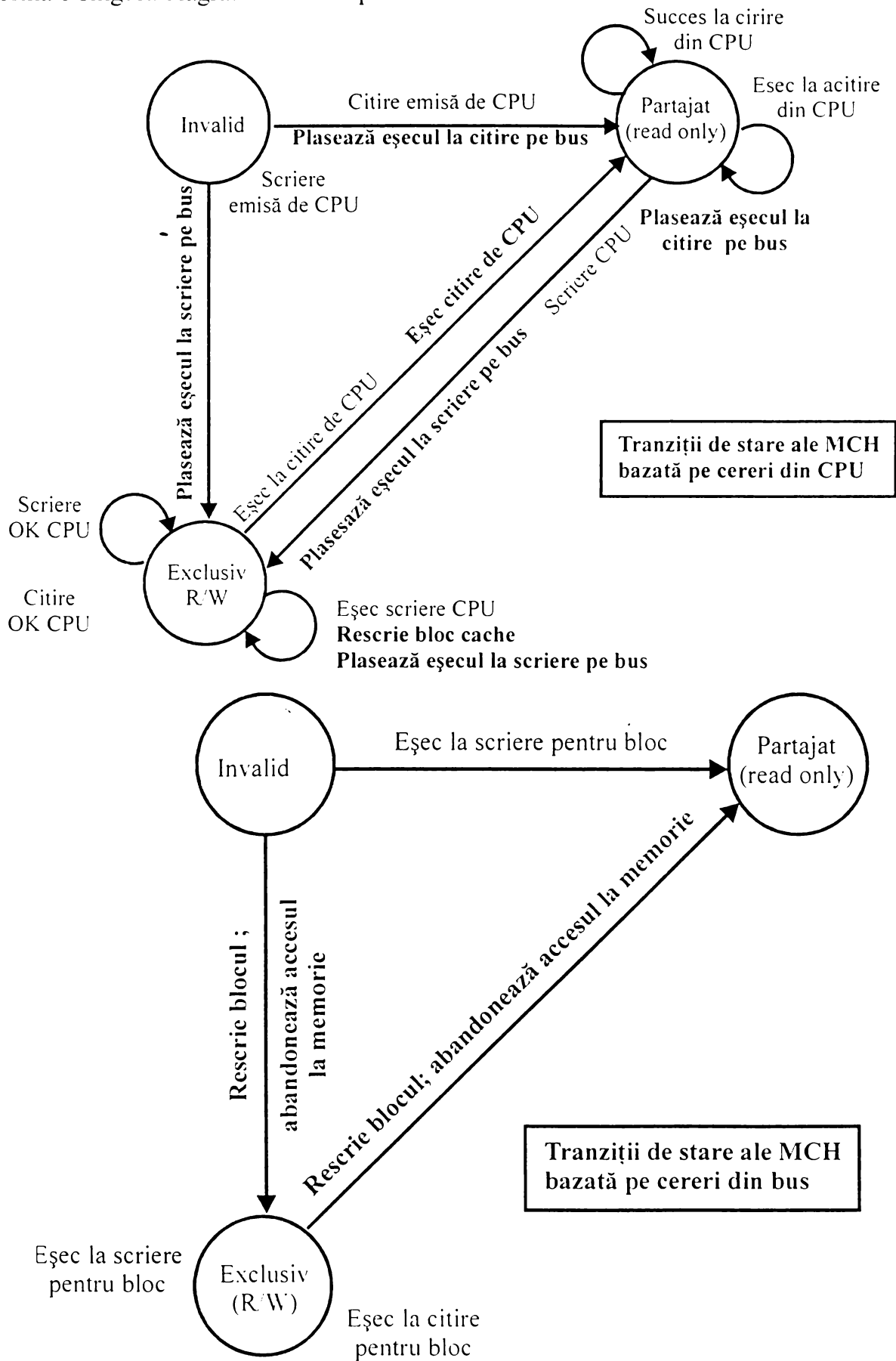


Fig.3.19 Un protocol cu invalidare la scriere pentru un cache cu rescriere prezentat ca și un automat finit

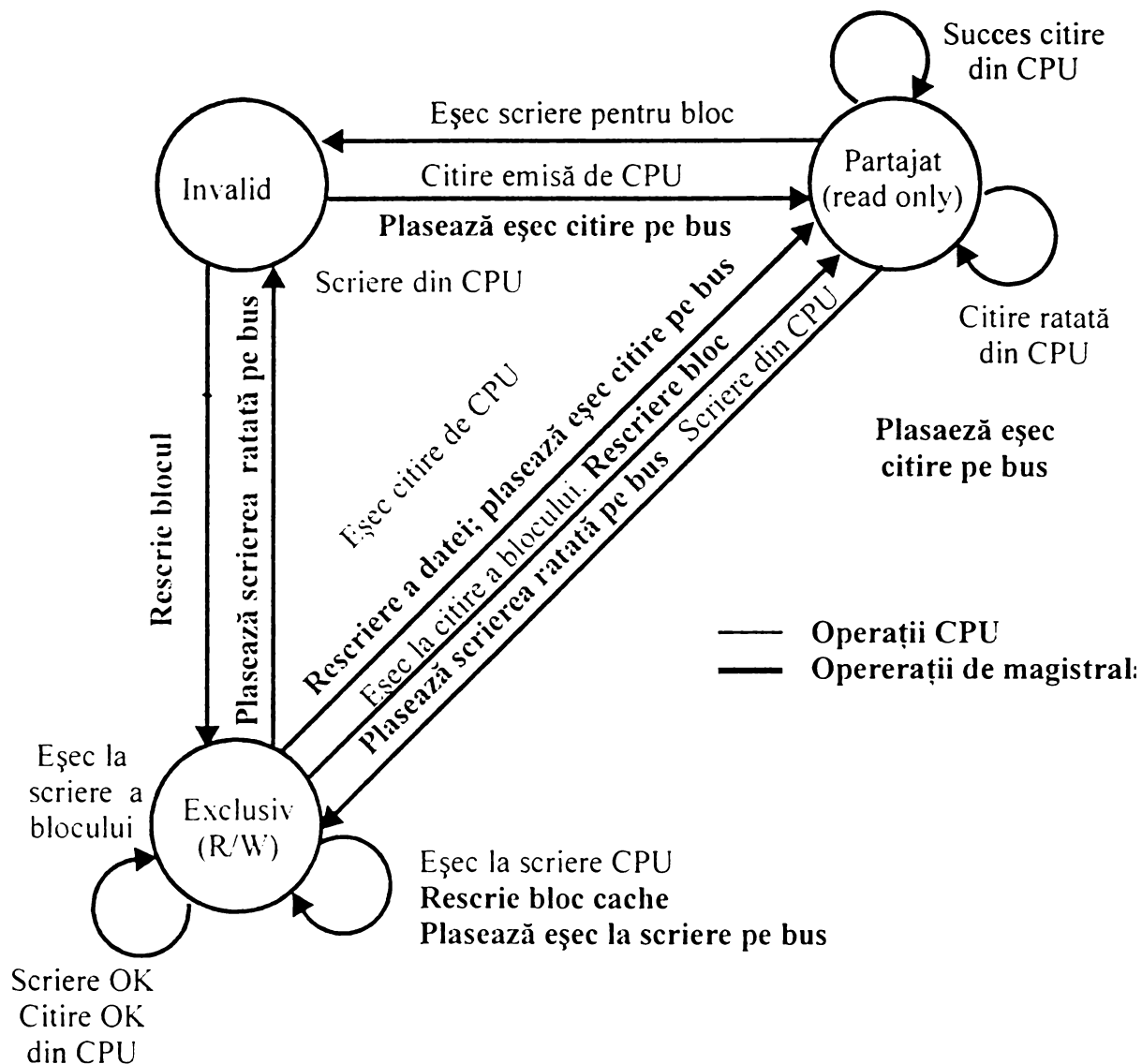


Fig. 3.20 Diagrama de stare a coerenței MCH cu tranziții de stare induse de procesorul local și activitățile magistralei

În implementarea acestor protocoale se presupune - ca ipoteză simplificatoare - că operațiile sunt *atomice*. De exemplu se consideră că ratările la scriere pot fi detectate, magistrala câștigată și un răspuns obținut toate fiind operații indivizibile. În realitate nu este așa. Modificând protocolul pentru a prezerva corectitudinea, acțiunile non atomice introduc posibilitatea ca protocolul să se poată *bloca (deadlock)*, adică automatul finit ce-l implementează să atingă o stare de unde să nu mai poată continua. În literatură [HP 96] se studiază performanțele MCH ce contribuie masiv la performanțele globale ale SMP. Se utilizează simulatoare pe care rulează *programe de etalonare (benchmark programs)*.

Se evaluează *ratele de eșecuri (miss rates)* ca și funcții de *numărul de procesoare (processor count)*, *talia MCH (cache size)* și *talia blocului de date (block size)*.

Rata totală de eșec se descompune în *eșecuri datorate coerenței (coherence misses)* și *eșecuri de uniprocessor normale (normal uniprocessor misses)*.

Eșecurile de uniprocessor normale constau din *eșecuri datorită capacității (capacity misses)*, *eșecuri datorită conflictului (conflict misses)* și *eșecuri obligatorii (compulsory misses)*.

De obicei aceste eșecuri se etichetează ca fiind *eșecuri datorate capacității* deoarece aceasta este dominantă în programele de etalonare.

În literatură [HP 96] se prezintă comparativ fiecare din acești parametri luând în considerare mai multe tipuri de aplicații pentru un parametru.

Se va prezenta doar o aplicație - analiză Fourier rapidă (FFT) - ce constă în prelucrarea a 1 milion de puncte de date reprezentate în complex.

Se observă că :

- rata de eșec datorită coerenței crește cu numărul de procesoare; (fig. 3.21);
- ratele de eșec datorită coerenței și capacității scad atunci când talia memoriei cache crește; (fig.3.22);
- ratele de eșec scad semnificativ atunci când talia blocului de memorie crește. (fig.3.23).

În fig. 3.24 se prezintă variația traficului pe magistrală în funcție de talia blocului. El crește odată cu talia blocului.

De asemenea rata de eșec scade cu creșterea dimensiunii memoriei cache.

În fig. 3.25 se prezintă evoluția componentelor ratei de eșec în funcție de talia memoriei cache și, respectiv, în funcție de talia blocului.

C. Concluzii referitoare la performanțele schemelor de asigurare a coerenței MCH utilizând monitorizarea

Traficul datorat coerenței poate introduce comportamente noi în sistemul de memorie ce nu răspund ușor la modificările taliei memoriei cache sau a taliei blocului ce sunt normal utilizate pentru a ameliora performanțele MCH la unipresoare.

În cadrul procesării paralele cerințele de menținere a coerenței sunt semnificative dar nu în mod exagerat.

În cadrul multiprogramării, utilizatorul și sistemul de operare (SO) evoluează foarte diferit.

În porțiunea de SO, combinațiile componentelor obligatorii și respectiv datorate capacității în ceea ce privește rata de eșec sunt mult mai semnificative. Aceste rate de eșec afectează performanțele CPU depinzând de restul sistemului de memorie, incluzând latența și lățimea de bandă a magistralei și a memoriei.

Graficele sunt prezentate în fig.3.28, 3.29, 3.30.

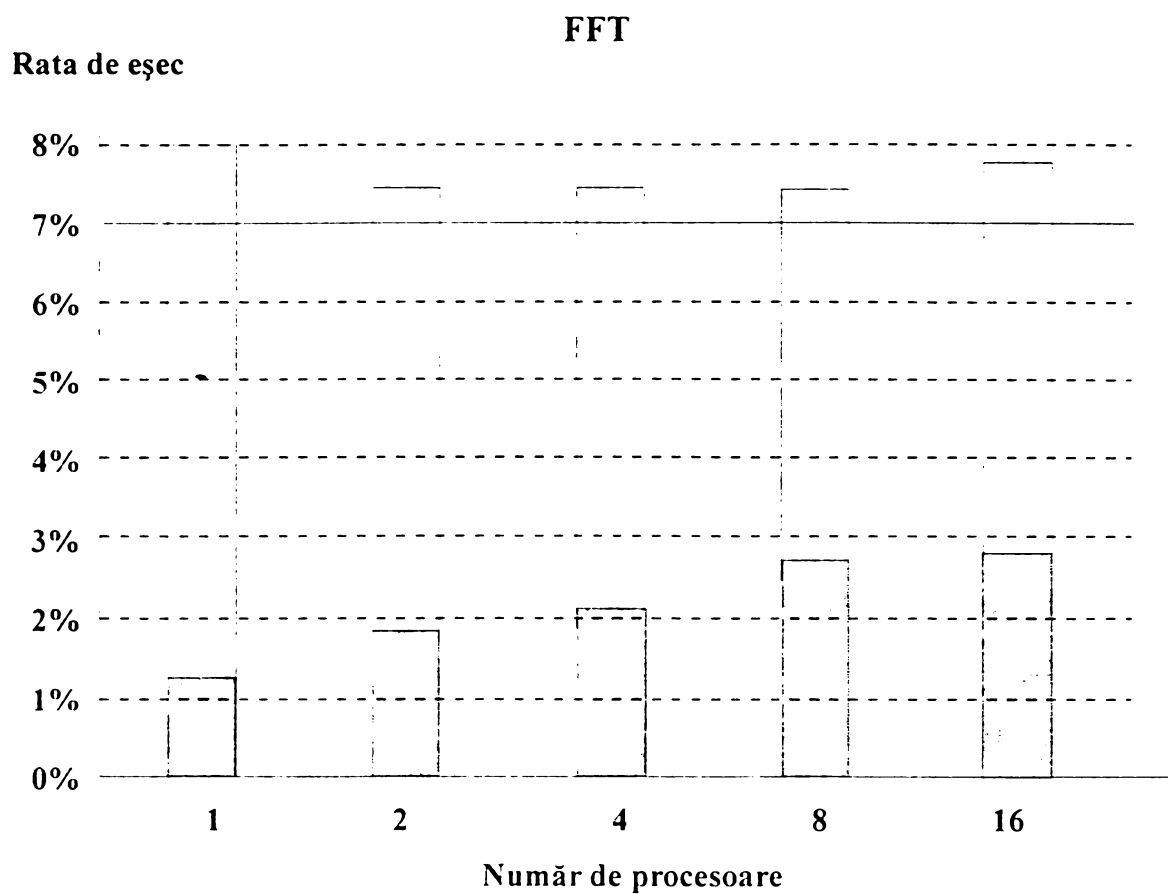


Fig. 3.21 Variația ratelor de eșec în funcție de numărul de procesoare

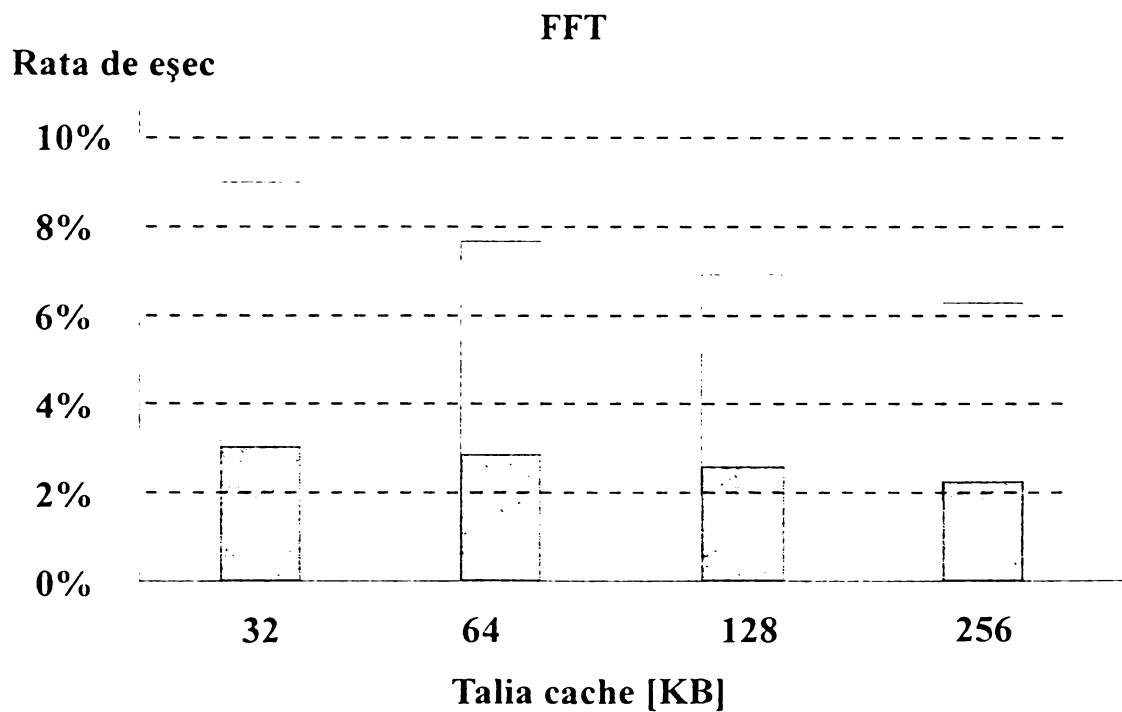


Fig. 3.22 Variația ratei de eșec funcție de talia MCH

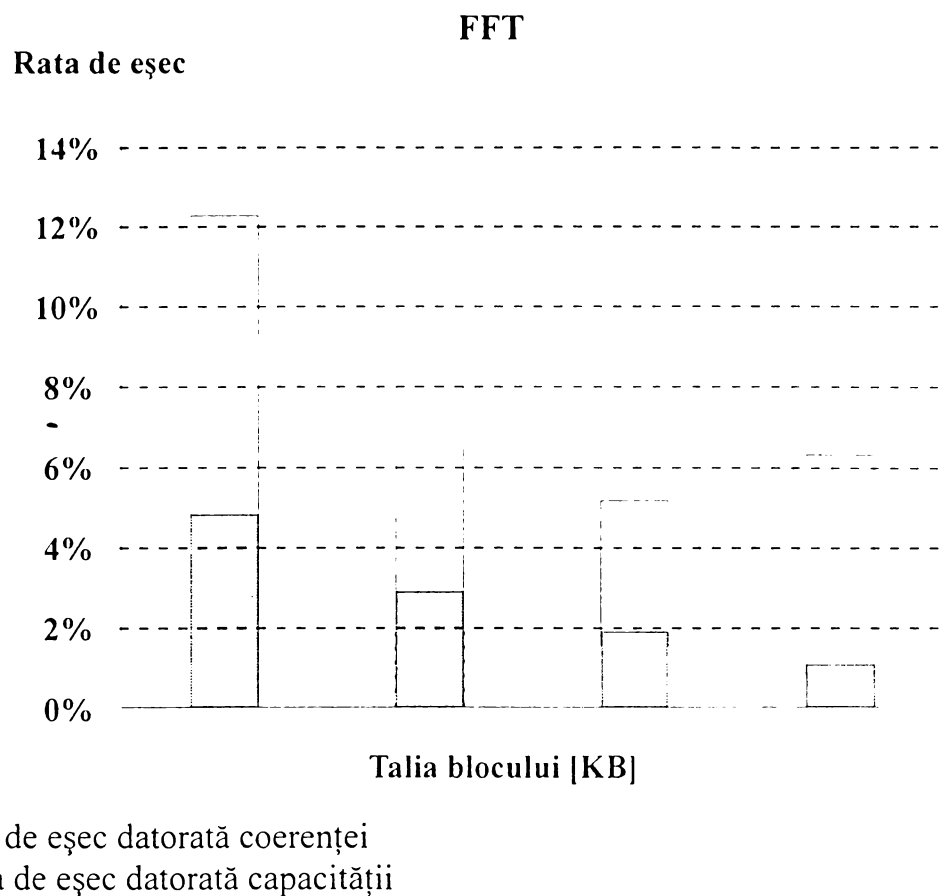


Fig. 3.23 Rata de eșec în funcție de talia blocului

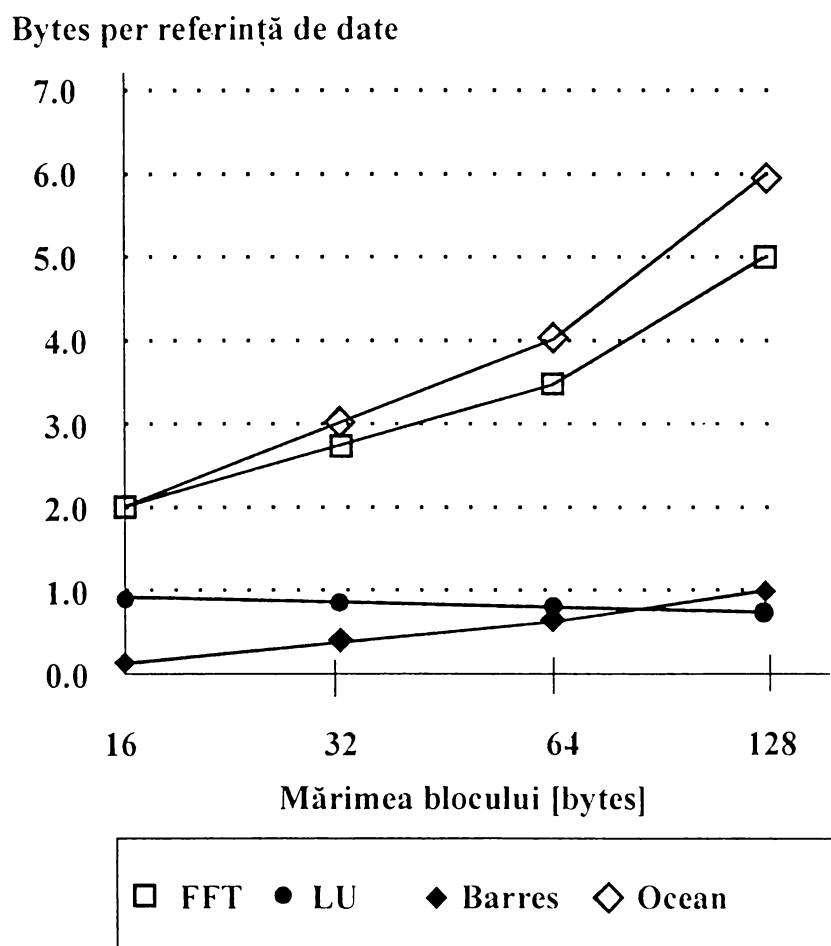


Fig. 3.24 Traficul de magistrală în raport cu talia blocului

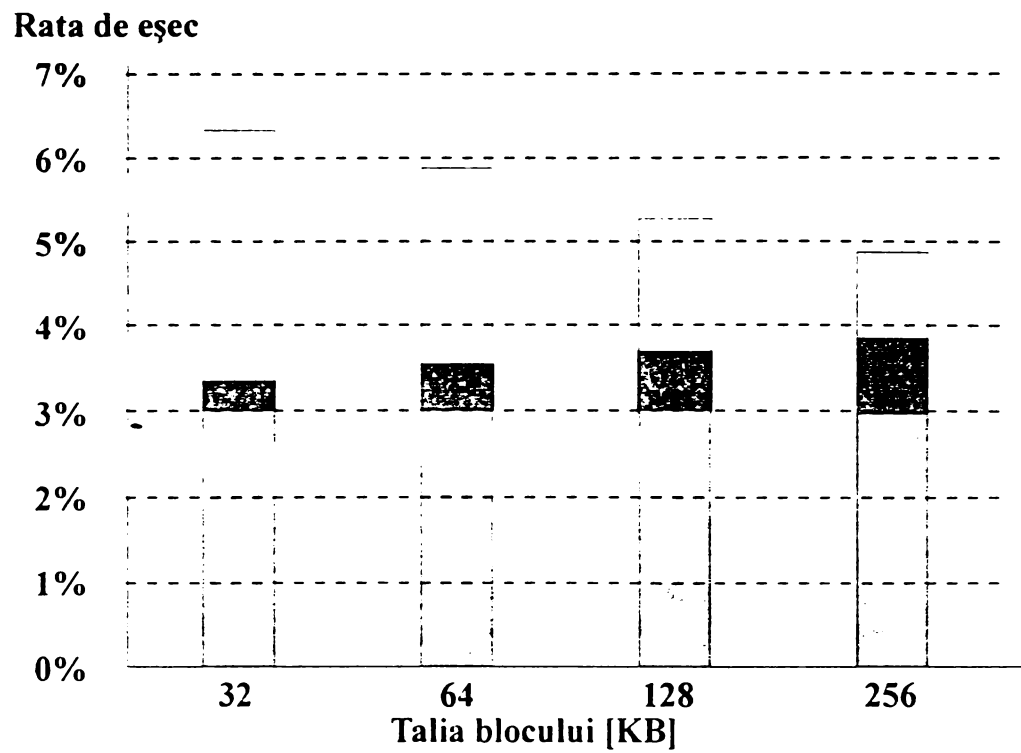


Fig. 3.25 Componentele ratei de eșec a datelor funcție de talia MCH

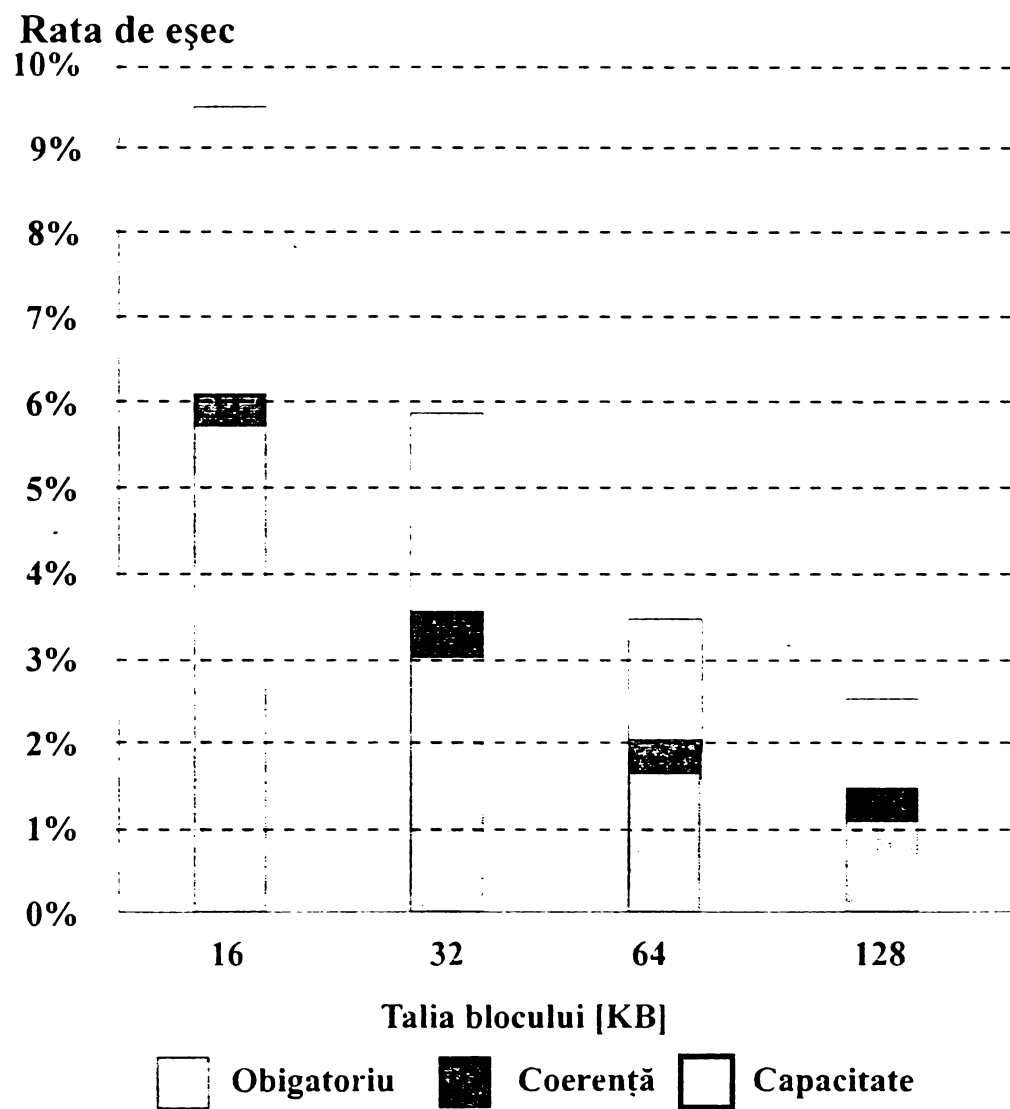


Fig. 3.26 Componentele ratei de eșec funcție de talia cache-ului și a blocului

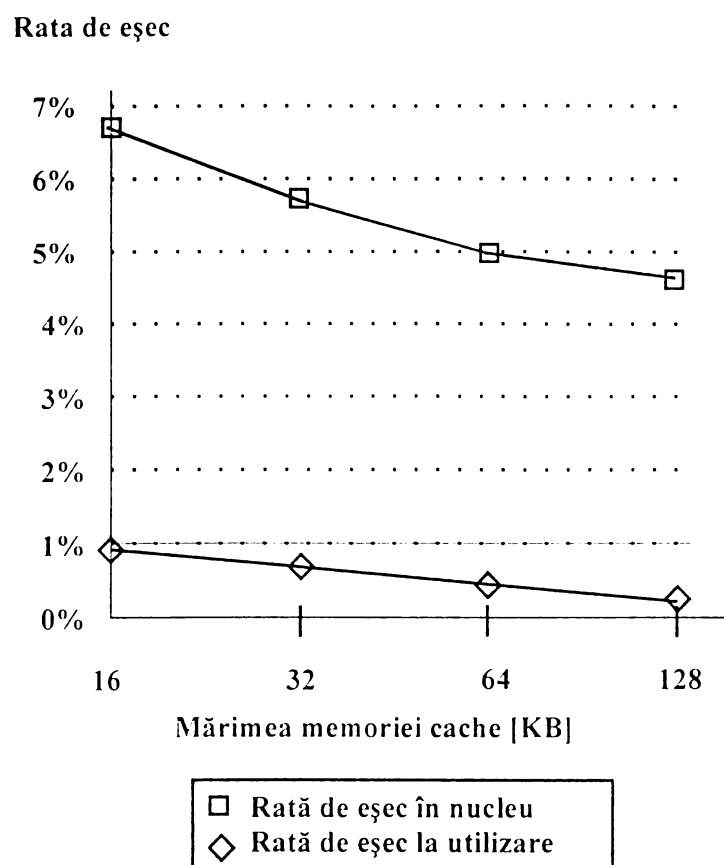


Fig. 3.27 Rata de eșec funcție de talia cache

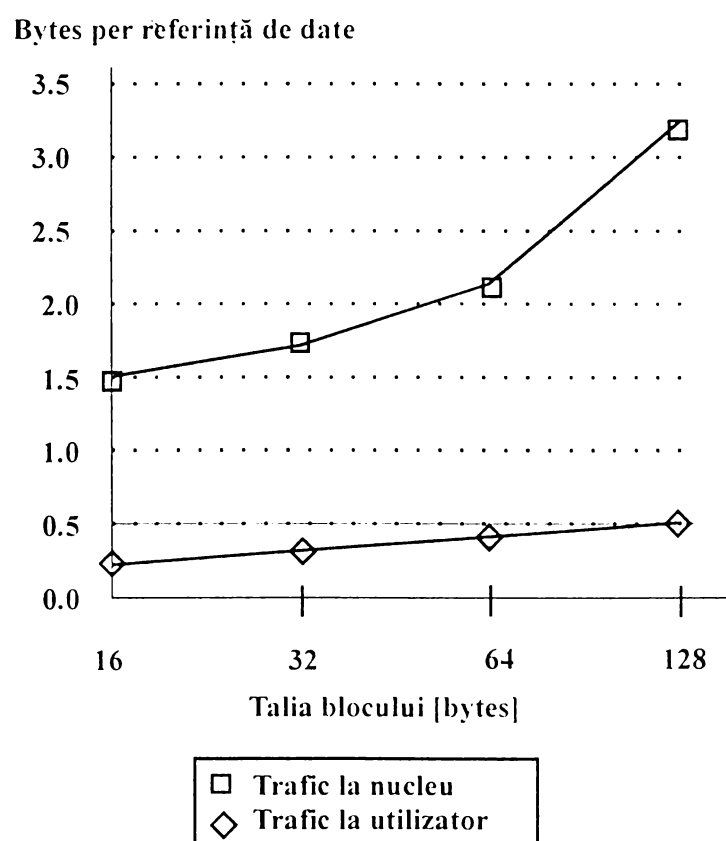


Fig. 3.28 Evoluția traficului în porțiunea de SO și utilizator în cazul programelor paralele

3.2.2 Arhitecturi cu memorie partajată distribuită

O mașină scalabilă ce suportă o memorie partajată ar putea să aibă exclusiv sau inclusiv coerența memoriei cache. Din punct de vedere hardware, cea mai simplă situație este aceea în care coerența memoriei cache este exclusivă. În astfel de sisteme memoria este distribuită printre noduri și toate nodurile sunt interconectate printr-o RIN. Accesul poate fi fie local fie la distanță, strategia fiind decisă de un controler dispus în fiecare nod, pe baza adresei, și care se referă la faptul că data este rezidentă în memoria locală sau în memoria de la distanță.

În ultimul caz, se trimite un mesaj controlerului dispus în memoria de la distanță, pentru accesarea datei.

Aceste sisteme sunt dotate cu MCH, dar pentru a preveni probleme legate de coerență, data partajată este marcată ca fiind indispensabilă pentru plasarea în MCH și dacă data privată este menținută în MCH.

Evident, sistemul de operare poate explicit să cache-eze valoare datei partajate copiind data din porțiunea partajată a spațiului de adresă în porțiunea privată locală a spațiului de adresă cache-at. Coerența este apoi controlată prin tehnici software.

Avantajul unui astfel de mecanism este că necesită un suport hardware redus.

Există câteva dezavantaje în această abordare.

În primul rând, mecanismele de compilare pentru asigurarea coerenței cache prin mijloace software, de-o manieră transparentă, sunt foarte limitate.

Dificultatea majoră o constituie faptul că algoritmi de menținere a coerenței, bazați pe tehnici software, trebuie să fie conservativi : fiecare bloc ce *poate* fi partajat trebuie să fie tratat ca și când *este* partajat. Această abordare conduce la cheltuieli excesive de asigurare a coerenței deoarece compilatorul nu poate prezice partajarea actuală cu acuratețe.

În al doilea rând, fără asigurarea coerenței MCH, sistemele pierd avantajul de-a fi capabile să găsească și să utilizeze cuvinte multiple într-un singur bloc cache.

În al treilea rând, mecanismele pentru tolerarea latenței ca și readucerea din memorie, sunt mult mai utile atunci când pot aduce cuvinte multiple, ca și un bloc cache, și unde data adusă rămâne coerentă.

Datorită motivelor prezentate mai sus, coerența MCH este o cerință acceptată în SMP de talie mică. Pentru arhitecturi de-o talie superioară, există diferite tehnici de extindere a modelului de memorie cache partajată și coerentă.

Prin urmare magistrala poate fi înlocuită cu o RIN mai scalabilă și memoria poate fi distribuită astfel încât lățimea de bandă să fie și ea scalată.

Schemele de asigurare a coerenței prin monitorizare (snooping coherence scheme) se caracterizează prin absența scalabilității. Un protocol de monitorizare necesită comunicarea cu toate cache-urile la fiecare eșec cache, incluzând scrieri a datelor partajate potențial. Absența oricărei structuri de date centralizate referitoare la starea cache-urilor este avantajul fundamental al schemei bazate pe monitorizare, deoarece îi permite să fie ieftină.

Se pot construi arhitecturi scalabile cu memorii partajate ce includ coerența cache. Un protocol alternativ protocolului bazat pe monitorizare îl reprezintă cel *bazat pe director (directory - based)*. Un *director* păstrează starea fiecărui bloc ce poate fi cachat. Informația conținută în director include marcarea situațiilor de tipul : există copii ale blocului în cache, sunt viciate, etc.

Implementările cu directori existente asociază o intrare în director cu fiecare bloc de memorie, în protocoale tipice, cantitatea de informație este proporțională cu produsul dintre numărul de blocuri de memorie și numărul procesoarelor.

Pentru evitarea situației când directorul devine un factor de strangulare a traficului, intrările în director pot fi distribuite în modulele de memorie astfel încât diferite accese la director pot fi distribuite la diferite locații, la fel ca și în situația în care diferite solicitări de acces la memorie sunt direcționate la diferite memorii.

Un director distribuit se caracterizează prin faptul că starea partajată a unui bloc este întotdeauna într-o singură locație cunoscută. Tocmai această proprietate este aceea care permite protocolului de menținere a coerenței să evite *difuzarea (broadcast)*.

În fig.3.29 se prezintă arhitectura unui SMP cu memorie distribuită cu directori atașați fiecărui nod. Este schema din fig. 3.14 detaliată!

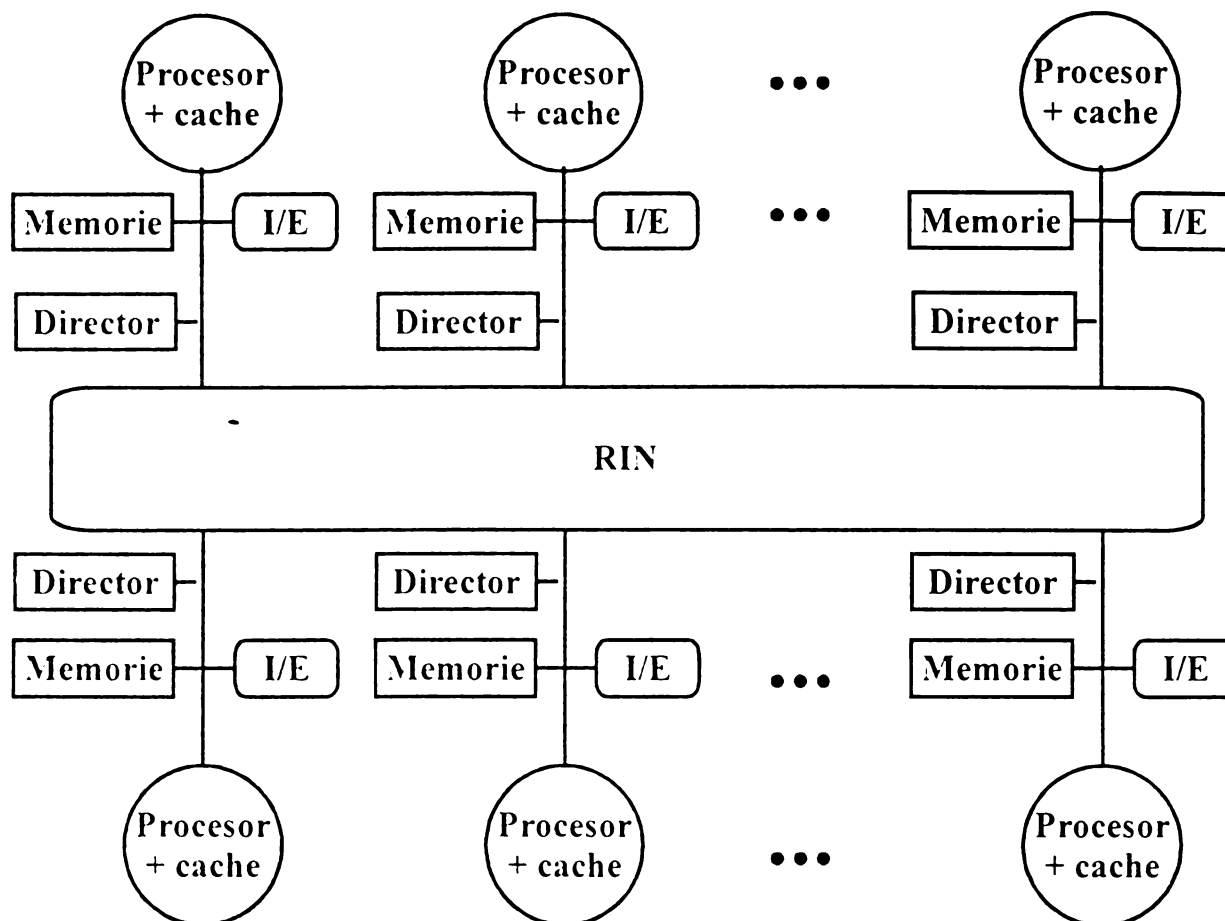


Fig. 3.29 Arhitectura SMP - TM cu directori atașați fiecărui nod pentru a implementa coerența cache

Din fig. 3.29 se observă că fiecare director este responsabil cu marcarea memoriilor cache ce partajează adresele de memorie ale porțiunii de memorie din nod. Directorul poate comunica cu procesorul și cu memoria prin intermediul unei magistrale comune, ca și în figură, sau poate avea un port separat la memorie, sau poate fi o parte a unui controler central de nod prin care trec toate comunicațiile intranod și internod.

A. Protocoale de asigurare a coerenței bazate pe director. Fundamente

Există două operații primare pe care un protocol bazat pe director trebuie să le implementeze:

- manipularea unui eșec la citire;
- manipularea unei scrieri la un bloc cache partajat.

Manipularea unui eșec la scriere într-un bloc partajat este o combinație dintre cele două prezentate anterior. Pentru a implementa aceste operații, un director trebuie să

traseze starea fiecărui bloc cache. Într-un protocol simplu, aceste stări ar putea fi următoarele:

- *Partajată (shared)* - Unul sau mai multe procesoare au blocul partajat și valoarea în memorie trebuie să fie actualizată. La fel în toate memoriile cache.
- *Necache-ată (unchached)* - Nici un procesor nu are o copie a blocului cache.
- *Exclusivă (exclusive)* - Doar un procesor - identificat exact - are o copie a blocului cache și blocul a fost scris, astfel că copie memoriei este perimată. Procesorul este denumit *proprietarul* blocului.

Suplimentar cu trasarea stării fiecărui bloc cache, trebuie marcate procesoarele ce au copii ale blocului când acesta este partajat, deoarece acestea trebuie invalidate la o scriere. Cea mai simplă cale pentru efectuarea acestei operații este de-a menține un vector binar (de unul sau mai mulți biți) pentru fiecare bloc de memorie. Când blocul este partajat, fiecare bit al vectorului indică dacă procesorul corespunzător are o copie a acestui bloc. Se poate de asemenea utiliza vectorul pentru menținerea traseului proprietarului blocului când blocul este într-o stare exclusivă. Din rațiuni de eficiență, se trasează starea fiecărui bloc cache la memoriile cache individuale.

Stăriile și tranzițiile pentru un SMP sunt identice la fiecare cache cu cele utilizate pentru monitorizarea cache-urilor. Acțiunile la o tranziție sunt ușor diferite.

Se fac următoarele ipoteze simplificatoare :

- încercările de scriere a unor date non exclusive întotdeauna generează eșecuri la scriere și procesoarele se blochează până când un acces se finalizează;
- interconectarea nu se poate utiliza ca și un singur punct de arbitrare. o funcție realizată de magistrală în cazul monitorizării.
- toate mesajele trebuie să aibă răspunsuri explicite, deoarece RIN este orientată pe mesaje (magistrala este orientată pe tranzacții).

În TABEL V se prezintă tipurile de mesaje trimise între noduri. Nodul *local* este nodul ce generează o solicitare. Nodul *rezident (home)* este nodul unde sunt rezidente locația de memorie și intrarea în director a unei adrese. Spațiul de adresă fizică este static distribuit astfel încât nodul ce conține atât memoria cât și directorul pentru o adresă fizică dată este identificat.

Nodul *îndepărtat (remote)* este nodul ce are o copie a unui bloc cache, fie exclusiv sau partajat. Nodul *local* poate fi nod *rezident* și vice versa. În ambele cazuri protocolul este același, prin urmare mesajele între noduri pot fi înlocuite prin tranzații în interiorul nodului ce sunt mai rapide.

În acest paragraf se presupune un model simplu de consistență a memoriei. Pentru a minimiza tipul mesajelor și complexitatea protocolului se presupune că mesajele vor fi recepționate și procesate în aceeași ordine în care ele au venit. Această presupunere, greu de îndeplinit în practică, ne asigură că invalidările transmise de un procesor sunt imediat onorate.

TABEL V. Mesajele posibile transmise între noduri pentru a menține coerența.

Tip	Sursă	Destinație	Conținut	Funcție
EC	CH-L	P A	P, A	P are un eșec la citire în A; Solicită data și îl face pe P partener la partajarea la citire.
ES	CH-L	P A	P, A	P are un eșec la scriere în A; Solicită data și îl face pe P proprietar exclusiv.
I	DR	CH-D	A	Invalidează o copie partajată a datei la adresa A.
A	DR	CH-D	A	Aduce blocul la adresa A și-l trimite la directorul lui de acasă; invalidează blocul în cache.
A/I	DR	CH-D	A	Aduce blocul la adresa A și-l trimite la directorul lui de acasă; invalidează blocul în cache.
RVD	DR	CH-L	Data	Returnează o valoare a datei din memoria de acasă.
RD	CH-D	DR	A, Data	Rescrie o valoare de dată pentru adresa A.

Abrevierile sunt următoarele :

- EC - eșec la citire (read miss);
- ES - eșec la scriere (write miss);
- I - invalidare;
- A - aducere;
- A/I - aducere/invalidare;
- RVD - returnează validarea datei;
- RD - rescrie data;
- CH-L - cache local;
- CH-D - cache la distanță;
- DR - director "rezident";
- P - procesor.

Un exemplu de protocol bazat pe director.

În literatură [HP 96] se prezintă mai multe exemple de protocoale bazate pe director. Stările de bază ale unui bloc cache într-un protocol bazat pe director sunt aceleași ca și într-un protocol de monitorizare și stările într-un director sunt similare cu acelea prezentate anterior.

În continuare se analizează diagrame simple de stare ce prezintă tranziții de stare pentru un bloc cache individual și apoi se examinează diagrama de stare pentru intrarea din director corespunzătoare fiecărui bloc din memorie. Aceste diagrame de tranziție de stare nu reprezintă toate detaliile unui protocol de asigurare a coerenței.

Fig. 3.30 prezintă acțiunile protocolului la care răspunde un cache individual.

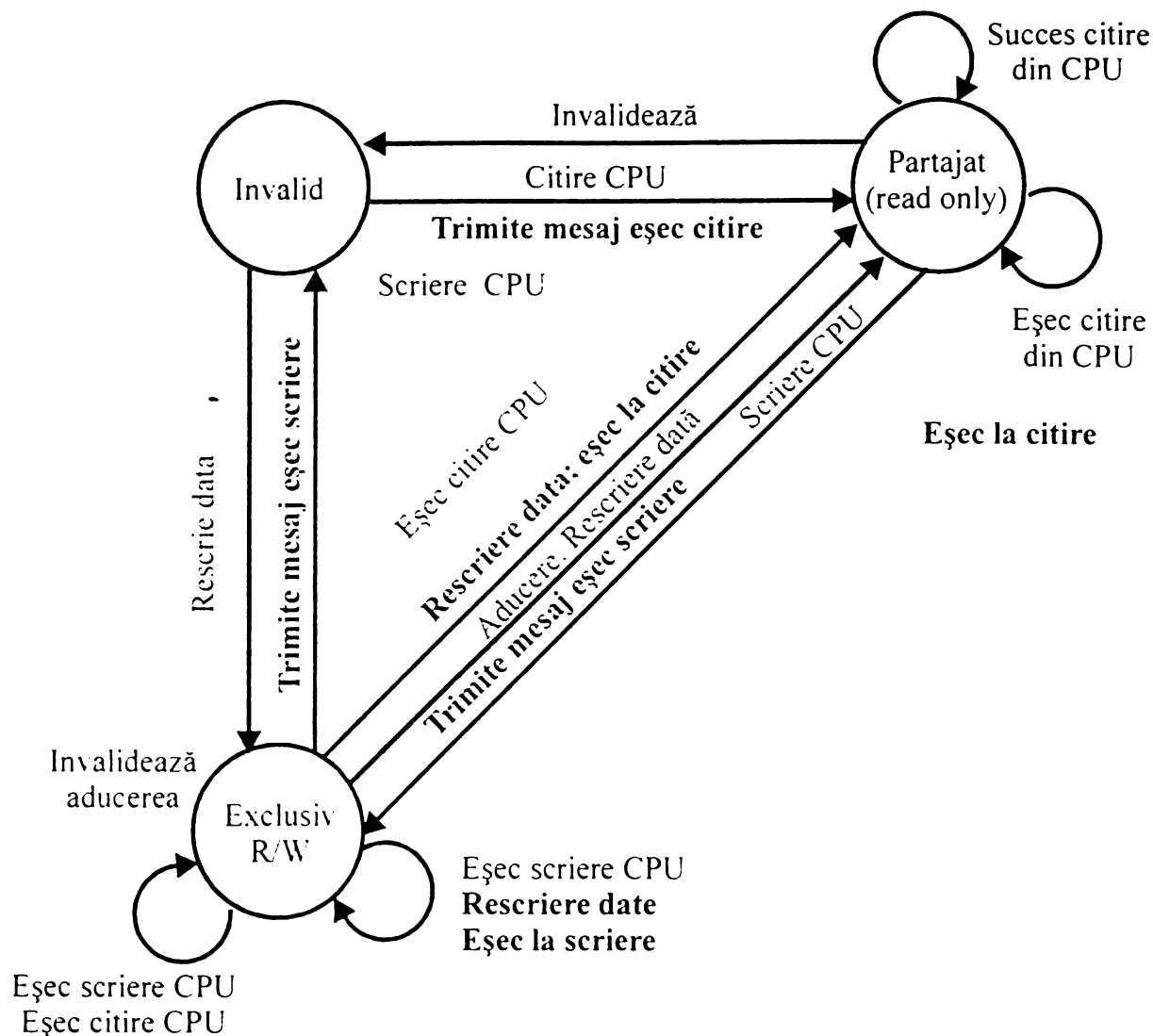


Fig. 3.30 Diagrama tranzițiilor de stare pentru un bloc cache individual într-un sistem bazat pe director

Câteva remarci sunt necesare. Stările sunt identice cu acelea din protocolul de monitorizare și tranzațiile sunt foarte asemănătoare cu solicitări de invalidare și de rescriere ce înlocuiesc eșecurile la scriere ce au fost difuzate pe magistrală. Se presupune că o încercare de a scrie un bloc cache partajat este tratată ca și un eșec; în practică o astfel de tranzație poate fi tratată ca și o solicitare a proprietarului și poate fi livrată proprietarului fără solicitarea ca blocul cache să fie adus. Prin urmare, operația de eșec la scriere, ce a fost difuzată pe magistrală în schema cu monitorizare, este realizată cu operații de căutare a datei și de invalidare ce sunt trimise selectiv de către controlerul directorului. Orice bloc cache trebuie să fie în starea exclusivă când este scris și orice bloc partajat trebuie actualizat în memorie.

În protocolul bazat pe director, directorul implementează cealaltă jumătate a protocolului de coerență. Un mesaj transmis spre un director provoacă două tipuri diferite de acțiuni :

- actualizează starea directorului;
- trimite mesaje adiționale pentru satisfacerea solicitării.

Stările într-un director reprezintă cele trei stări standard pentru un bloc, dar pentru toate copiile cache ale unui bloc de memorie nu numai pentru un singur bloc cache.

Blocul de memorie poate fi necache-at de către orice nod, cache-at în noduri multiple și partajat sau cache-at în mod exclusiv cu posibilitatea de a fi scris într-un nod anume.

Figura 3.31 prezintă acțiunile luate la nivel de director în răspuns la mesajele recepționate. Directorul recepționează trei solicitări diferite : *eșec la citire*, *eșec la scriere* și *rescriere a datei*.

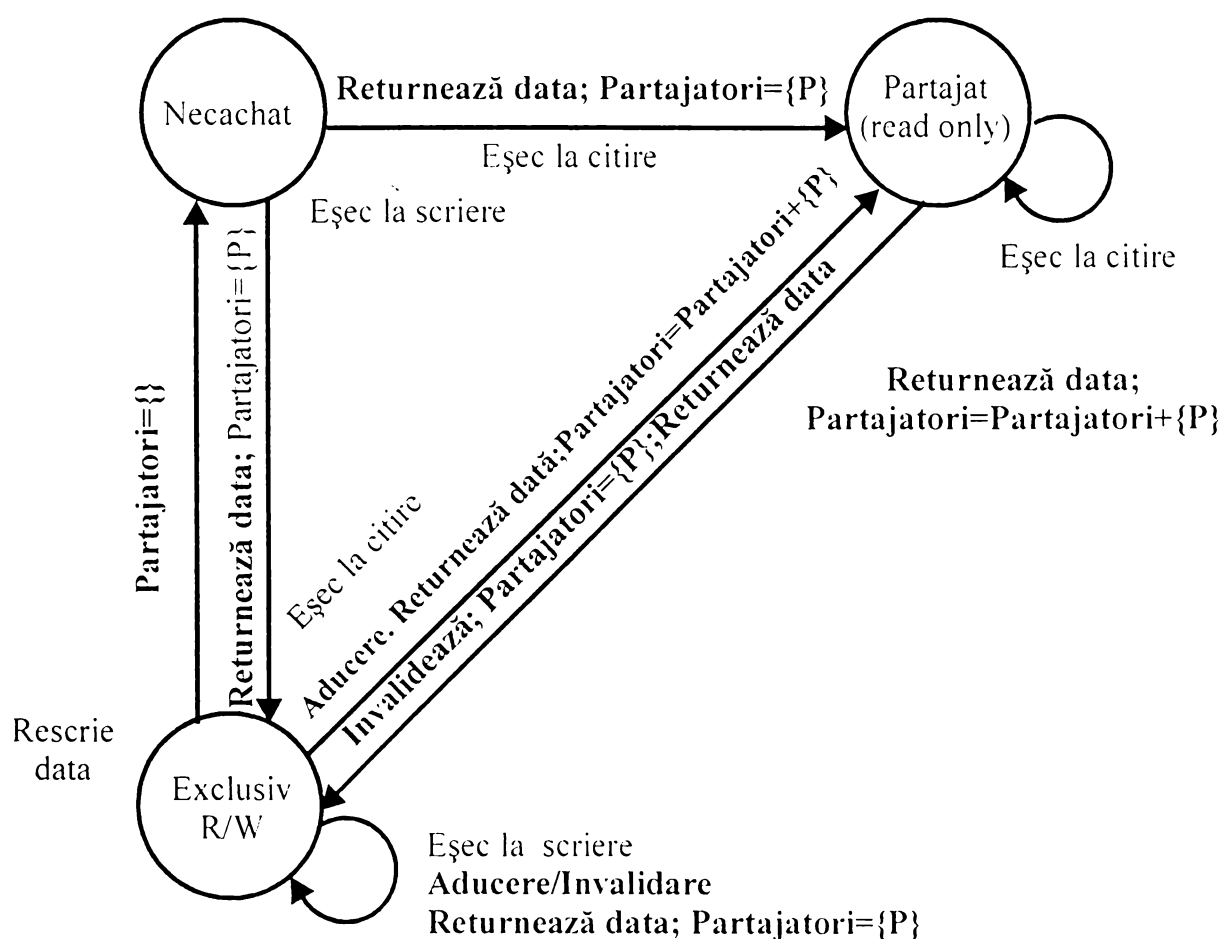


Fig. 3.31 Diagrama tranzițiilor de stare pentru director are aceleași stări și structură ca și diagrama de tranziție pentru o memorie cache individuală

Se remarcă existența așa numiților *partajatori* (*sharers*).

Aceștia sunt vectori binari ce efectuează trasarea seturilor de procesare ce au o copie a unui bloc.

Protocolul simplificat prezentat presupune că unele acțiuni sunt atomice, ca de exemplu solicitarea unei valori și transmiterea ei altui nod.

Analizând solicitările primite și acțiunile luate stare cu stare se constată că atunci când un bloc este într-o stare necachată următoarele solicitări sunt posibile :

- *Eșec la citire* - Procesorul solicitant a transmis data solicitată din memorie și solicitatorul este declarat singurul nod ce partajează. Starea blocului este declarată *partajată*.
- *Eșec la scriere* - Procesorul solicitat a transmis valoarea și devine *nod partajator*. Blocul este făcut *exclusiv* pentru a indica că doar copia validă este cachată. Partajatorii indică identitatea deținătorului. Când un bloc este în starea "*partajată*" valoarea din memorie este de actualizat, așa pot apare două cereri :
 - *Eșec la citire* - Procesorul solicitant a transmis data solicitată din memorie și procesorul solicitant este adăugat setului partajator.
 - *Eșec la scriere* - Procesorul solicitant a transmis valoarea. Toate procesoarele din setul Partajatorilor au transmis mesaje de invalidare, și setul Partajatorilor conține identitatea procesorului solicitant. Starea blocului este făcută "*exclusiv*".

Când blocul este în starea "*exclusiv*", valoarea curentă a blocului este menținută în memoria cache a procesorului identificat de către deținător astfel încât există trei posibilități de solicitări de director :

- *Eșec la citire* - Procesorul deținător a transmis un mesaj de aducere a datei ce provoacă tranziția stării blocului din cache-ul deținătorului în "*partajat*" și determină deținătorul să transmită data către director, unde este scrisă în memorie și retransmisă procesorului solicitant. Identitatea procesorului solicitant este adăugată la partajatori care încă mai conțin identitatea procesorului ce a fost deținător, el mai posedând o copie ce poate fi citită.

- *Rescrierea datei* - Procesorul solicitant înlocuiește blocul și prin urmare trebuie să-l rescrie. Copia memoriei se poate actualiza, blocul devine necach-at, și setul de partajatori se videază.
- *Eșec la scriere* - blocul are un nou deținător. A fost transmis un mesaj către vechiul deținător, ce provoacă cache-ul să trimită valoarea blocului către director, de la care este transmis procesorul solicitant, ce devine noul deținător. Partajatorii sunt poziționați să identifice noul deținător și starea blocului rămâne "*exclusiv*".

Diagrama de tranziție a stărilor din fig. 3.34 este simplificată.

În sistemele reale cele mai multe protocoale transmit data de la nodul deținător către nodul solicitant în mod direct. Astfel de optimizări nu duc la creșterea complexității dar repercursiunile se manifestă în alte părți ale sistemului.

B. Performanțe ale protocoalelor de asigurare a coerenței bazate pe director

Performanța unui sistem bazat pe director depinde de mulți dintre factorii ce influențează performanța mașinilor orientate pe magistrală (de ex. talia cache-ului, numărul de procesoare, talia blocului) cât și de distribuția eșecurilor la diferite locații în ierarhia de memorie. Locația unei date solicitate depinde atât de alocarea inițială cât și configurațiile partajate.

Analiza performanțelor se va efectua pe o memorie cache de 128 KB și asupra unui bloc de 64 bytes.

În arhitecturile cu memorie distribuită, distribuția solicitărilor de acces la memorie între memoria locală și cea de la distanță, reprezintă cheia performanței. Prin urmare eșecurile la memoria cache se vor separa în solicitări locale și la distanță.

În literatură [HP 96] se analizează performanțele SMP-TM, bazate pe director, observând evoluția mai multor programe paralele specifice SMP. De exemplu programele FFT, LU, Barnes și Ocean.

În acest paragraf vom analiza doar comportamentul schemelor bazate pe director în cazul evoluției unui program FFT.

În fig. 3.32 și 3.33 se prezintă dependența ratei de eșec de numărul de procesoare, respectiv de talia memoriei cache.

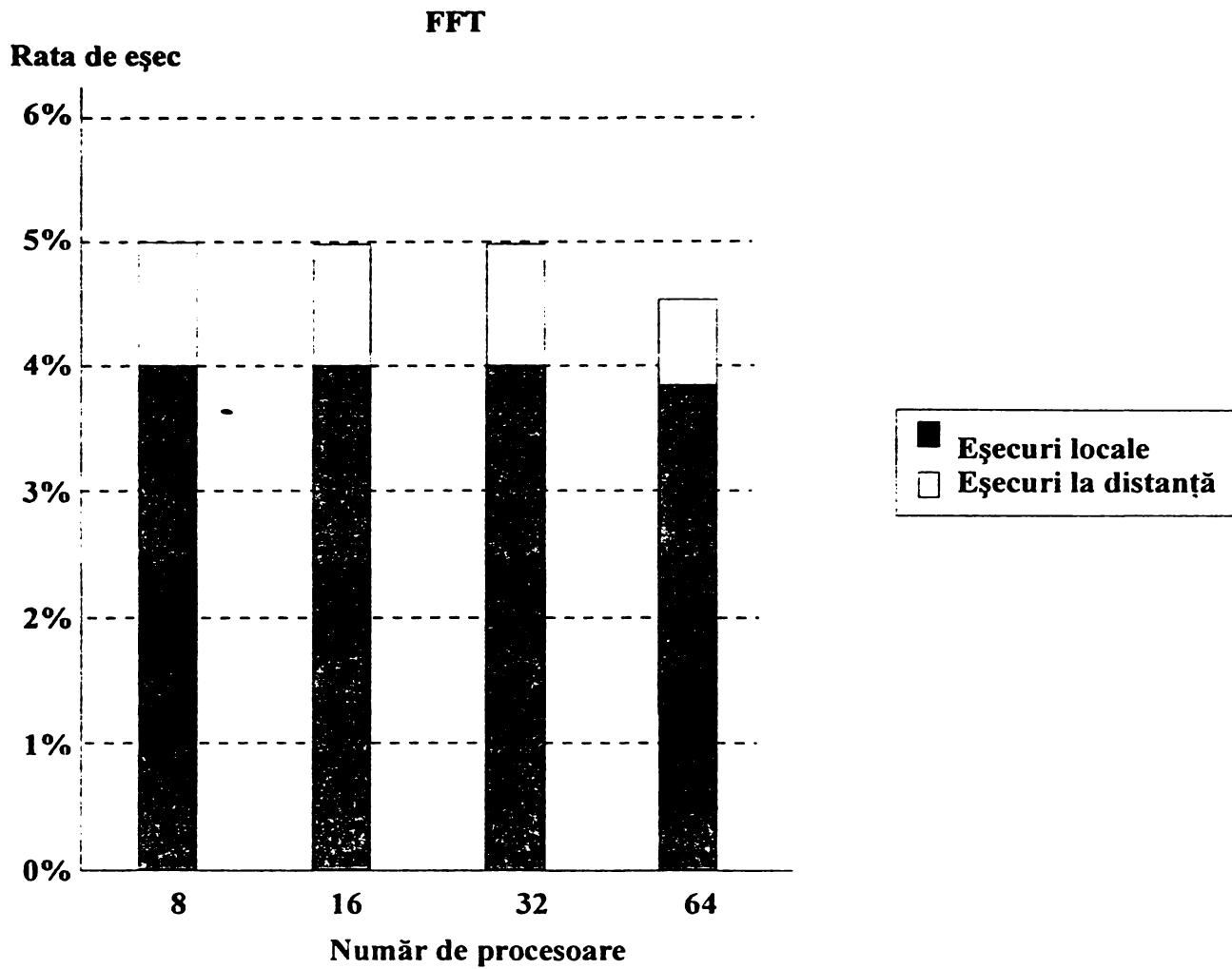


Fig. 3.32 Rata de eșec în funcție de numărul de procesoare

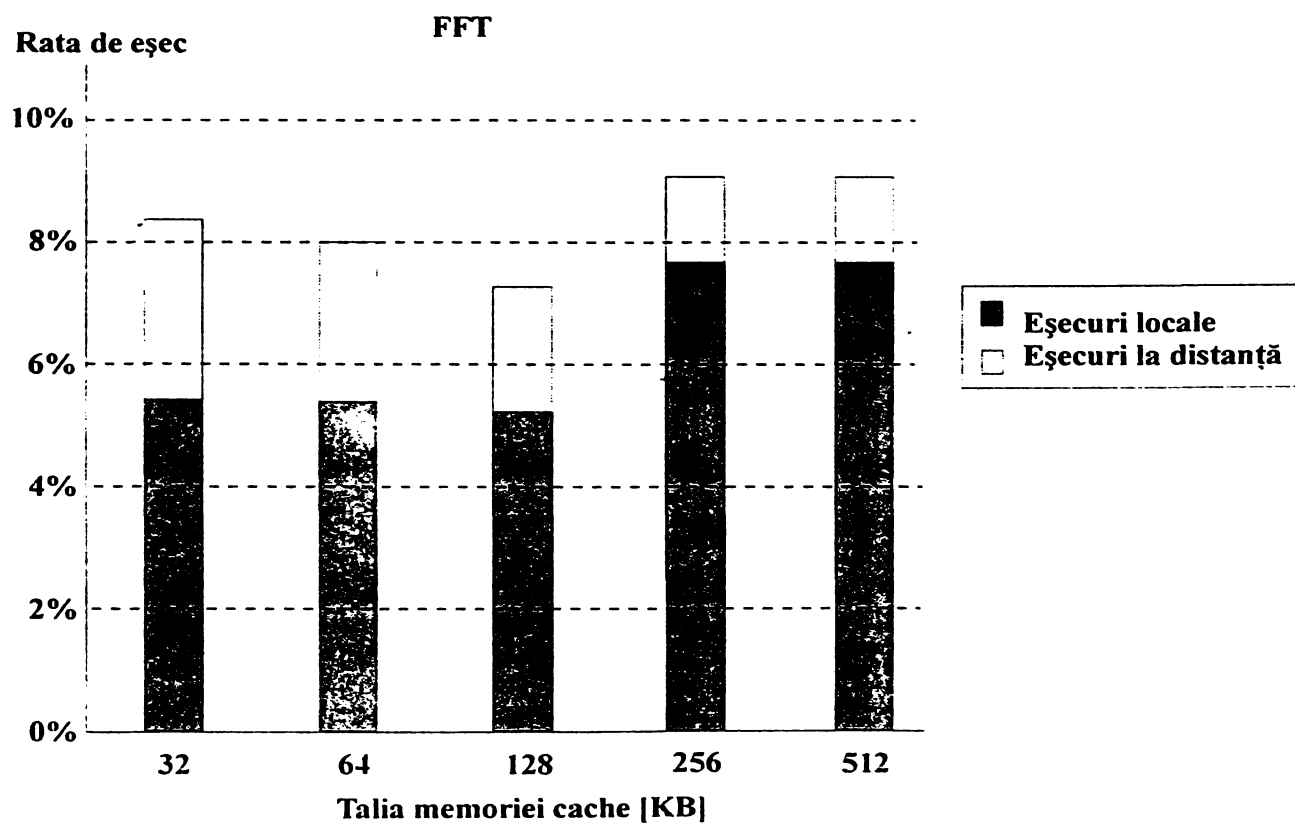


Fig. 3.33 Ratele de eșec scad pe măsură ce crește talia cache

În fig.3.34 se prezintă efectul schimbării taliei blocului din memoria cache. Deoarece aplicația FFT are o bună localitate spațială, creșterea taliei blocului reduce rata de eșec chiar și pentru blocuri mari și prin urmare nu este rentabilă utilizarea acestora.

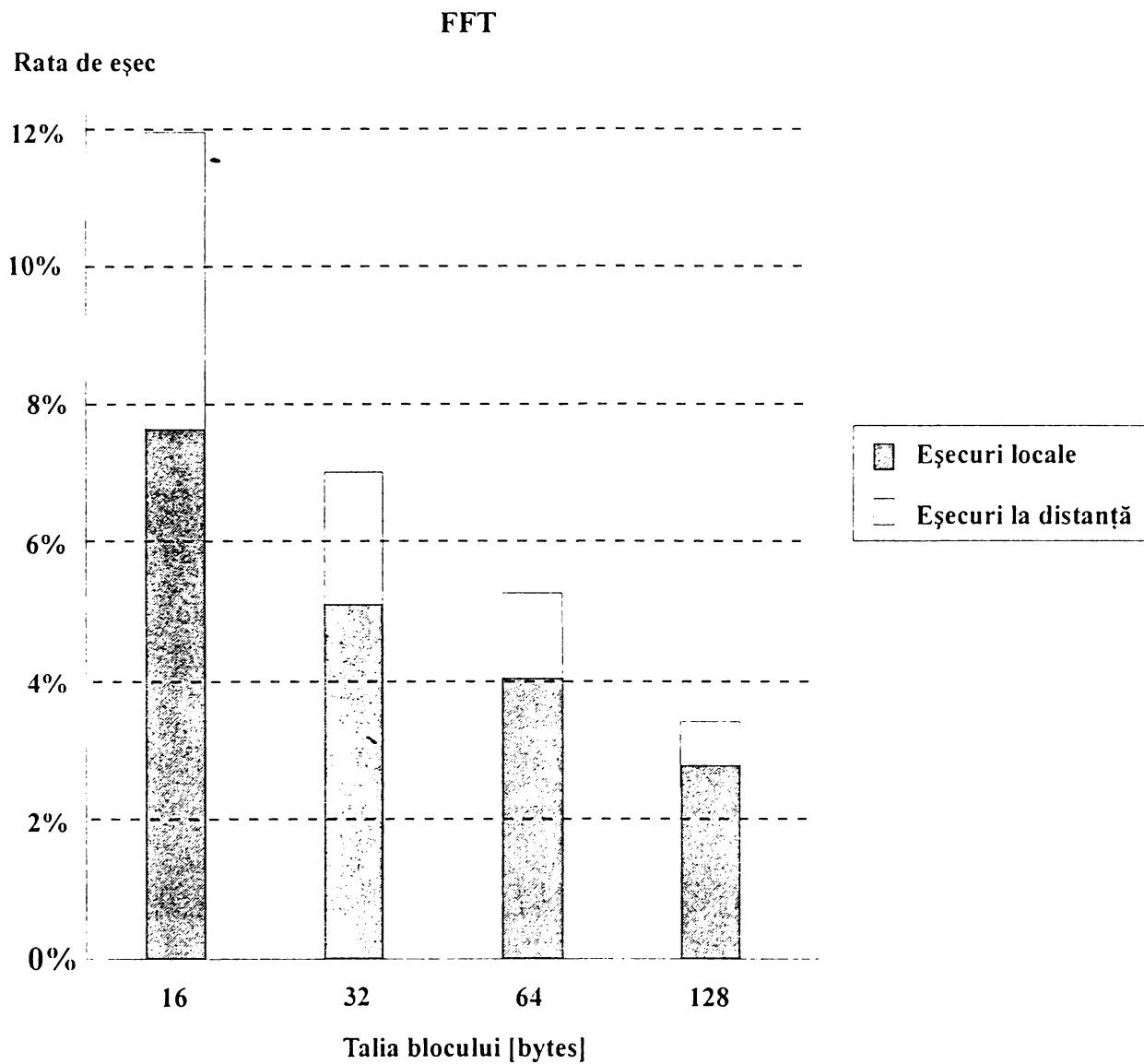


Fig. 3.34 Rata de eșec în funcție de talia blocului într-un sistem cu 64 de procesoare și 128 KB memorie cache

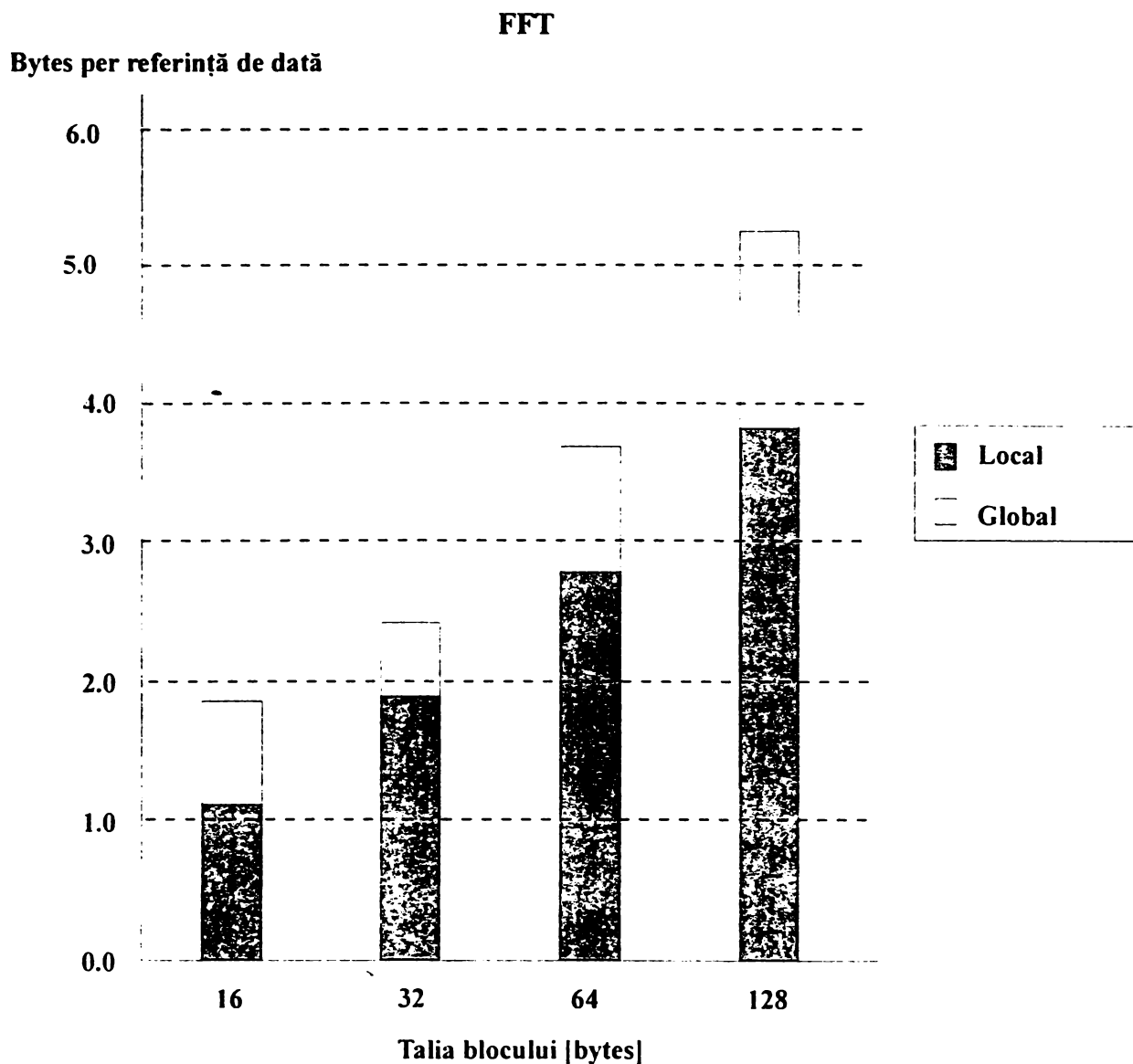


Fig. 3.35 Numărul de bytes per referință de dată crește cu talia blocului

3.2.3 Modele de consistență a memoriei

Coerența memoriei cache asigură procesoarelor multiple o viziune consistentă asupra memoriei.

Se pune problema gradului de consistență ce trebuie asigurat. De exemplu : când un procesor trebuie să "vadă" ce a fost actualizat de alt procesor, în ce ordine trebuie un procesor să observe scrierile de date ale altui procesor ?

Un exemplu simplu poate fi edificator asupra complexității procedurii de stabilire a gradului de consistență.

Fie două segmente de cod aparținând proceselor $P1$ și $P2$.

```

P1 :   A = 0 ;           P2 :   B = 0 ;
      A = 1 ;           B = 1 ;
L1 :   if (B == 0) ...   L2 :   if (A == 0) ...

```

Se presupune că procesele rulează pe procesoare diferite și că locațiile A și B sunt cachate inițial de ambele procesoare cu valoarea inițială 0.

Dacă scrierile au efect imediat și sunt văzute imediat de alte procesoare, va fi imposibil pentru *ambele* instrucții **if** (etichetate L1 și L2) să evalueze condițiile ca adevărate, deoarece ajungând la instrucția **if** înseamnă că fie A fie B trebuie să fi avut asigurată valoarea 1.

Dar presupunând că invalidarea la scriere este întârziată și că procesorului îi este permis să continue, atunci este posibil ca atât P1 cât și P2 să nu fi văzut invalidările pentru B și A înainte de-a încerca să citească valorile. Se ridică problema dacă un astfel de comportament este permis și dacă da, în ce condiții.

Cel mai direct model pentru consistența memoriei este denumit *consistență secvențială*.

Consistența secvențială pretinde ca rezultatul oricărei execuții să fie același ca și în cazul când accesele executate de fiecare procesor au fost ținute în ordine și accesele printre procesoare diferite au fost întrețesute.

Aceasta elimină posibilitatea unor execuții ambigue ca și în exemplul anterior, deoarece asignările trebuie să fie completate înaintea efectuării instrucțiilor **if**.

În fig. 3.36 se prezintă procedura de interzicere de către consistența secvențială a unei execuții când ambele instrucții **if** întâlnesc condiția "true".

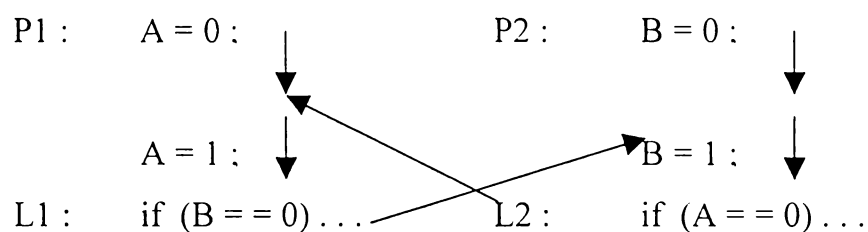


Fig. 3. 36 Un exemplu de consistență secvențială

Cea mai simplă metodă de implementare a consistenței secvențiale constă în solicitarea adresată unui procesor să-și întârzie efectuarea oricărui acces la memorie până când toate invalidările cauzate de acest acces sunt efectuate.

De reținut că consistența memoriei implică operații asupra diferitelor variabile. Cele două accese ce trebuie să fie ordonate se referă la locații diferite de memorie.

În exemplul anterior trebuie întârziată citirea lui A sau B ($A = 0$ sau $B = 0$) până când scrierea anterioară a fost efectuată ($B = 1$ sau $A = 1$). În consistența secvențială nu se poate plasa, de-o manieră simplistă, operația de scriere într-un tampon de scriere și să se continue cu citirea.

Consistența secvențială oferă o paradigmă simplă în programare și reduce performanța potențială îndeosebi la SMP cu un număr mare de procesoare sau cu întârzieri datorate interconectărilor la distanță.

Din punct de vedere al programatorului, acest model de consistență este cel mai agreat datorită simplității sale.

Programele trebuie să fie însă *sincronizate*, adică toate accesele la data partajată sunt ordonate de către operații de sincronizare. O referință la o dată este ordonată de către o operație de sincronizare dacă o scriere a unei variabile de către un procesor și un acces (scriere sau citire) a acelei variabile de către alt procesor sunt separate de către o pereche de operații de sincronizare, una executată *după* scriere de către procesorul ce scrie și alta executată *înaintea* accesului efectuat de cel de al doilea procesor.

Cursele de date apar când variabilele pot fi actualizate fără ordonarea prin sincronizare.

În literatură [HP 96] se prezintă mai multe tehnici de sincronizare.

În afară de operații de sincronizare, este necesară definirea ordonării operațiilor de memorie. Există două tipuri de restricții în ordonare : *bariere la scriere* (*write fences*) și *bariere la citire* (*read fences*).

Barierile sunt puncte fixe în cadrul unui calcul ce asigură că nici o citire sau scriere nu pot trece peste ele.

De exemplu, o barieră la scriere executată de către un procesor P asigură că :

- toate scrierile efectuate de către P, ce apar înainte ca P să fi executat operația de barieră la scriere, sunt efectuate și
- nu se inițiază nici o scriere ce apare după bariera din P, înainte de barieră.

În cadrul consistenței secvențiale, toate citirile sunt citiri cu barieră și toate scrierile sunt scrieri cu barieră. Această proprietate limitează abilitatea structurii hardware de-a optimiza accesele, deoarece ordinea trebuie strict menținută.

Din punct de vedere al performanței, procesorul ar trebui să execute citirile cât se poate de devreme și să efectueze scrierile cât de târziu posibil. Barierele acționează ca și frontiere, forțând procesorul să ordoneze citirile și scrierile în raport cu bariera.

O *barieră de memorie* (*memory fence*) reprezintă o operație ce acționează atât ca și o barieră de citire cât și ca una de scriere.

Barierele de memorie accentuează ordonarea printre accesele diferitelor procese. În cadrul unui singur proces este necesar ca ordinea din program să fie întotdeauna păstrată, astfel că citirile și scrierile la aceeași locație nu pot fi interschimbate.

A. Modele relaxate pentru consistența memoriei

Aceste modele permit implementări cu performanță mai ridicată și în același timp mai păstrează un model simplu de programare pentru programe sincronizate.

Modelele relaxate variază în ceea ce privește gradul de constrângere impus la implementare [CBZ95], [GLL+90], [GW88], [HP96].

Modelele pot fi definite în funcție de tipul de ordonări efectuate de un singur procesor, asupra citirilor și scrierilor.

Există patru astfel de ordonări :

1. $R \rightarrow R$ o citire urmată de o citire;
2. $R \rightarrow W$ o citire urmată de o scriere, care este întotdeauna păstrată dacă operațiile sunt la aceeași adresă, aceasta fiind o antidependență;
3. $W \rightarrow W$ o scriere urmată de o scriere, care este întotdeauna prezervată dacă ele sunt la aceeași adresă, aceasta fiind o dependență de ieșire;
4. $W \rightarrow R$ o scriere urmată de o citire, care este întotdeauna păstrată dacă ele sunt la aceeași adresă, deoarece aceasta reprezintă o adevărată dependență.

Dacă există o dependență între citire și scriere, atunci semanticile programului uniprocessor pretind ca operațiile să fie ordonate. Dacă nu există dependență, modelul de consistență, a memoriei determină ce ordonări trebuie păstrate. Un model de consistență secvențială pretinde ca toate cele patru ordonări să fie păstrate și aceasta

este echivalent cu asumarea unui singur modul centralizat de memorie ce serializează toate operațiile procesorului sau cu presupunerea că toate citirile și scrierile sunt bariere de memorie.

O ordonare este *relaxată*, atunci când se permite ca o operație ce trebuie executată mai târziu de către un procesor să fie efectuată prima.

În această situație, de exemplu, relaxarea ordonării $W \rightarrow R$ înseamnă că citirea care este în urma scrierii să se termine înainte ca scrierea să fie efectuată.

În realitate un model de consistență nu restricționează ordonarea evenimentelor. Ne interesează gradul de observabilitate a unor posibile ordonări. Cu titlu de exemplu, cităm cazul consistenței secvențiale în care sistemul trebuie să păstreze cele patru ordonări, descrise anterior. În practică este permisă reordonarea fără ca ea să fie observabilă.

Modelul de consistență trebuie de asemenea să definească ordonările impuse între accesese la variabilele de sincronizare, ce acționează ca și bariere, și toate celelalte accesese.

Când un sistem implementează consistența secvențială, toate citirile și scrierile - inclusiv accesesele de sincronizare - sunt bariere și prin urmare ținute în ordine.

În cazul modelelor mai slabe este necesară specificarea restricțiilor de ordonare impuse de accesesele de sincronizare la fel ca și restricțiile de ordonare implicând variabile ordinare. Cea mai simplă restricție de ordonare este aceea că fiecare acces de sincronizare este o barieră de memorie.

Fie S un acces la o variabilă de sincronizare.

Utilizând notația de ordonare putem scrie :

$S \rightarrow W, S \rightarrow R, W \rightarrow S$, și $R \rightarrow S$.

De remarcat că un acces de sincronizare este fie R fie W și ordonarea sa este afectată de alte accesese de sincronizare, ceea ce implică existența unei ordonări $S \rightarrow S$.

În literatură [HP 96] se citează mai multe exemple de modele relaxate, caracterizate prin eliminarea unui ordin.

De exemplu dacă o operație de sincronizare este executată înaintea unei citiri, atunci ordonările $W \rightarrow S$ și $S \rightarrow R$ asigură că scrierea se termină înaintea citirii.

Pentru acest model se utilizează nume ca și *consistență a procesorului* (*processor consistency*) și *ordonare de înmagazinare totală* (*total store ordering - TSO*).

Acest model este echivalent cu a face scrierile prin bariere de scriere. Prezentarea acestor modele este rezumată în TABEL VI iar un exemplu este prezentat în fig. 3.37.

TABEL VI. Ordonările impuse de diferite modele de
consistență pentru accese ordinare și accese de sincronizare

Model	Utilizat în	Ordonări ordinare	Ordonări de sincronizare
Consistență secvențială	Cele mai multe sisteme	$R \rightarrow R, R \rightarrow W,$ $W \rightarrow R, W \rightarrow W$	$S \rightarrow W, S \rightarrow R,$ $R \rightarrow S, W \rightarrow S,$ $S \rightarrow S$
TSO sau consistență de procesor	IBM 370, DEC VAX, SPARC	$R \rightarrow R, R \rightarrow W,$ $W \rightarrow W$	$S \rightarrow W, S \rightarrow R,$ $R \rightarrow S, W \rightarrow S,$ $S \rightarrow S$
PSO	SPARC	$R \rightarrow R, R \rightarrow W$	$S \rightarrow W, S \rightarrow R,$ $R \rightarrow S, W \rightarrow S,$ $S \rightarrow S$
Ordonare slabă	Power PC		$S \rightarrow W, S \rightarrow R,$ $R \rightarrow S, W \rightarrow S,$ $S \rightarrow S$
Consistență eliberată	Alpha, MIPS		$S_A \rightarrow W, S_A \rightarrow R,$ $R \rightarrow S_R, W \rightarrow S_R,$ $S_A \rightarrow S_A, S \rightarrow S_R$ $S_R \rightarrow S_A, S_R \rightarrow S_R$

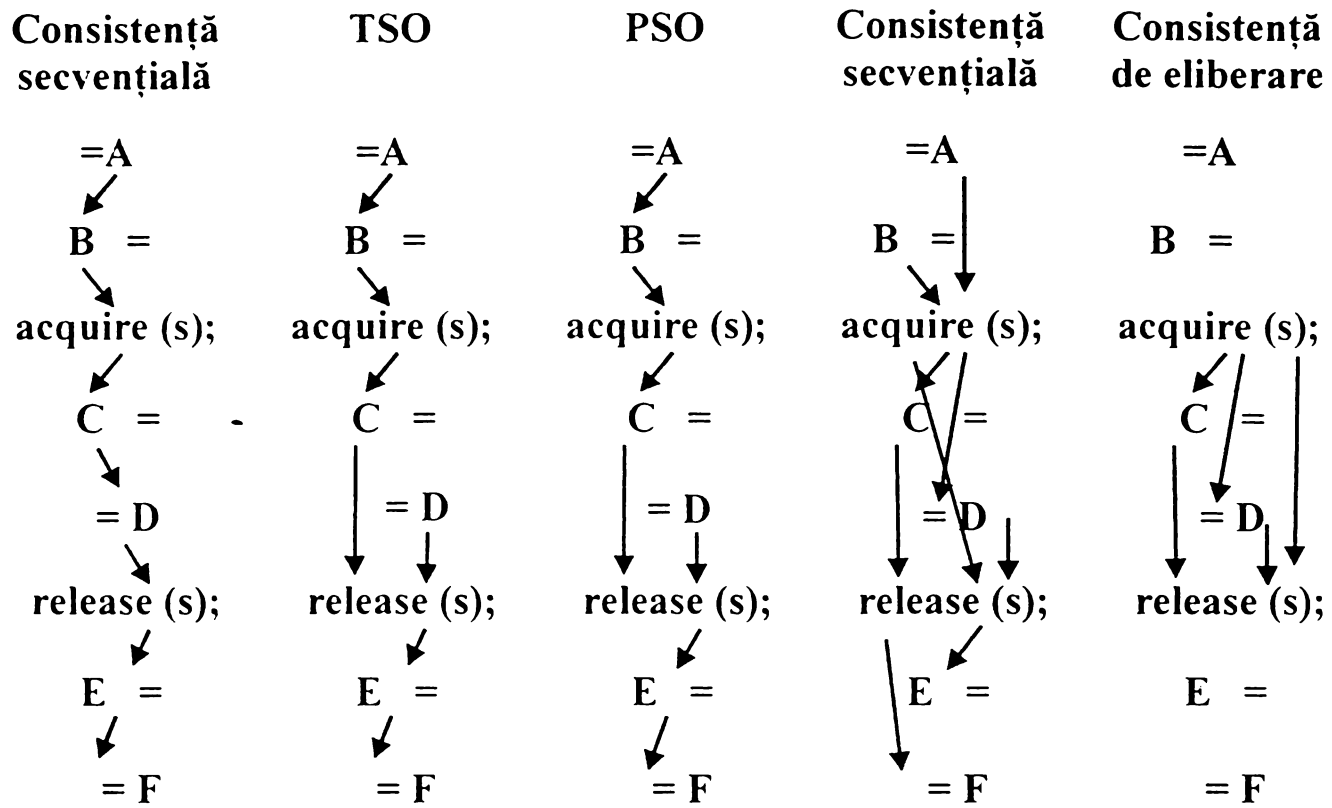


Fig.3.37 Reducerea numărului de ordine impuse pe măsură ce modelele devin mai relaxate. Ordinele minime sunt prezentate cu săgeți

B. Performanța modelelor relaxate

Potențialul de performanță a modelelor relaxate depinde atât de capacitățile sistemului cât și de aplicație. La examinarea performanței modelului de consistență a memoriei se recurge în prealabil la definirea mediului hardware. Configurațiile hardware luate în considerare trebuie să aibă câteva caracteristici comune. Acestea sunt:

- Erorile cache consumă 50 de cicluri de ceas;
- CPU include un tampon de scriere cu o adâncime de 16;
- Memoriile cache au 64 KB și au linii de 16 byte.

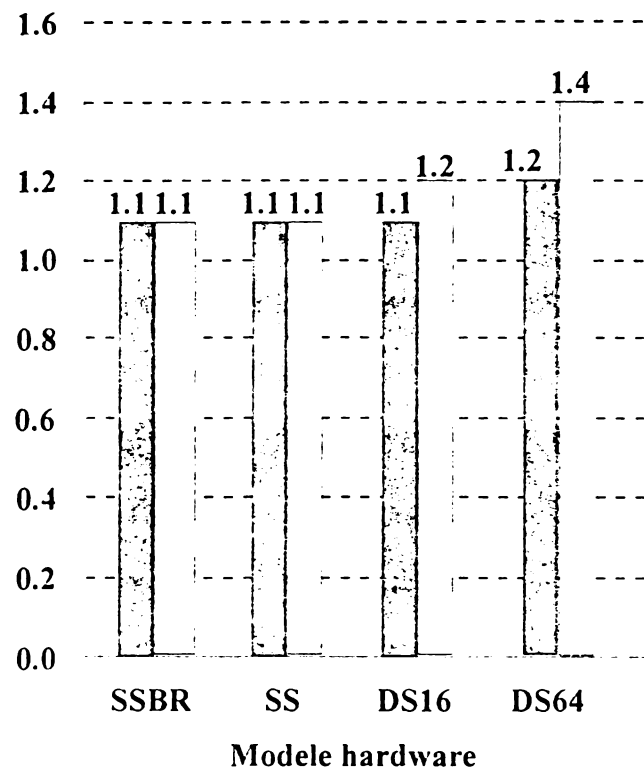
Se pot considera [HP96] patru modele hardware:

1. *SSBR (Statically Scheduled with Blocking Reads)* - Procesorul este planificat static și citirile ce eșuează în cache sunt imediat blocate.
2. *SS (Statically Scheduled)*- Procesorul este planificat static dar citirile nu cauzează blocarea procesorului până când nu se utilizează rezultatul.
3. *DS16 (Dynamically Scheduled with a 16-entry reorder buffer)*- Procesorul este planificat dinamic și are un tampon de reordonare ce permite înmagazinarea a până a 16 instrucții de orice tip, incluzând 16 instrucții de acces la memorie.
4. *DS64 (Dynamically Scheduled with a 64-entry reorder buffer)*- Procesorul este planificat dinamic și are un tampon de reordonare ce permite manevrarea a până la 64 de instrucții de oricare tip. Tamponul de reordonare este suficient de mare pentru a masca latența totală a eșecurilor cache.

În fig. 3.38 se prezintă performanțele pentru două dintre programele paralele de evaluare, LU și Ocean, menționate în paragrafele anterioare. Au fost luate în considerare cele patru metode hardware și două modele de consistență - TSO (Total Store Order) și RC (Release Consistency).

LU

Performanța relativ
la consistența secvențială



Ocean

Performanța relativ
la consistența secvențială

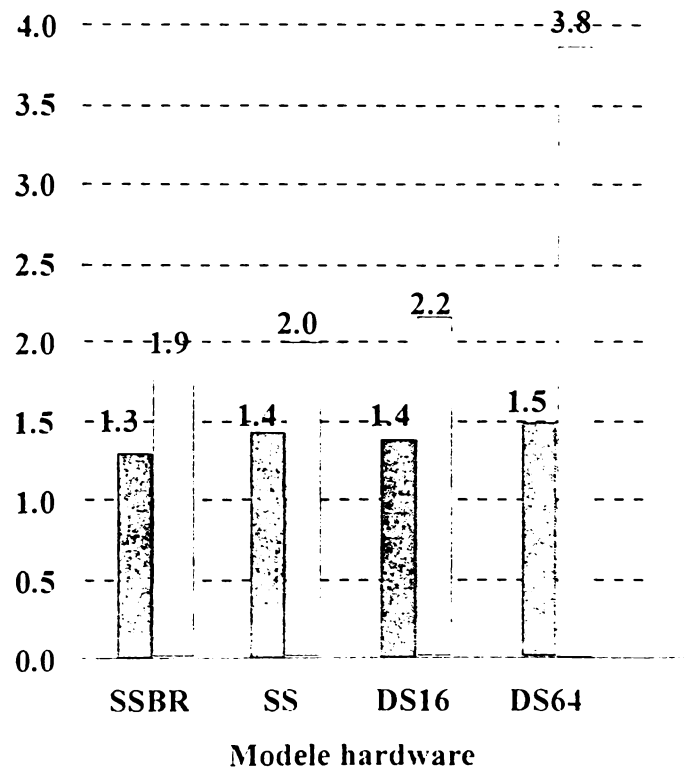


Fig. 3.38 Performanța modelelor de consistență relaxate prelevată la mai multe mecanisme hardware

Se remarcă că modelele relaxate oferă o performanță superioară la aplicații din clasa Ocean comparativ cu LU, deoarece Ocean are o mai mare rată de eșec și are o fracțiune semnificativă a eșecurilor la scriere.

În proiectare se preferă mărirea taliei memoriei cache înainte de-a include citiri fără blocare. Se reduce astfel semnificativ rata de eșec iar modelele relaxate devin atractive.

3.3 CONCLUZII

Utilizarea MCH private în cadrul SMP-MP poate reduce semnificativ timpul mediu de acces la memorie, dar aceste memorii introduc problematica complexă a menținerii coerenței acestora. Caracteristicile arhitecturale ce afectează mecanismul de menținere a coerenței memoriilor cache sunt:

- *strategia de detectare a coerenței* ce determină când și cum sunt definite în mod clar referirile la memorie pentru a se detecta existența unei posibile incoerențe între memoriile cache pentru date și memoria centrală. Mecanismele de detectare a coerenței în mod *dinamic* examinează adresele de memorie generate în timpul execuției. Eliminarea ambiguităților referitoare la memorie provoacă o rată redusă a eșecurilor în acces, dar mecanismele dinamice tind să producă un trafic de rețea ridicat datorită mesajelor necesare pentru a menține coerența. Schemele *statice* de detectare a coerenței, în contrast cu cele dinamice, examinează referirile la memorie în timpul compilării. *Autoinvalidarea* utilizată în mecanismele statice reduce traficul de rețea în comparație cu acela produs de strategiile dinamice de detectare a coerenței.
- *strategia de întărire a coerenței* ce trebuie să prevină referirea de către un procesor a unei date neactualizate. Acest deziderat se realizează prin forțarea procesoarelor ca să invalideze blocurile în memoriile cache aparținătoare. Dacă blocul este referit din nou, se va genera un eșec care va obliga procesorul să caute valoarea curentă a blocului, fie în memoria principală, fie în memoria aparținătoare altui procesor. Prin intermediul unei strategii de *actualizare*, noua valoare a unui bloc creat printr-o operațiune de scriere este automat distribuită către toate procesoarele împreună cu o copie cache-ată a blocului. Atunci când aceste procesoare fac din nou o referire a blocului, ele nu

generează altă cerere de deservire a eșecului în acces. Prin urmare, strategia de actualizare tinde să producă rate reduse de eșec în comparație cu strategia de invalidare. Această performanță se realizează însă prin creșterea semnificativă a traficului în rețea.

- *dimensiunea informației referitoare la partajarea blocului* care este întreținută de mecanismul de asigurare a coerenței. Această dimensiune are un impact direct asupra costului implementării mecanismului în termeni de număr de biți necesari pentru a înmagazina informația de partajare și la fel un impact direct asupra performanței sistemului de memorie. Pentru a reduce necesitățile de capacitate de memorie, mecanismul de menținere a coerenței poate să înmagazineze o cantitate relativ redusă de informație referitoare la procesoarele care au o copie a blocului cache-at. Mecanismul trebuie în continuare să transmită mesaje de invalidare atunci când numărul de procesoare ce partajează un bloc depășește resursele disponibile. Această abordare reduce necesitatea de memorie suplimentară dar crește traficul în rețea în comparație cu schemele bazate pe director care înmagazinează o informație completă despre modul de partajare. Recent, schemele propuse care conțin *tabel director etichetat* pot realiza o foarte mică creștere a capacității de memorare prin înmagazinarea informațiilor de partajare referitoare doar la acele blocuri care sunt cacheate la un moment dat. Acești directori pot să mențină un nivel redus al traficului în rețea deoarece sunt capabili să înmagazineze o informația suficientă referitoare la partajare pentru fiecare bloc cacheat.
- *dimensiunea blocului cache* reprezintă un factor important ce afectează performanța sistemului de memorie. Această dimensiune reprezintă numărul de cuvinte înmagazinate în memoria cache, ca și o singură unitate. Utilizarea blocurilor cache mai mari decât un singur cuvânt pot să permită procesoarelor să exploateze localitatea spațială tipică unui comportament ce face referire la memorie. Referirile la memorie într-un SMP tind să se disperseze printre procesoarele care reduc localitatea spațială disponibilă în comparație cu un sistem uniprosesor. În plus, blocurile mai mari decât un singur cuvânt introduc problema falsei partajări care tinde să provoace preferința pentru dimensiunile mici de blocuri cache. În câteva programe de aplicații este posibil să se reducă

rata de eșecuri în acces prin utilizarea blocurilor de cuvinte multiple dar studiile de simulare sugerează că blocurile de un singur cuvânt minimizează traficul în rețea prin reducerea atât a traficului de deservire a eșecurilor în acces cât și a traficului de invalidare.

În final, este important să se sublinieze că în arhitecturile moderne de SMP se pot incorpora diferite mecanisme de menținere a coerenței în memoriile cache. De exemplu prototipul DASH are încorporat în mecanismul de menținere a coerenței atât un mecanism de menținere a coerenței bazat pe monitorizarea magistralei cât și un mecanism de menținere a coerenței bazat pe director. Se poate oferi programatorului posibilitatea de a alege între strategii de întărire a coerenței bazate atât pe actualizare cât și pe invalidare. Suplimentar este posibil să se utilizeze informația obținută în timpul de compilare pentru a crește performanța mecanismului de menținere a coerenței: de exemplu de a reduce talia directorului prin reducerea numărului de pointeri de coerență ce sunt necesari ca să fie alocați și prin reducerea timpului necesar pentru a fi activați. Deoarece fiecare din factorii ce afectează mecanismul de menținere a coerenței produce facilități în termeni de rată de eșec în acces și de trafic în circuit, este clar că abordările hibride vor oferi cea mai bună oportunitate pentru creșterea performanței și reducerea costului de implementare a mecanismului de menținere a coerenței în SMP-MP de capacitate mare.

4. ANALIZA ARHITECTURILOR SISTEMELOR DE CALCUL MIMD TOLERANTE LA DEFECT

4.1 TERMINOLOGIE

În acest paragraf vom introduce câteva noțiuni de bază necesare în analiza sistemelor tolerante la defecte.

Un *sistem tolerant la defect* (*fault-tolerant system*) este atributul ce permite unui sistem să devină tolerant la defect.

Calcul tolerant la defect (*fault-tolerant computing*) este procesul prin care se efectuează calcule într-un sistem de calcul în manieră tolerantă la defect.

Un *defect* (*fault*) este o imperfecțiune fizică ce apare în cadrul unor componente hard/soft.

O *eroare* (*error*) reprezintă manifestarea defectului.

Defectarea (*failure*) reprezintă procesul prin care nu se efectuează acțiuni așteptate.

Reprezintă o performanță a unei funcții, în cantitate și în calitate, sub normalul admis.

În literatură [Prad 96] se citează *modelul triunivers* (*three universe model*), o adaptare a modelului *în patru universuri* (*four universe model*) dezvoltat de Avizienis (1982), în care noțiunile prezentate anterior pot fi mai bine înțelese.

Un defect apare într-un *univers fizic* (semiconductori, elemente mecanice, surse de alimentare, etc) și se manifestă printr-o eroare ce apare în *universul informațional* (cuvinte binare, imagini, voce digitală, etc).

Comportamentul eronat apare în *universul extern* sistemului. Se pot defini doi parametri importanți:

- *latența defectului* (*fault latency*) – timpul scurs între apariția unui defect și apariția erorii cauzate de defect
- *latența erorii* (*error latency*) – timpul scurs între apariția unei erori și apariția comportamentului eronat al sistemului.

Bazat pe modelul tri-univers, timpul total consumat între apariția defectiunii și apariția comportamentului eronat al sistemului va fi suma latenței defectiunii și a latenței erorii.

Pentru descrierea defectelor în mod adecvat, se utilizează patru atribute majore: *natura*, *durata*, *extinderea* și *valoarea Nelson și Carroll* [Prad 96].

Natura defectului specifică tipul defectului: hardware, software, în circuite analogice, etc.

Durata specifică durata manifestării defectului. Există *defect permanent*, *tranzitoriu* și *intermitent*.

Extinderea specifică dacă un defect este localizat la un modul dat sau afectează global modulele hard, soft sau le afectează pe ambele.

Valoarea unui defect poate fi *determinată* sau *nedeterminată*. Un defect determinat este unul a cărei stare rămâne nemodificată în timp până la apariția unei acțiuni externe. Un defect nedeterminat este unul a cărui stare la un moment T poate fi diferită de starea sa la $T + \Delta T$.

Există trei tehnici primare utilizate pentru a menține performanța normală a unui sistem într-un mediu în care defectele pot apare:

- *evitarea defectului (fault avoidance)*;
- *mascarea defectului (fault masking)*;
- *toleranța la defect (fault tolerance)*.

Toleranța la defect poate fi obținută prin mai multe tehnici:

Reconfigurarea – procesul de eliminare a entității cu funcționare defectuoasă dintr-un sistem și restaurarea sistemului la o stare operațională. Dacă se utilizează această tehnică atunci proiectantul trebuie să aibă în vedere *detectarea defectului*, *localizarea defectului*, *izolarea defectului* și *recuperarea din defect*.

Detectarea defectului (fault detection) este procesul de recunoaștere a apariției unui defect.

Localizarea defectului (fault location) este procesul de determinare a locului unde a apărut defectul.

Izolarea defectului (fault containment) este procesul de izolare a unui defect și de prevenire a extinderii lui prin sistem.

Recuperarea din defect (fault recovery) este procesul prin care se menține starea operațională a sistemului sau se recâștigă această stare prin reconfigurare.

Definiții echivalente pot fi date și în universul informațional operând asupra noțiunii de eroare în loc de defect.

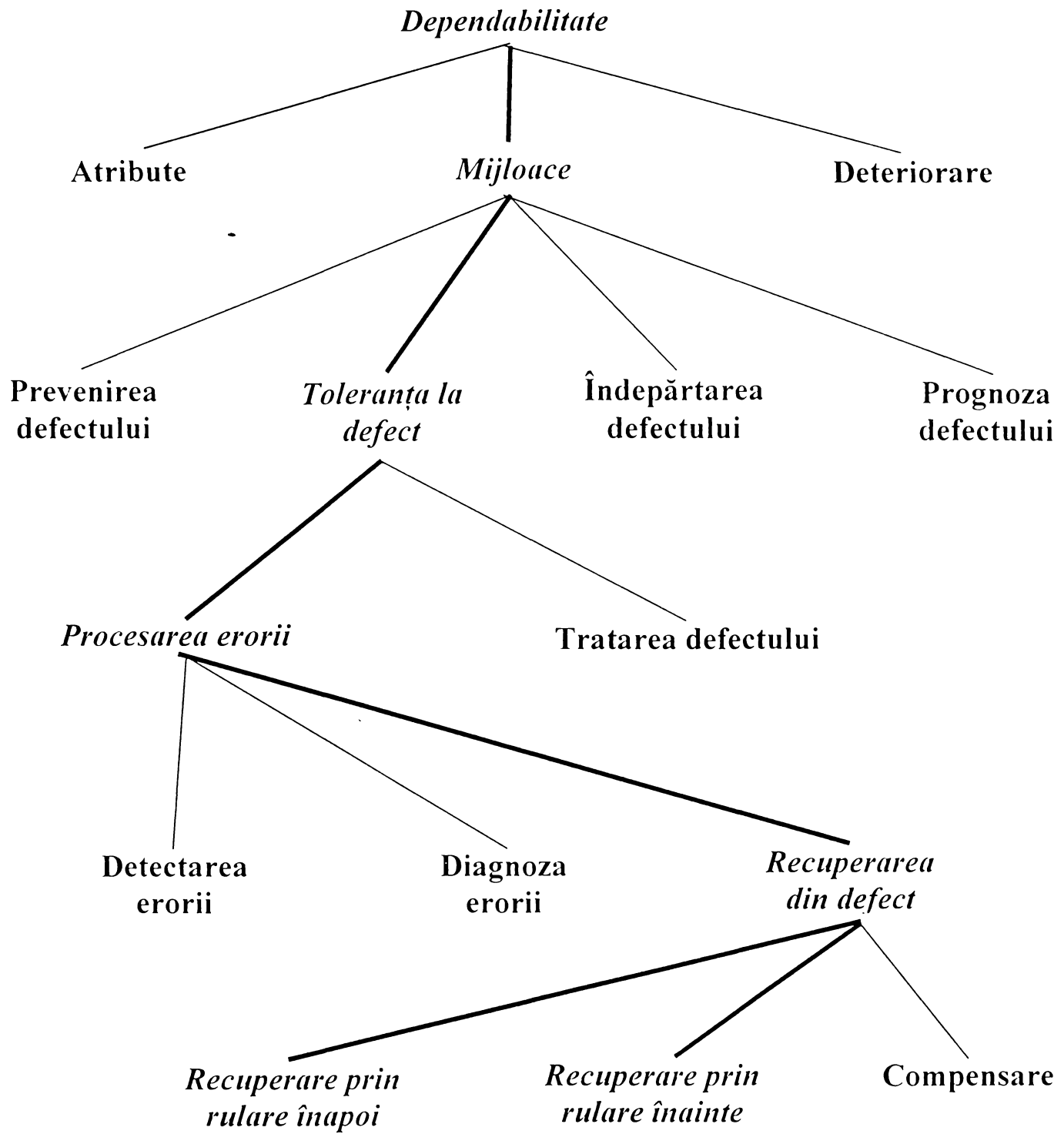


Fig. 4.1 Arborele dependabilității

4.1.1 Obiectivele toleranței la defect la sistemele multiprocesor

Toleranța la defect este un atribut structural în cadrul SMP pentru a atinge anumite cerințe de proiectare.

Cele mai proeminente cerințe adiționale, pe lângă cele vizând funcționalitatea și performanța, sunt :

dependabilitatea, fiabilitatea, disponibilitatea, siguranța, performabilitatea, mentenabilitatea și testabilitatea [Prad 96].

Toleranța la defect este unul dintre atributele sistemului capabil să satisfacă astfel de cerințe.

Dependabilitatea. (*Dependability*) Acest termen este utilizat pentru a încorpora conceptele de *fiabilitate, disponibilitate, siguranță, mentenabilitate, performabilitate și testabilitate*.

Dependabilitatea se poate defini simplu [Prad 96] ca fiind calitatea serviciilor furnizate de un anumit sistem de calcul. Componentele dependabilității sunt exemple de măsuri ale cantității de dependabilitate ale unui sistem. Arborele dependabilității este prezentat în fig.4.1.

Fiabilitatea (*Reliability*) unui sistem este o funcție de timp, $R(t)$, definită ca și probabilitatea condițională ca sistemul să funcționeze corect în intervalul de timp, $[t_0, t]$, dat fiind faptul că sistemul a funcționat corect la t_0 .

Cu alte cuvinte, *fiabilitatea* reprezintă probabilitatea ca sistemul să opereze corect într-un interval complet de timp.

Fiabilitatea este o probabilitate condițională ce depinde de starea operațională a sistemului la începutul intervalului de timp ales.

Nesiguranța în funcționare, opusul fiabilității, este o funcție de timp, $F(t)$, definită ca și probabilitatea condițională ca sistemul să evolueze incorect în intervalul, $[t_0, t]$, cunoscându-se faptul că la t_0 el funcționează corect.

Disponibilitatea. (*Availability*). *Disponibilitatea* este o funcție de timp, $A(t)$, definită ca și probabilitatea ca un sistem să fie operat corect și să fie disponibil la orice moment t .

Disponibilitatea diferă de siguranța în funcționare în sensul că trebuie să se manifeste în orice moment și nu în intervale temporale.

Siguranța. (*Safety*). *Siguranța* reprezintă probabilitatea, $S(t)$, ca un sistem să-și efectueze funcțiile corect sau să-și întrerupă funcționarea într-o manieră în care nu afectează operarea întregului sistem sau în care să compromită securitatea

personalului asociat cu sistemul. Siguranța este un parametru prin care se evidențiază modalitatea prin care - dacă nu funcționează corect - se ajunge sigur și controlat la starea de defecțiune.

Fiabilitatea și siguranța diferă deoarece fiabilitatea reprezintă probabilitatea ca un sistem să-și efectueze funcțiile corect. În timp ce siguranța este probabilitatea ca un sistem fie să-și efectueze funcțiile corect fie să-și întrerupă funcționarea într-o manieră care să nu dăuneze.

Performabilitatea. (*Performability*). *Performabilitatea* unui sistem este o funcție de timp, $P(L, t)$, definită ca și probabilitatea ca performanța sistemului în orice moment t să aibă același nivel L .

Performabilitatea diferă de fiabilitate în sensul că în timp ce fiabilitatea este o măsură globală ce se referă la faptul că toate funcțiile sunt efectuate corect, performabilitatea indică că un subset de funcții este efectuat corect.

Degradarea grațioasă este strâns legată de performabilitate și reprezintă abilitatea sistemului de a-și coborî nivelul de performanță pentru a compensa defecțiuni hard sau software.

Mentenabilitatea. (*Maintanability*). *Mentenabilitatea* reprezintă o măsură prin care se indică că un sistem poate fi reparat odată ce el a căzut.

$M(t)$ reprezintă probabilitatea ca un sistem defect să fie readus la o stare operațională într-o perioadă de timp t .

Testabilitatea. (*Testability*). *Testabilitatea* reprezintă abilitatea de-a testa anumite atribute ale sistemului. Ea este strâns corelată cu mentenabilitatea datorită importanței minimizării timpului cerut pentru identificarea și localizarea problemelor specifice.

4.1.2 Parametri de evaluare a dependabilității

Este necesară o abordare cantitativă a dependabilității prin introducerea unor parametri ca și:

Rata de defectare și funcția de fiabilitate (*Failure rate and the reliability function*)

Reprezintă numărul așteptat de defectări ale unui sistem într-o perioadă dată de timp.

Când are o valoare constantă se notează cu λ .

În literatură [Prad 96] se menționează definirea acestui parametru. Se consideră că $R(t)$ – fiabilitatea unui sistem – reprezintă probabilitatea condițională ca elementul

component să opereze corect în intervalul $[t_0, t]$, presupunând că a operat corect la t_0 . Presupunem că se rulează un test pe N componente identice plasând cele N componente în operare la t_0 și înregistrând numărul de componente căzute și a celor care sunt active la timpul t .

Fie $N_f(t)$ numărul de componente defecte la momentul t și $N_o(t)$ numărul de componente ce au operat corect la momentul t . Se presupune că odată ce un element s-a defectat el rămâne în această stare un timp nedefinit.

Fiabilitatea componentelor la momentul t este dată de

$$R(t) = \frac{N_o(t)}{N} = \frac{N_o(t)}{N_o(t) + N_f(t)} \quad (4.1)$$

și reprezintă probabilitatea ca un element component să fi supraviețuit în intervalul $[t_0, t]$.

Probabilitatea ca un element să nu fi supraviețuit în intervalul $[t_0, t]$ se numește *nefiabilitate* și este dată de

$$Q(t) = \frac{N_f(t)}{N} = \frac{N_f(t)}{N_o(t) + N_f(t)}. \quad (4.2)$$

În orice moment t există relația $R(t) = 1.0 - Q(t)$.

Prin urmare:

$$\frac{dR(t)}{dt} = \left(-\frac{1}{N}\right) \frac{dN_f(t)}{dt}$$

sau

$$\frac{dN_f(t)}{dt} = (-N) \frac{dR(t)}{dt}.$$

Funcția obținută

$$Z(t) = \frac{1}{N_o(t)} \frac{dN_f(t)}{dt} \quad (4.3)$$

se numește *funcție hazard, rata hazardului sau funcția rată a defectării*.

Exprimarea funcției $Z(t)$ în funcție de $R(t)$, respectiv $Q(t)$

$$Z(t) = \frac{\frac{dR(t)}{dt}}{R(t)} \quad \text{și} \quad Z(t) = \frac{\frac{dQ(t)}{dt}}{1 - Q(t)} \quad (4.4)$$

unde:

$\frac{dQ(t)}{dt}$ - funcția de densitate a defectării.

În literatură [Prad 96] se definesc $R(t) = e^{-\lambda t}$ ca fiind *legea defectării exponențiale* și *funcția rată a defectării* asociată unei distribuții Weibull

$$Z(t) = \alpha \lambda (\lambda t)^{\alpha-1} \quad (4.5)$$

unde α și λ reprezintă constante ce controlează variația funcției rată a defectării în timp.

Timpul mediu până la defectare (Mean Time To Failure)

Reprezintă un parametru util pentru a specifica calitatea unui sistem. *MTTF* este timpul așteptat ca un sistem să opereze înainte ca prima defectare să apară.

Dacă avem N sisteme identice plasate în operare la $t=0$ și dacă măsurăm timpul în care fiecare sistem operează înainte de a cădea, atunci timpul mediu este *MTTF*.

Dacă fiecare sistem i operează pentru un timp t_i , înainte de a întâlni prima defectare, *MTTF* este:

$$MTTF = \sum_{i=1}^N \frac{t_i}{N}. \quad (4.6)$$

Relația între *MTTF* și $R(t)$ este:

$$MTTF = \int_0^{\infty} R(t) dt \quad (4.7)$$

valabilă pentru $R(\infty) = 0$.

Timpul mediu de reparare (Mean Time To Repair)

MTTR reprezintă timpul mediu necesar pentru a repara un sistem. Dacă al i -lea defect din cele N necesită un timp t_i pentru a-l repara, atunci *MTTR* se estimează astfel:

$$MTTR = \frac{\sum_{i=1}^N T_i}{N}. \quad (4.8)$$

Dacă considerăm că μ - *rata de reparare* - este numărul de reparații pe o perioadă de timp atunci

$$MTTR = \frac{1}{\mu}. \quad (4.9)$$

Timpul mediu între defectări (Mean Time Between Failure)

Trebuie de la început să menționăm semnificațiile lui *MTTF* și *MTBF*.

MTTF reprezintă timpul mediu scurs până la prima defectare a unui sistem în timp ce *MTBF* este timpul mediu consumat între defectările unui sistem.

MTBF este calculat prin medierea timpului dintre defectări, incluzând orice timp necesar pentru a repara sistemul și a-l plasa apoi în starea operațională. Cu alte cuvinte, fiecare din cele N sisteme este operat pentru un timp oarecare T și numărul de defectări întâlnit de către al i -lea sistem se înregistrează ca n_i , media numărului de defectări este: -

$$n_{med} = \sum_{i=1}^N \frac{n_i}{N}$$

iar în final

$$MTBF = \frac{T}{n_{med}} . \quad (4.10)$$

Cu alte cuvinte, *MTBF* este timpul de operare total T , împărțit cu numărul mediu de defectări produse în timpul T .

Experiența arată că:

$$MTBF = MTTF + MTTR . \quad (4.11)$$

Acoperirea defectului (Fault coverage)

Există două definiții ale acoperirii defectului: una intuitivă și cealaltă matematic formalizată. Definiția intuitivă este aceea că acoperirea reprezintă o măsură a abilității sistemului de a efectua detectarea defectului, localizarea lui, izolarea și/sau acoperirea lui.

Primele patru tipuri de acoperire a defectului sunt: *acoperirea detectării defectului*, *acoperirea localizării defectului*, *acoperirea izolării defectului* și *acoperirea recuperării din defect*.

Acoperirea detectării defectului este o măsură a abilității sistemului de a detecta defecte. De exemplu, o cerință a sistemului poate fi că o anumită fracțiune din toate defectele trebuie să fie detectată; acoperirea detectării defectului este o măsură a sistemului de a îndeplini o astfel de cerință. *Acoperirea izolării defectului* este o măsură de a izola defecte. *Acoperirea recuperării din defect* reprezintă o măsură a abilității sistemului de a reveni din defecte și de a menține o stare operațională.

Acoperirea la defect este matematic defintă ca și probabilitatea ca, existând un defect, sistemul să-și revină. Problema fundamentală la acoperirea la defect este că se

calculează foarte laborios. Aici nu se va insista asupra tehnicilor de calcul al acestui parametru.

4.2 TOLERANȚA LA DEFECT PRIN UTILIZAREA REDUNDANȚEI

4.2.1 Toleranța la defect prin utilizarea *redundanței statice*

Există trei aplicații distincte ale *redundanței statice* în mediul multiprocesor: *redundanța pentru fiabilitate și disponibilitate*, *redundanța pentru siguranță* și *redundanța pentru tolerarea defectelor non-clasice*.

4.2.1a Redundanța pentru *fiabilitate și disponibilitate*

Prima abordare a toleranței la defect în SMP este la nivel de sistem, implicând utilizarea redundanței statice unde N copii ale unui modul efectuează calcule simultan pentru a fi verificate și validate. Există și alte abordări la alte nivele. Fiecare procesor poate fi replicat și rezultatul calculului fiecărui procesor validat sau întregul SMP poate fi replicat și rezultatul combinat validat. O a treia opțiune împarte cele P procesoare ale SMP în P/N grupuri de N procesoare, fiecare grup validându-și rezultatele înainte de comunicarea cu alte grupuri.

De asemenea, pentru a furniza comunicări robuste, toate tranzacțiile critice între grupuri pot fi replicate și ulterior validate. Replicarea și validarea în diverse grade a fost utilizată în mai multe proiecte din perioada anilor '60 [Prad96].

4.2.1b Redundanța pentru *siguranță*

Fiabilitatea se referă la probabilitatea ca sistemul să producă o ieșire corectă. Siguranța este definită ca fiind probabilitatea ieșirii sistemului de a fi corectă sau ca eroarea apărută să fie detectabilă. O siguranță ridicată este asigurată făcând neglijabilă probabilitatea apariției unei erori nedetectate la ieșire. Când se sesizează la ieșire o eroare ce nu poate fi detectată, se poate declanșa o procedură de recuperare. O schemă de toleranță la defect trebuie să fie aleasă atunci când îndeplinește cerințele de fiabilitate – siguranță. Trebuie menționat că atunci când un sistem este fiabil este de asemenea și sigur. Reciproca nu este întotdeauna valabilă. Un sistem poate fi sigur într-un grad înalt doar dacă are un timp mediu până la defectare mic. În general siguranța unui sistem este cel puțin egală cu fiabilitatea.

Strategiile de proiectare pot atinge atât o fiabilitate cât și o siguranță ridicate, utilizând sisteme formate din module identice multiple ce efectuează toate aceleași funcție.

Modelul generic este ilustrat în figura 4.2. Sistemul este format utilizând module identice multiple, toate executând aceeași funcție. Ieșirile modulelor sunt arbitrate iar

ieșirea arbitrului constituie ieșirea sistemului. Aceasta are două componente: *ieșirea de date* și *fanionul de indicare a nesiguranței* ce marchează eroare la datele de ieșire. Fanionul de indicare a nesiguranței, atunci când este poziționat, indică apariția unei erori la ieșirea de date. În această situație se apelează la un mecanism de recuperare. Astfel de sisteme se numesc SMR (Safe Modular Redundant). Un sistem SMR conținând n module este un sistem n SMR. Fiabilitatea și siguranța unui sistem n SMR depinde atât de fiabilitatea modulului individual cât și de strategia particulară de arbitrare selecționată.

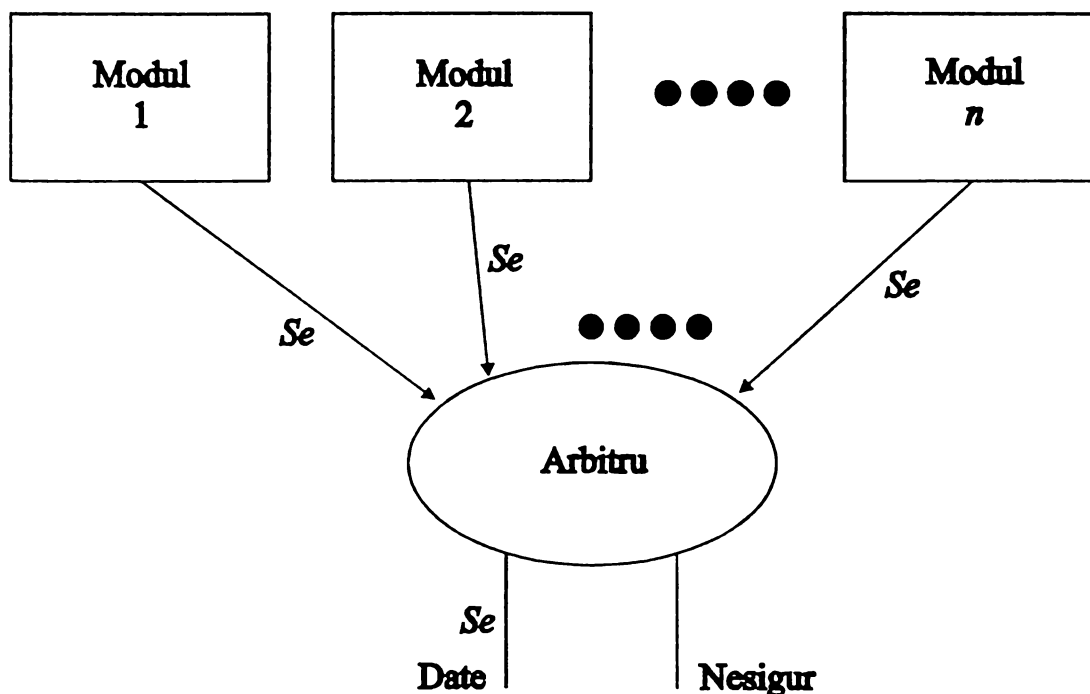


Fig. 4.2 Un sistem SMR

O strategie de arbitrare reprezintă funcția implementată de către un arbitru pentru a decide ce ieșire corectă există și când erorile de la ieșirile de modul exced capacitatea de corectare de o astfel de manieră încât ieșirea corectă nu poate fi furnizată. Când ieșirea corectă nu poate fi furnizată, fanionul de indicare a nesiguranței poate fi poziționat de către arbitru și astfel se poate renunța la evaluarea ieșirii.

În general în cadrul unei abordări uzuale se presupune că fiabilitatea și siguranța în funcționare ar putea să fie ameliorate prin adăugarea unor capacități suplimentare de redundanță. În literatură [Prad96] se arată că o redundanță augmentată nu trebuie în mod necesar să contribuie la o îmbunătățire atât a fiabilității cât și a siguranței în funcționare. Mai exact, pentru un model de eroare simetric, care modelează erori în circuitele VLSI trebuie să se rețină următoarele: în sisteme cu module non-redundante

(module cu ieșiri necodate) există un sistem n SMR în comparație cu care nici un sistem $(n+1)$ SMR nu poate atinge o fiabilitate mai ridicată și nici o siguranță în funcționare mai bună, ceea ce înseamnă că adăugarea unui modul nu ameliorează în mod simultan fiabilitatea și siguranța. Oricum, fiind dat un sistem n SMR întotdeauna există un sistem $(n+2)$ SMR cu o fiabilitate mai ridicată (siguranță) și cel puțin același nivel de siguranță (fiabilitate). Astfel, adăugarea unui minim de două module poate fi întotdeauna efectivă în ameliorarea atât a fiabilității cât și a siguranței. Prin urmare se poate demonstra că dacă există module individuale având o redundanță inclusă, întotdeauna există un sistem $(n+1)$ SMR cu o fiabilitate mai mare și cu o siguranță în funcționare comparabilă cu oricare alt sistem n SMR dat. Astfel, dacă se utilizează module cu capacitate inclusă de detectare, adăugarea unui singur modul poate fi suficientă. Dacă modulele nu au capacitate de detectare inclusă atunci adăugarea unor capacități suplimentare de redundanță nu garantează întotdeauna ameliorarea siguranței. Prin urmare, adăugarea unui singur modul poate să întărească atât fiabilitatea cât și siguranța doar în condițiile existenței unor capacități incluse de detectare a erorii. Acest lucru e valabil doar în ipoteza simplificatoare că toate configurațiile de eroare sunt aproximativ simplificatoare. În sistemele critice din punct de vedere al siguranței, una dintre problemele majore este reprezentată de acuratețea modelului de defect. Prin urmare, printre modelele de defect permanente tranzitorii și intermitente se utilizează un model cu grad mare de generalitate denumit *modelul Bizantin*.

4.2.1c Redundanța pentru tolerarea defectelor *non-clasice*

În timp ce cea mai simplă abordare pentru mascarea redundanței implică calcule triplicate, unde votarea asupra rezultatelor are loc pentru a tolera o singură eroare permanentă sau una intermitentă, votarea devine mult mai complicată în prezența defectelor arbitrare non-clasice. În mediile ce necesită o extrem de înaltă fiabilitate, metodele de votare trebuie să fie furnizate cu posibilitatea de a oferi toleranță împotriva defectelor arbitrare. Chiar și defectele maligne, unde două sau mai multe noduri defecte pot coopera, trebuie să fie tolerate.

În acest scop, a fost propus *modelul de defect al Acordului Bizantin* în 1982 de către Lamport, Shostak și Pease [LSP82]. Acest model permite o comportare arbitrară pentru modulul ce defectează. Un astfel de modul poate să producă nu numai valori incorecte dar de asemenea să transmită valori diferite la destinații diferite în loc de valori identice cum ar fi de așteptat. Sub acest model, două sau mai multe module ce

defectează pot accidental sau deliberat să coopereze de o astfel de manieră încât să producă un efect mult mai dezastruos asupra sistemului, chiar dacă defectele maligne la nivel de modul au fost tolerate.

Agreementul Bizantin, reprezintă o analogie tradițională utilizată pentru prezentarea dificultăților întâlnite în găsirea unui agreement între procesoarele ce operează corect în prezența unor procesoare ce defectează, chiar dacă aceste procesoare manifestă un comportament ce poate fi interpretat ca fiind intenționat malign.

În SMP critice din puncte de vedere al siguranței, unii proiectanți susțin acest model pentru a depăși efectul defectelor complexe dependente de timp ce sunt mult mai ușor acoperite de acest model.

În raport cu nivelul de fiabilitate cerut, proiectarea unui sistem bazat pe agreementul bizantin este totuși simplu. Un astfel de sistem conține doar un număr minim, în prealabil specificat, de procesoare și interconexiuni, furnizează sincronizarea lor și utilizează protocoale simple de schimb de informații.

Cu titlu de exemplu, amintim că pentru un protocol ce tolerează f defecte, bazat pe un sistem Bizantin f -rezilient, trebuie îndeplinite următoarele cerințe:

- Cel puțin $(3f+1)$ membri trebuie să participe
- Cel puțin $(2f+1)$ căi disjuncte de comunicație trebuie să existe între membri
- Cel puțin $(f+1)$ runde de comunicare trebuie să aibă loc

Astfel, un protocol pentru un sistem rezilient de tip 1-Bizantin trebuie să execute două runde de comunicare între cel puțin patru membri sincronizați conectați la cel puțin trei căi disjuncte de comunicație.

Pentru a specifica cerințele ce trebuie satisfăcute de un protocol bazat pe Agreementul Bizantin, în literatură [Prad96] se utilizează exemplul clasic al unui sistem de calcul compus din noduri ce comunică între ele. Nodul principal trebuie să trimită o comandă de tip ATAC sau RETRAGERE către toate nodurile subordonate. Nodul principal și nodurile subordonate utilizează un protocol de tip Agreement Bizantin pentru a conveni asupra comenzilor ce emană de la nodul principal. Nodurile trebuie să comunice prin mesaje și să ajungă la o decizie finală. Un nod subordonat care este defect poate, prin urmare, să transmită mesaje incorecte. Cerințele ce trebuie să fie satisfăcute de un algoritm bazat pe acest tip de agreement poate fi rezumat astfel:

C1: toate nodurile subordonate corecte agreează aceeași comandă

C2: dacă nodul principal funcționează corect toate nodurile subordonate ce funcționează corect sunt de acord asupra comenzii furnizate de către nodul principal

De exemplu se consideră patru noduri A , B , C și D . În ipoteza că nodul D dorește să transmită un mesaj către A , B și C Acordul Bizantin impune ca toate nodurile ce funcționează corect într-un set $\{A, B, C \text{ și } D\}$ să fie de acord asupra aceluiași mesaj. Dacă D este un nod principal și A , B și C noduri subordonate, D poate trimite un mesaj fie ATAC, fie RETRAGERE, situație prezentată în fig. 4.3.

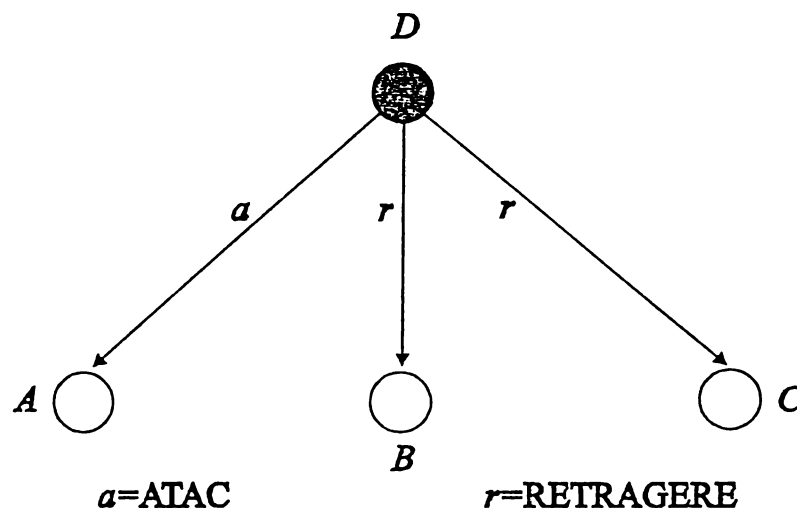


Fig.4.3 Acordul Bizantin

Nodurile subordonate libere de defect pot să fie de acord asupra comenzii trimise de D dacă D este liber de defect. Dacă D este defect atunci nodurile subordonate libere de defect sunt solicitate să fie de acord asupra unei comenzi identice, indiferent de comanda recepționată de la D . Se poate considera următorul scenariu: D defectează și A , B și C sunt libere de defect. D trimite comanda ATAC spre A și retragere spre B și C . Pentru a detecta un astfel de caz nodurile subordonate trebuie să schimbe între ele mesajul recepționat de la D . După ce schimbul a fost efectuat, fiecare dintre subordonați are trei copii ale acestui mesaj. Subordonatul acceptă majoritatea acestor copii ca fiind valoare agreată. Astfel, fiecare dintre nodurile subordonate înțelege comanda RETRAGERE ca fiind o valoare agreată, situație reprezentată în fig. 4.4.

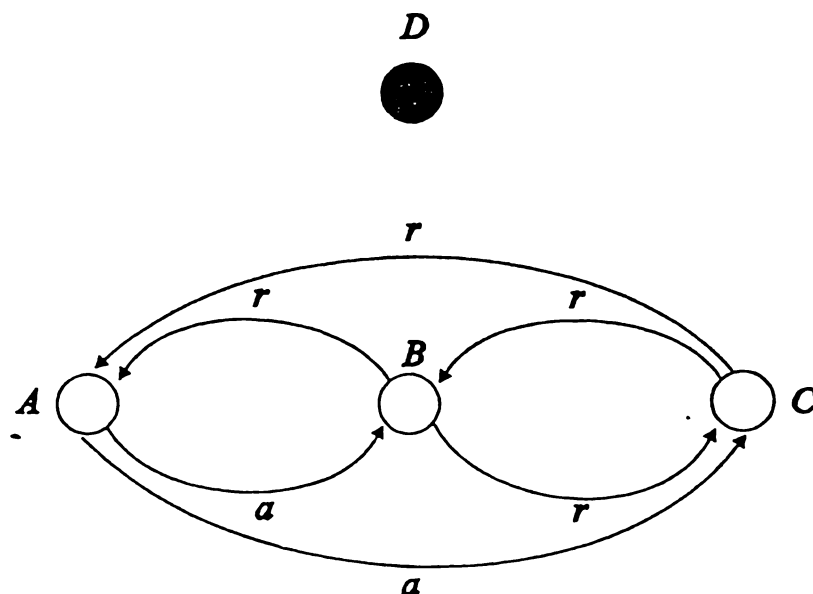


Fig. 4.4 Acordul Bizantin obținut după aplicarea protocolului

Un alt scenariu posibil este acela în care două dintre nodurile subordonate, A și B sunt defecte și nodul principal este liber de defect. Nodul principal emite comanda ATAC către toate cele trei noduri subordonate. Nodurile A și B se comportă în mod defectuos și trimit comanda RETRAGERE nodului C , situație prezentată în figura 4.5.

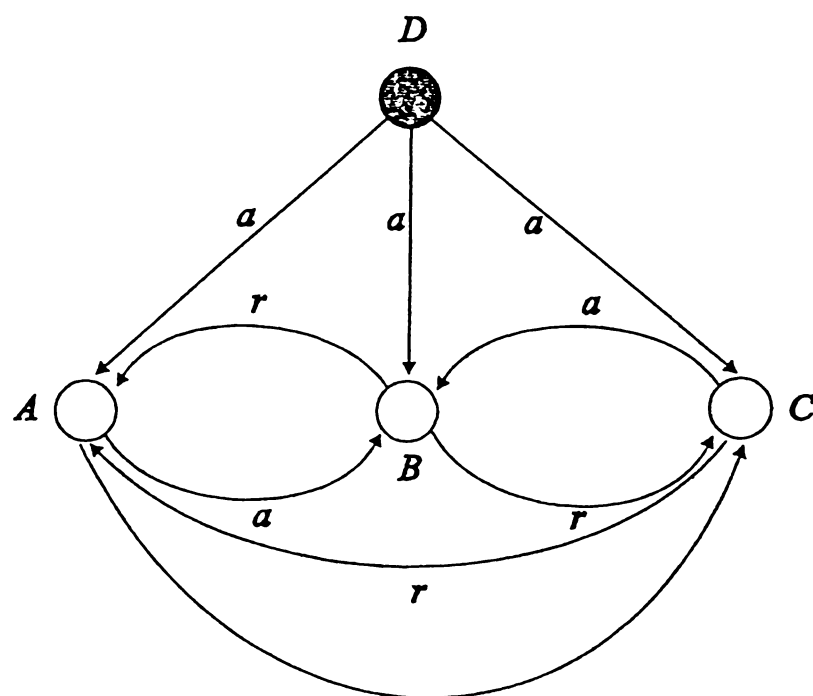


Fig. 4.5 Un acord Bizantin imposibil

Prin urmare, C va obține comanda retragere ca o majoritate a celor trei copii ale comenzii recepționate. Astfel, două noduri subordonate defecte pot încălca acest algoritm bizantin. Este demonstrat că orice sistem ce poate tolera t defectări după

modelul bizantin trebuie să conțină cel puțin $3t+1$ noduri. Aceasta înseamnă că pentru a tolera defectarea celor două noduri subordonate, sistemul prezentat ca și exemplu ar trebui să aibe un nod principal și șase noduri subordonate. Aceasta înseamnă că un model de reziliență pentru m defecte necesită cel puțin $(3m+1)$ noduri ce comunică între ele. În literatură se prezintă un algoritm $OM(m)$ bazat pe algoritmul Bizantin ce tolerează m defecte de tip bizantin. Numărul total n de noduri trebuie să satisfacă relația $n \geq 3m+1$. Termenii ATAC și retragere pot fi interpretați ca fiind mesaje binare.

Algoritm A. $OM(0)$

1. Valoarea emisă de nodul principal este transmisă fiecărui nod subordonat.
2. Fiecare nod subordonat utilizează valoarea recepționată de la nodul principal sau utilizează valoarea RETRAGERE dacă nu a primit nici un mesaj.

Algoritm B. $OM(t>0)$

1. Valoarea provenind de la nodul principal este transmisă la fiecare nod subordonat.
2. Pentru fiecare i , fie v_i valoarea pe care nodul subordonat unic a recepționat-o de la nodul principal sau să fie RETRAGERE dacă nu s-a recepționat nimic. Subordonat lui i acționează ca și nod principal în algoritmul $OM(t-1)$ pentru trimite valoarea v_i la fiecare din celelalte $(n-2)$ noduri subordonate.
3. Pentru fiecare i și pentru fiecare $j \neq i$ fie v_j valoarea nodului subordonat i transmisă către nodul subordonat j în pasul 2 (utilizând algoritmul $OM(t-1)$ sau RETRAGERE dacă nu s-a recepționat nimic. Nodul subordonat i utilizează valoarea *majoritate* $(v_1, v_2, \dots, v_{n-1})$.

În literatură [Prad96] se propune o formă diferită a protocolului Agreementului Bizantin denumită *Protocol de Agreement Degradabil*. În acest protocol se definesc doi parametri m și u unde $u \geq m$. Parametrul m reprezintă numărul de defecte până la care protocoalele definite anterior pot fi efective, iar u este cel mai mare număr de defecte până la care protocolul degradat poate fi util. Acest protocol definit de parametrul m și u este definit ca și *Protocolul de Agreement Degradabil m/u* . În continuare prezentăm acest protocol.

Protocolul de Agreement Degradabil

Fie f numărul de noduri defecte.

- Dacă $f \leq m$, atunci condițiile D.1 și D.2 enunțate mai jos trebuie să fie satisfăcute.

- Dacă $m < f \leq u$ atunci condițiile D.3 și D.4 enunțate mai jos trebuie să fie satisfăcute.
- (D.1) Dacă transmițătorul este liber de defect atunci toate receptoarele libere de defect trebuie să fie de acord asupra valorii transmise de transmițător.
- (D.2) Dacă transmițătorul este defect, atunci toți receptorii liberi de defect trebuie să fie de acord asupra unei valori identice.
- (D.3) Dacă transmițătorul este liber de defect, atunci receptorii liberi de defect trebuie să fie partiționați în cel mult două clase. Receptorii liberi de defect dintr-o clasă trebuie să fie de acord asupra valorii transmise de transmițător, și receptorii liberi de defect din cealaltă clasă trebuie să accepte valoarea de defect.
- (D.4) Dacă transmițătorul este defect, atunci receptorii liberi de defect trebuie să fie partiționați în cel mult două clase. Receptorii liberi de defect dintr-o clasă trebuie să fie de acord asupra valorii de defect și receptorii liberi de defect din cealaltă clasă trebuie să agreeze asupra unei valori identice.

Fie N numărul de noduri din sistem. Dacă $N > 2m + u$ atunci agreementul degradabil m/u definit anterior asigură că cel puțin $m+1$ noduri libere de defect agreează o valoare identică chiar dacă numărul de defecte este mai mare decât m dar cel mult u . Astfel, se poate obține o degradare grațioasă. De notat că degradarea grațioasă este posibilă până la u defecte chiar când $u \geq N/3$.

În mod specific, pentru a obține un agreement degradabil m/u sistemul trebuie să consistă din cel puțin $2m+u+1$ noduri, inclusiv transmițătorul, și aceste $2m+u+1$ noduri sunt suficiente.

Mulți algoritmi bazați pe agreementul bizantin au fost propuși în diverse ipoteze referitoare la sistemul luat în considerare. De exemplu sistemele FTTP (*Fault Tolerant Parallel Processor*) utilizează astfel de protocoale în diverse circumstanțe.

4.2.2 Toleranța la defect prin *redundanța dinamică*

O tehnică la toleranță la defect cu costuri scăzute în SMP, utilizează *redundanța dinamică*. În această tehnică se prevăd capacități interne de detectare a defectelor în procesoarele ce defectează. Când se detectează un defect, sistemul este reconfigurat prin activarea unui procesor de rezervă sau a unei capacități de procesare alternative. Această abordare furnizează un mecanism ce indică apariția unei erori în timp ce

sistemul operează. Acest lucru poate fi obținut utilizând fie resurse hardware, fie redundanță temporală. Odată ce prezența unui procesor defect a fost detectată, este demarată o procedură de localizare și de diagnosticare a erorii. Procesorul defect este înlocuit cu un procesor de rezervă prin proceduri de reconfigurare. În final, se efectuează o recuperare din defect, unde procesorul de rezervă, utilizând informații tipice de exploatare a PV-urilor, preia calculele procesorului defect din locul de unde acesta a ieșit din sistem.

4.3 DETECTAREA DEFECTELOR ÎN SISTEMELE MULTIPROCESOR

În vederea detectării defectelor în sistemele multiprocesor există câteva abordări [Prad96], [MP89], [PF90]. Acestea sunt:

- *testare off-line* planificată pentru depistarea defectelor permanente;
- *duplicarea și compararea* utilizând redundanța de spațiu și de timp;
- aplicarea unor *tehnici de codificare și de diagnosticare*.

4.3.1 Detectarea defectelor prin *duplicare și comparare*

Operațiile de *duplicare și comparare* pot fi efectuate în mai multe moduri. Fiecare procesor a SMP poate fi duplicat și rezultatele comparate înainte de a fi comunicate perechilor de procesoare. Această tehnică a fost aplicată la blocurile constitutive ale lui Intel 432.

Un mod echivalent de a trata această opțiune constă în divizarea celor P procesoare ale unui SMP în $P/2$ perechi. Memoria globală comună care constă din M module de memorie este și ea divizată în $M/2$ perechi. Duplicarea modulelor de memorie nu se utilizează deoarece codificarea poate oferi tehnici mult mai eficiente de detecție și corecție pentru aceste module. Comparatoarele pot fi încapsulate în interiorul fiecărui modul procesor și de memorie și rezultatul calculelor poate fi realizat hardware la fiecare ciclu de mașină în timpul fiecărei operații cu referire la memorie.

Dacă se detectează o eroare de către o pereche de procesoare, ambele sunt deconectate de la sursa de alimentare și calculele se efectuează pe cele $P - 2$ procesoare ce rămân și care sunt configurate ca și $(P - 2)/2$ perechi.

Intr-o altă abordare, duplicarea și compararea sunt utilizate cu *redundanță temporală*. O astfel de tehnică se utilizează când constructorul nu-și poate permite duplicarea

datorită restricțiilor de cost, greutate, putere, etc. ca de exemplu în aplicațiile spațiale. În principiu, cele P procesoare dublează același calcul paralel în timp asupra unui set diferit de procesoare, astfel încât fiecare procesor afectează doar o copie a rezultatului. Acest rezultat se poate obține prin divizarea celor P procesoare în $P/2$ grupe de câte două procesoare, fiecare din ele comparând rezultatele perechii înainte de a comunica cu alte procesoare.

În figura 4.6 se prezintă un exemplu de *graf de sarcină (task graph)* care este duplicat și comparat pentru detectarea defectelor. Cazul (a) prezintă graful original ce este mapat pe opt procesoare. Cazul (b) prezintă duplicarea grafului și maparea fiecărei copii pe un set disjunct de patru procesoare. Rezultatele celor două copii sunt comparate de două ori, de fiecare dată de câte un procesor din fiecare subgrupă. În cea mai defavorabilă situație procedura de mai sus poate să consume 100% un timp suplimentar.

Dacă graful nu are dependențe, sarcinile pot fi mai eficient mapate pe toate procesoarele unui SMP și toate procesoarele sunt ocupate tot timpul. În prezența dependențelor, unele procesoare vor fi în stare de așteptare activă deoarece există sarcini nepregătite pentru prelucrare. În astfel de situații se poate mapa graful original pe $P/2$ procesoare (ca și în fig. 4.6 b), dându-se o mai bună utilizare a procesoarelor, iar cele $P/2$ procesoare ce rămân sunt utilizate pentru calculul dublat al grafului. Compararea poate fi efectuată și prin tehnici software, cu mențiunea că în această situație granularitatea calculului, înainte de comparare, poate fi mai mare decât un ciclu mașină, deoarece compararea software durează mai mult decât cea hardware. Operația de comparare poate fi efectuată eficient prin utilizarea *punctelor de control (check point)*.

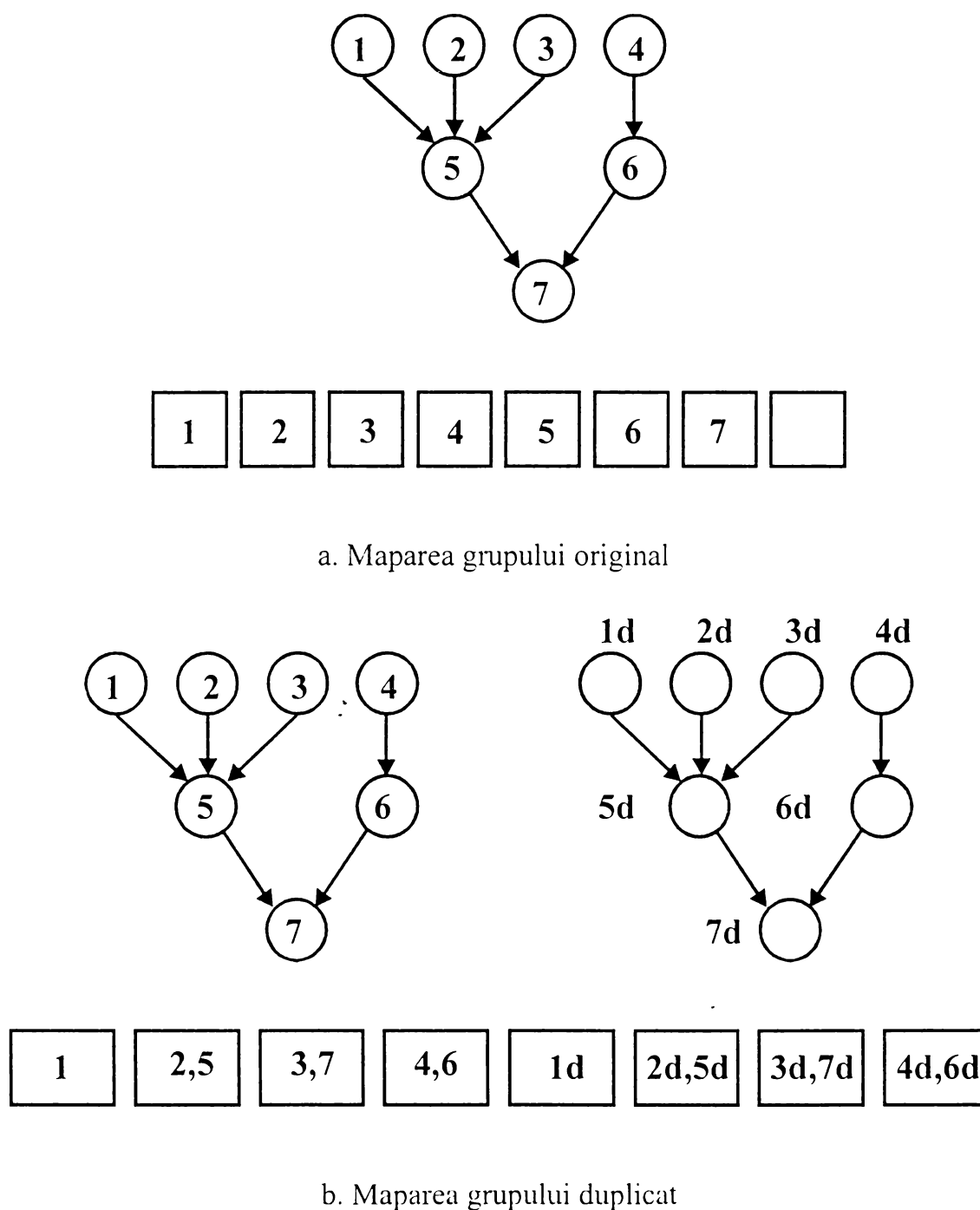


Fig. 4.6 Exemplu de mapare a grafurilor duplicate pe seturi disjuncte de procesoare.

4.3.2 Detectarea defectului utilizând *diagnosticarea și tehnici de codificare*

Diagnoza utilizează programe de diagnosticare ce pot fi utilizate pentru detectarea defecțiunilor hard în procesoare, memorii și RIN. Această metodă detectează defecte permanente dar nu poate fi aplicată pentru defecte tranzitorii și intermitente. Pentru detectarea defectelor tranzitorii și intermitente se utilizează tehnici concurente de

detectare a defectelor ce utilizează metode de *autotestare* [Prad 96]. *Tehnicile de codificare* pot furniza detectări de defect la preț scăzut pentru magistrale, memorii și registre. Verificări simple de paritate în magistrale și registre sunt un exemplu de utilizare efectivă a codificării. Violarea parității provoacă situația de excepție. Revenirea este manipulată prin reîncercarea efectuării instrucției (instruction retry) sau prin proceduri similare. Arhitecturile cele mai moderne oferă capabilități de reîncercare pentru aproape toate instrucțiile. Un exemplu îl constituie sistemele IBM unde verificări intensive de paritate sunt prevăzute intern pentru a obține o diagnoză rapidă.

Implementarea toleranței la defect a decodificatoarelor și a codificatoarelor este o chestiune de maximă importanță.

În codificarea utilizată în tehnicile de calcul, se poate avea uneori încredere în informația apriori referitoare la localizarea erorii.

4.3.3. Controlul erorii în memoriile de mare viteză.

Codurile corectoare de erori se utilizează atât la memoriile cache de mare viteză cât și la memoriile principale datorită faptului că acestea sunt vulnerabile la erori soft. Pentru ca un cod să fie util pentru memoriile de mare viteză, structura sa trebuie să permită o codificare și o decodificare paralelă, rapidă, complexitatea circuitului de verificare a parității utilizată în componența decodificatoarelor și a codificatoarelor poate fi un factor major în determinarea vitezei. Examinarea structurii *matricii de verificare la paritate*, cunoscută sub numele de matrice \mathbf{H} , dă informații referitoare la complexitatea circuitului de verificare a parității.

De exemplu, se consideră un cod de lungime 6, $n = 6$, cu trei biți informaționali, $k = 3$ și trei biți de verificare, $r = 3$. Cele două matrici \mathbf{H} prezentate în continuare au aceeași capabilitate de corectare a erorii, d_0, d_1 și d_2 reprezentând data iar c_0, c_1 și c_2 , verificarea de paritate. Bitul de verificare a parității c_i reprezintă verificarea de paritate efectuată asupra lui d_i dacă coloana i a matricii \mathbf{H} conține un 1 în coloana d_i . Deoarece toate coloanele în matricea \mathbf{H} sunt disjuncte, codul poate corecta toate erorile singulare, dar circuitul de verificare la paritate pentru \mathbf{H}_1 este mai puțin complex decât pentru \mathbf{H}_2 . Se consideră ecuațiile pentru c_0, c_1 și c_2 - biții de control la paritate – separat pentru \mathbf{H}_1 și \mathbf{H}_2 . \mathbf{H}_1 necesită doar trei porți SAU-EXCLUSIV

cu 2 intrări pentru calculul lui c_0, c_1 și o poartă SAU-EXCLUSIV cu trei intrări pentru calculul lui c_2 . Datorită porții SAU-EXCLUSIV cu trei intrări, codificatorul și decodificatorul utilizând \mathbf{H}_2 vor fi mai lente și mai complexe.

$$\mathbf{H}_1 = \begin{bmatrix} & d_0 & d_1 & d_2 & c_0 & c_1 & c_2 \\ 1 & 0 & 1 & 1 & 0 & 0 & \\ 1 & 1 & 1 & 0 & 1 & 0 & \\ 1 & 1 & 0 & 0 & 0 & 1 & \end{bmatrix} \quad \begin{aligned} c_0 &= \overline{d_0 \oplus d_2} \\ c_1 &= \overline{d_1 \oplus d_2} \\ c_2 &= \overline{d_0 \oplus d_1} \end{aligned} \quad (4.12)$$

$$\mathbf{H}_2 = \begin{bmatrix} & d_0 & d_1 & d_2 & c_0 & c_1 & c_2 \\ 1 & 0 & 1 & 1 & 0 & 0 & \\ 1 & 1 & 1 & 0 & 1 & 0 & \\ 1 & 1 & 0 & 0 & 0 & 1 & \end{bmatrix} \quad \begin{aligned} c_0 &= \overline{d_0 \oplus d_1} \\ c_1 &= \overline{d_0 \oplus d_2} \\ c_2 &= \overline{d_0 \oplus d_1 \oplus d_2} \end{aligned} \quad (4.13)$$

Numărul de 1-uri din matricea \mathbf{H} determină complexitatea totală a circuitului de verificare a parității. Cu cât mai puțini de 1, cu atât mai simplu circuitul. Prin urmare, circuitul de verificare a parității cel mai lent corespunde cu linia cu numărul maxim de 1-uri și numărul de 1-uri în orice linie la un minim.

Coduri detectoare/corectoare de erori pentru un bit. În memoriile de mare viteză *coduri corectoare de erori pentru un singur bit (single-bit error-correcting codes)* și *coduri detectoare de erori pentru doi bit (double-bit error-detecting codes)*- coduri SEC-DED- sunt cele mai utilizate. Aceasta pentru că multe memorii RAM sunt organizate pe un bit; prin urmare orice defect într-un circuit integrat se manifestă ca și o eroare pe un bit.

Se consideră două r -tuple de pondere impară. Ponderea se referă la numărul de 1-uri. De notat că suma a două r -tuple de pondere impară este o r -tuplă de pondere pară. Datorită acestei proprietăți, un cod SEC-DED cu r biți de verificare poate fi construit după cum urmează: aici matricea \mathbf{H} constă din vectori coloană de pondere impară. Astfel, sindromul de eroare pe doi-bit este o r -tuplă de pondere pară. Acest cod este

diferit de *codul SEC-DED Hamming* original, a cărui matrice \mathbf{H} are în plus la matricea \mathbf{H} a codului SEC adăugat, un vector coloană de 1-uri.

Pentru a obține un circuit de codificare/decodificare de mare viteză, numărul de 1-uri în fiecare linie este egal sau cât mai apropiat de numărul mediu. Lungimea de cod maximă este egală cu numărul maxim de r -tuple de pondere pară și prin urmare $n = 2^{r-1}$. În continuare se prezintă un exemplu din această clasă de coduri:

$$\mathbf{H} = \left[\begin{array}{cccc|cccc} d_0 & d_1 & d_2 & d_3 & c_0 & c_1 & c_2 & c_3 \\ 1 & 1 & 1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 & 0 & 1 \end{array} \right] \quad (4.14)$$

$$\text{Sindrom: } S = D \cdot \mathbf{H}^T [S_0, S_1, S_2, S_3]$$

$$S_0 = \overline{d_0 \oplus d_1 \oplus d_2 \oplus c_0}$$

$$S_1 = \overline{d_0 \oplus d_1 \oplus d_3 \oplus c_1}$$

$$S_2 = \overline{d_0 \oplus d_2 \oplus d_3 \oplus c_2}$$

$$S_3 = \overline{d_1 \oplus d_2 \oplus d_3 \oplus c_3}$$

Pentru orice cod SEC-DED, probabilitatea de funcționare eronată la apariția unor erori triple sau pe mai mulți biți trebuie să fie minimală. O funcționare eronată se referă la o decodificare eronată ce se transformă într-un cod corector eronat.

În concluzie, codurile SEC-DED cu coloane de pondere impară au două avantaje practice: simplitate constructivă a codicatorului/decodicatorului și posibilitate scăzută a unei decodificări eronate.

Coduri detectoare-corectoare de erori pe byte. Unele memorii sunt organizate în capsule de mare densitate, având o lățime de b -bit. Dacă apare o defecțiune, cuvântul extras se comportă de parcă ar avea un bloc de b -bit (byte) defect. În acest tip de aplicații este de dorit să existe un cod corector de erori capabil de detectare/corectare a erorilor la nivel de byte.

Coduri corectoare de erori la nivel de byte. În literatură [FP 90] se face o prezentare exhaustivă a acestei clase de coduri. Vom prezenta modul de construcția matricii \mathbf{H} pentru un cod corector de eroare la nivel de byte, astfel: Se aleg ca și coloane ale

matricii \mathbf{H} toate r -tuplele nonzero de elemente dintr-un câmp finit F , în particular dintr-un câmp Galois $GF(2^b)$, astfel că nici o coloană din \mathbf{H} să fie multiplu al altei coloane. Astfel, fiecare pereche de coloane este linear independentă și prin urmare există o distanță Hamming minimă de 3 pentru cod. Codul cunoscut sub acronimul de SbEC, este capabil să corecteze toate erorile singulare la nivel de cuvânt de b -bit.

Implementarea acestui tip de cod necesită transformarea matricii \mathbf{H} din $GF(2^b)$ într-o formă binară. Prin utilizarea unor funcții polinomiale primitive binare, $g(x)$ de grad b , se poate defini o matrice nesingulară \mathbf{T} exprimată ca și o matrice binară ($b \times b$). Setul acestor matrici reprezintă un câmp ce este izomorfic la $GF(2^b)$. Prin urmare, elementele lui $GF(2^b)$ pot fi exprimate ca și $\{0, \mathbf{T}, \mathbf{T}^2, \mathbf{T}^3, \dots, \mathbf{T}^{2b-2}, \mathbf{T}^{2b-1} = \mathbf{I}\}$, unde \mathbf{I} este matricea ($b \times b$) de valoare 0.

Simbolurile n și k denotă lungimile de cod și lungimile de informație ale acestui cod peste $GF(2^b)$. Codul derivat SbEC este un cod (N, K) în formă binară, unde $N = n \cdot b$ și $K = k \cdot b$. Similar, numărul de biți de verificare este $R = r \cdot b = n - k$.

Lungimea maximă (în bit) a acestui cod de tip Hamming este dată de $N_H = b \cdot n = b \cdot ((2^{br} - 1)/(2^b - 1))$.

Deoarece codurile SbEC nu garantează detectarea erorilor aleatorii pe doi bit, respectiv pe doi bytes, aceste coduri nu sunt utilizate în sistemele de calcul. În schimb, calculatoarele utilizează coduri corectoare de erori pe cuvânt de b -bit și coduri detectoare de erori duble pe cuvânt de b -bit, denumite SbEC-DbED.

4.4.STRATEGII DE RECUPERARE PENTRU SISTEME MULTIPROCESOR

Recuperarea unui SMP tolerant la defect, reprezintă ultimul pas în procedura de revenire a sistemului după ce defectul a fost detectat, izolat și sistemul reconfigurat în jurul procesorului (procesoarelor) defect(e) [Sita93], [VP93], [WFP90].

Tehnicile de recuperare sunt diferite pentru SMP-MP respectiv SMP-TM.

Soluția cea mai trivială de recuperare în cazul apariției unui defect în unul sau mai multe procesoare în cadrul unui multiprocesor o reprezintă terminarea programului în lucru și reexecutarea întregului program, pentru toate procesoarele, de la început. Este procedura *restartului global* (*global restart*). Evident această soluție este costisitoare

datorită timpului consumat. Practic timpul de rulare se dublează. Din această cauză această soluție nu este considerată ca și o procedură standard de recuperare din defect.

Recuperarea din defecțiune poate avea trei forme [Prad 96]:

- *recuperare prin rulare înapoi (rollback recovery)*: Transformarea stării eronate constă din aducerea înapoi a sistemului la o stare ocupată tot timpul înainte de apariția erorii;
- *recuperare prin rulare-înainte (forward recovery)*: Transformarea stării eronate constă din găsirea unei noi stări, din care sistemul să poată opera, de obicei într-un mod degradat;
- *compensarea (compensation)*: Starea eronată conține suficientă redundanță pentru a permite transformarea sa într-o stare liberă de eroare.

Problema recuperării prin *rulare înapoi* este complexă în cazul SMP deoarece procese multiple pot accesa memoria și pot avea copii diferite sau eronate ale aceluiași variabile și prin urmare producând o stare inconsistentă când eroarea este detectată. În consecință schemele în care se aplică tehnica recuperării prin *rulare înapoi* trebuie să asigure memorarea informațiilor referitoare la stările libere de erori într-un loc sigur de unde să poată fi apoi regăsite și utilizate pentru restartarea programului dintr-o stare consistentă. În esență recuperarea prin *rulare înapoi* constă din salvarea periodică a unui proces în timpul evoluției corecte a unui program de o astfel de manieră încât restartarea să fie posibilă dintr-o astfel de stare salvată. Stările salvate ale procesului se numesc *puncte de verificare (checkpoints)* -PV- iar tehnica de creare și de exploatare a astfel de puncte se numește *checkpointing*.

Starea unui sistem, în general, implică setul de registre ale procesorului, numărătorul de program, starea memoriei cache și uneori chiar memoria cel puțin până la ultimul punct de verificare – PV.

Aceste PV-uri sunt memorate în unități de memorie sigure, discuri sau memorii protejate de coduri corectoare de erori.

A. Tehnica recuperării prin rulare înapoi utilizând puncte de verificare

Această tehnică ce utilizează PV-uri este o metodă de a asigura toleranță la defect împotriva defectelor tranzitorii și a celor intermitente. În continuare se vor descrie câteva metode de implementare a acestei tehnici [BP92], [BP93], [Ng91].

- **Puncte de verificare bazate pe memoria cache a procesului**

Datele punctului de verificare (*checkpoint data*) reprezintă starea celui mai recent PV, iar *datele active* reprezintă ceea ce este accesat după verificarea PV-ului.

Ierarhia memoriei oferă o soluție naturală pentru această problemă în sensul că datele active pot fi plasate în registrele procesorului și în memoria cache în timp de datele PV pot fi plasate în memoria principală a sistemului. Astfel operația de bază pentru starea unui PV constă în a rămâne în memoria principală, iar procesorul avansează actualizând doar memoria cache și registrele sale. Un PV este inițiat prin o *tentativă eșuată de acces la cache (cache miss)* ce forțează o linie cache modificată spre memoria principală. Un PV este obținut prin copierea tuturor registrelor modificate și a liniilor cache în memoria principală. O *rulare înapoi (rollback)* se obține prin îndepărtarea tuturor registrelor modificate și a liniilor cache și prin restaurarea stării anterioare utilizând datele PV-ului. Această schemă devine operațională prin utilizarea unei politici de actualizare de tip “*copiază-înapoi*”. Astfel memorările sunt efectuate în cache și nu în memoria principală. Memoria principală este actualizată când apare un acces ratat, concretizându-se în înlocuirea unui linii viciate (*dirty-line*). O linie sau o pagină se consideră *viciată (dirty)* dacă ea a fost modificată de către deținătorul ei de la ultimul său PV. Dacă acest lucru nu s-a produs, atunci linia (pagina) este *curată (clean)*. Prin urmare un acces ratat la cache forțează un PV.

În fig. 4.7 se prezintă separarea stării active și a stării PV-ului.

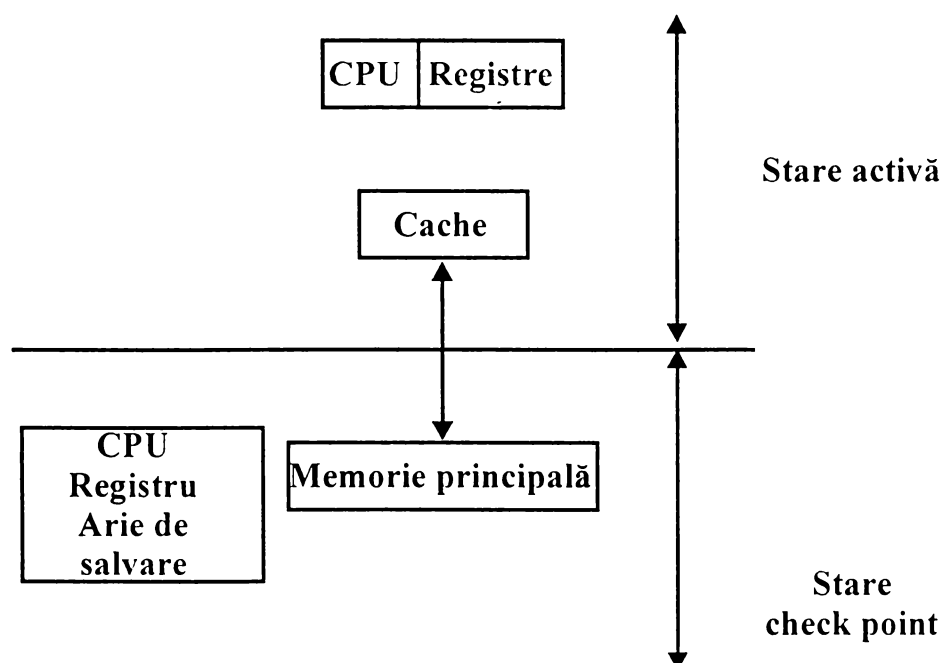


Fig. 4.7 Crearea de PV bazate pe procesor și recuperarea prin revenire

Un controler ce controlează PV-ul și procesul de revenire trebuie să aibă următoarele abilități care să:

- salveze registrele procesorului în memorie (sau în cache);
- aducă toate liniile cache modificate în memoria principală;
- încarce registrele CPU din memorie (sau din cache);
- invalideze toate liniile modificate în cache.

Dacă sistemul este într-o stare oprită cu memoria principală conținând un PV valid, cu registrele CPU pentru această stare salvate în memoria principală (aria de salvare) atunci controlerul încarcă registrele din aria de salvare și sistemul începe să ruleze. Evident inițial registrele și cache-ul sunt goale. Datele încep să fie încărcate la cerere în cache. Cache-ul utilizează o politică *copiază înapoi*, așa că scrierile în memorie nu depășesc spațiul cache. Orice schimbări în stare de la ultimul PV se produc doar în cache și în registre.

Aceste două seturi de locații înmagazinează starea activă. Dacă se detectează o eroare tranzitorie (de exemplu o eroare de paritate pe o magistrală internă), sistemul trebuie readus până la un PV anterior. Deoarece acesta este conținut în întregime în memoria principală a sistemului, registrele CPU și memoria cache trebuie restaurate la starea lor de la începutul intervalului. Memoria cache distinge între liniile nemodificate și cele modificate așa că doar liniile modificate sunt luate în considerare, ele reprezentând starea activă. Prin urmare, liniile cache modificate trebuie să fie selectiv invalidate.

În mod alternativ, întreaga memorie cache ar putea fi inițializată. Următorul pas îl constituie reîncărcarea registrelor CPU din starea salvată în memorie. Acum sistemul a fost readus înapoi la PV anterior și rularea programului poate reîncepe.

În orice schemă ce utilizează PV-uri este esențial să existe posibilitatea manipulării defectelor în timpul exploatarei PV-urilor. În literatură [Prad 96] se descriu mai multe scheme de principiu ce sunt utilizate în calculatoare de firmă și în care se regăsesc principiile enunțate anterior. În fig 4.8 se prezintă șase pași pentru redirectionarea informației din cache către memoria principală. Această tehnică poate tolera o singură defecțiune în orice punct.

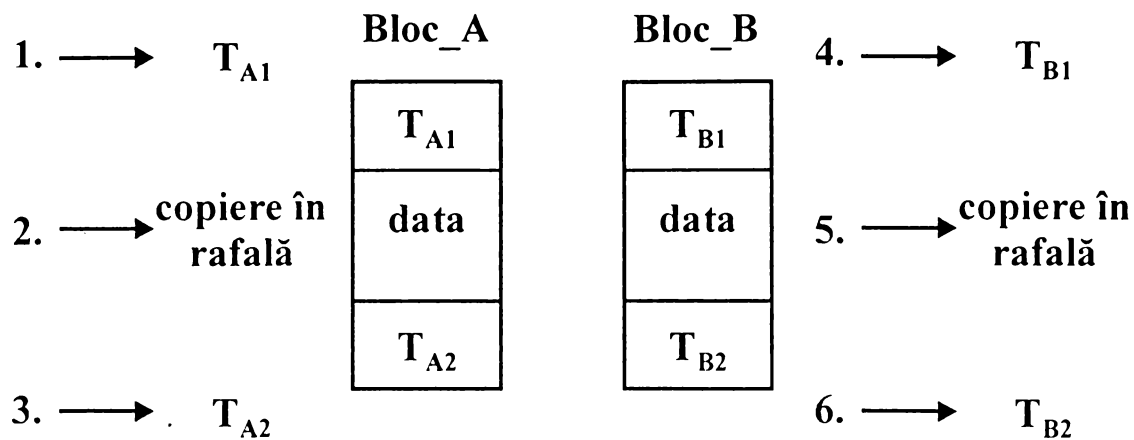


Fig. 4.8 Tehnici tolerante la defect pentru copierea în cascadă a cache-urilor

Conținutul cache-ului este scris secvențial în cele două blocuri de memorie. O *marcă ce marchează timpul (time stamp)* este scrisă în blocul de memorie înainte de *transferul în rafale (flush)*. Aceasta poate fi T_{A1} sau T_{B1} . După transferul în rafale eticheta va fi T_{A2} sau T_{B2} . Deoarece cele patru mărci de timp sunt scrise secvențial, o defecțiune poate fi descrisă de aceste etichete ce sunt diferite. Utilizând aceste patru mărci de timp, se poate corecta o defecțiune o procesului de checkpointing. TABEL VII descrie diferitele condiții ale mărcilor de timp găsite după defectare. Deși, din punct de vedere logic, există patru mărci de timp distincte, în practică doar trei mărci sunt necesare deoarece $T_{A2} = T_{B1}$. În consecință, pașii de recuperare subliniați în liniile 3 și 4 din TABEL VII sunt aceiași. În mod clar, dacă ei sunt egali, atunci procesul de creare și gestionare a PV-urilor a fost efectuat cu succes. Dacă T_{A1} este singura marcă de timp scrisă, atunci defecțiunea a apărut în timpul transferului spre blocul A.

Noul PV (în blocul A) este incomplet și PV anterior (blocul B) trebuie să fie utilizat. Astfel, procedura de recuperare va reveni înapoi la PV anterior prin copierea blocului A în blocul B.

Dacă ambii T_{A1} și T_{A2} sunt scriși, $T_{A1} = T_{A2}$, atunci noul PV a fost scris cu succes. Această situație derivă în două cazuri distincte de defectare ambele necesitând aceeași acțiune. În primul rând, dacă T_{B1} nu a fost scris, atunci defecțiunea a apărut între două transferuri. Oricum, doar dacă T_{B2} nu a fost scris, eroarea apare în timpul transferului spre B. În celălalt caz, noul PV din blocul A este copiat în blocul B.

TABEL VII. Condiții de defectare

Condiția	Defecțiune	Acțiune
$T_{A1} = T_{A2} = T_{B1} = T_{B2}$	Nici una	Nici una
$T_{A1} > T_{A2} = T_{B1} = T_{B2}$	Transfer în A	Copiază B în A
$T_{A1} = T_{A2} > T_{B1} = T_{B2}$	Intre	Copiază A în B
$T_{A1} = T_{A2} = T_{B1} > T_{B2}$	Transfer în B	Copiază A în B

Talia memoriilor cache afectează direct frecvența creerii și gestionării PV-urilor. Cu cât mai mare este talia unui cache, cu atât mai mică va fi necesitatea de a forța un PV din cauza condiției de memorie cache plină. Pe de altă parte, frecvența utilizării PV-urilor este direct proporțională cu rata de defectare și cu abilitatea de recuperare. În literatură [Prad 96] se oferă mai multe soluții de ameliorarea tehnicilor de creare și gestionare a PV-urilor.

- **Puncte de verificare virtuale**

O deficiență majoră a tehnicilor de creare și exploatare a PV-urilor bazate pe memoria cache este aceea că frecvența de creare și testare a PV-urilor, ce depinde de talia cache-ului, poate fi foarte ridicată afectând negativ performanțele.

Când frecvența de checkpointare este ridicată, o alternativă o constituie plasarea PV-urilor și a recuperării prin *rulare înapoi* într-o memorie virtuală, prin extragerea lor din cache-ul procesorului.

Metoda *punctelor de verificare virtuale (virtual checkpoints)* este utilizată pentru a reduce penalitățile în performanțe ale unei testări a PV-urilor cu frecvență ridicată, prin utilizarea unui mecanism de translație și a unei memorii adiționale. Metoda combină concepte provenite de la două tehnici de recuperare utilizate la baze de date [Prad 96]: *paginarea din umbră (shadow paging)* și *paginarea geamână (twin paging)*.

Termenul de *date ale punctului de verificare* se referă la starea în cel mai recent PV iar cel de *date active* la acelea ce sunt accesate după PV. Conceptul de *suportare a datelor active* este implementat prin alocarea dinamică a unei copii secundare a paginii virtuale. Paginile active pot fi identificate prin utilizarea unui contor de PV-uri asociat fiecărei pagini.

Un aspect fundamental al acestei tehnici este acela al suportării a două clase de date în sistemul memoriei virtuale (date active și PV-uri). Fiecare clasă suportă înmagazinarea tradițională pe două niveluri a memoriei virtuale –memoria reală și discul de paginare. Există două cerințe importante ce se impun pentru această tehnică. Prima este abilitatea de a detecta toate paginile active. Aceasta deoarece o *rulare înapoi* la un PV anterior se obține prin purjarea tuturor paginilor active. A doua cerință este abilitatea de a face permanente toate paginile active la momentul verificării PV-urilor. Aceste pretenții sunt satisfăcute având un contor de PV global (V) ce acoperă toate contoarele de date și de PV-uri locale (v) pentru pagini individuale. În esență, contorul global este copiat în contorul local la fiecare referire. Prin urmare prima cerință este îndeplinită deoarece toate datele active au contorul local v egal cu contorul global V . Contorul V este incrementat când este preluat un PV. Astfel, a doua cerință este îndeplinită deoarece toate paginile active versiuni PV când V este incrementat. Figura 4.9 ilustrează aceste concepte de bază. Pagina virtuală k nu a fost referită de la PV anterior. Pagina j a fost accesată în intervalul curent și are atât o versiune activă cât și una PV.

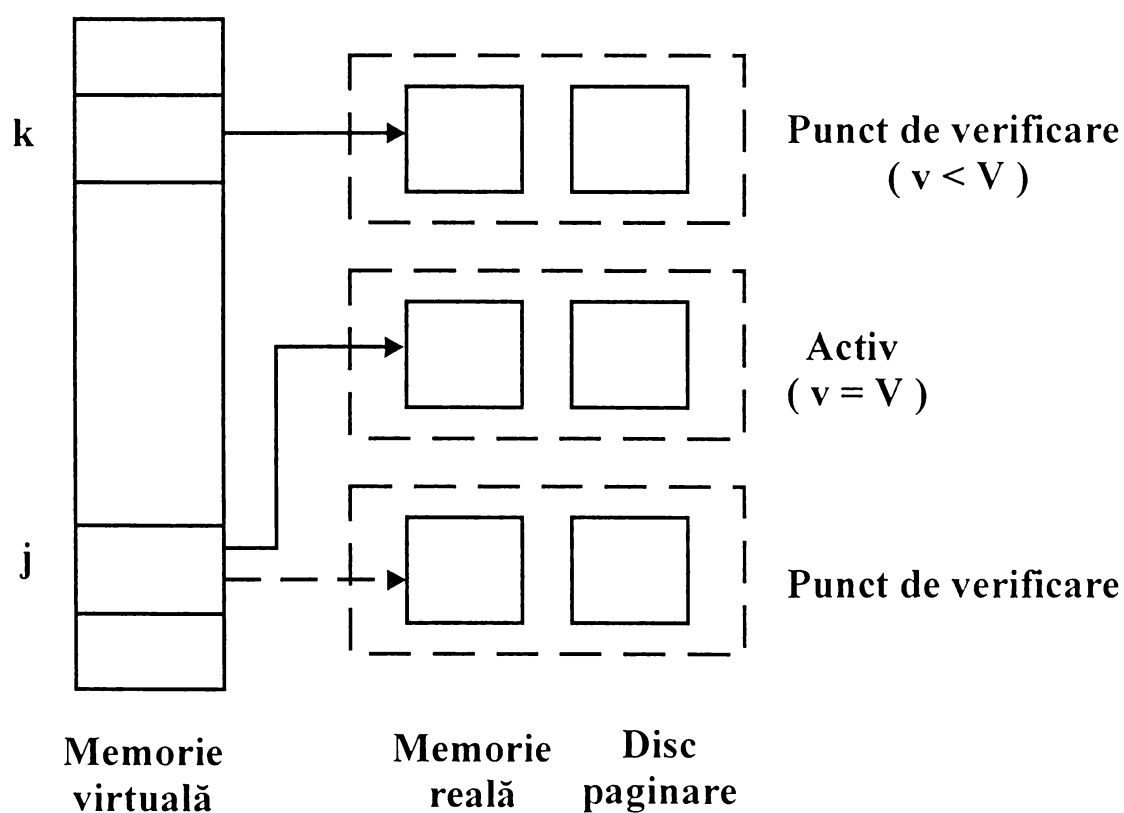


Fig.4.9 Conceptul de bază

Pentru a implementa această schemă, se prevăd câteva câmpuri care trebuie să fie adăugate la mecanismul de translatăre virtuală. Trebuie să existe informații pentru a distinge între o versiune activă și una de PV a fiecărei pagini. În continuare vom descrie translatărea unei singure pagini virtuale. Aceasta înseamnă că notarea va exclude orice referință la numărul paginii virtuale. Pentru fiecare pagină virtuală, mapările sunt replicate și referite ca și m_0 și m_1 . O mapare m_l conține mapări în cadrul real (r_l) și copia disc (d_l) ca și în fig. 4.10.

Fiecare pagină are un câmp de un bit, l , care poate fi imaginat ca și un comutator ce indică cea mai recentă mapare utilizată (de exemplu m_0 și m_1). Astfel, notarea m_l se referă la maparea ce a fost utilizată ultima. În plus, fiecare pagină are un număr (v) a unui PV local de k -bit ce conține o copie a unui număr (V) de PV global în timpul celei mai recente referiri. Numărul (V) este o valoare globală care este incrementată la fiecare PV. Fiecare cadru real, desemnat de r_l , conține un bit de schimb (c_l) care indică dacă rata din cadrul real a fost schimbată și este diferită de copia disc la d_l .

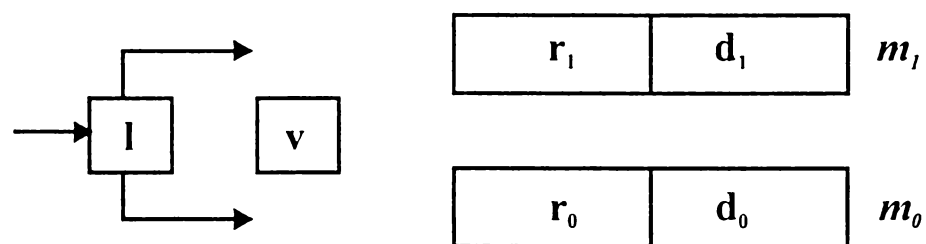


Fig. 4.10 Maparea unei singure pagini

O caracteristică importantă a acestei scheme este că acțiunile de preluare ale unui PV nu sunt concentrate la timpul actual ale PV ci sunt distribuite în timpul ce urmează PV-ului. Procesarea pentru paginile individuale este amânată până la ultima referire la pagina ce urmează după PV. Pentru a determina dacă procesarea amânată trebuie să apară, valorile V și v trebuie să fie comparate la fiecare referire. Astfel, când o pagină este referită, apare fie cazul unde procesarea PV-ului trebuie să apară ($v \neq V$) sau un acces la pagina activă ($v = V$). Fig. 4.11 prezintă o situație unde PV-urile au fost prelevate la timpii t_{c1} și t_{c2} .

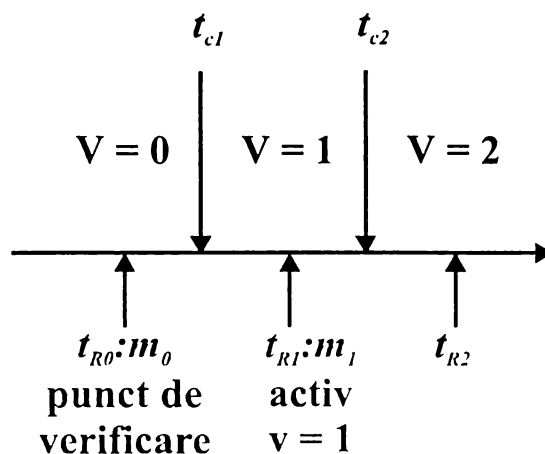


Fig.4.11 Cazul 1: Prima referire după PV

Să considerăm evenimentele la t_{R2} . Pagina activă adresată de m_1 ($l=1$) a fost ultima referită la t_{R1} (astfel, $v=1$). Referința la t_{R2} este prima referință după PV (deoarece $v \neq V$) și conținutul paginii m_1 trebuie să fie rezervat ca și pagină PV. Mai departe, conținutul paginii m_1 trebuie să fie utilizată în timp ce *resursele* vechii pagini m_0 de tip PV sunt utilizate. Odată ce data validă a fost copiată în m_0 , bitul l este inversat și devine $l=0$ astfel că m_0 devine pagina activă și m_1 devine PV. În final, V este copiat în numărul PV-ului local pentru această pagină așa că la următorul acces în intervalul de procesare al PV-ului, apare o translată normală. Fig. 4.12 prezintă situația la următoarea referire în acest interval de preprocesare al PV-ului. O referire la timpul t_{R3} urmează după data la m_0 , deoarece $V=v$.

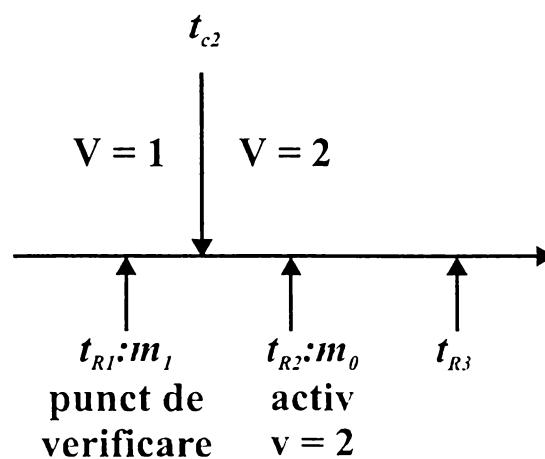


Fig. 4.12 Cazul 2: Pagină referită anterior

O rulare înapoi este efectuată prin îndepărtarea oricărei date ce a fost modificată după PV-ul anterior. Dacă pagina nu a fost încă referită după PV-ul anterior, atunci pagina este într-o stare obținută prin rularea înapoi și nu trebuie efectuată nici o operație (de exemplu cazul 1 din fig. 4.11. Dacă pagina a fost referită după PV-ul anterior, atunci există o pagină activă ce trebuie îndepărtată.

De exemplu, dacă o defecțiune apare la t_{R3} , în fig. 4.12, se dorește îndepărtarea lui m_0 și restuararea lui m_1 . Astfel pentru toate paginile cu $V=v$, valoarea v este decrementată și bitul l este inversat.

4.4.1 Recuperarea din defect prin *rulare înapoi* în sistemele multiprocesor

Recuperarea prin rulare înapoi reprezintă un subiect complex în multiprocesoarele ce comunică între ele deoarece orice recuperare prin *rulare înapoi* trebuie să țină cont de toate comunicațiile și de efectul lor după ultimul punct de verificare.

4.4.1.a Recuperarea prin *rulare înapoi* la SMP - MP

Memoriile cache utilizate în SMP - MP, atașate fiecărui procesor, utilizează o strategie de *scriere - înapoi* (*write - back*) pentru a reduce traficul la memoria principală. Această strategie de menținere a coerenței admite existența la un moment dat a mai multor copii ale aceleiași date în diferite memorii cache. Aceasta presupune utilizarea unor protocoale de menținere a coerenței în timpul modificării unei anumite copii. De exemplu orice procesor ce dorește să modifice o copie poate transmite un mesaj astfel încât toate procesoarele ce dețin copii trebuie în prealabil să și le invalideze. De asemenea, copia din memoria principală trebuie să fie actualizată în cel mai scurt timp posibil după ce procesorul modifică copia din cache. Utilizând tehnica *scrie - prin* (*write - through*), modificarea se obține instantaneu, în același timp cu modificarea copiei cache. Utilizând procedura *scrie - înapoi* (*write - back*), această actualizare este obținută când conținutul cache-ului este copiat înapoi datorită unei tentative eșuate de acces. Cele mai multe scheme utilizează această proprietate a procesoarelor cu cache-uri actualizate prin scriere înapoi, pentru a face o recuperare rapidă. Această optimizare nu se poate obține utilizând tehnica *scrie - prin* [Prad 96]. Prin urmare cele mai multe scheme de recuperare din defect pentru SMP - MP utilizează protocolul *scrie - înapoi* pentru menținerea coerenței.

Prezența cache-urilor implică faptul că nu numai un procesor poate avea copii ale aceleiași date, unde o modificare subsecventă poate să conducă la situația în care mai

multe procesoare să aibă copii diferite ale aceleași date, deci să producă incoerența memoriilor cache și principală. În consecință se utilizează schema ce etichetează datele în cache-uri pentru menținerea coerenței. Se utilizează [Prad96], [AFM90], [BP93], [HM87], [JF91] scheme bazate pe magistrală singulară și bazate pe director. Cele mai multe scheme bazate pe magistrală utilizează politicile *invalidare de scriere* și *actualizare de scriere*. În alte lucrări aceste protocoale sunt denumite *scrie – invalidează*, respectiv *scrie – actualizează*.

Reamintim că esența schemei *invalidare de scriere* presupune că ori de câte ori într-o locație de memorie trebuie scris, se trimite un semnal ce invalidează toate copiile preexistente ale aceleași date dispuse în alte cache-uri înainte de actualizarea copieii locale din cache.

A doua categorie de scheme bazate pe magistrală, cele care utilizează tehnica *actualizare de scriere*, este mai bună dacă scrierile sunt mai puțin frecvente [Prad 96], deoarece ele implică livrarea datelor de scris la toate cache-urile simultan ori de câte ori un acces de scriere la un bloc cache partajat este generat. Schemele bazate pe director utilizează un *tabel (directory)* pentru a păstra evidența stărilor fiecărui bloc de memorie. Diferențele specifice între scheme rezidă în modul în care tabelul este înmagazinat și actualizat.

Tehnica de creare și exploatare a PV-urilor utilizând memoriile cache, este potrivită pentru construcția unei mașini ce poate tolera defecte tranzitorii.

Utilizând această structură ca și o bază, se poate focaliza un efort suplimentar pe detectarea erorilor mai eficient decât obținerea unor capacități combinate de detectare și corectare ale lor.

Problemele ce apar în astfel de scheme se referă la situațiile în care procese multiple, rulând pe procesoare diferite, crează dependența de date între ele. Apare problema *propagării rulării înapoi (rollback propagation)* care constă în posibilitatea ca un procesor care rulează înapoi poate necesita rularea înapoi a altui proces și așa mai departe. O tehnică de evitare constă în îndepărtarea dependențelor de date prin menținerea unui singur PV global.

Acest PV poate fi adresat prin metode variate ce vor fi prezentate în continuare. În [Prad 96] se prezintă cazul unui SMP - MP, unde fiecare procesor are un cache, și o magistrală comună partajată. Consistența cache-urilor este asigurată printr-un protocol ce reclamă existența a trei linii adiționale cu următoarele funcții:

- *partajarea unui bloc pe magistrală;*

- *stabilirea unui punct de rulare înapoi;*
- *rularea înapoi până la punctul anterior.*

O altă abordare a fost propusă în cazul SMP - MP cu memorii cache private, unde memoria partajată este atașată procesoarelor prin intermediul unei RIN. Pentru a se împiedica propagarea rulării înapoi, se generează PV-uri ori de câte ori alte procesoare citesc o linie de cache modificată de la ultimul PV. Se utilizează termeni ca și: *identificatori de puncte de valorificare, stivă de recuperare, numărător de lungime k.* etc.

Un avantaj semnificativ al SMP - MP îl constituie faptul că îndepărtarea unui proces defectat de la un procesor defect spre unul bun, este o operațiune relativ ușoară. Când un procesor cade el poate fi îndepărtat din operarea întregului SMP, sarcina fiind preluată de celelalte.

Dacă defecțiunea procesorului este una tranzitorie, atunci el poate rula înapoi până la PV anterior sau până la începutul procesorului. O chestiune interesantă este aceea a procesorului pe care să se reexecute procesul, în sensul că pot exista și alte procesoare care să efectueze această operație. Uzual, procesul eronat poate fi reexecutat pe același procesor și o defecțiune repetată poate indica un defect hard și necesitatea de a deplasa procesul la alt procesor.

În SMP - MP, această situație conduce la activarea stării PV a procesorului defect într-un procesor diferit. Aceasta poate fi obținută ușor deoarece într-un SMP - MP, tot ceea ce trebuie făcut este să se încarce starea PV - ului din procesul defect într-un procesor corespunzător. Pe de altă parte, într-un SMP - TM, neexistând memorie globală partajată pentru comunicare, rularea înapoi poate constitui o procedură complexă.

4.4.1.b Recuperarea prin *rulare înapoi* în SMP - TM

Sistemele multiprocesoare cu memorie distribuită (SMP – TM) constau dintr-un set interconectat de procesoare cu memorii locale ce comunică între ele prin intermediul transferului de mesaj. Nu există o memorie globală partajată și atunci când un procesor se defectează, informația este comunicată la alte procesoare din sistem într-o perioadă finită de timp după apariția defectului.

În literatură [Prad 96], [Gros 99] etc, se descriu mai multe scheme pentru crearea și exploatarea PV-urilor în cadrul SMP - TM.

Un SMP - TM constă dintr-un set finit de N procese secvențiale P_1, P_2, \dots, P_N ce comunică și se sincronizează doar prin mesaje. Fiecare proces execută un algoritm local și rulează pe un procesor.

Apariția de acțiuni efectuate de algoritmi locali formează *evenimente*. Există trei tipuri de evenimente: *transmisii (send)*, *recepții (receive)* și *interne*. S-a definit următoarea relație între două evenimente: evenimentul a se produce în mod direct înaintea evenimentului b dacă și numai dacă:

1. a și b sunt evenimente în același proces și a apare înainte de b ; sau
2. a este transmisia mesajului m de către un proces și b este recepționarea lui m de către alt proces.

Închiderea tranzitivă a relației *se întâmplă direct înainte* este relația *se întâmplă înainte*. Dacă evenimentul a se produce înainte de b , b se produce după a .

Diagramele de timp reprezintă o metodă convenabilă pentru reprezentarea calculului distribuit.

În fig. 4.13 este reprezentată o diagramă de timp pentru un sistem de trei procese. Liniile orizontale sunt axele de timp ale proceselor, punctele sunt evenimente iar săgețile reprezintă mesaje de la procesul transmițător la cel receptor.

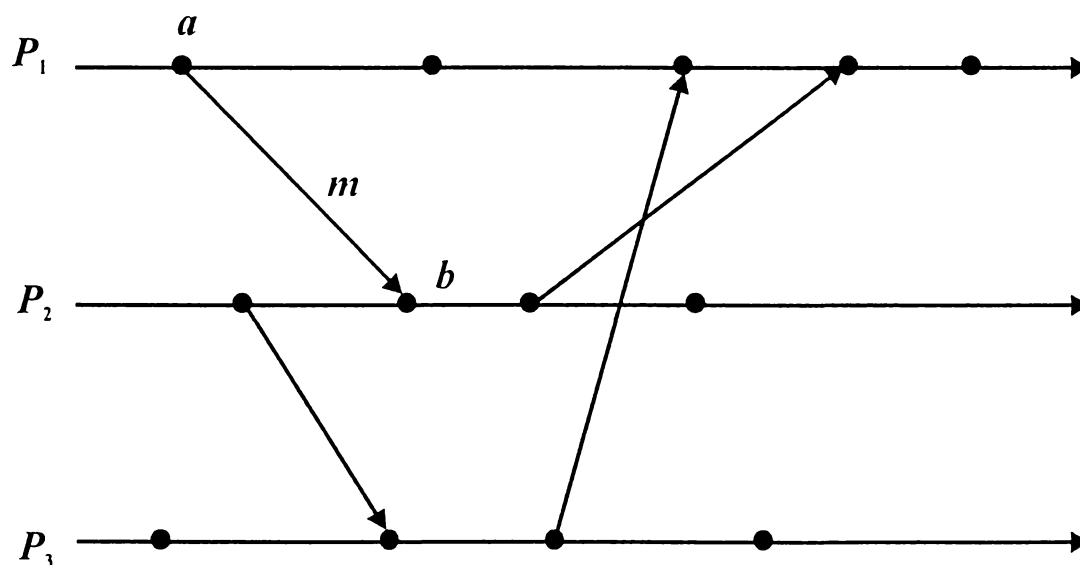


Fig. 4.13 O diagramă de timp pentru un calcul distribuit

O *stare globală* a unui SMP - TM este un set de stări locale, una pentru fiecare proces. *Starea locală* a unui proces P este definită de către starea sa inițială și de către secvența de evenimente ce a apărut la P . Starea canalelor corespunzând unei stări globale S este setul mesajelor transmise dar încă nerecepționate în S . Se poate

reprezenta o stare globală pe o diagramă de timp printr-o tăietură în diagramă. O *stare a sistemului consistentă* (*consistent system state*) este una în care fiecare mesaj ce a fost recepționat este de asemenea prezentat ca fiind transmis în starea transmițătorului, adică evenimentul de transmisie s-a produs înainte de evenimentul recepție. Informal, o stare globală (tăietură) este consistentă dacă nici o săgeată nu punctează de la dreapta spre stânga de-a lungul tăieturii. În fig. 4.14, C este o tăietură consistentă și C' este una *inconsistentă*.

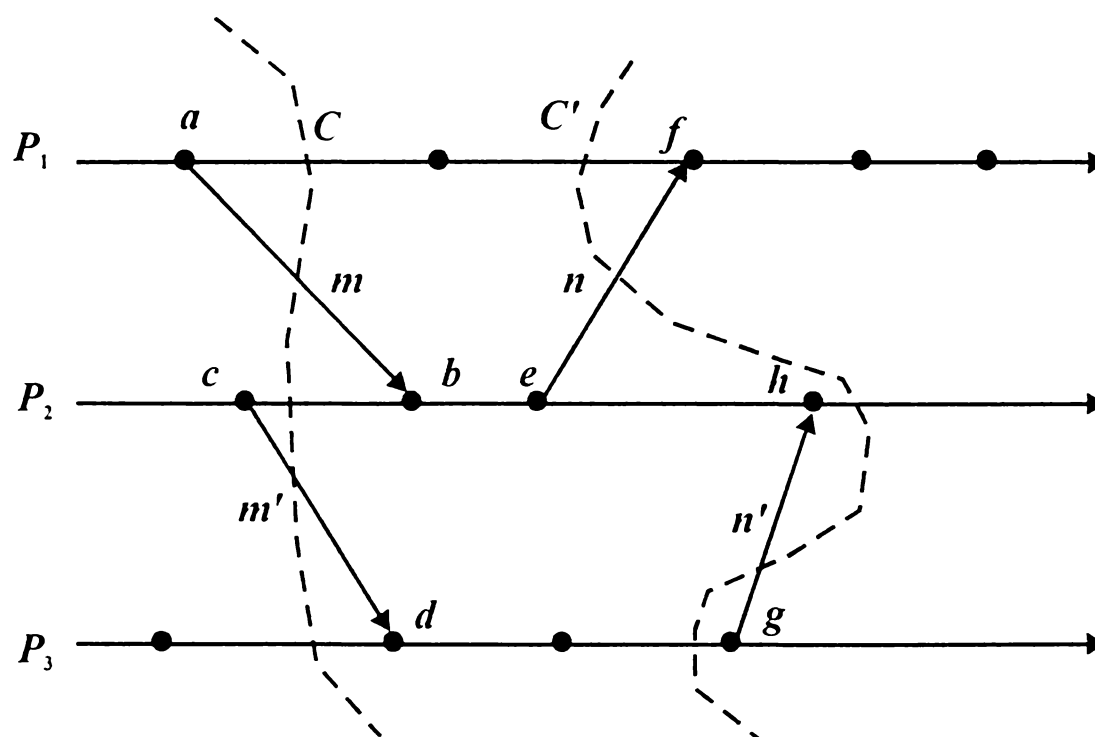


Fig. 4.14 Tăieturi consistente și inconsistente

Un punct de verificare (PV) reprezintă starea locală salvată a unui proces. Un set de PV - uri, unul per proces, este consistent dacă stările salvate formează o stare globală consistentă [Gros 99]. Orice PV global consistent poate fi utilizat pentru restaurare după apariția unei defecțiuni. O stare globală consistentă se numește *linie de recuperare* (*recovery line*). În scopul de-a minimiza timpul pierdut datorită unei defecțiuni, sistemul trebuie să fie rulat înapoi spre cea mai recentă linie de recuperare.

Cu ocazia recuperării pot apare următoarele patru probleme:

- *mesaje orfane* (*orphan messages*)
- *mesaje lipsă* (*missing messages*)

- *mesaje duplicat (duplicate messages)*
- *blocaj activ (livelock)*

Mesajele orfane pot apare în următoarea situație. Fie două procese P_i și P_j . P_i este restartat din punctul de verificare C_i^s și P_j din C_j^r , situație prezentată în figura 4.15.

Dacă procesul P_i nu regenerează mesajul m_1 după recuperare, datorită nondeterminismului lui P_i atunci mesajul m_1 , este un *mesaj orfan*. Procesul P_j depinzând de un astfel de mesaj orfan, devine un *proces orfan*. În mod normal, un mesaj m_1 trimis de P_i spre P_j este un mesaj orfan în corespondență cu perechea ordonată de puncte de verificare locale (C_i^s, C_j^r) dacă evenimentul său de recepție aparțin lui C_j^r în timp ce evenimentul de transmisie nu aparțin lui C_i^s .

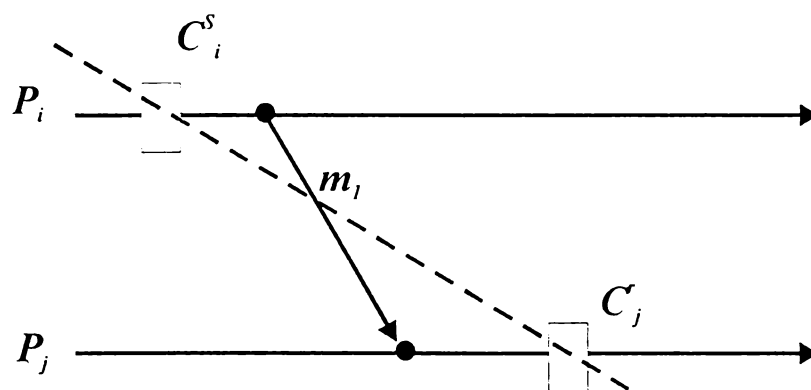


Fig. 4.15 Mesaj orfan

Pentru a preveni apariția mesajelor orfane au fost propuse câteva soluții [Prad 96], [Gros 99].

- Mesaje orfane pot fi evitate dacă procesele sunt deterministice și sunt executate pe procesoare ce se blochează la defect (fail - stop processors).
- Fiecare proces orfan trebuie să fie rulat înapoi la un PV înainte ca mesajul ce a cauzat această stare a procesului, să fie recepționat. Dacă nu se ia această precauție mai multe procese pot deveni orfane, care de asemenea trebuie rulate înapoi în timpul recuperării, producând *efectul domino-ului*.

Mesaje pierdute pot apare în următoarea situație. Fie două procese P_i și P_j . P_i este restartat din PV-ul notat C_i^s și P_j din C_j^r (fig. 4.16). Dacă procesul P_j nu execută un eveniment de recepție a mesajului m_2 , acest mesaj devine pierdut. Formal, un mesaj m_2 trimis de către P_i spre P_j este denumit *pierdut* în raport cu perechea ordonată (C_i^s, C_j^r) dacă evenimentul său de transmisie aparține lui C_i^s iar recepția nu aparține lui C_j^r .

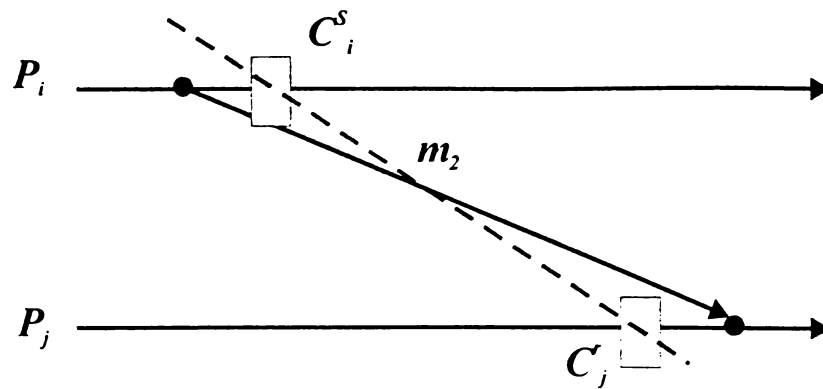


Fig. 4.16 Mesaj pierdut

Apariția mesajelor duplicate poate fi evidențiată în figura ce reprezintă cazul mesajelor orfane. Dacă evenimentul de transmisie este reexecutat în timpul calculului, mesajul m_1 este livrat de două ori spre P_j afectând apoi starea locală a lui P_j de două ori. Acest mesaj duplicat trebuie să fie manipulat cu atenție. Soluții posibile sunt: Suprimarea regenerării mesajelor duplicate, procesul receptor recunoaște mesajul ca fiind un duplicat etc.

Blocajul activ (livelock) în situația în care procesul de recuperare este asincron iar procesele nu se pot recupera simultan. Această recuperare asincronă poate introduce *blocaje active* adică situații în care o singură defectare poate cauza un număr infinit de rulări înapoi, împiedicând sistemul să progreseze. Se consideră următoarea situație prezentată în fig. 4.17. Procesul P_i este rulat înapoi spre C_i^s . Această rulare înapoi poate implica rularea înapoi a lui P_j spre C_j^r datorită mesajului m_1^1 .

Notăția m_i^j se referă la a j -a situație a transmiterii lui m_i . Să presupunem că P_j a transmis mesajul m_2 înainte de a fi rulat înapoi. Atunci m_2 poate fi recepționat în situația de rulat înapoi a lui P_i . Apoi P_i poate rula înapoi a doua oară pentru a restaura consistența. Dacă m_1 a fost recepționat în situația de rulat înapoi a lui P_j , aceasta presupune o a doua rulare înapoi a lui P_j spre C_j^r . Prin urmare, P_i și P_j pot fi forțate să ruleze înapoi pentru totdeauna, chiar dacă nu au apărut defectiuni suplimentare.

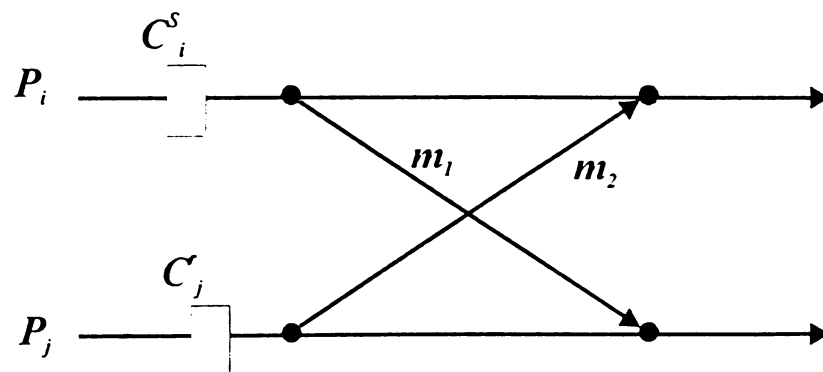


Fig.4.17 Situația blocajului activ

Se pot lua în considerare următoarele soluții pentru evitarea blocajului activ: ignorarea mesajelor de la un moment anterior al unui proces sau utilizarea unui protocol în două faze (procesele să ruleze înapoi la același moment).

Tehnicile de recuperare din defect prin rulare înapoi pentru SMP - TM pot fi clasificate în două clase principale [Gros 99]. Această taxonomie ia în considerare modalitățile prin care problemele prezentate anterior pot fi rezolvate.

Tehnicile existente de creare și gestionare a PV-urilor pot fi clasificate în următoarele categorii:

- *gestionarea sincronă a punctelor de verificare (synchronous checkpointing)* - procesele sincronizează gestionarea PV-urilor de-o astfel de manieră că întotdeauna este menținut un set consistent de PV-uri în sistem;
- *gestionarea asincronă a punctelor de verificare (asynchronous checkpointing)* - fiecare proces își ia PV-urile sale în mod independent existând riscul apariției efectului domino;
- *gestionarea quasi-sincronă a punctelor de verificare (quasi - synchronous checkpointing)* - procesele iau atât PV-urile induse cât și cele luate independent.

Tehnicile de recuperare din defect prin rulare înapoi, bazate pe înregistrare (log - based rollback recovery techniques), utilizează înregistrarea mesajelor suplimentar la gestionarea PV-urilor.

Aceste tehnici pot fi clasificate în trei categorii în concordanță cu modalitățile în care sunt evitate procese orfane: înregistrare pesimistă, optimistă și cauzală.

- în *înregistrarea pesimistă (pessimistic logging)* crearea proceselor orfane este evitată prin efectuarea unor acțiuni corespunzătoare în timpul unei execuții;
- în *înregistrarea optimistă (optimistic logging)* crearea unor procese orfane este evitată prin efectuarea de acțiuni în timpul recuperării;
- *înregistrarea cauzală (causal logging)* combină tehnicile anterioare.

4.4.2 Recuperarea din defect prin *rulare-înainte* în sistemele multiprocesor.

Caracteristica comună la toate schemele de *recuperare prin rulare-înainte* este aceea că, după detectarea unei defecțiuni, sistemul îndepărtează starea eronată curentă și determină starea corectă fără nici o pierdere în calcul. Există două metode diferite de abordare: una se bazează pe redundanță hardware, alta pe cea software. Abordările utilizând redundanța hardware pot fi apoi clasificate în două categorii :

- abordări utilizând redundanța statică;
- abordări utilizând redundanța dinamică.

4.4.2.a Abordări utilizând redundanța statică

Tehnicile bazate pe redundanța statică sunt foarte similare în principiu: toate utilizează acest tip de redundanță pentru mascarea defectelor.

- ***Redundanța cu mascare activă***

Se utilizează un nivel adecvat de replicare pentru a tolera defecțiunile, utilizând sisteme *TMR (triple-modular-redundant)* ce maschează un defect singular fără a pierde din performanțe.

- ***Redundanța activă utilizând module cu blocare pe defect***

Module multiple ale fiecărui procesor execută activ fiecare proces. Fiecare proces este presupus ca fiind cu *blocare pe defect (fail-stop)*. Astfel, dacă un procesor se defectează, el se oprește și celelalte continuă execuția task-ului fără a se produce degradări în performanțe.

În literatură [Prad 96], se citează sistemul Stratus în care fiecare subsistem este duplicat, formând o pereche iar una din replici este identificată ca fiind de rezervă. Fiecare subsistem și rezerva sa sunt cu autotestare prin duplicare. Circuitele sunt astfel replicate de patru ori. Toate copiile hardware sunt strâns sincronizate.

Când se detectează un defect într-un subsistem prin mecanismul propriu de autotestare, se deconectează subsistemul iar rezerva sa pornește fără nici o întrerupere sau rulare înapoi.

- ***Redundanța activă utilizând autodiagnoze***

Identificarea procesorului defect se face prin utilizarea task-urilor de autodiagnoză în locul mecanismelor de autoverificare.

Mecanismul de *duplicare reconfigurabilă* este descris în [Prad 96]. În duplicarea reconfigurărilor, procesul este replicat pe două procesoare. Se compară apoi în mod continuu ieșirile, orice nepotrivire semnaland o defecțiune la cel puțin un procesor din

pereche. După ce defectul a fost detectat, fiecare procesor rulează task-uri de auto-diagnoză pentru a localiza defectul. Odată ce procesorul defect este identificat, ieșirea procesorului liber de defect poate fi acceptată ca fiind corectă.

Utilizarea task-urilor de autodiagnoză în locul autoverificării concurente conduce la un calcul suplimentar pentru determinarea procesorului defectat.

4.4.2.b Abordări utilizând redundanța dinamică

Scheme de recuperare prin *rulare-înainte* bazate pe redundanță dinamică și pe crearea și gestionarea de PV-uri, prin care se încearcă rularea înapoi chiar și în prezența defecțiunilor, au fost prezentate în [Prad 96].

Pentru a ilustra principiul de bază utilizat în recuperarea prin *rulare-înainte* se consideră un sistem duplex ce detectează defecte prin crearea și verificarea PV-urilor din două module din sistem, în mod periodic, și apoi comparând stările lor. Când este detectată o defecțiune, o schemă cu *rulare-înapoi* ar încerca să restabilească starea celor două module din PV-ul anterior. Schemele cu *rulare-înainte* încearcă să determine care din cele două module procesoare (PM) este liber de eroare.

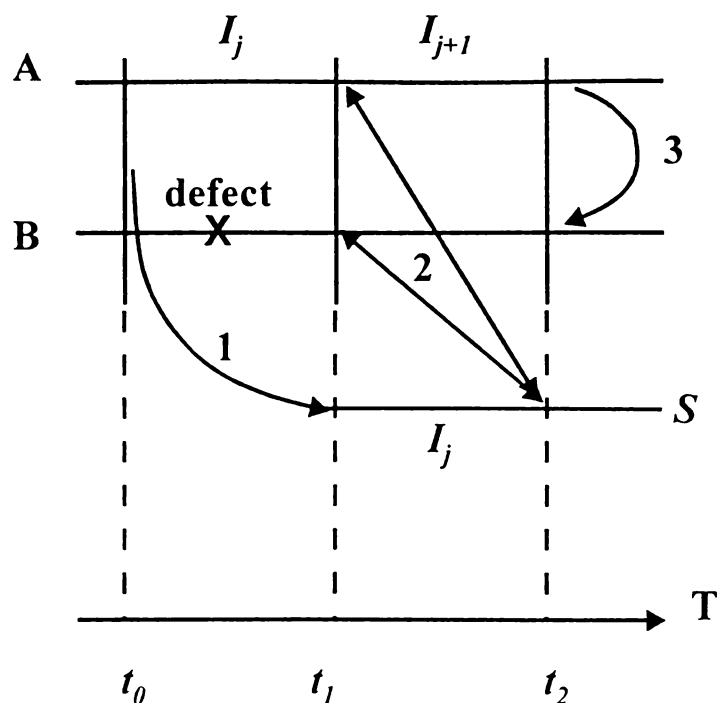
În literatură [Prad96], [LFA90], [LFA91], [PV94a], [PV94] se prezintă mai multe *scheme de creare și gestionare a punctelor de verificare prin rulare-înainte (Roll-Forward Checkpointing Scheme) - RFCS*.

O schemă RFCS este un mecanism pentru identificarea modulului de procesare defect în cadrul unui sistem duplex prin *relansarea executării programului (computation retrying)* pe un *modul partener de rezervă (spare module)*.

Odată ce PM defect este identificat, starea acestuia este făcută consistentă în raport cu starea PM-ului liber de eroare din cadrul sistemului duplex.

Calculul se reîncearcă în modulul de rezervă, concurent cu execuția în continuare a task-ului pe cele două module din cadrul sistemului duplex.

În fig. 4.18 se prezintă executarea a două copii ale unui task, denumite A și B. În situația în care B defectează în intervalul de checkpointare I_j se produc următoarele acțiuni. Semnăturile lui A și B se vor compara la sfârșitul intervalului I_j la timpul t_1 . Această comparare la t_1 va activa relansarea în manieră concurentă a intervalului de verificare I_j , pe modulul de rezervă, după cum urmează. În timpul acestei relansări, atât A cât și B continuă să avanseze spre următorul interval de checkpoint I_{j-1} .



- 1 : Starea de copiere în modulul de rezervă
- 2 : Starea de comparare a modulului de rezervă cu starea lui A și B
- 3 : Starea de copiere din A și B

Fig. 4.18 Relansarea concurrentă în cadrul unei scheme RFCS

La începutul relansării concurrente, PV-ul anterior (prelevat la t_0) este încărcat într-un modul de rezervă. Intervalul de verificare I_j , în care a apărut defectul, este relansat în modulul de rezervă. În mod concurrent, A și B continuă execuția următorului interval I_{j-1} . Cum A și B nu au fost în aceeași stare la t_1 , este puțin probabil că vor fi în aceeași stare la sfârșitul lui I_{j-1} , care a fost prelevată la t_2 .

După ce modulul S completează intervalul I_j la timpul t_2 , starea lui S este comparată cu starea lui A și B la t_1 .

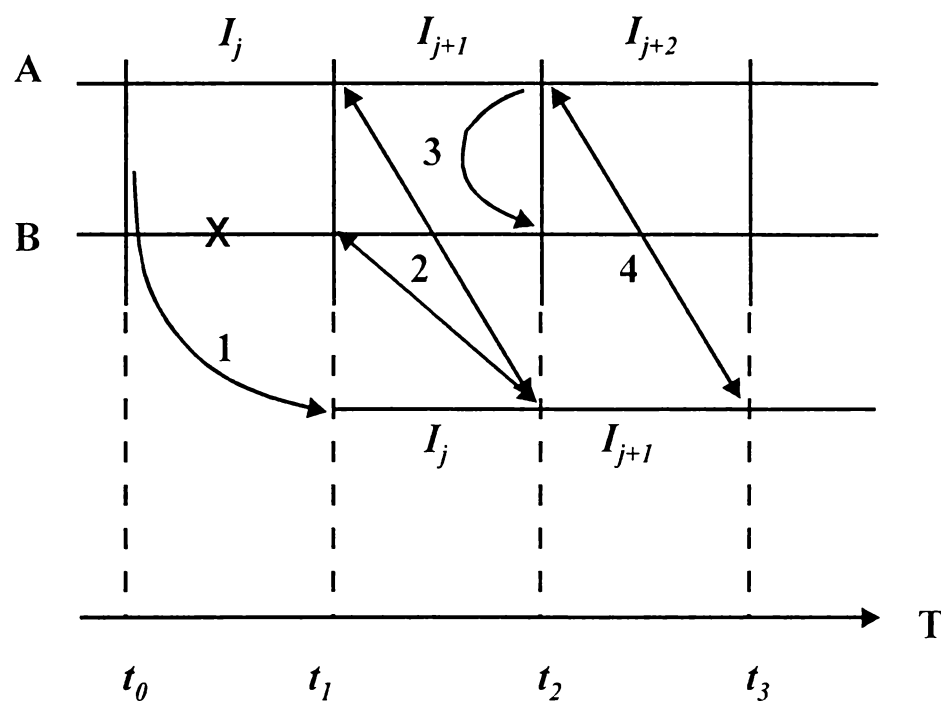
Dacă A și S nu se defectează în I_j , atunci starea lui S la t_2 se va potrivi cu starea lui A la t_1 . Când apare această similitudine, s-a presupus că A este liber de eroare la t_1 .

Acum, starea lui B este făcută identică cu starea lui A la t_2 . Astfel, dacă A nu defectează în intervalul I_{j-1} , atunci A și B vor fi ambele în starea corectă la momentul t_2 .

Pentru a determina dacă A a defectat în I_{j-1} , modulul de rezervă trece la executarea intervalului I_{j-1} . Acest caz este prezentat în fig. 4.19

După ce modulul de rezervă parcurge acest interval, starea sa este comparată cu starea lui A la sfârșitul lui I_{j-1} .

Dacă cele două stări se potrivesc, atunci se presupune că a fost liber de defect în I_{j+1} ; astfel sistemul duplex este rulat înapoi. O analiză detaliată a funcționării schemelor RFCS este prezentată în [Prad 96]. În toate schemele RFCS se presupune că există înglobate capacități de detectare a defectului în modulele procesoare.



4. Starea de comparare a modulului A și a celui de rezervă

Fig. 4.19 Tehnica de relansare concurrentă în schemele RFCS.

Schemele RFCS evită rularea înapoi chiar și în prezența unui singur defect. Defectul este astfel tolerat fără o penalizare în performanță.

Relativ recent (1994) s-au produs variante ale schemei clasice RFCS. Presupunând că fiecare modul are un anumit număr de capacități de detectare a defectului conținute în module - verificări la paritate, capacități de detectare a situațiilor de excepție, etc. - se pot conceptualiza patru scenarii diferite, prezentate în TABEL VIII.

TABEL VIII. Patru scheme diferite cu *rulare înainte*

Strategie de recuperare	Resurse utilizate	
	Cu modul de rezervă	Fără modul de rezervă
Optimistă	<i>Rulare înainte I</i>	<i>Rulare înapoi I</i>
Pesimistă	<i>Rulare înainte II</i>	<i>Rulare înapoi II</i>

În cadrul unei strategii de recuperare *optimistă*, se acordă încredere în cel mai mare grad capacităților de detectare conținute în module.

De exemplu se consideră că un defect a fost detectat în modulul B de către un astfel de mecanism intern, situație prezentată în fig. 4.20.

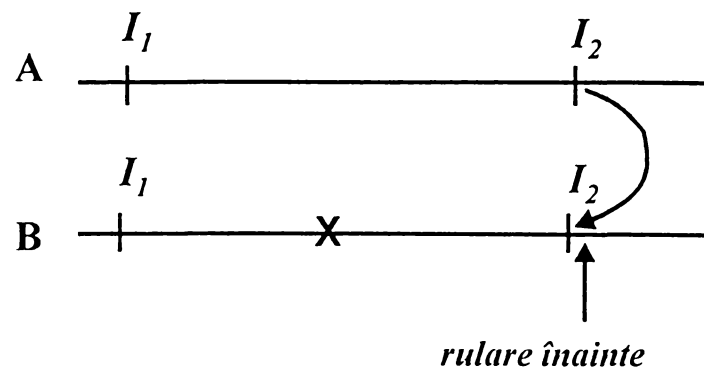


Fig. 4.20 Schemă optimistă cu sau fără modulul de rezervă

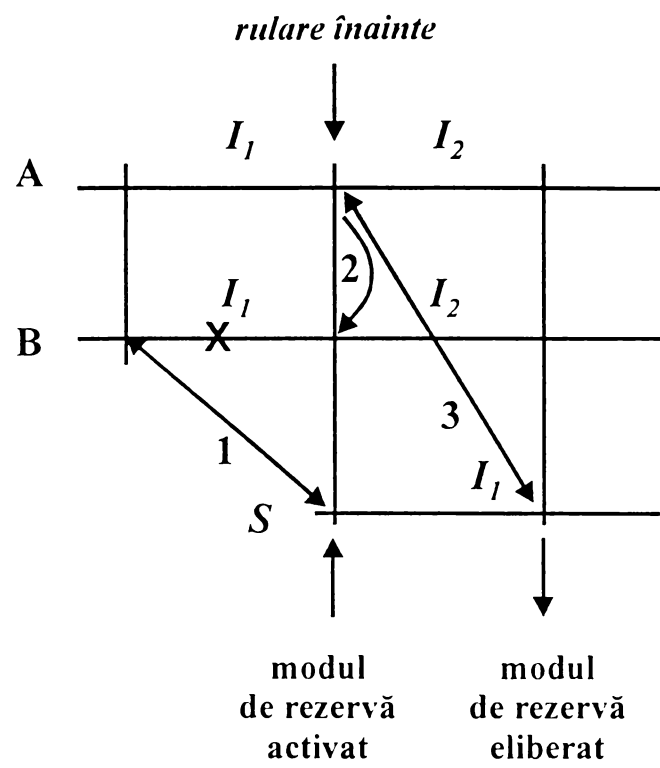
La sfârșitul intervalului de checkpointare I_1 , starea lui A poate fi copiată în B și procesul ar putea rula înainte. Aceasta nu va solicita utilizarea unui modul de rezervă, chiar dacă el este disponibil. Această schemă este denumită *Rulare - înainte I*.

În schema *pesimistă*, modulul de rezervă este activat în mod automat și se execută relansarea concurentă și recuperarea ca și în fig. 4.21 a. Această schemă este denumită *Rulare - înainte II*. Se observă că deoarece modulul B a fost tot timpul suspectat ca fiind defect, doar în situația în care A a prezentat o defecțiune ce a scăpat capabilității de detectare conținută în modul (a lui A în intervalul I_1) s-a acționat de-o manieră mai conservatoare. În acest caz, s-ar putea ca A să se fi defectat în intervalul I_1 (ceea ce ar corespunde la o dublă eroare în timpul lui I_1) și procesul să ruleze înapoi când modulul de rezervă nu are aceeași stare cu a după reexecutarea lui I_1 .

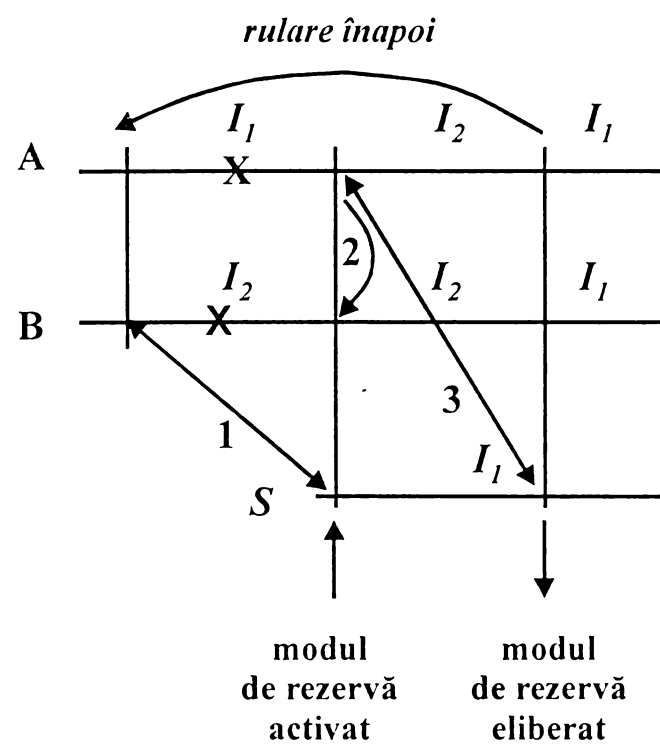
Oricum, în ipoteza că nu sunt disponibile module de rezervă, schema *optimistă* (*Rulare-înapoi I*) va rula înainte, ca și în fig. 4.20, în timp ce schema *pesimistă* (*Rulare înapoi II*) va rula înapoi atât A cât și B pentru a reexecuta I_1 , ca și în cazul schemelor standard cu rulare înapoi.

O analiză prezentată de Pradhan și Vaidya în [Prad 96] arată că un câștig semnificativ în performanță poate fi obținut cu un sacrificiu de talie redusă în ceea ce privește siguranța în funcționare.

În fig. 4.22 se prezintă comparativ trei sisteme în spațiul siguranță în funcționare - performanță. Se observă că este convenabil un sistem ale cărui caracteristici sunt descrise de curba 1. Acest sistem care permite o mică reducere în siguranța în funcționare va avea însă un semnificativ câștig în performanță. Schemele optimiste descrise oferă potențial pentru scheme ca și cele descrise în fig. 4.23 și fig. 4.24.



- a. Schemă pesimistă cu modulul de rezervă rulând înainte, cu toate defectele singulare



- b. Schemă pesimistă cu modulul de rezervă rulând înapoi, cu defecte duble

Fig. 4.21 Scheme pesimiste.

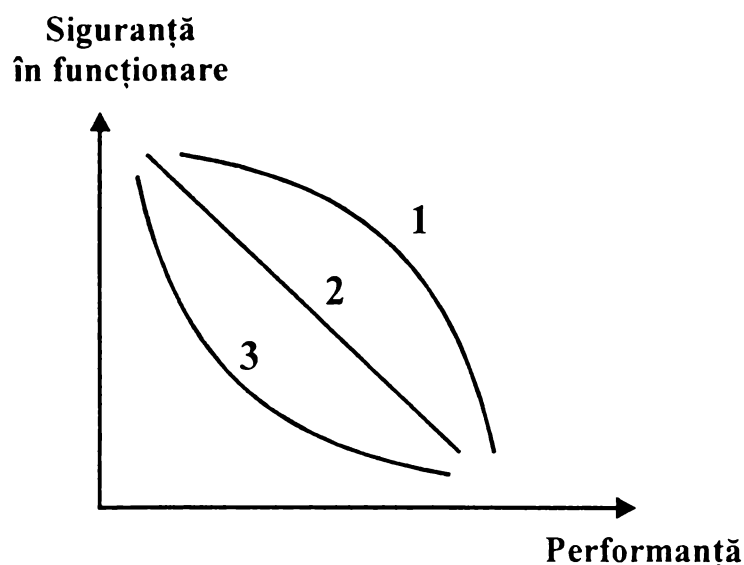


Fig. 4.22 Trei scheme diferite cu *rulare înainte*

În fig. 4.24 se ilustrează o comparație tipică de performanță indicată de o formulare analitică. Aici $\lambda = 10^{-6}$ reprezintă rata asumată de defectare și $c = 0.8$ este gradul de acoperire al detectării defectului. În general, timpul mediu de tratare și de reușire din eroare când a apărut o defecțiune, este mai mic pentru schemele cu *rulare-înainte* în ambele scenarii : optimistă și pesimistă. Acest parametru reprezintă o metrică importantă deoarece fără nici o defecțiune, toate sistemele operează similar. Prin urmare este important să se cunoască comportamentul sistemului după ce acesta a experimentat o situație în care a fost nevoit să se recupereze din defect.

În ceea ce privește siguranța în funcționare, se pot accepta valori diferite pentru capabilitățile de detecție conținute în module. Parametrul c poate lua valori de la 0 - nu există nici o capabilitate de detecție - până la 1 - capabilitate perfectă de detecție. Se observă că atunci când $c = 0$, *schemele Rulează - înainte I = Rulează - înainte II și Rulează - înapoi I = Rulează - înapoi II* din punct de vedere al nesiguranței în funcționare. Cu alte cuvinte când nu există capabilități de detecție înglobate, schema optimistă și corespondenta sa pesimistă au siguranțe în funcționare identice. Deoarece nu există detectare încapsulată nu există nici o modalitate de-a identifica modulul defect fără a efectua compararea între module.

Pe măsură ce c crește apare o diferență între ambele abordări ca și în fig. 4.24. Aici rata de defectare este presupusă ca fiind 10^{-3} și $n = 10$ reprezintă intervalul de creare și gestionare a PV-urilor. Schema pesimistă suferă mai puțin de degradarea siguranței în funcționare de cât cea optimistă.

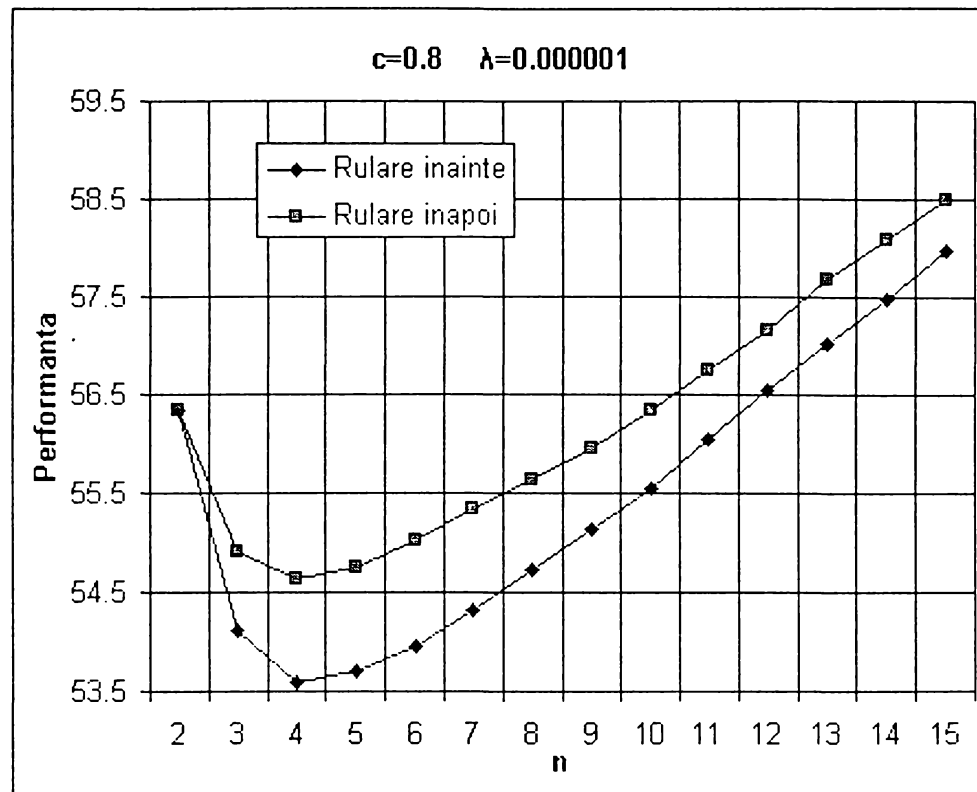


Fig. 4.23 Comparația în performanță dintre schemele optimiste și pesimiste
Timpul mediu de completare în prezența unui defect

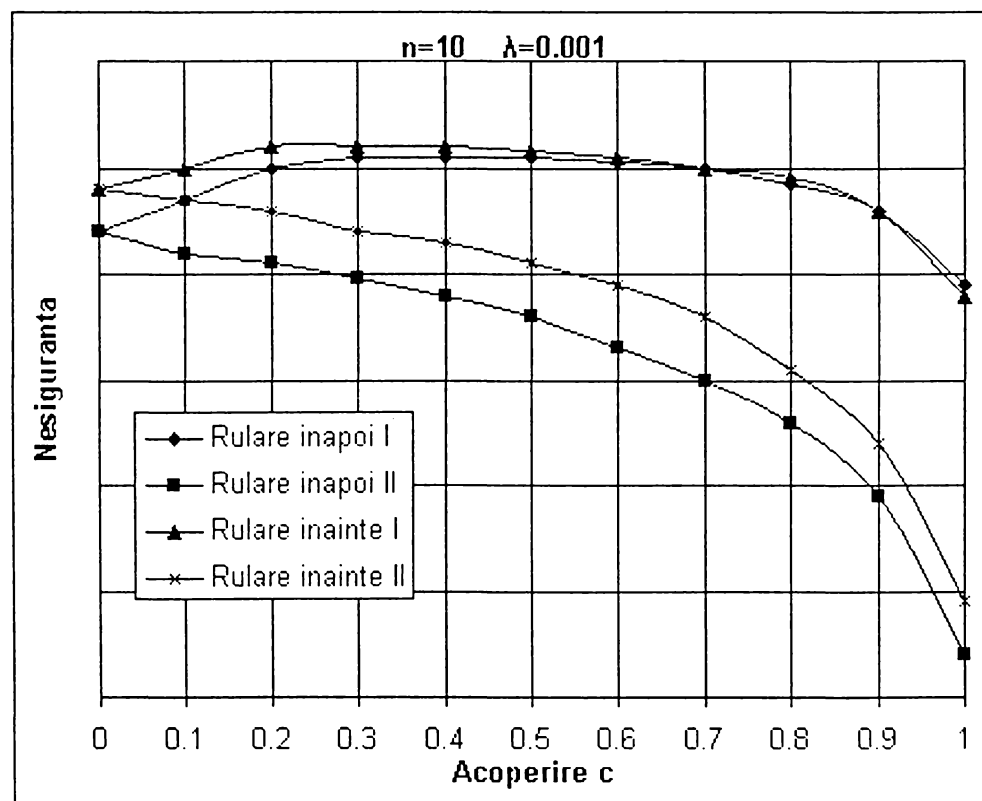


Fig. 4.24 Comparații ale siguranței în funcționare între schemele optimiste și pesimiste

Dar, în realitate, diferența este mică în termeni de nesiguranță în funcționare. Presupunând însă că se poate tolera această diferență, se preferă schema optimistă deoarece are performanțe superioare în situația apariției unui defect. Pe de altă parte, dacă siguranța în funcționare este pe primul loc atunci se pot alege scheme pesimiste care sunt mai conservatoare. Trebuie notat că efectul utilizării modulelor de rezervă poate să reducă siguranța în funcționare, dată fiind o strategie fixată de recuperare, dar pot ameliora performanțele.

În fig. 4.25 se prezintă comportamentul schemelor cu *rulare - înapoi* în cazul apariției unui defect.

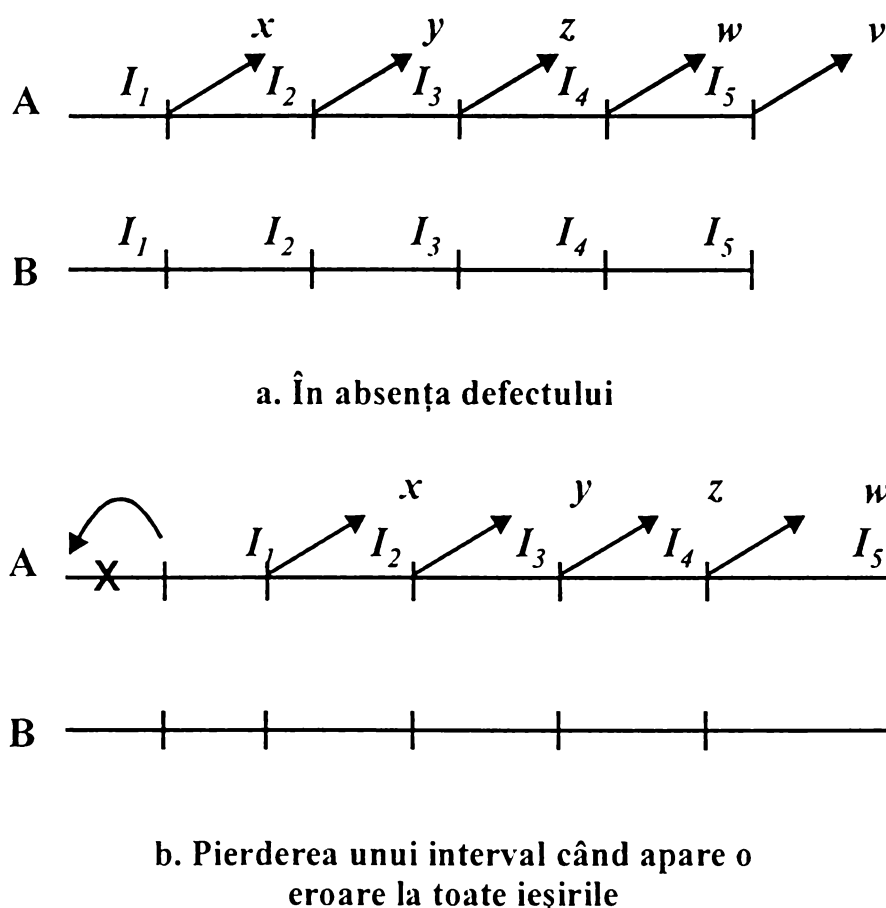


Fig. 4.25 Întârzierea permanentă la ieșirile schemelor cu *rulare-înapoi* în cazul apariției unui defect.

Unul dintre avantajele importante ale schemelor cu *rulare - înainte* este degradarea minimă în performanțele de I/E. Se consideră schema standard cu *rulare - înapoi*. În mod normal, ieșirile se vor poziționa la sfârșitul unui interval după ce o comparație indică că nu a apărut nici o eroare în timpul intervalului. Se consideră un proces care are ieșirea x la sfârșitul lui I_1 , y la sfârșitul lui I_2 , z la sfârșitul lui I_3 și w la sfârșitul lui I_4 .

Dacă apare un defect în I_1 , intervalul I_2 va fi utilizat pentru reexecutarea lui I_1 . Astfel, toate ieșirile după I_1 vor experimenta o întârziere de un interval de checkpointare, caz prezentat în fig.4.26.

În cazul schemei cu *rulare - înainte*, în cazul aceluiași scenariu, ieșirile vor fi întârziate temporar ca și în fig. 4.25. Ieșirile x și y sunt singurele întârziate și toate celelalte ieșiri vor apărea la intervalul de timp planificat, în mod regulat. Mai exact, ieșirea x va fi întârziată două intervale, în timp ce y doar unul. Toate ieșirile subsecvente se vor poziționa fără nici o întârziere, dacă nu există defect în sistem.

De exemplu ieșirile z , w și v vor apărea la sfârșitul lui I_3 , I_4 și I_5 ca și în fig. 4.26.

Acesta este același timp ce se așteaptă să apară când nu sunt defecte, ca și în fig. 4.25 a.

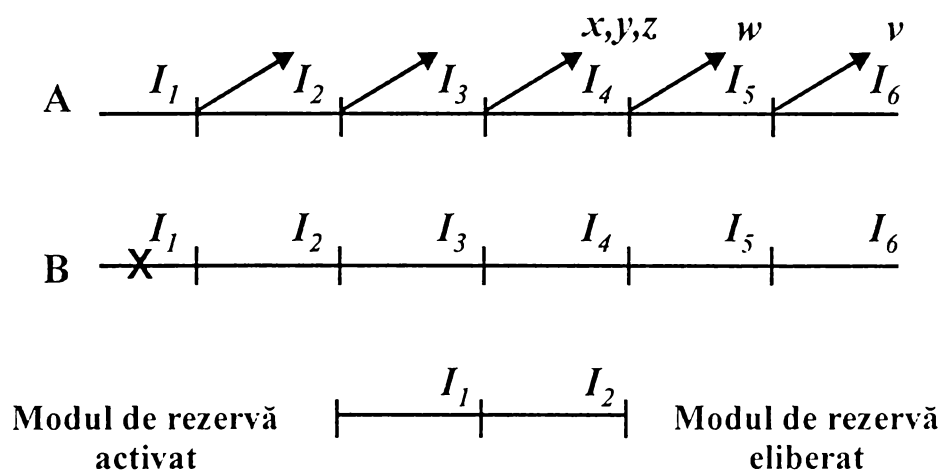


Fig. 4.26 Întârziere temporară în ieșirile schemelor cu *rulare - înainte* în situația unui defect

4.5 RECONFIGURAREA ÎN SISTEMELE MULTIPROCESOR

Sistemele multiprocesor tolerante la defect trebuie să aibă capacitatea de *reconfigurare*.

Într-o primă abordare, scopul este obținerea unui nivel controlabil de performanță degradată după ce defectele au fost identificate și izolate în tendința lor de propagare. Aceasta implică posibilitatea de a identifica un subset de procesoare active, conectate în aceeași topologie ca și cea originală, dar de o talie mai redusă. De exemplu, dacă topologia originală fără defecte este o *plasă* $N \times N$, după apariția defectului se încearcă identificarea unui set de $M \times M$ procesoare, unde $M < N$. În astfel de scheme, aplicațiile de calcul paralel care rulează pe aceste multiprocesoare utilizează proprietatea topologiei multiprocesorului și pot evolua pe orice măsură a topologiei. Cea mai simplă abordare a reconfigurării software în *hipercuburi*, identifică subcubul liber de erori de cea mai mare talie și rulează aplicațiile pe el. Dificultatea în această abordare este că degradarea în performanță poate fi până la 100% pentru primul defect în *hipercub* deoarece se va utiliza un subcub la jumătate din dimensiunea inițială.

A doua abordare în scopul obținerii toleranței la defect în RIN statice este aceea prin care interconexiunea este privită ca și o *cale* de a conecta un mare număr de procesoare ce pot comunica prin transmiterea de mesaje de-a lungul legăturilor. Într-o astfel de situație, toleranța la defect poate fi ușor obținută printr-o *degradare grațioasă* (*graceful degradation*), prin *redistribuirea sarcinii de calcul* (*computational load*) în jurul procesoarelor defecte și *rerutarea* mesajelor. În această ipostază este necesară investigarea algoritmilor de rutare ce vor transmite mesaje ce trebuie să fie rutate spre procesoarele active, evitând procesoarele defecte. Această abordare este adecvată pentru scenariile unde aplicațiile ce se execută pe aceste multiprocesoare nu necesită o topologie specifică de interconectare dar reclamă ca procesoarele să fie conectate.

A treia abordare necesită proiectarea unei *topologii augmentate* (*augmented topology*) prin adăugarea de procesoare suplimentare și conexiuni adiționale structurii originale, astfel că atunci când apar defecte în procesoare sau în conexiuni, topologia originală să fie identificată. În această abordare, reconfigurarea este obținută prin substituirea elementelor defecte cu cele suplimentare pentru a se menține nivelul de performanță nemodificat.

În continuare se vor detalia aceste abordări pentru mai multe topologii de multiprocesoare. Se va accentua asupra reconfigurării rețelelor de interconectare, acestea având un rol major în crearea statutului de sistem tolerant la defect.

4.5.1 Reconfigurarea în SMP – MP bazate pe *magistrală*

Magistrala este cea mai comună RIN în care toate procesoarele și toate modulele de memorie sunt conectate la o resursă comună. SMP – MP bazate pe magistrală sunt cel mai vulnerabile la defecte.

SMP – MP bazate pe magistrală, tolerante la defect, au avantajul că perechile de module procesor și module de memorie pot fi extrase ușor din sistem atunci când au fost identificate ca fiind defecte. Defectele sunt gestionate prin utilizarea *magistralelor redundante*. Într-un sistem redundant toate procesoarele și toate memoriile sunt conectate la un număr de magistrale ca și în fig.3.4. Funcționarea SMP – MP cu *magistrală multiplă* a fost descrisă *în extenso* în § 3.1.3.

Când o magistrală se defectează, sistemul operează ca și un sistem cu $B-1$ magistrale, din cele B existente. Arbitrii trebuie să fie reconfigurați pentru a opera cu $B-1$ magistrale.

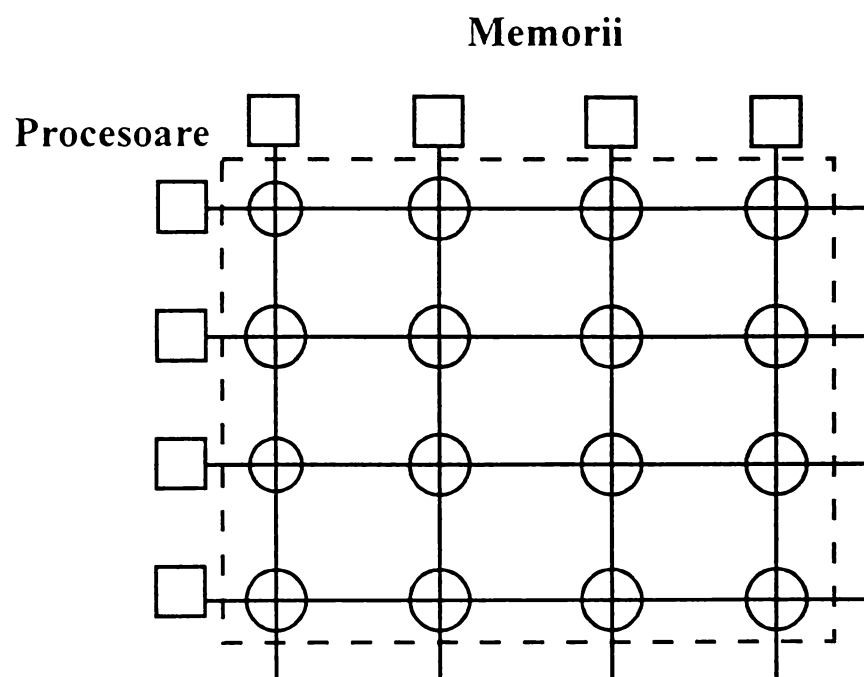
În literatură [Prad 96] se prezintă mai multe metode de reconfigurare a acestor sisteme. În principiu toate permit interconectarea a N procesoare utilizând un număr fixat de procesoare / magistrală, cu fiecare procesor conectat la două magistrale. Aceste topologii permit un anumit grad de toleranță la defect.

4.5.2 Reconfigurarea rețelelor de tip *crossbar*

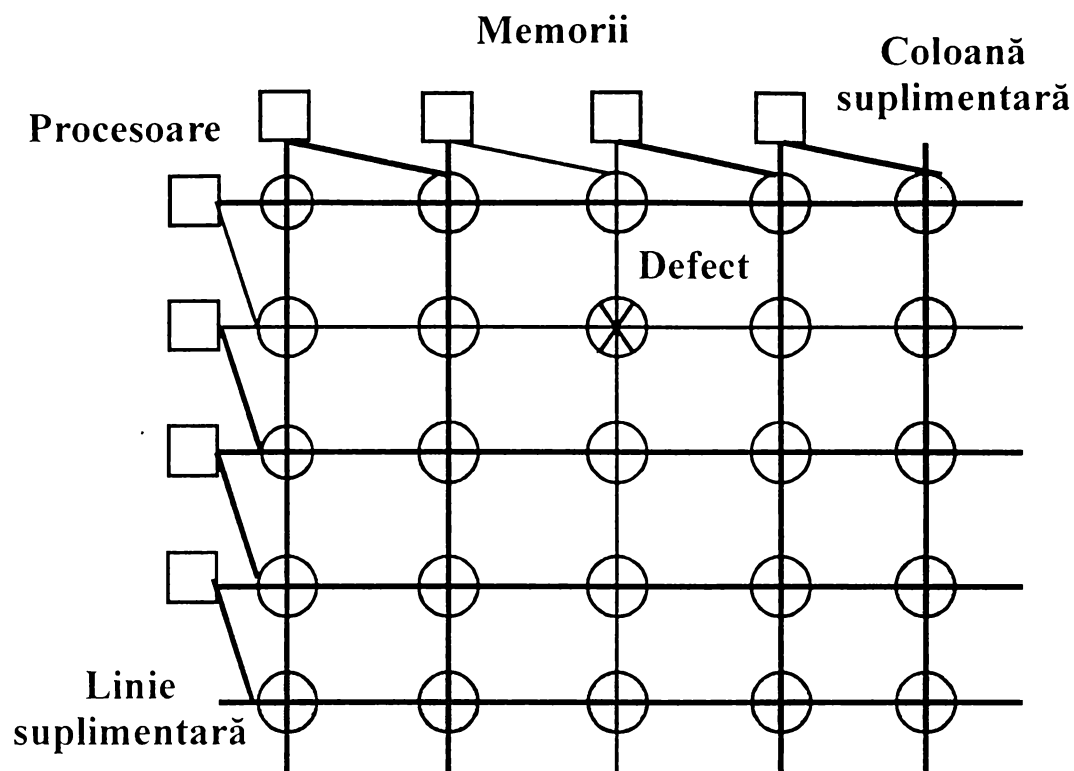
RIN *crossbar* a fost prezentată în §2.4.4. Reamintim că o altfel de conexiune se caracterizează prin prezența a N^2 comutatoare ce conectează N surse cu N destinații. Este permisă conectarea simultană între toate perechile procesor-memorie. Logica de arbitraj este încorporată în comutator. Un mod echivalent de-a construi un *crossbar* $N \times N$ este acela prin care este privit ca și o colecție de N multiplexoare $N \times 1$, cu cele N intrări legate la fiecare din multiplexoare. Fiecare din cele $N \times 1$ MUX-uri poate fi implementat ca și un arbore de 2×1 MUX-uri având $N - 1$ astfel de comutatoare.

O metodă de-a obține o mai mare fiabilitate la astfel de rețele este *replicarea* lor în întregime, adică păstrarea a R copii a fiecărei din cele N MUX-uri $N \times 1$.

O altă metodă utilizează o *linie* și o *coloană* suplimentare astfel că se utilizează un *crossbar* $(N + 1) \times (N + 1)$ cu fiecare port conectat la două linii sau coloane de comutatoare. La defectarea unui comutator, se deconectează întreaga linie și coloană de comutatoare și se activează coloana și linia suplimentară. Schema este prezentată în fig. 4.27 cu linie îngroșată fiind prezentată reconfigurarea *crossbar-ului* în caz de defect a unui comutator.



a) RIN *crossbar* inițial

b) RIN *crossbar* redundantFig. 4.27 RIN *crossbar* tolerantă la defect.

4.5.3 Reconfigurarea rețelelor *hipercub*

Multicalculatoarele de tip *hipercub*, apărute relativ recent, oferă o abordare spre calculatoarele masiv paralele, relativ ieftine și fezabile. Ele conectează un număr mare de procesoare $P = 2^d$, ce au memorie locală, utilizând conexiuni directe în conformitate cu topologia de interconectare de tip *d-cub* (*binary d-cube*).

În literatură [Prad 96], [Ban90] se prezintă mai multe soluții de rezolvare a problemei toleranței la defect pentru acest tip de multicalculatoare.

O soluție constă în reconfigurarea *hipercubului* d -dimensional prin augmentarea cu un port adițional în dimensiunea $d + 1$. Extra portul este utilizat pentru a conecta un *procesor de extensie* (*spare processor*). Soluția este prezentată în fig. 4.28a. Procesorul de extensie se poate conecta la portul adițional al fiecărui procesor din *hipercub*, astfel că atunci când un procesor se defectează, se poate trece peste acesta activându-se conexiunile corespunzătoare spre vecinii procesului defect din topologie. Această aparent simplă abordare implică un grad $N = 2^d$ pentru procesorul de extensie.

Soluția problemei constă în utilizarea unui set de comutatoare *crossbar* într-o încapsulare VLSI ce implementează această conectare. Se consideră un *hipercub* cu 16 procesoare astfel că fiecare procesor să aibă în mod normal patru porturi.

Se adaugă un extra port pentru asigurarea toleranței la defect. Cele 16 procesoare sunt divizate în două grupe de câte 8 procesoare. Cele opt conexiuni ce ies din fiecare procesor al subcubului sunt conectate la un comutator *crossbar* 8×5 . Fiecare subcub de 8 procesoare are un comutator *crossbar* corespunzător (în acest caz, două comutatoare). Ieșirile celor două comutatoare *crossbar* sunt conectate împreună cu cele cinci porturi ale procesorului de extensie. Când un procesor dintr-un subcub se defectează (de ex. procesorul 4), procesorul de extensie activează cele patru conexiuni ale subcubului defect prin comutatorul *crossbar* corespunzător (de ex. procesoarele 0, 5, 6) și o conexiune a altor comutatoare *crossbar* (de ex. procesor 12) în alte subcuburi la care procesorul defect a fost conectat, ca și în fig. 4.28b.

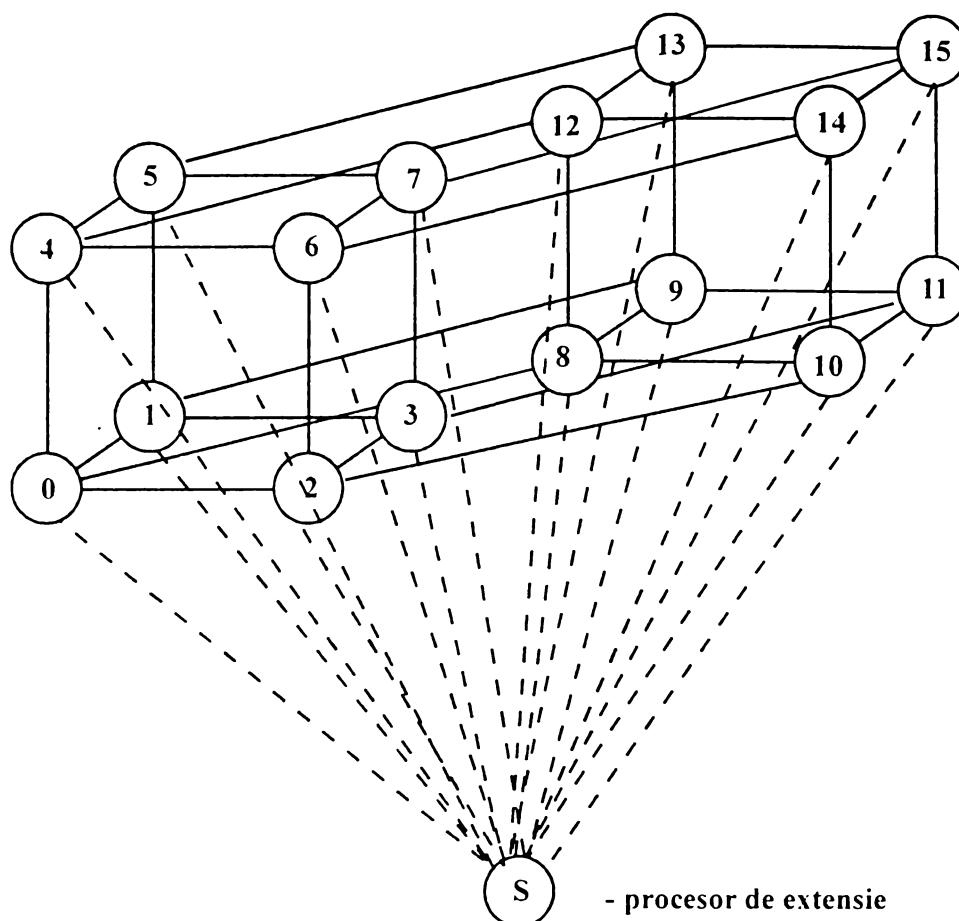
Schema prezentată poate fi extinsă pentru utilizarea mai multor procesoare de extensie unde un astfel de procesor a fost afectat la un subcub de procesoare, în loc de întreg *hipercubul*.

Alte metode hardware pentru reconfigurarea unui *hipercub* sunt descrise în [Prad 96]. În prima schemă hardware pentru reconfigurare este presupus că procesoarele de extensie pot fi atașate la procesoare specifice ale *hipercubului*. În această abordare, două tipuri de noduri (noduri P și noduri S) sunt proiectate în *hipercub*. În fig. 4.29a se prezintă un *hipercub* cu 16 procesoare unde patru din noduri (0, 7, 9 și 14) sunt *noduri* S, restul fiind *noduri* P. Un nod P constă dintr-un procesor de calcul CPU (computation processor) conectat printr-o magistrală internă la o memorie locală, și o logică de rutare a mesajului constând dintr-o unitate DMA și un comutator *crossbar* $(d + 1) \times (d + 1)$ pentru *hipercubul* cu 2^d procesoare. Un nod S constă din două copii ale CPU și ale memoriei locale conectate la două magistrale interne. Unitatea DMA și logica de rutare a mesajului sunt partajate între cele două unități de procesare, dintre care una este activă în condiții normale iar alta este de extensie. În caz de defectare a oricărui element de procesare (CPU sau memoria locală), fie prin nodul S sau prin nodul P, cuplul procesor de extensie - memorie de la nodul S corespunzător este conectat la linie. Canalele de rutare și logica de rutare a mesajului ce aparțin nodului S sunt partajate între procesoarele active și cele de extensie prin intermediul magistralelor interne.

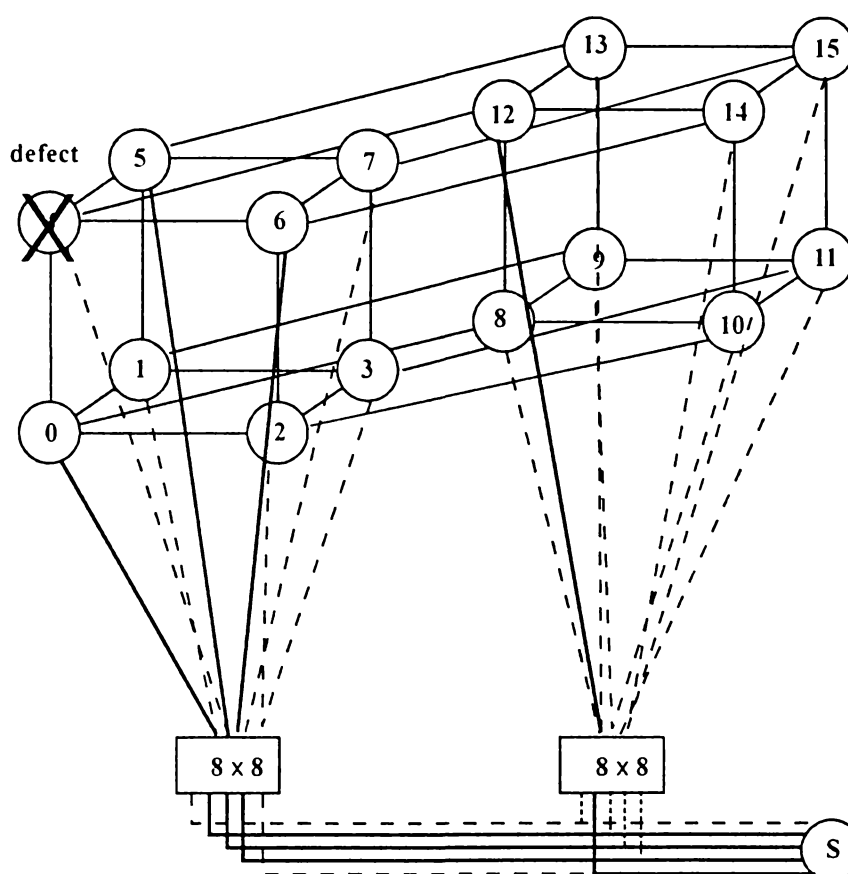
O chestiune importantă o reprezintă modalitatea optimă de alocare a procesoarelor de extensie la un *hipercub*. Evident, alocând un procesor de extensie pentru fiecare procesor, fiecare nod din *hipercub* devine un nod S. Dacă un procesor de extensie este asignat unui set de procesoare active, atunci costul hardware-ului adițional poate fi minimizat.

Trebuie rezolvată problema mapării nodurilor S în topologia *hipercubului* astfel încât fiecare procesor din sistem să fie adiacent la cel puțin un procesor de extensie. S-au dezvoltat algoritmi adecvați pentru rezolvarea acestei probleme. Fig. 4.29a arată un exemplu de alocare a patru procesoare de extensie (noduri S) la un *hipercub* de 16 procesoare.

Fiind dat un set de defecte multiple într-un astfel de sistem, s-a dezvoltat o strategie de reconfigurare bazată pe un model de graf bipartit. Nodurile acestui graf corespund nodurilor defecte și nodurilor de extensie disponibile, iar muchiile dintre ele măsoară distanța dintre noduri ce trebuie să fie traversate între ele. Algoritmul de reconfigurare oferă o mapare optimă a nodurilor defecte pe nodurile disponibile prin minimizarea dilatării totale a conexiunilor logice de mapare de-a lungul multiplelor conexiuni fizice. Figura 4.29b prezintă un *hipercub* cu patru defecte localizate la nodurile 0, 5, 6, 15. Figura 4.29c prezintă reconfigurarea unui *hipercub* cu 16 procesoare supus la patru defecte. Defectul în nodul 0 este înlocuit printr-un nod de extensie 4, defectul din nodul 3 este înlocuit prin nodul 12, defectul din nodul 6 prin nodul 2 și defectul din nodul 15 este înlocuit prin nodul de extensie 10.

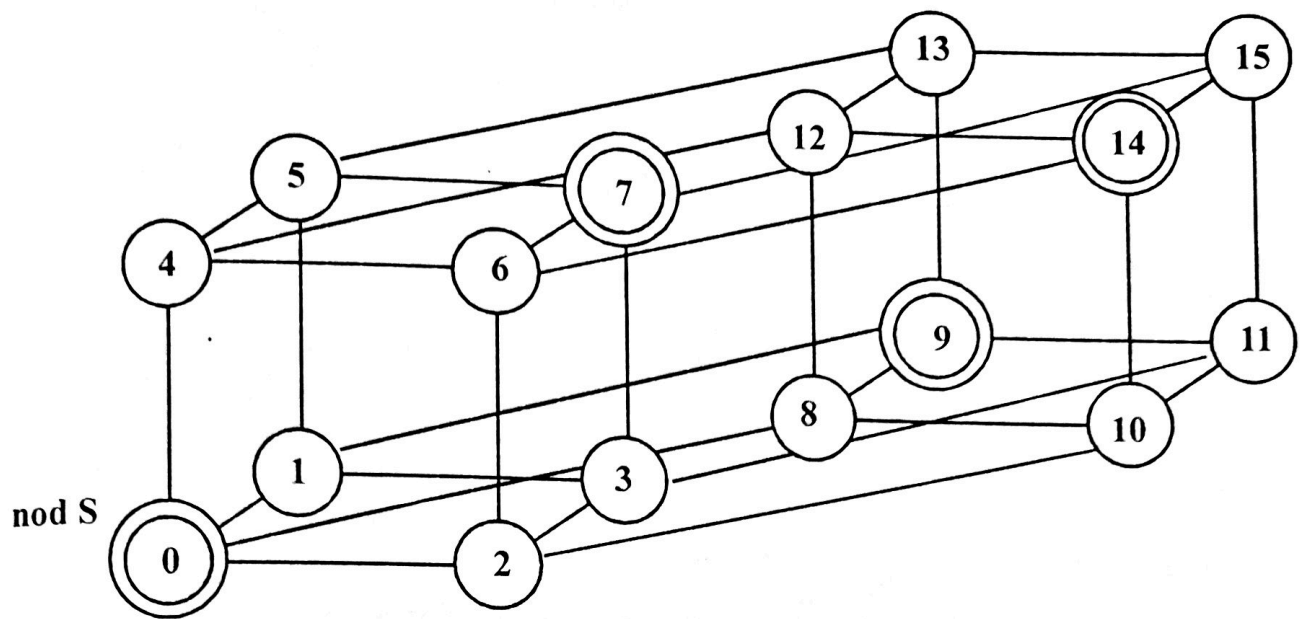


a. Conceptul de adăugare a unui procesor de extensie prin intermediul extra portului

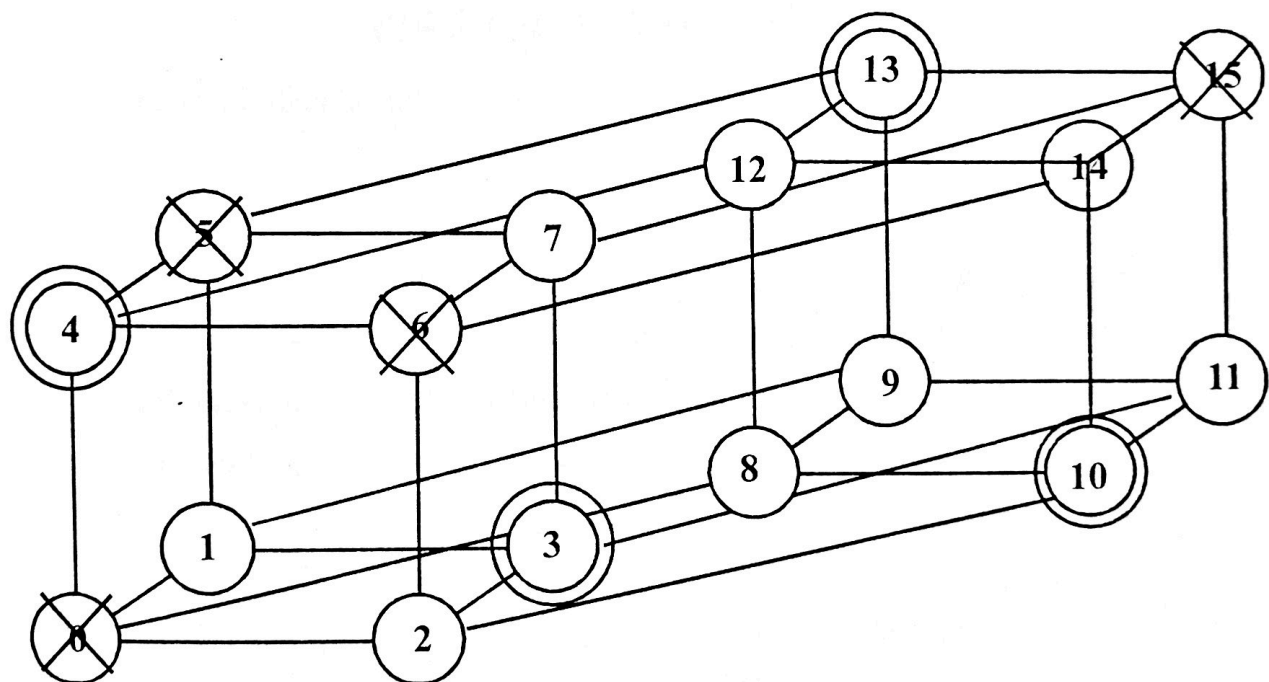


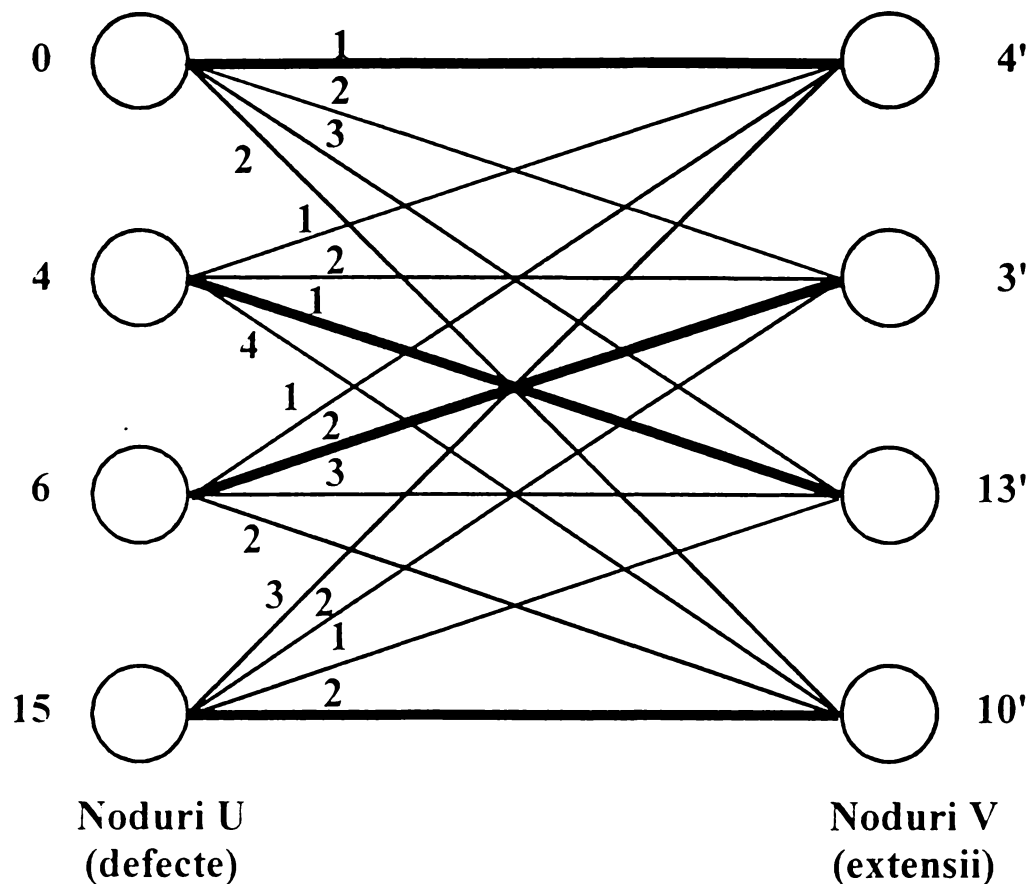
b. Implementarea unui extraport utilizând comutatoare *crossbar*

Fig. 4.28 Reconfigurarea de *hipercub* prin extraporturi



a. Alocarea de noduri de extensie

b. Localizarea a patru defecte în *hipercub*



c. Grafic de reconfigurare

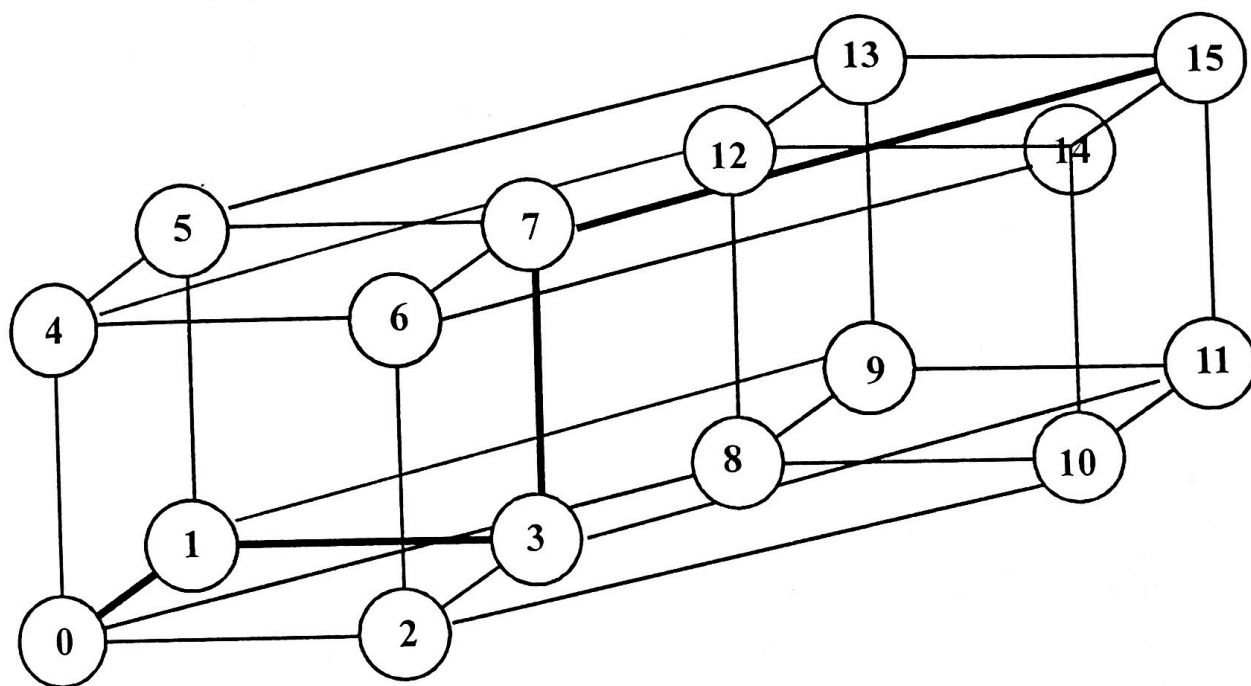
Fig. 4.29 Reconfigurare a hiper-cubului utilizând noduri de extensie

Rutarea tolerantă la defect în rețele hiper-cub.

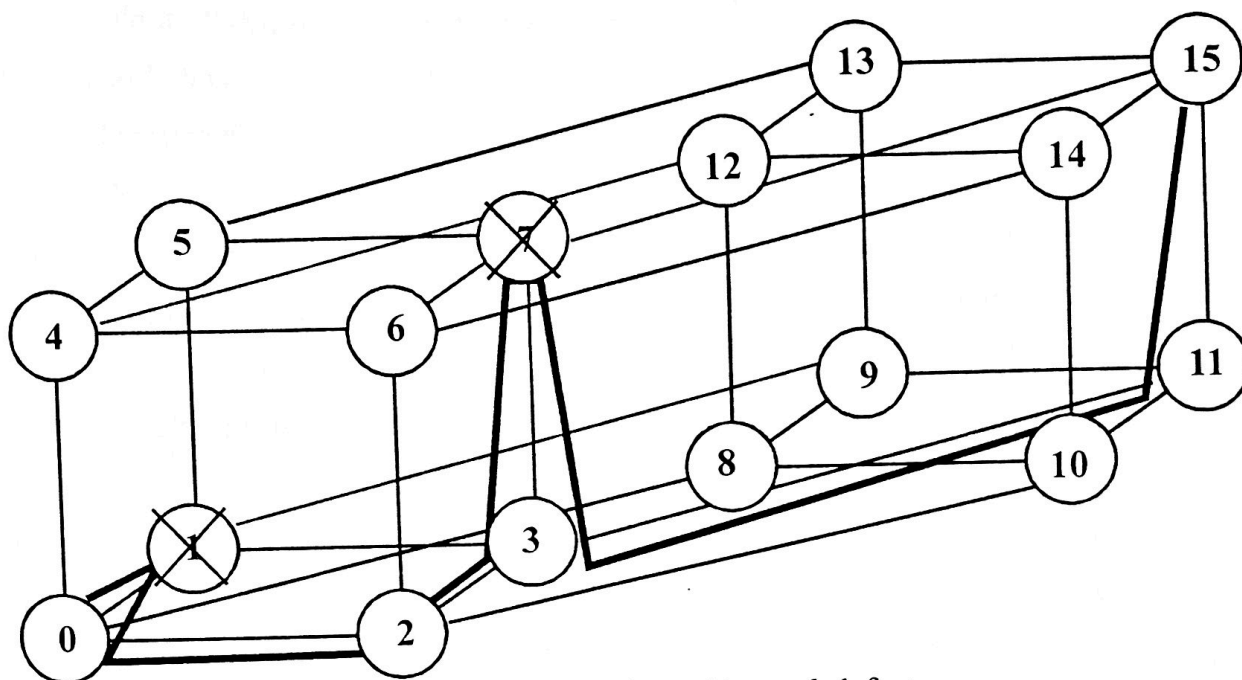
Procedura de rutare fixă standard se numește *rutarea e-cub*. În algoritmul *e-cub* diferențele în bit dintre sursa și destinația unei căi se rezolvă de la dimensiuni reduse la dimensiuni mai mari. Figura 4.30a prezintă o cale de rutare *e-cub* într-un hiper-cub de 16 procesoare. În caz de defecțiune în hiper-cub, mesajele trebuie să fie rerutate utilizând o cale ca cea prezentată în fig. 4.30b.

Există și *algoritmi de rutare adaptivi*. Într-un astfel de algoritm, fiecare mesaj poartă un *antet* - o listă ordonată a dimensiunilor de traversat - și o *etichetă* de n bit. Lista este inițializată cu ordinul dimensiunilor ce vor fi traversate de o rutare *e-cub*. La fiecare nod din cale, conexiunile sunt încercate în ordinea dată în listă pînă când este întâlnită una fără defect. Conexiunea bună este reținută și aceea dimensiune este extrasă din listă. Dacă toate dimensiunile din listă poziționează calea pe conexiuni defecte, atunci o dimensiune de extindere este aleasă din exteriorul listei. Conexiunea de această dimensiune este luată în considerare și lista se mărește prin adăugarea acestei dimensiuni la sfârșit.

În [Prad 96] se descriu și alți algoritmi de rutare. Este de notat că generalizări ale *hipercuburilor* și algoritmi de rutare au fost propuși de Agrawal și Bhuyan în 1988 [Agra88].



a. Direcționare *E-cub* la un *hipercub* tolerant la defect



b. Direcționare adaptivă la un *hipercub* defect.

Fig. 4.30 Direcționare în *hipercuburi* cu funcționare corectă și în cele cu defecte

4.5.4 Reconfigurarea în rețele *plasă*

RIN *plasă* bidimensionale au fost propuse pentru calitățile lor de scalabilitate, modularitate și expandabilitate. Fig. 4.31 prezintă 16 procesoare dispuse într-o *plasă* bidimensională. O modificare a *plasei* o reprezintă *torul* unde de asemenea există conexiuni finale (end-around) pentru procesoarele de pe direcțiile externe X și Y.

Există mai multe strategii pentru a face aceste RIN *tolerante la defect*.

O primă strategie se bazează pe ideea de-a face aceste RIN capabile ca, în caz de defecțiune, să sufere o degradare grațioasă. Pentru arii bidimensionale de calculatoare, se propune *ștergerea unei întregi linii sau coloane de procesoare* la care aparține un procesor defect. În fig. 4.31 se prezintă organizarea de bază a unei scheme de reconfigurare în care se utilizează comutatoare speciale.

Dându-se un set de procesoare defecte, un set de coloane și linii este șters din topologia *plasei* în vederea obținerii unei *subplase fără defecte (fault-free submesh)*.

Din nefericire o astfel de schemă degradează performanțele în mod semnificativ pentru un singur defect. Acest lucru se întâmplă datorită faptului că se exclude un număr mare de procesoare active (o linie sau o coloană) în timpul reconfigurării.

O altă strategie ce a fost propusă constă în utilizarea unui set de *tehnici de extensie (sparing techniques)* pentru ariile bidimensionale prin utilizarea unor circuite complexe de comutare a nodului.

În metoda de reconfigurare directă, se adaugă câteva coloane și linii în plus la aria $N \times N$ de procesoare. Reconfigurarea este obținută de-a lungul unei direcții specifice prin efectuarea unei redenumiri globale a tuturor procesoarelor în care indicii logici ai procesoarelor (i, j) sunt mapați în indici fizici (k, h) a procesoarelor active. Se utilizează o *funcție de mapare*. Dându-se un procesor defect (i, j) , este necesar să se decidă asupra unei direcții (orizontale sau verticale) pentru reconfigurare ; aceasta înseamnă că celula defectă va fi înlocuită cu celula de extensie în direcție orizontală sau verticală printr-o *conexiune de ocolire (by-pass link)* în direcție orizontală sau verticală.

Ordinea pașilor de reconfigurare este critică. Strategia recomandată a fost să se clasifice primul defect într-o coloană când se scanează de jos în sus, ca fiind un defect vertical, și toate celelalte defecte ca fiind orizontale. Dacă nici o linie nu are mai mult de un defect orizontal, atunci reconfigurarea este posibilă. În fig. 4.32 se prezintă un

exemplu de reconfigurare directă a unei arii de 5×5 augmentate cu o coloană și cu o linie suplimentară. Sunt necesare conexiuni și comutatoare auxiliare.

Schema prezentată poate tolera doar câteva configurații de defect. În orice linie doar un defect poate fi marcat ca și orizontal și în orice coloană doar un defect poate fi vertical.

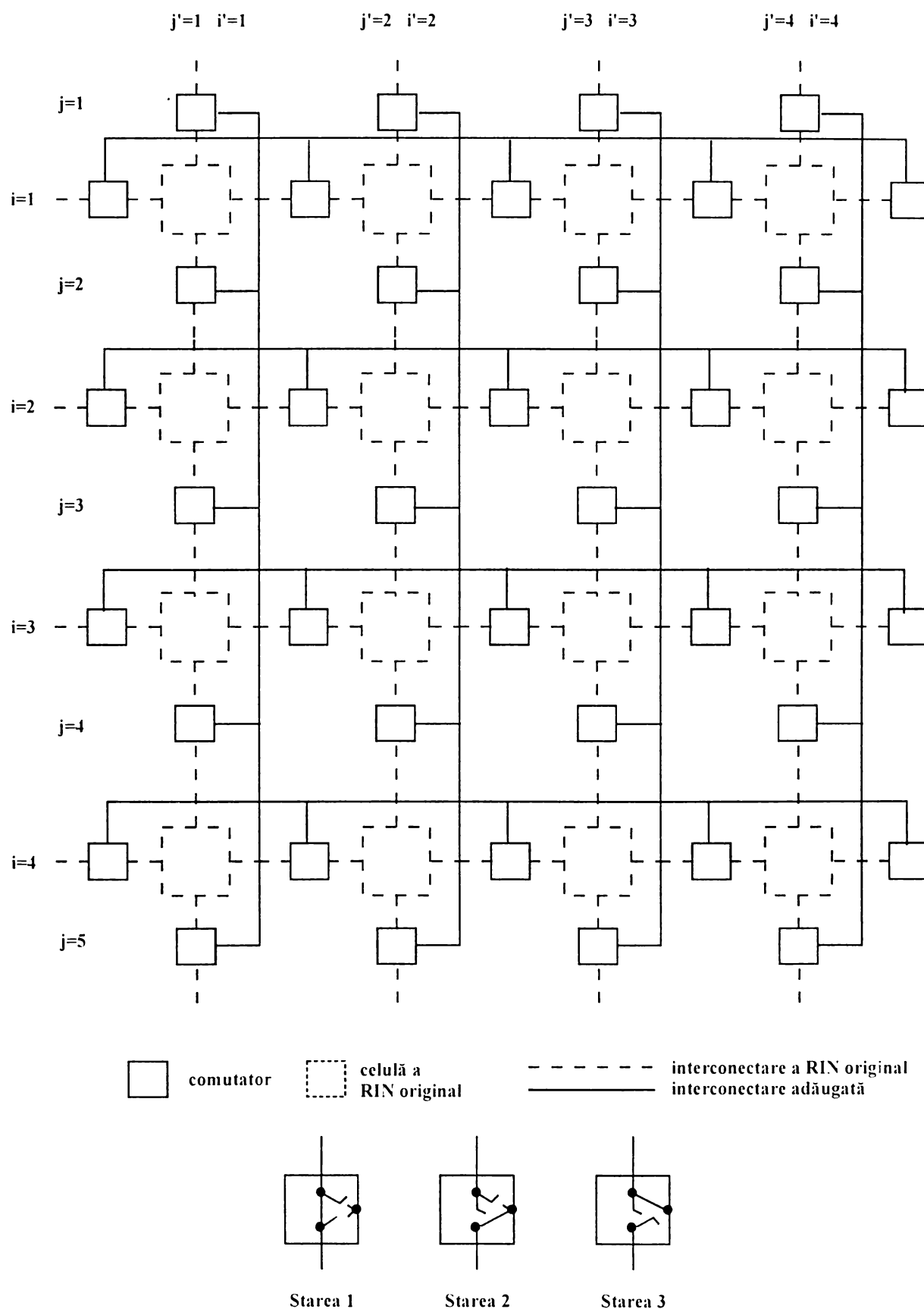


Fig. 4.31 Reconfigurarea unei RIN *plasă* prin eliminarea unei linii sau coloane

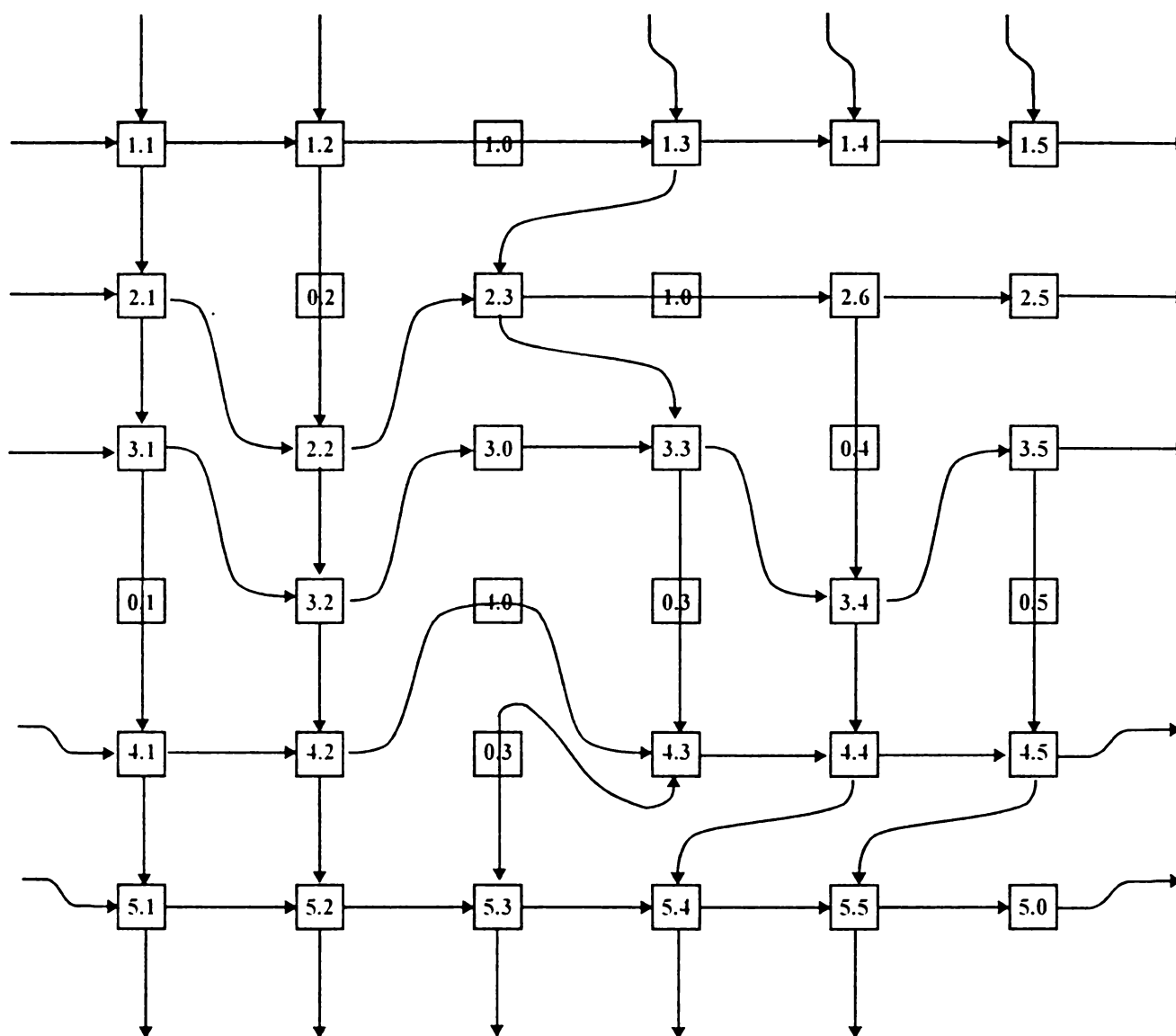


Fig. 4.32 Reconfigurarea unei RIN *plasă* prin mapare logică

4.5.5 Reconfigurarea în rețele *arbore*

RIN *hipercub* pot conecta P procesoare utilizând un grad de $\log(P)$ și o distanță maximă de $\log(P)$ iar *plasa* poate conecta P procesoare utilizând un grad de 4 și o distanță maximă de \sqrt{P} . RIN *arbore* pot conecta $P - 1$ procesoare utilizând un grad de 3 și o distanță maximă de $2(\log(P) - 1)$. Sistemele bazate pe acest tip de rețea sunt configurate ca și arbori binari unde fiecare nod este conectat la doi "copii" și la "părintele" său într-o manieră recursivă. Astfel de arhitecturi sunt excelente pentru prelucrarea algoritmilor paraleli recursivi.

Există mai multe tehnici de transformare a RIN *arbore* în RIN *arbore tolerant la defect*.

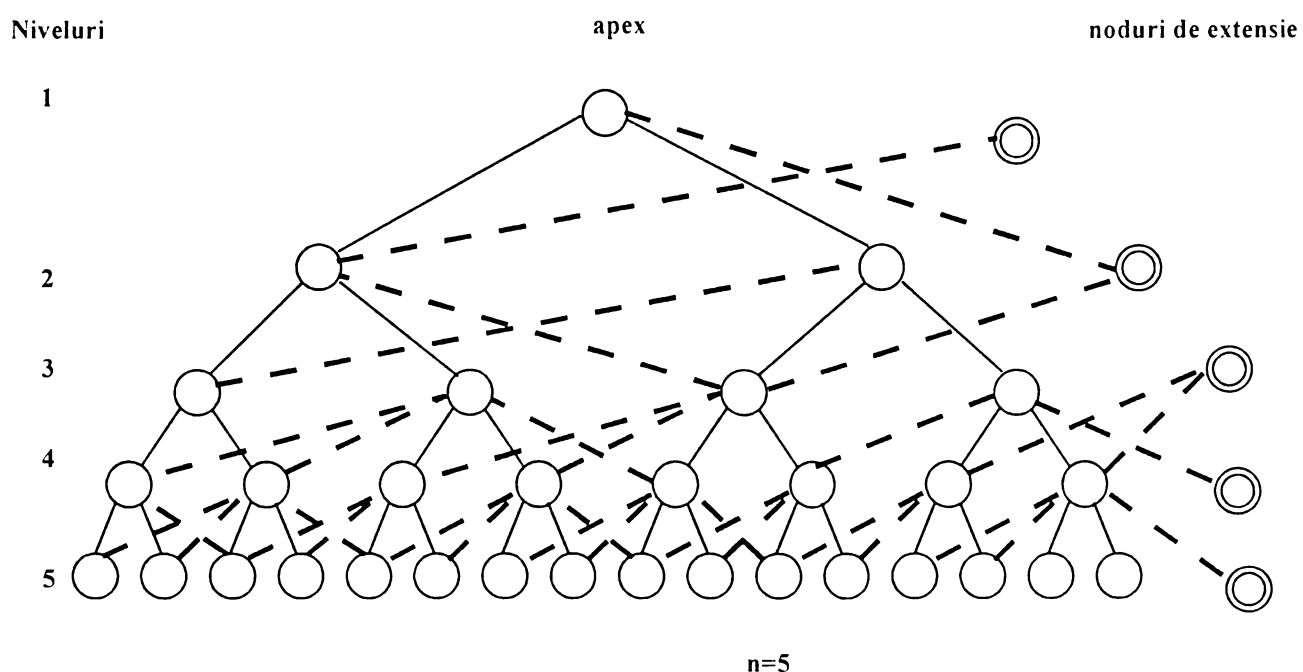
Una dintre cele mai vechi metode de reconfigurare dezvoltă arbori toleranți la un defect cu structură optimă .

Abordarea începe cu un arbore binar complet de m niveluri, utilizând m extensii, una pe nivel. Câteva conexiuni suplimentare sunt prevăzute de la nodurile de extensie ale unui anumit nivel spre nodurile arborelui binar ale următorului nivel superior astfel că în caz de defect, nodul defect și conexiunile sale sunt înlocuite de conexiunile suplimentare.

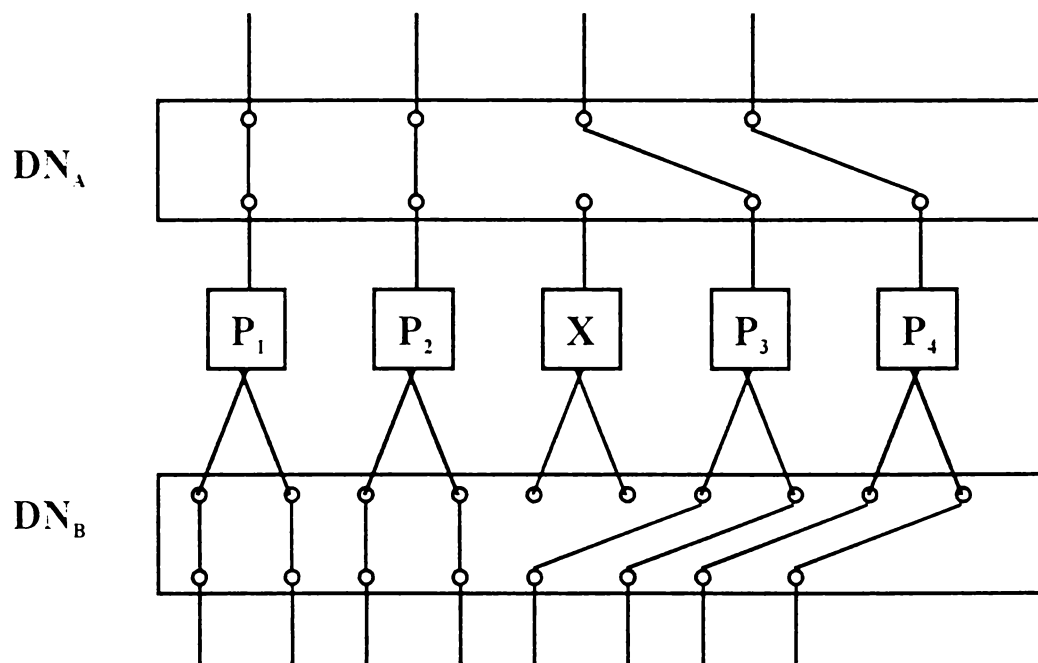
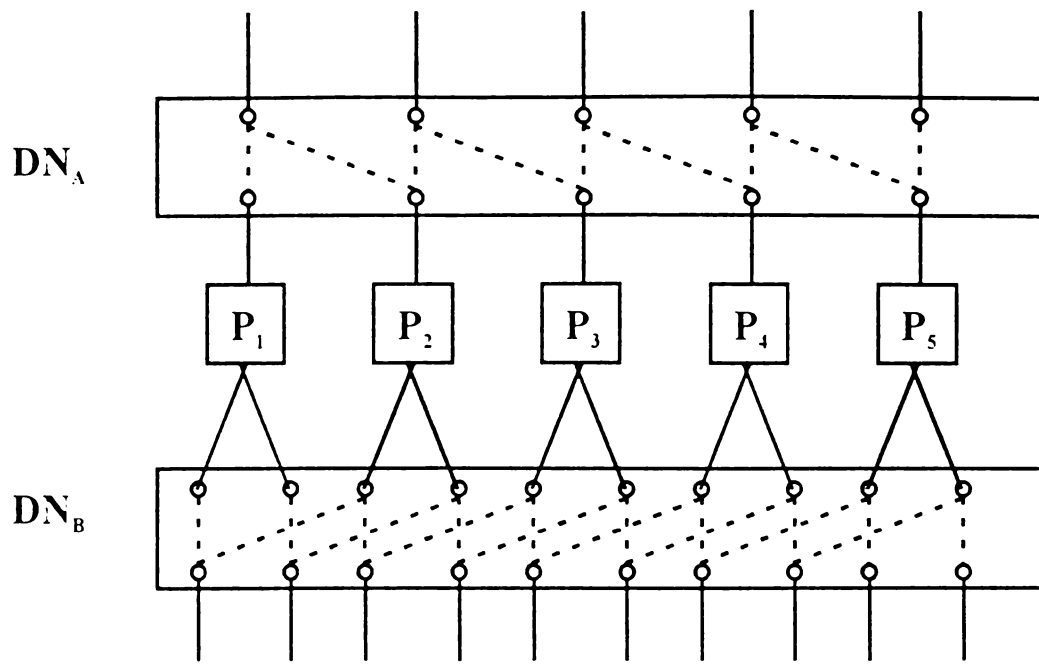
Schema poate fi ameliorată prin *adăugarea de conexiuni redundante*, suficiente pentru a tolera defecte multiple.

Un nod suplimentar este utilizat pe nivel în partea dreaptă a arborelui. Fiecare nod al arborelui include extensiile; se conectează la cei doi "copii" ai săi și de asemenea la cei doi "copii" ai vecinului său din stânga și la unul din "copii" vecinului din dreapta. Figura 4.33a prezintă schema de bază pentru reconfigurare. Când un procesor se defectează, toate conexiunile unui procesor spre dreapta sa sunt reajustate spre vecinii din dreapta. Pentru a reduce complexitatea conexiunilor redundante, pot fi utilizate două rețele de decuplare, situație prezentată în fig. 4.33b.

Schemele au fost extinse cu noduri de extensie multiple per nivel, unde un nod de extensie este poziționat pentru 2^i noduri.



a. Reconfigurarea în RIN arbore.



b. Rețea de decuplare în RIN arbore

Fig. 4.33 Reconfigurare în RIN arbore

4.5.6 Reconfigurarea RIN *multinivel*

Terminologie.

În acest paragraf vom detalia terminologia de bază utilizată în descrierea RIN-MN *tolerante la defect*.

În fig. 4.34 se prezintă structura hardware a unei RIN-MN.

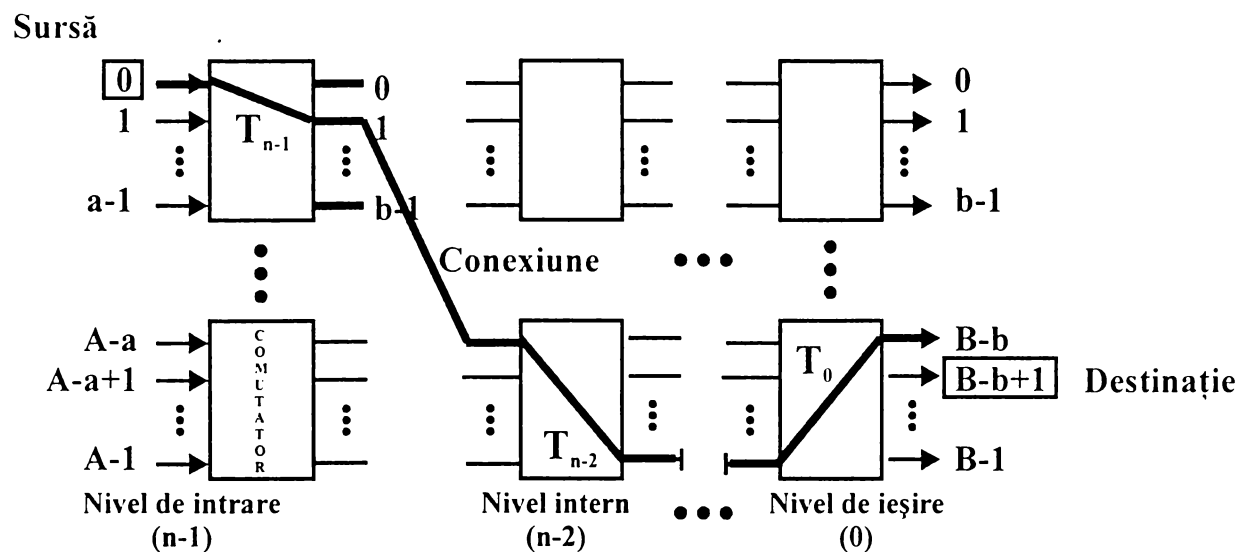


Fig. 4.34 O diagramă generică a RIN-MN evidențiată o cale

O RIN-MN este o colecție de *comutatoare* și *conexiuni* între comutatoare. Un semnal intră sau părăsește o rețea printr-un *port*. O rețea cu A porturi de intrare și B porturi de ieșire este o rețea $A \times B$.

Comutatoarele sunt dispozitive *multiport*. În RIN-MN sunt utilizate comutatoare de tip *selector* care la un moment dat conectează doar una din intrările sale la una din ieșirile sale.

O RIN-MN *tolerantă la defect* este una ce oferă servicii chiar când conține una sau mai multe componente defecte. Un defect poate fi fie *permanent* fie *tranzitoriu*. Vom considera în continuare doar *defecte permanente*. Toleranța la defect este definită doar în raport cu un *model ales de toleranță la defect (fault-tolerance model)* care are două părți. *Modelul de defect (fault-model)* caracterizează toate defectele ce se presupune că pot apare, declarând modurile de defectare. Pentru fiecare componentă a rețelei *criteriul de toleranță la defect (fault-tolerance criterion)* reprezintă condiția ce trebuie întâlnită pentru ca rețeaua să poată fi considerată că a tolerat un defect(e) dat. Un model de defect poate fi ales cu caracteristici ce simplifică analiza de fiabilitate, chiar

dacă acele caracteristici se depărtează semnificativ de realitate. În timp ce criteriile de toleranță la defect reflectă capacitățile operaționale normale (fault-free) ale rețelei, în realitate situațiile sunt mai complexe.

O rețea este *tolerantă la un singur defect (single fault tolerance)* dacă ea poate funcționa așa cum a fost specificată de criteriul de toleranță la defect. În general, dacă orice set de i defecte poate fi tolerat, atunci rețeaua este *tolerantă la i defecte (i -fault tolerant)*.

RIN-MN pot fi descrise atât *topologic (grafic)* cât și *algebric*. Topologia unei rețele reprezintă configurația de conexiuni din structura sa. Configurația poate fi reprezentată printr-un graf. Topologia este determinată de proiectarea comutatoarelor și de configurația conexiunilor. Nodurile din graful unei RIN-MN pot fi numerotate și atunci rețeaua este descrisă în termeni de relații algebrice între noduri. Modelul algebric este util în discutarea strategiilor de control și comunicare.

Reamintim că există trei forme de bază pentru conexiuni într-o rețea.

O *conexiune unu-la-unu (one-to-one connection)* trece informația de la un port al RIN, *sursa*, la alt port al RIN, *destinația*. Ruta exactă parcursă de informație constituie *calea* sa. Conexiuni multiple unu-la-unu pot fi active simultan.

O *conexiune de permutare* este un set de conexiuni unu-la-unu astfel încât două conexiuni unu-la-unu să nu aibă aceeași sursă sau destinație. Această clasă de conexiuni sunt utile doar în contextul rețelelor cu un număr egal de surse și de destinații. Fluxul de informație ce circulă simultan de la o sursă spre două sau mai multe destinații este suportată de o *conexiune de difuzare* și calea urmată se numește *cale de difuzare (broadcast path)*.

Etichetele de rutare (routing tags) reprezintă un mod de a descrie o cale printr-o rețea și de a fi oferite pentru un control distribuit al rețelei. Controlul este *distribuit* dacă dispozitivele ce utilizează rețeaua generează propriile etichete de direcționare și comutatoarele rețelei se pot poziționa bazându-se pe informația conținută în etichetă. În fig. 4.34 se prezintă un comutator în etajele $n-1$, $n-2$, și 0 , ce este poziționat de biții T_{n-1} , T_{n-2} și T_0 din etichetă. Etichetele de rutare sunt importante pentru RIN-MN *tolerante la defect* deoarece ele sunt capabile să specifice o cale funcțională, dacă ea există; limitările în numărul și dimensiunile etichetelor se reflectă în limitări ale toleranței la defect. Există trei metode utilizate de surse pentru a genera etichete de rutare ce specifică o cale liberă de defect. Cu o *rutare nonadaptivă* o sursă realizează existența unui defect doar când calea ce trebuie stabilită traversează componenta

defecții din RIN. Elemente de identificare al defectului sunt transmise spre sursă, care încearcă următoarea cale alternativă. Această abordare necesită un număr redus de circuite dar și performanțele obținute sunt reduse.

O particularizare constă în faptul că uzual $m = 2$ și atunci vor exista două variante de rutare adaptivă. Cu modificarea la sursă (modification on demand) a sursă menține un tabel cu defectele pe care le-a întâlnit în încercarea de-a stabili o cale și utilizează această informație pentru a ghida o nouă rutare.

Cu RIN cu $\log N$ niveluri de două proprietăți importante erori, toate sursele sunt notificate asupra componentelor defecte pe măsură ce ele sunt diagnosticate.

O cale liberă de defect într-un RIN este specificată de ieșire, dar dacă o conexiune rutare pot fi modificate ca răspuns la defectele întâlnite pe măsură ce o cale este

Apărând avantajele datorită acestei topologii poate fi realizată în RIN-MN construite cu comutatoare 2×2 realizată de către o conexiune stabilă la o dată de efect

În literatură [Prad 87], [Kro 90], [Moo 96], [Moo 98] se face o analiză sistematică și o comparație a RIN-MN non-reinversiuni sau al unui insistent, de conexiune tehnic și căi de comunicație conducând la o lipsă a toleranței la defect și la o

fiabilitate scăzută.

În fig. 4.35 se ilustrează această situație pentru o RIN omega având 16 intrări și 16 ieșiri construită cu comutatoare 2×2 suplimentare.

Figura evidențiază toate căile care utilizează conexiunea 1×1 defectiune la nivel de l deconectează aceste căi. Situația devine dramatică atunci când dimensiunea RIN crește deoarece numărul de căi ce traversează o conexiune dată, crește liniar cu N .

Tehnicile pentru obținerea toleranței la defect în RIN-MN pot fi clasificate după modul în care ele implică modificarea topologiei sistemului.

Există metode cunoscute de creare a toleranței la defect în RIN-MN care nu modifică topologia sistemului. Aceste tehnici utilizează coduri detectoare de erori, *bit-slice* etc.

Ele sunt însă limitate la o singură defectiune. Concluziile desprinse demonstrează că domeniul RIN-MN tolerante la defect este Orice schemă de toleranță la defect care utilizează notiunea de modificare a topologiei încă în stadiul incipient, având largi posibilități de dezvoltare.

RIN implică detectarea și localizarea defectelor în rețea, urmată de o reconfigurare. Aplicarea unor metrice fundamentale ca și modelul de defect și criteriul de toleranță corespunzătoare. În literatură, [Prad 96] se menționează mai multe metode de detectarea defectului în RIN-MN. De exemplu se utilizează tehnici de testare *off-line* de RIN-MN.

sau *on-line* în care se aplică un număr redus de teste independent de talia rețelei. Se RIN-MN au fost descrise în §2.4.4. Ele se caracterizează prin următorii parametri poziționează toate comutatoarele aplicând la intrări configurații 01 sau 10. Dacă nu constructivi:

există defecte se poate demonstra că porturile de ieșire vor furniza configurații 01 sau $N=m$ intrări

10. Apoi se poziționează toate comutatoarele în poziția “*încrucișat*” (*cross-over*) și se aplică aceleași configurații 01,10 iar la ieșiri se testează dacă apar 10, 01.

Dacă apare o defecțiune, câteva dintre ieșiri vor furniza o configurație ilegală adică 00 sau 11. Configurațiile pot fi utilizate pentru localizarea defectului.

Odată ce un defect a fost detectat în RIN-MN există numeroase strategii pentru reconfigurarea rețelei. Acestea includ replicarea întregii rețele, adăugarea unui nivel suplimentar de comutatoare, varierea dimensiunii comutatorului, adăugarea de conexiuni suplimentare între niveluri, adăugarea de porturi suplimentare, etc.

În continuare se va prezenta o metodă de transformare a unui RIN-MN cu o *singură cale* într-un RIN-MN cu *căi multiple* între fiecare pereche de I/E a rețelei.

Factorii care contribuie la *complexitatea comutatorului* (*switch complexity*) într-un RIN-MN sunt: numărul de nivele de comutare, numărul de comutatoare/nivel și talia comutatoarelor (numărul de contacte).

O altă măsură a costului unui RIN-MN este *complexitatea conexiunii* (*link complexity*) care depinde de numărul de conexiuni în cadrul unui nivel și de numărul de niveluri. Dintre cele două tipuri de conexiuni în RIN-MN, cele între niveluri sunt mai scumpe decât cele în cadrul unui nivel.

Un *model de defect* (*fault-model*) are în vedere efectele defectelor fizice în operarea unui sistem. Se utilizează trei modele: *stuck-at fault*, *modelul defectului de conexiune* și *modelul defectului de comutator*.

În modelul *stuck-at*, o defecțiune cauzează blocarea comutatorului într-o stare particulară indiferent de starea intrării de comandă. Astfel se afectează capacitatea de a stabili conexiunile necesare. Comutatorul poate fi folosit pentru stabilirea căii de comunicație dacă starea blocată a comutatorului este de asemenea starea solicitată. În modelul *defectului de conexiune*, o defecțiune afectează o conexiune individuală a unui comutator, lasând partea care rămâne a comutatorului, operațională.

În modelul de *defect de comutator* – cel mai consistent dintre modelele prezentate – o defecțiune face comutatorul total inutilizabil. Exemplul prezentat a fost proiectat să tolereze astfel de defecte de comutator.

Topologia RIN-MN *cu căi multiple* poate permite ca rerutarea să fie făcută doar la sursă sau în anumite puncte fixe din rețea. În acest caz o linie ocupată, o linie defectă, sau un comutator defect, întâlnite în timp ce trebuie stabilită o cale, pot necesita o *parcurgere inversă a traseului* (*back tracking*) până la un nivel la care există o bifurcație și de acolo să se încerce stabilirea altei căi de comunicație. Dacă căile între

fiecare pereche I/E într-o RIN-MN *cu căi multiple* are o bifurcație la fiecare nivel, atunci redirecționarea poate fi făcută din mers fără procedura de parcurgere inversă.

Un *graf de redundanță* oferă o cale convenabilă pentru a studia proprietățile RIN-MN *cu căi multiple* ca și numărul de defecte tolerate sau tipul posibil de redirecționare. Acest graf descrie toate căile posibile între o sursă și o destinație în cadrul unei rețele. El constă din două noduri distincte – sursa S și destinația D – și restul nodurilor corespunde cu comutatoarele ce apar pe căile între S și D.

Condițiile de proiectare pentru un RIN-MN *omega* tolerant la defect sunt:

- *toleranță la defect de comutator* în toate nivelurile, inclusiv primul și ultimul
- posibilitatea de *rerutare* fără parcurgerea inversă a traseului
- *reparabilitatea on-line* a componentelor defecte
- *complexitatea scăzută* a comutatoarelor și a conexiunilor.

RIN-MN *omega* cu nivel suplimentar

O RIN-MN *omega* căreia i se adaugă conexiuni și comutatoare suplimentare pentru a o face tolerantă la defect se numește o RIN-MN *omega* *cu nivel suplimentar* (*extra-stage shuffle-exchange network*).

În literatură [BT88] ,[BT89], RIN-MN *omega* și RIN-MN *omega* *mărită* se notează, din rațiuni de simplitate, cu SEN (Shuffle Exchange Network), respectiv SEN+ Vom utiliza următoarele convenții de reprezentare:

- nivelurile sunt numărate de la 0 la $n-1$
- comutatoarele în fiecare nivel sunt numărate de la 0 la $N/2-1$
- comutatorul i de pe nivelul j este referit ca și comutatorul (i,j)
- conexiunile superioare (denumite conexiuni 0) de comutatoare i sunt numărate $2i$
- conexiunile inferioare (denumite conexiuni 1) sunt numărate $2i+1$.

Rerutarea se face într-o manieră distribuită utilizând eticheta binară a destinației ca și *etichetă de rutare* (*routing tag*).

Fie d_0, d_1, \dots, d_{n-1} adresa binară a destinației dorite.

Comutatoarele de-a lungul căii din nivelul i utilizează bitul d_i și fac o conectare spre ieșirea superioară dacă bitul este 0 și spre ieșirea inferioară dacă bitul este 1. Generarea etichetei de direcționare este sarcina nodului sursă.

Pentru a obține toleranța la defect se exploatează faptul că există subseturi de comutatoare în fiecare nivel care se găsesc pe căile ce conduc la același subset de destinații. De exemplu, orice comutator de pe nivelul 0 conduce la orice destinație; toate comutatoarele *pare* de pe nivelul 1 conduc la destinații de la 0 la $N/2-1$; toate comutatoarele *impare* de pe nivelul 1 conduc la destinații numerotate de la $N/2$ la $N-1$. Toate comutatoarele de pe un nivel dat care conduc la același subset de destinații cuprind un *subset conjugat* de comutatoare. Comutatoarele de pe fiecare nivel sunt partiționate în câteva subseturi conjugate.

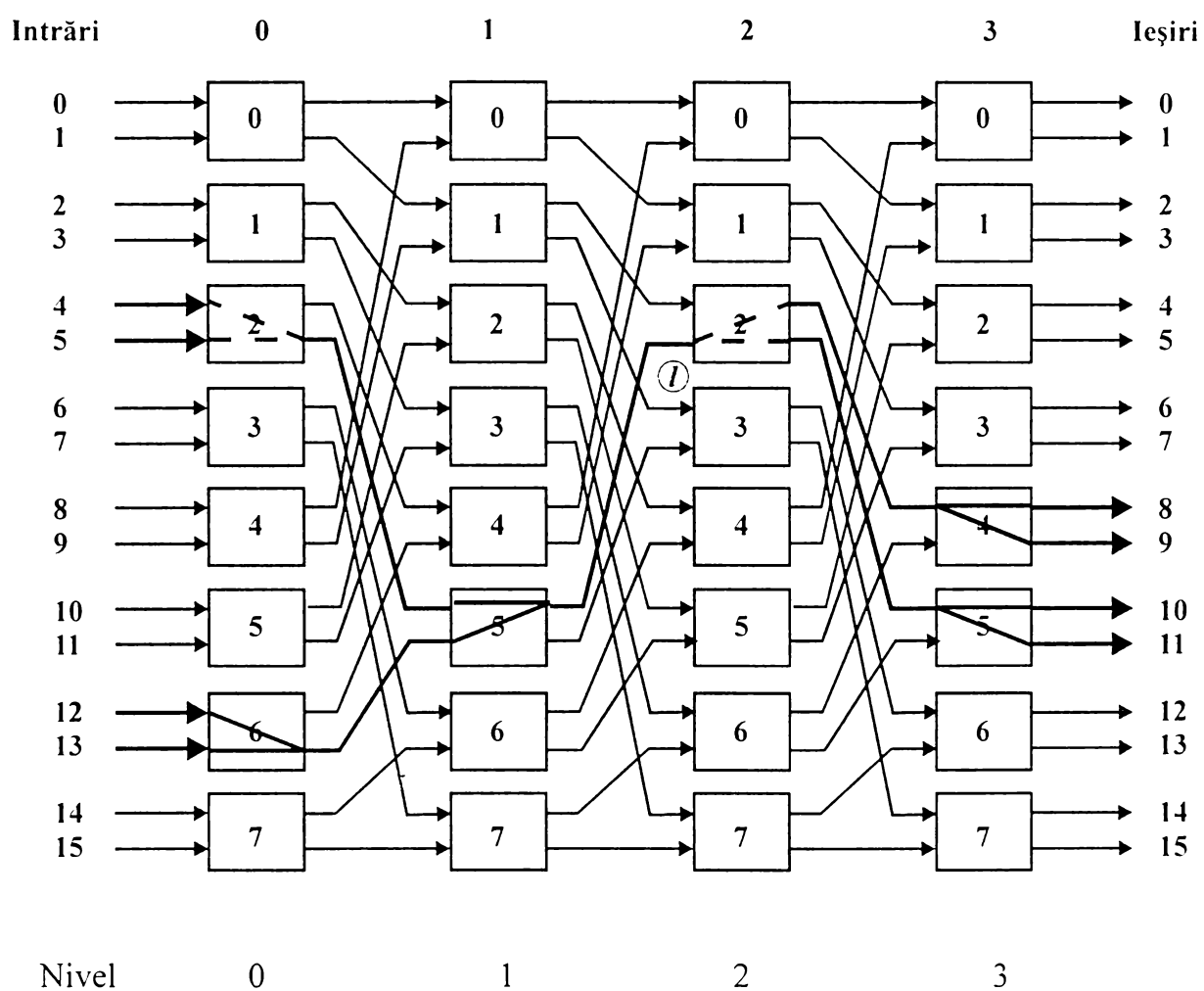


Fig. 4.35 RIN omega 16×16 . Există două căi ce utilizează conexiunea 1

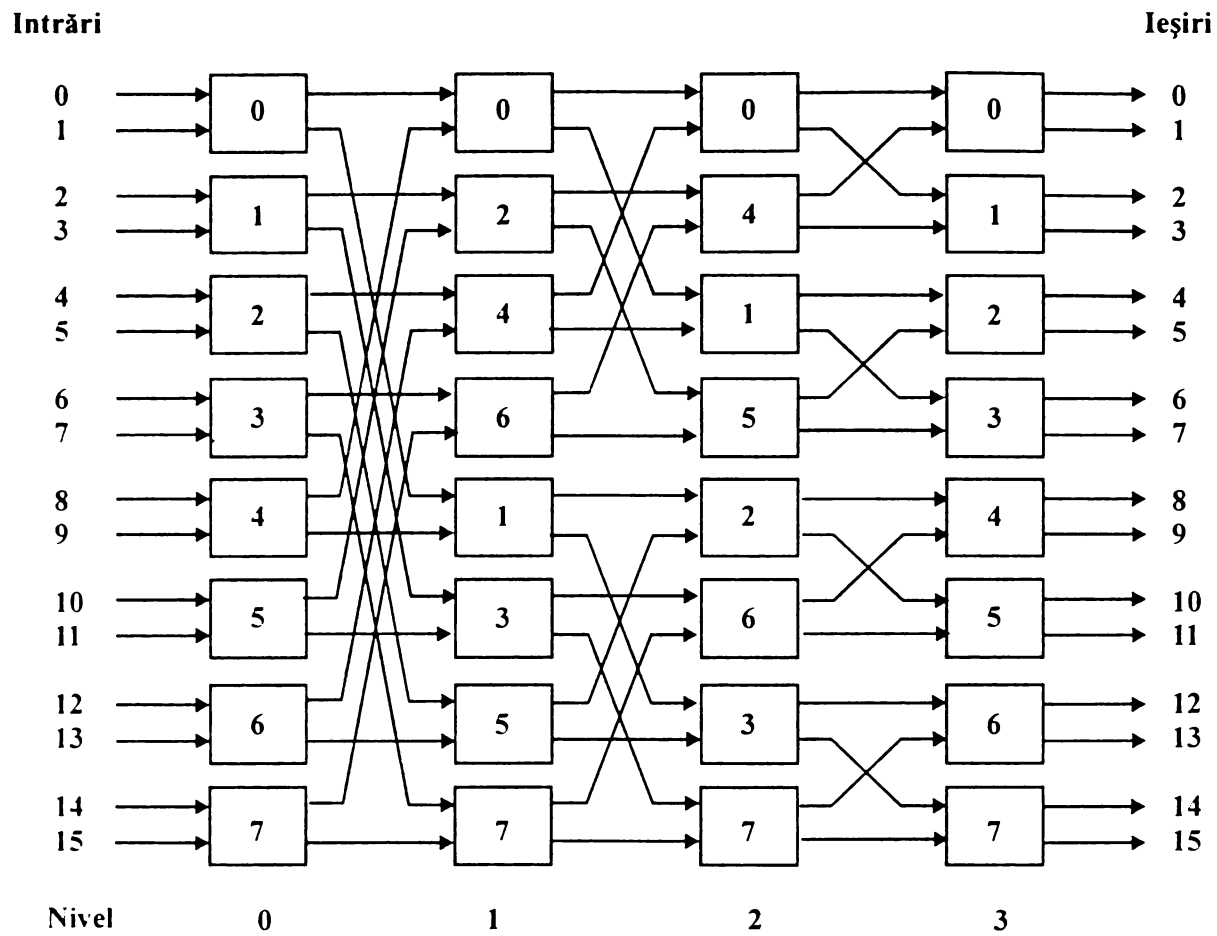
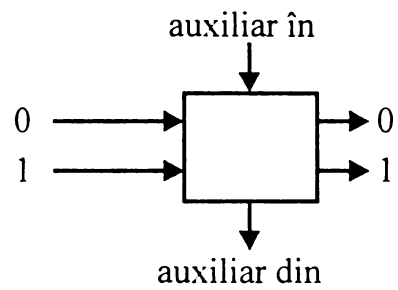


Fig. 4.36 RIN-MN redesenată, cu identificarea subseturilor conjugate de comutatoare.



a) Comutator necesar pentru formarea buclelor de comutare

Fig. 4.37 Comutatorul defect (1, 2) este tolerat cînd intrarea 9 este conectată la ieșirea 4

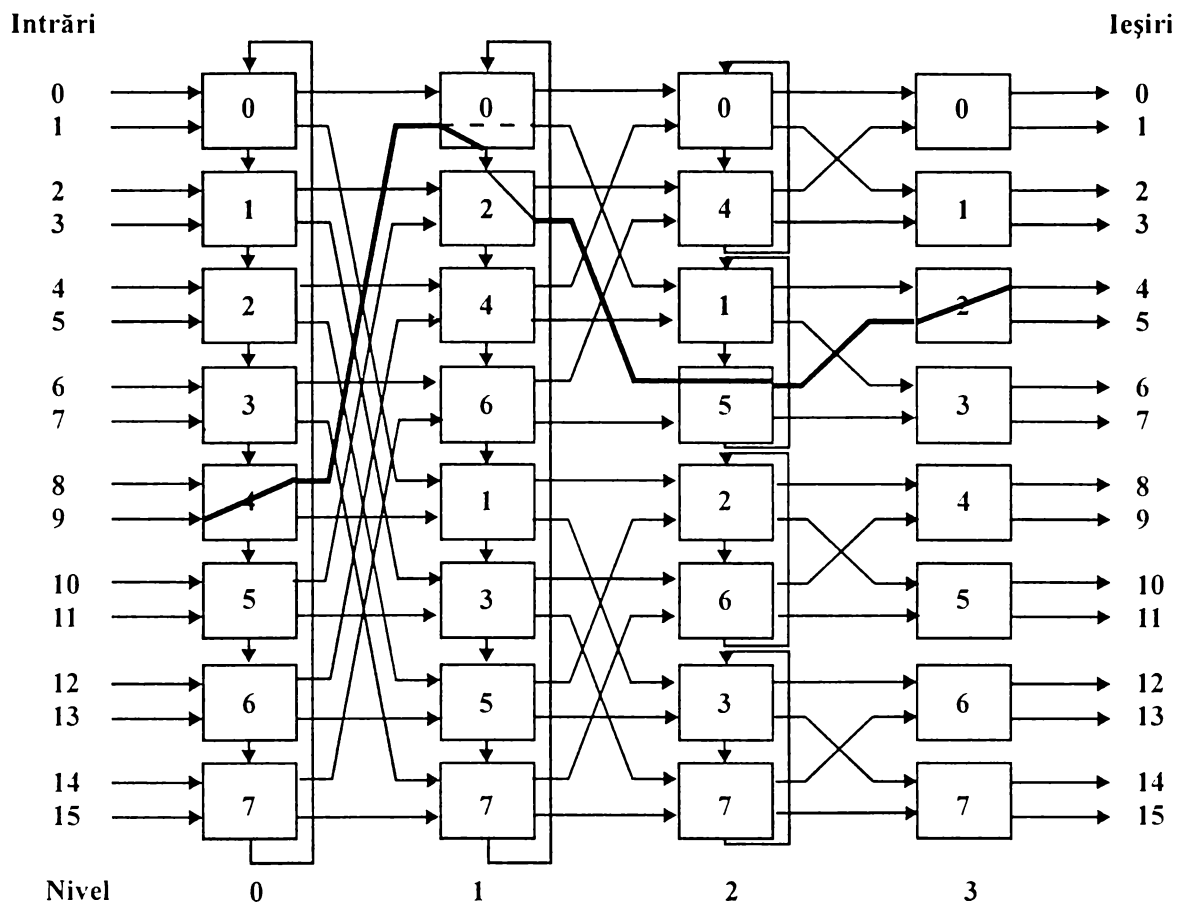
b) RIN_MN *omega augmentat*

Fig. 4.37 Comutatorul defect (1, 2) este tolerat cînd intrarea 9 este conectată la ieșirea 4

În fig. 4.36 se prezintă RIN-MN din fig. 4.35 redesenată pentru facilitarea identificării subseturilor conjugate de comutatoare. De exemplu pe nivelul i , două comutatoare k și l aparțin aceluiași subset conjugat dacă $k \equiv l \pmod{2^i}$. Eticheta de rutare parțială necesară pentru stabilirea unei conectări pentru o destinație, fie de la comutatorul k fie de la l , este aceeași.

În fiecare subset conjugat de comutatoare există câteva perechi de comutatoare denumite *perechi conjugate de comutatoare* (*conjugate pairs of switches*). Comutatoarele într-o astfel de pereche sunt conectate la aceleași comutatoare din nivelul următor. De exemplu în fig. 4.35, comutatoarele (0,0) și (4,0) formează o pereche conjugată deoarece ambele sunt conectate la comutatoarele (0,1) și (1,1). În general, comutatoarele (i,j) și $(i+N/4,j)$ unde $0 \leq i < N/4$, $0 \leq j < n-1$, sunt ambele conectate la comutatoarele $(2i,j+1)$ și $(2i+1,j+1)$. Subseturile conjugate și perechile conjugate de comutatoare joacă un rol fundamental în acest tip de scheme.

Schema pentru *crearea de căi multiple* se bazează pe conectarea comutatoarelor aparținând unui subset conjugat prin utilizarea conexiunilor adiționale pentru a forma bucle. Pentru aceasta este necesară înlocuirea comutatoarelor 2×2 cu comutatoare 3×3 de tipul prezentat în fig. 4.37a. Conexiunile verticale utilizate pentru conectarea comutatoarelor pe același nivel sunt denumite *conexiuni auxiliare (auxiliary links)*. RIN-MN din fig. 4.35, augmentată în această manieră, este prezentată în fig. 4.37b. Această figură prezintă de asemenea cum poate fi obținută *rutarea tolerantă la defect*. Dacă un comutator nu este capabil să proceseze o solicitare de conectare datorită unui comutator defect dispus pe nivelul următor sau datorită unei conexiuni ocupate, el poate ruta solicitarea prin conexiunea auxiliară de ieșire spre următorul comutator din buclă. Următorul comutator poate apoi să facă o conectare la un comutator diferit (corect) dispus pe nivelul următor.

Dacă se presupune că trebuie făcută o conectare între conexiunea de intrare 9 spre conexiunea de ieșire 4, calea de date în condiții normale ar fi prin comutatoarele (4,0), (0,1), (1,2) și (2,3). În ipoteza defectării comutatorului (1,2), traseul inițial va fi schimbat utilizându-se calea (4,0), (0,1), (2,1), (5,2) și (2,3).

De fapt se observă că utilizând o conexiune auxiliară ori de câte ori se întâlnește un defect, se permite oricărei surse să fie conectată la orice destinație în timp ce se tolerează orice comutator singular defect în orice nivel cu excepția celui inițial și a celui final.

O defectare a unui comutator în nivelul inițial deconectează sursele atașate acelui comutator de restul rețelei. Similar orice defectare a unui comutator de pe nivelul final deconectează destinațiile atașate comutatorului. Pentru a extinde toleranța la defect în aceste niveluri, fiecare sursă și fiecare destinație trebuie să fie prevăzută cu cel puțin un port adițional de I/E, astfel ca ele să poată fi conectate la cel puțin două comutatoare distincte. În fig. 4.38 se prezintă schema pentru conectarea surselor și destinațiilor echipate cu două porturi fiecare. Un multiplexor 2×1 este plasat la fiecare conexiune de intrare a nivelului 0 și fiecare sursă este conectată la două MUX-uri distincte. Pentru a prevedea coexiuni adiționale spre destinații, cele $N/2$ comutatoare din nivelul final al RIN-MN sunt înlocuite cu $N \times 2$ DMUX-uri. Fiecare destinație este direct conectată la două DMUX.

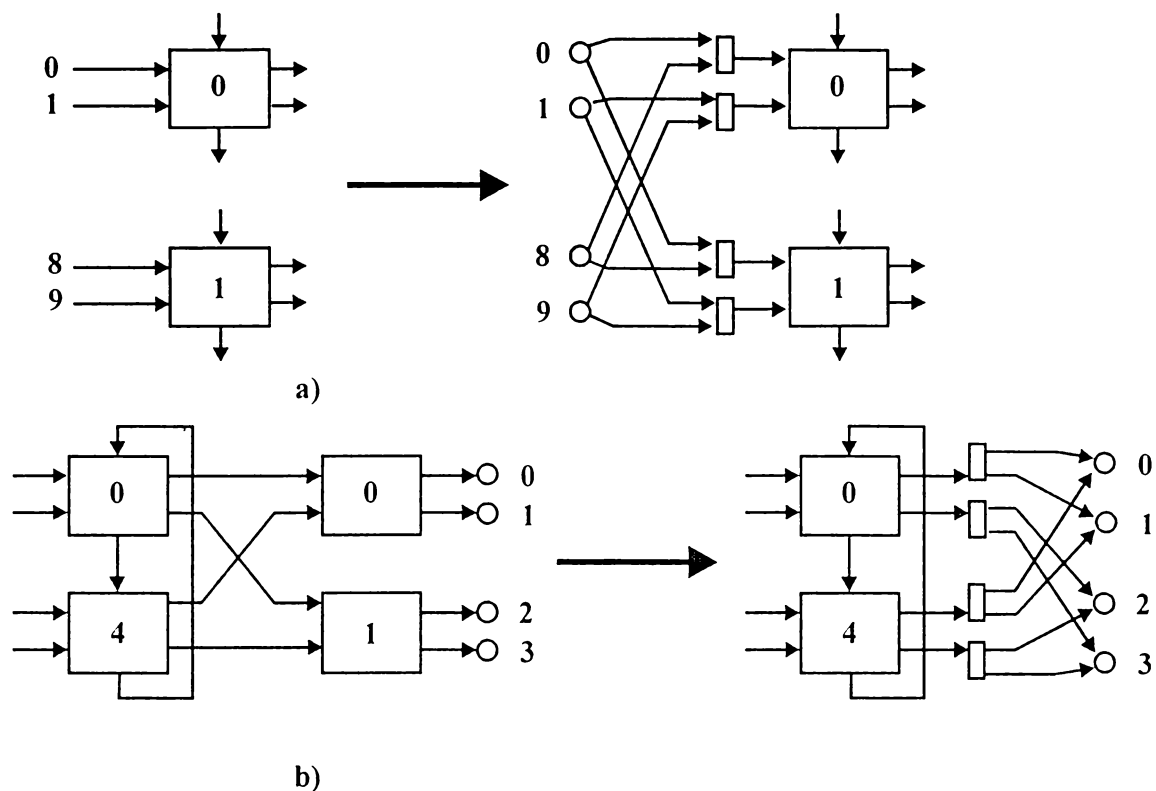


Fig. 4.38 Scheme de realizare unor conexiuni multiple

pentru surse și destinații - adăugarea a MUX-uri 2×1 la nivelul de intrare (a);

înlocuirea comutatoarelor din nivelul final cu DMUX-uri 1×2 (b)

Dacă fiecare buclă din fig. 4.38 este divizată în două bucle, cu condiția ca să nu existe două comutatoare într-o buclă dată ce să formeze o pereche conjugată, atunci se obține o rețea care are nu numai toleranță la defect singular dar și câteva proprietăți ce să faciliteze reparația on-line și întreținerea.

O astfel de rețea este denumită ASEN (Augmented Shuffle-Exchange Network). În fig. 4.39a se reprezintă un ASEN pentru $N=16$ iar graful său de redundanță este dat în fig. 4.39b. Această rețea se numește ASEN-Max deoarece ea are numărul maxim de comutatoare pe buclă cu restricțiile că comutatoarele dintr-o buclă aparțin aceluiași subset conjugat și că nu există două comutatoare într-o buclă ce să formeze o pereche conjugată. Numărul de comutatoare dintr-o buclă de pe nivelul i a lui ASEN-Max este 2^{n-2-i} . De notat că în nivelul final (nivelul $n-2$) există doar un comutator/buclă.

Este posibil să se proiecteze ASEN-Max și să se rețină proprietatea toleranței la un singur defect. De exemplu ASEN din fig. 4.40 are o buclă de talia a două comutatoare în fiecare nivel cu excepția ultimului. În general, un ASEN în care numărul de comutatoare dintr-o buclă pe nivel i este egal cu $\text{Min}(2^{n-2-i}, k)$ este denumit ASEN- k .

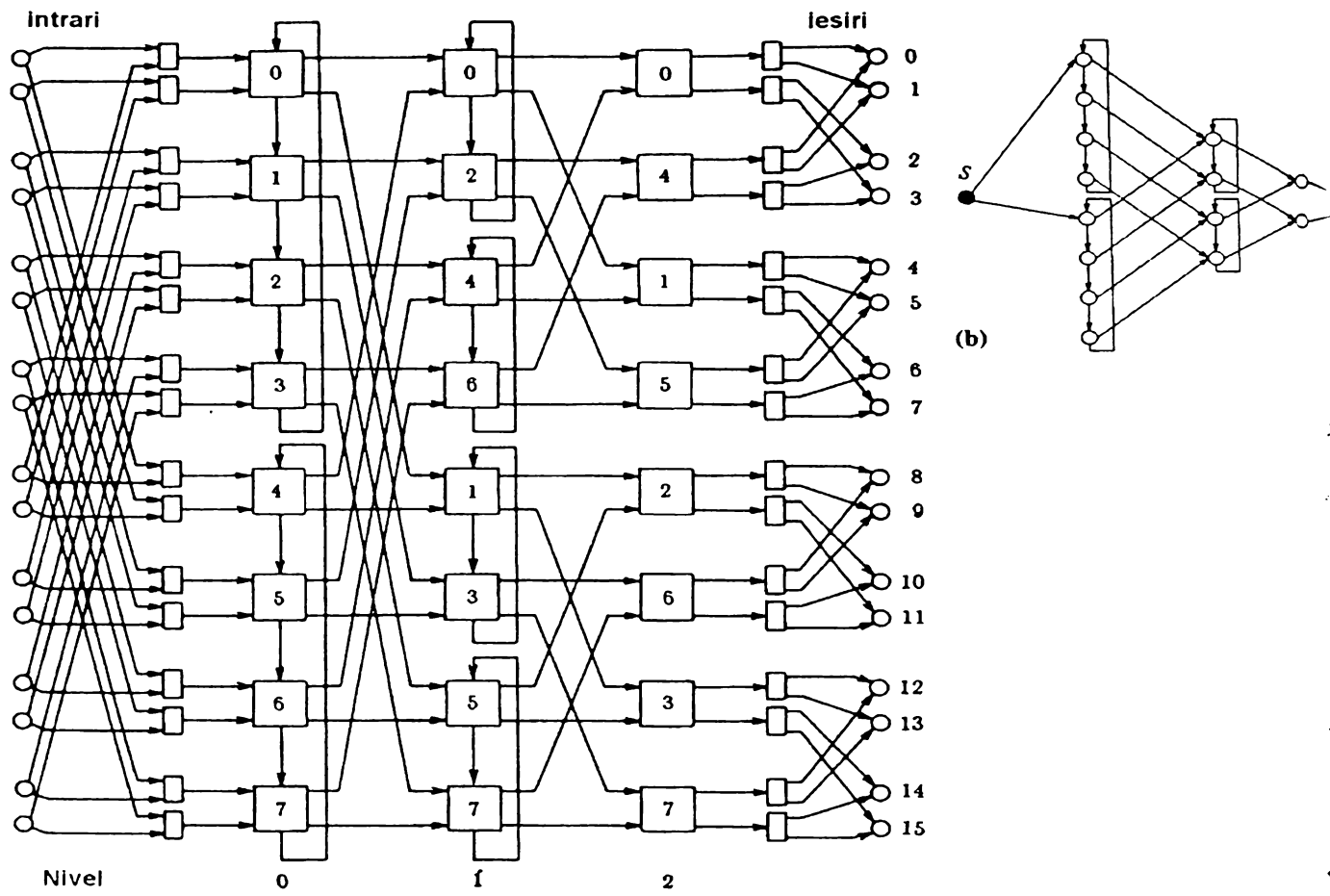


Fig. 4.39 Un ASEN – Max pentru $N=16$ (a); graful de redundanță al acestui circuit(b)

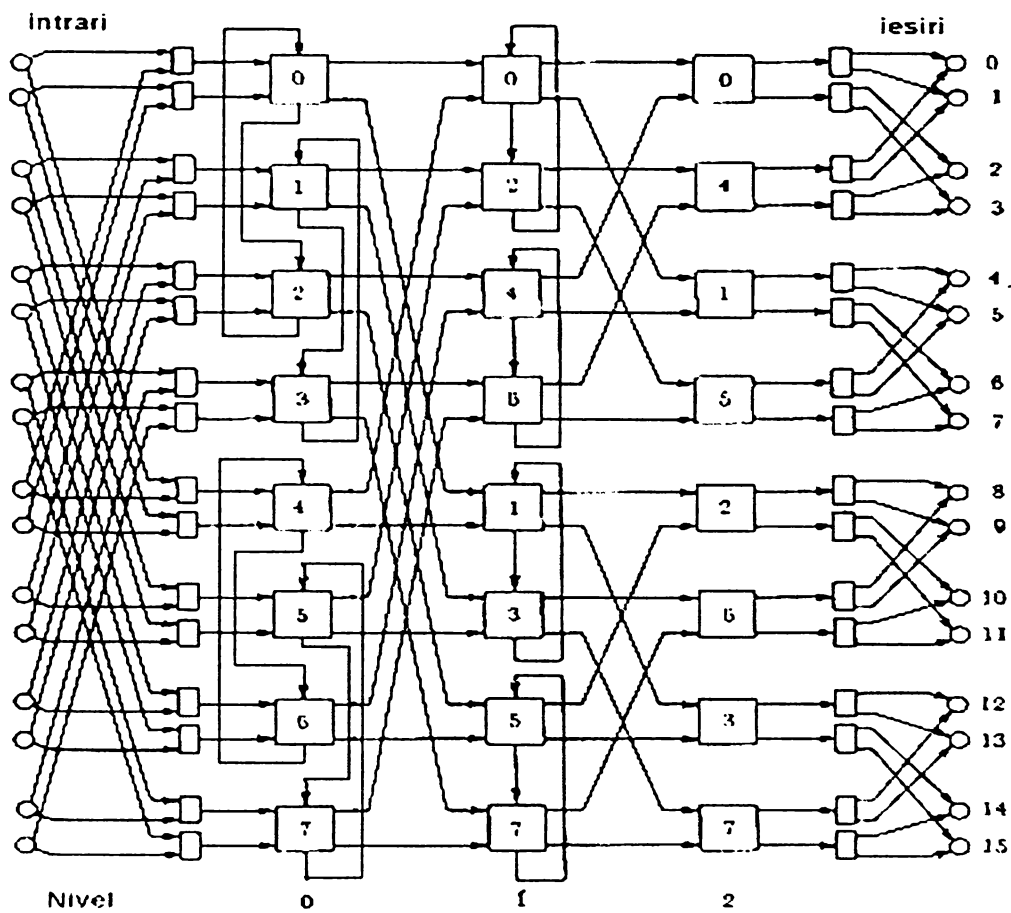


Fig. 4.40 Un ASEN-2 pentru $N=16$

Toate buclele într-un nivel dat trebuie să conțină un număr egal de comutatoare, ceea ce restrânge pe k la o putere a lui 2.

Grafurile de redundanță pentru ASEN-2 și ASEN-4 sunt reprezentate în fig. 4.41. De notat că buclele formate în toate nivelurile, cu excepția nivelului final, sunt astfel formate încât pentru fiecare buclă există altă buclă care este conectată la același set de comutatoare pe următorul nivel. Astfel de perechi de bucle sunt denumite *bucle conjugate*.

Rețelele ASEN au $\log N-1$ niveluri și $N/2$ comutatoare de talie $2 \times 3/\text{nivel}$, plus N MUX-uri 2×1 și N DMUX-uri 1×2 .

A. Controlul adaptiv al rutării la RIN – MN ASEN.

Siguranța în funcționare și ameliorarea performanței ce se obțin de la un RIN-MN, depind de cât de efectiv căile alternative sunt utilizate de către algoritmul de rutare. Algoritmii bazați pe *parcurgerea inversă a traseului (back-tracking)* sunt costisitori în termeni de resurse hardware și de timp afectat pentru stabilirea conexiunilor. Un algoritm ce nu se bazează pe parcurgerea inversă este ușor de implementat. El presupune să sursele și comutatoarele au abilitatea de a detecta defecte în comutatoarele la care sunt conectate.

Algoritmul pentru rutare în ASEN este descris în literatură [KR87]. El este următorul:

- **Pentru fiecare sursă:**

Transmite cererea de conectare comutatorului primar, împreună cu eticheta binară a destinației solicitate. Dacă comutatorul primar sau multiplexorul la intrarea comutatorului primar este defect, atunci transmite cererea comutatorului secundar. Dacă comutatorul secundar este de asemenea defect, abandonează cererea.

- **Pentru fiecare comutator de pe nivelul i ($i < n - 2$):**

Cererile pentru conectare pot sosi la oricare din cele trei conexiuni de intrare.

Pentru fiecare cerere, utilizează bitul i al etichetei destinație asociată și conectează la conexiunea de ieșire etichetează corespunzător. Dacă conexiunea de ieșire solicitată este ocupată sau nu poate fi utilizată pentru că există un defect în nivelul următor,

rutează cererea prin conexiunea auxiliară de ieșire spre următorul comutator din buclă. Dacă conexiunea auxiliară este de asemenea inutilizabilă deoarece este ocupată sau deoarece este defectă, atunci abandonează cererea. Un defect în DMUX la ieșirea unui comutator de pe nivelul $n - 2$ este privit ca și un defect în acel comutator.

- **Pentru fiecare comutator de pe nivelul $n - 2$:**

Cererile pot sosi la oricare din cele două conexiuni de intrare.

Pentru fiecare cerere, utilizează bitul $n - 2$ a etichetei de rutare și rutează cererea prin conexiunea de ieșire. Dacă ieșirea cerută este ocupată, abandonează cererea.

- **Pentru fiecare DMUX de la ieșirea nivelului $n - 2$:**

Poate primi maximum o cerere.

Dacă sosește o cerere, atunci utilizează bitul $n - 1$ a rutinei de rutare și fă o conectare spre conexiunea de ieșire superioară sau inferioară în concordanță cu bitul de direcționare 0 sau 1.

- **Pentru fiecare destinație:**

Pot sosi maxim două solicitări.

Acceptă cererile.

Acest algoritm este ușor încorporat în proiectare hardware a comutatoarelor, care este uniformă pentru toate nivelurile.

Algoritmul de rutare este independent de numărul de comutatoare dintr-o buclă.

B. Toleranța la defect și fiabilitatea rețelelor ASEN.

Un criteriu general de evaluare a robusteții RIN-MN se referă la faptul că fiecare membru a unui subset de surse trebuie să aibă căi către fiecare membru a unui subset de destinații. Probabilitatea ca și criteriul menționat anterior să fie satisfăcut se numește *fiabilitate de tip multiterminal (multiterminal reliability)*. Ne interesează două cazuri speciale ale acestui criteriu.

Primul caz este acela în care subseturile de surse și destinații conțin exact un element fiecare. Acest caz special conduce la o măsură denumită *fiabilitate de tip terminal*

(*terminal reliability*) care este probabilitatea ca o pereche sursă/destinație dată să aibă cel puțin o cale fără defect între ele, fiind cunoscut că fiecare comutator are o oarecare fiabilitate.

Fiabilitatea unui comutator este probabilitatea ca el să fie *fără defect (fault – free)*. Fiabilitatea de tip terminal a unei RIN –MN poate fi evaluată din grupul său de redundanță.

Al doilea caz special al criteriului de fiabilitate de tip multiterminal este acela al subseturilor de surse și destinații conținând toate sursele și toate destinațiile. Acest caz conduce la supoziția că RIN – MN este defect (sau deconectat) ori de câte ori toate căile sunt deconectate între unele perechi sursă / destinație și oferă *fiabilitatea $R(t)$* a RIN – MN. Funcția $R(t)$ poate conduce la calculul *MTTF – timpul mediu până la defectare*.

Deoarece algoritmul de rutare prezentat anterior pentru ASEN nu explorează toate căile disponibile, se introduce funcția ce reprezintă probabilitatea de succes a algoritmului de rutare – *fiabilitatea de tip terminal efectivă – ETR (effective terminal reliability)*. Ea este aceeași pentru toate perechile dacă defectele apar uniform printre comutatoare, deoarece rețelele ASEN sunt simetrice.

În fig. 4.42 se prezintă funcțiile *ETR* pentru fiecare din rețelele ASEN – 2, ASEN – 4 și ASEN – MAX, obținute analitic și verificate prin simulare Monte Carlo. Se presupune că probabilitatea unui element de comutare este $p = 0.9$

C. Performanțe ale rețelelor ASEN.

În literatură [KR 87] și [Prad 96] se prezintă câteva performanțe ale ASEN într-un mediu tipic comutării de circuite și se arată că ameliorarea față de RIN – MN *cu cale unică* este semnificativă. Ameliorarea se datorează utilizării de conexiuni în interiorul nivelului, care rezultă în lungimi de cale mai mari decât acelea din RIN – MN *cu cale unică*. Oricum creșterea lungimii nu este semnificativă.

Probabilitatea de acceptanță P_A , în acest caz, se referă la faptul că o cerere transmisă de o sursă este acceptată de o destinație fără să fie blocată de alte cereri sau conectări în RIN. P_A se evaluează presupunând că toate celelalte surse generează simultan cererile lor pentru conectare cu o probabilitate p și țintesc destinații alese aleator, la

începutul unui ciclu. Când două sau mai multe cereri sosesc la comutator solicitând aceeași conexiune de ieșire, cereile ce sunt prelucrate sunt alese pe o bază aleatoare iar celelalte sunt blocate sau abandonate.

În cazul ASEN-urilor, un număr de maximum trei cereri pot să sosească la un comutator. Se pot procesa cu același bit al etichetei de rutare până la două cereri, prin conectarea uneia la conexiunea de ieșire cerută spre nivelul următor și alta la conexiunea auxiliară de ieșire ce conduce la următorul comutator din buclă.

Destinațiile acceptă până la două cereri / ciclu, deoarece ele au două conexiuni la rețea. P_A este definită ca și raportul dintre numărul așteptat de cereri ce au reușit la numărul așteptat de cereri emise de către surse.

Rezultatul analizei P_A în rețelele ASEN sunt prezentate în fig. 4.42. S-a considerat probabilitatea de generare a cererilor ca fiind 1,0 respectiv 0,5. Se observă că la rata maximă de cereri generate, ASEN – 4 se comportă mai prost decât un *crossbar* dar mult mai bine decât RIN – MN *cu cale unică*. Interesant este faptul că, pentru valori mici ale lui N , dacă rata de generare a cereilor este în limite rezonabile scăzută, P_A este mai ridicată pentru ASEN decât pentru *crossbar*. Există două explicații: în primul rând, destinațiile în ASEN sunt echipate cu două porturi I/E și pot accepta până la două cereri în loc de una cât poate accepta un comutator *crossbar*; în al doilea rând, gradul de blocare, care este întotdeauna coborât pentru un N mic și un trafic scăzut, scade prin utilizarea conexiunilor auxiliare.

Prețul plătit în ASEN pentru un P_A mai mare este o mică creștere în lungimea căilor de comunicație. În fig. 4.44 se prezintă distribuția lui P_A în funcție de lungimea căii. Masa probabilității este concentrată în jurul valorilor mici ale lungimii de cale.

Prin combinarea avantajelor oferite de RIN – MN *omega cu o singură cale* cu facilitățile prezentate mai sus, rețelele ASEN devin foarte atractive pentru multiprocesoare și vor intra în circuitul industrial.

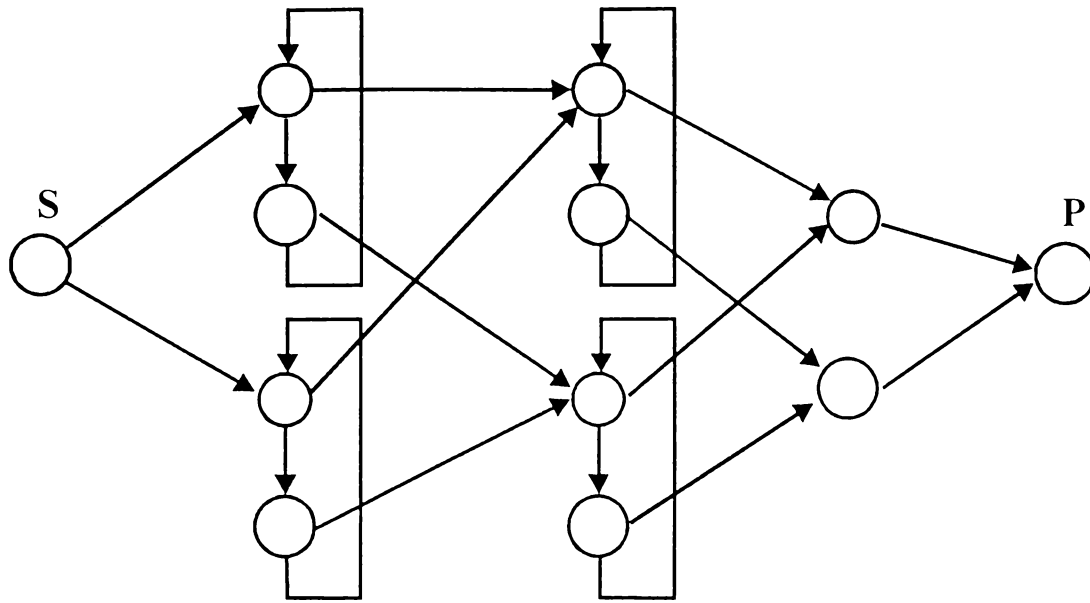
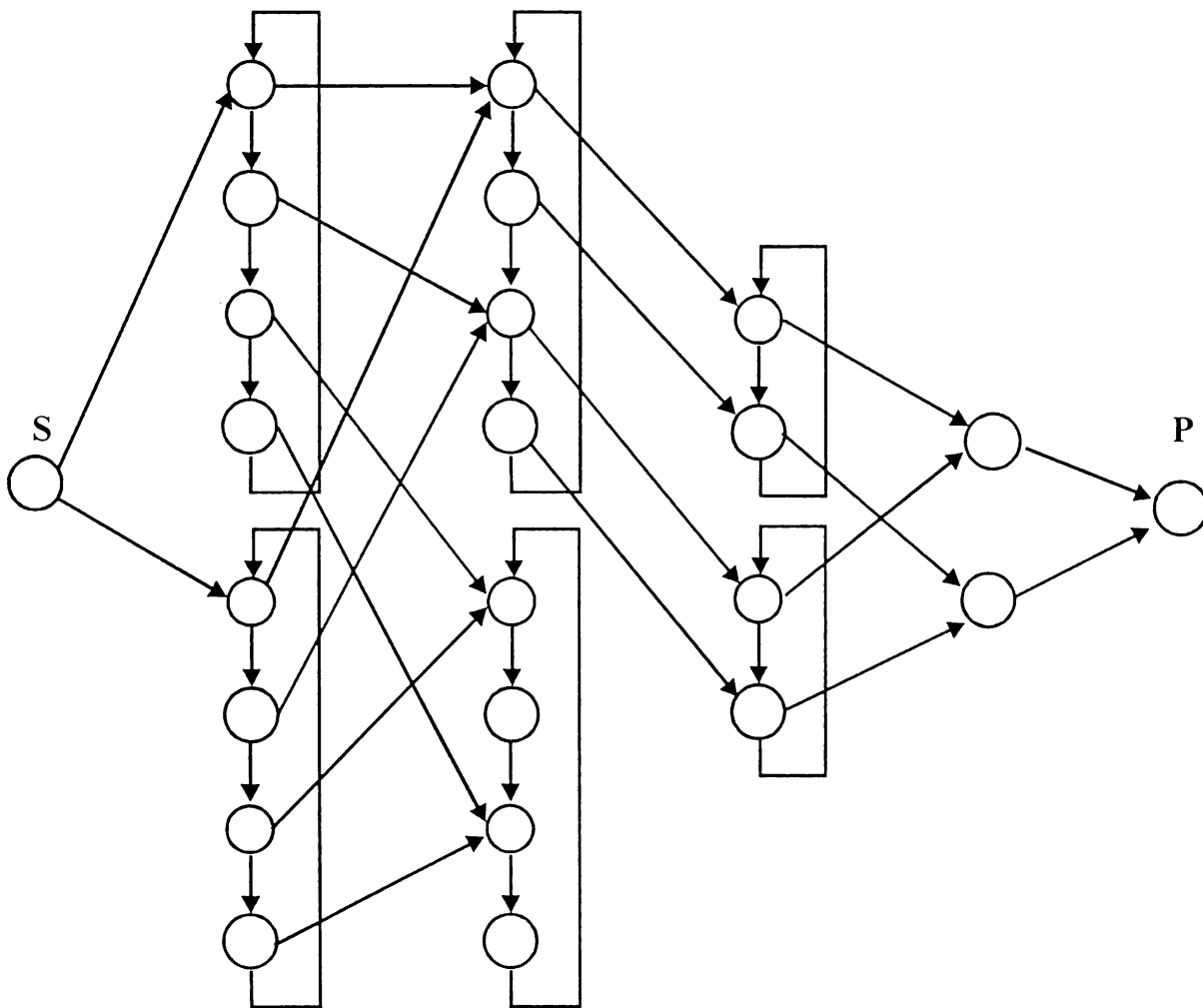
a) Graful de redundanță pentru *ASEN-2*b) Graful de redundanță pentru *ASEN-4*

Fig. 4.41 Grafurile de redundanță pentru RIN -MN ASEN.

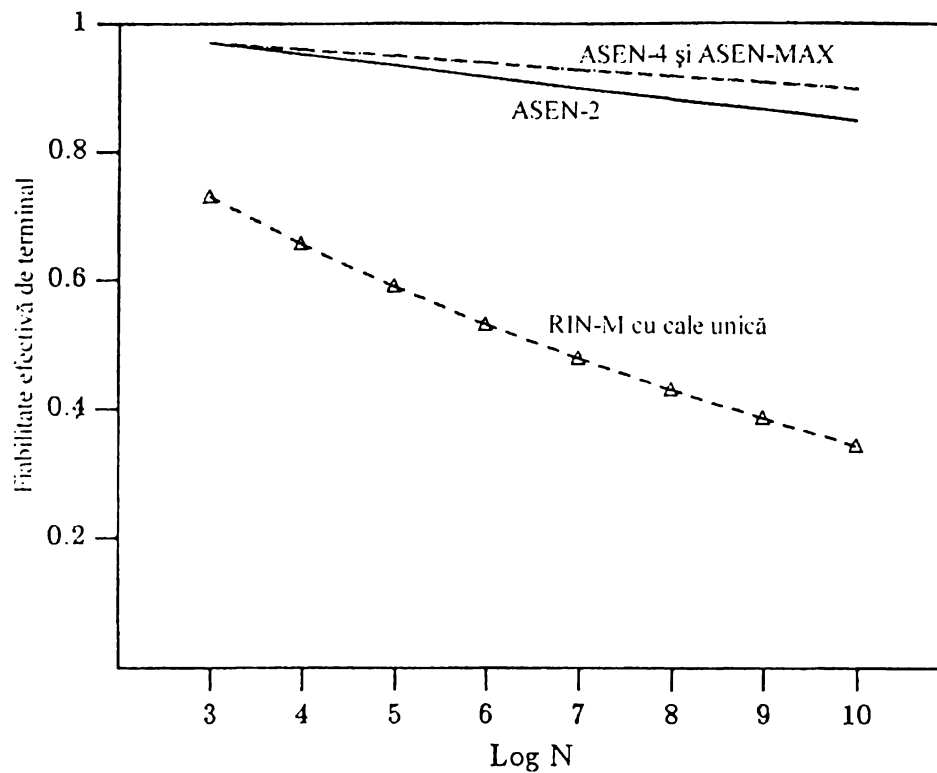


Fig. 4.42 Probabilitatea de terminal efectivă pentru ASEN și RIN-MN

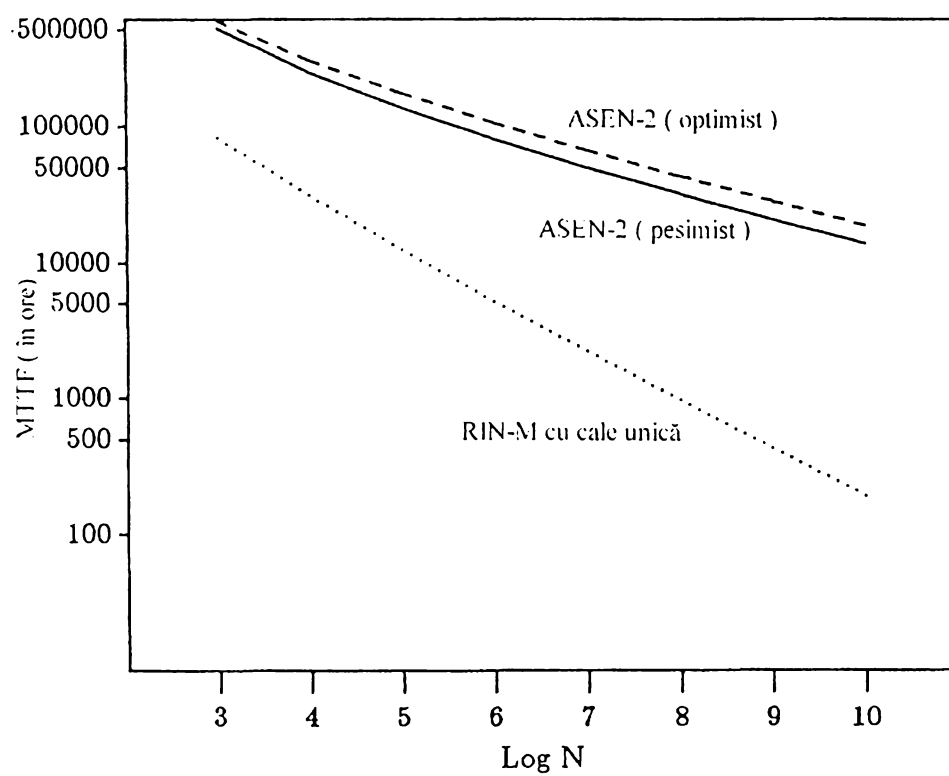


Fig. 4.43 Parametrul $MTTF$ pentru ASEN-2 și pentru RIN-MN cu *cale unică*, la o rată de defectare a comutatorului de 10^{-6} pe oră

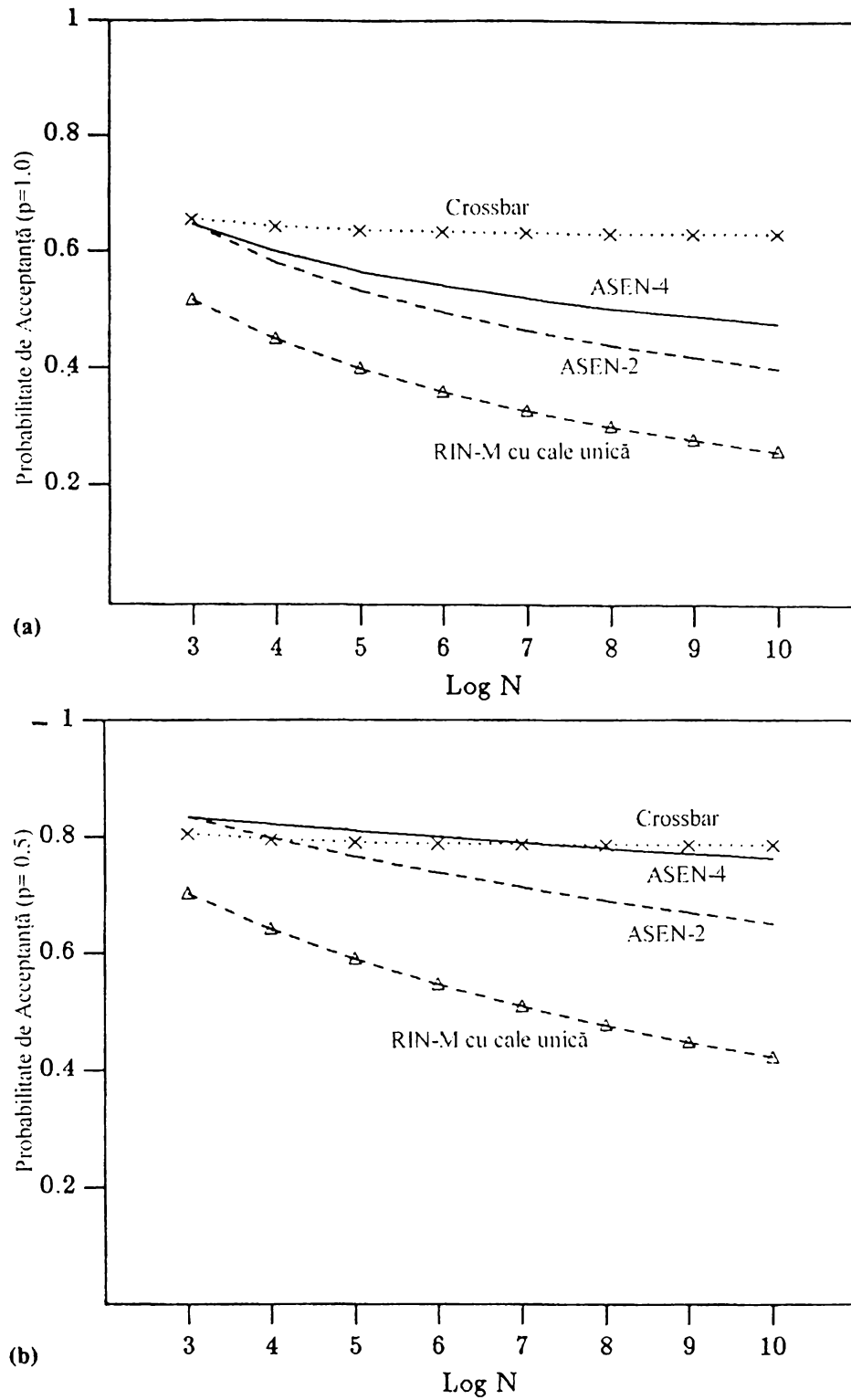


Fig. 4.44 Probabilitatea de acceptanță pentru rețele ASEN, pentru RIN-MN *cu cale unică* și pentru comutator *crossbar*. Probabilitatea de generare a unei erori este de 1.0 pentru (a) și de 0.5 pentru (b)

4.6 CONCLUZII

Capitolul 4 este dedicat analizei unor arhitecturi MIMD *tolerante la defect*. Se analizează mai multe abordări propuse pentru parcurgerea celor trei faze ale toleranței la defect dinamice: *detectarea defectului și localizarea lui, recuperarea din defect și reconfigurarea sistemului în jurul procesorului defect*. Dacă defectul este tranzitoriu, recuperarea se face fără reconfigurare. În cazul unui defect permanent, reconfigurarea sistemului devine esențială, evident, acesta suferind o degradare în performanță.

Referitor la detectarea defectului în SMP, se pot concluziona următoarele:

- detectarea prin *duplicare și comparare*, utilizând redundanța spațială, reprezintă o tehnică costisitoare deoarece un SMP cu P procesoare este configurat ca având $P/2$ procesoare. Acest lucru este valabil și pentru cele M module de memorie configurate în $M/2$ perechi. Un defect detectat într-o pereche de procesoare, dezactivează două procesoare, rămânând active $P-2$ procesoare configurate ca și $(P-2)/2$.
- *compararea* prin tehnici software este o operațiune consumatoare de timp dar se poate eficientiza prin utilizarea PV-urilor
- *detectarea* defectelor tranzitorii și a celor intermitente se produce utilizând tehnici de autodiagnoză a căror grad de complexitate crește odată cu finețea și acuratețea diagnosticării
- proiectarea și implementarea circuitelor de codificare și decodificare în cadrul circuitelor VLSI trebuie să asigure existența unor capacități interne de detectare a propriilor defecte.

Asigurarea toleranței la defect în cazul defectelor tranzitorii sau a celor intermitente se efectuează prin proceduri de recuperare ce asigură reexecutarea calculului. Există penalizări în performanță și de aceea procedurile de recuperare trebuie reduse atât ca număr cât și ca durată. De asemenea recuperarea din defect reprezintă pasul final după ce defectul din cadrul unui sistem a fost detectat și localizat iar sistemul reconfigurat în jurul procesorului defect.

Câteva concluzii se impun:

- tehnicile de recuperare sunt diferite pentru SMP–MP, SMP–TM sau SMP–MPD
- atât *rularea înapoi* cât și *rularea înainte*, utilizează tehnici de creare și de exploatare a PV-urilor. Ierarhia de memorie oferă soluții de amplasare atât pentru *datele active* cât și pentru *datele punctelor de verificare*.
- pentru a evita apariția unei dependențe între frecvența de verificare a PV-urilor și talia memoriei cache, se pot utiliza PV-uri *virtuale*. Ambele tipuri

de date – active și PV-uri – sunt conținute într-un sistem de memorii virtuale

- recuperarea din defect prin *rulare înapoi* în cazul multiprocesoarelor ce comunică între ele reprezintă un proces complex deoarece este necesară luarea în considerare a tuturor comunicațiilor și efectul lor până la cel din urmă PV
- crearea și verificarea PV-urilor în mod *sincron* evită apariția efectului *domino* și calculul poate fi repornit după ultimul PV global. Recuperarea este foarte simplă și timpul de recuperare este redus. Dezavantajul major îl constituie o suprasarcină semnificativă atât la crearea și verificarea PV-urilor cât și în procesul de comunicare
- crearea și verificarea PV-urilor în mod *asincron* are cea mai redusă sarcină dar pentru a putea fi efectiv utilizată în practică trebuie să fie combinată cu tehnica de rulare înapoi bazată pe evidența operațiunilor de rulare (log – based rollback). În schemele recente se preferă implementarea variantei sincrone în loc de cea asincronă. Sarcina de coordonare a creerii și exploatării PV-urilor în cazul acestor sisteme este neglijabilă în comparație cu sarcina de salvare a stărilor. Astfel sarcina în cazul variantei sincrone este aproximativ aceeași ca și în cazul variantei asincrone dar beneficiul în cazul sincron îl reprezintă o schemă mult mai simplă
- crearea și verificarea PV-urilor în mod *quasi-sincron* combină ușurința și sarcina redusă a variantei asincrone cu avantajele timpului redus de recuperare de la varianta sincronă. Reprezintă o abordare atractivă în sistemele moderne ca de exemplu în sistemele de calcul mobile.

Referitor la proiectarea și utilizarea RIN-MN *tolerante la defect* sunt necesare câteva concluzii:

- RIN-MN *tolerante la defect* au caracteristici ce includ *modelul de defect*, *criteriul de toleranță la defect*, *metoda de creare a toleranței la defect*, *toleranța la defecte singulare sau multiple*, *complexitatea rutării și complexitatea hardware*. Modelul de defect și criteriul de toleranță la defect sunt metrici esențiale pentru orice comparație a capabilităților de toleranță la defect conținute de rețelele studiate.
- codurile corectoare de erori, ce conectează subsisteme redundante în paralel, și ce împart un subsistem în părți identice conectate în paralel, reprezintă tehnici consacrate pentru obținerea toleranței la defect

- tehnicile de obținere a toleranței la defect, independente de topologia rețelei, sunt:
 - transmisia multiplă cu aceeași informație pe căi distincte
 - parcurgerea de mai multe ori a RIN-MN
- tehnicile ce modifică topologia unei RIN-MN pentru a obține o mai bună toleranță la defect prevăzute în faza de proiectare, sunt:
 - adăugarea unui nivel de comutatoare suplimentar cu sau fără conexiuni by-pass
 - adăugarea unor legături suplimentare
 - creșterea numărului de porturi
 - modificarea dinamică a dimensiunii comutatorului și a numărului de niveluri de comutatoare
 - augmentarea dimensiunii comutatorului prin adăugarea unor conexiuni corespunzătoare
 - replicarea *in integrum* a unei RIN-MN

În literatură [AAS 87] se descriu pe larg și alte tehnici. În ultimă instanță importantă este stabilirea unui raport corect performanță/preț în ceea ce privește gradul de acoperire la defect.

- RIN-MN se compară alegând un model de defect, valabil pentru toate rețelele, și un criteriu de toleranță la defect bazat pe o reprezentare realistă a defectelor hardware ce pot apare. Modelul comun de defect presupune că orice componentă a rețelei ar putea să se defecteze și ca acest defect să facă rețeaua inutilizabilă. Criteriul de toleranță la defect, specifică că o RIN-MN tolerează un defect doar dacă mai poate încă conecta orice intrare la oricare ieșire.

Doar câteva dintre RIN-MN sunt tolerante la defect singular în condițiile respectării tuturor cerințelor impuse anterior. Dacă se relaxează modelul de defect comun și se presupune o perfectă siguranță în funcționare a nivelurilor de comutatoare, de la intrare și de la ieșire, atunci vor exista mai multe tipuri de RIN-MN ce pot fi considerate tolerante la defect singular. Sunt necesare însă resurse hardware suplimentare ca și circuite de evitare a comutatorului defect sau includerea unor conexiuni multiple între fiecare intrare a rețelei și portul de ieșire, respectiv comutatoarele din rețea. Aria de cercetare a RIN-MN *tolerante la defect* este departe de a fi epuizată iar abordarea noilor tehnologii hardware și software în contextul mai larg al dependibilității vor asigura ameliorarea semnificativă a performanței acestora.

5. MODELAREA DEPENDABILITĂȚII PENTRU SISTEMELE MULTIPROCESOR CU MEMORIE PARTAJATĂ

SMP-MP trebuie să ofere un nivel ridicat de performanță și, în plus, să fie *dependabile* în prezența defectelor. Proiectarea SMP-MP în condițiile obținerii unei puteri de calcul ridicate și a unei dependabilități superioare, este foarte complexă. Prin modelare se pot identifica cerințele unei bune dependabilități, identifica sursele de limitare a performanțelor și, prin urmare, selecta o arhitectură optimă.

În cap. 4 s-au prezentat atributele unui SMP tolerant la defect.

Reamintim că *fiabilitatea* unui sistem este o funcție $R(t)$ definită ca fiind probabilitatea condițională ca sistemul să funcționeze corect în intervalul $[t_0, t]$, dat fiind faptul că a fost operațional la t_0 .

Disponibilitatea la momentul t este o funcție $A(t)$ ce reprezintă probabilitatea ca sistemul să fie operațional la momentul t . Evaluarea lui $A(t)$, comparativ cu cea a lui $R(t)$, nu necesită o operare continuă a sistemului în intervalul $[t_0, t]$, deoarece un sistem poate traversa un ciclu *defectare-reparare* (*fail-repair cycle*) de un număr de ori în acest ciclu.

În literatură [DKT90], [BRL94], [BT93] se indică patru modele operaționale pentru analiza dependabilității. Fiecare model definește cerințe minimale pentru ca un sistem să fie considerat operațional.

La proiectarea SMP trebuie luate în considerare măsuri de asigurare a *dependabilității raportată la performanță* (*performance-related dependability*) în plus față de cele de asigurare doar a dependabilității [DKT90], [RA90], [RG87].

Modelarea cu acuratețe a RIN este extrem de complexă și, de aceea, diferite topologii de RIN pot necesita tehnici de modelare diferite.

5.1 ETAPELE PRELIMINARE ÎN DESCRIEREA MODELULUI

Un model este o reprezentare abstractă a sistemului și crearea lui presupune trei faze: *definirea*, *parametrizarea* și *evaluarea*.

Definirea

Se definește sistemul și obiectivul modelării. Acesta poate fi: *fiabilitatea*, *disponibilitatea* sau ambele în funcție de *performanță*. Identificarea sistemului ca fiind operațional se face în termeni de conectivitate disponibilă între noduri, conectate prin RIN.

Conectivitatea minimă pentru un sistem operațional se poate baza pe:

- o conexiune între o pereche anume de noduri de intrare și noduri de ieșire (dependabilitate de tip terminal)
- o conexiune între un set cunoscut de noduri de intrare și de ieșire (dependabilitate de tip multiterminal)
- conexiuni între un procentaj dat de noduri într-un sistem (dependabilitate bazată pe activitate)
- conexiuni între toate nodurile de intrare și de ieșire (dependabilitate de tip rețea)

Parametrizarea

Parametrii de intrare includ *tipuri de componente, caracteristici de defectare* ca și *distribuția defectului, rata de defectare și de recuperare și caracteristici de reparare* ca și *distribuția timpului de reparare, rata de reparare și probabilitatea unei reparări cu succes.*

Rata de defectare este λ_i , iar cea de reparare μ_i .

Au o distribuție exponențială cu rate constante.

Parametri de ieșire sunt: *fiabilitatea, disponibilitatea stării stabile sau tranzitorii, MTTF, factorul de ameliorare a siguranței în funcționare.* De asemenea, se pot obține măsuri ale dependabilității în raport cu performanța ca și *performanța așteptată la timpul t sau performanța acumulată într-un interval de timp.*

Evaluarea

Evaluarea modelului se face utilizând fie tehnici analitice, fie simularea.

Tehnicile analitice sunt:

- diagrame bloc ale fiabilității
- arbori de defect
- modele Markov
- rețele petri

În literatură [DKT90] sunt descrise în amănunt aceste tehnici, dar cea mai puternică tehnică pentru analiza de dependabilitate este *modelarea Markov.*

În cadrul *modelelor Markov* comportarea sistemului este modelată ca și un lanț Markov cu stări finite și în spațiul timpului continuu.

5.2 MODELE DE FIABILITATE

Există patru tipuri de tehnici de evaluare a fiabilității: *de tip terminal*, *de tip multiterminal*, *bazată pe sarcină* și *de tip rețea* [DKT90].

5.2.1 Fiabilitatea de tip *terminal*

Reprezintă probabilitatea ce există cel puțin o cale între două terminale specifice. Căderea apare când un terminal se defectează sau când nu există nici o cale între două terminale.

Un terminal poate fi oricare dintre o gamă de entități.

De exemplu, fiecare terminal ar putea fi o pereche procesor-memorie, în care caz fiabilitatea de tip terminal ar putea determina fiabilitatea a două procese ce comunică. În mod alternativ, un terminal ar putea fi un procesor și celălalt un modul de memorie.

Pentru evaluarea fiabilității de tip *terminal*, în primul rând se stabilește siguranța componentelor individuale și se enumeră conectivitatea lor relativă. Defectele de legătură între noduri sunt ignorate. Ele pot fi înglobate în defecțiunile de comutator sau de nod.

5.2.1.a Sisteme orientate pe *crossbar*

Fie SMP-MP orientat pe *crossbar* de dimensiuni $M \times N$.

Elementul de procesare PE_i se conectează la modulul de memorie MM_j prin intermediul comutatorului SW_{ij} .

Fiabilitatea terminalului este:

$$TR(t) = R_p(t) \cdot R_m(t) \cdot R_{sw}(t)$$

unde:

$R_p(t)$ – fiabilitatea procesorului

$R_m(t)$ – fiabilitatea memoriei

$R_{sw}(t)$ – fiabilitatea comutatorului

În cazul unei distribuții exponențiale cu ratele de defectare λ_p , respectiv λ_m se poate scrie:

$$R_p(t) = e^{-\lambda_p t}; \quad R_m(t) = e^{-\lambda_m t}; \quad R_{sw}(t) = e^{-\lambda_{sw} t}$$

5.2.1.b Sisteme bazate pe magistrală

Un sistem cu B magistrale asigură comunicarea dintre orice element de procesare și orice modul de memorie atât timp cât cel puțin o magistrală este disponibilă.

Dacă fiabilitatea fiecărei magistrale este:

$$R_b(t) = e^{-\lambda_b t}$$

atunci

$$TR(t) = R_p(t) \cdot R_m(t) \cdot [1 - (1 - R_b(t))^B]$$

Ultimul termen reprezintă probabilitatea ca cel puțin o magistrală să fie operațională la timpul t .

5.2.1.c Sisteme orientate pe RIN-MN

Fie o RIN-MN cu o *singură cale* cu $n = \log_a N$ etaje de comutatoare.

Atunci fiabilitatea de terminal este:

$$TR(t) = R_p(t) \cdot R_m(t) \cdot R_{SW}^n(t)$$

În cazul RIN-MN cu *cale multiplă* pot exista una sau mai multe căi ce conectează un procesor la o memorie.

Analiza unui astfel de sistem se face utilizând o diagramă bloc a fiabilității pentru a reprezenta sistemul ca și o structură serie-paralel.

În ipoteza că într-un RIN-MN cu *cale multiplă* analizăm starea a două comutatoare A și B în raport cu o anumită cale, atunci pot exista două posibilități:

Ambele A și B sunt necesare pentru o conectare validă.

Comutatoarele sunt modelate ca și un sistem serie cu:

$$R_s(t) = R_A(t) \cdot R_B(t)$$

Dacă fiecare comutator poate asigura o conectare validă atunci ele constituie un sistem în paralel cu:

$$R_p(t) = [1 - (1 - R_A(t)) \cdot (1 - R_B(t))]$$

Metoda se complică în cazul sistemelor mari.

5.2.1.d Sisteme *hipercub*

Analiza sistemelor *hipercub* este similară cu cea a sistemelor orientate pe RIN-MN. Analiza unui n -cub cu 2^n noduri se face enumerând căile și prezentându-le ca și o diagramă bloc. Există $n!$ căi alternative non disjuncte de la fiecare sursă la fiecare destinație.

În fig. 5.1. se prezintă diagrama bloc a fiabilității pentru un *hipercub* unde se analizează fiabilitatea de tip terminal între nodurile 0 și 7.

Există șase căi posibile non disjuncte în nodurile 0 și 7.

În ipoteza că toate nodurile sunt la fel de sigure, $R_n(t) = R_i(t)$ pentru $0 \leq i \leq 7$ atunci

$$TR(t) = 6 \cdot R_n^1(t) - 6 \cdot R_n^5(t) - 3 \cdot R_n^6(t) + 6 \cdot R_n^7(t) - 2 \cdot R_n^8(t)$$

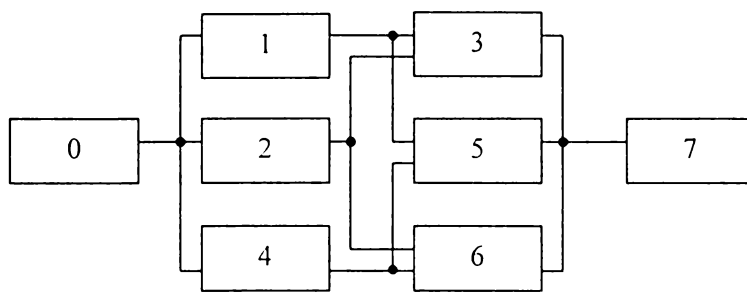


Fig. 5.1. O diagramă bloc a fiabilității reprezentând diverse căi între nodul 0 și nodul 7

5.2.2 Fiabilitatea de tip *multiterminal*

Reprezintă o extensie logică a fiabilității de tip terminal.

În locul unei perechi specifice de terminale se ia în considerare un subset specific de terminale. Prin urmare, ea reprezintă probabilitatea ca un subset dat de terminale să fie complet conectat pentru comunicații simultane. De exemplu, fiabilitatea de tip *k-terminal*, reprezintă probabilitatea ca k noduri ($k > 1$) să poată comunica prin RIN.

5.2.2.a Sisteme orientate pe *crossbar*

Un SMP *crossbar* de talie $M \times N$ necesită $x \cdot y$ comutatoare *crossbar* operaționale pentru a conecta x procesoare la y memorii. Fiabilitatea este:

$$MTR(t) = R_p^x(t) \cdot R_m^y(t) \cdot R_{SMP}^{xy}(t) \text{ pentru } 1 \leq x \leq M \text{ și } 1 \leq y \leq N$$

5.2.2.b Sisteme orientate pe magistrală

Pentru un SMP $M \times N \times B$ fiabilitatea de tip *multiterminal* este aceeași ca și cea de tip *monoterminal*, deoarece la un moment dat este nevoie de cel puțin o magistrală. Dacă pentru o conectare validă sunt necesare un set de x procesoare și y memorii, atunci

$$MTR(t) = R_p^x(t) \cdot R_m^y(t) \cdot [1 - (1 - R_b(t))^B]$$

5.2.2.c Sisteme orientate pe RIN-MN

Fiabilitatea de tip *multiterminal* a unui sistem orientat pe RIN-MN *cu o singură cale* se raportează la S , numărul de comutatoare necesar pentru conectarea a x elemente de procesare la y memorii. Acest S depinde de localizarea procesoarelor solicitate, cât și a memoriilor în cadrul RIN-MN. O dată ce S este cunoscut, fiabilitatea multiterminal devine:

$$MTR(t) = R_p^x(t) \cdot R_m^y(t) \cdot R_{SMP}^S(t)$$

Analiza lui MTR în cazul RIN-MN *cu cale multiplă* este complexă și necesită un studiu mai aprofundat.

5.2.2.d Sisteme *hipercub*

Un sistem este reprezentat ca și un graf probabilistic $G(N, E)$ unde N și E sunt seturile de noduri, respectiv de conexiuni.

Fiecărei componente i îi este asigurată o fiabilitate. Algoritmii de stabilire a fiabilității de tip *multiterminal* utilizează un *identificator de cale* de n bit pentru a reprezenta fiecare cale dintre o sursă și destinație, unde n este numărul de componente ce se pot defecta.

Prin urmare, selecționând $n=N$ se consideră doar defecțiunile de nod, iar $n=E$ doar cele de conexiune. În situația $n=N+E$ se consideră ambele.

Un bit din identificatorul de cale este 1 dacă el este solicitat pentru conexiune și x în cealaltă. Toți identificatorii de cale pentru conexiunile de intrare și de ieșire sunt apoi prelucrați, utilizând tehnici din algebra booleană pentru a găsi o expresie simbolică finală, care este convertită într-o formă probabilistică pentru a determina fiabilitatea de tip *multiterminal*.

5.2.3 Fiabilitate bazată pe sarcină

Modelele de fiabilitate bazată pe sarcină (task-based) sunt mai generale decât fiabilitatea de terminal sau multiterminal, deoarece nu se specifică explicit nodurile sursă și destinație. Un sistem rămâne operațional atât de mult timp cât sarcina poate fi executată cu resursele disponibile. Prin urmare, dacă o sarcină necesită cel puțin I elemente de procesare și J module de memorie într-un SMP-MP, atunci sistemul este considerat operațional atâta timp cât sunt disponibile aceste resurse. Fiabilitatea bazată pe sarcină permite să se amelioreze toleranța la defect prin încorporarea *degradării grațioase (graceful degradation)* în sistem.

5.2.3.a Sisteme orientate pe *crossbar*

Probabilitatea ca i din x componente să funcționeze la momentul t este:

$${}_x C_i(t) = \binom{n}{i} R^i(t) \cdot (1 - R(t))^{n-i} \cdot C^{x-1}$$

unde:

$R(t)$ – fiabilitatea componentei $M \times N$

C – factorul de acoperire

Utilizând această notație, fiabilitatea unui SMP-MP *crossbar* de talie $M \times N$, de exemplu $C.mmp$ ce necesită cel puțin I elemente de procesare și J module de memorie este:

$$TBR(t) = \sum_{i=1}^M \sum_{j=1}^N {}_M C_i(t) \cdot {}_N C_j(t) \cdot R_{x-bar}(t) \quad (5.1.)$$

$R_{x-bar}(t)$ – fiabilitatea unei matrici *crossbar* de 16×16 din sistemul C.mmp.

$M = N = 16$ pentru sistemul C.mmp

În acest model nu este luată în calcul degradarea performanței comutatorului. Pentru a lua în calcul acest deziderat se construiește un nou model astfel: se presupune că se dorește un grup conectat de i elemente de procesare și j module de memorie într-un sistem *crossbar* de $M \times N$. Deoarece un comutator este necesar pentru a conecta fiecare memorie, sunt necesare $i \cdot j$ elemente de comutare pentru conectarea cerută.

Fie $\gamma = i \cdot j$ și $\varepsilon = M \cdot N - \gamma$, comutatoarele neutilizate. Există două modalități [DKT90] prin care un sistem *crossbar* $M \times N$ ar putea asigura conexiunea $i \times j$ solicitată:

- când exact i elemente de procesare și j memorii sunt funcționale, și elementele de comutare γ solicitate sunt disponibile să asigure conexiunile cerute. În acest prim caz (1) probabilitatea este $P_{i,j(1)}(t)$.
- Când mai mult decât i procesoare sau j memorii sunt funcționale, dar doar γ comutatoare sunt disponibile pentru a realiza o conexiune $i \times j$. În acest caz (2) probabilitatea este $P_{i,j(2)}(t)$

Se poate scrie:

$$P_{i,j(1)}(t) = \prod_{x=1}^M C_x(t) \cdot \prod_{y=1}^N C_y(t) \cdot R_{SW}^\gamma(t) \quad (5.2.a.)$$

De notat că ε – condiția de comutatoare suplimentare – este irelevantă când acoperirea C este perfectă în ceea ce privește comutatoarele. Acoperirea imperfectă în comutatoare se face aproximând al doilea termen $P_{i,j(2)}(t)$ ca fiind:

$$P_{i,j(2)}(t) = \sum_x^M \sum_y^N \prod_{x \neq x'} C_{x'}(t) \cdot \prod_{y \neq y'} C_{y'}(t) \cdot R_{SW}^\gamma(t) \cdot (1 - R_{SW}(t))^\varepsilon \quad (5.2.b.)$$

unde:

$$x > i \quad y > j$$

Combinând cei doi termeni se definește $TBR(t)$ astfel:

$$TBR(t) = \sum_{i=1}^M \sum_{j=1}^N P_{i,j(1)}(t) + P_{i,j(2)}(t) \quad (5.2.c.)$$

5.2.3.b Sisteme orientate pe magistrală

La aceste sisteme, analiza este simplă deoarece există întotdeauna o cale între orice procesor și orice memorie, atâta timp cât cel puțin o magistrală este disponibilă. Contribuția structurii de

magistrală este $1 - (1 - R_b(t))^b$. Intruducând acest termen în ecuația (5.1) în locul lui $R_{x-bar}(t)$ rezultă:

$$TBR(t) = \sum_{i=1}^M \sum_{j=1}^N C_i(t) \cdot C_j(t) \cdot [1 - (1 - R_b(t))^b] \quad (5.2.d.)$$

ce reprezintă fiabilitatea bazată pe sarcină pentru SMP-MP cu magistrală multiplă.

5.2.3.c Sisteme orientate pe RIN-MN

Analiza fiabilității bazată pe sarcină pentru această clasă de sisteme este extrem de dificilă, în principal datorită complexității RIN-MN. Un model Markov cu un singur nivel (incluzând procesoare, memorii, comutatoare) este aproape imposibil de construit datorită numărului mare de stări. De aceea, se utilizează tehnici ca și *simularea* și *descompunerea de sistem* în analiza fiabilității bazată pe sarcină.

Descompunerea ierarhică este o tehnică specială pentru analiza sistemelor orientate pe RIN-MN. Este descrisă în amănunt în [DKT90], [DTB90] și [DTB93]. Această schemă presupune descompunerea unui sistem mai mare în subsisteme mai mici pentru o analiză mai ușoară.

În literatură [DKT90] se citează exemplul unui sistem bazat pe o rețea *Butterfly* ce utilizează comutatoare *crossbar* 4×4 ca și în fig. 5.2.

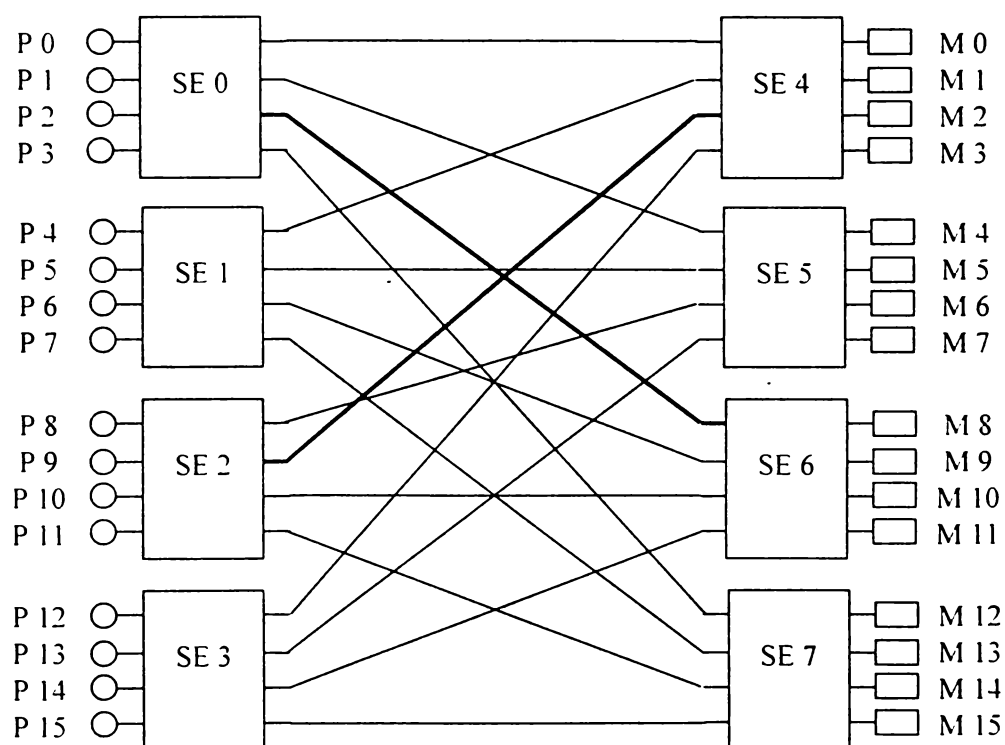


Fig. 5.2. Un SMP-MP cu RIN-MP *butterfly* de 16×16

Numărul de noduri din sistem este 4^n și fiecare nod constă dintr-un procesor și o memorie. Este un sistem $4^n \times 4^n$. Configurația $4^n \times 4^n$ poate fi partiționată pe 4 subsisteme $4^{n-1} \times 4^{n-1}$, fără a

perturba conexiunile. Prin urmare, fiabilitatea unui SMP-MP de rang 4^n depinde de 4 subsisteme 4^{n-1} și de configurația de conectare dintre ele. Analiza pornește cu expresiile probabilităților pentru o configurație de patru elemente de procesare, patru memorii și două comutatoare.

Fie $P_4(i,j)^{(t)}$ probabilitatea ca i elemente de procesare și j module de memorie ($0 \leq i, j \leq 4$) să lucreze la momentul t . Analiza are drept scop să găsească probabilitatea ca I elemente de procesare și J module de memorie să lucreze într-un SMP-MP de talie 16×16 . Aceasta se exprimă folosind patru grupe de subsisteme 4×4 astfel:

$$P_{16(i,j)}(t) = P_{4(i_0,j_0)}(t) \cdot P_{4(i_1,j_1)}(t) \cdot P_{4(i_2,j_2)}(t) \cdot P_{4(i_3,j_3)}(t)$$

$$\text{Astfel ca: } (i_0 + i_1 + i_2 + i_3) = I \text{ și } (j_0 + j_1 + j_2 + j_3) = J$$

Toate combinațiile posibile ale celor I procesoare și ale celor J memorii sunt generate pentru expresiile de mai sus.

În final, fiabilitatea sistemului de 16 noduri cu cel puțin I procesoare și J memorii active este:

$$TBR(t) = \sum_{i=I}^{16} \sum_{j=J}^{16} P_{16(i,j)}(t) \quad (5.3)$$

Fiabilitatea unui sistem cu 64 de noduri se obține din configurații de 16 noduri, utilizând tehnica de mai sus. Un sistem cu 64 noduri are trei niveluri cu câte 16 comutatoare pe nivel.

Primul și ultimul nivel sunt incluse în analiza de bază 4×4 . Prin urmare, se poate determina numărul cerut de comutatoare din nivelul mediu, pentru orice grup activ utilizând un calcul de aproximare. Fiabilitatea acestor comutatoare este multiplicată cu $P_{64(i,j)}(t)$ pentru a garanta că comutatoarele din nivelul median asigură $i \times j$ conexiuni pentru cele două cazuri posibile discutate la analiza sistemelor crossbar.

5.2.3.d Sisteme hiper cub

Modelul de fiabilitate bazat pe sarcină dezvoltat este utilizat pentru a găsi probabilitatea ca un n -cub să aibă cel puțin I noduri active conectate. Se utilizează tehnica descompunerii și se poate lucra cu orice nivel de degradare a sistemului.

Deoarece un n -cub poate fi divizat în două $(n-1)$ -cuburi, această schemă descompune recursiv un cub mai mare în cuburi mai mici, până când dimensiunea cubului obținut este suficient de simplă pentru model. Cel mai mic cub este denumit „modelul de bază” al analizei.

Un 2-cub (patru noduri) sau un 3-cub (opt noduri) pot constitui modelul de bază. Apoi, în mod recursiv, se obține fiabilitatea unor cuburi de dimensiuni mai mari.

De exemplu, pentru a calcula probabilitatea a j noduri conectate într-un n -cub, îl partajăm în două $(n-1)$ -cuburi. Există două cazuri în care un n -cub va avea j noduri conectate.

În primul caz, exact j noduri sunt active în n -cub.

În al doilea caz, mai mult de j noduri, de exemplu x , sunt active, dar conectivitatea este doar j .

Se poate diviza j sau n noduri între două $(n-1)$ -cuburi ca și k și $j-k$ sau $x-k$, respectiv. Cinci distribuții ale nodurilor active pot fi analizate pornind de la aceste două cazuri pentru a obține $P_j(t)$, probabilitatea ca să existe j noduri active conectate în două $(n-1)$ -cuburi.

Fiabilitatea n -cubului este:

$$TBR(t) = \sum_{j=1}^{2^n} P_j(t) \quad (5.4)$$

Această abordare poate fi extinsă pentru a găsi un subcub de dimensiune m , unde $m \leq n$, într-un n -cub, presupunând că un program în sistem necesită o alocare cubică în loc de l noduri conectate.

5.2.4 Fiabilitatea de tip rețea

Fiabilitatea de tip rețea (network reliability) este probabilitatea de-a menține o rețea cu acces total. Dacă se consideră două rețele cu o singură cale, probabilitatea rețelei reprezintă probabilitatea ca toate procesoarele și memoriile sistemului să fie conectate. O rețea *cu cale multiplă* poate avea acces total chiar dacă componentele rețelei se defectează.

Analiza RIN *cu cale multiplă* este extrem de dificilă. În ultimă instanță, acest tip de analiză reprezintă un caz special a fiabilității bazate pe sarcină, unde sistemul nu permite degradarea resurselor de calcul.

5.3 MODELE DE DISPONIBILITATE

În cadrul evaluării fiabilității nu se iau în considerare activitățile de reparare a sistemului. În sistemele cu degradare grațioasă, repararea componentelor ce s-au defectat contribuie la creșterea duratei operaționale a sistemului și oferă o disponibilitate mai mare.

Evaluarea disponibilității unui sistem complex cu o singură facilitate de reparare este mult mai complexă decât modelarea de fiabilitate. Această complexitate derivă din luarea în considerare a doi factori.

În primul rând, trebuie generat un lanț Markov al sistemului pentru a analiza disponibilitatea sa. Acest proces poate fi complex pentru un sistem mare. Numărul de stări din model poate, de asemenea, să fie prea mare pentru a fi ușor manipulat.

În al doilea rând, găsirea unei soluții pentru disponibilitatea în cazul stărilor tranzitorii sau a celor stabile este dificilă. Din această cauză, se utilizează tehnici numerice pentru a calcula disponibilitatea.

5.3.1 Sisteme orientate pe *crossbar*

În literatură [DKT90] se prezintă un lanț Markov dezvoltat pentru un sistem *crossbar* (C.mmp) de dimensiuni 16×16 . Întreaga rețea este tratată ca și un singur element. Modelul este simplu deoarece o defecțiune în *crossbar* se transformă într-o defecțiune de sistem. Modelul este un lanț Markov bidimensional, unde o dimensiune reprezintă defecțiunea procesorului, iar a doua reprezintă defecțiunea de memorie.

5.3.2 Sisteme orientate pe magistrală

SMP-MP cu magistrală multiplă se modelează ca și un lanț Markov tri-dimensional, unde fiecare dimensiune reprezintă degradarea unui tip de componentă - procesor, memorie și magistrală. Presupunând că procesoarele, memoriile și magistralele au fiecare facilități de reparare, lanțul Markov poate include rețele de reparare pentru a ilustra un sistem reparabil.

5.3.3 Sisteme orientate pe RIN-MN

Modelul de disponibilitate pentru SMP-MP orientate pe RIN-MN propus în [DTB90] poate manipula sisteme mari și orice nivel de degradare. Această schemă utilizează o abordare a descompunerii sistemului pentru a evita generarea unui lanț Markov mononivel detaliat. Acest model presupune o comportare independentă la defecțiune și la reparare a procesoarelor, memoriilor și comutatoarelor și, din punct de vedere conceptual, împarte sistemul în trei subsisteme conectate în serie. Se pot dezvolta lanțuri Markov independente pentru cele trei subsisteme.

Fiecare diagramă Markov a subsistemului este un simplu lanț unidimensional. În fig. 5.3. se prezintă un lanț Markov pentru un subsistem procesor ce necesită cel puțin I elemente de procesare, λ_p și μ_p reprezentând *rata de cădere*, respectiv *rata de reparare ale procesorului*.

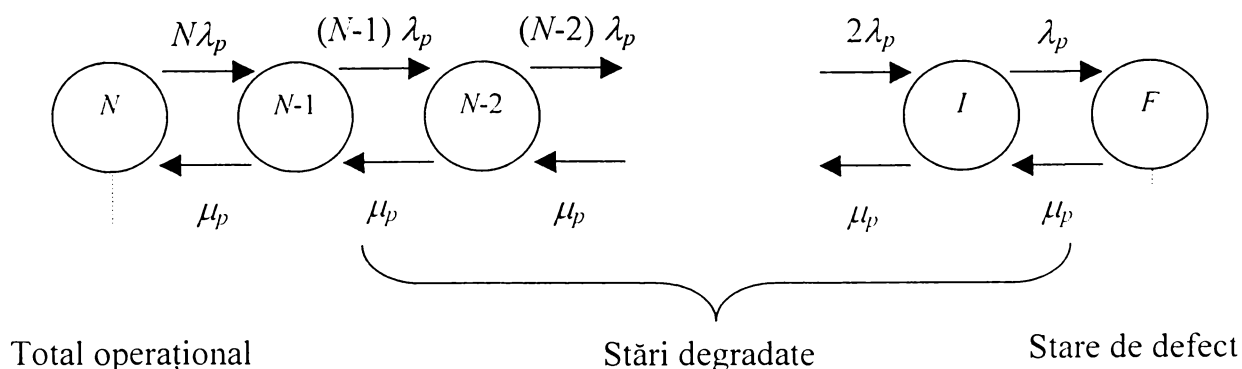


Fig. 5.3. Un lanț Markov pentru un subsistem procesor

Sistemul evoluează de la starea i la starea vecină $i-1$ cu o rată de cădere $i \cdot \lambda_p$ sau la starea $i+1$ cu o rată de reparare μ_p . Din rațiuni de simplificare a analizei se omit acoperirea imperfectă și parametrii de reparare. Un astfel de model poate fi rezolvat prin tehnici de analiză numerică pentru a găsi probabilitatea a i procesoare active $P_i(t)$ pentru $1 \leq i \leq N$.

Modelul Markov pentru subsistemul de memorie este aproape identic cu excepția că sistemul tolerează până la J defectări. *Ratele de defectare ale memoriei și ratele de reparare* sunt λ_m și μ_m . Rezolvarea acestui model oferă $P_j(t)$, probabilitatea ca să existe j module de memorie active la orice timp t .

Lanțul Markov pentru elementele de comutare se bazează pe identificarea numărului de comutatoare necesare pentru o conexiune $i \times j$. De exemplu, calculul probabilității ca opt procesoare și opt memorii să fie conectate într-un sistem 16×16 ca în fig. 5.2, se obține parționând sisteme în patru grupe, fiecare grup cuprinzând patru procesoare și patru memorii. Se poate distribui cele opt procesoare și opt memorii în două, trei sau patru grupuri.

Există două cazuri ce ar putea apare într-o conexiune 8×8 . În primul caz, exact opt procesoare și opt memorii sunt active și sunt conectate printr-un număr necesar de comutatoare. Celelalte comutatoare nu afectează conectivitatea. De exemplu, dacă componentele active sunt distribuite în două grupuri, trebuie să funcționeze patru comutatoare pentru ca această conexiune să fie validă. Aceste două comutatoare aparțin la două grupuri de lucru și pot fi separate de celelalte patru comutatoare. Conectivitatea este 8×8 , indiferent de care alte patru comutatoare sunt active.

Lanțul Markov din fig. 5.4 prezintă acest caz, cu λ_{SE} și μ_{SE} reprezentând *ratele de defectare*, respectiv *ratele de reparare ale comutatoarelor*.

În fig. 5.4, (x, y) reprezintă comutatoarele (x) cerute și comutatoarele suplimentare (y) , astfel că $x+y$ este numărul total de comutatoare din RIN-MN. Tranzițiile orizontale arată că sistemul cade (starea F_i , unde $0 \leq i \leq y$) când oricare dintre cele x comutatoare solicitate cade. În mod consecvent, fiabilitatea RIN-MN este:

$$P_{SE(i)}(t) = \sum_{j=0}^x P_{x,j}(t)$$

Unde (i) indică primul caz, adică cazul în care exact i procesoare și j memorii sunt active. În fig. 5.4 cu $x=y=4$, $P_{SE(i)}(t) = (P_{4,4}(t) + P_{4,3}(t) + P_{4,2}(t) + P_{4,1}(t) + P_{4,0}(t))$.

Probabilitatea conexiunii $i \times j$ pentru crossbar este:

$$P_{i,j(i)}(t) = P_i(t) \cdot P_j(t) \cdot P_{SE(i)}(t)$$

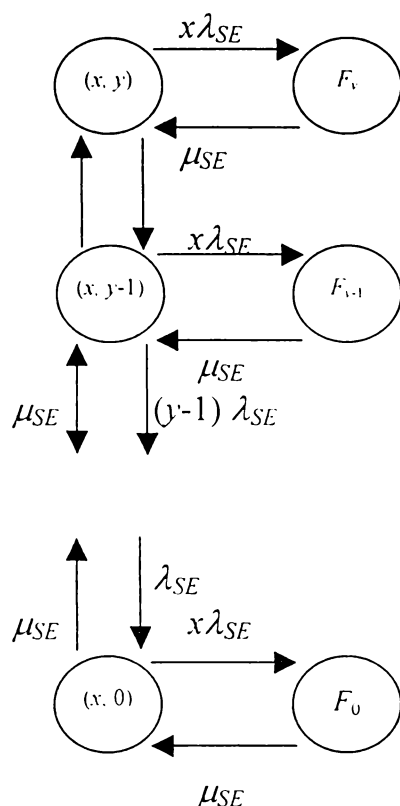


Fig. 5.4. Un lanț Markov reprezentând o RIN-MN

În al doilea caz, mai mult decât i procesoare, fie ele α , sau mai mult decât j memorii, fie ele β , pot lucra în timp ce conectivitatea rămâne $i \times j$. De exemplu, mai mult decât opt procesoare și de opt memorii pot lucra în fig. 5.2, cu conectivitatea rămânând 8×8 , atâta timp cât doar comutatoarele cerute pentru opt procesoare și opt memorii sunt active. Dacă conexiunea 8×8 este distribuită în două grupuri, sunt necesare patru comutatoare. Celelalte patru comutatoare ar putea cădea (făcând $y=0$) astfel că procesoarele suplimentare și memoriile din restul sistemului să nu contribuie la fiabilitate. Prin urmare, $P_{SE(2)}(t) = P_{x,0}(t)$ reprezintă probabilitatea ca doar x comutatoare să fie active. Se obține $P_{i,j(2)}(t)$ prin multiplicarea lui $P_{x,0}(t)$ cu probabilitățile ca α procesoare și β memorii să fie active. Aceste probabilități sunt $P_\alpha(t)$ și $P_\beta(t)$. Valorile pentru $P_\alpha(t)$ și $P_\beta(t)$ se obțin rezolvând lanțurile Markov pentru procesor și memorie. Se determină disponibilitatea sistemului adăugând toate probabilitățile configurației active exprimate de $TBR(t)$.

5.4 MODELE DE DEPENDABILITATE RAPORTATĂ LA PERFORMANȚĂ

Măsura dependabilității raportată la performanță combină aspecte ale performanței și ale dependabilității. Este utilă analiza din această perspectivă, mai ales pentru sistemele cu degradare grațioasă.

Utilizând notația unui model lanț Markov în spațiu continuu, un SMP generic se poate reprezenta ca și un proces cu stări finite $Z(t)$, cu $t \geq 0$ și cu spațiul de stări $0, 1, \dots, n$. Starea 0 reprezintă starea defectă a sistemului, iar stările 1 până la n , configurațiile de lucru. Lanțul Markov ar putea fi *ciclic* sau *aciclic*, depinzând de faptul că sistemul este *reparabil* sau *nu*.

Rezolvând lanțul Markov se poate obține un vector coloană $P(t) = p_1(t); p_2(t), \dots, p_n(t)$ de unde se obține probabilitatea ca sistemul să fie în diferite stări de lucru la momentul t . De exemplu, $p_i(t)$ este probabilitatea ca sistemul să fie în starea i la momentul t .

Se asociază stării i o *rată de recompensă* r_i pentru $1 \leq i \leq n$. Dacă sistemul rămâne în starea i pentru o durată τ , atunci recompensa acumulată în starea i este τr_i . *Funcția de recompensă acumulată* $Y(t)$ în timpul t este:

$$Y(t) = \int_0^t x(\tau) \cdot d\tau$$

unde $x(\tau)$ este rata de recompensă la momentul τ .

La rândul său, $x(\tau) = r Z(\tau)$ unde r este un *vector de recompensă* asociat. Se poate defini r cu parametrii ca și: *lățime de bandă*, *rată de execuție a instrucției* sau alte metrici de performanță relevante.

Studiile asupra acestui tip de dependabilitate au produs un număr interesant de măsuri. Astfel, pentru SMP degradabile, în literatură [DKT90] se citează: *fiabilitate a calculului* și *disponibilitate a calculului*. Utilizând conceptul de *probabilitate a unei stări de lucru* $p_i(t)$ cât și a *ratei de recompensă asociate* r_i , se poate exprima *puterea de calcul la timpul t* ca fiind:

$$E[x(t)] = \sum_{i=1}^n p_i(t) \cdot r_i \quad (5.5)$$

În analiza dependabilității raportată la performanța SMP-MP se introduce un parametru denumit *disponibilitatea lățimii de bandă (band-width availability)*, $BA(t)$. Ea este definită ca fiind lățimea de bandă așteptată a unui SMP la momentul t . Lățimea de bandă este o măsură a performanței pentru sistemele sincrone, definit ca și numărul mediu de perechi procesor-memorie conectate într-un ciclu.

Utilizând o expresie de fiabilitate bazată pe sarcină, *disponibilitatea lățimii de bandă* este:

$$BA(t) = \sum_{i=1}^M \sum_{j=J}^N P_{i,j}(t) \cdot BW_{i,j} \quad (5.6)$$

unde :

$P_{i,j}(t)$ – probabilitatea ca sistemul să lucreze ca i procesoare și j memorii

$BW_{i,j}$ – lățimea de bandă cu factorul de recompensă r_i

Performanța acumulată la momentul t este $E[Y(t)]$ care reprezintă cantitatea totală de activitate utilă, exprimată ca și valoare medie, pe care un SMP o poate efectua integrând ecuația 5.5 pe intervalul de timp observat:

$$E[Y(t)] = \sum_{i=1}^n r_i \cdot \int_0^t p_i(\tau) \cdot d\tau \quad (5.7)$$

Un alt parametru introdus este *performabilitatea* ce reprezintă funcția de distribuție a probabilității performanței de sistem acumulată. Ea oferă probabilitatea ca un sistem să efectueze o cantitate dată de activități x în timpul t . Utilizând funcția $Y(t)$, performabilitatea este:

$$y(x,t) = P[Y(t) \geq x] \quad (5.8)$$

Calculul expresiei este complex. Pentru un sistem cu n stări active, fiecare distribuție a recompensei stării implică o integrare.

Evaluarea performabilității la SMP a fost limitată deoarece modelele de dependabilitate nu au fost suficient dezvoltate pentru unele circuite și determinarea distribuției este complexă pentru un număr mare de stări.

Cele mai multe dintre SMP-MP *crossbar*, *magistrală multiplă*, *RIN-MN* utilizează lățimea de bandă BW ca și indicator de performanță.

Compararea dependabilității la SMP de talie mică se face luând în considerare conexiunile *crossbar*, *magistrala partajată* și *RIN-MN*.

Fig. 5.5 prezintă fiabilitatea sistemului pentru un SMP-MP de dimensiuni 16×16 în trei ipostaze derivate din tipurile de conexiuni. Se presupune că au aceeași rată de defectare. Toate cele trei sisteme operează aproximativ cu aceeași fiabilitate când nu se permite degradarea.

Oricum, ele au în mod distinct diferite fiabilități, dacă se permite o degradare de 25% ($I=J=12$). SMP-MP cu *magistrală partajată* operează cel mai bine deoarece are opt căi alternative, cele mai multe, decât celelalte sisteme. Diferența în fiabilitate între *crossbar* și *RIN-MN* devine notabilă când există degradare de sistem datorită diferitelor capacități de conectare. Un defect de comutator în *RIN Butterfly* deconectează patru perechi procesor-memorie, în timp ce în *crossbar* doar o pereche este deconectată. Pe de altă parte, complexitatea *RIN-MN* este mai mică deoarece există mai puține comutatoare. Prin urmare, se poate ameliora fiabilitatea *RIN-MN* adăugând niveluri suplimentare de comutatoare.

Cuantificarea *dependabilității raportată la performanță* este prezentată în fig. 5.6 iar în fig. 5.7 se prezintă *distribuția mediată în timp a lățimii de bandă* pentru un sistem de tip C.mmp.

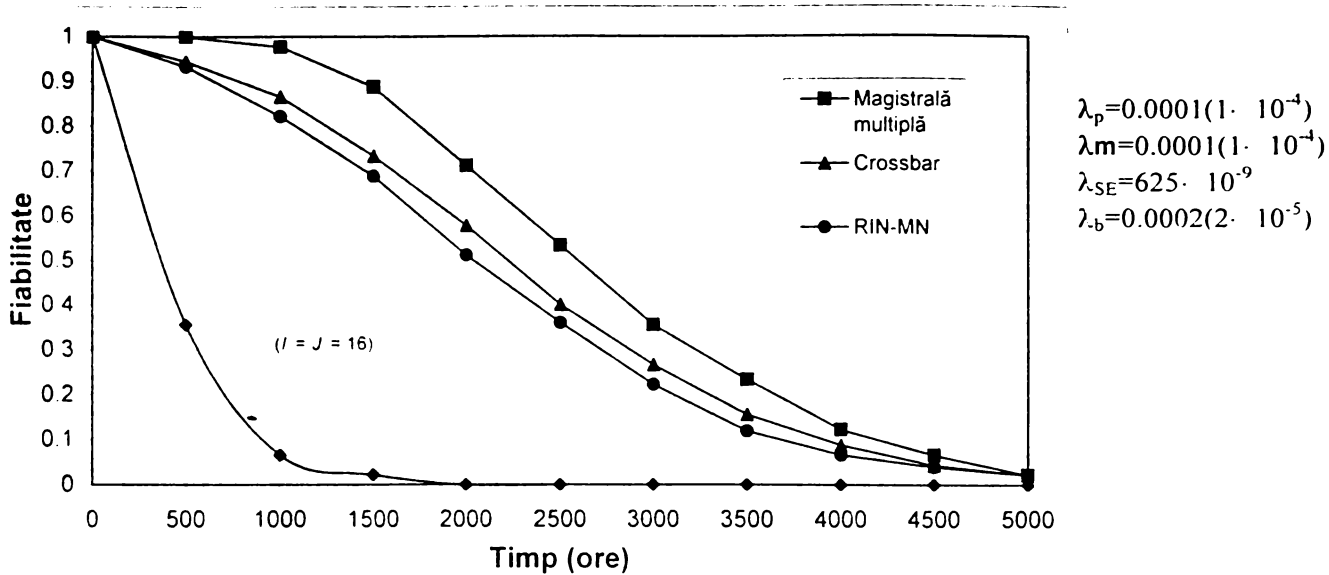


Fig. 5.5 Fiabilitatea unui SMP - MP în funcție de natura RIN

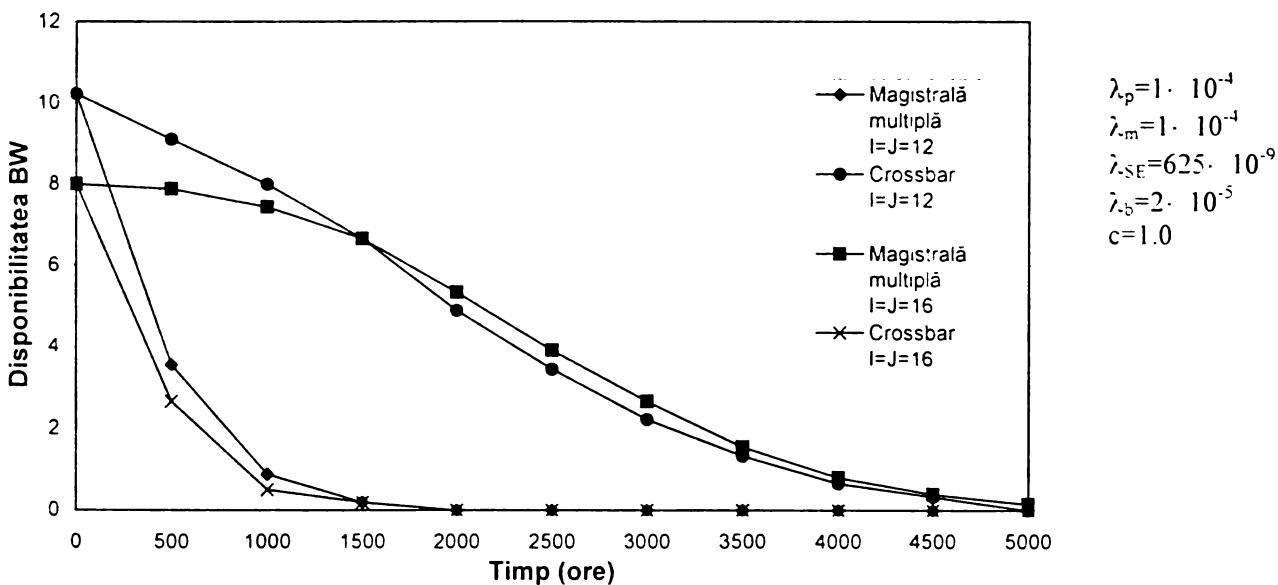


Fig. 5.6 Disponibilitatea unui SMP - MP în funcție de natura RIN

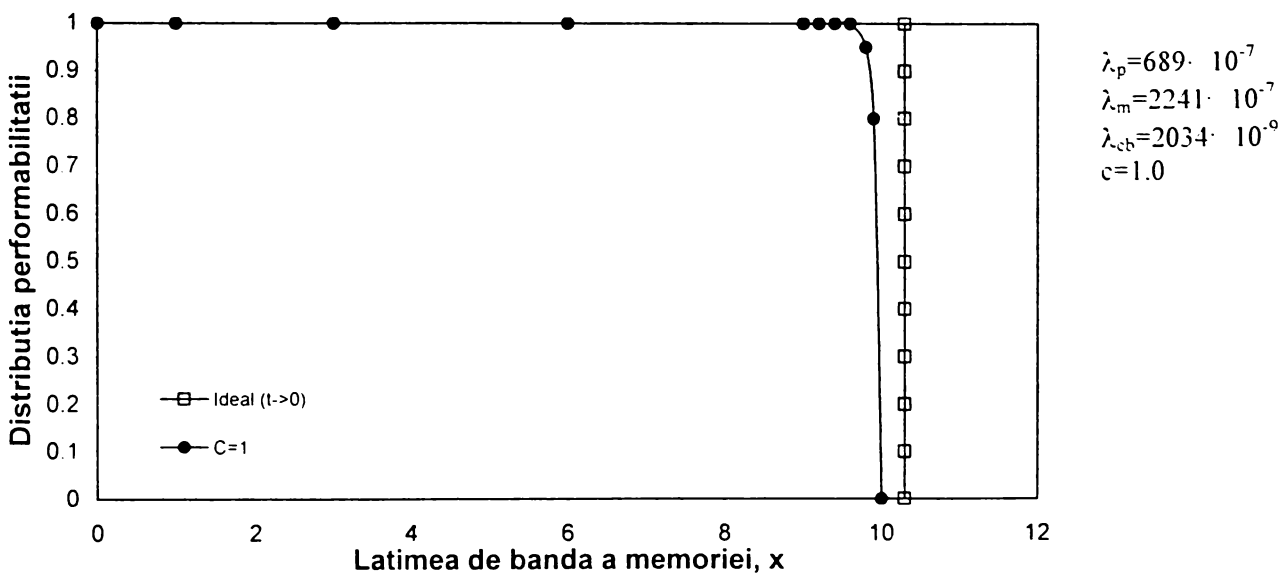


Fig. 5.7 Distribuția mediată a performabilității

5.5 MODELAREA DISPONIBILITĂȚII UTILIZÂND DESCOMPUNEREA DE SISTEM PENTRU SMP-MP, ORIENTATE PE RIN-MN

În paragraful § 5.3.3. s-a prezentat succint modul de abordare a construcției unui model de dependabilitate pentru un SMP-MP orientat pe RIN-MN.

În acest paragraf vom detalia *descompunere de sistem* ca și tehnică de modelare a dependabilității unui sistem complex fără a construi un *lanț Markov mononivel (single-level Markov Chain - MK)*.

În literatură [DTB 93] se prezintă metoda *descompunerii de sistem* în cadrul clasei de modele de disponibilitate bazate pe sarcină. În cadrul unor astfel de modele, un sistem este considerat operațional atâta timp cât cerințele sarcinii sunt îndeplinite. Modelele celor două simple MK ce pot fi dezvoltate - subsistemul procesor-memorie și subsistemul RIN-MN - se rezolvă apoi prin pachete de programe specializate. În literatură [DKT 90], [DTB 93] sunt citate pachete de programe **Hybrid Automated Reliability Predictor- HARP**, **Computer-Aided Reliability Estimation III- Care III** și **Symbolic Hierarchical Automated Reliability and Performance Evaluation- SHARPE**.

În cadrul modelului se calculează probabilitățile ca i elemente de procesare PE și j module de memorie MM, $P_i(t)$ respectiv $P_j(t)$, să fie operaționale în orice moment t prin găsirea numărului de comutatoare cerute pentru conectarea celor i procesoare în cele j memorii.

5.5.1 Descrierea sistemului

Sistemul analizat este un SMP-MP de talie $N \times N$ în care N elemente de procesare PE sunt conectate cu N module de memorie MM, prin intermediul unei RIN-MN de tip *Butterfly* având $n(n = \log_4 N)$ niveluri de comutatoare-SE- de dimensiuni 4×4 . Pe un nivel sunt plasate $N/4$ SE.

Particularizarea unui astfel de sistem pentru $N=16$ este prezentată în fig. 5.2.

Se presupune că un element de procesare PE_i și modulul său de memorie corespunzător MM_i sunt localizate pe o singură placă denumită *nodul i*. Prin urmare, sistemul are o configurație de N noduri. O referire la memoria locală nu trebuie să parcurgă RIN-MN. O referire la o memorie externă modulului necesită două treceri prin RIN-MN. În prima trecere, procesorul solicitant transmite cererea către memoria externă modulului prin RIN-MN. Răspunsul este transmis de memorie către procesor la a doua trecere. De exemplu, o parcurgere completă a circuitului de la PE_2 la MM_8 este reprezentată prin linie groasă în fig. 5.2. În acest tip de schemă de comunicare

sunt necesare patru SE într-un sistem (16×16) în raport cu două SE necesare în cadrul unui protocol cu o singură trecere.

În construcția modelului se fac următoarele ipoteze simplificatoare:

1. Componentele subsistemului sunt toate independente, omogene și au distribuții de defectare exponențiale.

Acestea sunt $\lambda_p, \lambda_m, \lambda_{se}$ și reprezintă *ratele de defectare* ale PE, MM și SE.

2. Există o singură facilitate de reparare pentru fiecare tip de componentă cu o distribuție exponențială a timpului de reparare.

Ratele de reparare sunt μ_p, μ_m, μ_{se} .

Este interesant de remarcat că se admite repararea on-line și a comutatoarelor SE. Această abordare presupune existența unor algoritmi de detectare a defectului pentru RIN-MN care pot localiza defectul în comutatoare. Izolarea și remedierea defectului sunt posibile fără oprirea activității. De asemenea, se admite că toleranța la defect se manifestă și prin degradarea grațioasă. Datorită complexității RIN-MN, pentru a se putea gestiona defectele din SE, acestea sunt prevăzute cu memorii cache incluse.

5.5.2 Tehnica de modelare

Analiza se bazează pe *metoda descompunerii* în care SMP este divizat în trei subsisteme ca și în fig. 5.8. Fiecare sistem este analizat independent. La un moment dat pot fi conectate I procesoare și J memorii pentru satisfacerea cerințelor impuse de sarcina la nivel de sistem.

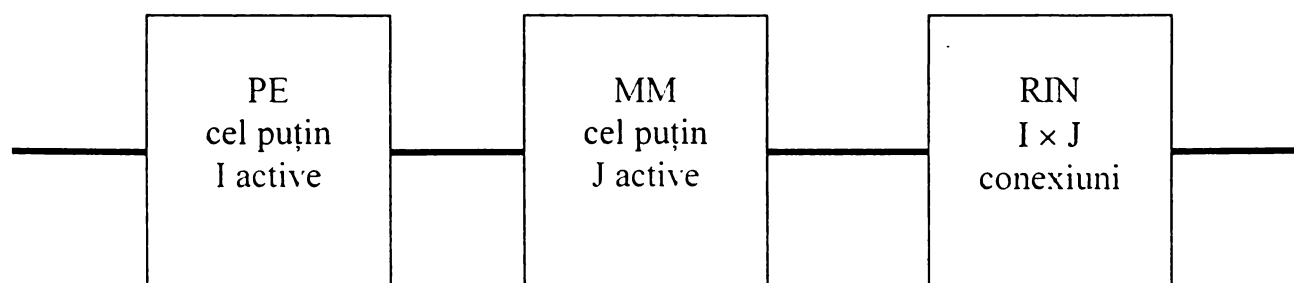


Fig. 5.8. Descompunerea unui sistem

A. Subsistemul Procesor/Memorie (SPM)

Configurația inițială a SMP are N procesoare și cel puțin I dintre acestea trebuie să fie active pentru ca sistemul să fie operațional. MK-ul pentru SPM este prezentat în fig. 5.9.

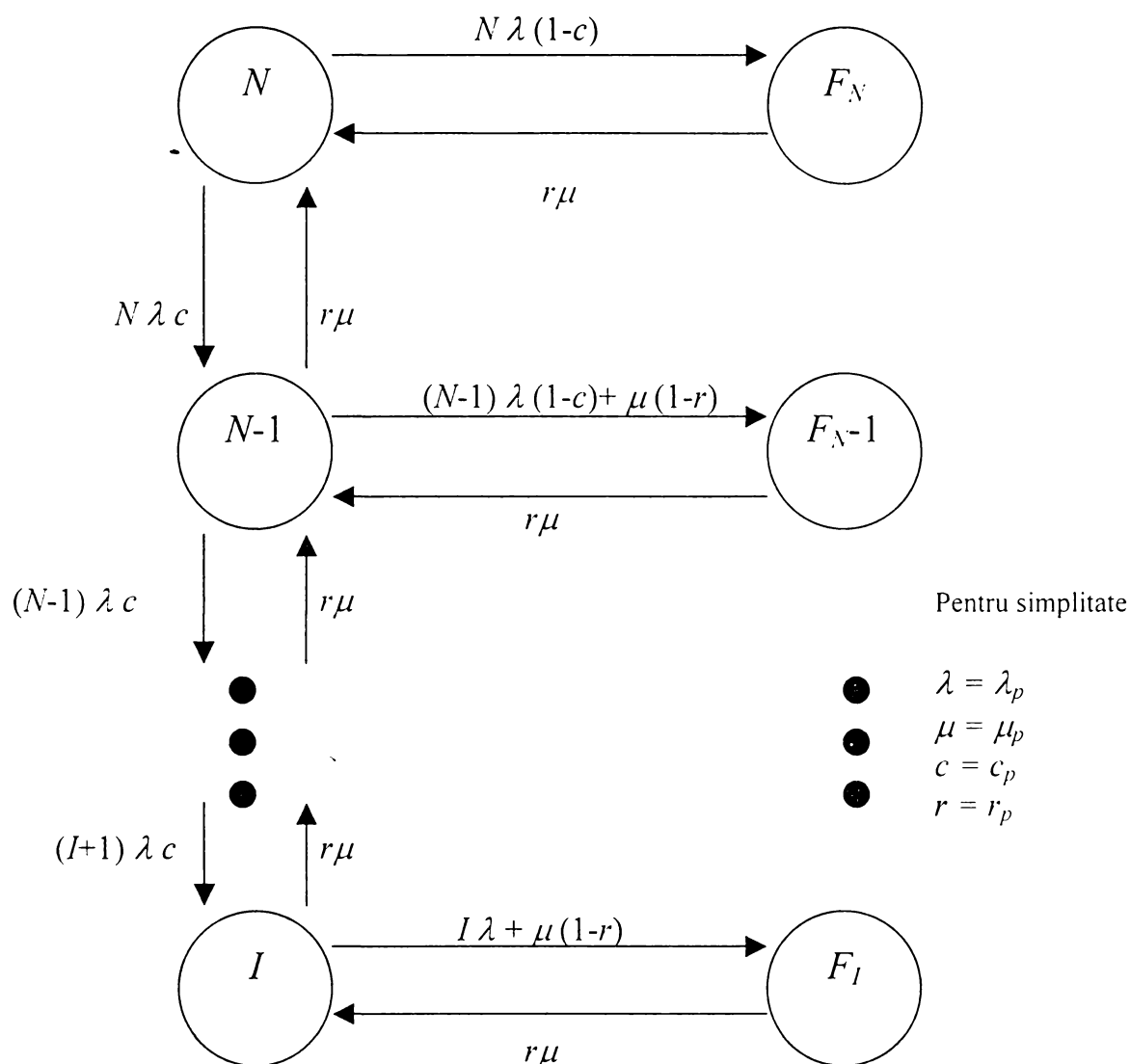


Fig. 5.9. Lanț Markov pentru subsistemul SPM

Modelul include o acoperire imperfectă și o reparare imprecisă presupuse ca putând apare.

Notând:

C_p –acoperirea unui procesor

R_p –factorul de reparare cu succes

se observă că SPM evoluează de la o stare de lucru i , pentru $I \leq i \leq N$ spre o stare defectă F_i fie unei acoperiri imperfecte $(1-c_p)$ sau datorită unei reparări imprecise $(1-r_p)$. Dacă acoperirea este integrală și procesul de reparare este perfect, atunci MK se reduce la un ciclu de reparare.

Soluția pentru MK prezentat se obține utilizând un pachet HARP ce furnizează probabilitatea ca i procesoare să fie active în sistem la momentul t . Aceasta este $P_i(t)$ pentru $1 \leq i \leq N$.

MK pentru subsistemul memorie este asemănător cu cel procesor, cu mențiunea că sistemul tolerează până la $N - J$ defecte. Rata de defectare, rata de reparare, factorul de acoperire și factorul de reparare cu succes sunt $\lambda_m, \mu_m, c_m, r_m$. Rezolvând MK-ul memoriei, se obține $P_j(t)$, probabilitatea ca să existe j memorii active la momentul t , pentru $1 \leq j \leq N$.

B. Subsistemul RIN-MN

O RIN-MN degradabilă poate asigura conectarea dintre i PE și j MM. Astfel, modelarea unui RIN-MN poate garanta că numărul solicitat de SE trebuie să funcționeze în scopul satisfacerii conexiunii. Găsirea SE potrivite pentru conectare este dificilă deoarece PE și MM ce se conectează sunt distribuite de o manieră arbitrară. Modelul exact pentru RIN-MN este cunoscut ca fiind o problemă cu NP-complet. Aici se dezvoltă însă un model aproximativ.

Un grup conectat de i procesoare și j memorii, adică $i \times j$ conexiuni, poate fi obținut în două moduri.

În primul caz, i procesoare și j memorii sunt active și cel puțin numărul cerut de SE este disponibil pentru asigurarea conectării.

În al doilea caz, mai mult de i procesoare și j memorii sunt active, dar conexiunea este încă $(i \times j)$. Acest caz apare când SE active sunt suficiente să conecteze doar i procesoare și j memorii.

O RIN-MN de dimensiuni $(N \times N)$ cu 4×4 SE are un număr total de $N/4 \cdot \log_4 N$ comutatoare. Pentru o distribuție dată $(i \times j)$, fie x numărul de SE-uri cerut pentru conectare și $y = N/4 \cdot \log_4 N - x$ comutatoarele SE adiționale. Starea acestor SE adiționale nu afectează conectivitatea. Un MK a unei RIN-MN cu până la y defectări poate încă furniza conectarea cerută pentru un sistem $(i \times j)$. În al doilea caz, când mai mult decât $(i \times j)$ elemente sunt active, dacă se utilizează probabilitatea ca doar x SE să fie perfecte, talia sistemului încă va fi $(i \times j)$.

Un al doilea nivel de partiționare al SMP este utilizat pentru a-l găsi pe x , numărul cerut de SE. Talia sistemului este uzual dată de $(4' \times 4')$ deoarece RIN-MN este construită din SE de talie 4×4 . O configurație de $(4' \times 4')$ poate fi partiționată în patru subsisteme $(4'^{-1} \times 4'^{-1})$ fără a perturba conexiunea.

Metodologia de modelare nu este limitată la RIN-MN cu SE 4×4 . Tehnica de descompunere poate fi aplicată la orice nivel $\log_4 N$ proiectat cu SE $a \times a$. Rețeaua în acest caz ar putea fi

divizată logic în a grupuri. Precizia rezultatelor va depinde de valorile x și y . Este foarte dificil să se producă expresii generice pentru a calcula aceste valori comutatoare 4×4 .

C. Configurația cu 16 noduri

Cele 16 PE sunt partiționate în patru grupe unde fiecare grup are patru PE și comutatoarele corespunzătoare. Aceste grupuri sunt notate (PG_0, PG_1, PG_2, PG_3). Similar, cele 16 MM sunt partiționate în (MG_0, MG_1, MG_2, MG_3). Fie N_p numărul minim de grupuri de procesoare cerute pentru i PE active $1 \leq i \leq N$.

Atunci:

$$N_p = \left\lceil \frac{i}{N \text{div} 4} \right\rceil \quad (5.9.a)$$

unde: $\lceil \rceil$ reprezintă limitarea superioară

N = talia sistemului

Similar, pentru N_m , numărul minim al grupurilor de memorie necesare pentru j MM active ($1 \leq j \leq N$):

$$N_m = \left\lceil \frac{j}{N \text{div} 4} \right\rceil \quad (5.9.b)$$

Cele i PE pot fi distribuite în p grupuri de procesare pentru $1 \leq p \leq 4$ și j MM pot fi distribuite peste m grupuri de memorie, $1 \leq m \leq 4$. Numărul total de căi din care se alege p și m este

$\binom{4}{p} \binom{4}{m}$. Fiecare din aceste distribuții necesită SE specifice pentru conectare.

De exemplu, pentru $p = m = 2$, cele două grupuri de procesoare și două grupuri de memorii pot fi distribuite în 36 de căi.

Dacă două grupuri de procesare sunt conectate la două grupuri corespunzătoare de memorii, având același număr de grup, $\{PG_0, PG_1\} \leftrightarrow \{MG_0, MG_1\}$, atunci aceasta este o conexiune ($2 \leftrightarrow 2$). Sunt necesare două SE de intrare și două SE de ieșire pentru această conexiune. Acest caz reprezintă activitatea lui SE 0, SE 1, SE 4 și SE 5 din fig. 5.2.

Dacă în loc de selectarea grupurilor de memorie 0 și 1 se selectează 0 și 2, conexiunea devine $\{PG_0, PG_1\} \leftrightarrow \{MG_0, MG_2\}$. Aceasta reclamă o conexiune ($3 \leftrightarrow 3$). Sunt implicate SE 0, 1, 2, 4, 5 și 6. Există 24 de conexiuni ($3 \leftrightarrow 3$).

În fine, dacă conexiunea este între $\{PG_0, PG_1\} \leftrightarrow \{MG_2, MG_3\}$ toate cele patru SE de la intrare și cele patru SE de la ieșire sunt necesare în exemplul din fig. 5.2.

O astfel de conexiune este denumită $(4 \leftrightarrow 4)$ și există șase astfel de conexiuni. Conexiunea $(k \leftrightarrow k)$ este utilizată pentru a găsi numărul de SE de intrare și de ieșire necesar pentru a conecta p și m grupuri. Pentru un sistem (16×16) , $x = 2k$ pentru $1 \leq k \leq 4$ și $y = 8 - x$.

Pentru a exprima formal aceste combinații se utilizează un parametru de distribuție $D_g(p, m, \delta)$ pentru un p și m dat astfel:

$$D_g(p, m, \delta) = \binom{4}{\alpha} \binom{\alpha}{p - \delta} \binom{4 - \alpha}{\delta} \quad (5.10)$$

unde:

$$\alpha = \max(p, m)$$

$$\beta = \min(p, m)$$

δ - numărul de grupuri de memorie (procesoare) selecționate în exteriorul grupurilor de procesoare (memorii).

Pentru $\delta = 0$, $m \subset p$, dacă $p > m$ sau $p \subset m$, dacă $m > p$. Această situație produce o conexiune $(\alpha \leftrightarrow \alpha)$. Pentru $\delta = 1$, apare cazul când fie unul dintre grupurile de procesoare sau unul dintre grupurile de memorie este selectat din exteriorul numerelor corespunzătoare pentru a da o conexiune $((\alpha+1) \leftrightarrow (\alpha+1))$. În mod similar, $\delta = 2$ reprezintă două grupuri selectate în exteriorul numerelor de grup corespunzătoare pentru a da o conexiune $((\alpha+2) \leftrightarrow (\alpha+2))$. Teoretic, α poate varia între 0 și β .

Fie $p = 3$ și grupurile de procesoare 0, 1 și 2. Dacă se aleg toate grupurile de memorii ce lucrează dintre 0, 1 și 2, atunci aceasta este o selecție corespundentă și $m \subset p$. În exemplu prezentat, selecția grupului de memorie 3 nu este o selecție corespundentă și $\delta = 1$.

Primul termen în (5.10) arată cât de multe căi avem pentru a selecta α din patru grupuri.

Al doilea termen reprezintă numărul de căi posibile pentru a selecta $(\beta - \delta)$ din α pentru o valoare dată a lui δ .

Al treilea termen arată selecția lui δ din $(4 - \alpha)$ grupuri.

Ecuția (5.10) generează următoarele numere pentru $p = m = 2$

$$D_g(2,2,0) = 6 \Rightarrow \text{conectare } 2 \leftrightarrow 2$$

$$D_g(2,2,1) = 24 \Rightarrow \text{conectare } 3 \leftrightarrow 3$$

$$D_g(2,2,2) = 6 \Rightarrow \text{conectare } 4 \leftrightarrow 4$$

Pentru un i și j dat, p și m pot varia între N_p și 4 și N_m și 4, respectiv.

Se reprezintă numărul total de conexiuni ($k \leftrightarrow k$) ca fiind $sum(i, j)$, unde:

$$sum(i, j) = \sum_{p=N_p}^4 \sum_{m=N_m}^4 \sum_{\delta=0}^{\beta} D_g(p, m, \delta) \quad (5.11)$$

Ecuția (5.11) generează diferite conexiuni ($k \leftrightarrow k$) pentru fiecare p și m posibile, care sunt funcții ale lui δ . Deoarece există patru grupuri, $k \in \{\alpha, 4\}$.

Fiind dată $sum(i, j)$ se pune problema câte din acestea sunt conexiuni ($k \leftrightarrow k$). Numărul de conexiuni ce sunt de tipul ($k \leftrightarrow k$) se notează cu $g(k, k)$ pentru $\alpha \leq k \leq 4$ și este necesar pentru a calcula probabilitatea ca o conexiune ($k \leftrightarrow k$) să fie cerută.

$$g(k, k) = \sum_{p=N_p}^k \sum_{m=N_m}^k D_g(p, m, k - \alpha) \text{ pentru } \alpha \leq k \leq 4 \quad (5.12)$$

Utilizând (5.10) – (5.12) sunt generate mai multe conexiuni ($k \leftrightarrow k$) cerute pentru un număr dat i de PE și j de MM.

Rezultatele sunt prezentate în TABEL IX.

Pentru un p și m dați, fiecare căsuță din tabel prezintă numărul diferit de grupe ($k \leftrightarrow k$). Dacă elementele pot fi distribuite peste toate cele patru grupuri, atunci $sum(i, j)$ este obținută prin adăugarea de intrări a tuturor căsuțelor între p și 4 și m și 4.

Funcția $g(k, k)$ este obținută prin adăugarea a ($k \leftrightarrow k$) intrări tuturor căsuțelor utilizate în calcularea lui $sum(i, j)$.

Pentru $p = m = 2$, numărul total de căsuțe ce trebuie incluse în calcul este prezentat cu linii îngroșate în TABEL IX. El dă $sum(i, j)=121$, $g(2,2)=6$, $g(3,3)=52$ și $g(4,4)=63$. Atunci $g(k, k)/sum(i, j)$ este probabilitatea că o conexiune ($k \leftrightarrow k$) este cerută.

TABEL IX. Tabelul de căutare a conexiunilor

$\begin{matrix} p \\ m \end{matrix}$	1	2	3	4
1	4(1-1) 12(2-2)	12(2-2) 12(3-3)	12(3-3) 4(4-4)	4(4-4)
2	12(2-2) 12(3-3)	6(2-2) 24(3-3) 6(4-4)	12(3-3) 12(4-4)	6(4-4)
3	12(3-3) 4(4-4)	12(3-3) 12(4-4)	4(3-3) 12(4-4)	4(4-4)
4	4(4-4)	6(4-4)	4(4-4)	1(4-4)

Reamintim că o conexiune ($k \leftrightarrow k$) determină numărul x de comutatoare solicitate ce trebuie să fie active pentru grupurile de p procesoare, în vederea conectării la grupurile de m memorii. MK pentru RIN-MN este prezentat în fig. 5.10.

Se notează cu (x, y) numărul de x SE-uri cerute și cu y numărul de SE-uri suplimentare în RIN-MN, astfel ca $x + y = N/4 \log N/4 = S_i$ – numărul total de comutatoare din sistem.

Se notează cu λ_{se} , μ_{se} , C_{se} și r_{se} ratele de defectare, de reparare, factorul de acoperire și factorul de reparare cu succes pentru comutatoarele SE.

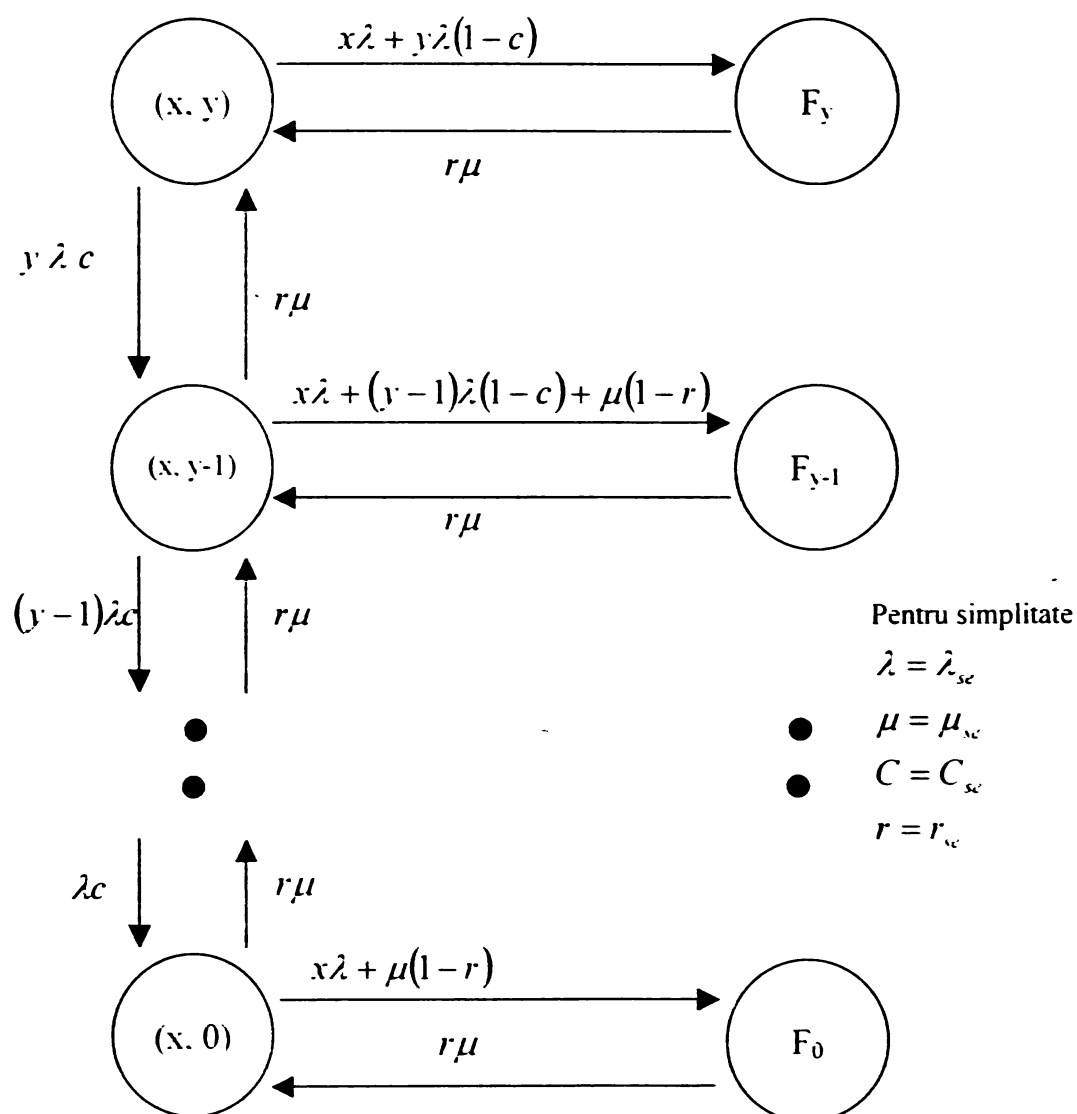


Fig. 5.10. Lanț Markov pentru RIN-MN

Tranziția orizontală cu rata de $x \lambda_{se} C_{se}$ arată că sistemul se defectează când oricare dintre comutatoarele solicitate se defectează. A doua parte a tranziției orizontale cu termenul $(1 - C_{se})$ reprezintă defectarea sistemului datorită acoperirii imperfecte a SE-urilor defecte ce aparțin grupului suplimentar, de extensie.

Deși un comutator suplimentar, un extra-SE, nu contribuie la nici o conexiune utilă, defectarea sa poate fi detectată de o astfel de manieră încât reconfigurarea sistemului este posibilă. A treia parte reprezintă repararea imprecisă a SE-urilor.

Tranzițiile verticale reflectă defectarea și repararea a celor y extra SE-uri. Se observă că dintr-o stare (x, y) , subsistemul RIN-MN ar putea evolua fie spre o stare defectă – tranziție orizontală – sau spre o stare $(x, y-1)$ – tranziție verticală – deoarece doar un singur SE ar putea să se defecteze la un moment dat. Prin urmare, rata de reparare este întotdeauna μ_{se} în MK prezentat în fig. 5.10. Repararea unui SE necesar în conexiune se consideră că este prioritară față de cea a unui extra-SE. Această strategie evită tranziția între stările defecte și ameliorează disponibilitatea RIN-MN.

Modelul prezentat în fig. 5.10 poate modela orice degradare în performanțe cu valori adecvate ale lui x și y . De exemplu, dacă procesoarele și memoriile active au nevoie de cel puțin două grupuri de procesoare și două grupuri de memorii, RIN-MN din fig. 5.10 trebuie să furnizeze o conexiune $2 \leftrightarrow 2$, una $3 \leftrightarrow 3$ sau una $4 \leftrightarrow 4$, depinzând de distribuția procesoarelor și memoriilor. O conexiune $2 \leftrightarrow 2$ furnizează $x=4, y=4$, o conexiune $3 \leftrightarrow 3$ dă $x=6, y=2$ și o conexiune $4 \leftrightarrow 4$ înseamnă $x=8, y=0$ (nu există degradare a RIN-MN). Astfel, pentru evaluarea unui sistem cu o degradare de 50%, MK-ul din fig. 5.10 ar putea fi rezolvat de trei ori, cu trei seturi diferite de valori (x, y) .

Fie $P_{(x, l)}(t)$ probabilitatea ca sistemul să fie în starea (x, l) la momentul t . Soluția MK-ului din fig. 5.10, utilizând pachetul HARP, va furniza $P_{(x, l)}(t)$ pentru $0 \leq l \leq y$.

Probabilitatea ca RIN-MN să asigure o conexiune $k \leftrightarrow k$ cu x SE active este $P_{(k \leftrightarrow k)}(t)$ și are expresia :

$$P_{(k \leftrightarrow k)}(t) = \sum_{l=0}^y P_{(x, l)}(t) \quad (5.13)$$

Pentru a recapitula, pentru orice subsistem $(i \times j)$, numărul p de grupuri de procesoare și numărul m de grupuri de memorii se situează în domeniul $N_p \leq p \leq 4$ și $N_m \leq m \leq 4$, respectiv. Valoarea lui k pentru o conexiune $k \leftrightarrow k$ se situează în domeniul $\max(p, m) \leq k \leq 4$. În final,

pentru o conexiune $k \leftrightarrow k$ cunoscută, valorile x și y ale unei configurații (16×16) este obținută ca fiind $2k$ și respectiv $(8-2k)$.

Probabilitatea că există un grup $(i \times j)$ conectat într-un sistem (16×16) atunci când i procesoare și j memorii sunt active este :

$$P_{(i,j)}^{(1)}(t) = \sum_{p=N_p}^4 \sum_{m=N_m}^4 \left(\sum_{k=\max(p,m)}^4 \frac{g(k,k)}{\text{sum}(i,j)} P_{(k \leftrightarrow k)}(t) \right) P_i(t) P_j(t) \quad (5.14)$$

În această expresie, indicele superior 1 reprezintă primul caz, acela în care doar i PE și j MM sunt active. Primele două sume denotă diversele combinații ale grupurilor de procesoare și de memorii conținând i PE și j MM. Al treilea termen de însumare arată că k poate varia până la 4. De exemplu, dacă $p = 2$, $m = 3$, atunci $k = 3$. Aceasta înseamnă că dacă i PE și j MM sunt distribuite în două grupuri de procesoare și în trei grupuri de memorii, depinzând de distribuție, atunci sunt necesare șase SE-uri de tip $(P_{(3 \leftrightarrow 3)}(t))$ sau toate cele opt SE-uri $(P_{(4 \leftrightarrow 4)}(t)) \cdot g(k,k) / \text{sum}(i,j)$ ce dau probabilitatea de a solicita o conexiune $k \leftrightarrow k$.

Pentru $p=m=3$, $g(3,3)=16$, $g(4,4)=39$ și $\text{sum}(i,j)=55$ pentru TABEL IX. Multiplicând probabilitatea subsistemului de comutatoare cu $P_i(t)$ și $P_j(t)$, obținute din subsistemele procesor și de memorie, se obține probabilitatea din primul caz care este $P_{(i,j)}^{(1)}(t)$.

În al doilea caz, conexiunea este restricționată la $(i \times j)$ când α PE și β MM sunt active în sistem, pentru $\alpha \geq i$, $\beta \geq j$. Aceasta este posibil doar dacă comutatoarele solicitate i și j sunt active în grupuri separate și $(\alpha - i)$ și $(\beta - j)$ elemente funcționează în grupuri separate și nu există conexiuni între subsistemul $(i \times j)$ activ și restul componentelor active. Să considerăm un exemplu în care i PE și j MM lucrează utilizând comutatoarele 1, 2, 3, 5, 6 și 7 din fig. 5.2. Dacă oricare altă componentă lucrează în aceste trei grupuri de procesoare și în aceste trei grupuri de memorii, atunci conectivitatea crește.

Procesoarele suplimentare și/sau memoriile trebuie să funcționeze în grupul de procesoare 0 și/sau în grupul de memorii 0, și în grupul SE 0 sau 1. Aceasta implică că dacă i PE și j MM solicitate sunt distribuite în cele patru grupuri, nu poate exista orice altă componentă suplimentară.

Probabilitatea ce este specifică cazului doi, devine nulă în această situație.

Această situație poate fi enunțată astfel : «Care este probabilitatea ca i procesoare și j memorii să necesite o conexiune $k \leftrightarrow k$ pentru $1 \leq k \leq 3$, și ca $(\alpha - i)$ și $(\beta - j)$ elemente să fie distribuite în restul celor $(4-k)$ grupuri ? »

Acest enunț poate fi exprimat astfel :

$$P_{cx}(t) = \frac{\sum_{k=\max(N_p, N_m)}^3 \binom{4}{k} \binom{k \times 4}{i} \binom{N-4 \times k}{\alpha-i} \binom{k \times 4}{j} \binom{N-4 \times k}{\beta-j}}{\binom{N}{\alpha} \binom{N}{\beta} \binom{\alpha}{i} \binom{\beta}{j}} \quad (5.15)$$

În această expresie, numărătorul oferă numărul de căi pentru a alege i PE din $k \times 4$ procesoare aparținând la k grupuri (fiecare grup într-un sistem 16×16 are 4 PE sau 4 MM), j MM dintre $k \times 4$ memorii din k grupuri, $(\alpha - i)$ procesoare din restul a $(N - 4 \times k)$ procesoare, distribuite în $(4 - k)$ grupuri, și $(\beta - j)$ memorii din restul celor $(N - 4 \times k)$ memorii. $\binom{4}{k}$ reprezintă numărul de căi prin care se poate selecta un grup $k \leftrightarrow k$ procesor-memorie ce conține un subsistem $(i \times j)$. Numitorul oferă numărul total de căi prin care i PE și j MM pot fi selectate. Limita superioară a lui k este 3 deoarece distribuția a i PE și j MM în cele patru grupuri face ca probabilitatea cazului al doilea să fie nulă.

Starea $(x,0)$ în MK prezentat în fig. 5.10 reprezintă situația în care doar comutatoarele solicitate lucrează și toate celelalte SE adiționale s-au defectat.

Includerea acestei probabilități de stare ar garanta că procesoarele adiționale active și/sau memoriile corespunzătoare nu pot să crească talia sistemului activ mai mult de $(i \times j)$. De notat că $P_{(x,0)}(t)$ poate fi obținută pentru un grup $(k \leftrightarrow k)$ dat.

Probabilitatea pentru cel de al doilea caz unde se poate obține o conexiune $(i \times j)$ din partea a α PE active și β MM active se notează cu $P_{(\alpha,\beta)}(t)$.

$$P_{(\alpha,\beta)}(t) = P_{cx}(t) \cdot P_{(x,0)}(t) \cdot P_{\alpha}(t) \cdot P_{\beta}(t) \quad (5.16)$$

unde $P_{\alpha}(t)$ și $P_{\beta}(t)$ dau probabilitatea ca α PE și β MM să fie active la momentul t . Deoarece α și β pot fi în domeniul $i \leq \alpha \leq N$ și $j \leq \beta \leq N$ respectiv, probabilitatea cazului al doilea este :

$$P_{(i,j)}^{(2)}(t) = \sum_{\alpha=i+1}^N \sum_{\beta=j+1}^N P_{(\alpha,\beta)}(t) + \sum_{\alpha=i+1}^N P_{(\alpha,i)}(t) + \sum_{\beta=j+1}^N P_{(i,\beta)}(t) \quad (5.17)$$

Al doilea și al treilea termen în (5.17) au o contribuție neglijabilă și pot fi ignorate pentru a face calculul mai rapid.

Finalmente, disponibilitatea sistemului devine :

$$A(t) = \sum_{i=1}^N \sum_{j=1}^N (P_{(i,j)}^{(1)}(t) + P_{(i,j)}^{(2)}(t)) \quad (5.18)$$

5.5.3 Un model generalizat

Pentru a generaliza modelul prezentat anterior în [DTB93] sunt prezentate ca fiind necesare următoarele.

În primul rând, metodologia de modelare trebuie să fie extinsă pentru a se putea analiza sisteme mai mari cu o acuratețe rezolvabilă.

Al doilea obiectiv este de a include analiza sistemelor a căror talie nu este o putere a lui 4. Această necesitate apare datorită faptului că multe SMP pot fi configurate pentru diferite talii. De exemplu, configurația Butterfly poate fi proiectată pentru orice talie până la 256 de noduri. Extensia modelului trebuie să includă și alte clase de RIN-MN și de SMP-MP.

Ideea principală este de a descompune sistemul în diferite grupuri. Distribuția a i PE și j MM poate fi făcută într-o manieră similară cu cea de la sistemele (16×16) .

În funcție de talia sistemului și cu ajutorul lanțului Markov prezentat în fig. 5.9, se poate calcula $P_i(t)$ și $P_j(t)$. Problema ce trebuie detaliată o reprezintă determinarea numărului x de SE necesare. Numărul y de SE suplimentare poate fi obținut calculând $y = S_i - x$ unde S_i este numărul total de SE din RIN-MN.

Cunoașterea lui x și y este esențială pentru aflarea soluției modelului MK prezentat în fig. 5.10.

A. Modelarea într-o configurație de noduri $N=4^i$

Se consideră un sistem $(N \times N)$ unde $N = 4^i$. El are $n = \log_4 N$ niveluri conținând $N/4$ SE în fiecare nivel. Se numerează nivelurile de la 0 la $(n-1)$ dinspre partea procesorului.

Prin urmare, $S_i = n \cdot (N/4)$ reprezintă numărul total de comutatoare. Procesoarele sunt grupate în patru grupuri de procesare – PG – și modulele de memorie în patru grupe de memorii – MG. Fiecare grup, indiferent dacă este PG sau MG, are 4^{i-1} elemente. Un PE întotdeauna are nevoie de un nivel specific 0 de SE în scopul de accesa oricare MM. Astfel, comutatoarele din nivelul 0 sunt incluse cu PG-urile corespunzătoare. Pentru a simplifica calculul, SE din fiecare nivel de la 2 la $(n - 1)$ sunt, de asemenea, divizate în patru grupuri. Fiecare grup conține $N/16$ SE / nivel. Există o relație între diferitele grupuri de $N/16$ SE/nivel și grupurile MG asociate. Accesul la grupurile MG poate fi asigurat prin includerea acestor SE cu MG-urile lor respective. Fiecare MG conține, astfel, $(n - 2) \cdot (N/16)$ comutatoare. În acest mod, doar comutatoarele din nivelul 1 nu sunt incluse, nici în grupurile PG, nici în cele MG.

Comutatoarele din nivelul 1 sunt acelea ce realizează conectarea dintre grupurile PG și cele MG. Numărul de SE din nivelul 1 depinde de distribuția grupurilor PG și MG cerute. Bazat pe funcționalitate, se grupează comutatoarele nivelului 1 în grupuri de comutatoare – SG. Un grup

de comutatoare SG_{pq} asigură conexiunea de la PG_p la MG_q . Nodurile de comutare sunt prezentate în TABEL X și acesta este util pentru formalizarea conceptului.

TABEL X. Tabelul nodurilor de comutare

PG \ MG	0	1	2	3
0	00	01	02	03
1	10	11	12	13
2	20	21	22	23
3	30	31	32	33

Fiecare locație din tabel se referă la un SG și, prin urmare, există 16 astfel de locații. Numărul de SE în fiecare grup este egal cu $N / 64$. Din descrierea conexiunii ($k \leftrightarrow k$) oferită ca și exemplu pentru analiza unui sistem cu 64 de noduri, doar un SG este cerut pentru o conexiune ($1 \leftrightarrow 1$) în timp ce pentru conectări ($2 \leftrightarrow 2$) sau ($3 \leftrightarrow 3$) sunt necesare patru sau, respectiv nouă SG.

Atât SG_{pq} și SG_{qp} sunt necesare cu scopul de a avea o conexiune între PG_p și MG_q . Numărul de SG-uri necesar pentru o conexiune ($k \leftrightarrow k$) este, prin urmare, egal cu k^2 . De exemplu, o configurație cu 1024 noduri are 256 SE în nivelul 1 și fiecare din cele 16 grupuri are 16 SE-uri. Fiecare locație din TABEL X în acest caz reprezintă un grup de 16 SE-uri.

Numărul total de SE-uri necesare pentru o conexiune de grup ($k \leftrightarrow k$), pentru $1 \leq k \leq 4$, se obțin astfel :

$x = (\text{numărul de SE în } k \text{ PG-uri}) + (\text{numărul de SE în } k \text{ MG-uri}) + (\text{numărul de comutatoare de pe nivelul 1 necesare pentru a conecta } k \text{ PG-uri la } k \text{ MG-uri})$

$$\begin{aligned}
 &= k \cdot \frac{N}{16} + k(n-2) \cdot \frac{N}{16} + k^2 \cdot \frac{N}{64} \\
 &= \frac{Nk}{16} \left[n-1 + \frac{k}{4} \right] \tag{5.19}
 \end{aligned}$$

Modelul MK din fig. 5.10 este rezolvat utilizând valoarea lui x și valoarea corespunzătoare a lui y .

Ecuțiile (5.13) și (5.14) se utilizează, în continuare, pentru a găsi probabilitatea ca să existe un grup $(i \times j)$ conectat când sunt active i PE și j MM. $P_{ct}(t)$ este obținut din (5.12) înlocuind $(k \times 4)$ cu $(k \times N/4)$

Probabilitatea unei conexiuni $(i \times j)$ când mai mult decât i PE și j MM sunt active se obține rezolvând (5.16) și (5.17). În final, (5.18) oferă disponibilitatea sistemului.

B. Modelarea într-o configurație de noduri $N \neq 4^i$

Se consideră un sistem $N \times N$ unde $4^{i-1} < N < 4^i$ pentru un întreg i oarecare. Numărul n de niveluri necesar pentru a configura RIN-MN este i . Următorul parametru de configurare este numărul de procesoare și de grupuri de memorie. Divizarea SMP în patru grupe ar putea duce la o proiectare ineficientă și scumpă. De exemplu, să analizăm situația unui sistem cu 100 de noduri. Se pot repartiza aceste noduri în toate cele patru grupe. Multe dintre comutatoarele (4×4) ar putea fi subutilizate sau chiar neutilizate. Este, prin urmare, esențial ca aceste noduri să fie grupate în cât mai puține grupe posibile. O grupare densă înseamnă a selecta nodurile de o astfel de manieră încât să se maximizeze utilizarea comutatorului. Aceasta conduce la formarea de grupe de 4^{n-1} noduri. Astfel, numărul de grupe va depinde de talia N a sistemului.

Fie G numărul de grupe unde $1 \leq G \leq 4$.

Atunci

$$G = \left\lceil \frac{N}{4^{n-1}} \right\rceil \quad (5.20)$$

Prin urmare, exemplul sistemului cu 100 noduri necesită două grupe. Distribuția a N noduri în G grupe poate fi realizată în multe moduri dacă $N \neq G \cdot 4^{n-1}$.

În exemplul anterior, distribuția a 100 noduri în două grupe este (50, 50), (64, 34), etc.

Se presupune că N este distribuit în mod egal în cele G grupe, astfel că numărul de SE necesare pentru fiecare grup este același.

Numărul minim de PG-uri și MG-uri pentru analiza a i PE și j MM este :

$$N_p = \left\lceil \frac{i}{N \text{div} G} \right\rceil \text{ și } N_m = \left\lceil \frac{j}{N \text{div} G} \right\rceil \quad (5.21)$$

Cele i PE pot fi distribuite peste p PG pentru $N_p \leq p \leq G$ și j MM pot fi distribuite peste m MG pentru $N_m \leq m \leq G$. Fiecare din aceste distribuții necesită un set specific de SE pentru

interconectare. Înlocuind valoarea 4 în G în ecuațiile (5.10), (5.11) și (5.12), se pot calcula valorile pentru $g(k,k)$ și $sum(i,j)$ unde $\alpha \leq k \leq G$.

Calcularea lui x este similară ca și tehnică cu cea descrisă în 5.5.3.A. Comutatoarele de pe nivelul 0 sunt incluse în grupurile PG și astfel, fiecare PG are N/G PE-uri și $\lceil (1/G)/(N/4) \rceil$ comutatoare pe fiecare nivel. Toate comutatoarele situate pe nivelul 2 până la cele de pe nivelul $(n - 1)$ sunt incluse în MG-urile corespunzătoare. Prin urmare, fiecare MG constă din M/G grupuri MM și $(n - 2) \cdot \lceil N/4G \rceil$ comutatoare. Pasul următor constă în determinarea numărului solicitat de comutatoare de pe nivelul 1 necesare pentru o conexiune ($k \leftrightarrow k$). Tabelul nodurilor de comutare pentru o astfel de conectare este TABEL X, dar numărul de linii și coloane este restrâns la G .

Deoarece există G grupe PG și G grupe MG în TABEL X, se descompun comutatoarele de pe nivelul 1 în G^2 comutatoare. Fiecare SG conectează un PG specific la un MG. Numărul de comutatoare în fiecare SG este dat de $(N/4) \times (1/G^2) = N/4G^2 \cdot N/4G^2$. El este un număr întreg în anumite cazuri. Aceasta se datorează situației când talia sistemului nu este o putere a lui 4. Atunci, câteva dintre comutatoarele de pe nivelul 1 nu sunt utilizate pentru conectare.

Numărul total x de comutatoare pentru o conexiune ($k \leftrightarrow k$) pentru $1 \leq k \leq G$ se obține din :

$$\begin{aligned} x &= k \cdot \frac{N}{4G} + k(n-2) \cdot \frac{N}{4G} + k^2 \cdot \left\lfloor \frac{N}{4G^2} \right\rfloor \\ &= \frac{Nk}{4G}(n-1) + k^2 \left\lfloor \frac{N}{4G^2} \right\rfloor \end{aligned} \quad (5.22)$$

Valorile x și y astfel obținute sunt utilizate în modelul MK din fig. 5.10.

C. Extensia modelului

Tehnica de modelare propusă poate fi utilizată pentru evaluarea RIN-MN ce utilizează comutatoare de altă mărime decât 4×4 . Analiza grupurilor PE și MM rămâne aceeași. Calculul diferiților parametrii utilizând p grupuri de PE active și m grupuri de MM active se pot utiliza aplicând expresiile derivate în 5.5.2.

Dacă talia comutatorului este $a \times a$ atunci numărul de grupuri trebuie să fie a în loc de 4.

O a doua considerație trebuie făcută în legătură cu x ce trebuie utilizat în modelul MK din fig. 5.10.

Se pot utiliza diferite tehnici de aproximare pentru aflarea valorii lui x . O abordare constă în utilizarea tuturor comutatoarelor unui grup, așa cum s-a prezentat anterior. Această metodă este

rezonabilă pentru $a \geq 4$, deoarece cu un număr mai mare de comutatoare, numărul de niveluri scade. Pe de altă parte, pentru $a = 2$ sau 3 , este necesară o aproximare mai exactă a lui x .

Următoarea extensie a modelului se referă la SMP-MP. Un astfel de sistem are toate procesoarele PE dispuse pe o parte, iar memoriile MM pe cealaltă parte a RIN-MN. Este necesară doar o cale de la un PE la un MM, în vederea conectării.

Utilizând notațiile din 5.5.3.B se vor obține $p \leftrightarrow m$ conexiuni în loc de $(k \leftrightarrow k)$ conexiuni, unde $N_p \leq p \leq G$ și $N_m \leq m \leq G$. Numărul necesar de comutatoare este prezentat în TABEL X.

Toleranța la defect poate fi ameliorată prin utilizarea unor tehnici specifice. Aceste tehnici includ adăugarea unor niveluri suplimentare, crescând numărul de porturi, replicând RIN-MN, etc. etc. Deoarece redundanța și topologia a RIN-MN tolerante la defect variază, este dificil de a formaliza o tehnică generală pentru modelare și pentru calcularea disponibilității.

Tehnica de modelare a unui subsistem RIN-MN variază în funcție de topologie. Conceptul prezentat poate fi aplicat la RIN-MN tolerante la defect, dacă comutatoarele pot grupate în două grupuri: numărul x cerut de SE și numărul y de comutatoare suplimentare. Este dificil de identificat x și y într-o RIN-MN *cu căi multiple*. Validitatea modelului unei RIN-MN tolerantă la defect depinde de ușurința cu care se pot găsi valorile corespunzătoare pentru x și y .

5.6 CONCLUZII

Analiza de dependabilitate în arhitecturile paralele se referă atât la *analiza de fiabilitate* cât și la cea de *disponibilitate* [RA90], [RG87], [TD89], [TZ94].

Fiabilitatea este una dintre cele mai importante attribute ale sistemelor de calcul. Cele mai utilizate tehnici de analiză ale fiabilității sunt abordările analitice. Dintre acestea am analizat *modelarea combinatorie* și *modelarea Markov*.

Modelarea combinatorie utilizează tehnici probabilistice ce enumeră modurile diferite în care un sistem poate rămâne operațional. Probabilitățile evenimentelor ce conduc la starea în care un sistem devine operațional sunt calculate ca să formeze o estimare a fiabilității sistemului. Fiabilitatea sistemului este în general exprimată în termeni de fiabilități ale componentelor individuale. În practică se utilizează două modele, serie și paralel. În *sistemul serie*, fiecărui element al sistemului i se pretinde să funcționeze corect pentru ca sistemul în ansamblul lui să funcționeze corect. Într-un *sistem paralel*, doar câteva din elemente trebuie să fie operaționale pentru ca sistemul să funcționeze corect.

Ca și dezavantaj major al acestei metode se remarcă faptul că sistemele complexe nu pot fi modelate convenabil deoarece expresiile fiabilității devin foarte complexe. În plus, acoperirea defectului este dificil de incorporat în expresii combinaționale.

Modelele Markov utilizează concepte de *stare* și de *tranziție a stării*. Pentru modelele de fiabilitate fiecare stare a modelului Markov reprezintă o combinație distinctă de module ce defectează și de module fără defect. Tranzițiile de stare sunt caracterizate de probabilități, ca și *probabilitatea de defectare*, *de acoperire a defectului* și *de reparare*. Principalul beneficiu al utilizării MK în modelarea dependabilității îl reprezintă abilitatea de a manipula acoperirea defectului într-o manieră sistematică. Incorporarea acoperirii C în analiza MK permite modelarea siguranței sistemului.

Complexitatea din punct de vedere al calculului al unei analize de dependabilitate a fost redusă drastic prin fragmentarea în module mai mici, utilizând un model de rețea ierarhică mai apropiat de modelul de sistem și care favorizează comunicările locale. Proiectanții combină metode algebrice Booleene binecunoscute cu un model de rețea ierarhică pentru a evalua *fiabilitatea de terminal*, *fiabilitatea de rețea* și *fiabilitatea bazată pe sarcină*.

S-au dezvoltat mai mulți algoritmi pentru evaluarea fiabilității SMP. În acești algoritmi, fiabilitatea de tip terminal este utilizată în mod comun ca și o măsură a conectivității. A fost studiată de mulți cercetători și există diferite metode de evaluare a acestui parametru. Amintim, cu titlu informativ, că cea mai simplă dintre aceste tehnici este de a enumera stările favorabile prin examinarea evenimentelor posibile. Evident, acest lucru este imposibil în rețelele de dimensiuni mari.

Modelarea disponibilității presupune încorporarea reparării. Lanțurile Markov au proprietatea de a fi ciclice. Prin introducerea noțiunii de *reparare* se permite reîntoarcerea de la o stare mai puțin operațională la o stare mai operațională sau chiar complet operațională. Pentru aceasta se specifică o probabilitate de întoarcere la nivel de tranziție. Evaluarea experimentală a disponibilității nu este întotdeauna posibilă din cauza constrângerilor de timp și de cost. Din această cauză se utilizează frecvent parametri *MTTF* sau *MTTR* ceea ce conduce la disponibilitatea stării ferme, A_{SS} .

Disponibilitatea sistemului este evaluată în sistemele de dimensiuni mari și prin alte tehnici dintre care amintim strategiile de rutare ierarhice. Aceasta înseamnă că fiabilitatea este asociată cu existența unor căi de comunicare între toate perechile de noduri. Procesul de obținere a expresiilor de fiabilitate pornind de la un set minimal este un proces consumator de timp. De exemplu calculul fiabilității de tip *k-terminal* prin considerarea tuturor combinațiilor posibile de noduri și legături funcționale și nefuncționale a unei rețele ierarhice este exponențial în natura sa

și puternic restrictiv în rețelele ierarhice. În literatură [DA93] se indică mai multe metode de evitare prin furnizarea unor expresii aproximative ale fiabilității pentru o rețea oricât de mare, prin modelarea acestora într-o formă ierarhică și prin calcularea diverșilor parametri ale fiecărui grup ce conține un număr relativ redus de noduri. Ulterior se combină aceste expresii pentru a obține o valoare generală a fiabilității.

Metodele de analiză și evaluare a dependabilității sunt numeroase, înregistrându-se o direcție fermă de cercetare în acest sens. Tehnici de evaluare exacte presupun însă un volum mare de calcul și de aceea proiectanții se bazează pe soluții de proiectare devenite clasice, ale căror fiabilitate și disponibilitate sunt cunoscute, concentrându-se însă asupra unor expresii statistice ale dependabilității întregului sistem. Tendința firească este însă de a utiliza simulatoare și programe de test pentru a se obține o privire de ansamblu cu un grad de generalitate acceptat asupra comportării unui sistem din punct de vedere al dependabilității.

6. MODELUL EXPERIMENTAL

6.1 ARHITECTURA SMP – MP UTILIZAT CA ȘI MODEL

Sistemul de calcul multiprocesor, utilizat ca și model experimental, reprezintă o arhitectură clasică pentru un mediu de calcul cu memorie partajată distribuită.

Memoria este distribuită printre procesoare, dar toate memoriile partajează același spațiu de adresă.

Nodurile N sunt conectate între ele prin intermediul unei rețele de interconectare de bazată pe *magistrală multinivel (Multistage Bus Network)*. Vom utiliza în continuare pentru acest tip de rețea acronimul RMM.

Arhitectura sistemului multiprocesor este prezentată în fig.6.1.

Fiecare nod N_i conține :

- procesorul P_i
- memoria M_i , locală pentru P_i , și controlerul ei CM_i
- memoria cache C_i și controlerul ei CC_i
- memoria cache utilizată pentru director D_i și controlerul ei CD_i
- interfața cu RMM

Un procesor P_i din nodul N_i este conectat la un modul de memorie M_i atât prin intermediul RMM cât și direct. M_i este o *memorie locală* pentru P_i .

Ierarhia de memorie este constituită din *memorii cache* situate pe stratul (layer) superior L_1 , din *memorii cache* situate pe stratul L_2 utilizat de procesoarele ce aparțin aceluiași *grup (cluster processors)* și din *memoriile externe*, locale sau situate la distanță. Nivelul L_1 conține memorii cache de capacitate redusă (8 – 32 Kbytes) ce servesc o cerere de acces în cadrul unui singur ciclu de procesor. Nivelul L_2 conține memorii cache de capacitate mult mai mare (1 – 4 Mbytes) și onorează o solicitare în mai multe cicluri de procesor. Deoarece întregul spațiu de adresă este partajat, dar memoria este fizic distribuită, eșecurile de acces la L_1 și la L_2 pot fi tratate atât în memoria locală cât și în cea situată la distanță. Accesul la memorie este neuniform și prin urmare organizarea este denumită cu *acces neuniform la memorie (Non-Uniform Memory Access)* - NUMA. Accesul la memoria locală necesită transferul de date prin magistrala internă din nod și

acces la memoria DRAM. Accesul la memoria situată la distanță necesită transferuri de date prin magistrală, transferuri de mesaje prin RIN, și acces la DRAM. Astfel, timpul de deservire a unui acces ratat la memoria cache, în mod special atunci când data trebuie adusă de la o memorie situată la distanță, poate fi de câteva ordine de mărime mai mare ca și timpul de acces la memoria cache.

Arhitectura sistemului de calcul este prezentată extins, cu evidențierea ierarhiei de memorie, în fig. 6.2.

De remarcat că în timp ce latențele în accesul la memoria locală pot fi tolerate, accesele la memoria de la distanță, generate în timpul execuției, pot reduce drastic performanța.

Pentru a reduce impactul latențelor în accesul memoriilor de la distanță în literatură [SLM89], [MBL89], [ADT90], [NB93], [IB00], [IB99], [BINK98] s-au propus mai multe soluții. Amintim câteva dintre ele :

- creșterea eficienței memoriilor cache private
- ameliorare performanțelor RIN
- managementul memoriilor ierarhizate

O atenție specială trebuie acordată proiectării comutatorului în cadrul unei RIN-MN deoarece există un impact direct asupra performanței SMP-MPD.

În cadrul modelului experimental am luat în considerare un comutator 2 X 2 bazat pe o magistrală internă. Am utilizat comutarea de pachete ca și tehnică de comutare prin rețea. În ipoteza că se utilizează tehnici de comutare mai avansate [BWK98], [BINK97], [BIK00] ca și *wormhole*, *canale virtuale* și *rutare adaptivă* se pot aborda arhitecturi de comutatoare mult mai sofisticate. De exemplu SGI Spider [IB00], similar cu Intel Cavallino, sunt comutatoare *crossbar* ce utilizează canale virtuale. În literatură [BWK98], [Iyer99] se prezintă diferite arhitecturi de comutatoare 2 X 2 ce utilizează canale virtuale cu sau fără tamponare. Este interesant de remarcat că prin incorporarea unei memorii cache rapide SRAM în cadrul fiecărui comutator dintr-o RIN-MN se pot colecta datele partajate așa cum evoluează ele prin rețea și se pot astfel asigura accesele viitoare de la procesoarele ce reutilizează aceste date.

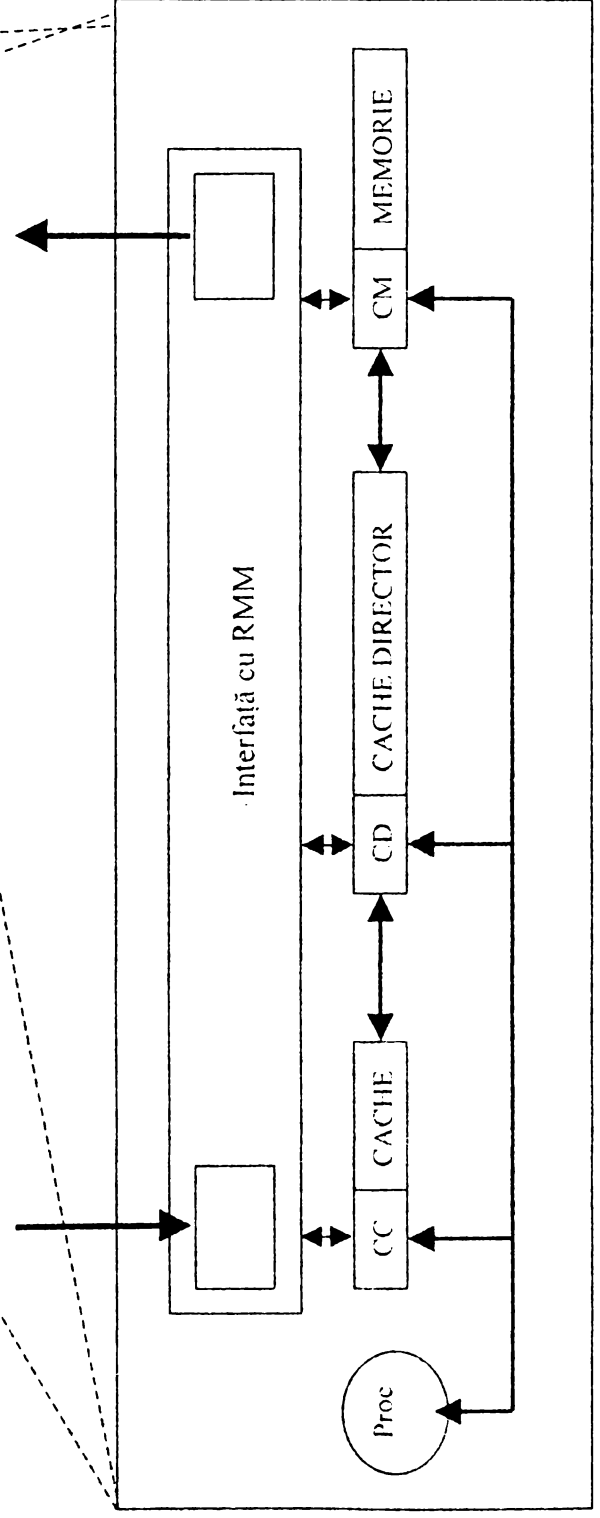
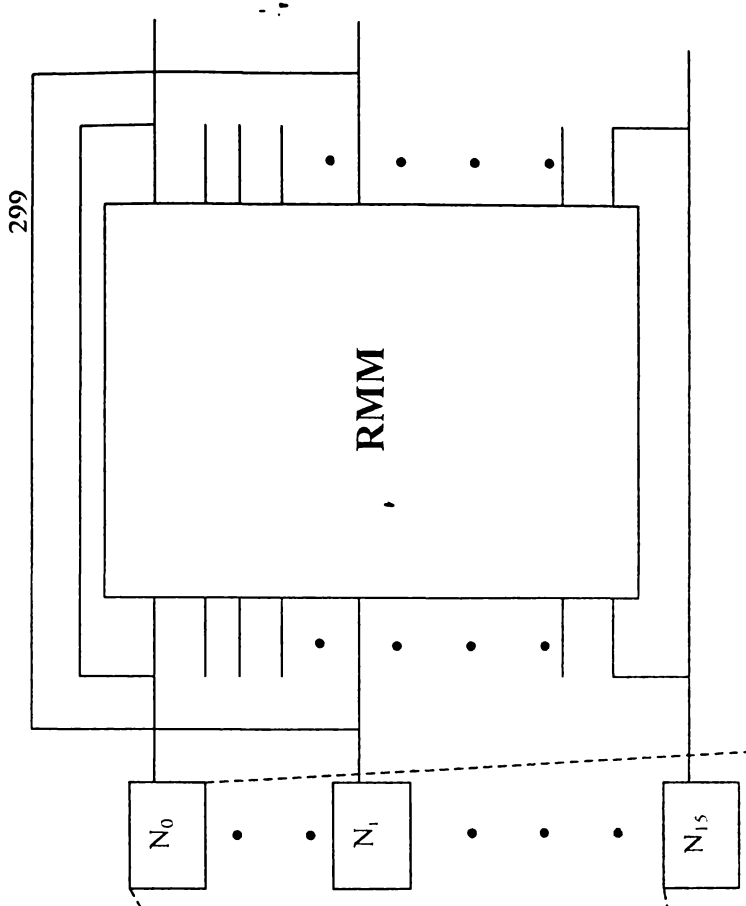


Fig. 6.1 Arhitectura SMP propusă ca și model experimental

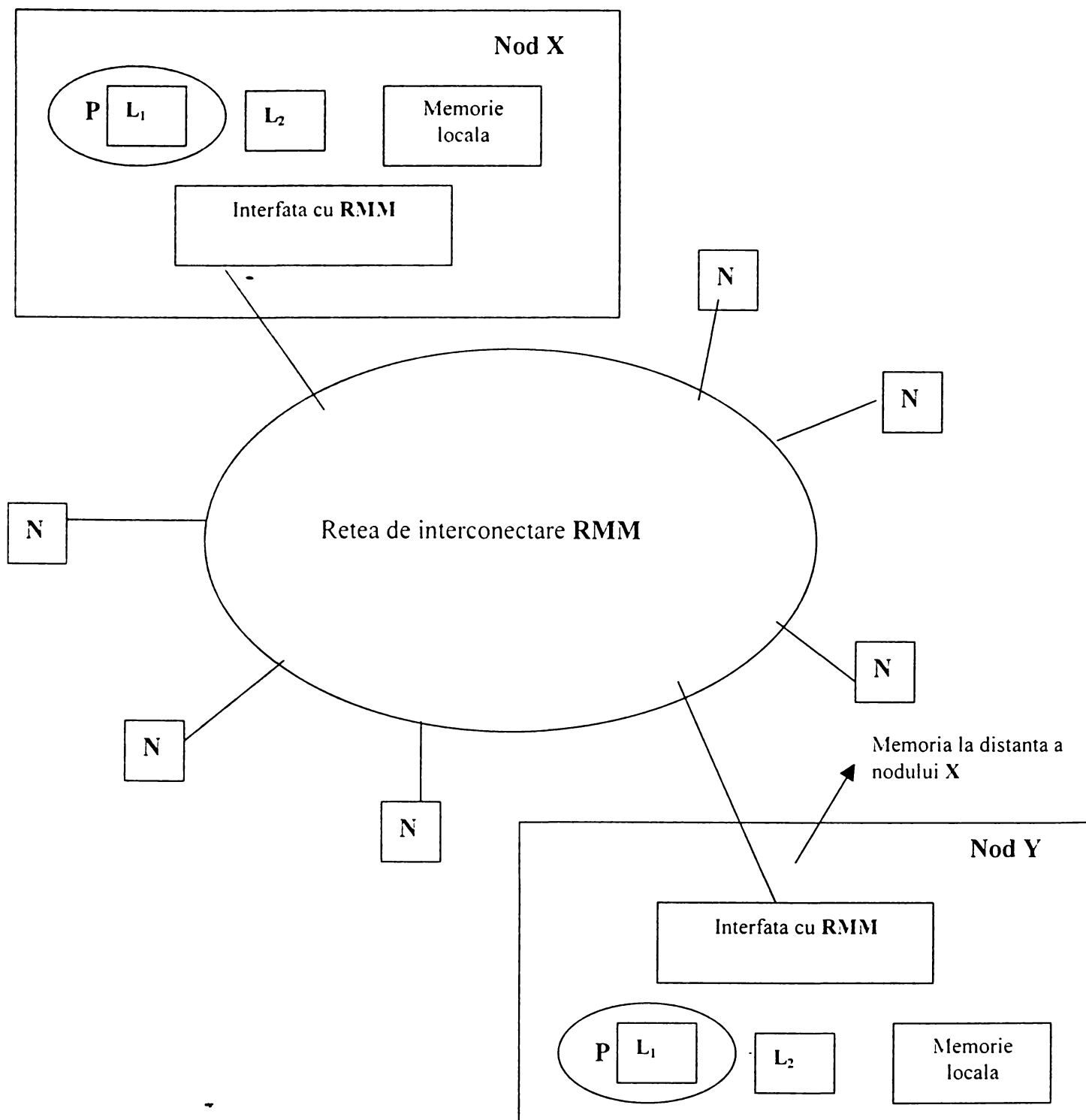


Fig. 6.2 Arhitectura SMP cu evidențierea ierarhiei de memorie

Extinderea ierarhiei memoriei se poate efectua prin plasarea unor capacități suplimentare de memorare (DRAM, SRAM) în fiecare dintre comutatoarele componente ale RIN. Se obține o rețea de interconectare *multinivel cu ierarhie de memorie*, RIN – MIM. Costul unei astfel de soluții constă în complicarea mecanismelor de păstrare a coerenței cache ceea ce impune soluții hardware suplimentare.

Din aceste rațiuni am optat pentru o soluție de compromis între o ierarhie echilibrată de memorie, în termeni de straturi, și o rețea RMM ce permite circulația bidirecțională a mesajelor și în același timp existența unor mecanisme simple de menținere a coerenței.

6.2 STRUCTURA REȚELEI DE INTERCONECTARE DE TIP RMM

În cadrul sistemelor cu arhitectura bazată pe *memorie partajată distribuită* (MPD) modulele de memorie sunt conectate direct la procesoarele corespondente dar spațiul de adresă este partajat. În situația în care RIN este de tip *magistrală ierarhică*, magistrala de pe nivelul superior constituie un potențial de limitare a performanțelor sistemului [BN91], [NB93], [BINK97], [BNA94], [BNT94], [Coif01].

Pentru a ameliora performanța, un număr de magistrale trebuie să fie conectate la nivelul superior utilizând tehnici de interțesere a memoriei. În literatură [BINK 97] se indică o astfel de conectare prezentată în fig. 6.3. Fiecare magistrală împreună cu controlerul asociat se plasează într-un comutator similar cu cel utilizat într-o RIN-MN clasică. Se obține o rețea de interconectare de tip *Rețea Magistrală Multinivel* (RMM).

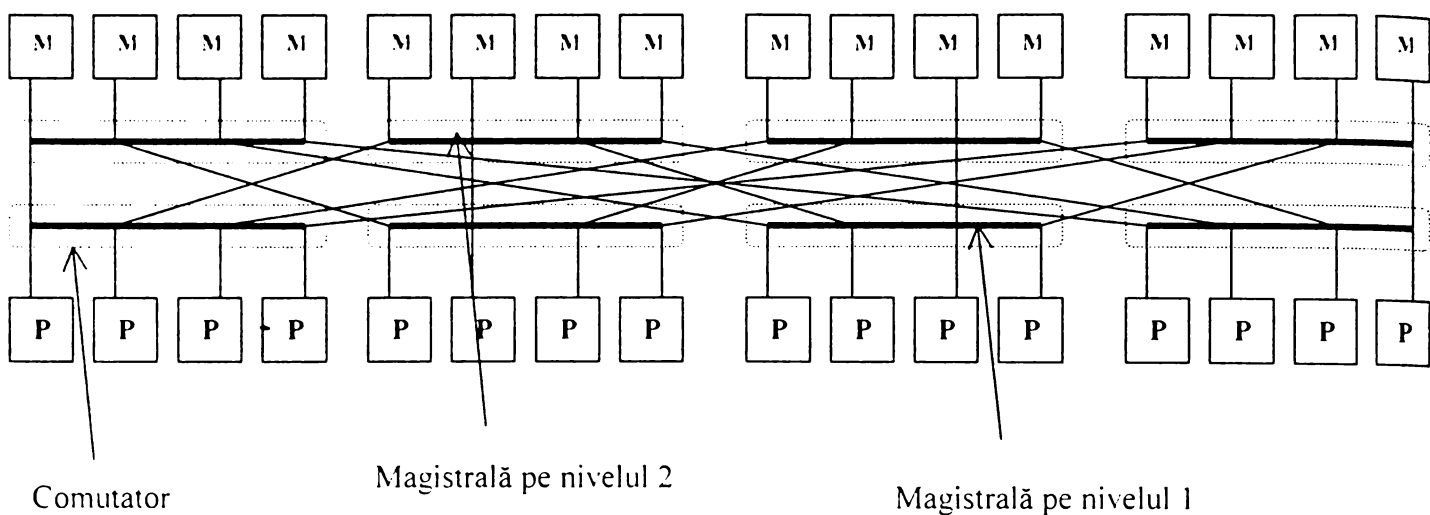


Fig.6.3 Un SMP 16×16 bazat pe RMM utilizând comutatoare 4×4

În general o RMM de dimensiune $N \times N$, utilizând comutatoare $k \times k$, are $l = \log_k N$ niveluri de comutatoare numerotate de la 0 la $l-1$. Pe fiecare nivel sunt plasate N/k comutatoare. Fiecare comutator are k porturi superioare și k porturi inferioare. Interconexiunea dintre două niveluri de magistrale implementează funcția "*k-amestec*".

O astfel de rețea pentru SMP din fig. 6.1 este prezentată în fig. 6.4

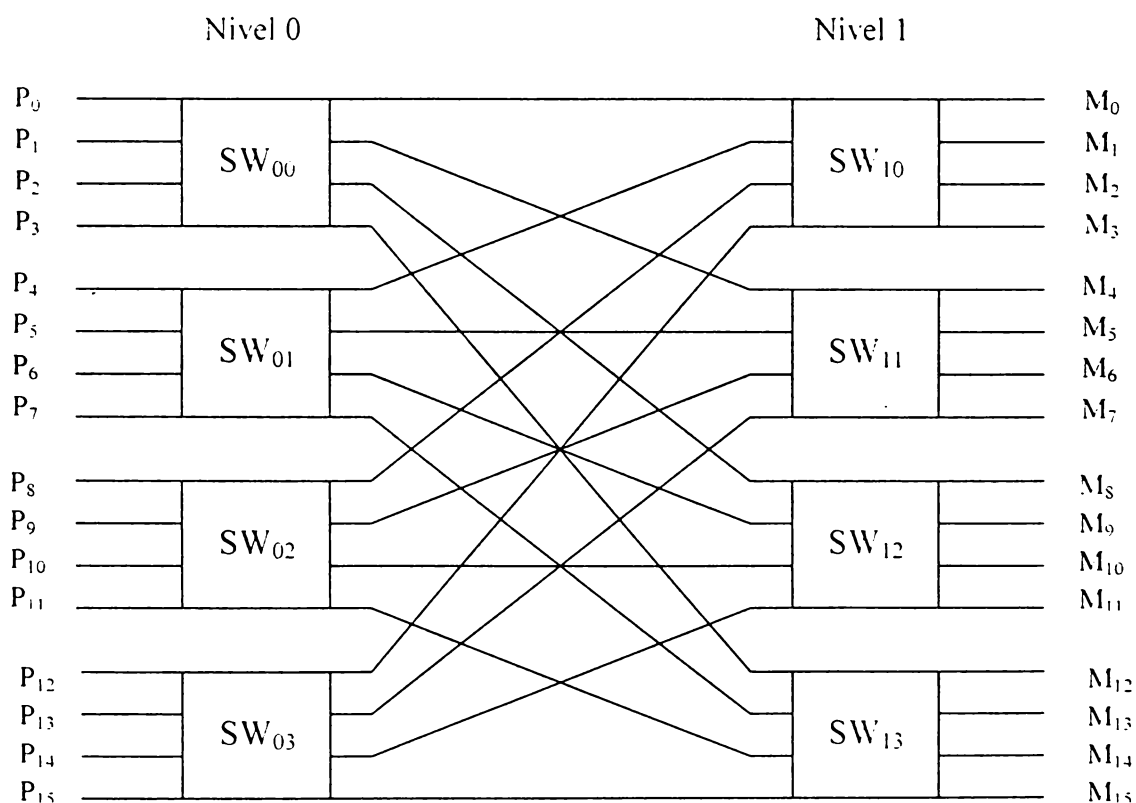


Fig.6.4 RMM 16×16 , utilizând comutatoare 4×4

În acest capitol se analizează performanțele unui SMP ce utilizează o RMM 16×16 , construită cu comutatoare 2×2 . Am optat pentru această soluție deoarece compararea performanțelor acestei rețele față de ale altor RIN-MN este mai facilă atunci când se utilizează comutatoare 2×2 , utilizate în majoritatea sistemelor comerciale.

Prin urmare RIN de tip RMM 16×16 este prezentat în fig 6.5

La acest tip de rețea poate să existe sau poate să nu existe o interconectare de tip “*amestec perfect*” înainte de primul nivel de comutatoare. În primul caz rețeaua este de tip *omega* și nu face obiectul analizei prezentate în acest capitol. Algoritmii de rutare sunt dezvoltati pe o rețea unde nu există implementată funcția “*amestec perfect*” înainte de primul nivel. Prin urmare, un set de procesoare cu memoriile aferente sunt conectate la un comutator din primul nivel și la alt comutator dispus pe ultimul nivel.

Dacă există o conectare “*k - amestec perfect*” înainte de primul nivel, atunci la comutatoarele dispuse pe primul nivel și pe ultimul vor fi conectate un set diferit de procesoare. Indiferent însă de amplasarea conecticii de interconectare sunt posibile patru tipuri de rutări ce vor fi descrise în continuare.

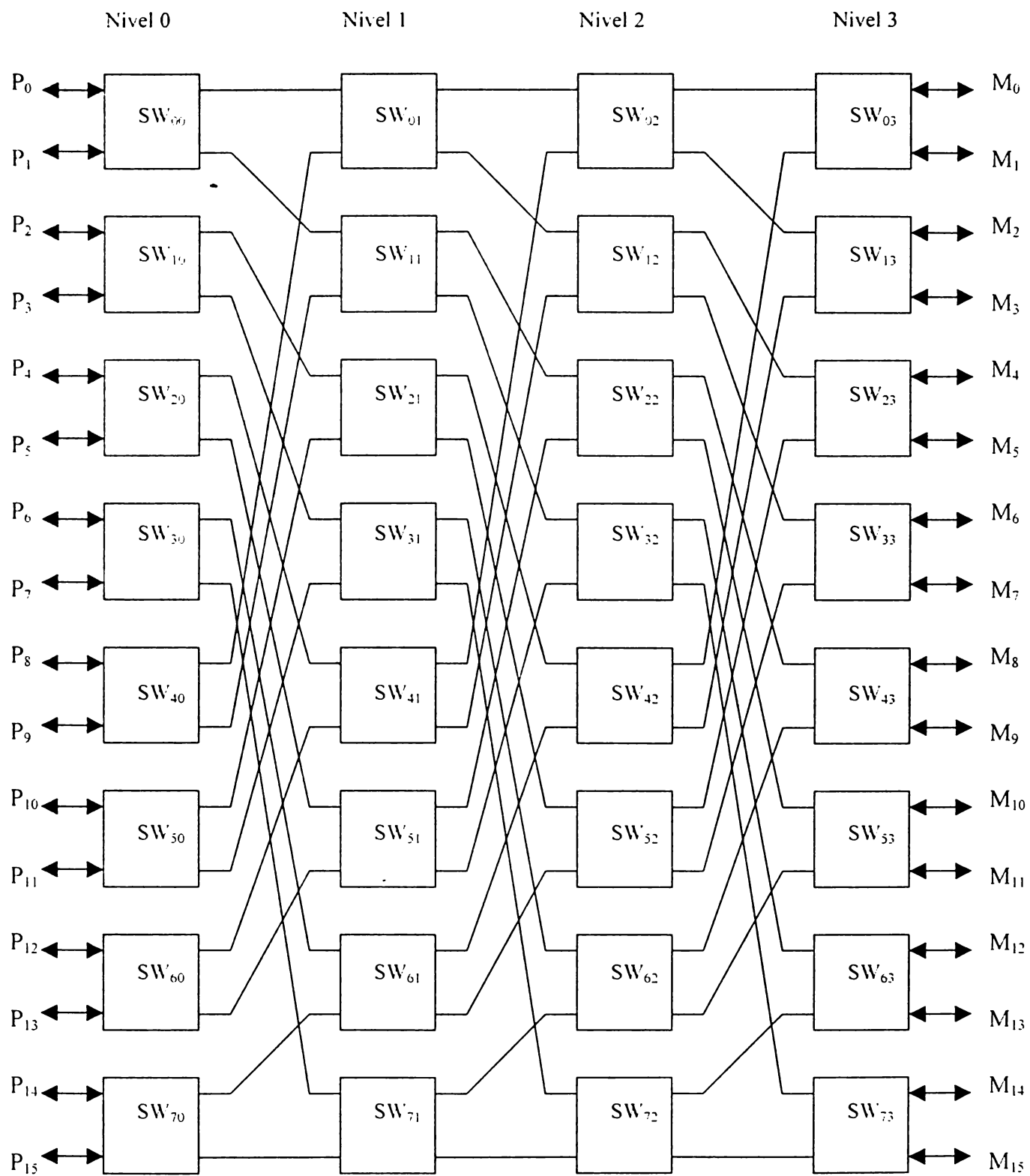


Fig.6.5 RMM 16x16, utilizând comutatoare 2x2

RMM este o rețea bidirecțională, permițind trecerea datelor atât în sens direct cât și în sens invers. Fiecare comutator are un set de conexiuni plasate în stânga mai aproape de procesoare și un set de conexiuni plasate în dreapta mai aproape de memorii.

Structura internă a unui comutator 2x2, utilizat în RMM, este prezentată în fig.6.6

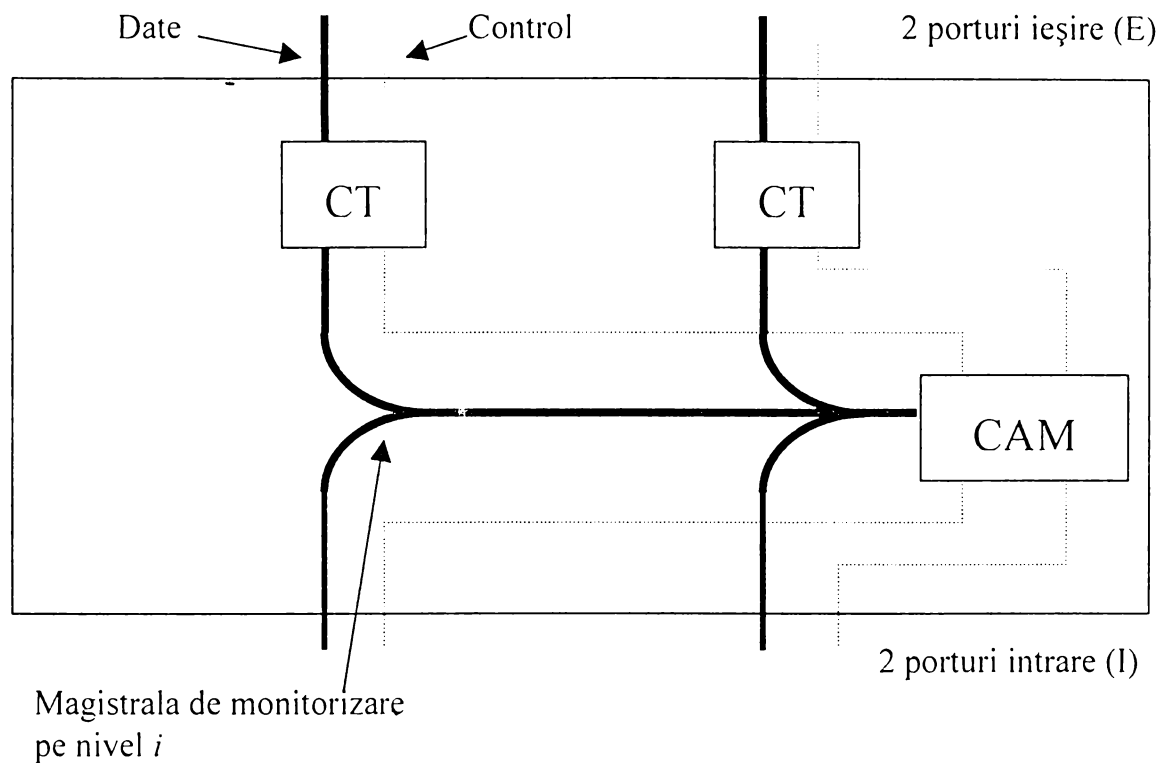


Fig.6.6 Diagrama bloc a unui comutator RMM 2x2

Comutatorul constă din următoarele:

- o magistrală de monitorizare (MM)
- un controler de acces la magistrală (CAM)
- un controler de trafic (CT)

Diagrama bloc a unui CT este prezentată în fig. 6.7. Blocurile CT din componența comutatoarelor RMM păstrează un director al informației de stare asupra tuturor blocurilor partajate din sistem care l-ar putea parcurge în concordanță cu conexiunile comutatoarelor. Un director nu conține însă nici o informație referitoare la blocurile private.

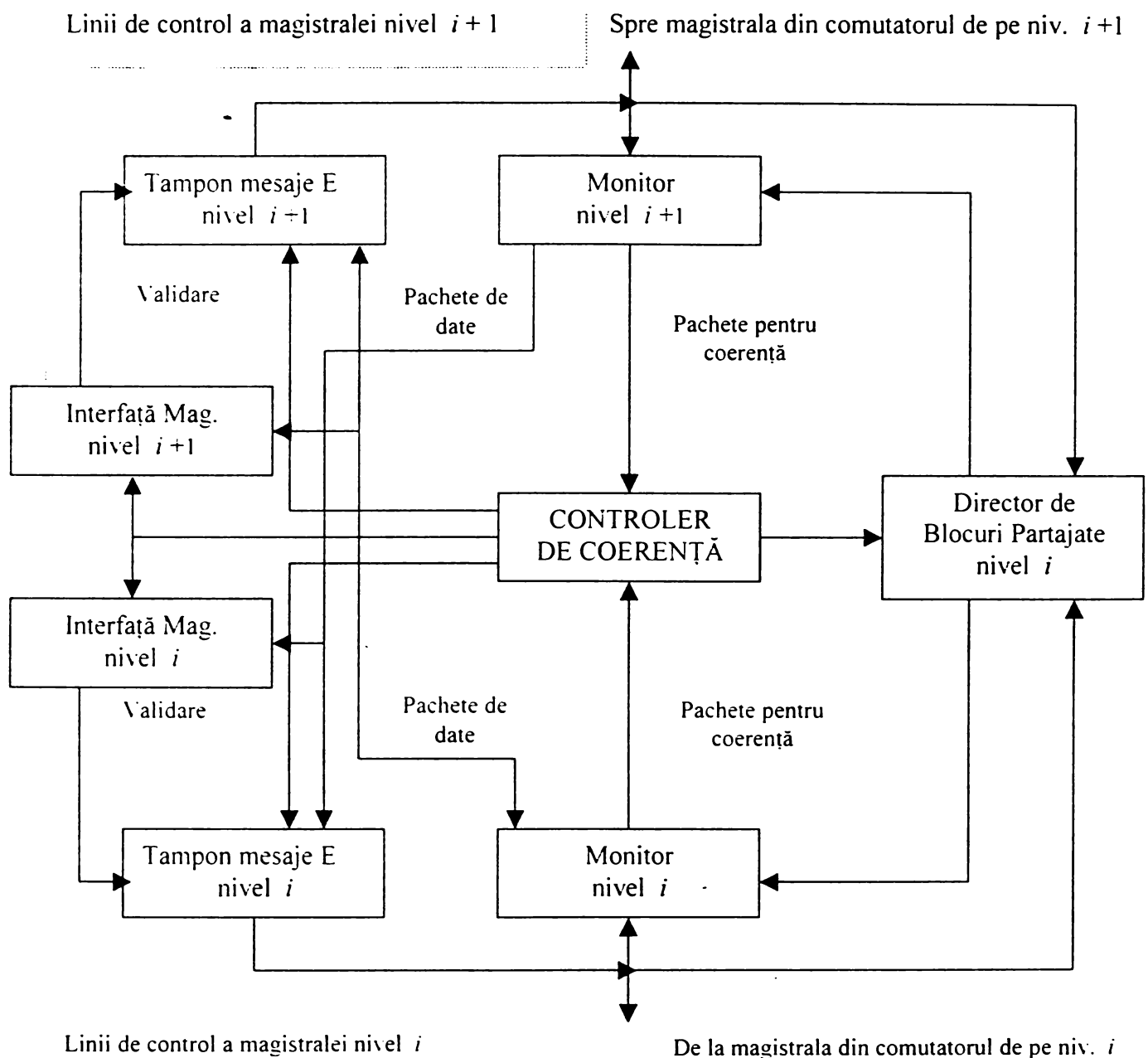


Fig.6.7 Diagrama bloc a unui controler de trafic (CT)

Blocul CT constă dintr-un comutator plasat pe nivelul i ce utilizează informația din director pentru a filtra cererile provenind de la magistrala plasată pe nivelul $i+1$ spre magistrala de pe nivelul i și vice versa.

Leșirile unui CT sunt conectate la magistralele de pe următorul nivel superior utilizând aceeași configurație de conectare ca și o RIN-MN convențională. Accesul la magistrala de monitorizare este comandat de un *controler de acces la magistrala* (CAM). CAM

primește semnale reprezentând cereri de acces la magistrală din parte blocurilor CT conectate la magistrală și validează semnale de garantare a accesului către solicitanți. Există mai multe politici de arbitrare. În general un grup de magistrale plasate pe nivelul i este conectat la o magistrală de pe nivelul $i+1$ prin intermediul unor CT individuale plasate pe nivelul i . Blocurile CT situate pe nivelul cel mai de sus sunt conectate direct la modulele de memorie corespunzătoare. Controlerele memoriilor cache din nodurile procesoare monitorizează magistralele de pe nivelul i . Un procesor lansează o cerere de acces la memorie spre cache-ul aparținător. Dacă este o solicitare pentru o dată privată și dacă este un succes, atunci blocul de date este furnizat imediat. Dacă este un eșec, atunci solicitarea este avansată spre modulul de memorie corespunzător prin intermediul RMM. Blocul CT manipulează solicitările pentru date private în mod separat așa că acestea nu fac subiectul acțiunilor de menținere a coerenței.

Pentru blocurile partajate, este permisă coexistența unor copii multiple în diferite memorii cache locale. Coerența acestora este menținută prin protocoale pentru copii multiple. Un astfel de protocol va fi descris ulterior.

Funcționarea diferitelor module dintr-un CT este următoarea. Fiecare CT conține două *module de monitorizare (Snoopers)* (MM). Cele situate pe nivelul i și pe nivelul $i+1$ monitorizează magistralele corespondente. Sarcina lor este supravegherea tuturor tranzacțiilor de pe magistrală și luarea de măsuri adecvate. Dacă cererea este pentru un bloc privat ce trebuie să parcurgă printr-un anume CT determinat de etichetă, atunci MM plasează cererea *tamponului de mesaje de ieșire* superior. Dacă însă cererea este pentru un bloc partajat atunci MM citește directorul pentru a determina dacă trebuie luată o acțiune de menținere a coerenței. Apoi plasează cererea *controlerului de coerență* (CCH), dacă ea este solicitată. Acțiunea specifică de menținere a coerenței depinde de starea blocului din comutatorul respectiv. Un CCH execută acțiunile necesare pentru a implementa un anumit protocol de menținere a coerenței. El modifică directorul CT-ului cum este necesar, corespunzător mesajului de intrare recepționat și regulilor protocolului. Apoi transmite un mesaj următorului nivel, dacă este solicitat.

Interfețele cu magistrala (IM) gestionează accesele la magistralele respective. Când există un mesaj în coada de așteptare de ieșire ce trebuie transmis unei magistrale, atunci interfața trimite un semnal de solicitare a magistralei spre CAM și după ce câștigă accesul, validează tamponul de ieșire pentru a plasa mesajul pe magistrală.

Din punct de vedere al modelării prin lanțuri Markov, un comutator RMM este reprezentat ca și în fig. 6. 8.

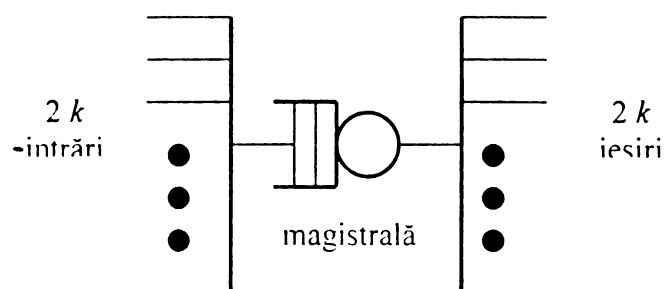


Fig.6.8 Înlănțuirea pachetelor într-un comutator RMM

Comutatoarele RMM pot efectua conectări în sens *direct (înainte)*, *invers (înapoi)* și cu *buclare*.

În general procesoarele P_0, P_1, \dots, P_{N-1} sunt conectate la conexiunile din stânga comutatoarelor RMM de pe nivelul 0. Modulele de memorie M_0, M_1, \dots, M_{N-1} sunt conectate la conexiunile din dreapta ale comutatoarelor dispuse pe nivelul $l-1$. Un modul de memorie M_i este de asemenea conectat la un procesor P_i și este o *memorie locală* pentru P_i .

O *sursă* este identificată printr-o etichetă sursă $S = s_0 s_1 \dots s_j \dots s_{l-1}$, iar o *destinație* printr-o etichetă de destinație $D = d_0 d_1 \dots d_j \dots d_{l-1}$ unde s_j și d_j sunt cifre în sistemul k de numerație.

Cifrele s_0 și d_0 sunt cele mai semnificative și s_{l-1} și d_{l-1} sunt cel mai puțin semnificative.

Conectica între niveluri în cadrul unei RMM implementează funcția "*k-amestec*", ceea ce înseamnă că terminalul din dreapta de la poziția $a_0 a_1 \dots a_{l-1}$ din nivelul i este conectat cu terminalul din stânga la poziția $a_l \dots a_{l-1} a_0$ de pe nivelul $i+1$ pentru $i=0, 1, \dots, l-2$.

O cerere de acces la memorie este satisfăcută intern de către memoria locală când eticheta sursă și eticheta destinație a unei cereri sunt identice. Dacă etichetele sunt diferite, solicitarea parcurge rețeaua spre o memorie situată la distanță. În fig. 6.5 se prezintă ca și exemplu o RMM 16×16 cu comutatoare 2×2 . Înainte de nivelul 0 nu există implementată funcția "*k-amestec*". Astfel un set de procesoare cu memoriile aferente sunt conectate la un comutator pe primul nivel și la alt comutator în aceeași poziție pe ultimul nivel.

În fig.6.9 o cerere traversează în *direcția înainte* sau în *sens direct* când ea pornește de la un procesor și parcurge nivelurile $0, 1, \dots, (l-1)$ în această ordine.

În fig. 6.10 se prezintă *rutarea înapoi* sau *sensul invers* de parcurgere. Cererea pornește de la memorie și parcurge nivelurile $(l - 1), \dots, 1, 0$ în această ordine. Un pachet poate de asemenea să traverseze de la stânga la dreapta și să facă o buclă, denumită în formă de U, într-un nivel intermediar ca și în fig. 6.11. Aceasta este o *buclare în sens direct*. Similar în fig. 6.12 se prezintă o *buclare în sens invers*.

Cele patru tehnici de rutare prezentate anterior asigură patru căi distincte între o sursă și o destinație în RMM. Ca și rezultat, toleranța la defect și fiabilitatea RMM sunt mult mai bune decât în cazul unei RIN-MN convenționale. În cazul rețelelor convenționale ca și Omega, Delta, etc., eticheta destinație este utilizată în scopul autorutării unei cereri doar în direcția înainte - sensul direct. În cazul RMM, eticheta destinație poate fi de asemenea utilizată pentru autorutare în sensul direct. Deoarece conexiunile de pe nivelul 0 sunt directe în loc să implementeze funcția "*k-amestec*", eticheta destinație propriu-zisă nu poate fi utilizată pentru autorutare în sensul invers. Eticheta de rutare în cazul rutării în sens invers -RI- se obține prin aplicarea funcției "*amestec invers*" asupra etichetei destinație, pe o cifră.

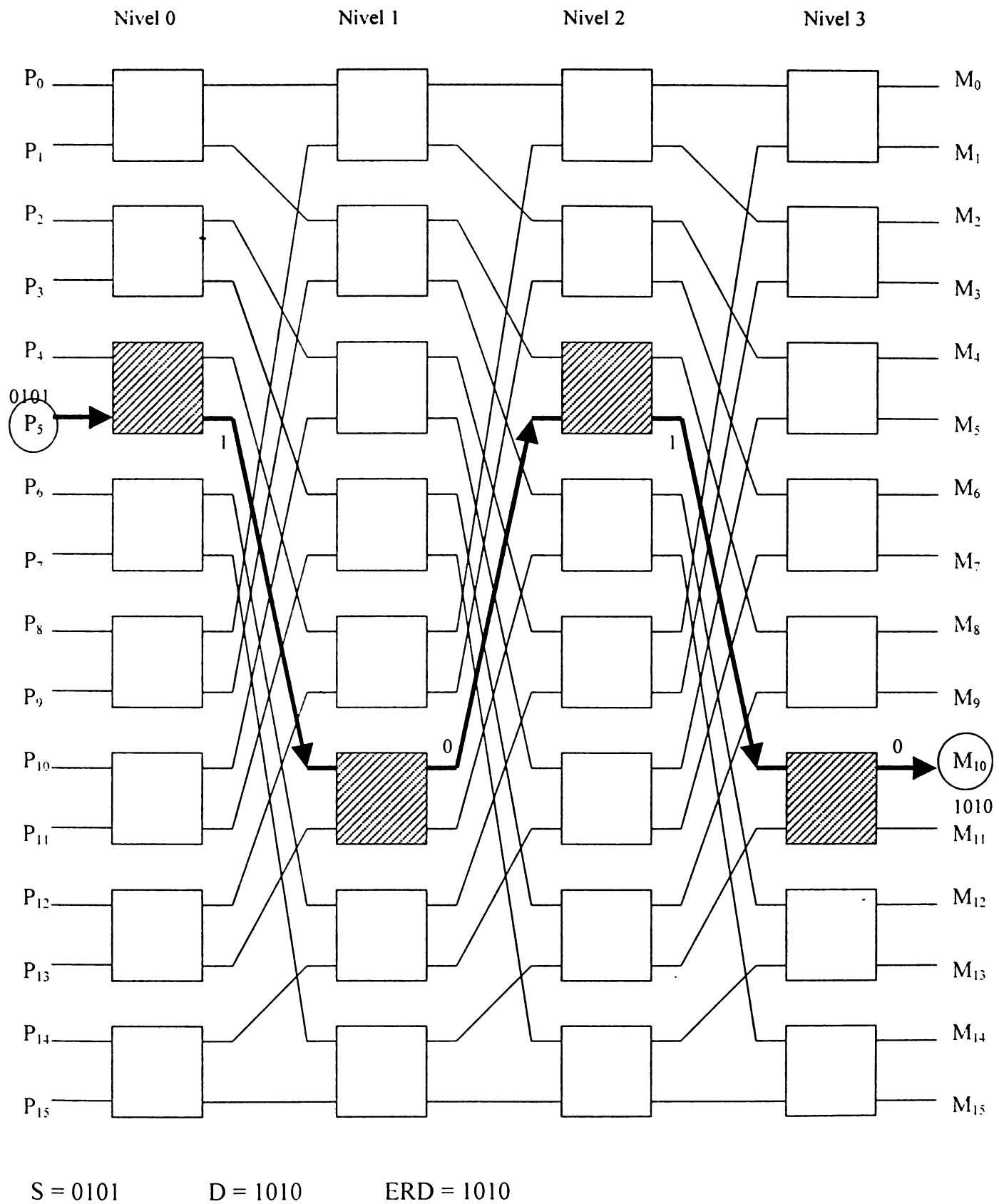


Fig.6.9 Rutarea în sens Direct în RMM

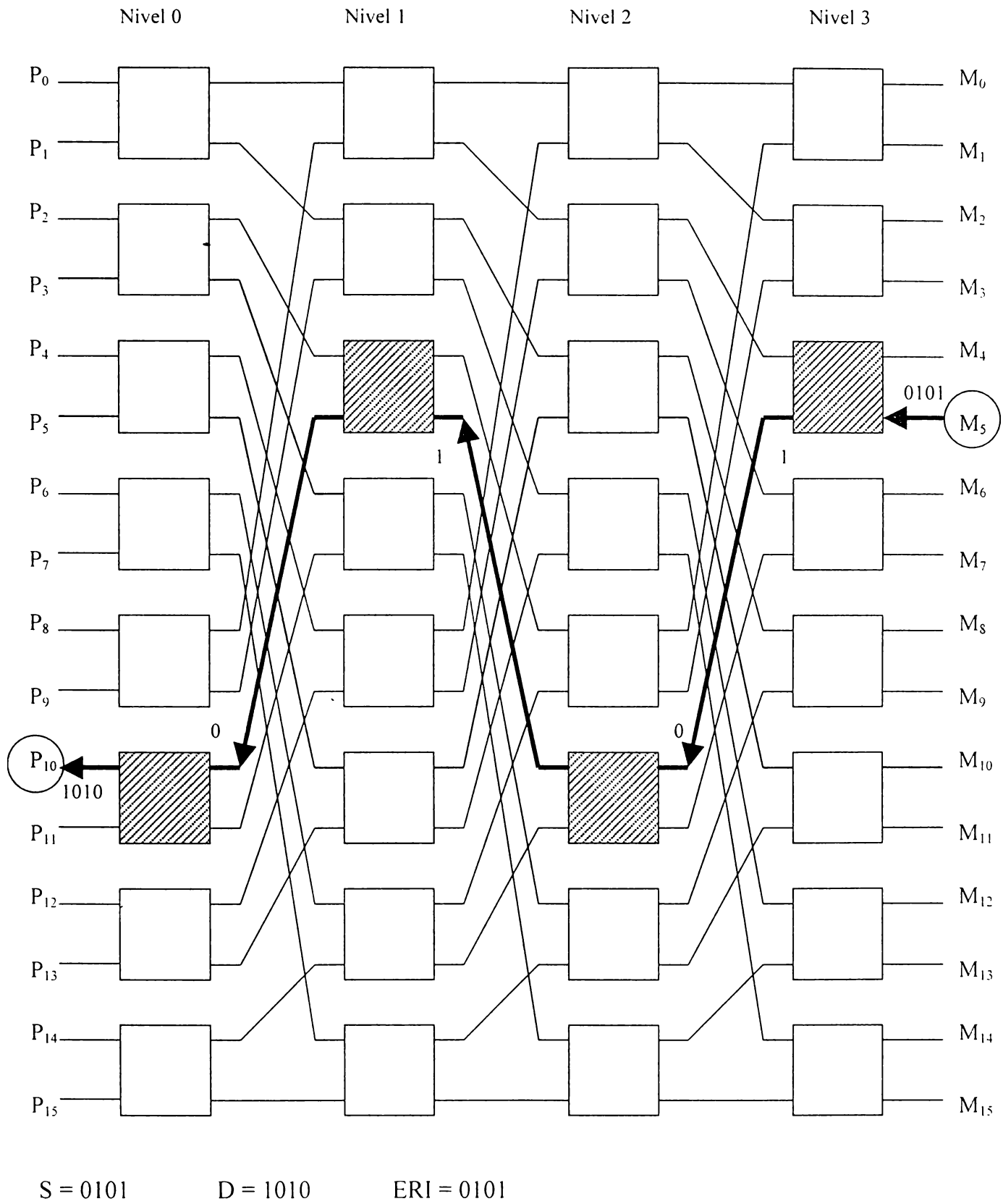


Fig.6.10 Rutarea în sens Invers în RMM

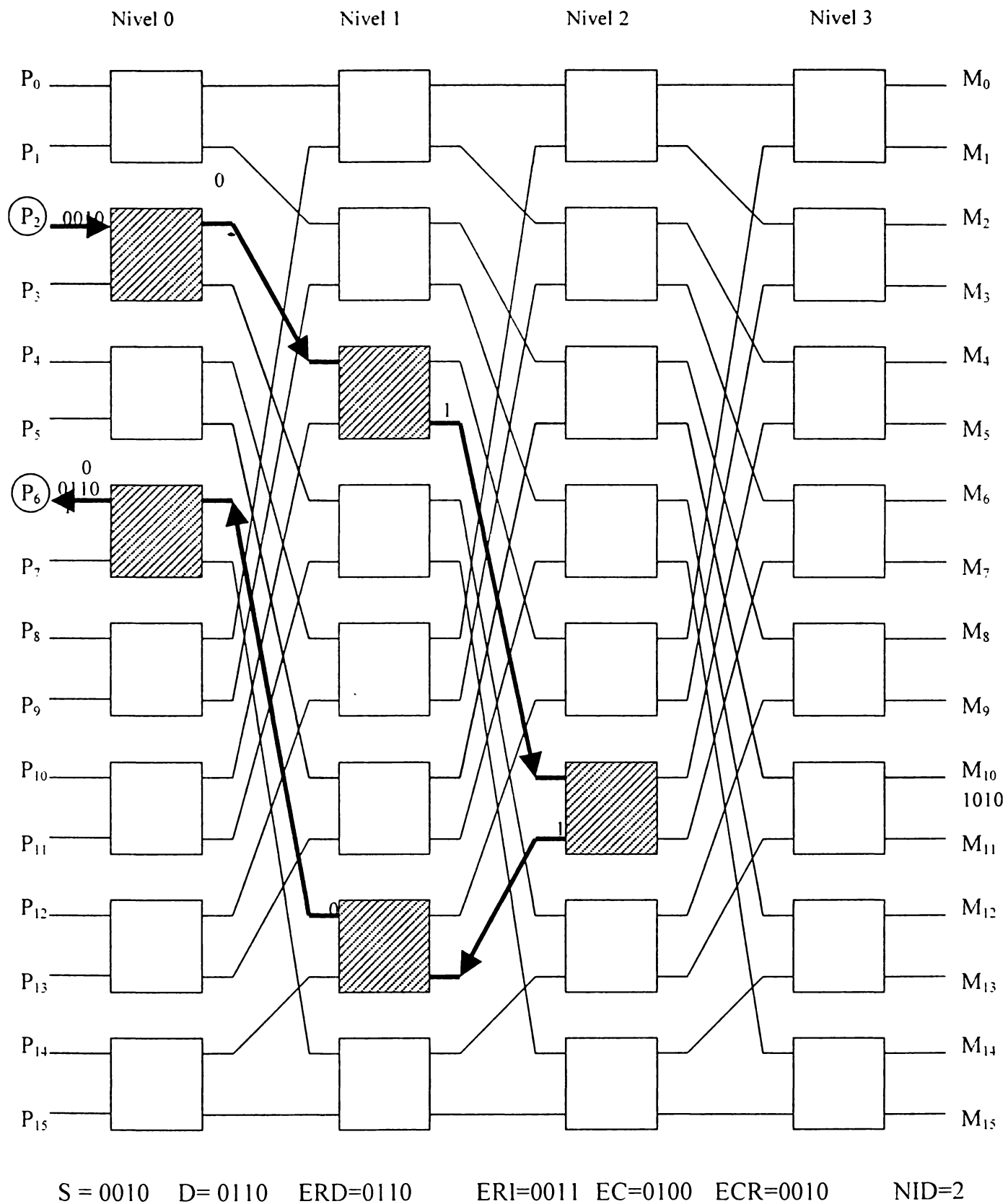


Fig.6.11 Buclare în sens direct (BD)

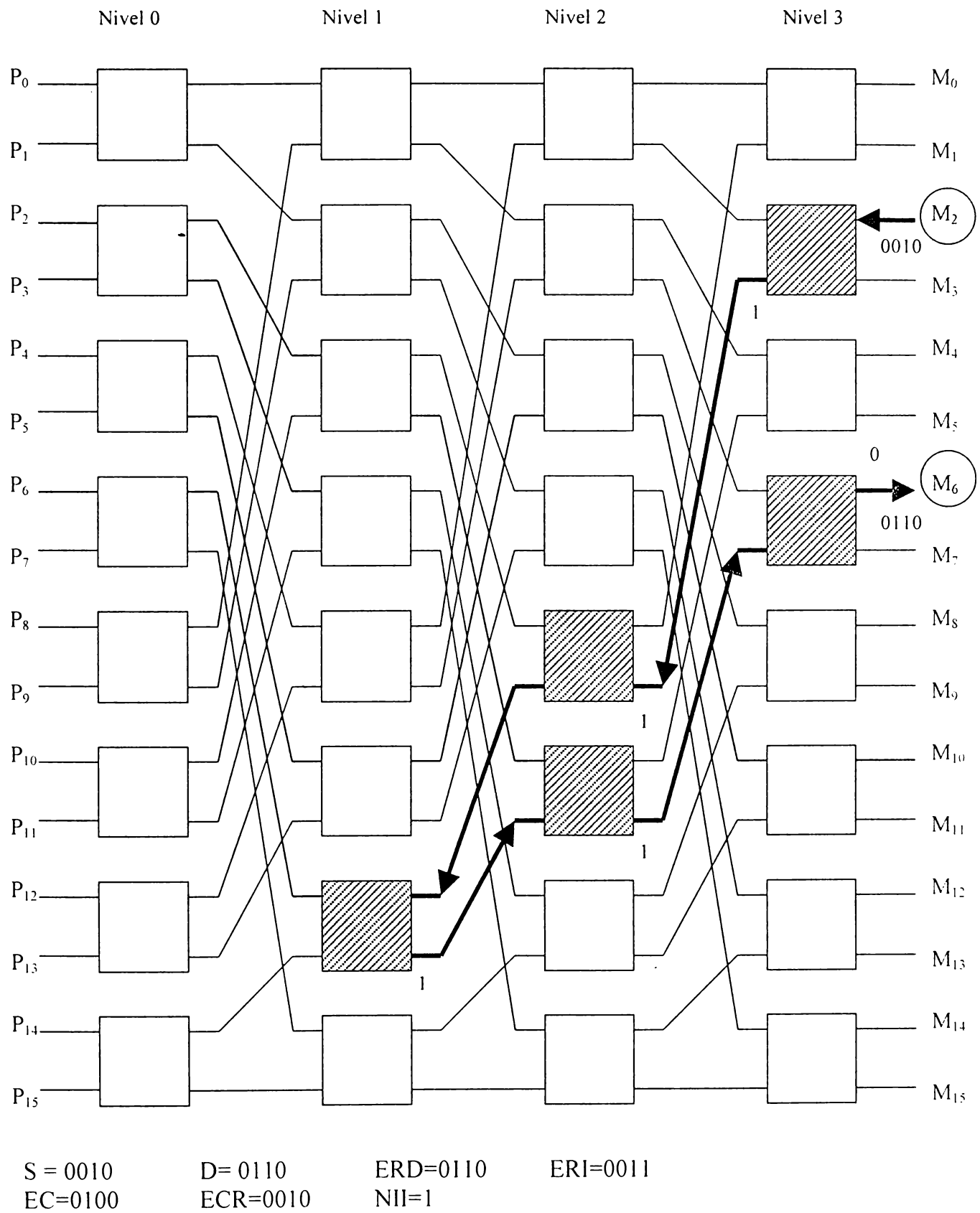


Fig.6.12 Buclare în sens invers (BI)

DEFINIȚIA 6.1 (ERD). *Eticheta de Rutare în sens Direct (ERD) este aceeași ca și eticheta destinație a unei cereri de acces la memorie.*

$$ERD = d_0d_1 \dots d_{l-1}$$

DEFINIȚIA 6.2 (ERI). *Eticheta de Rutare în sens Invers (ERI) este eticheta destinație asupra căreia a operat funcția "amestec" pe o cifră.*

$$ERI = b_0b_1 \dots b_{l-1} = d_{l-1}d_0d_1 \dots d_{l-2} \text{ unde: } b_j = d_{(j-1) \bmod l}$$

DEFINIȚIA 6.3 (EC). *Eticheta Combinată (EC) este eticheta obținută prin aplicarea funcției SAU-EXCLUSIV între cifrele sursă și a celei destinație*

$$EC = c_0, c_1, \dots, c_j, \dots, c_{l-1} \quad \text{unde } c_j = s_j \oplus d_j$$

De notat că în timp ce cifrele în etichetele sursă (S) respectiv destinație (D) sunt în baza k , cifrele din EC sunt în baza 2.

DEFINIȚIA 6.4 (ECR). *Eticheta Combinată Rotită (ECR) este eticheta combinată (EC) asupra căreia s-a aplicat funcția "amestec invers" sau rotită spre dreapta cu o cifră.*

$$ECR = r_0r_1, \dots, r_{l-1} = c_{l-1}c_0, c_1, \dots, c_j, \dots, c_{l-2} \quad \text{unde } r_j = c_{(j-1) \bmod l}$$

DEFINIȚIA 6.5 (NID). *Nivelul de Întoarcere spre sensul Direct (NID) - înainte - este poziția non-zero.*

Aceasta este $NID = m$, astfel că $r_m = 1$ și $r_j = 0$ pentru $m < j \leq l-1$.

DEFINIȚIA 6.6 (NII). *Nivelul de Întoarcere spre sensul Invers (NII) - înapoi - este poziția non-zero cea mai din stânga în Eticheta Combinată (EC). Aceasta este $NII = n$, astfel ca $c_n = 1$ și $c_j = 0$ pentru $0 \leq j < n$.*

Etichetele de rutare ERD și ERI sunt utilizate pentru autorutare în cazul direcțiilor *înainte*, respectiv *înapoi*. Etichetele EC și ECR sunt utilizate pentru a identifica nivelurile U de întoarcere de tip NID sau NII. Nivelurile NID și NII sunt utilizate pentru a determina unde se produc întoarcerile *înainte* sau *înapoi* în timpul parcurgerilor de tip U.

6.2.1 Algoritmi de rutare pentru RMM.

Există mai multe tehnici de rutare pentru RMM. În literatură [BINK97], [CWH01], [Coif01] se prezintă algoritmi de selectare a *căii cu distanță minimă* în cadrul unei RMM.

6.2.1.1 Tehnici de rutare

A. Rutare în sens direct (RD)

În această rutare, o solicitare de acces la memorie generată de procesorul S traversează rețeaua de la nivelul 0 până la nivelul $l-1$ spre memoria de destinație D.

Un exemplu de RD de la sursa 0101 la destinația 1010 este prezentat în fig. 6.9.

Cifra de indice j a etichetei ERD este utilizată de către un comutator amplasat pe nivelul j , pentru autorutare. De obicei se utilizează comutatoare 2×2 și atunci cifra are reprezentare binară.

Astfel o cerere de acces ce provine din poziția $S = s_0s_1\dots s_j\dots s_{l-1}$ din stânga nivelului 0 este comutată spre poziția $s_0s_1\dots s_j\dots s_{l-2}d_0$ spre dreapta nivelului 0. apoi traversează o conecție " k -amestec" și ajunge la poziția $s_1s_2\dots s_j\dots s_{l-2}d_0s_0$ la intrarea nivelului 1 și este comutată la $s_1s_2\dots s_j\dots s_{l-2}d_0d_1$ de la ieșirea nivelului 1.

În general, când o cerere sosește în poziția $s_j s_{j+1} \dots s_{l-2} d_0 d_1 \dots d_{j-1} s_{j-1}$ la stânga nivelului j , ea este comutată spre poziția $s_j \dots s_{l-2} d_0 d_1 \dots d_{j-1} d_j$ din dreapta nivelului j , parcurge o conecție " k -amestec" (excepție în cazul ultimului nivel, $j = l - 1$) și sosește în poziția

$s_{j-1} \dots s_{l-2} s_{j-1} \dots s_{l-2} d_0 d_1 \dots d_j s_j$ din stânga nivelului $j+1$. În final, ea atinge destinația $d_0 d_1 \dots d_j \dots d_{l-1}$ la ieșirea ultimului nivel al RMM.

B. Rutare în sens Invers (RI).

În această tehnică o cerere provenind de la nodul sursă S se deplasează înapoi de la nivelul $l-1$ prin nivelul 0 spre nodul de destinație D. Un exemplu de RI de la sursa 0101 spre destinația 1010 este prezentat în fig. 6.10. Cifra de indice j din eticheta ERI este utilizată de către un comutator de pe nivelul j pentru autorutare. Uzual se operează cu bitul j în cazul comutatoarelor 2×2 .

O cerere ce a fost lansată din poziția $S = s_0s_1 \dots s_j \dots s_{l-1}$ la dreapta nivelului $l-1$, este comutată spre poziția $s_0s_1 \dots s_j \dots s_{l-2}d_{l-2}$ la stânga nivelului $l-1$, apoi parcurge o conecție care implementează funcție "*k-amestec invers*" și atinge poziția $d_{l-2}s_0s_1 \dots s_j \dots s_{l-2}$ la dreapta nivelului $l-2$. În general când o cerere sosește în poziția $d_jd_{j-1} \dots d_{l-2}s_0s_1 \dots s_{j-1}s_j$ la dreapta nivelului j , ea este comutată spre poziția $d_jd_{j-1} \dots d_{l-2}s_0s_1 \dots s_{j-1}d_{(j-1) \bmod l}$. Apoi parcurge o conecție "*amestec invers*", cu excepția ultimului nivel $j=0$, și sosește în poziția $d_{j-1}d_jd_{j-1} \dots d_{l-2}s_0s_1 \dots s_{j-1}$ din stânga nivelului $j-1$.

C. Rutare cu Buclare în sens Direct (BD).

Solicitarea pornește de la sursa S dispusă pe nivelul 0, urmează o rutare RD (utilizând ERD) până la nivelul $NID-1$ și atinge partea stângă a nivelului NID în poziția $s_{NID}s_{NID-1} \dots s_{l-2}d_0d_1 \dots d_{NID-1}s_{NID-1}$.

La NID, schimbă de sens formând o buclă în formă de U, în loc de-a comuta spre dreapta nivelului NID. Solicitarea este comutată spre stânga nivelului NID în poziția $s_{NID}s_{NID-1} \dots s_{l-2}d_0d_1 \dots d_{NID-1}d_{(NID-1) \bmod l}$ și urmează o rutare RI (utilizând ERI) până la nivelul 0. În final atinge poziția $d_0d_1 \dots d_{NID-1}s_{NID} \dots s_{l-2}d_{l-1}$ în stânga nivelului 0. În fig. 6.11 se indică o astfel de rutare pentru 0010 la 0110.

D. Rutare cu Buclare în sens Invers (BI)

În BI, solicitarea este lansată de sursa S plasată pe nivelul $l-1$, urmează RI (utilizând ERI) până la nivelul $NII+1$ și atinge poziția $d_{NII}d_{NII-1} \dots d_{l-2}s_0s_1 \dots s_{NII-1}s_{NII}$ la dreapta nivelului NII. La nivelul NII, formează o buclă U, în loc să comute spre stânga lui NII.

Solicitarea este comutată spre dreapta nivelului NII în poziția $d_{NII}d_{NII-1} \dots d_{l-2}s_0s_1 \dots s_{NII-1}d_{NII}$ și urmează un traseu RD (utilizând ERD) până la nivelul $l-1$. În final, solicitarea atinge poziția $s_0s_1 \dots s_{NII-1}d_{NII}d_{NII-1} \dots d_{l-1}$ în dreapta nivelului $l-1$. Un exemplu de astfel de rutare de la 0010 la 0110 este arătat în fig.6.12.

6.2.1.2a Algoritmul de determinare a numărului de căi din fiecare tip de rutare între surse și destinații

Algoritmul pentru determinarea numărului de căi din fiecare tip de rutare între orice sursă și destinație este un algoritm de tip *backtracking*. Scopul este verificarea și măsurarea unor proprietăți și nu realizarea unei unelte de simulare. Pentru fiecare nod sursă se determină toate căile posibile la fiecare destinație folosind RD, RI, BD, BI. Pentru aceasta se explorează întreg spațiul de stări, selectându-se doar stările valide. Concomitent se determină și lungimea minimă a fiecărei căi pentru toate perechile sursă-destinație. Vectorul de stare este reprezentat printr-o stivă.

initializare stiva cu nodul sursa

k=1

LOOP WHILE (k >= 0)

LOOP

 are_succesor = succesor(k)

IF (are_succesor) **THEN**

 este_valid = valid(k)

UNTIL (NOT are_succesor OR (are_succesor AND este_valid))

IF (are_succesor) **THEN**

IF (solutie(k)) **THEN**

 inregistrare

ELSE

 k = k + 1

ELSE

 k = k - 1

succesor(k) determină succesorul nodului din poziția *k*. Succesorul poate fi unul din cele patru noduri, corespunzătoare direcțiilor posibile de deplasare. Când un nod nu mai are succesor se coboară în stivă la nodul precedent și se reia explorarea ($k=k-1$).

valid(k) verifică validitatea vectorului parțial de stare curent. Vectorul este valid dacă nu s-au depășit marginile, dacă nu sunt cicluri, dacă nodul $k+1$ nu coincide cu nodul $k-1$, dacă sunt mai multe bucle.

solutie(k) verifică dacă s-a ajuns la o soluție, adică dacă s-a ajuns la un nod destinație. În caz afirmativ se va incrementa contorul de căi în funcție de tipul căii și totodată se va reține lungimea căii minime. Dacă nu s-a ajuns încă la o soluție se continuă algoritmul pentru următorul nod ($k=k+1$).

Programul de implementare al acestui algoritm și rezultatele sunt prezentate în ANEXA 1, ANEXA 2 și în [CWH01].

6.2.1.2 b Algoritm de stabilire a căii optime

Distanța dintre o sursă și o destinație în RMM se definește ca fiind numărul minim de comutatoare prin care un pachet de date trebuie să traverseze rețeaua.

Pentru o RIN-MN convențională, distanța este întotdeauna l . În cazul unei RMM, distanța poate fi mai mică decât l dacă se aleg tehnici de rutare BD sau BI. Acestea sunt utilizate atunci când nivelul din rețea unde se produce întoarcerea - buclarea - este mai mic decât nivelul central al rețelei. Prin urmare se obțin beneficii nete în termeni de distanță între o sursă dată și toate destinațiile.

Un algoritm pentru a selecta rutarea optimă pentru o pereche sursă - destinație dată, este următorul:

1. $S = s_0s_1\dots s_{l-1}$
2. $D = d_0d_1\dots d_{l-1}$
3. $EC = S \oplus D = c_0, c_1, \dots, c_j, \dots, c_{l-1}$
4. $ECR = c_{l-1}, c_0, c_1, \dots, c_j, \dots, c_{l-2}$
5. $d_l = \lfloor l/2 \rfloor, d_u = \lceil l/2 \rceil$
6. **IF** (sursa = destinația)
7. **THEN** cererea de acces este pentru memoria locală
8. **ELSE**
9. Caută NID și NII, în funcție de ECR și EC
10. **IF** ($NID = (l-1 - NII) = 0$)
11. **THEN** selectează rutarea BD sau rutarea BI

12. **ELSE IF** ($NID < d_l$)
13. **THEN** selectează rutarea BD
14. **ELSE IF** ($NID \geq d_u$)
15. **THEN** selectează rutarea BI
16. **ELSE**
17. selectează rutarea RD sau RI

Algoritmul de stabilire a căii optime selectează un traseu care are o lungime minimă de cale în termeni de număr de comutatoare traversate.

Fiind dată o sursă $S = s_0s_1\dots s_{l-1}$ și o destinație $D = d_0d_1\dots d_{l-1}$ acest algoritm calculează etichetele descrise anterior. Se compară apoi etichetele pentru a decide care din cele patru tehnici de rutare va oferi lungimea minimă prin rețea. În acest algoritm d este definit ca și nivelul central a RMM. Trebuie remarcat că rutarea optimă între două noduri este fixată într-o rețea dată. Astfel, calea optimă poate fi precalculată și memorată într-un tabel ce poate fi citit când se generează o cerere. Nu este necesar ca acest algoritm să fie executat ori de câte ori un mesaj este extras.

Dacă sursa și destinația sunt identice, atunci cererea este pentru memoria locală. În acest caz, nu se solicită nici o traversare prin RMM. Toate celelalte cereri traversează cel puțin un nivel din RMM.

Memoriile ce sunt conectate la un procesor prin primul sau ultimul nivel al RMM sunt denumite *memorii de grup (cluster memories)*. Similar, procesoarele ce au cel puțin un comutator între ele și memoriile formează *procesoare de grup (cluster processors)* al acestor memorii. Un grup de procesoare conectat la un comutator de pe nivelul 1, îndeplinește această condiție. Cererile de acces la memoriile unui grup necesită ca doar un comutator să fie traversat. Astfel când ($NID = (l-1-NII) = 0$), se optează pentru rutări BD sau BI.

Dacă $NID < \lfloor l/2 \rfloor$, sau $NII \geq \lceil l/2 \rceil$, există niveluri de întoarcere în RMM înaintea sau după nivelul din centrul RMM.

Aceasta va reduce lungimea totală a căii la mai puțin de l și astfel rutările BD sau BI sunt selectate. Dacă nici una din aceste condiții nu sunt îndeplinite, atunci $NID \geq \lceil l/2 \rceil$, și $NII < \lfloor l/2 \rfloor$. În acest caz rutările RD sau RI sunt singurele opțiuni posibile.

Lungimile de cale, în termeni de număr de comutatoare traversate, sunt următoarele:

- Memorie locală : 0 comutatoare (RMM nu e traversată)
- Rutare RD sau RI : l comutatoare
- Rutare BD:
 - memorii grupate : 1 comutator
 - alte memorii : $2 \times NID + 1$ comutatoare
- Rutare BI :
 - memorii grupate : 1 comutator
 - alte memorii : $2 \times (l - 1 - NID) + 1$ comutatoare

În ANEXA 1 este prezentat un program scris în limbajul de programare C ce determină numărul de căi pentru fiecare tip de rutare între toate sursele și toate destinațiile. De asemenea acest program calculează lungimile minime ale căilor stabilite între surse și destinații. S-au luat în considerare succesiv RMM 16X16, 256X256 și, respectiv, 1024X1024.

În ANEXA 2 sunt prezentate tabele ce conțin numărul de căi între o sursă dată și toate destinațiile pentru o RMM 16 X 16. luând în considerare toate cele patru tipuri de rutări. Sunt prezentate de asemenea tabele ce afișează lungimile minime ale căilor pentru rutările RD, RI, BD și BI.

În ANEXA 3 este prezentat programul de determinare a căii optime între o sursă și o destinație dată în RMM.

În TABEL XI sunt prezentate aceste lungimi pentru o sursă și o destinație dată. S-a considerat o RMM 1024x1024[BINK97]. [CWH01].

TABEL XI

Lungimile de cale pentru fiecare din rutări, considerând $s = 0$ și diferite destinații

Destinația i	Rutări RD/RI	Rutare BD	Rutare BI
514	10	19	19
256	10	5	17
2	10	19	3
72	10	19	13

6.2.2 Analiza de performanță a RMM

Se analizează comparativ comportamentul din punct de vedere al performanțelor a două rețele de interconectare, RMM și respectiv RIN-MN *bidirecțional*, abreviat RIN-MN(B). Ele sunt plasate în cadrul aceluiași mediu cu memorie partajată distribuită, adică SMP din fig. 6.1.

În ambele cazuri, modulul de memorie M_i este direct conectat la procesorul P_i și este denumit *memorie locală* a lui P_i . Cererile de la un procesor către memoria sa locală sunt denumite *cereri interne* și sunt extrase prin magistrala internă dintre procesor și memoria sa locală. O memorie poate deasemenea recepționa *cereri externe* ce sunt generate de alte procesoare și traversează rețeaua RMM sau RIN-MN (B).

Un *grup de procesoare (cluster processors)* asociate unui procesor P sunt acele procesoare situate la distanță de un comutator în raport cu P. Într-un SMP-MPD, există $k - 1$ procesoare ce pot fi accesate prin comutatorul de pe primul nivel sau de pe ultimul la care P este conectat.

Prin urmare o *cerere externă* destinată unui procesor dintr-un grup sau unei memorii dintr-un grup se întoarce din primul nivel sau din ultimul, fără a parcurge întreaga RMM. Cu alte cuvinte *cererile externe* lansate de procesoare parcurg RMM de la stânga spre dreapta în timp ce răspunsurile de la memorii traversează rețeaua de la dreapta spre stânga.

Ambele *cereri interne* și *externe* sosesc la coada de așteptare formată la memorie. Disciplina de deservire a acestor cereri este FIFO.

6.2.2.1 Funcționarea rețelei

În analiza funcționării rețelei se fac următoarele ipoteze de lucru:

- Sistemul este *sincron și comutat prin pachete*
- Tamponul pentru pachete are *capacitatea infinită*
- *Timpul de serviciu al magistralei*, în cazul RMM, sau *timpul de serviciu al legăturii*, în cazul RIN-MN(B), pentru a transfera un mesaj, formează un singur ciclu sistem
- Un procesor este reprezentat printr-un *centru de întârziere*, într-un ciclu dat

- *Timpul de serviciu al modulelor de memorie reprezintă multiplii întregi ai timpului ciclului de sistem.*

Odată ce transmite cererea de acces la memorie, procesorul așteaptă până când obține pachetul de răspuns de memorie (în caz de citire) sau de confirmare (în caz de scriere)

Se definesc următorii parametrii de sistem:

$k \times k$	talia comutatoarelor RMM sau RIN-MN
N	număr de procesoare sau memorii în SMP-MPD
L	$\log_k N$ nivele în RIN
t_s	timp de serviciu al comutatorului
t_m	timp de serviciu al memoriei
p	probabilitatea ca un procesor să plaseze o cerere de acces la memorie
m	probabilitatea ca un procesor să solicite memoria sa locală când s-a făcut o cerere de acces
m_1	probabilitatea ca un pachet să fie destinat unui procesor din grup.
p_i	probabilitatea ca o cerere să traverseze prin nivelul i
r_i	timp mediu de răspuns a unui comutator de pe nivelul i
q_i	număr mediu de cereri locale de la un procesor, în cadrul unui ciclu
q_e	număr mediu de cereri la distanță de la un procesor, în cadrul unui ciclu
d_n	întârziere totală în rețea considerând toate nivelurile
l_m	lungime medie a firului de așteptare în modulul de memorie
d_m	întârziere medie într-un modul de memorie
P_u	utilizare a procesorului (fracțiune din timpul în care procesorul este ocupat)

Analiza performanței atât a RMM cât și a RIN-MN(B) se bazează pe următoarele supoziții :

- Pachetele sunt generate în fiecare nod sursă de către procese aleatoare independente și identic distribuite.
- În orice moment, un procesor este fie ocupat executând calcule interne, fie este în așteptare pentru un răspuns în urma unei cereri de acces la memorie.
- Fiecare procesor ocupat, dacă nu există cereri ce așteaptă să fie rezolvate, generează un pachet cu probabilitatea p la începutul fiecărui ciclu.

- Probabilitatea ca această cerere să fie adresată unei memorii locale (*cerere internă*) este m și probabilitatea ca ea să fie adresată altui modul (*cerere externă*) este $(1-m)$
- O *cerere externă* se va adresa unei memorii dintr-un grup cu probabilitatea m_1 și unui alt modul de memorie, exterior grupului, cu probabilitatea $1-m_1$
- Un răspuns de la memorie parcurge rețeaua în direcția opusă, pe același traseu, în cazul rețelelor RMM, respectiv RIN-MN(B).

Mesajele de la procesor spre memorie sunt generate utilizând probabilități specificate astfel :

- *Probabilitatea de cerere a accesului (p)* : Probabilitatea de cerere a accesului este utilizată ca și un mijloc de estimare a comportamentului procesorului în termeni de cereri de acces la memorie. Când un procesor este ocupat în procesul de calcul, el poate genera o cerere. La fiecare ciclu procesorul decide dacă un mesaj poate fi sau nu transmis cu această probabilitate. În medie, sunt necesari $1/p$ cicluri pentru a extrage o cerere din procesor.
- *Probabilitatea de cerere a accesului la memoria locală (m)* : Fiind cunoscut că o cerere trebuie lansată către memorie, se utilizează o probabilitate m pentru a se decide dacă această cerere este pentru memoria locală sau pentru cea externă.

Deși simple, aceste probabilități au un rol important fiind singurele intrări în analiza de performanță. După fiecare cerere de acces la memorie, procesorul așteaptă pentru o confirmare. Odata ce această confirmare a fost recepționată, procesorul efectuează calcule pe durata unui ciclu și apoi, bazat pe probabilitățile de mai sus, decide dacă va genera sau nu altă cerere de acces la memorie.

Utilizarea procesorului : *Utilizarea procesorului, P_u* , definită ca și fracțiunea de timp în care un procesor este ocupat, va fi determinată de *timpul de așteptare* și de *timpul de serviciu* necesare satisfacerii unei cereri de acces în diverse puncte.

În cazul în care rutările de tip BD și BI vor fi permise doar în primul și ultimul etaj și toate celelalte cereri vor fi rutate direct, utilizarea procesorului este :

$$P_u = 1 / \{ 1 + p(1-m)(1-m_1)(d_n + d_m) + pm d_m + p(1-m)m_1(2r_0 + d_m) \} \quad (6.1)$$

În cazul studiat, un mesaj în RMM sau în RIN-MM(B) va fi expedit de-a lungul unui traseu minim. În acest caz

$$P_u = \frac{1}{1 + \alpha + \beta + \gamma + \eta} \quad (6.2)$$

unde :

- α : întârzierea așteptată pentru ca o cerere locală să fie deservită
- β : întârzierea așteptată pentru ca o cerere la grupul de memorii să fie deservită
- γ : întârzierea așteptată pentru a deservi toate cererile, cu excepția memoriilor grupate ce urmează rutarea RD sau RI
- η : întârzierea așteptată pentru a deservi toate cererile ce urmează o rutare RD sau RI

Termenii $\alpha, \beta, \gamma, \eta$ depind de :

- Probabilitățile de rutare de-a lungul fiecărei căi
- Intensitatea traficului de rețea
- Solicitarea de prelucrare a cererii în fiecare centru de deservire

Puterea de procesare, definită ca și numărul mediu de procesoare ce rămân ocupate, este NP_u

În scopul de a compara eficacitatea comutatorului RMM sau RIN-MN(B) se definește un parametru de cost. O soluție simplă de a defini *costul* unui comutator este reprezentată de *numărul de conexiuni* în comutator. Numărul de conexiuni este k^2 pentru un *crossbar* în timp ce pentru un RMM este $2k$. Într-un SMP-MPD costul total al RIN-MN(B) este $kN \log_k N$ iar pentru RMM este $2N \log_k N$.

Eficacitatea costului rețelei este definită ca fiind raportul dintre *puterea de procesare* și *costul* rețelei. Eficacitatea costului va indica performanța ce poate fi atinsă pe unitatea de cost a rețelei.

6.2.2.2 Probabilități de rutare și întârzierile în parcurgerea căii.

Probabilitățile de rutare și întârzierile în parcurgerea unei căi, atât pentru RMM cât și pentru RIN-MN(B), se obțin presupunând că toate memoriile non-locale sunt adresate în mod egal de către un procesor [BINK97]. Aceste ecuații pot fi modificate în cazul referirilor neuniforme la memorii situate la distanță [IB00]. [BIK00]. Deoarece lungimea căii în rutarea R1 este aceeași ca și rutarea RD, se extrage η bazat pe RD și se multiplică cu 2 pentru a include și rutarea RI. Similar se procedează și pentru rutările BD și BI.

În continuare definim :

Cereri de acces la memoria locală (α) : O cerere de acces la memoria locală nu implică traversarea comutatorului. Prin urmare singura întârziere apare în deservirea cererii în modulul de memorie și ea este d_m .

Prin urmare fiind date p și m , anterior definite. se deduce :

$$\alpha = p \times m \times d_m \quad (6.3)$$

Rutarea de grup (β) : Cererile la procesoarele asociate într-un grup traversează rețeaua spre un comutator din primul sau din ultimul nivel și sunt rutate BD sau BI spre procesorul de destinație. Toate perechile sursă-destinație de acest tip, unde toți biții cu excepția celor mai puțin semnificativi $\log_2 k$ ai etichetei EC sunt 0, atrag acest tip de rutare. Astfel numărul memoriilor asociate în grup pentru o sursă dată este $k-1$ deoarece $k \times k$ este dimensiunea unui comutator din RMM sau din RIN-MN(B). Comutatorul dispus pe nivelul 0 este traversat odată pentru accesarea memoriei din grup și odată pentru a retransmite confirmarea. Când o memorie externă este solicitată, probabilitatea ca memoriile unui grup să fie solicitate este :

$$\beta_p = \frac{k-1}{N-1} \quad (6.4)$$

Considerând $2 \times r_0$ este întârzierea / comutator și că d_m este întârzierea în deservirea memoriei. se obține :

$$\beta = p \times (1-m) \times \frac{k-1}{N-1} \times (2r_0 + d_m) \quad (6.5)$$

Rutarea BD sau BI pentru memoriile negrupate (γ) : În rutarea BD și BI, cererea parcurge rețeaua într-o direcție până la un nivel anume și apoi face o buclă de întoarcere pentru a ajunge la procesorul de destinație. Cunoscând nivelul NID unde se produce buclarea, lungimea căii este $2 \times NID + 1$.

NID este traversat doar odată în timp ce toate nivelurile din stânga lui NID sunt traversate de două ori, dar nu în mod necesar prin același comutator. Deoarece întotdeauna există memorii grupate acoperite ($NID=0, NII=l+1$) se va începe procedura de optimizare cu $NID \geq 1$ și $NII \leq l+2$.

Se consideră $NID < \lfloor l/2 \rfloor$. O dezvoltare similară se poate face pentru $NII \geq \lceil l/2 \rceil$

Se cunoaște că numărul total de destinații este $N-1$. Pentru un nivel dat de întoarcere $1 \leq i < d$, deoarece NID este definit ca și cel mai din dreapta bit din etichetă, toți biții din stânga acestei poziții pot fi 1 sau 0. Această constatare conduce la concluzia că pot exista k' căi de parcurgere în rețea. Numărul de căi în care un bit din ECR poate fi 1 este $k-1$. Atunci când o memorie externă este solicitată, ecuația pentru probabilitatea unei rutări BD sau BI pentru memorii grupate este :

$$\gamma_p = 2 \times \left(\sum_{i=1}^{d-1} \frac{k-1}{N-1} \times k^i \right) \quad (6.6)$$

unde $d = \lfloor l/2 \rfloor$.

Întârzierea într-o astfel de rutare este dependentă de nivelul în care se produce buclarea. Întârzierea pentru toate comutatoarele de pe un traseu dat, cu excepția celei din nivelul de buclare este :

$$\Delta_1 = 2 \times \sum_{j=0}^{i-1} 2r_j$$

Deoarece trebuie luat în considerare faptul că nivelul de întoarcere este traversat de asemenea și de cererea de acces și de confirmare cu întârzierea $r_i + d_m$ se obține :

$$\Delta_2 = 2 \times \left(\sum_{j=0}^{i-1} 2r_j + r_i \right) + d_m$$

Ecuția de determinare a lui γ este :

$$\gamma = p \times (1-m) \times 2 \times \left\{ \sum_{i=1}^{d-1} \frac{k-1}{N-1} \times k \times \left[2 \times \left(\sum_{j=0}^{i-1} 2r_j + r_i \right) + d_m \right] \right\} \quad (6.7)$$

Rutarea RD (η) : Această tehnică se utilizează în cazul tuturor acelor perechi sursă - destinație care nu cad sub incidența categoriilor descrise anterior. Deoarece atât rutarea RD cât și cea RI reprezintă o ultimă opțiune pentru orice alt tip de perechi sursă - destinație, se poate exprima :

$$\eta = (1 - \beta - \gamma) + d_n$$

În acest tip de rutare toate comutatoarele sunt traversate, ceea ce înseamnă :

$$d_n = \sum_{i=0}^{l-1} r_i$$

Prin urmare, întârzierea așteptată pentru toate aceste rutări este :

$$\eta = p \times (1 - m) \times \eta_p \times \left[\left(\sum_{i=0}^{l-1} r_i \right) + d_m \right] \quad (6.8)$$

unde:

$$\eta_p = 1 - \beta_p - \gamma_p \quad (6.9)$$

iar β_p , γ_p sunt date de (6.4) și respectiv de (6.6)

Ecuțiile de la (6.1) la (6.9) sunt valide când memoria locală este accesată cu o probabilitate m și toate celelalte memorii sunt adresate cu probabilități egale, adică $(1-m)/(N-1)$.

TABEL XII prezintă numărul de destinații ce pot fi atinse pornind de la P_0 cu fiecare dintre rutări ca și o funcție de talia rețelei. Talia unui comutator este 2×2 . Analizând rezultatele din tabel, se observă că :

- un număr semnificativ de conexiuni beneficiază de alte tipuri de rutări decât cele de tip RD sau RI, utilizate în mod curent

- aceleași număr de procesoare utilizează rutări BD și BI în două rețele succesive ca și dimensiuni.

De exemplu se consideră dimensiunile 64, respectiv 128, ale rețelei. Aceasta implică $l=6$ și $l=7$. Rețelele, deși au dimensiuni diferite, au aceleași număr de destinație pentru rutare BD deoarece adăugarea unui nivel introduce un nivel central autentic în cazul $l = 7$. Acesta este nivelul 3. În situația $l = 6$ nu există un centru autentic al rețelei. Adăugarea unui nivel central nu va crește numărul posibil de rutări BD sau BI.

TABEL XII
Numărul de destinații pentru diferite dimensiuni ale rețelei utilizând diferite rutări.

<i>Rutare</i>	<i>Dimensiune RMM</i>				
	64	128	256	512	1024
Rutare BD sau BI $(N-1) \times (\beta_p + \gamma_p)$	13	13	29	29	61
Rutări RD sau RI $(N-1) + (\eta_p)$	50	114	226	482	962

Întârzierile r_0 , r_l , d_n și d_m vor depinde de :

- valoarea traficului în rețea, care este o funcție de P_u
- cererea de deservire la centrele individuale de deservire

6.2.2.3 Întârzierile în formarea cozilor de așteptare din comutatoare.

Pentru a simplifica analiza, fiecare nivel din rețea este considerat izolat de celelalte nivele. Se consideră un centru de formare a cozilor de așteptare cu n intrări.

Fie probabilitatea q că există cel puțin un pachet la una dintre intrări în fiecare ciclu dat și că cererea de deservire a unui pachet este t cicluri. Numărul de cereri ce vin în coadă în timpul de deservire a oricărei cereri anterioare va forma o distribuție binomială având următoarele caracteristici :

nt - numărul de încercări

q - probabilitatea de succes

Numărul mediu al cererilor ce sosesc este $E = ntq$ iar varianța este $V = ntq(1 - q)$

În literatură [BINK97], se prezintă metoda de calcul a *lungimii medii a cozii de așteptare* Q în cadrul unei stații de deservire.

$$Q = \frac{E}{2} + \frac{V}{2(1-E)} \quad (6.10)$$

Capacitatea de trecere a acestor cereri este E / t . Timpul de răspuns mediu al centrului este r , calculat prin aplicarea legii lui Little, devine:

$$r = \frac{Q.t}{E} = \frac{1}{2} \left(1 + \frac{1-q}{1-ntq} \right) t \quad (6.11)$$

Modelele de formare a cozilor de așteptare a unui comutator RMM și a unui comutator RIN - MN (B) sunt prezentate în fig. 6.13

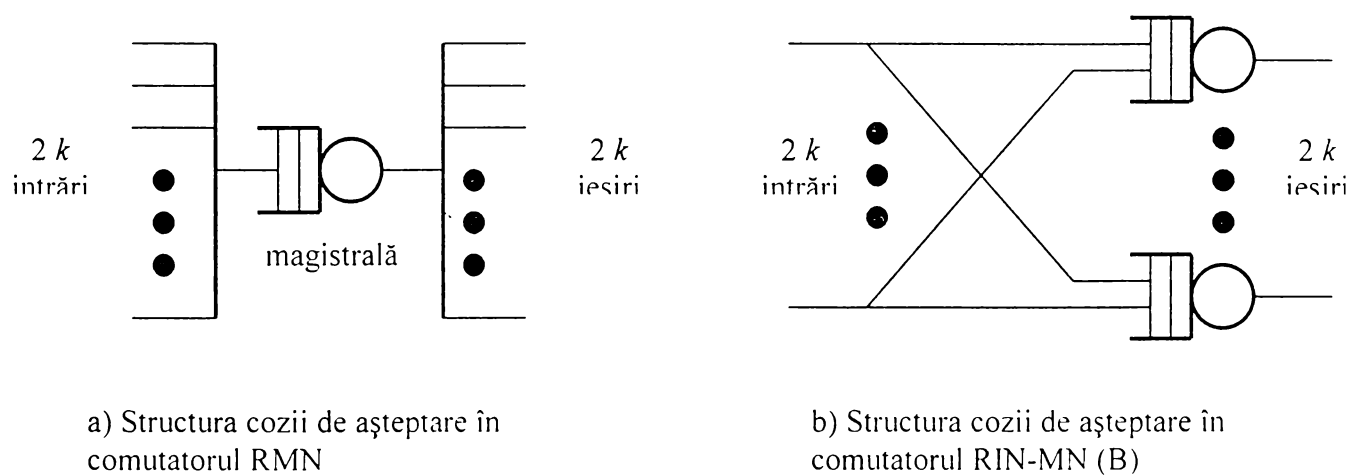


Fig. 6.13 Cozi de așteptare în comutatoarele RMN și RIN-MN (B)

În cadrul unui comutator RMM există o competiție între pachete pentru câștigarea accesului la magistrală. Aceste pachete provin, pe de o parte, de la cele k porturi din dreapta comutatorului, iar pe de altă parte, de la k porturi din stânga. Pentru comutatoare dispuse pe niveluri, altele decât primul și ultimul, $n=2k$, $t=t_s$ și $q=q_n = P_u p (1-m)(1-m_1)$. Se poate calcula *timpul mediu de răspuns* al comutatorului, r_{RMMs} , pentru oricare dintre aceste comutatoare RMM.

Acesta este :

$$r_{RMMs} = \frac{1}{2} \left(1 + \frac{1 - q_s}{1 - 2kt_s q_s} \right) t_s \quad (6.12)$$

Pentru comutatoare dispuse pe primul și pe ultimul nivel, $n=2k$, $t=t_s$ și

$$q = q_s = P_u p (1 - m) (2 - m_1) / 2$$

Timpul mediu de răspuns, r_{RMMi} , este:

$$r_{RMMi} = \frac{1}{2} \left(1 + \frac{1 - q_s}{1 - 2kt_s q_s} \right) t_s \quad (6.13)$$

Întârzierea rețelei RMM, d_n , va fi suma timpilor de răspuns a nivelurilor traversate de un pachet în timp ce este rutat prin rețea.

$$d_n = 2((\log_k N - 2)r_{RMMs} + 2r_{RMMi}) \quad (6.14)$$

În cazul rețelei RIN-MN(B), într-un comutator există $2k$ intrări și $2k$ ieșiri. Probabilitatea de apariție a unei cereri la o intrare sau la o ieșire a unui comutator RIN-MN(B) dispus pe nivelul i va fi $P_u p (1 - m) p_i$.

Timpul de răspuns a unei RIN-MN(B) se calculează utilizând $n=2k$, $t=t_s$ și

$$q = P_u p (1 - m) p_i / n.$$

Acesta este:

$$r_{RMBi} = \frac{1}{2} \left(1 + \frac{1 - \frac{q}{n}}{1 - 2kt_s \frac{q}{n}} \right) t_s \quad (6.15)$$

Întârzierea totală din rețea este d_n și reprezintă suma timpilor de răspuns ale comutatoarelor din diferite niveluri.

În ambele tipuri de rețele, numărul mediu de cereri sosind la un modul de memorie este $E_m = q_i + t_m q_e$, unde q_i și q_e sunt cereri interne și externe pentru acel modul de memorie.

Varianța $V_m = q_i(1 - q_i) + t_m q_e(1 - q_e)$.

Astfel lungimea medie a cozii de așteptare la memorie este :

$$l_m = \frac{q_i + t_m q_e}{2} + \frac{q_i(1 - q_i) + t_m q_e(1 - q_e)}{2(1 - q_i + t_m q_e)} \quad (6.16)$$

iar timpul de răspuns mediu al memoriei este :

$$d_m = \frac{l_m t_m}{E_m} \quad (6.17)$$

Ecuția din care se determină P_u este neliniară. Ea se rezolvă utilizând tehnici de iterație.

O tehnică de iterație pentru a determina P_u este următoarea :

1. Se inițializează P_u cu o valoare ce reprezintă utilizarea presupusă a procesorului.
Cu cât valoarea este mai bine aleasă cu atât vor fi mai puține iterații pentru calcul.
2. Se calculează probabilitățile de cerere la fiecare nivel al rețelei și la fiecare modul de memorie. Un pas intermediar ar putea fi calculul valorii statice pentru p_i ce reprezintă probabilitatea ca un nivel să fie traversat.
3. Se calculează timpii medii de răspuns ai comutatorului și timpul mediu de răspuns al memoriei, adică r_i respectiv r_m .
4. Se calculează întârzierea rețelei și a memoriei considerând valorile anterior calculate și utilizând ecuațiile pentru determinarea lui α , β , γ și η .
5. Se calculează un nou P_u , considerând valorile anterioare.
6. Se repetă pașii 2 - 5 până când noul P_u este în limita unei toleranțe admise comparabil cu ultimul P_u . Uzual se lucrează cu o valoare inițială de 0.5 pentru P_u și cu o precizie de 0.00001.

6.2.3 Rezultate numerice și discuții.

În această secțiune se prezintă performanța relativă a RMM, RIN-MN(B) și RIN-MN convențională. RIN-MN(B) este o rețea RIN-MN *bidirecțională*, echivalentă cu RMM, dar constituită din comutatoare *crossbar*. Pentru simetria abrevierii vom utiliza acronimul RMB. Această rețea permite cele patru tipuri de rutări propuse pentru RMM. RIN-MN convențională – abreviat RMC – este similară cu rețeaua utilizată în sistemul Butterfly [BINK 97] și se caracterizează prin faptul că pachetele reprezentând cereri și răspunsuri traversează rețeaua într-o singură direcție, de la procesoare către memorii.

Se analizează *utilizarea procesorului și timpul de răspuns* pentru cele trei tipuri de rețele luând în considerare valori diferite ale lui m .

S-au efectuat simulări ciclu cu ciclu pentru a verifica că rutările propuse funcționează și s-au măsurat probabilitățile și întârzierile rețelei [Iyer96].

Simularea a fost făcută utilizând un mediu descentralizat distribuit cu comutare de pachete în mod sincron. Comutatoarele au tamponare de ieșire cu lungime infinită. Toate evenimentele apar în sistem la începutul ciclului mașină. Fiecare ciclu, utilizat în simulare, reprezintă timpul cerut pentru transmisia unui pachet de la un tampon de ieșire al unui comutator sau al unui procesor (memorie) spre tamponul de ieșire al destinației sale. Astfel se include transmisia pachetului prin canal și timpul necesar comutatorului să-l ruteze spre tamponul de destinație corespunzător. Timpul minim necesar pentru ca un pachet să ajungă la memoria de destinație se bazează pe numărul de comutatoare dispuse în calea sa. Întârzierea introdusă de memorie, timpul de deservire a cererii, a fost fixată la patru cicluri.

Rutările descrise în §6.2.1 au fost simulate. Simularea compară fiecare sursă și destinație, rulând algoritmul de rutare optimă și apoi selectând rutarea. Deciziile referitoare la tipul de rutare se bazează exclusiv pe etichetele generate de către algoritmul de rutare optimă.

Mesajele de la procesor la memorie sunt generate utilizând probabilitățile p , m , m_1 prezentate în §6.2.2.

S-au comparat performanțele RMM și RMB obținute pe cale analitică, respectiv prin simulare, în termeni de utilizare a procesorului și respectiv timpul de răspuns.

În fig. 6.14 se prezintă o comparație între rezultatele analizei și cele ale simulării pentru utilizările procesorului în cazul RMM, pentru două valori ale lui m .

În fig. 6.15 este prezentată comparația pentru timpul de răspuns.

În ambele figuri se observă că rezultatele analitice se aproprie foarte mult de cele obținute din simulare. Aceasta indică că independența cozilor de așteptare, presupusă în timpul analizei, nu a fost una grosieră și pentru valori mai mari ale lui m sunt satisfăcute mai multe cereri de acces fără a se parcurge RMM. Astfel, pentru $m=0.9$, P_u este mult mai mare în fig. 6.14 iar timpul de răspuns mult mai mic în fig. 6.15. Cu cât p crește, sunt generate mai multe cereri și crește timpul de răspuns iar P_u se reduce.

În fig. 6.16 se prezintă utilizarea procesorului la un sistem 1024×1024 pentru RMM, RMB și RMC ca și funcția probabilității cererii de acces la memorie, p , cu $m_1 = (1-m)(k-1)/(N-1)$. Se observă că timpul de răspuns este același pentru toate rețelele atunci când $m=0.9$. Pentru $m=0.1$, RMB se comportă mai bine ca și RMM, care la rândul lui se comportă mai bine ca și RMC. Pentru valori mai mari ale probabilității de cerere de acces la memoria locală, m , sunt satisfăcute mai multe cereri fără a se traversa RIN, cu cât p este mai mare, mai multe cereri sunt generate. Concluzia este că apare o diferență semnificativă în performanță între RMM și RMC pentru valori coborâte ale lui m și valori mari ale lui p . Diferența este de aproximativ 3% pentru valori superioare ale lui m și aproximativ 10% pentru valori inferioare ale lui m .

În fig. 6.17 se prezintă timpii de răspuns a celor trei rețele pentru valori diferite ale lui m ca și o funcție de p cu $m_1 = (1-m)(k-1)/(N-1)$. *Timpul de răspuns* este definit ca și diferența dintre momentul când un procesor generează o cerere de acces și momentul când el primește răspunsul din partea memoriei. Timpul de răspuns pentru toate RIN este mai mare când m are valori inferioare. Acest lucru se datorează traficului intens prin RIN. Din nou se observă că diferența maximă este de aproximativ 3% pentru valori ridicate ale lui m și de aproximativ 10% pentru valori inferioare ale lui m .

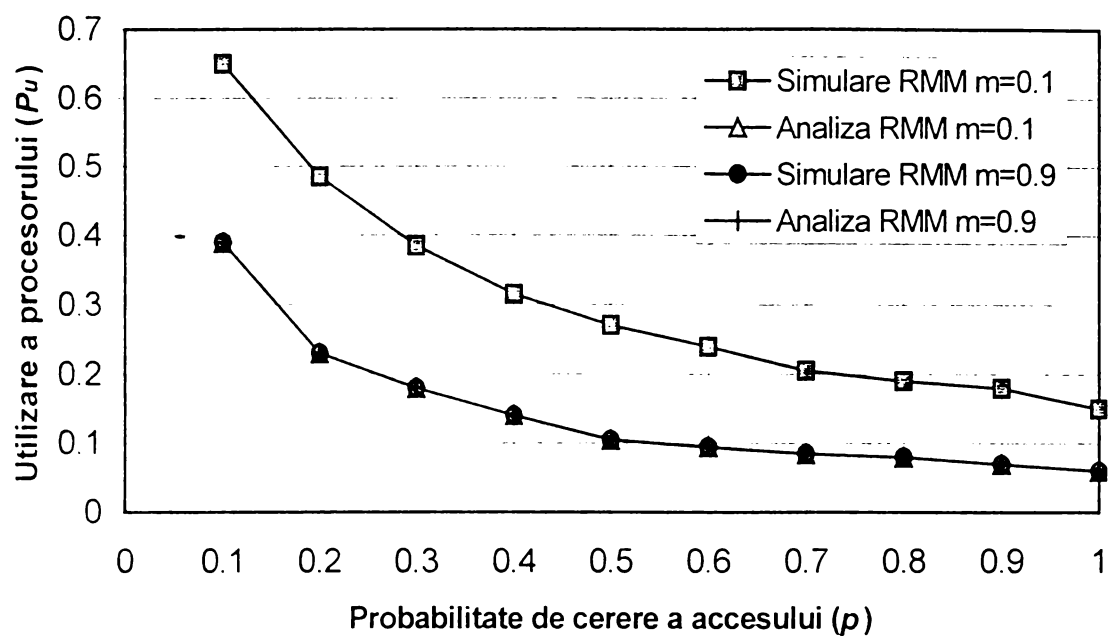


Fig. 6.14 Comparație a rezultatelor analizei și simulării pentru utilizările procesorului la RMM, variind m

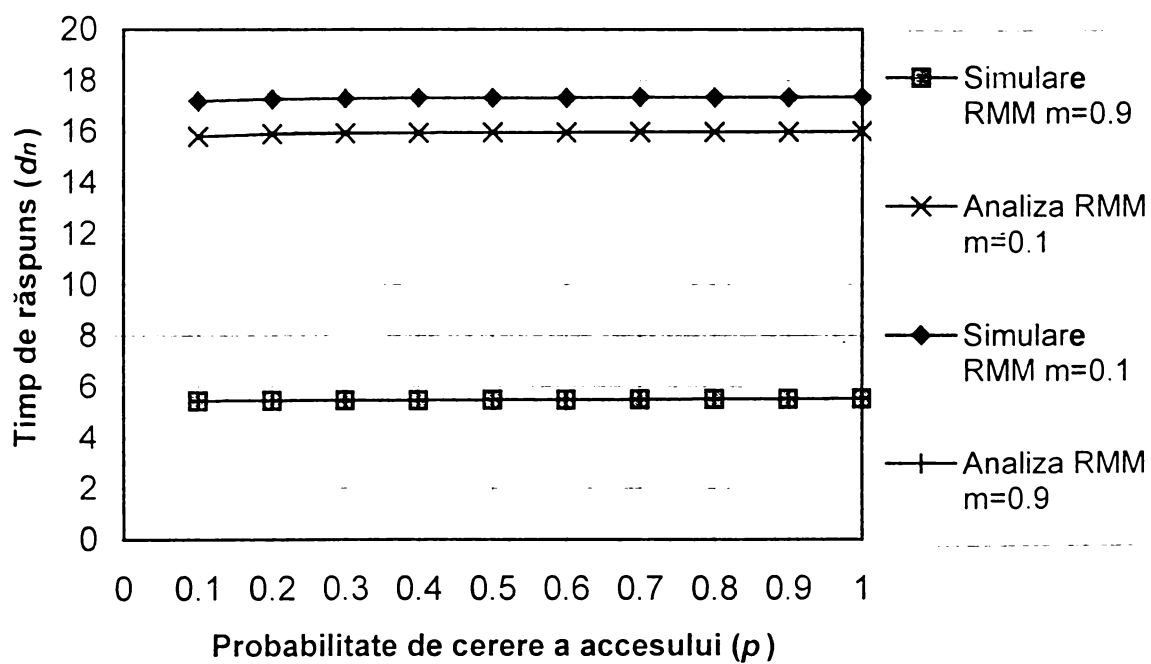


Fig. 6.15 Comparație a rezultatelor analizei și simulării pentru timpul de răspuns al RMM, variind m

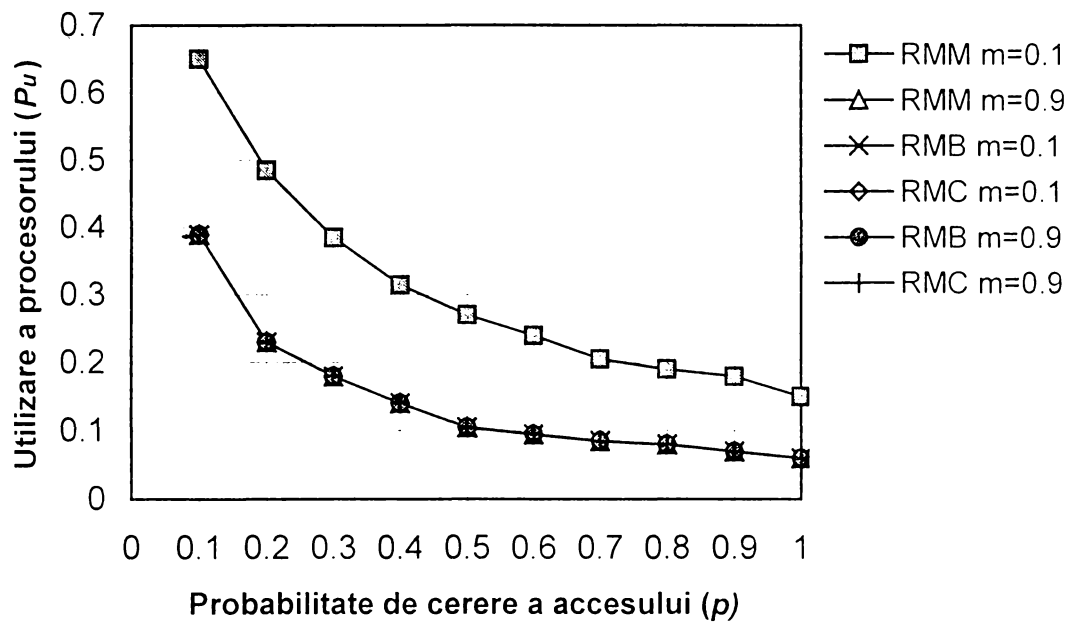


Fig. 6.16 Comparația între utilizările procesorului la diferite RIN, variind m

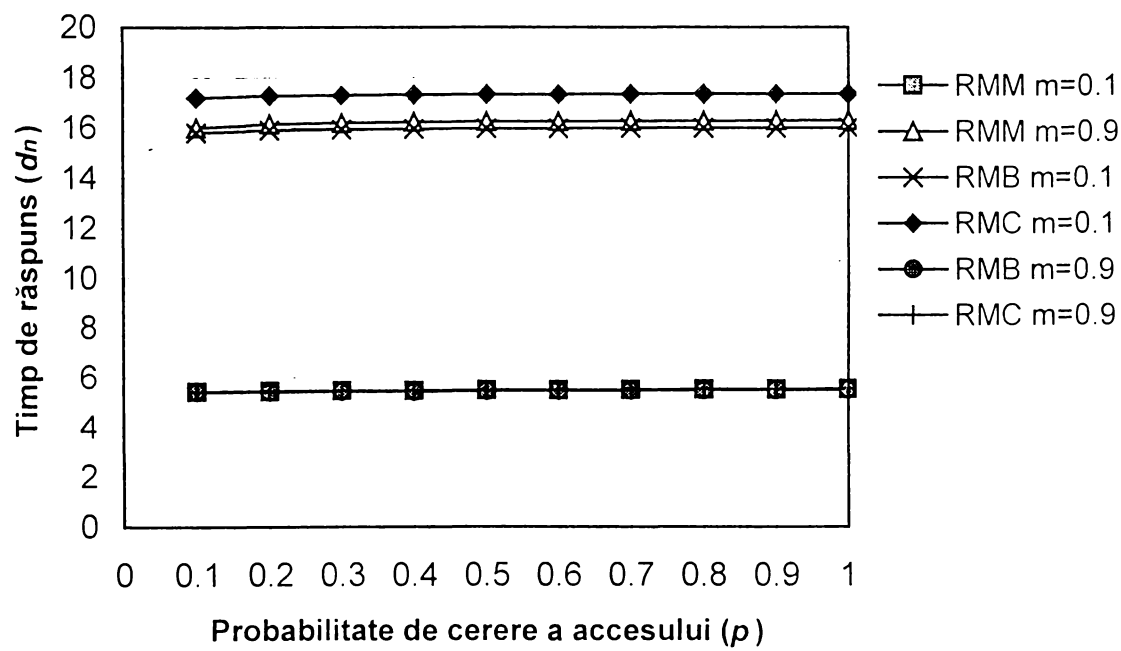


Fig. 6.17 Comparația între timpii de răspuns la diferite RIN, variind m

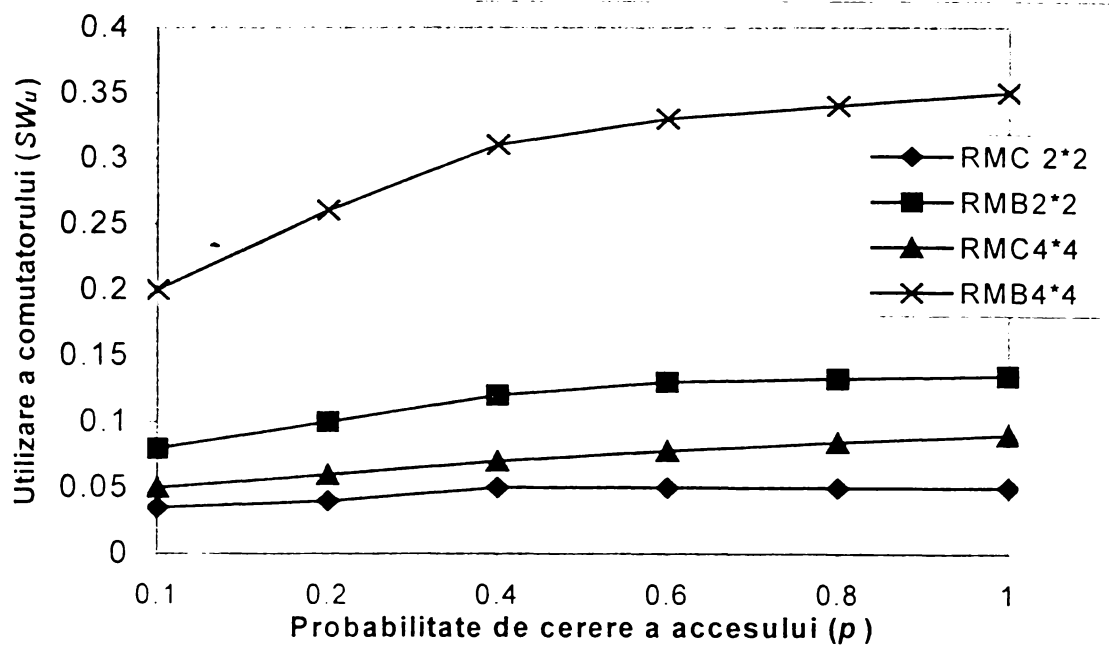


Fig. 6.18 Utilizarea comutatorului în RMM și RMC
variind dimensiunea comutatorului

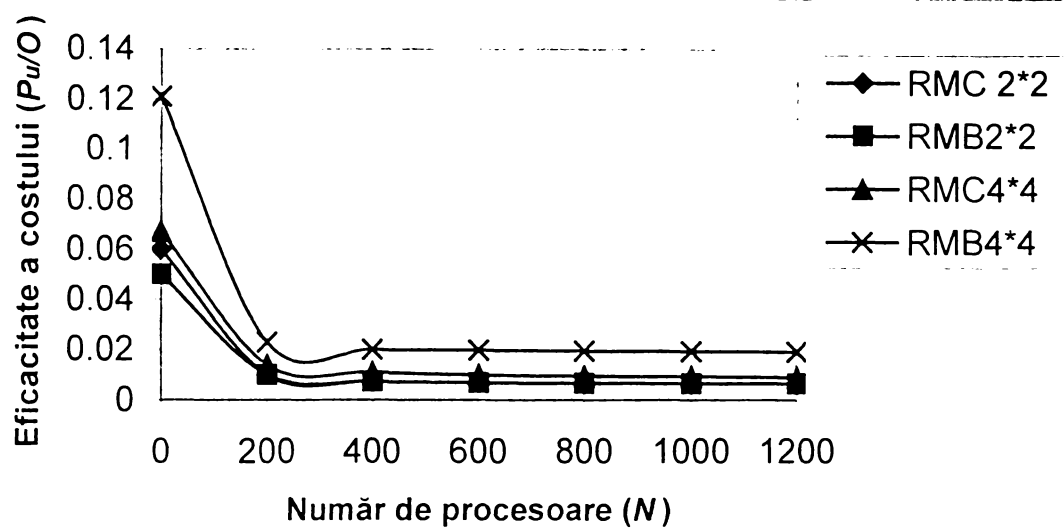


Fig. 6.19 Eficacitatea costului în RMM și RMC
variind dimensiunea comutatorului

În fig.6.18 se compară *utilizarea magistralei* într-un comutator RMM și *utilizarea legăturii* într-un comutator *crossbar* din RMC, pentru 2 dimensiuni diferite ale comutatoarelor într-un SMP 1024×1024 . Utilizarea comutatorului RMC este mai mică de 10%. Utilizarea unui comutator RMM 2×2 este mai puțin de 15% iar a unui comutator RMM 4×4 mai puțin de 35%. Pentru valori inferioare ale lui p utilizarea comutatorului este redusă și eă crește cu p pentru toate cazurile. Se confirmă teoria conform căreia conexiunile din RMC sunt neutilizate în mod masiv în comparație cu utilizarea comutatoarelor RMM.

Eficacitatea costului sistemelor bazate pe RMM și RMC este prezentată în fig. 6.19 ca și o funcție de talia sistemului. Deoarece *costul* unei RMM comparativ cu a unei RMC este identic atunci când se utilizează comutatoare 2×2 , eficacitatea costului este aceeași pentru cele două sisteme. O rețea RMM cu comutatoare 4×4 este mai eficientă din punct de vedere a costului decât o rețea RMC echivalentă. Diferența este de aproximativ 100% pentru $p=0.5$. În fig.6.19 se observă clar superioritatea unei RMM în raport cu RMC.

În fig.6.20 și 6.21 se prezintă *utilizarea procesorului* și *timpii de răspuns* pentru diferite talii ale sistemelor. Rezultatele au fost obținute cu o probabilitate de cerere de acces a memoriei locale, m , fixată la 0.5 și pentru două valori diferite ale lui p , adică 0.1 respectiv 0.5. Se poate constata că performanța rețelei RMM rămâne apropiată de cea a RMB chiar și în situația când talia sistemului crește. Se observă că în toate situațiile RMM se comportă mai bine decât RMC. Rețeaua RMM este scalabilă la un nivel ridicat în condițiile unui trafic impus.

În TABEL XIII se prezintă dimensiunile optime ale unui comutator pentru o valoare dată a lui p și pentru diferite talii de sistem ($m=1/N$ și $m_1=(k-1)/N$). Pentru $p=0.5$ toate sistemele sunt optime pentru comutatoare 4×4 . În acest caz dimensiuni superioare ale comutatorului vor conduce la un trafic crescut și prin urmare mai puțină performanță și mai puțină eficacitate a costului. Comutatoarele 2×2 vor avea un număr sporit de etaje și prin urmare o întârziere mai mare și vor fi mai puțin eficace din punct de vedere a costului. O excepție este sistemul de talie 512×512 . Oricum, pentru $p=0.1$ situația se schimbă. Talia comutatorului 8×8 este optimă oriunde acesta este utilizat. Pentru dimensiuni ale comutatorului egale sau mai mari de 16×16 magistralele sunt saturate și prin urmare performanța este redusă iar rețeaua are o eficacitate a costului diminuată.

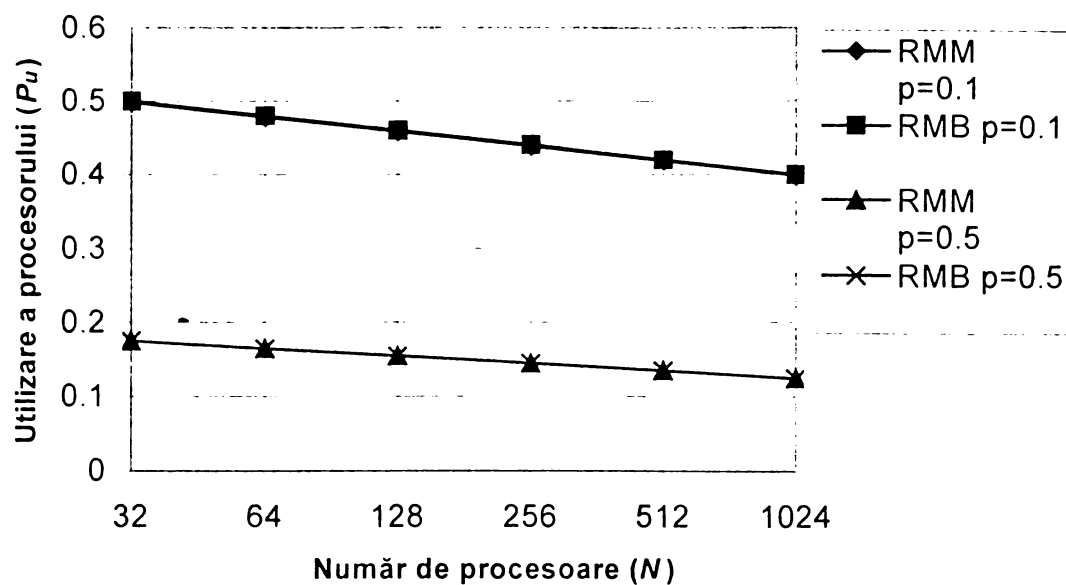


Fig. 6.20 Utilizarea procesorului: Scalabilitatea RMM în raport cu RMB, variind p

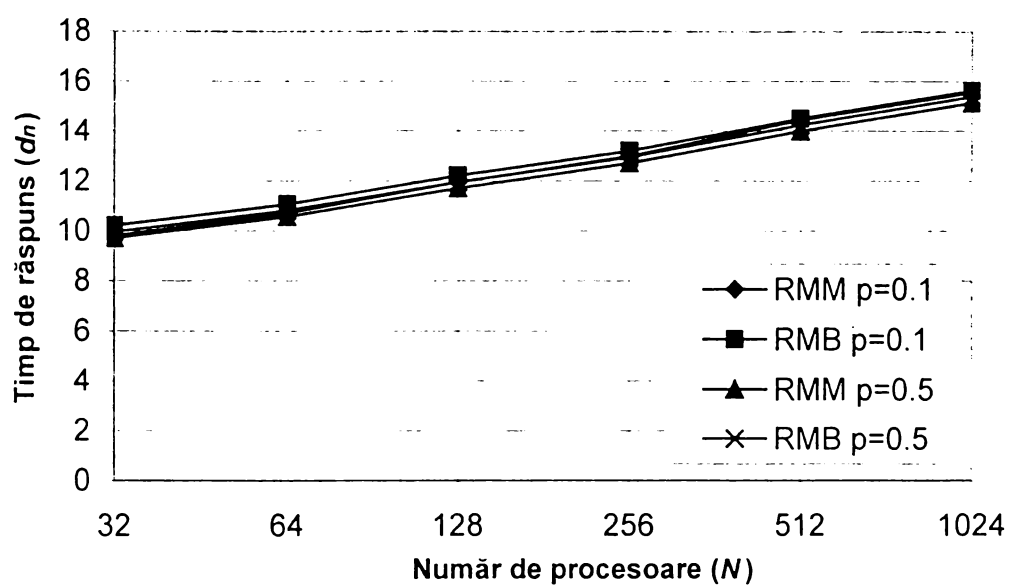


Fig. 6.21 Timp de răspuns : Scalabilitatea RMM în raport cu RMB, variind p

În TABEL XIV se prezintă utilizarea procesorului și timpul de răspuns ale RMM pentru diferite dimensiuni ale comutatorului și număr diferit de procesoare. În tabel

câteva locuri sunt necompletate deoarece o RMM $N \times N$ nu poate fi construită utilizând doar comutatoare $k \times k$. Atât p cât și m sunt fixate la 0.5. Pentru $N=64$, se observă o scădere a lui P_u când k crește de la 4 la 8. Aceasta datorită faptului că RMM este mai puțin eficientă în condițiile apariției unui conflict mărit și a unei întârzieri augmentate în cazul unui comutator RMM 8×8 . Pe de altă parte se constată o ameliorare a performanței într-un sistem 512×512 , unde k crește de la 2 la 8. Numărul de comutatoare din întreaga rețea este încă mare, menținând redusă rata conflictelor ceea ce este suficient pentru câștigul în performanță.

Comparând performanța RMM cu cea a RMB se poate constata că pe măsură ce dimensiunea comutatorului crește, RMB oferă o utilizare mai mare a procesorului și un timp de răspuns mai mic. Acest câștig în performanță se datorează unui număr mai mic de conflicte în comutatorul *crossbar*. Oricum RMB oferă o performanță crescută în detrimentul costului. Luând în considerare numărul de conexiuni într-un comutator $k \times k$ care este k^2 , că pentru un comutator bazat pe magistrală el este $2k$, iar costurile totale ale RMB și RMM sunt $kN \log_k N$ și $2N \log_k N$ iar apoi coroborând aceste date cu *utilizarea procesorului* și cu *timpul de răspuns*, rezultă că *eficacitatea costului* unei RMM este mai mare decât a unei RMB. Se constată de asemenea că un comutator 4×4 este cel mai eficient din punct de vedere al costului pentru diferite dimensiuni de rețea.

TABEL XIII

Dimensiunea optimă a unui comutator pentru eficacitatea costului.

N	Dimensiuni posibile ale comutatorului	Dimensiune optimă	
		$p=0.5$	$p=0.1$
16	$2 \times 2, 4 \times 4, 16 \times 16$	4×4	4×4
64	$2 \times 2, 4 \times 4, 8 \times 8, 64 \times 64$	4×4	8×8
256	$2 \times 2, 4 \times 4, 16 \times 16, 256 \times 256$	4×4	4×4
512	$2 \times 2, 8 \times 8, 512 \times 512$	2×2	8×8
1024	$2 \times 2, 4 \times 4, 32 \times 32, 1024 \times 1024$	4×4	4×4

TABEL XIV

Utilizarea procesorului și timpul de răspuns al RMM și RMB variind k și n

Nr.proc N	Diferite talii ale comutatorului											
	k=2				k=4				k=8			
	P_u		T_r		P_u		T_r		P_u		T_r	
	RMM	RMB	RMM	RMB	RMM	RMB	RMM	RMB	RMM	RMB	RMM	RMB
64	0.150	0.153	10.96	10.73	0.181	0.191	9.06	8.17	0.179	0.210	9.19	7.51
256	0.130	0.132	13.17	12.88	0.166	0.176	9.95	9.30	-	-	-	-
512	0.121	0.123	14.40	14.07	-	-	-	-	0.158	0.190	10.69	8.54
1024	0.114	0.116	15.37	15.03	0.151	0.161	11.20	10.41	-	-	-	-

6.2.4 Toleranța la defect și fiabilitatea RMM

6.2.4.1 Toleranța la defect

În cadrul analizei fiabilității la RMM, modelul de defect de bază este *defectul de comutator*. La apariția unei defecțiuni la nivel de comutator, acesta devine total inutilizabil și, pe cale de consecință, toate legăturile conectate de către acest comutator devin inutilizabile [BNA94], [BNT94].

În cazul RIN-MN convenționale există doar o singură cale între fiecare pereche sursă-destinație. Prin urmare când un comutator se defectează în oricare nivel al rețelei, rețeaua devine inutilizabilă. În cazul RMM există căi multiple între fiecare pereche sursă-destinație ceea ce conferă rețelei capacitatea de a tolera toate defecțiunile de tip *comutator singular* și, în cazuri bine determinate, defecțiunile de tip *comutatoare multiple*.

În capitolul 5, se arată că toleranța la defect în cazul defectării unui singur comutator, implică existența a cel puțin două căi distincte între fiecare pereche sursă-destinație. Aceasta înseamnă că toate comutatoarele dintr-o cale trebuie să fie distincte de comutatoarele din cealaltă.

În cazul RMM există trei cazuri ce trebuie să fie luate în considerare depinzând de etichetele sursă și destinație. $S=s_0s_1\dots s_{l-1}$ și $D=d_0d_1\dots d_{l-1}$

Cazul 1: $s_i=d_i$ pentru toți i , adică $S=D$.

În acest caz cererea de acces parcurge doar magistrala internă și nu traversează nici un comutator. Prin urmare nu sunt relevante defectele comutatoarelor.

Cazul 2: $s_i=d_i$ pentru toți i , cu excepția $i=l-1$, adică $S \neq D$ doar pentru bitul din extrema dreaptă.

În acest caz atât sursa cât și destinația sunt conectate la același comutator în nivelul 0 și de asemenea la același comutator în nivelul $l-1$. Cererea poate utiliza oricare dintre aceste comutatoare pentru a efectua o buclare și pentru a ajunge la destinație. Prin urmare, când unul dintre aceste comutatoare se defectează, RMM poate încă funcționa normal.

Cazul 3: $s_i \neq d_i$, adică S și D diferă cel puțin într-un bit i , dar $i \neq l-1$.

În acest caz cererea traversează rețeaua fie în sens Direct, fie în sens Invers, utilizând rutările RD sau RI.

În cazul rutării RD, o cerere sosește în poziția $s_j s_{j+1} \dots s_{l-2} d_0 d_1 \dots d_{j-1} s_{j-1}$ la stânga nivelului j . Această poziție corespunde comutatorului $SW1 = s_j s_{j+1} \dots s_{l-2} d_0 d_1 \dots d_{j-1}$ în nivelul i .

În cazul rutării RI, o cerere sosește în poziția $d_j d_{j+1} \dots d_{l-2} s_0 s_1 \dots s_{j-1} s_j$ în dreapta nivelului j . Această poziție corespunde comutatorului $SW2 = d_j d_{j+1} \dots d_{l-2} s_0 s_1 \dots s_{j-1}$ în nivelul j .

Deoarece există cel puțin un i , altul decât $i=l-1$, astfel ca $s_i \neq d_i$, comutatoarele SW1 și SW2 diferă unul de celălalt în fiecare nivel j . Prin urmare calea în sens Direct și calea în sens Invers sunt *disjuncte* în termeni de comutatoare.

Prin urmare, RMM este tolerant la defectele de tip *un comutator defect*.

6.2.4.2 Fiabilitatea

Anduranța RMM la defecțiuni multiple de comutator poate fi exprimată în termeni de *fiabilitate de terminal medie* (*Average Terminal Reliability*) – FTM

Fiabilitatea de terminal (FT) a unei perechi sursă – destinație date, este definită ca fiind probabilitatea ca să existe cel puțin o cale liberă de defect între sursă și destinație, cunoscându-se că fiecare comutator are o anumită fiabilitate.

FTM a rețelei este media FT -urilor a tuturor perechilor sursă-destinație. Pentru o RIN-MN cu l niveluri, există doar o cale de lungime l între fiecare pereche sursă-destinație cu excepția situației când sursa și destinația sunt identice. Prin urmare FTM pentru o RIN-MN este:

$$FTM = R^l$$

RMM este bidirecțională și există mai multe căi între o pereche sursă-destinație. Câteva comutatoare sunt partajate de către două sau mai multe căi. *Graficul de redundanță* pentru RMM ce descrie toate căile disponibile între o pereche sursă-destinație, este prezentat în fig.6.22. El constă din două noduri speciale, nodul sursă S și nodul destinație D , și un set de noduri reprezentând comutatoarele ce se găsesc pe diferitele trasee între S și D .

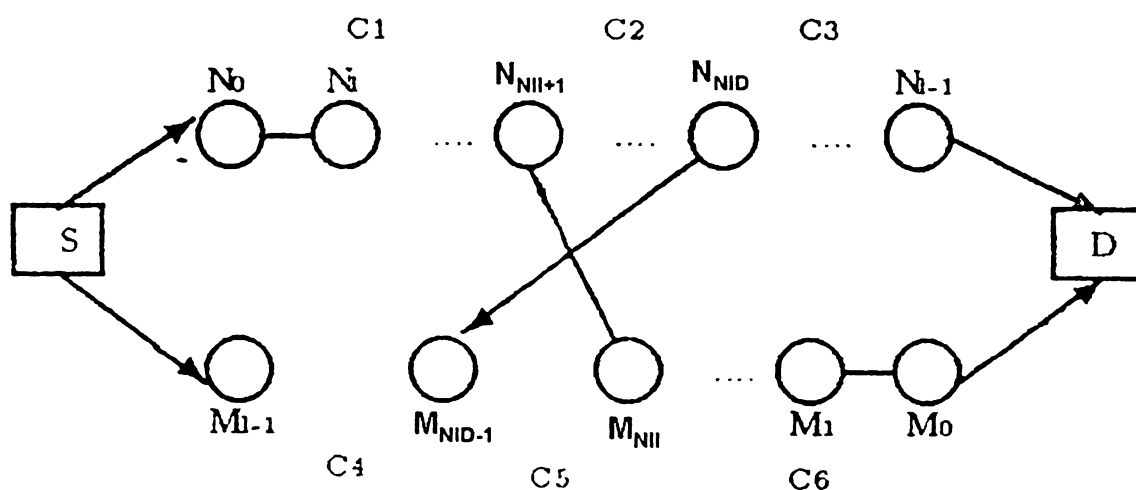


Fig.6.22 Graf de redundanță în RMM

$$C_1 = N_0 \dots N_{NII}$$

$$C_2 = N_{NII+1} \dots N_{NID}$$

$$C_3 = N_{NID-1} \dots N_{i-1}$$

$$C_4 = M_{i-1} \dots M_{NID}$$

$$C_5 = M_{NID-1} \dots M_{NII}$$

$$C_6 = M_{NII-1} \dots M_0$$

Calea parcursă în cazul RD constă din componente $C_1 \cdot C_2 \cdot C_3$ și calea în cazul RI din componente $C_4 \cdot C_5 \cdot C_6$. În mod similar căile parcurse în cazul BD și BI constau din componente $C_1 \cdot C_2 \cdot C_5 \cdot C_6$ și, respectiv, din $C_4 \cdot C_5 \cdot C_2 \cdot C_3$. Fiabilitatea acestor căi va depinde de fiabilitatea componentelor individuale ce sunt funcții de NID și NII.

Fiabilitatea de terminal medie (FTM) a RMM, respectiv a RMC, este prezentată în fig.6.23. S-au utilizat comutatoare 2×2 și diferite valori ale *fiabilității comutatorului* – *FSM*. În fig.6.24 sunt prezentate diagramele pentru comutatoare 4×4 .

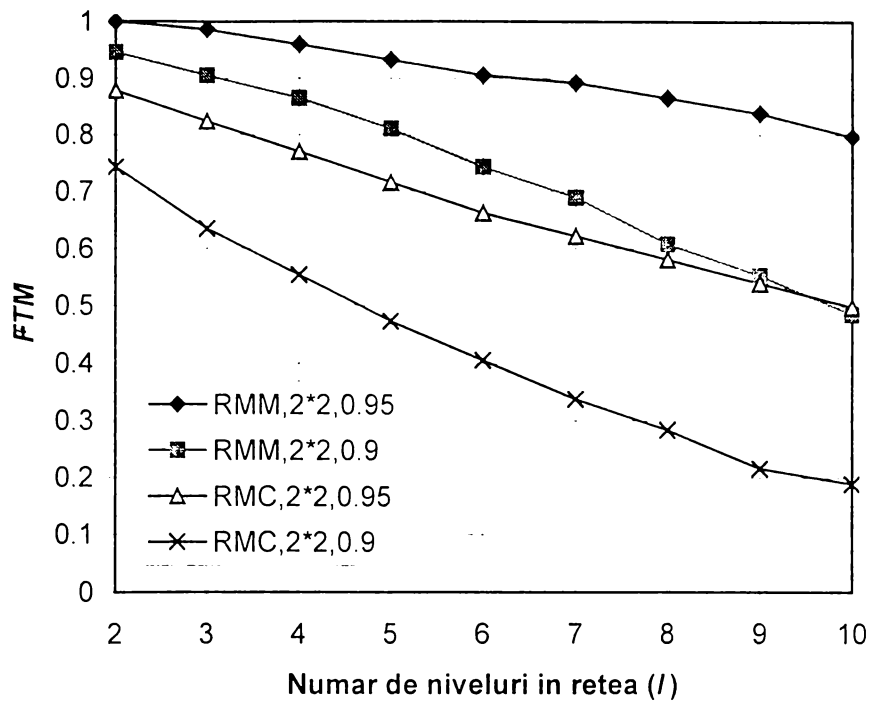


Fig. 6.23 Fiabilitatea rețelelor în cazul comutatoarelor 2×2

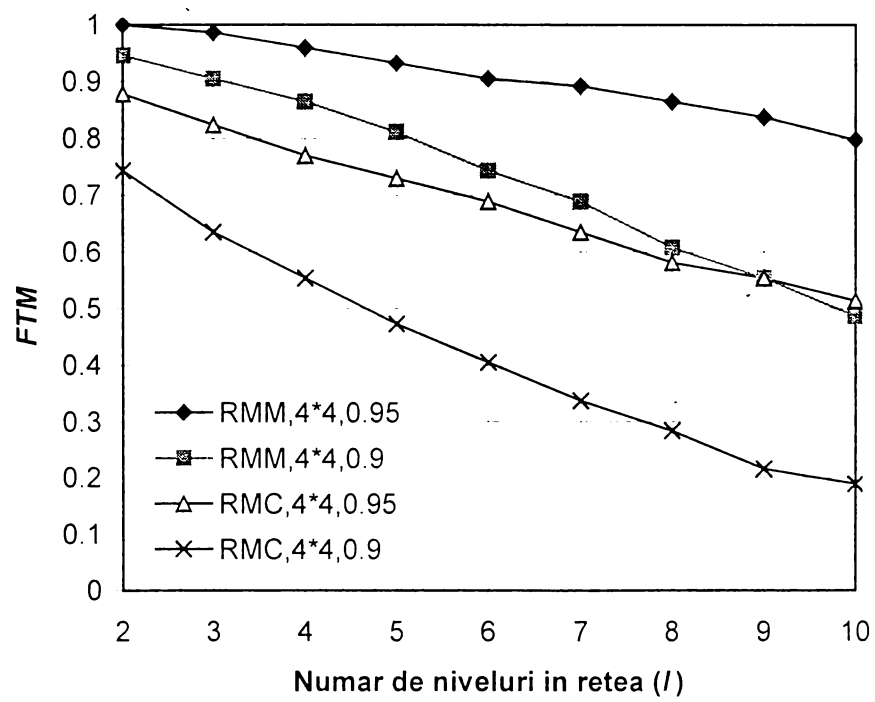


Fig. 6.24 Fiabilitatea rețelelor în cazul comutatoarelor 4 × 4

Fiabilitatea ambelor rețele scade ca și o funcție de l , deoarece numărul de comutatoare în fiecare cale crește cu l . Se observă că RMM are o mai bună fiabilitate decât RMC în întreg spectrul de valori ale parametrilor. Chiar și în situația în care $FSM=0,9$ pentru RMM, această rețea este mai fiabilă decât RMC cu $FSM=0,95$. Fiabilitatea mai ridicată a RMM se datorează în mod evident căilor multiple dintre perechile de noduri sursă-destinație, în contrast cu căile unice într-o RMC.

6.3 PROTOCOALE DE MENȚINERE A COERENȚEI

Un modul de memorie este conectat la toate procesoarele într-o structură de *arbore* ca și în fig. 6.25. Modulul formează *apex*-ul, comutatoarele sunt *noduri* intermediare și procesoarele sunt *frunze*.

Graful arborelui este identic cu talia I/E a unui comutator. Un bloc este localizat într-un modul de memorie, fie intern, fie extern. Toate cererile de acces sau răspunsurile pentru acel bloc vor parcurge arborele având ca apex modulul de memorie. Prin urmare protocoalele de menținere a coerenței pot fi descrise în termeni de ierarhie de nivele ale unui arbore, ceea ce corespunde etajelor din RMM.

Procesoarele și memoriile cache asociate sunt dispuse pe *nivelul 0*. Primul etaj de comutatoare din RMM este dispus pe *nivelul 1*, al doilea pe *nivelul 2*, etc.

Porturile superioare ale unui comutator de pe nivelul i sunt conectate cu porturile inferioare ale comutatoarelor de pe nivelul $i+1$.

Un comutator este un *predecesor* al unui procesor dacă o cerere de acces la memorie lansată de procesor traversează acel comutator.

Un procesor sau o cache este *în aval* față de un comutator, dacă comutatorul este un predecesor al aceluși procesor. Dacă un procesor nu este *în aval față de* un comutator, atunci el este *în amonte* față de acel comutator. Comutatoarele de pe *nivelul i* care sunt conectate la un comutator de pe *nivelul i+1* se numesc *descendenți direcți* ale comutatorului dispus pe *nivelul i+1*.

Descrierea mișcării unui bloc partajat se va face în continuare în termenii mai sus menționați.

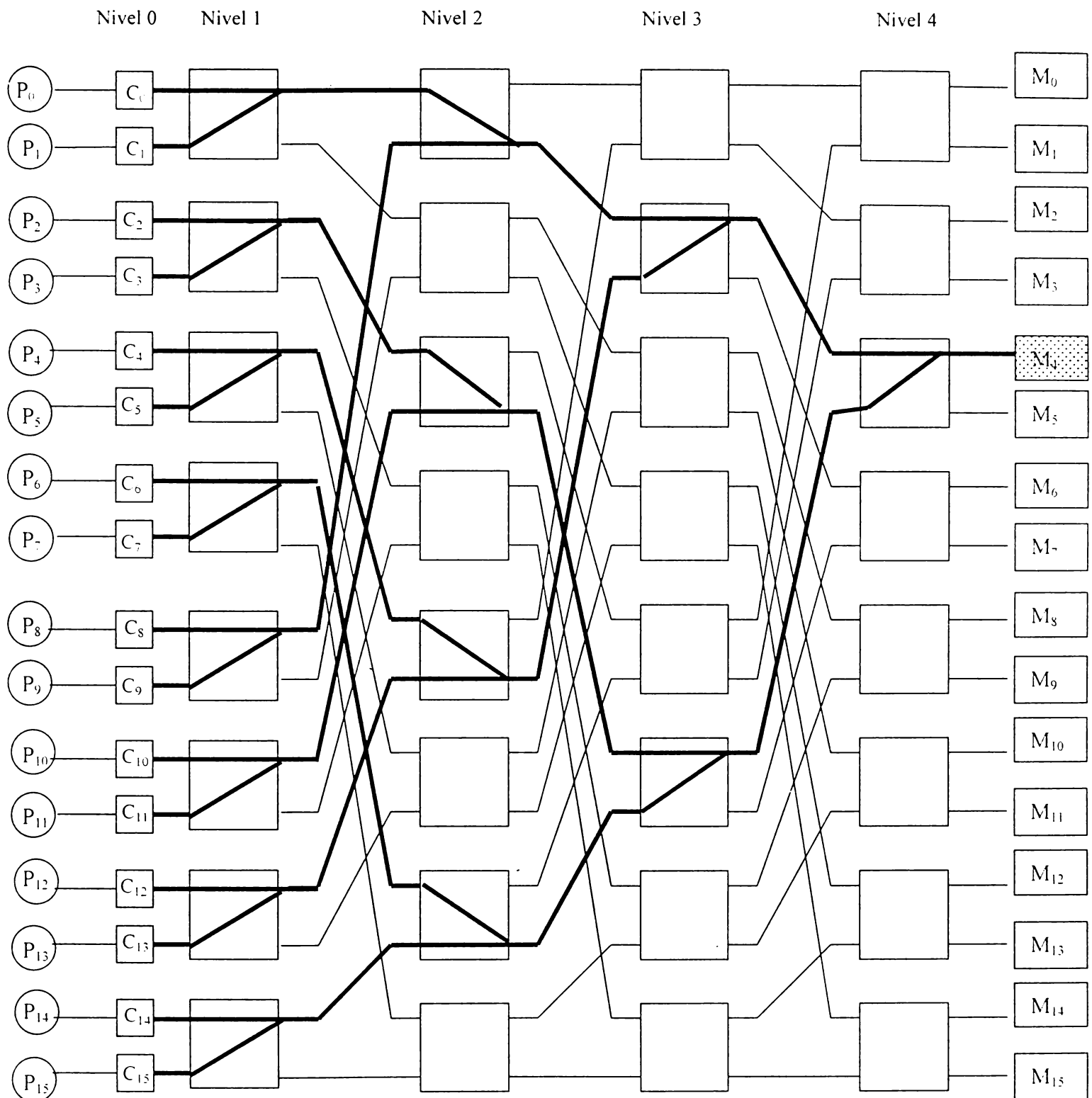


Fig. 6.25 Un arbore ce conectează M_4 cu toate procesoarele

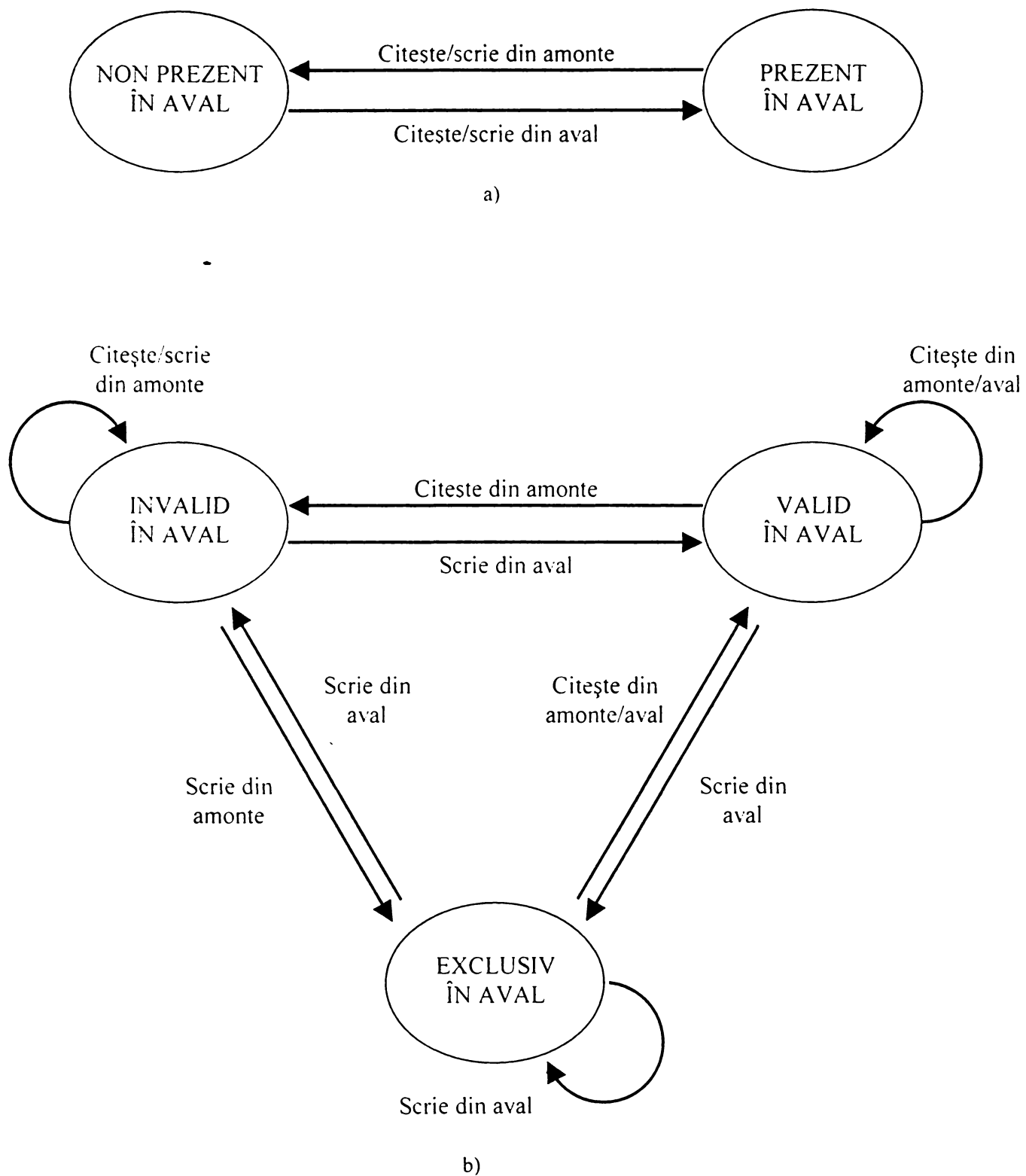


Fig.6.26 Diagrame ale tranzițiilor de stare în protocol

a) Protocol pentru o singură copie

b) Protocol pentru copii multiple

A. Protocol pentru o singură copie.

Acest protocol permite prezența în sistem a unei singure copii a unui bloc partajat la un moment dat. Un bloc partajat poate fi în starea *prezent în aval* sau în starea *non-prezent în aval* în raport cu un comutator la un nivel dat. Această situație este prezentată în fig. 6.26(a). Starea *prezent în aval* înseamnă că un procesor este localizat sub comutatorul ce are o copie a blocului. Stare *non - prezent în aval* înseamnă că blocul este fie prezent în alt procesor, fie în memorie. Dacă este *prezent în aval* directorul din comutator păstrează informația referitoare la descendentul ce trebuie să fie căutat pentru a obține o copie a blocului partajat. Dacă *nu este prezent în aval* cererea traversează RMM spre memorie așa cum se indică în eticheta de rutare. În final, procesorul ajunge la deținătorul curent al blocului prin intermediul directorilor din comutatoare, încarcă blocul și devine noul deținător. Astfel fiecare referire de citire sau de scriere a blocului de la un alt procesor implică o mișcare a datelor de la deținătorul curent la acel procesor.

B. Protocol pentru copii multiple.

Acest protocol permite existența unor copii multiple a blocurilor partajate din sistem și se bazează pe *invalidarea la scriere*.

Un bloc din memoria cache poate fi în una din cele trei stări : *invalid în aval*, *valid în aval* și *exclusiv în aval* în raport cu o memorie cache locală sau cu comutatoare de pe un nivel intermediar.

Semnificația stărilor în raport cu comutatorul de pe nivelul *i* sau cu memoriile cache locale este după cum urmează.

- *Invalid în aval* înseamnă că blocul nu este prezent în oricare dintre memoriile cache care sunt în aval față de comutator.
- *Valid în aval* înseamnă că blocul este prezent în starea *valid* în una sau mai multe cache *în aval* de comutator.
- *Exclusiv în aval* înseamnă că blocul este în starea *exclusiv* într-o cache *în aval* de comutator.

Diagrama de tranziții de stare pentru un bloc partajat în raport cu comutatorul de pe un nivel *i* sau cu o memorie cache locală este prezentată în fig. 6.26 (b) Acțiunile efectuate

de controlerele de coerență la diferite nivele pot fi descrise în termeni de *eșec la citire*, *succes la scriere* și *eșec la scriere* într-un bloc partajat.

În cazul unui *succes la citire*, cache-ul local furnizează blocul și nu mai este necesară nici o acțiune. Toate blocurile sunt marcate *invalid* în cache-uri și comutatoarele de pe toate nivelurile sunt aduse în stare inițială.

În cazul unui *eșec la citire* în cazul unei cereri de acces la un bloc partajat din cache-ul local, cererea este transferată comutatorului de pe nivelul 1 care este antecesorul aceluși cache în raport cu blocul ratat. Dacă comutatorul de pe nivelul 1 arată că blocul este în starea *valid în aval* sau în *exclusiv în aval*, se extrage o copie din cache-ul plasat pe nivelul 0 ce deține această copie și ea este furnizată solicitantului. Dacă blocul este în starea *exclusiv în aval*, memoria este actualizată în raport cu politica de rescriere de către cache-ul ce furnizează răspunsul.

Dacă blocul este în starea *invalid în aval*, cererea traversează rețeaua spre comutatoarele de pe nivel superior până când starea blocului este identificată ca fiind fie în *valid în aval* sau în *exclusiv în aval* la un anumit nivel. Dacă blocul este în starea *valid* într-un comutator de pe un nivel superior dincolo de nivelul 1, blocul este extras din memoria principală prin comutatoarele de pe nivel superior. În cazul *exclusiv*, blocul este întotdeauna furnizat de cache-ul ce deține copia exclusivă. Dacă blocul este în starea *invalid în aval* la toate nivelurile, el este extras din memoria principală. Comutatoarele de pe traseul de răspuns schimbă starea blocului în *valid în aval*, în directoarele lor locale.

Dacă există un *succes la scriere* la un bloc aflat în starea *exclusiv în aval* în cache-ul local, blocul este scris imediat și nu este necesară nici o acțiune de menținere a coerenței.

În cazul unui *succes la scriere* asupra unui bloc valid, se transmite un semnal de invalidare spre comutatoarele plasate pe un nivel superior. Semnalul de invalidare traversează rețeaua în sens direct spre cel mai de sus nivel pentru a schimba starea blocului în *exclusiv în aval* în comutatoarele dispuse de-a lungul căii. În același timp, comutatorul de pe nivelul i transmite mai departe, în aval, spre comutatorul dispus pe nivelul $i-1$, un semnal de invalidare, doar dacă blocul este în starea *valid în aval* în acel comutator.

În final, toate cache-urile de pe nivelul 0 ce au o copie validă a blocului invalidează copiile lor și transmit un semnal de validare a invalidării înapoi spre comutatoarele dispuse pe nivelul 1, urmând aceeași cale.

Comutatoarele de pe nivelurile intermediare colectează aceste semnale și după ce toate semnalele sunt recepționate, transmit un semnal de validare a acestora.

Toate comutatoarele de pe calea de invalidare schimbă starea blocului în *invalid în aval*. În final solicitantul recepționează un semnal de validare, schimbă starea blocului în *exclusiv în aval* și finalizează operațiunea. Comutatoarele predecesoare ale solicitantului schimbă starea în *exclusiv în aval*.

În cazul unui *eșec la scriere* în cache-ul local, cererea traversează în sus rețeaua până când ajunge la un comutator ce conține blocul fie în starea *valid în aval* fie în *exclusiv în aval*. Dacă blocul este în starea *valid în aval* în cadrul unui comutator, se transmite un semnal de invalidare spre comutatoarele dispuse pe nivelurile superioare și inferioare.

Semnalul se propagă ca și în cazul *succesului la scriere*. Cache-urile având o copie validă schimbă starea blocului în *invalid*. Blocul este extras din memoria principală prin comutatoarele de pe nivelul superior și este încărcat în stare *exclusiv în aval*. Dacă blocul este în starea *exclusiv în aval* în comutator atunci cache-ul având copia, transmite blocul spre solicitant. Cache-ul furnizor și predecesorii săi schimbă starea blocului în *invalid în aval*. Solicitantul și toți predecesorii săi schimbă starea blocului în *exclusiv în aval*.

C. Organizarea tabelului director.

În SMP există două tipuri de referiri la memorie :

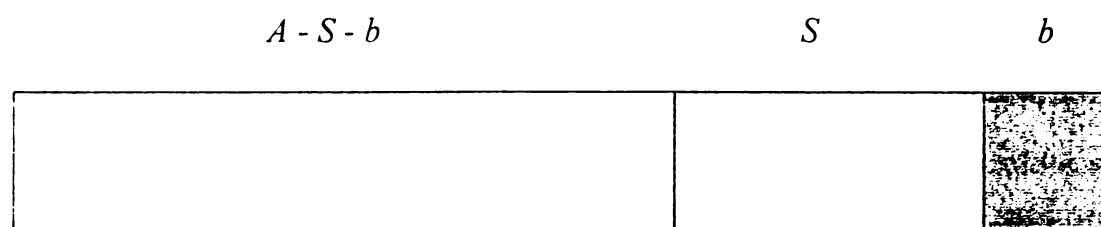
- *private* - instrucții, date private, date partajate de tip citește-numai
- *partajate* - date din locații partajate de tip R/W

Se presupune că în cursul compilării, compilatorul distinge între datele private și cele partajate astfel că ele sunt menținute în modulele de memorie diferite [MN99], [MNLS97].

Presupunem că numărul de blocuri partajate, N_{sb} , în cadrul sistemului este fix și că blocurile partajate sunt distribuite în modulele de memorie. Compilatorul detectează aceste structuri de date și le plasează în locații corespunzătoare.

Fie A numărul total de biți din adresa fizică a unei locații de memorie. Sunt necesari b biți pentru a adresa un cuvânt în interiorul unui bloc de memorie. Prin urmare $A-b$ biți sunt utilizați pentru a adresa un bloc de memorie, privat sau partajat. Deoarece N_{sb} reprezintă numărul de blocuri partajate, sunt necesari $S = \log_2 N_{sb}$ biți pentru a distinge printre blocurile partajate.

O configurație tipică de adresă de memorie este prezentată în fig. 6.25



A - număr de biți în adresa fizică

S - număr de biți pentru adresarea unui bloc partajat

b - număr de biți pentru adresarea unui cuvânt

Fig.6.27 Configurarea adresei fizice

Configurarea se face astfel :

Când c.m.s. biți ($A-S-b$) sunt în intervalul 000...00 la 111...10, adresa este a unui bloc *privat*. Când acești biți sunt toți 1, adresa este a unui bloc *partajat*.

Controlerul memoriei cache poate distinge un bloc partajat prin decodificarea doar a bitului ($S+b+1$). Controlerul validează o linie de comandă indicând că această adresă este a unui bloc partajat. Această linie face parte din liniile de adresă și se transmite comutatoarelor din RMM. Cei S biți sunt utilizați pentru adresarea în mod unic a unui bloc partajat în director.

Directorul prezent într-un comutator RMM este organizat ca și în fig. 6.28.

Talia directorului este fixă și el poate conține informații referitoare la N_{sb} blocuri partajate. Câmpul *Stare* în cadrul unei intrări a directorului memorează starea blocului

partajat. Pentru protocolul prezentat în fig. 6.26(b) sunt necesari doar 2 biți pentru codificarea informațiilor de stare. Câmpul *Prezență* înmagazinează informația asupra blocului în raport cu cei k descendenți ai comutatorului (k este numărul de intrări în comutator).

	<i>Stare</i>	<i>Biți de prezență</i>			
		Ch1	Ch2	...	Ch k
Bloc 1					
Bloc 2					
.			.		
.			.		
.			.		
Bloc N_{sb}					

Fig. 6.28 Organizarea directorului într-un comutator RMM

Dacă câmpul *Stare* indică că blocul este *invalid în aval*, toți biții de prezență sunt pe 0.

Dacă câmpul *Stare* indică că blocul este *exclusiv în aval* doar unul și numai unul din biții de prezență - corespunzător descendentului unde blocul este în această stare - se poziționează pe 1.

Dacă se indică că blocul este *valid in aval*, toți biții de prezență corespunzători descendenților unde blocul este valid sunt 1, ceilalți sunt 0.

Se poate observa că directorul păstrează informația de stare doar referitoare la descendenții săi, și nu la toate procesoarele situate în aval. Aceasta permite ca lungimea cuvântului a unei intrări a directorului să fie redusă și uniformă de-a lungul tuturor nivelurilor. Deoarece organizarea directorului dintr-un comutator RMM este similară cu cea a unei memorii cache și pentru că talia directorului, în termeni de biți, este mult mai mică decât cea a unei cache, directorul încorporat în comutator este cel puțin la fel de rapid ca și un controler al unei memorii cache.

6.4. SIMULAREA PILOTATĂ DE EXECUȚIE

În ultimă instanță, timpul de execuție al unei aplicații rulate pe o arhitectură multiprocesor reprezintă un indicator final al preformanței. În scopul verificării modelului analitic al RMM și a comparării lui cu modelele analitice ale RMC și ale RMB, au fost efectuate simulări ale mai multor aplicații. Simulatorul utilizat [BINK97], [BD91], [Iyer96], [Iyer99] se bazează pe simulatorul Proteus [BD91] dezvoltat la MIT. Acest simulator este un simulator pilotat de execuție ce simulează SMP-MP și SMP-TM. El poate de asemenea simula câteva tipuri de RIN cât și configurații de memorii.

6.4.1 Construcția simulării

Simularea pilotată de execuție constă din două părți de bază : *codul aplicației* și *codul simulatorului* [Iyer96] [Kum96], [KBI96].

Aplicația este scrisă într-un superset C. Se definesc cuvinte-cheie ce sunt utilizate pentru declararea de variabile partajate și operatori pentru accesarea variabilelor partajate. Simulatorul este apelat după un acces la o variabilă partajată sub forma unei citiri sau a unei scrieri. Simulatorul oferă o bibliotecă de rutine pentru transferul de mesaje, managementul memoriei, sincronizare, etc. Sarcina principală în construcția unei simulări ce modelează rețeaua luată în considerare se realizează prin adăugarea sa la codul simulatorului. În fig. 6.29 se prezintă etapele necesare în construcția unei simulări.

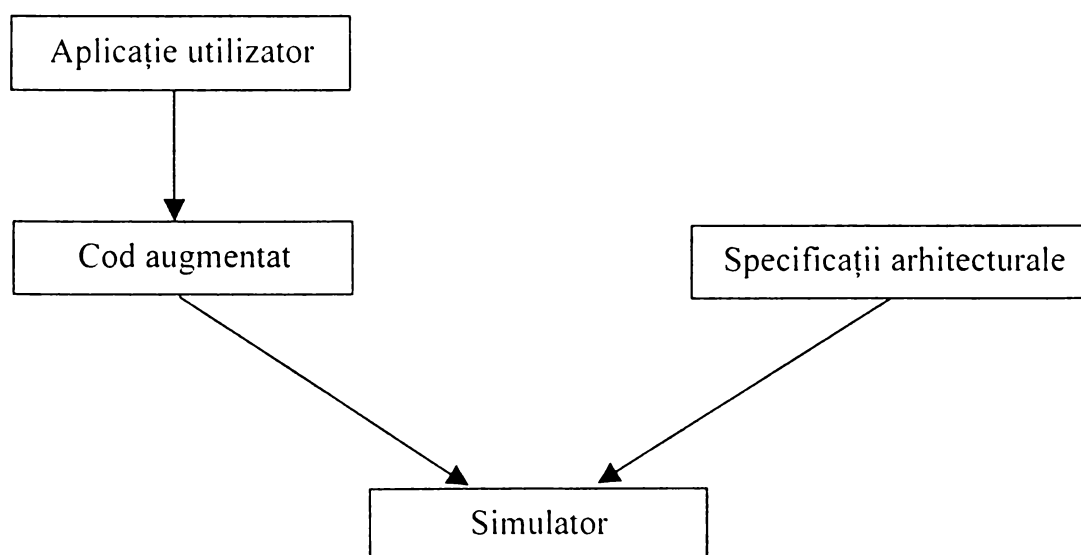


Fig. 6.29 Simulator pilotat prin execuție

6.4.2 Interfața cu rețeaua

Configurația nodului și interfața cu rețeaua în cadrul simulatorului sunt prezentate în fig. 6.30.

Atât controlerul memoriei cache (CC), cât și controlerul memoriei (CM), sunt conectate la interfață, așa că ele pot comunica direct cu alte noduri. Se presupune că procesoarele au un singur flux și nu se efectuează nici o aducere prealabilă a blocurilor de date cache. Sistemul este presupus secvențial consistent, adică nu există tamponare de scriere, și eșecurile de acces cache la operații de tip LOAD/STORE blochează procesorul.

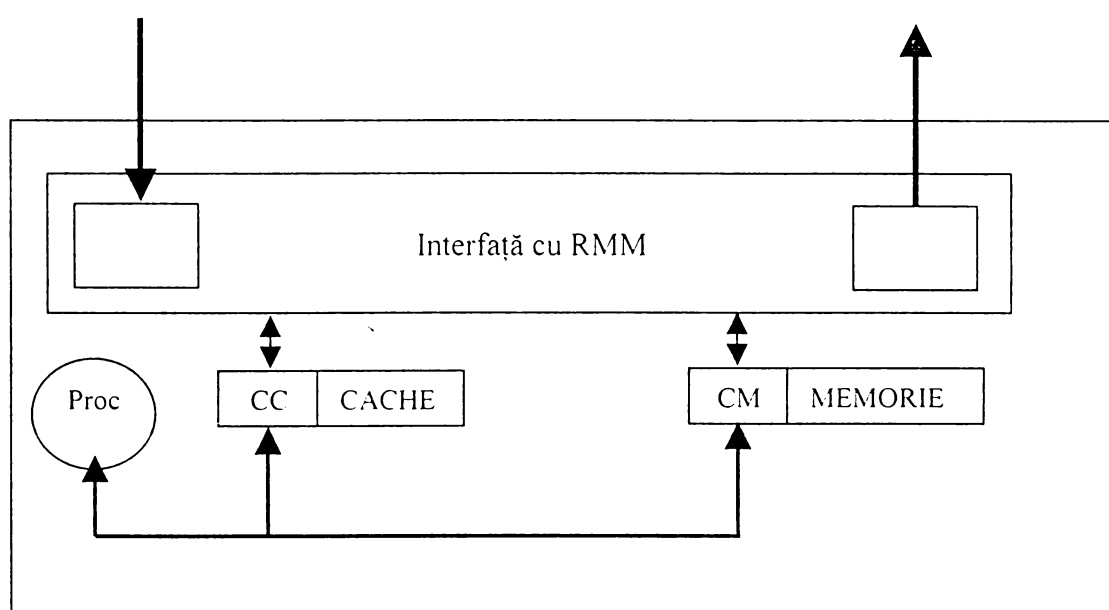


Fig. 6.30 Configurația unui nod și interfața cu rețeaua

Interfața cu rețeaua (IR), furnizează spațiul tampon pentru toate mesajele ce vin sau ce pleacă în/din nod. Se presupune că IR are suficient spațiu pentru a înmagazina toate mesajele. Interfața oferă de asemenea servicii constând în divizarea mesajelor în pachete, inițializarea header-ului sau a oricărui pachet cu informații necesare pentru rutare, etc.

Comutatoarele conțin un număr limitat de tamponare la fiecare intrare destinate pentru recepționarea mesajelor. Spațiul tampon este organizat în termeni de pachete de dimensiune fixată. Mesajele sunt divizate în pachete de dimensiune fixată, iar fiecare pachet conține informația de rutare. Pachetul este transferat de la un comutator la celălalt

ca și un întreg. Un pachet este avansat către comutatorul din etajul următor doar dacă are suficient spațiu pentru a înmagazina întregul pachet. În caz de conflict la ieșire se apelează la disciplina de ordonare *round-robin*.

6.4.3 Parametrii de simulare

Parametrii de sistem utilizați în simulare sunt prezentați în TABEL XV. Am simulat un SMP-MPD cu 16 noduri ce utilizează o rețea cu 4 niveluri, conținând comutatoare 2 X 2. S-a presupus că timpul de întârziere datorat comutatorului este de 1 ciclu, suficient pentru a produce decizia de rutare în ceea ce priveștea header-ul pachetului sau pentru primul flit al mesajului în cazul rutării de tip wormhole. Dimensiunea flit-ului și dimensiunea legăturii este de 16 bit. În acest mod se poate considera că transferul unui flit prin legătură durează un ciclu. În sistemele cache-coerente, mesajele de date conținând blocuri de memorii sunt mai lungi. în timp ce mesajele de cerere de acces și de menținere a coerenței conținând doar informații referitoare la protocolul transmisiei și adrese, sunt mai scurte.. S-au luat în considerare mesaje de lungime de 8 byte pentru menținerea coerenței și de 40 bit pentru date. Aceste dimensiuni se regăsesc și în experimentele citate în literatură [KB93], [BINK97]. În cazul comutării de pachete, pachetele au o lungime de 8 byte (1 – header, 7 – date). Un pachet consumă un timp de întârziere format din 1 + 4 cicli pentru comutare și respectiv pentru transmisia spre nivelul următor din RMM.

TABEL XV
Parametrii pentru simulare

Parametru	Văloare
Număr de procesoare	16
Dimensiunea memoriei partajate în nod	32 Kbytes
Dimensiunea memoriei cache	8 Kbytes
Dimensiunea liniei memoriei cache	32 bytes
Dimensiunea setului	2
Timpul de acces la memoria cache	1
Timpul de acces la memorie	8
Întârzierea de comutare	1
Lățimea conexiunii	16 biți
Lungimea unității de transmis din pachet	16 biți
Dimensiunea pachetului	8 bytes

6.4.4 Protocolul de menținere a coerenței și sincronizarea

Se utilizează un protocol de menținere a coerenței bazat pe director. În această schemă fiecărui nod îi este asignat un bloc de memorie partajată, denumit *nod local* ce menține intrările directorului pentru acel bloc. Fiecare intrare în director este un vector de aceeași lungime ca și numărul de noduri. Directorul menține de asemenea informația despre starea blocurilor. Ori de câte ori o copie a unui bloc de memorie este transmis spre un cache, bitul corespunzător acelui nod este poziționat. Un protocol de invalidare a fost implementat în cadrul căruia toate copiile cache-ate ale unui bloc sunt invalidate la o operație de scriere.

Sosirea unor mesaje în afara ordinii prestabilite provoacă probleme dacă protocolul de menținere a coerenței nu este modificat. De exemplu, o situație poate apărea acolo unde un mesaj de invalidare neplanificat ajunge la un nod înainte ca data solicitată de o cerere de citire să fi ajuns. Dacă invalidarea este confirmată, ea poate conduce la o stare inconsistentă unde un nod crede că are un acces exclusiv la un bloc, în timp ce alt nod are o copie a blocului. Din această cauză se utilizează un protocol de menținere a coerenței modificat unde CC detectează dacă un mesaj a sosit în afara ordinii și îl reține pentru a fi prelucrat ulterior.

Metoda de sincronizare utilizată în cadrul sincronizării se bazează pe utilizarea operației TAS [MCS91]. Au fost implementate bariere utilizându-se un numărător partajat. Sarcina de calcul datorită existenței unui conflict între variabilele de sincronizare este semnificativă în unele aplicații.

6.4.5 Programe de evaluare a performanțelor

Au fost selectate câteva aplicații numerice pentru evaluarea performanței rețelei în SMP-MPD. Aceste aplicații sunt:

- multiplicarea a două matrici 2D – MATMUL
- algoritmul de selectare a căii cele mai scurte dintr-un graf –Floyd–Warshall–FWA
- factorizarea LU a unei matrici dense 2D – LU
- transformarea Fourier rapidă – FFT

Sunt aplicații clasice citate în [BINK97], [BIK00], [HP96].

MATMUL operează asupra a două matrici 128×128 în dublă precizie. Structurile de date principale sunt pentru tabloul 2D de numere reale: două matrici de intrare, o matrice transpusă și una de ieșire. Datele partajate ocupă aproximativ 512 Kbytes.

Pentru FWA se utilizează un graf cu 128 de noduri cu ponderi aleatoare atribuite muchiilor. Structurile principale de date sunt două tablouri 2D de întregi: o matrice distanță și altă matrice predecesor. Datele partajate ocupă aproximativ 128 Kbytes. Programul parcurge atâtea iterații câte vertice există și în timpul unei iterații se citește de către toate procesoarele o linie anume din matricile distanță respectiv predecesor. Fiecare iterație este urmată de o barieră.

Aplicația de descompunere LU utilizează o matrice 256×256 structurată în blocuri 8×8 . Structura principală de date este un tabel 2D având prima dimensiune blocul iar a doua conținând toate punctele de date din acel bloc. Astfel toate punctele, operate de același procesor, sunt alocate în mod contiguu.

Pentru aplicația FFT, la intrare se dau 2^{14} puncte. Structurile principale de date sunt două tablouri 2D de numere complexe. Nu există o partajare a datelor în primele $\log_2(N-P)$ etaje, unde N este numărul de puncte de date și P este numărul de procesoare. În restul celor $\log_2 P$ etaje, fiecare punct de dată este partajat de două procesoare. În aceste etaje se întretes tablouri de intrare și de ieșire pentru a se evita erorile de conflict.

6.4.6 Rezultate ale simulării

Caracteristicile aplicațiilor, luate în considerare în cadrul simulării, sunt prezentate în TABEL XVI. Aceste caracteristici sunt: numărul de referiri la memoria partajată, eșecurile în referirile la memoria partajată, numărul total de mesaje generate în timpul execuției.

TABEL XVI

Caracteristici ale aplicațiilor utilizate în evaluare

Aplicația	Referințe la memorie	Eșecuri acces cache	Mesaje generate
MATMUL	9,215,571	743,258	1,474,220
FWA	11,050,133	323,159	739,624
FFT	2,146,864	394,972	1,093,376
LU	112,270,761	706,579	1,671,854

Pentru început se prezintă latențele medii ale mesajelor în cele două rețele comparate – RMM și RMB – când aplicațiile menționate anterior sunt rulate.

Valorile prezentate în TABEL XVII demonstrează că timpul de răspuns al mesajelor nu diferă semnificativ atunci când se utilizează RMM în loc de RMB.

TABEL XVII

Latențele medii ale mesajelor utilizând RMM și RMB

Aplicația	Latența RMB	Latența RMM
MATMUL	224.48	231.25
FWA	491.24	506.74
FFT	183.86	187.26
LU	162.00	165.93

Performanța RMM într-un mediu mai apropiat de realitate este dată de o evaluare bazată pe timpul de execuție al acestor aplicații.

Simulatorul oferă timpul de execuție al unei aplicații în milioane de cicluri.

În fig. 6.31 se prezintă timpii de execuție pentru diferitele programe de aplicație anterior prezentate.

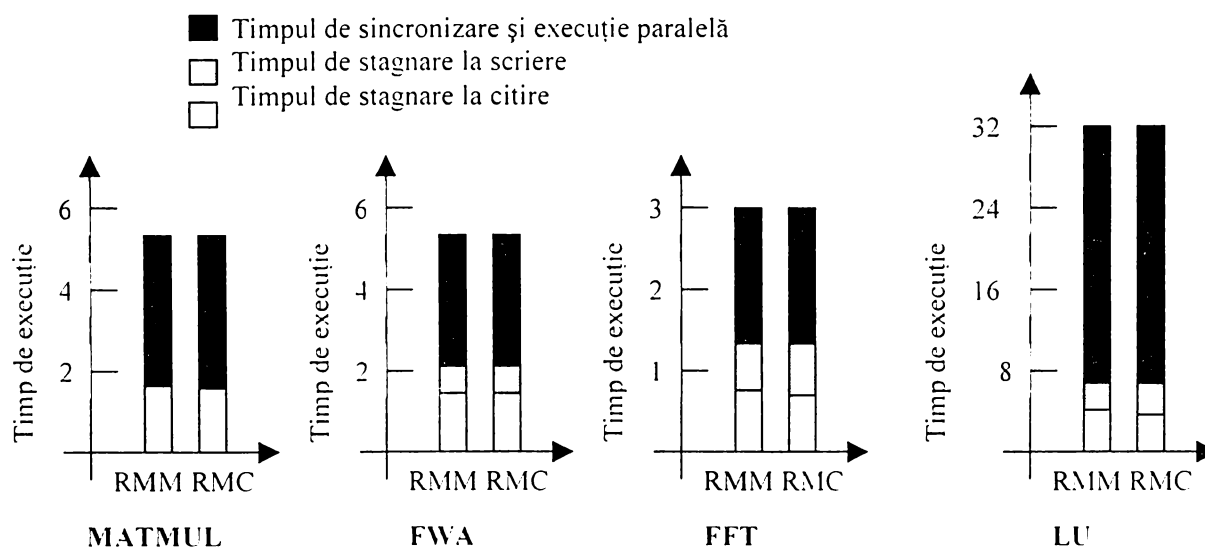


Fig. 6.31 Rezultate ale simulării bazate pe execuție

Timpii sunt repartizați astfel:

- timpul consumat pentru calcul și sincronizare
- timpul de staționare în citire
- timpul de staționare în scriere

Graful din fig. 6.31 este divizat în patru seturi de câte două bare fiecare. Fiecare set este alocat unei anumite aplicații și prezintă timpul de execuție utilizând RMM în comparație cu RMB. Se observă că performanța RMM este foarte apropiată de cea a RMB. Ameliorarea timpului de execuție în general se ameliorează datorită timpului de staționare în citire și a timpului de staționare în scriere ce este direct proporțional cu latența rețelei. La RMM timpul de staționare în scriere este mai coborât decât cel de citire datorită unui număr mai mic de eșecuri la scriere.

6.5 EFECTUL POLITICILOR DE MANAGEMENT AL MEMORIEI

Non-uniformitatea în accesele la memorie reprezintă un subiect ce are un impact semnificativ asupra performanțelor sistemelor multiprocesor cu memorie partajată distribuită. Alocarea și plasarea datelor în astfel de sisteme tinde să creeze puncte critice, din punct de vedere al acumulării de mesaje, în rețea și în modulele de memorie. În literatură [BR90], [RN91], [BIK00] se remarcă efortul cercetătorilor de a dezvolta abordări experimentale în reducerea sau chiar în eliminarea unor astfel de puncte critice. Toate aceste politici se referă la mișcarea *dinamică* a datelor utilizând politici de replicare și de migrare a lor. Aceste abordări măresc semnificativ costul total deoarece necesită resurse hardware suplimentare. În acest paragraf ne vom referi la organizări ale ierarhiei memoriilor și la politici de plasare în mod *static* a paginilor ce pot fi utilizate pentru a ameliora în mod semnificativ performanța.

Aceste politici sunt:

1. Întreșeserea în Grad Înalt (High-Order Interleaving) – IGI
2. Întreșeserea în Grad Redus (Low-Order Interleaving) – IGR
3. Plasarea Paginii cu Întreșeserea în Grad Înalt (Page Placement with High-Order Interleaving) – PPIGI

4. Plasarea Paginii cu Întreșeserea în Grad Redus (Page Placement with Low-Order Interleaving) – PPIGR
5. Plasarea Paginii bazat pe Tipul de Acces (Access Pattern based Page Placement with High-Order Interleaving) – PPAIGI

În primul tip de organizare – IGI, atât memoria privată cât și cea partajată sunt strâns întreșesute.

A doua organizare reprezintă politica de întreșesere în grad redus. În figura 6.32 se prezintă organizarea modelului experimental în situația existenței a 4 procesoare. Fiecare procesor P_i are o memorie cache C_i ce este conectată la RIN prin interfața I_i . Spațiul de adresă al memoriei este divizat în două părți: *memoria privată* PM și *memoria partajată* SM. Memoria PM este uzual întreșesută într-un grad ridicat, în schimb SM, în termeni de blocuri, este slab întreșesută. Rezultatul constă în maniera distribuită în care blocurile de date partajate sunt accesate. Se evită astfel apariția punctelor critice de acumulare a acceselor.

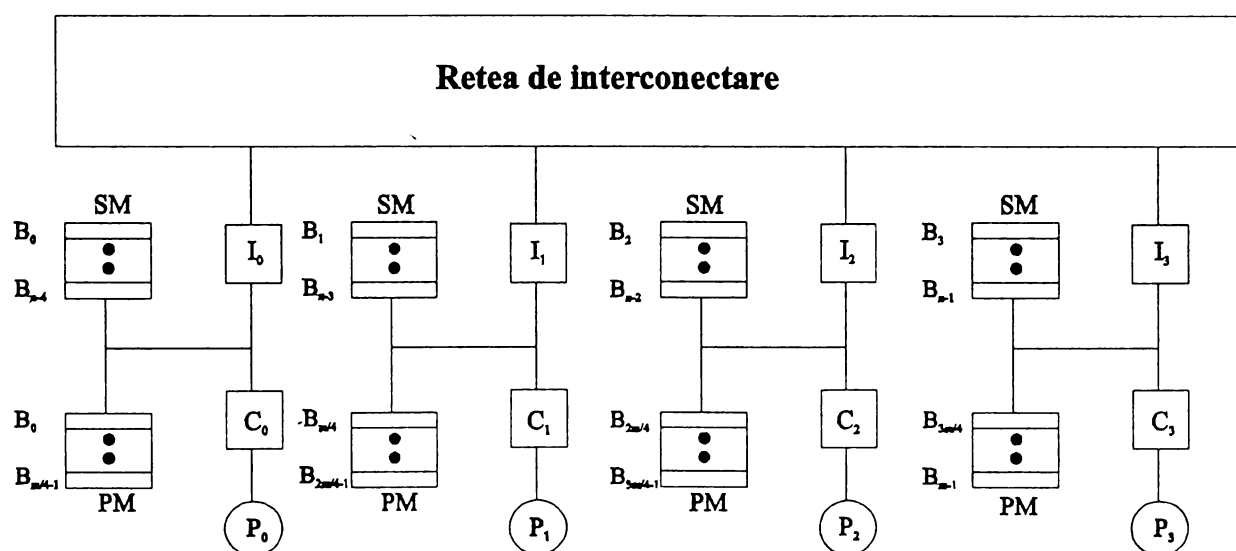


Fig. 6.32 Organizare IGR a unui SMP – MPD

În a treia, respectiv a patra organizare, se utilizează diferite politici de întreșesere pentru implementarea amplasării paginilor. În exemplul prezentat, amplasarea paginii se bazează pe apariția erorii în procesor. Când un procesor se defectează la un acces al unui bloc, pagina corespunzătoare este extrasă dintr-o memorie secundară și plasată în memoria locală a procesorului. Dacă spațiul memoriei locale partajate a procesorului unde apare

eroarea nu este disponibil pentru plasarea paginii, atunci pagina este plasată la cel mai apropiat procesor.

A cincea politică este una în care aplicațiile sunt rulate cu o tehnică PPIGI și numărul paginii accesate de către fiecare procesor este salvat ca și informație de configurație de acces. Cu ajutorul acestor configurații de acces a paginii se definește un algoritm pentru determinarea celei mai bune amplasări atunci când informația este dată. Această politică însă este dificil de implementat deoarece pretinde cunoașterea configurațiilor de acces ale paginii. Accesele la pagină ar putea fi estimate cu unele tehnici asistate de compilator [Iyer 96].

În literatură [Iyer 96] studiul acestor politici s-a efectuat prin incorporarea tehnicilor mai sus amintite în simulatoare pilotate de aplicație și derivate din simulatorul Proteus. S-a introdus paramentru *Page_Size* și s-a optat pentru o valoare de 1024 a dimensiunii paginii. Atunci când dimensiunea tamponului se mărește de la un pachet la două pachete, performanța se ameliorează semnificativ, dar în continuare, orice augmentare a taliei tamponului nu mai produce îmbunătățiri relevante.

6.5.1 Algoritm pentru determinarea amplasării paginilor bazată pe acces

În literatură [Iyer 96], se precizează că politica PPAIGI de gestionare a memoriei implică determinarea celei mai bune amplasări a datelor când sunt cunoscute *a priori* câteva informații de către alocaturul de memorie. Cea mai utilizată informație o reprezintă numărul de accese la oricare pagină ce sunt efectuate din partea oricărui procesor. Algoritmul prezentat parcurge o serie de iterații, având această informație, pentru a determina cea mai bună amplasare a fiecărei pagini. Se utilizează un tabel bidimensional P_{to_p} ce reprezintă harta de alocare *procesor la pagină*.

Algoritmul este următorul:

1. Pentru fiecare pagină, se extrage nodul al cărui procesor efectuează numărul maxim de accese și se introduce această valoare în tabelul P_{to_p} .
2. Pentru fiecare nod, se extrag atât de multe pagini cât se pot potrivi din linia respectivă în tabelul P_{to_p} cu scopul de a reduce numărul de accese.
3. Se verifică dacă toate paginile au fost alocate.
4. Dacă pasul 3 este *adevărat* atunci *exit*.

5. Dacă pasul 3 este *fals*, atunci:
 - a) Se extrag paginile anterior alocate din lista de alocare
 - b) Se extrag procesoarele selectate de către paginile nealocate din listele lor individuale și
 - c) Se pornește din nou de la pasul 1 până când toate paginile disponibile devin memorie alocată

Algoritmul este astfel proiectat încât paginile și procesoarele sunt ambele beneficiare ale alocării. Aceasta ar putea reduce numărul de accese la un nod îndepărtat deoarece procesoarele selectează o pagină să fie în memoria locală a nodului său dacă frecvența de acces este ridicată.

Deficiența majoră a acestui algoritm constă în faptul că informația este reprezentată doar de *numărul de accese* și nu de *perioada acceselor*. Dacă o etichetă de marcarea a timpului este asociată fiecărui acces, s-ar putea ameliora semnificativ algoritmul, astfel încât pagina să se plaseze perfect. Cea mai potrivită informație de care este nevoie este calea critică în execuție.

6.6 CONCLUZII

Sistemul multiprocesor utilizat ca și model experimental este un sistem cu memorie partajată distribuită de tip CC-NUMA.

Rețeaua de interconectare este o Rețea cu Magistrală Multinivel. S-a analizat RMM și s-a demonstrat potențialul său în raport cu RIN existente. Inițial RMM a fost comparată cu o RIN convențională bidirecțională (RMB). S-au prezentat patru sisteme de autorutare pentru RMM ce sunt valabile și pentru RMB. A fost prezentat un algoritm pentru selectarea căii de distanță minimă dintre o sursă și o destinație. Acest algoritm reduce numărul de niveluri traversate în RIN.

Pentru evaluarea performanțelor RMM și ale RMB s-au utilizat analize probabilistice și de formare a cozilor de așteptare. Simulările au confirmat că RMM se comportă mai bine în comparație cu RMC și similar în raport cu RMB. A rezultat că scalabilitatea unui RMM este aproape identică cu cea a unui RMB.

Performanța unui RMM este similară cu cea a unui RMB în termeni de timp de răspuns și de utilizare a procesorului. Ea este însă mai bună decât a unui RMC.

Pentru a reliefa potențialul de performanță al unei RMM, a fost prezentată o evaluare bazată pe execuție într-un mediu multiprocesor cu memorie partajată, coerent din punct de vedere cache. Mecanismul de asigurare a coerenței cache se bazează pe director.

Simulatorul bazat pe execuție utilizat este PROTEUS. S-a evidențiat că pentru un număr de patru aplicații performanța RMM este aproximativ egală cu cea a unei RMB în termeni de latență a mesajului și timpi de execuție. Întrucât RMM oferă o simplitate din punct de vedere hardware, acest tip de rețea reprezintă o alternativă viabilă la RIN existente. În plus, există facilități de toleranță la defect incorporate. Studiul politicilor de gestionare a memoriei scot în evidență faptul că modul de plasare a paginii depinde puternic de configurațiile de acces ale aplicațiilor și sunt necesare tehnici de investigare în termeni de politici de mișcare în mod dinamic a paginii. Se evidențiază că politica de plasare a paginii cu un grad coborât de întrețesere reprezintă cea mai bună alternativă deoarece ameliorează întrețeserea cu grad redus prin plasarea primului bloc la procesorul defect.

7. CONCLUZII GENERALE ȘI DIRECȚII VIITOARE DE CERCETARE

A. Concluzii generale

Lucrarea de față vizează un domeniu actual din știința calculatoarelor, acela al arhitecturilor paralele din clasa MIMD cu facilități de toleranță la defect.

Actualitatea problematicii abordate în teză rezidă din faptul că sistemele multiprocesor au devenit industriale și prin urmare, pe lângă performanța ce validează anumite concepte teoretice, dependabilitatea acestor sisteme devine un factor determinant. Pe cale de consecință, trebuie abordată proiectarea sistemelor multiprocesor la nivel de arhitectură, prin prisma unui raport corect între performanță și dependabilitate.

Lucrarea este concepută ca și o dezvoltare succesivă a temelor majore ce apar în *clasificarea, analiza, proiectarea și simularea SMP* din clasa MIMD. Fiecare capitol se finalizează cu concluzii și recomandări ce contribuie la abordarea mai eficientă a tematicilor conținute în capitolul următor.

Capitolul 1 este dedicat prezentării principalelor tendințe arhitecturale în dezvoltarea SMP cât și a unor noi concepte ce permit o abordare mult mai nuanțată a arhitecturilor paralele. Examinarea sistematică, în termeni calitativi și cantitativi, a diferitelor arhitecturi se face utilizând conceptul de *convergență* între paradigmele *memorie partajată* și *transfer de mesaje*. Prin prisma acestui concept sunt descrise în continuare patru teme majore pentru proiectare: *modelul de comunicare, modelul de programare, comunicarea și replicarea, performanța. Coerența și consistența* memoriilor cache reprezintă subiecte majore în abordarea corectă a reducerii latenței comunicației. Diferite modele de consistență a memoriei au fost tratate ulterior în Capitolul 2 și în Capitolul 3. *Arhitectura rețelelor de interconectare* contribuie determinant la reducerea latenței și la tolerarea unei creșteri semnificative a procesoarelor conectate, menținând scalabilitatea. Tot în acest capitol sunt prezentate principalele metode de obținere a tolerării latenței: *reducerea timpului de acces la fiecare nivel al unei ierarhii extinse de memorie, structurarea sistemului pentru a reduce frecvența acceselor cu latență ridicată, structurarea aplicației pentru a reduce frecvența acceselor cu latență ridicată*. În fine, se introduce *toleranța la defect* ca și un obiectiv major în proiectarea SMP și în evaluarea performanței sale.

În Capitolul 2 se prezintă o analiză a arhitecturilor SMP din clasa MIMD din perspectiva conceptelor prezentate anterior.

Se constată că în ceea ce privește SMP-MP simetrice elementele cheie din punct de vedere tehnic sunt organizarea și implementarea sistemului de memorie partajată, care este utilizată pentru comunicarea între procesoare, suplimentar la gestionarea tuturor acceselor la memorie. Cele mai multe dintre SMP-MP cu număr redus de procesoare utilizează magistrala sistemului ca și un mediu de interconectare pentru comunicare și problemele care se ridică sunt menținerea coerenței datelor partajate în MCH private ale procesoarelor și asigurarea toleranței la defect. Există o mare varietate de opțiuni disponibile arhitectului de sistem, incluzând seturi de stări asociate cu blocurile din MCH, acțiuni și tranzacții de magistrală, alegerea dimensiunii unui bloc cache și protocoale de actualizare și invalidare. Sarcina principală a arhitectului de sistem este de a face o corectă alegere între menținerea în limite acceptabile ale configurațiilor de date partajate și de a face implementarea lor cât mai ușoară.

Se poate aștepta ca SMP-MP de scară redusă să continue să rămână atractive din cel puțin trei motive. Primul îl reprezintă raportul cost/performanță atractiv. În al doilea rând microprocesoarele actuale sunt proiectate să fie integrate în medii multiprocesor, ele având tot setul și toate facilitățile necesare interfațării în astfel de sisteme. În fine, al treilea motiv îl reprezintă existența unui suport software performant pentru mașinile paralele.

Deși interconectarea prin magistrala partajată reprezintă un mecanism puternic, există în mod clar limitări în ceea ce privește lățimea de bandă pe măsură ce numărul de procesoare sau viteza procesorului cresc. O soluție generală în construcția de SMP-MP coerente din punct de vedere cache și scalabile, este aceea de a distribui fizic memoria printre noduri și de a utiliza o interconectare scalabilă împreună cu protocoale de menținere a coerenței ce nu se bazează pe monitorizare. Indiferent de evoluție, plasarea rețelelor de interconectare în cadrul ierarhiei de memorie, problema menținerii coerenței cache și diferite protocoale de menținere a coerenței la nivel de tranziție de stare sunt fundamentale pentru toate arhitecturile SMP-MP coerente din punct de vedere cache, independent de modul de interconectare utilizat.

Cele mai moderne sisteme construite pe scară largă sunt proiectate pornind de la noduri de procesare cu utilizare generală, dotate cu ierarhie de memorie, ce au în plus și facilități de interfațare, până la rețele scalabile. Este disponibilă o gamă largă de opțiuni de proiectare pentru asistarea facilităților de comunicare. Proiectarea este puternic influențată de nivelul unde are loc interfațarea dintre facilitatea de comunicare cu arhitectura nodului: nivelul procesorului, al controllerului MCH, al magistralei de memorie sau al magistralei I/E. Este de asemenea influențată de arhitectura de comunicație și de modelul de programare. Modelele de programare sunt implementate utilizând protocoale ce se bazează pe primitive de tip tranzacții de rețea.

Proiectarea SMP este strâns legată de revoluția tehnologică a anilor '80 – microprocesorul și comutatorul de rețea scalabilă încapsulate fiecare într-un singur chip. Ca și rezultat, proiectarea mașinilor paralele s-a divizat în două direcții: SMP orientate în jurul conceptelor de transfer de mesaj, de grupuri de procesoare, etc. și SMP ce integrează într-un grad ridicat RIN în sistemul de memorie pentru a furniza un acces coerent cache la un spațiu global de adresă fizică.

În Capitolul 3 se analizează performanța diferitelor arhitecturi de SMP. În urma analizei a trei tipuri de RIN și anume: RIN-MM, RIN-MN și RIN-CB se constată că RIN-CB manifestă cea mai bună probabilitate de acceptanță în raport cu creșterea numărului de procesoare având o lățime de bandă scalabilă și o bună utilizare a procesorului. Toleranța la defect este acceptabilă, cel mai bun comportament însă în acest sens fiind manifestat de RIN-MM. Sub raport performanță/cost la un număr mare de procesoare, RIN-MN rămâne însă soluția cea mai bună. Utilizarea unor tehnici moderne de transmisie a mesajelor ca și tehnica *wormhole*, cât și introducerea unor capacități redundante încorporate în rețea elimină în mare măsură deficiențele evidențiate anterior în urma analizei. În acest sens, rețelele ASEN devin atractive prin faptul că pot deveni tolerante la mai mult de un defect.

Reducerea semnificativă a latenței unui SMP poate fi obținută și prin utilizarea MCH private în cadrul unei ierarhii de memorie. Există însă o interdependență între caracteristicile arhitecturale și mecanismele de menținere a coerenței. Aceste caracteristici sunt: *strategia de detectare a coerenței*, *strategia de întărire a coerenței*.

dimensiunea informației referitoare la partajarea blocului, dimensiunea blocului cache cât și altele, mai puțin importante dar nu de neglijat.

În arhitecturile moderne de SMP se pot combina diferite mecanisme de menținere a coerenței bazate pe *snooping* împreună cu cele bazate pe *director*. Abordările hibride vor oferi cea mai bună oportunitate pentru creșterea performanței și pentru reducerea costului de implementare.-

Capitolul 4 este dedicat analizei unor arhitecturi MIMD tolerante la defect. Se analizează mai multe abordări propuse pentru parcurgerea celor trei faze ale toleranței la defect dinamice: *detectarea defectului și localizarea lui, recuperarea din defect și reconfigurarea sistemului în jurul procesorului defect*. Referitor la recuperarea din defect câteva concluzii se impun și anume: tehnicile de recuperare sunt diferite pentru SMP-MP, SMP-TM și SMP-MPD. Pe de altă parte atât *rularea înapoi* cât și *rularea înainte* utilizează tehnici de creare și de exploatare a PV-urilor. Ierarhia de memorie oferă soluții de amplasare atât pentru datele active cât și pentru datele PV-urilor. Atât recuperarea din defect prin *rulare înapoi* cât și prin *rulare înainte* reprezintă procese complexe în care crearea corectă și gestionarea PV-urilor sunt cele mai importante aspecte. Totuși, coordonarea creerii și exploatării PV-urilor este neglijabilă în comparație cu sarcina de salvare a stărilor.

Referitor la proiectarea și utilizarea RIN-MN tolerante la defect se constată existența unor caracteristici comune ce includ *modelul de defect, criteriul de toleranță la defect, metoda de creare a toleranței la defect, toleranța la defecte singulare sau multiple, complexitatea rutării și complexitatea hardware*. Tehnicile de obținere a toleranței la defect sunt în general independente de topologia rețelei. Acestea sunt: transmisia multiplă pe căi distincte a aceleiași informații, parcurgerea de mai multe ori a RIN-MN, etc.

În faza de proiectare însă se pot prevedea tehnici ce modifică topologia unei RIN-MN și care în general se caracterizează prin intervenții hardware la nivel de: comutator, niveluri de comutatoare, conexiuni suplimentare, etc. În ultimă instanță replicarea *in integrum* a unei RIN-MN reprezintă o soluție fiabilă dar costisitoare. Finalitatea ultimă însă o reprezintă stabilirea unui raport corect performanță/cost în ceea ce privește gradul de acoperire la defect.

Datorită unor condiții restrictive ce au fost expuse pe larg în acest capitol se constată că doar câteva din RIN-MN sunt tolerante la defect singular. Condiția creșterii numărului de rețele tolerante la acest tip de defect o reprezintă relaxarea modelului de defect comun și presupunerea unei siguranțe ridicate în funcționare a nivelurilor de comutatoare. Aria de cercetare a RIN-MN tolerante la defect este departe de a fi epuizată iar abordarea noilor tehnologii hardware și software în contextul mai larg al dependabilității vor ameliora semnificativ performanța acestora.

În Capitolul 5 se tratează modelarea dependabilității pentru SMP-MP. Sunt prezentate etapele preliminare în descrierea unui model care sunt: *definirea*, *parametrizarea* și *evaluarea*. Modelele de fiabilitate luate în considerare sunt cele de tip *terminal*, de tip *multiterminal* bazat pe sarcină și de tip *rețea*. Se consideră fiabilitatea ca fiind una dintre cele mai importante atribute ale SMP. Cele mai utilizate tehnici de analiză de fiabilitate sunt *modelarea combinatorie* și *modelarea prin lanțuri Markov*.

Modelarea combinatorie bazată pe tehnici din teoria probabilității reliefează modurile diferite în care un sistem poate rămâne operațional. Fiabilitatea sistemului este în general exprimată în termeni de fiabilități ale componentelor individuale. În practică se utilizează două abordări și anume *modelarea serie* respectiv *modelarea paralel*. Se remarcă un dezavantaj major al metodei combinatorii și anume faptul că sistemele complexe nu pot fi modelate convenabil deoarece expresiile fiabilității devin foarte complexe. De asemenea acoperirea defectului este dificil de încorporat în expresiile combinaționale.

Modelele Markov utilizează concepte de *stare* și de *tranziție a stării*. Principalul beneficiu al utilizării MK în modelarea dependabilității îl reprezintă abilitatea de a gestiona acoperirea defectului într-o manieră sistematică. Incorporarea acoperirii C în analiza MK permite analiza siguranței sistemului.

Modelarea disponibilității presupune încorporarea conceptului de *reparare*. Lanțurile Markov au proprietatea de a fi ciclice și, prin urmare, se permite reîntoarcerea de la o stare mai puțin operațională la o stare mai operațională sau chiar complet operațională. Evaluarea experimentală a disponibilității nu este întotdeauna posibilă din cauza constrângerilor de timp și de cost. Din această cauză se utilizează frecvent parametrii *MTTF* sau *MTTR* ceea ce conduce la evaluarea disponibilității stării ferme. A_{SS} .

În fine, se prezintă modele de dependabilitate raportată la performanță introducându-se noțiunea de *performabilitate* ce reprezintă funcția de distribuție a probabilității performanței de sistem acumulată în timp.

Se abordează în mod sistematic modelarea disponibilității utilizând *descompunerea de sistem* pentru SMP-MP orientate pe RIN-MN. Se tratează configurații cu 16 noduri și un model generalizat.

Arhitectura SMP-MP, utilizat ca și model experimental, este prezentată în Capitolul 6. Acest sistem este unul cu memorie partajată distribuită, coerent din punct de vedere cache, și cu acces la memorie de tip NUMA. Rețeaua de interconectare este o rețea de tip RMM. S-au considerat 16 noduri procesoare și o ierarhie de memorie structurată pe 3 straturi. S-a analizat RMM și s-a demonstrat potențialul său în raport cu RIN existente. Rețelele de interconectare luate în considerare la comparație au fost o RIN convențională bidirecțională (RMB) și o RIN convențională (RMC). S-au elaborat doi algoritmi pentru cercetarea comportării rețelei RMM. Primul algoritm determină numărul de căi posibile între o sursă și toate destinațiile, în timp ce al doilea determină lungimea minimă de cale între o sursă și o destinație dată. Rezultatele au fost prezentate în anexele 1, 2 și 3. S-au evidențiat tehnicile de rutare RD, RI, BD și BI.

Pentru evaluarea performanțelor RMM, RMB și RMC s-au utilizat analize probabilistice și de formare a cozilor de așteptare. S-a constatat că RMM se comportă mai bine în comparație cu RMC și similar în raport cu RMB. Scalabilitatea unui RMM este aproximativ similară cu cea a unui RMB. Performanța unui RMM este similară cu cea a unui RMB în termeni de timp de răspuns și de utilizare a procesorului. Avantajul RMM îl reprezintă însă o simplificare din punct de vedere hardware al comutatorului în comparație cu comutatorul clasic *crossbar* utilizat în RMB. Există și facilități de toleranță la defect incorporate. Simulatorul bazat pe execuție care a fost utilizat este Proteus elaborat în laboratoarele MIT.

B. Direcții de cercetare

În ultima perioadă de timp se constată o convergență semnificativă între arhitecturile SMP-MP și SMP-TM. Cele mai recente SMP-MPD utilizează arhitecturi bazate pe accesul neuniform la memorie și coerente cache (CC-NUMA).

Performanța unui calculator este primordial dictată de existența unei latențe reduse cât și de furnizarea unei mari lățimi de bandă. Performanța raportată la dependabilitate reprezintă un raport corect ce trebuie stabilit între o performanță ușor diminuată în avantajul unei toleranțe la defect confortabile. Evident, raportul cost/performanță în aceste condiții trebuie să fie unul cât mai acceptabil. Prin urmare, direcțiile viitoare de cercetare în proiectarea SMP-MPD tolerante la defect trebuie să vizeze *reducerea și tolerarea latenței, creșterea lățimii de bandă cu păstrarea scalabilității și creșterea toleranței la defect.*

O direcție de cercetare viitoare o reprezintă proiectarea comutatoarelor cu capacități de memorie cache încorporate. Procesoarele uzuale ca și Pentium Pro și MIPS R10000 sunt superscalare, utilizând niveluri multiple de memorii cache și tolerând mai multe eșecuri în acces decât microprocesoarele clasice. Proprietatea de superscalaritate are trei implicații în traficul prin rețea ce trebuie să fie studiate atent. În primul rând, intensitatea traficului poate fi redusă datorită nivelurilor multiple de cache-are. În al doilea rând, numărul de cereri de acces simultane poate crește semnificativ. În final, analizarea impactului latenței RIN asupra timpului de stagnare a aplicației va deveni dificilă datorită suprapunerii comunicării cu calculul.

Proiectarea comutatoarelor *crossbar* sau cu *magistrală inclusă* trebuie să adopte soluții hardware de cache-are. De exemplu înglobarea unor MCH de tip SRAM în comutatoarele *crossbar* ar permite captarea datelor partajate pe măsură ce ele parcurg RIN și apoi livrarea lor, în mod eficient, către următorii solicitanți.

În ceea ce privește politicile de management al memoriei, utilizarea unor politici de management dinamic poate reduce solicitările ce congestionează traficul prin RIN, generate de aplicații, printr-o mișcare optimă a datelor spre memoria locală.

Se cunoaște că utilizarea unor tehnici de mascare a latenței, ca de exemplu aducerea în avans a datelor din memorie, poate produce o degradare semnificativă în performanță.

Dacă însă aducerea în avans a datelor se produce în MCH utilizate ca și comutatoare, atunci aceste tehnici își recâștigă eficiența și ameliorează semnificativ performanța.

C. Contribuții personale ale autorului

Pornind de la obiectivele declarate ale acestei dizertații, în continuare se evidențiază principalele contribuții originale, structurate pe capitole:

- 1.1. Abordarea sistematică a arhitecturilor sistemelor multiprocesor din clasa MIMD, utilizând concepte moderne ca și: *convergența arhitecturilor paralele, evaluarea și analizarea arhitecturilor prin aplicații, metodologii sistematice pentru evaluarea performanței.*
- 1.2. Utilizarea în cadrul caracterizării și analizei diferitelor arhitecturi a unor concepte de proiectare fundamentale ca și: *modelul de comunicare, modelul de programare, comunicarea și replicarea, performanța.*
- 1.3. Utilizarea în analiza arhitecturilor sistemelor multiprocesor cu memorie partajată a unor concepte-cheie ca și: *coerența memoriilor cache, consistența memoriei, sincronizarea.*
- 2.1. Prezentarea unei taxonomii a schemelor de asigurare a coerenței memoriilor cache
- 2.2. Sistematizarea prezentării rețelelor de interconectare în funcție de: structura organizatorică, topologia de interconectare, mecanisme de rutare, structura comutatorului.
- 3.1. Prezentarea sistematică a metodelor de analiză de performanță la sistemele multiprocesor din clasa MIMD.
- 3.2. Analiza de performanță la rețelele de interconectare cu prezentarea comparativă a rezultatelor.
- 3.3. Sistematizarea unor scheme și a unor protocoale de asigurare a coerenței memoriilor cache la multiprocesoarele cu memorie partajată, respectiv la multiprocesoarele cu transfer de mesaj.
- 3.4. Evidențierea și sistematizarea caracteristicilor arhitecturale ce afectează mecanismele de menținere a coerenței memoriilor cache.

- 4.1. Sintetizarea unor parametri ce definesc *dependabilitatea* ca fiind calitatea serviciilor furnizate de un sistem de calcul.
- 4.2. Abordarea sistematică a metodelor de simulare a toleranței la defect, a strategiilor de recuperare din defect și a tehnicilor de reconfigurare în sistemele multiprocesor.
- 4.3. Sistemătizarea metodelor și tehnicilor de obținere a toleranței la defect pentru sistemele multiprocesor.
- 5.1. Abordarea într-o variantă unitară a modelelor de dependabilitate și raportarea lor la performanțe.
- 5.2. Abordarea modelării disponibilității, utilizând descompunerea de sistem pentru sistemele multiprocesor cu memorie partajată orientate pe rețele de interconectare multinivel.
- 6.1. Definirea unei arhitecturi de sistem multiprocesor cu memorie partajată distribuită având 16 noduri procesoare în vederea evaluării performanței.
- 6.2. Proiectarea unei rețele de interconectare cu magistrală multinivel (RMM) de dimensiuni 16 X 16.
- 6.3. Elaborarea unui algoritm și a programului de calcul pentru determinarea tuturor rutelor posibile într-o RMM cu $N=16, 64, 256, 1024$ de noduri cât și a numărului minim de comutatoare traversate.
- 6.4. Elaborarea unui algoritm și a programului de calcul pentru determinarea căii optime într-o RMM cu $N=16, 64, 256, 1024$ de noduri.
- 6.5. Evaluarea performanțelor rețelei cu magistrală multinivel (RMM), rețelei multinivel bidirecțională (RMB) și a rețelei multinivel convențională (RMC) pentru a evidenția comportamentul RMM.
- 6.6. Utilizarea simulatorului PROTEUS

ANEXA 1

Program de determinare a numărului de căi pentru fiecare tip de rutare între toate sursele și destinațiile cât și de calcul a lungimilor minime ale căilor.

```

#include <stdio.h>
#define D 256 //numar procesoare
#define MARGD 15 // log2(D)*2-1 - indicele marginii din dreapta

int stiva[2*D][3]; //stiva pentru determinarea rutelor
int stat[D][4]; //vector de numarare a rutelor
int min[D][4]; //vector de calcul a minimelor

// initializeaza vectorul de numarare a rutelor
void init_stat(void) {
    int i;

    for(i=0;i<D;i++) {
        stat[i][0]=0;
        stat[i][1]=0;
        stat[i][2]=0;
        stat[i][3]=0;
    }
}

// tipareste vectorul de numarare a rutelor
void print_stat(void) {
    int i;

    for(i=0;i<D;i++)
        if (i!=stiva[1][0])

printf("%d\t%d\t%d\t%d\t%d\n",i,stat[i][0],stat[i][1],stat[i][2],stat[i][3
]);
        else
            printf("%d\t-\t-\t-\t-\n",i);
    }

//initializeaza vectorul de calcul al minimelor
void init_min(void) {
    int i;

    for(i=0;i<D;i++) {
        min[i][0]=-1;
        min[i][1]=-1;
        min[i][2]=-1;
        min[i][3]=-1;
    }
}

//calculeaza minimul pe o pozitie si un tip de cale
void save_min(int k,int i) {
    int i;

```

```

    j=stiva[k+1][0];
    if ((min[j][i]==-1) || (min[j][i]>(k+1)/2))
        min[j][i]=(k+1)/2;
    }

//tipareste vectorul de calcul al minimelor
void print_min(void) {
    int i;

    for(i=0;i<D;i++)
        if (i!=stiva[1][0])

        printf("%d\t%d\t%d\t%d\t%d\n",i,min[i][0],min[i][1],min[i][2],min[i][3]);
        else
            printf("%d\t0\t0\t0\t0 \n",i);
    }

//initializeaza stiva pe nodul sursa
void init_st(int k, int s, int margin) {

    stiva[k][0] = s; stiva[k][1] = margin; stiva[k][2] = 0;

    }

//rotire stanga
int rol (int i) {
    int tmp, ret;

    tmp = i % 2;
    ret = i / 2;
    if (tmp)
        ret += D / 2;
    return ret;
    }

//rotire dreapta
int ror (int i) {
    int tmp=0, ret;

    if (i>=D/2) {
        tmp = 1;
        i -= D/2;
    }
    ret = i*2 + tmp;
    return ret;
    }

//genereaza succesorul unui nod
int succ(int k) {

    int i=stiva[k][0];
    int j=stiva[k][1];
    int dir=stiva[k][2];

    if (dir < 4) {

        //in functie de directia (0-3) selectata genereaza succesorul

```

```

switch (dir) {

    //continuarea sensului pe acelasi nivel orizontal
    case 0: stiva[k+1][0] = i;
           stiva[k+1][2] = 0;
           if (j % 2 == 0)
               stiva[k+1][1] = j+1;
           else
               stiva[k+1][1] = j-1;
           break;

    //continuarea sensului pe celalalt nivel orizontal
    case 1: stiva[k+1][2] = 0;
           if (j % 2 == 0) {
               stiva[k+1][1] = j+1;
               if (i % 2 == 0)
                   stiva[k+1][0] = i+1;
               else
                   stiva[k+1][0] = i-1;
           }
           else {
               stiva[k+1][1] = j-1;
               if (i % 2 == 0)
                   stiva[k+1][0] = i+1;
               else
                   stiva[k+1][0] = i-1;
           }
           break;

    //intoarcere
    case 2: stiva[k+1][2] = 0;
           stiva[k+1][1] = j;
           if (i % 2 == 0)
               stiva[k+1][0] = i+1;
           else
               stiva[k+1][0] = i-1;
           break;

    //cont sensului pe urmatorul nivel de comutatoare
    case 3: stiva[k+1][2] = 0;
           if (j % 2 == 0) {
               stiva[k+1][0] = rol(i);
               stiva[k+1][1] = j-1;
           }
           else {
               stiva[k+1][0] = ror(i);
               stiva[k+1][1] = j+1;
           }
           break;

}

stiva[k][2]++;
return 1;
}
else
return 0;

```

```

    }

// verificarea validitatii nodului urmator
int valid (int k) {

    int i;
    int bucle=0;

    //excludere marginilor
    if (stiva[k+1][0] < 0 || stiva[k+1][0] >= D)
        return 0;
    if (stiva[k+1][1] < 0 || stiva[k+1][1] > MARGD)
        return 0;

    //excluderea buclelor inchise
    for(i=1; i<=k; i++)
        if (stiva[i][0]==stiva[k+1][0] && stiva[i][1]==stiva[k+1][1])
            return 0;
    //excluderea repetarii unor segmente la infinit
    if ((stiva[k-1][2]-1==3) && (stiva[k][2]-1==3))
        return 0;
    if ((stiva[k][2]-1!=3) && (stiva[k-1][2]-1!=3))
        return 0;
    //excluderea cailor cu mai mult de o intoarcere
    for (i=1; i<k; i++)
        if (stiva[i][2]-1==2)
            bucle++;
    if (bucle>1)
        return 0;
    return 1;

}

//verificarea ajungerii la solutie
int solutie (int k) {

    if ((stiva[k+1][1]==MARGD) || (stiva[k+1][1]==0))
        return 1;
    return 0;
}

//inregistrarea solutiei
void inregistrare (int k) {

    int i;
    int j;

    j=stiva[k+1][0];
    for(i=0; i<=k; i++)
        if (stiva[i][2]-1==2) {
            if (stiva[1][1]==0) {
                stat[j][1]++;
                save_min(k,1); }
            else {
                stat[j][3]++;
                save_min(k,3); }
        }
    return;
}

```



```

    }
    if (stiva[1][1]==0) {
        stat[j][0]++;
        save_min(k,0); }
    else {
        stat[j][2]++;
        save_min(k,2); }
/*
for(i=1;i<=k;i++)
    printf("%d %d %d\t",stiva[i][0],stiva[i][1],stiva[i][2]-1);
printf("\n");*/
}

//generarea datelor pentru un nod sursa
void genereaza (int s, int m) {

    int k=1;
    int as, ev;

    stiva[0][2]=4;
    init_st(1,s,m);
    while (k >= 0) {

        do {
            as = succ(k);
            if (as)
                ev = valid(k);

        } while ((as) && !(as && ev));

        if (as)
            if (solutie(k))
                inregistrare(k);
            else {
                k++;
            }
        else
            k--;
    }
}

void main (void) {
    int i;
    //selectarea fiecarui nod sursa pt generare date
    for(i=0;i<D;i++) {
        init_stat();
        init_min();
        genereaza(i,0);
        genereaza(i,MARGD);
        printf("Nr cai Sursa: %d\n",i);
        printf("DEST\tRD\tBD\tRI\tBI\n");
        print_stat();
        printf("Nr minim comutatoare Sursa: %d\n",i);
        printf("DEST\tRD\tBD\tRI\tBI\n");
        print_min();
        printf("\n");
    }
}

```


Numărul de căi între o sursă dată și toate destinațiile și lungimile minime ale căilor pentru fiecare tip de rutare într-o RMM 16 X 16

Număr de căi		S=0					Număr minim de comutatoare		S=0					Număr de căi		S=1					Număr minim de comutatoare		S=1											
D	RD	BD	RI	BI	D	RD	BD	RI	BI	D	RD	BD	RI	BI	D	RD	BD	RI	BI	D	RD	BD	RI	BI	D	RD	BD	RI	BI	D	RD	BD	RI	BI
0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	0	4	1	4	1	0	0	0	0	0	1	0	0	0	0					
1	1	1	1	1	1	4	1	4	1	2	1	8	1	2	2	4	7	4	3	2	1	8	1	2	3	4	7	4	3					
2	1	8	1	2	2	4	7	4	3	4	1	4	1	4	4	4	5	4	5	4	1	4	1	4	4	4	5	4	5					
3	1	8	1	2	3	4	7	4	3	5	1	4	1	4	5	4	5	4	5	6	1	8	1	4	6	4	7	4	5					
4	1	4	1	4	4	4	5	4	5	7	1	8	1	4	7	4	7	4	5	8	1	2	1	8	8	4	3	4	7					
5	1	4	1	4	5	4	5	4	5	9	1	2	1	8	9	4	3	4	7	9	1	2	1	8	10	4	7	4	7					
6	1	8	1	4	6	4	7	4	5	10	1	8	1	8	10	4	7	4	7	11	1	8	1	8	11	4	7	4	7					
7	1	8	1	4	7	4	7	4	5	11	1	8	1	8	12	4	5	4	7	12	1	4	1	8	12	4	5	4	7					
8	1	2	1	8	8	4	3	4	7	13	1	4	1	8	13	4	5	4	7	13	1	4	1	8	13	4	5	4	7					
9	1	2	1	8	9	4	3	4	7	14	1	8	1	8	14	4	7	4	7	14	1	8	1	8	14	4	7	4	7					
10	1	8	1	8	10	4	7	4	7	15	1	8	1	8	15	4	7	4	7	15	1	8	1	8	15	4	7	4	7					
11	1	8	1	8	11	4	7	4	7																									
12	1	4	1	8	12	4	5	4	7																									
13	1	4	1	8	13	4	5	4	7																									
14	1	8	1	8	14	4	7	4	7																									
15	1	8	1	8	15	4	7	4	7																									

Număr de căi		S=2					Număr minim de comutatoare		S=2					Număr de căi		S=3					Număr minim de comutatoare		S=3											
D	RD	BD	RI	BI	D	RD	BD	RI	BI	D	RD	BD	RI	BI	D	RD	BD	RI	BI	D	RD	BD	RI	BI	D	RD	BD	RI	BI	D	RD	BD	RI	BI
0	1	8	1	2	0	4	7	4	3	0	1	8	1	2	0	4	7	4	3	0	1	8	1	2	0	4	7	4	3					
1	1	8	1	2	1	4	7	4	3	1	1	8	1	2	1	4	7	4	3	1	1	8	1	2	1	4	7	4	3					
2	0	0	0	0	2	0	0	0	0	2	1	1	1	1	2	4	1	4	1	2	1	1	1	1	2	4	1	4	1					
3	1	1	1	1	3	4	1	4	1	3	0	0	0	0	3	0	0	0	0	3	0	0	0	0	3	0	0	0	0					
4	1	8	1	4	4	4	7	4	5	4	1	8	1	4	4	4	7	4	5	4	1	8	1	4	4	4	7	4	5					
5	1	8	1	4	5	4	7	4	5	5	1	8	1	4	5	4	7	4	5	5	1	8	1	4	5	4	7	4	5					
6	1	4	1	4	6	4	5	4	5	6	1	4	1	4	6	4	5	4	5	6	1	4	1	4	6	4	5	4	5					
7	1	4	1	4	7	4	5	4	5	7	1	4	1	4	7	4	5	4	5	7	1	4	1	4	7	4	5	4	5					
8	1	8	1	8	8	4	7	4	7	8	1	8	1	8	8	4	7	4	7	8	1	8	1	8	8	4	7	4	7					
9	1	8	1	8	9	4	7	4	7	9	1	8	1	8	9	4	7	4	7	9	1	8	1	8	9	4	7	4	7					
10	1	2	1	8	10	4	3	4	7	10	1	2	1	8	10	4	3	4	7	10	1	2	1	8	10	4	3	4	7					
11	1	2	1	8	11	4	3	4	7	11	1	2	1	8	11	4	3	4	7	11	1	2	1	8	11	4	3	4	7					
12	1	8	1	8	12	4	7	4	7	12	1	8	1	8	12	4	7	4	7	12	1	8	1	8	12	4	7	4	7					
13	1	8	1	8	13	4	7	4	7	13	1	8	1	8	13	4	7	4	7	13	1	8	1	8	13	4	7	4	7					
14	1	4	1	8	14	4	5	4	7	14	1	4	1	8	14	4	5	4	7	14	1	4	1	8	14	4	5	4	7					
15	1	4	1	8	15	4	5	4	7	15	1	4	1	8	15	4	5	4	7	15	1	4	1	8	15	4	5	4	7					

Număr de
căi

S=4

D	RD	BD	RI	BI
0	1	4	1	4
1	1	4	1	4
2	1	8	1	4
3	1	8	1	4
4	0	0	0	0
5	1	1	1	1
6	1	8	1	2
7	1	8	1	2
8	1	4	1	8
9	1	4	1	8
10	1	8	1	8
11	1	8	1	8
12	1	2	1	8
13	1	2	1	8
14	1	8	1	8
15	1	8	1	8

Număr minim
de comutatoare

S=4

D	RD	BD	RI	BI
0	4	5	4	5
1	4	5	4	5
2	4	7	4	5
3	4	7	4	5
4	0	0	0	0
5	4	1	4	1
6	4	7	4	3
7	4	7	4	3
8	4	5	4	7
9	4	5	4	7
10	4	7	4	7
11	4	7	4	7
12	4	3	4	7
13	4	3	4	7
14	4	7	4	7
15	4	7	4	7

Număr de
căi

S=5

D	RD	BD	RI	BI
0	1	4	1	4
1	1	4	1	4
2	1	8	1	4
3	1	8	1	4
4	1	1	1	1
5	0	0	0	0
6	1	8	1	2
7	1	8	1	2
8	1	4	1	8
9	1	4	1	8
10	1	8	1	8
11	1	8	1	8
12	1	2	1	8
13	1	2	1	8
14	1	8	1	8
15	1	8	1	8

Număr minim
de comutatoare

S=5

D	RD	BD	RI	BI
0	4	5	4	5
1	4	5	4	5
2	4	7	4	5
3	4	7	4	5
4	4	1	4	1
5	0	0	0	0
6	4	7	4	3
7	4	7	4	3
8	4	5	4	7
9	4	5	4	7
10	4	7	4	7
11	4	7	4	7
12	4	3	4	7
13	4	3	4	7
14	4	7	4	7
15	4	7	4	7

Număr de
căi

S=6

D	RD	BD	RI	BI
0	1	8	1	4
1	1	8	1	4
2	1	4	1	4
3	1	4	1	4
4	1	8	1	2
5	1	8	1	2
6	0	0	0	0
7	1	1	1	1
8	1	8	1	8
9	1	8	1	8
10	1	4	1	8
11	1	4	1	8
12	1	8	1	8
13	1	8	1	8
14	1	2	1	8
15	1	2	1	8

Număr minim
de comutatoare

S=6

D	RD	BD	RI	BI
0	4	7	4	5
1	4	7	4	5
2	4	5	4	5
3	4	5	4	5
4	4	7	4	3
5	4	7	4	3
6	0	0	0	0
7	4	1	4	1
8	4	7	4	7
9	4	7	4	7
10	4	5	4	7
11	4	5	4	7
12	4	7	4	7
13	4	7	4	7
14	4	3	4	7
15	4	3	4	7

Număr de
căi

S=7

D	RD	BD	RI	BI
0	1	8	1	4
1	1	8	1	4
2	1	4	1	4
3	1	4	1	4
4	1	8	1	2
5	1	8	1	2
6	1	1	1	1
7	0	0	0	0
8	1	8	1	8
9	1	8	1	8
10	1	4	1	8
11	1	4	1	8
12	1	8	1	8
13	1	8	1	8
14	1	2	1	8
15	1	2	1	8

Număr minim
de comutatoare

S=7

D	RD	BD	RI	BI
0	4	7	4	5
1	4	7	4	5
2	4	5	4	5
3	4	5	4	5
4	4	7	4	3
5	4	7	4	3
6	4	1	4	1
7	0	0	0	0
8	4	7	4	7
9	4	7	4	7
10	4	5	4	7
11	4	5	4	7
12	4	7	4	7
13	4	7	4	7
14	4	3	4	7
15	4	3	4	7

Număr de
căi**S=8**

D	RD	BD	RI	BI
0	1	2	1	8
1	1	2	1	8
2	1	8	1	8
3	1	8	1	8
4	1	4	1	8
5	1	4	1	8
6	1	8	1	8
7	1	8	1	8
8	0	0	0	0
9	1	1	1	1
10	1	8	1	2
11	1	8	1	2
12	1	4	1	4
13	1	4	1	4
14	1	8	1	4
15	1	8	1	4

Număr minim
de comutatoare**S=8**

D	RD	BD	RI	BI
0	4	3	4	7
1	4	3	4	7
2	4	7	4	7
3	4	7	4	7
4	4	5	4	7
5	4	5	4	7
6	4	7	4	7
7	4	7	4	7
8	0	0	0	0
9	4	1	4	1
10	4	7	4	3
11	4	7	4	3
12	4	5	4	5
13	4	5	4	5
14	4	7	4	5
15	4	7	4	5

Număr de
căi**S=9**

D	RD	BD	RI	BI
0	1	2	1	8
1	1	2	1	8
2	1	8	1	8
3	1	8	1	8
4	1	4	1	8
5	1	4	1	8
6	1	8	1	8
7	1	8	1	8
8	1	1	1	1
9	0	0	0	0
10	1	8	1	2
11	1	8	1	2
12	1	4	1	4
13	1	4	1	4
14	1	8	1	4
15	1	8	1	4

Număr minim
de comutatoare**S=9**

D	RD	BD	RI	BI
0	4	3	4	7
1	4	3	4	7
2	4	7	4	7
3	4	7	4	7
4	4	5	4	7
5	4	5	4	7
6	4	7	4	7
7	4	7	4	7
8	4	1	4	1
9	0	0	0	0
10	4	7	4	3
11	4	7	4	3
12	4	5	4	5
13	4	5	4	5
14	4	7	4	5
15	4	7	4	5

Număr de
căi**S=10**

D	RD	BD	RI	BI
0	1	8	1	8
1	1	8	1	8
2	1	2	1	8
3	1	2	1	8
4	1	8	1	8
5	1	8	1	8
6	1	4	1	8
7	1	4	1	8
8	1	8	1	2
9	1	8	1	2
10	0	0	0	0
11	1	1	1	1
12	1	8	1	4
13	1	8	1	4
14	1	4	1	4
15	1	4	1	4

Număr minim
de comutatoare**S=10**

D	RD	BD	RI	BI
0	4	7	4	7
1	4	7	4	7
2	4	3	4	7
3	4	3	4	7
4	4	7	4	7
5	4	7	4	7
6	4	5	4	7
7	4	5	4	7
8	4	7	4	3
9	4	7	4	3
10	0	0	0	0
11	4	1	4	1
12	4	7	4	5
13	4	7	4	5
14	4	5	4	5
15	4	5	4	5

Număr de
căi**S=11**

D	RD	BD	RI	BI
0	1	8	1	8
1	1	8	1	8
2	1	2	1	8
3	1	2	1	8
4	1	8	1	8
5	1	8	1	8
6	1	4	1	8
7	1	4	1	8
8	1	8	1	2
9	1	8	1	2
10	1	1	1	1
11	0	0	0	0
12	1	8	1	4
13	1	8	1	4
14	1	4	1	4
15	1	4	1	4

Număr minim
de comutatoare**S=11**

D	RD	BD	RI	BI
0	4	7	4	7
1	4	7	4	7
2	4	3	4	7
3	4	3	4	7
4	4	7	4	7
5	4	7	4	7
6	4	5	4	7
7	4	5	4	7
8	4	7	4	3
9	4	7	4	3
10	4	1	4	1
11	0	0	0	0
12	4	7	4	5
13	4	7	4	5
14	4	5	4	5
15	4	5	4	5

Număr de
căi

S=12

D	RD	BD	RI	BI
0	1	4	1	8
1	1	4	1	8
2	1	8	1	8
3	1	8	1	8
4	1	2	1	8
5	1	2	1	8
6	1	8	1	8
7	1	8	1	8
8	1	4	1	4
9	1	4	1	4
10	1	8	1	4
11	1	8	1	4
12	0	0	0	0
13	1	1	1	1
14	1	8	1	2
15	1	8	1	2

Număr minim
de comutatoare

S=12

D	RD	BD	RI	BI
0	4	5	4	7
1	4	5	4	7
2	4	7	4	7
3	4	7	4	7
4	4	3	4	7
5	4	3	4	7
6	4	7	4	7
7	4	7	4	7
8	4	5	4	5
9	4	5	4	5
10	4	7	4	5
11	4	7	4	5
12	0	0	0	0
13	4	1	4	1
14	4	7	4	3
15	4	7	4	3

Număr de
căi

S=13

D	RD	BD	RI	BI
0	1	4	1	8
1	1	4	1	8
2	1	8	1	8
3	1	8	1	8
4	1	2	1	8
5	1	2	1	8
6	1	8	1	8
7	1	8	1	8
8	1	4	1	4
9	1	4	1	4
10	1	8	1	4
11	1	8	1	4
12	1	1	1	1
13	0	0	0	0
14	1	8	1	2
15	1	8	1	2

Număr minim
de comutatoare

S=13

D	RD	BD	RI	BI
0	4	5	4	7
1	4	5	4	7
2	4	7	4	7
3	4	7	4	7
4	4	3	4	7
5	4	3	4	7
6	4	7	4	7
7	4	7	4	7
8	4	5	4	5
9	4	5	4	5
10	4	7	4	5
11	4	7	4	5
12	4	1	4	1
13	0	0	0	0
14	4	7	4	3
15	4	7	4	3

Număr de
căi

S=14

D	RD	BD	RI	BI
0	1	8	1	8
1	1	8	1	8
2	1	4	1	8
3	1	4	1	8
4	1	8	1	8
5	1	8	1	8
6	1	2	1	8
7	1	2	1	8
8	1	8	1	4
9	1	8	1	4
10	1	4	1	4
11	1	4	1	4
12	1	8	1	2
13	1	8	1	2
14	0	0	0	0
15	1	1	1	1

Număr minim
de comutatoare

S=14

D	RD	BD	RI	BI
0	4	7	4	7
1	4	7	4	7
2	4	5	4	7
3	4	5	4	7
4	4	7	4	7
5	4	7	4	7
6	4	3	4	7
7	4	3	4	7
8	4	7	4	5
9	4	7	4	5
10	4	5	4	5
11	4	5	4	5
12	4	7	4	3
13	4	7	4	3
14	0	0	0	0
15	4	1	4	1

Număr de
căi

S=15

D	RD	BD	RI	BI
0	1	8	1	8
1	1	8	1	8
2	1	4	1	8
3	1	4	1	8
4	1	8	1	8
5	1	8	1	8
6	1	2	1	8
7	1	2	1	8
8	1	8	1	4
9	1	8	1	4
10	1	4	1	4
11	1	4	1	4
12	1	8	1	2
13	1	8	1	2
14	1	1	1	1
15	0	0	0	0

Număr minim
de comutatoare

S=15

D	RD	BD	RI	BI
0	4	7	4	7
1	4	7	4	7
2	4	5	4	7
3	4	5	4	7
4	4	7	4	7
5	4	7	4	7
6	4	3	4	7
7	4	3	4	7
8	4	7	4	5
9	4	7	4	5
10	4	5	4	5
11	4	5	4	5
12	4	7	4	3
13	4	7	4	3
14	4	1	4	1
15	0	0	0	0

ANEXA 3

Programul de determinare a căii optime între o sursă și o destinație date în RMM.

```

#define LMAX 255
#include <stdio.h>
#include <math.h>
#include <stdlib.h>

char S[LMAX], D[LMAX], EC[LMAX], ECR[LMAX];

short int NID, NII, l;

short int cai_directe = 1; //se selecteaza cai directe (1) sau cai
inverse (0)

void verifica (char * S) {

    int i;
    for(i=0; i<l; i++)
        if (S[i] != '0' && S[i] != '1') {
            printf("Eroare.\n");
            exit(2);
        }
}

void introducere (void) {

    printf("\n\nIntroduceti lungimea (l) [max %d]: ", LMAX);
    scanf("%d", &l);
    if (l>255) {
        printf("Eroare. l > %d\n", LMAX);
        exit(1);
    }
    gets(S);
    printf("Introduceti S: ");
    gets(S);
    verifica(S);
    printf("Introduceti D: ");
    gets(D);
    verifica(D);
}

int egalitate (char * S, char * D) {

    int i;
    for (i=0; i<l; i++)
        if (S[i] != D[i])
            return 0;
    return 1;
}

char xor (char s, char d) {
    if (s == d)
        return '0';
}

```

```

        else
            return '1';
    }

void calc_EC (void) {
    int i;
    for (i=0; i<l; i++)
        EC[i] = xor (S[i], D[i]);
}

void calc_ECR (void) {
    int tmp=EC[l-1];
    int i;
    for (i=1; i<l; i++)
        ECR[i] = EC[i-1];
    ECR[0] = tmp;
}

void calc_NID (void) {
    int i = l-1;
    while (i >= 0 && (ECR[i] == '0'))
        i--;
    NID=i;
}

void calc_NII (void) {
    int i = 0;
    while (i < l && (EC[i] == '0'))
        i++;
    NII=i;
}

void rotire_st (char * S) {
    char tmp = S[0];
    int i;
    for (i=0; i<l-1;i++)
        S[i] = S[i+1];
    S[l-1] = tmp;
}

void rotire_dr (char * S) {
    char tmp = S[l-1];
    int i;
    for (i=l-1; i>0; i--)
        S[i] = S[i-1];
}

```



```

    S[0] = tmp;
}

void cale_RD (char * S, char * D) {
    int j = 0;
    while (j < l) {
        printf("Nivel %d st: ",j);
        puts(S);
        S[l-1] = D[j];
        printf("Nivel %d dr: ",j);
        puts(S);
        rotire_st (S);
        j++;
    }
}

void cale_RI (char * S, char * D) {
    int j = l-1;
    while (j >= 0) {
        printf("Nivel %d dr: ",j);
        puts(S);
        if (j == 0)
            S[l-1] = D[l-1];
        else
            S[l-1] = D[j-1];
        printf("Nivel %d st: ",j);
        puts(S);
        rotire_dr (S);
        j--;
    }
}

void cale_BD (char * S, char * D) {
    int j = 0;
    while (j <= NID) {
        printf("Nivel %d st: ",j);
        puts(S);
        if (j == NID) {
            if (j==0)
                S[l-1] = D[l-1];
            else
                S[l-1] = D[NID-1];
            printf("Nivel %d st: ", j);
            puts(S);
            rotire_dr (S);
            j--;
            break;
        }
        S[l-1] = D[j];
        printf("Nivel %d dr: ",j);
        puts(S);
        rotire_st (S);
    }
}

```

```

        j++;
    }

    while (j >= 0) {
        printf("Nivel %d dr: ",j);
        puts(S);
        if (j == 0)
            S[l-1] = D[l-1];
        else
            S[l-1] = D[j-1];
        printf("Nivel %d st: ",j);
        puts(S);
        rotire_dr (S);
        j--;
    }
}

void cale_BI (char * S, char * D) {

    int j = l-1;
    while (j >= NII) {
        printf("Nivel %d dr: ",j);
        puts(S);
        if (j == NII) {
            S[l-1] = D[NII];
            printf("Nivel %d dr: ",j);
            puts(S);
            rotire_st (S);
            j++;
            break;
        }
        if (j == 0)
            S[l-1] = D[l-1];
        else
            S[l-1] = D[j-1];
        printf("Nivel %d st: ",j);
        puts(S);
        rotire_dr (S);
        j--;
    }

    while (j < l) {
        printf("Nivel %d st: ",j);
        puts(S);
        S[l-1] = D[j];
        printf("Nivel %d dr: ",j);
        puts(S);
        rotire_st (S);
        j++;
    }
}

void main (void) {

    short int dl,du;

```

```

introducere();

calc_EC();
calc_ECR();
dl = floor ((double)l/2.);
du = ceil ((double)l/2.);

printf("S: "); puts(S);
printf("D: "); puts(D);
printf("EC: "); puts(EC);
printf("ECR: "); puts(ECR);

if (egalitate (S, D)) {
    printf("Cerere de acces pentru memoria locala\n");
    exit(0);
}
else {
    calc_NID();
    calc_NII();
}

printf("NID=%d NII=%d\n", NID,NII);

if ((NID == (l-l-NII))==0) {

    if (cai_directe) {
        printf ("Directionare BD:\n");
        cale_BD (S, D);
    }
    else {
        printf ("Directionare BI:\n");
        cale_BI (S, D);
    }
}
else if (NID < dl) {
    printf ("Directionare BD:\n");
    cale_BD (S, D);
}
else if (NID >= du) {
    printf ("Directionare BI:\n");
    cale_BI (S, D);
}
else {
    if (cai_directe) {
        printf ("Directionare RD:\n");
        cale_RD (S, D);
    }
    else {
        printf ("Directionare RI:\n");
        cale_RI (S, D);
    }
}
}

```


B I B L I O G R A F I E

- [AAS87] G.B. Adams, D.P. Agrawal and H.J. Siegel, "A Survey and Comparison of Fault-Tolerant Multistage Interconnection Networks", *Computer*, vol.20 (June 1987), pp.14 - 27
- [ADT91] M.S.Algudady, C.R.Das and M.J. Thazhuthaveetil, "A Write Update Cache Coherence Protocol for MIN -Based Multiprocessors with Accessibility -Based Split Caches", *IEEE Trans On Computers*, pp.544-552, July 1990
- [AG89] G.S. Almasi and A. Gottlieb, *Highly Parallel Computing*, Redwood City, CA, Benjamin/Cummings, 1989
- [Agra88] P. Agrawal, "Fault Tolerance in Multiprocessor Systems Without Dedicated Redundancy", *IEEE Transactions on Computers*, vol.37 (3), pp. 358 – 362, 1988
- [AFM90] R. Ahmed, R. Frazier and P. Marinos, "CARER: Cache-Aided Rollback Error Recovery Algorithm for Shared Memory Multiprocessors", *Proc. 20th Int. Symp. Fault-Tolerant Computing*, pp.82 – 88, June 1990
- [Arch88] J. K. Archibald, "A Cache Coherence Approach for Large Multiprocessor Systems", *ACM International Conference on Supercomputing*, pp.337 – 345, 1988
- [Ban90] P. Banerjee. "Strategies for Reconfiguring Hypercubes under Faults", *Proc. 20th Int. Conf. Fault-Tolerant Computing Systems*, pp.210 – 215, June 1990
- [BD91] E. Brewer, C. Dellarocas, "PROTEUS: A High Performance Parallel Architecture Simulator", TR516 MIT/LCS, Mass. Inst. of Technology, Sept.1991
- [BIK00] L. Bhuyan, R.Iyer and A. Kumar, "Impact of CC-NUMA Memory Management Policies on the Application Performance of Multistage

- Switching Networks ", *IEEE Trans on Parallel and Distributed Systems*, 13(3), pp.230-246, March 2000
- [BINK97] L.Bhuyan, R. Iyer, A.Nanda and M. Kumar, "Performance of Multistage Bus Networks for a Distributed Shared Memory Multiprocessor ", *IEEE Trans. On Parallel and Distributed Systems*, 8(1),pp. 82-95, January 1997
- [BLA89] L.Bhuyan, B.Liu, and L.Ahmed "Analysis of MIN Based Multiprocessors with Private Cache Memories ", *Proc 1989 International Conf. On Parallel Processing* ", vol I. pp. 51-58, August 1989
- [BN91] L.Bhuyan and A. Nanda,"Multistage Bus Network (MBN) An Interconnection Network for Cache Coherent Multiprocessors ",*Proc Third IEEE Symp. Parallel and Distributed Processing*, Dec1991.
- [BNA94] L.Bhuyan,A.Nanda and T.Askar,"Performance and Reliability of the MBN",*Proc.1994 Int'l Conf. Parallel Processing*,pp. 26-33
- [BNT94] L. Bhuyan. A. Nanda and T. Askar, "Performance and Reliability of MBN". *Proc. 1994 Int'l Conf. Parallel Processing*. pp. 26-33
- [BP92] N. Bowen and D. Pradhan, "Virtual Checkpoints: Architecture and performance", *IEEE Transactions on computers*, vol.41, pp. 516 – 525, May 1992
- [BP93] N. Bowen and D. Pradhan. "Processor and Memory Based Checkpoint and Rollback Recovery", *Computer*, pp. 22 – 31, Feb. 1993
- [BRA89] C. Botting, S. Rai and D. Agrawal, "Reliability Computation of MINs", *IEEE Trans. on Reliability*, vol.38 (1), April 1989
- [BRL94] A. Bhattacharya, R. Rao and T. Lin, "Cumulative Performance Measure for Gracefully Degradable MINs", *Proc. 1st International Workshop on Parallel Processing*, pp. 234 – 239, December 1994
- [BT88] I.Blake and K.Trivedi, "Reliabilities of Two Fault - Tolerant INs ",*Proc.18th Intern. Symp. Fault - Tolerant Computing*. pp.300-305, June 1988

- [BT89] I.Blake and K.Trivedi,"Reliability Analysis of INs Using Hierarchical Composition ", *IEEE Trans. On Reliability*, 38(1), pp.111-119, April 1989
- [BT93] N. Benitez and K. Trivedi. "Multiprocessor Performability Analysis ", *IEEE Trans. On Reliability*, 42(4), pp. 579-587, Dec 1993
- [BWK98] L. Bhuyan, H.Wang and A.Kumar, "Impact of Switch Design on the Application Performance of Cache-Coherent Multiprocessors ", *International Parallel Processing Symposium*, pp.466-474, March 1998
- [BYA89] L.N. Bhuyan, Qing Yang, D.P. Agrawal, "Performance of Multiprocessor Interconnection Networks", *Computer*, vol.22 (February 1989). pp.25 – 37
- [CBZ95] J. Carter, J. Bennett and W. Zwaenepoel, "Techniques for Reducing Consistency-Related Information in Distributed Shared Memory Systems", *ACM Transaction on Computer Systems*, vol.13 (3). pp. 205 – 243, Aug. 1995
- [CF78] L.M.Censier and P.Feautrier, "A new solution to coherence problems in multicache systems", *IEEE Transactions on Computers*, C-27 (December 1978), pp.1112 – 1118
- [CFKA90] D. Chaiken. C. Fields, K.Kurihara and A. Agarwal, "Directory-Based Cache Coherence in Large Scale Multiprocessors", *IEEE Computer*. vol.23 (6), pp.49 – 58, June 1990
- [Choi96] L. Choi, *Compiler and Hardware Support for Cache Coherence in Large Scale Multiprocessors*, PhD Dissertation, University of Illinois at Urbana, Aug. 1996
- [CK92] F. Chong and T. Knight, jr., "Design and Performance of Multipath MIN Architectures", *Proc. 4th Annual ACM Symposium on Parallel Algorithms and Architectures*, pp.286 – 295, 1992
- [Coif94] V. Coifan. "Stadiul actual al sistemelor de calcul multiprocesor". referat la teza de doctorat, Universitatea Politehnica Timișoara, Fac. de Automatică și Calculatoare. Timișoara 1994

- [Coif01] V. Coifan, "Optimizarea performanțelor la sistemele multiprocesor în funcție de dependabilitate", *va apare în Editura Orizonturi Universitare*
- [CSG99] D. Culler, J. Singh and A. Gupta, *Parallel Computer Architecture: A Hardware/Software Approach*, San Francisco, Morgan Kaufmann Publishers, 1999
- [CWH01] V. Coifan, R. Wagmann and N. Hanigovszki, "A Routing Algorithm for Multistage Bus Networks for a Distributed Shared Memory Multiprocesor", *Buletinul Științific al Universității Politehnica*, Timișoara, 2001
- [DA93] T.Dahlberg and D. Agrawal, "Dependability Analysis for Large Systems : a Hierarchical Modeling Aproach ", *IEEE Trans on Reliability*, pp.450-459
- [Dal92] W. J. Dally. "Virtual-Channel Flow Control", *IEEE Transactions on Parallel and Distributed Systems*, vol. 3 (2), pp.194 – 205, Mar. 1992
- [Dasg89] S. Dasgupta. "*Computer Arhitecture - A Modern Synthesis*", vol.2, John Wiley and sons, 1989
- [DB94] J. Ding and L. Bhuyan, "Finite Buffers Analysis of MINs", *IEEE Transactions on Computers*, vol. 43(2), pp.243 – 247, Feb. 1994
- [DKT90] C.R.Das. J.T.Kreulen and M.J.Thazhuthaveetil. "Dependability Modeling for Multiprocessors ", *Computer*, October 1990, pp. 8-19
- [DSB88] M.Dubois, C.Scheurich and F.Briggs, "Synchronization Coherence and Event Ordering in Multiprocessors", *Computer*. February 1988, pp.9 – 21
- [DTB90] C.R. Das,L. Tien and L.N. Bhuyan, "Availability Evaluation of MIN - Connected Multiprocessors using Decomposition Technique ",
- [DTB93] C.R.Das. L.Tien and L.N. Bhuyan, "An Availability Model for MIN - Based Multiprocessors", *IEEE Transactions on Parallel and Distributed Systems*, vol.4, October 1993, pp. 1118-1129
- [FP90] E. Fujiwara and D. Pradhan, "Error Control Coding in Computers", *Computer*, pp. 63 – 72, July 1990

- [GGH91] K.Gharachorloo, A.Gupta and J. Hennessy, "Performance Evaluation of Memory Consistency Models for Shared Memory Multiprocessors", *Int. Conf. on Architectural Support for Programming Languages and Operating Systems*, pp. 245 – 257, 1991
- [GHGM91] A.Gupta, J.Hennessy, K.Gharachorloo and T.Mowry, "Comparative Evaluation of Latency Reducing and Tolerating Techniques", *Int. Symp. on Computer Architecture*, pp. 254 – 263, 1991
- [GLL+90] K. Gharachorloo, D. Lenoski, J. Laudon, P. Gibbons, A. Gupta et al., "Memory Consistency and Event Ordering in Scalable Shared Memory Multiprocessors", *Proceedings of the International Symposium on Computer Architecture*, pp.15 – 26, 1990
- [Good83] J. R. Goodman, "Using Cache Memory to Reduce Processor-Memory Traffic", *International Symposium on Computer Architecture*, pp. 124 – 131, 1983
- [Gros99] D. Grosu, "Rollback Recovery Techniques in Message Passing Systems : A Survey and Synthesis " Technical Report, Dept. of Computer Engineering, Technical University of Timisoara, June 1999
- [GVW89] J. Goodman, M. Vernon and P. Woest, "A Set of Efficient Synchronisation Primitives for a Large Scale Shared Memory Multiprocessors". *Int. Conf. on Architectural Support for Programming Languages and Operating Systems*, pp.64 – 73, 1989
- [GW88] J. Goodman and P. Woest, "The Wisconsin Multicube: A New Large Scale Cache Coherent Multiprocessor", *International Symposium on Computer Architecture*, pp. 422 – 431, 1988
- [GY93] P. Gawghan and S. Yalamanchi, "Adaptive Routing Protocols for Hypercube Interconnection Networks", *IEE Transactions on Computers*, vol. 26 (5). pp. 12 – 23, May 1993

- [HM87] D. Hunt and P. Marinos, "A General Purpose Cache-Aided Rollback Error Recovery (CARER) Tehnique", *17th Symp. On Fault Tolerant Computing*, pp. 170 – 175, June 1987
- [HP94] J.L. Hennessy and D.A. Patterson, "*Computer Organization and Design*", Morgan Kaufmann. San Francisco, CA. 1994
- [HP96] J.L. Hennessy and D.A. Patterson. "*Computer Organization - A Quantitative Aproach*", Morgan Kaufmann, San Francisco, CA, 1996
- [Hwan93] K. Hwang, "*Advanced Computer Architecture : Parallelism, Scalability, Programmability*", Mc Graw-Hill Inc., 1993
- [IB99] R. Iyer and L. Bhuyan, "Switch Cache: A Framework for Improving the Remote Memory Access Latency of CC-NUMA Multiprocessors", *Fifth International Conference on High Performance Computer Arhitecture (HPCA -5)*. pp.152 – 160, Jan. 1999
- [IB00] R. Iyer and L.Bhuyan, "Design and Evaluation of a Switch Cache Architecture for CC-NUMA Multiprocessors, *IEEE Trans . on Computers*. 49(8), pp. 779-797. August 2000
- [Iyer96] R. Iyer, "Distributed Shared Memory Multiprocessors Using Multistage Bus Networks", MS Thesis, Texas A&M Univ. Dept. of Comp. Science, Aug. 1996
- [Iyer99] R.Iyer, "High Performance Switch Architectures for CC-NUMA Multiprocessors", PhD Dissertation, Texas A&M.Univ. Dept. of Comp. Science, Aug. 1999
- [JF91] B. Janssens and W. Fuchs, "Experimental Evaluation of Multiprocessor Cache-Based Error Recovery", *1991 International Conference on Parallel Processing*. pp. 505 – 508, Aug. 1991
- [Kak96] S.H. Kakas. "*Contribuții la analiza și sinteza de fiabilitate a sistemelor multiprocesor*". teză de doctorat. Universitatea Politehnica din Timișoara. Fac. Automatică și Calculatoare. Timișoara. 1996

- [KB96] A. Kumar and L. Bhuyan, "Evaluating Virtual Channels for Cache Coherent Shared Memory Multiprocessors", *ACM International Conference on Supercomputing*, pp. 253 – 260, May 1996
- [KBI96] A. Kumar, L. Bhuyan and R. Iyer, "Execution Based Evaluation of MINs for Cache Coherent Multiprocessor", *International Conference on Advanced Computing (ADCOMP '96)*, India, May 1996
- [KR87] V.P. Kumar and S.M. Reddy, "Augmented Shuffle - Exchange Multistage Interconnection Networks". *Computer*, vol.20 (June 1987), pp.30 – 40
- [KS83] C. Kruskal and M. Snir, "The Performance of MINs for Multiprocessors", *IEEE Trans. on Computers*, vol.32 (12), pp.1091-1098, Dec.1983
- [Kum96] A. Kumar, *Execution-Driven Evaluation of Cache Coherent Shared Memory Multiprocessors*, PhD Dissertation, Dept. of Computer Science, Texas A&M University, Aug. 1996
- [Lamp79] L. Lamport, "How to Make a Multiprocessor Computer that Correctly Executes Multiprocess Programs", *IEEE Transactions on Computers*, vol. C-28 (9), pp. 690 – 691, September 1979
- [LFA90] J. Long, W. Fuchs and J. Abraham, "A Forward Recovery Using Checkpointing in Parallel Systems", *Proceedings of the 1990 International Conference on Parallel Processing*, pp.272 – 275, Aug. 1990
- [LFA91] J. Long, W. Fuchs and J. Abraham, "Implementing Forward Recovery Using Checkpoints in Distributed Systems", *IFIP 2nd Int. Working Conf. Dependable Computing for Critical Applications*, Feb.1991
- [Lil93] D. J. Lilja, "Cache Coherence in Large Scale Shared Memory Multiprocessors", *ACM Computing Surveys*, vol. 25 (3), pp. 303 – 338. Sept. 1993
- [LSP82] L.Lamport, R.Shoshtak and M.Pease, "The Byzantine General's Problem". *ACM Trans. Prog. Lang. Syst.*, July 1982, pp.382 – 401

- [Marq89] D. Marquardt, "C2MP: A Cache Coherent Distributed Memory Multiprocessor System", *Proceedings Supercomputing '89*, pp.466 – 475, 1989
- [MBL89] H.E.Mizrahi, J.L.Baer, E.D. Lazowska, "*Introducing Memory into the Switch Elements of Multiprocessor Interconnection Networks*", University of Washington, Dept. of Comp. Sci, Technical Report 88-11-03, Nov 1988
- [MD96] P.Mohapatra and C.Das, "Performance Analysis of Finite - Buffered Asynchronous MINs", *IEEE Trans. On Parallel and Distributed Systems*, 7(1), pp.18-25, January 1996.
- [MD98] A. Moga and M. Dubois, "The Effectiveness of SRAM Network Caches on Clustered DSMs", *Fourth International Symposium on High Performance Computer Architecture*, pp. 103 – 112, Feb. 1998
- [MHB86] T.N. Mudge. I.P. Hayes, G.D. Buzzard, "Analysis of Multiple-Bus Interconnection Networks". *Journal of Parallel and Distributed Computing*. (1986), pp.328 – 343
- [MHW87] T.N. Mudge. J.P. Hayes, D.C. Winsor, "Multiple-Bus Architectures", *Computer*, vol. 20 (June 1987), pp.42 – 48
- [MN99] M. Michael and A. Nanda. "Design and Performance of Directory Caches for Scalable Shared Memory Multiprocessors", *Fifth International Conference on High Performance Computer Architecture (HPCA – 5)*, pp.142 – 151, 1999
- [MNLS97] M. Michael, A. Nanda, B. Lim and M. Scott. "Coherence Controller Architectures for Shared Memory Multiprocessors", *Proceedings of the International Symposium on Computer Architecture (ISCA '97)*, pp. 219 – 228, 1997
- [Mold93] D.I. Moldovan, "*Parallel Processing : From Applications to Systems*", Morgan Kaufman, 1993

- [MP89] J. Meyer and D. Pradhan, "Dynamic Testing Strategy for Distributed Systems", *IEEE Transactions on Computers*, vol 38(3), pp. 356 – 365, March 1989
- [NB93] A.K.Nanda and L.N.Bhuyan, "Design and Analysis of Cache Coherent Multistage Interconnection Networks ",*IEEE Trans. On Computers*, vol.42,no.4, April 1993
- [Ng91] T. P. Ng. *Checkpointing Algorithms for Shared Memory Systems*. Technical Report, Dept. of Computer Science, University of Illinois, Urbana, 1991
- [NK93] L. Ni and P. McKinley, "A Survey of Wormhole Routing Tehniques in Direct Netoworks", *Computer*, Vol. 26 (2), pp. 62 – 76, Feb. 1993
- [NOS96] B. Nayfeh, C. Olukotun and J. Singh, "The Impact of Shared-Cache Clustering in Small-Scale Shared memory Multiprocessors", *Proceedings of the 2nd International Symposium on High Performance Computer Arhitecture*, pp.74 – 84, Feb. 1996
- [PA96] R. Potlapalli and D. Agrawal, "On the Design of Efficient MINs", *IEEE Proc. of MASCOTS '96*
- [Pat81] J. Patel, "Performance of Processor-Memory Interconnection for Multiprocessors", *IEEE Transactions on Computers*, vol. C-30 (10). pp. 771 – 780, Oct. 1991
- [PF90] D. Pradhan and E. Fujiwara, "Error Computing Codes in Fault Tolerant Computers". *Computer*, vol.23 (7), pp. 62 – 72, July 1990
- [Prad89] D. K. Pradhan, *Redundancy Schemes for Recovery*. TR-89-CSE-16. ECE Department. University of Massachusets, 1989
- [Prad96] D.K. Pradhan. "*Fault-tolerant Computer System Design*", Prentice Hall PTR, New Jersey. 1996
- [Popa96] M.Popa. "Intruducere în arhitecturi paralele și neconvenționale". *Editura AS Computer Press*, Timișoara 1992

- [PV94a] D. Pradhan and N. Vaydia, "Roll Forward Checkpointing Scheme: A Novel Fault-Tolerant Architecture", *IEEE Transactions on Computers*, vol. 43 (10), pp. 1163 – 1174, October 1994
- [PV94] D. Pradhan and N. Vaydia, "Performance and Reliability Trade-offs in Roll Forward Schemes", *Proc. FTCS-24*, Austin Texas, June 1994
- [RA90] S. Rai and D. Agrawal, *Distributed Computing Network Reliability*, Tutorial Textbook, IEEE CS Press, 1990
- [RG87] D.A. Reed and D.C. Grunwald, "The Performance of Multicomputer Interconnection Networks", *Computer*, vol. 20 (June 1987), pp.63 – 73
- [SD89] C. Scheurich and M. Dubois, "Dynamic Page Migration in Multiprocessors with Distributed Global Memory", *IEEE Transactions on Computers*, vol. 38 (8), pp. 1154 – 1163, Aug. 1989
- [Sita93] R. Sitamaran, "Communication and Fault Tolerance in Parallel Computers", PhD Dissertation CS-TR-414-93, Dept. of Comp.Science, Princeton Univ., Princeton NJ, 1993
- [SLM89] Sang Lyul Min, "*Memory Hierarchy Management Schemes in Large Scale Shared - memory Multiprocessor*", Technical Report # 89-08-07, Dept. of Computer Science and Engineering, FR-35. University of Washington, Seattle, WA 98195 USA, August 1989
- [Smit82] A. J. Smith, "Cache Memories", *ACM Computing Surveys*, vol. 14 (3), pp. 473 – 530, 1982
- [TD89] L. Tien and C. Das, "Reliability Evaluation of Butterfly Network-Based Multiprocessors", TR – 89 – 072, The Pennsylvania State University
- [TH90] J.Torrellas and J.Hennessy, "Estimating the Performance Advantages of Relaxing Consistency in a Shared Memory Multiprocessor". *Int. Conf. on Parallel Processing*, vol. I, pp.26 – 33, 1990
- [TZ94] J. Torrellas and Z. Zhang, "The Performance of the Cedar Multistage Switching Network", *IEEE Transactions on Parallel and Distributed Systems*, vol. 8 (4), pp. 321 – 336, Apr. 1994

- [Vaid86] A. V. Veidenbaum, "A Compiler-Assisted Cache Coherence Solution for Multiprocessors", *International Conference on Parallel Processing*, pp. 1029 – 1036, 1986
- [VP93] N. Vaidya and D. Pradhan, "Fault Tolerant Design Strategies for High Reliability and Safety", *IEEE Transactions on Computers*, vol. 43 (10), Oct. 1993
- [VSD97] A. Vaidya, A. Sivasubramaniam and C. Das, "Performance Benefits of Virtual Channels and Adaptive Routing: An Application-Driven Study", *Proceedings of the 11th International Conference on Supercomputing*, pp. 140 – 147, Jul. 1997
- [WFP90] K.Wu,K.Fuchs and I.Patel, "Error Recovery in Shared Memory Multiprocessors Using Private Caches ", *IEEE Trans. On Parallel and Distributed Systems*, 1(2), pp.231-240. April 1990
- [Wils87] A. W. Wilson, "Hierarchical Cache/Bus Architecture for Shared Memory Multiprocessors", *Proceedings of the 14th Annual International Symposium on Computer Architecture*, pp. 244 – 252, 1987
- [Witt81] L.D. Wittie, "Communication Structures for Large Networks of Microcomputers", *IEEE Transactions on Computers*, vol. 30 NO. 4 apr. 1981, pp.246 – 273
- [YB91] Q. Yang and L. Bhuyan, "Analysis of Packet - Switched Multiple-Bus Multiprocessors ", *IEEE Trans. on Computers*, 40(3) pp.352-357, March 1991
- [YLL90] H. Yoon, K. Lee and M. Liu, "Performance Analysis of Multibuffered Packet-Switching Networks in Multiprocessor Systems", *IEEE Transactions on Computers*, vol. 39 (3), pp. 319 – 327, Mar. 1990
- [ZB92] R.Zucker and J.-L.Baer, "A Performance Study of Memory Consistency Models", *Int. Symp. on Computer Architecture*. pp.2 – 12, 1992

LISTĂ SELECTIVĂ

A LUCRĂRILOR ȘTIINȚIFICE PUBLICATE, A INVENȚIILOR ȘI INOVAȚIILOR ELABORATE ȘI A CERCETĂRILOR CONTRACTUALE EFECTUATE

A. Lucrări științifice publicate

- [Coif75] V. Coifan. "Bazele logice ale calculatoarelor", *Îndrumător de lucrări*, IPTVT 1975
- [RPC77] Al. Rogojan, V. Pop și V. Coifan, "Microcalculator didactic MCD1", *Buletinul Științific și Tehnic al IPTVT*, tom 22, iulie – decembrie 1977, pp.383 – 388
- [Coif81] V. Coifan, "Microcalculator didactic MCD85", *Primul simpozion de microprocesoare, microcalculatoare și aplicații*, noiembrie 1981, pp. 68 – 76
- [Coif84] V. Coifan, "Sistem ierarhizat de microcalculatoare construit în jurul familiei de calculatoare M18", *Conferința națională a cercurilor științifice studențești*. Timișoara, noiembrie 1984
- [CC84] V. Crețu și V. Coifan, "Executiv în timp real implementat pe un sistem de calcul ECAROM 881", *Conferința națională a cercurilor științifice studențești*, Timișoara, noiembrie 1984
- [Coif85] V. Coifan. "An arithmetic processing unit APU – for the ECAROM 881 microcomputer", *Simpozionul de microprocesoare, microcalculatoare și aplicații în economie*, noiembrie 1985, pp.67 – 68
- [CP85] V. Coifan și O. Proștean, "A HD 64180 μ P-based SBC", *Simpozionul de microprocesoare, microcalculatoare și aplicații în economie*, noiembrie 1985, pp.67 – 68
- [CPDL+85] V. Coifan, O. Proștean, T. Dragomir ș.a., "Some aspects of automatic voltage regulation at hidrogenerator using a process computer", *Hydro-Turbo*, 1985. Brno, Cehoslovacia
- [BC85] N. Budișan și V. Coifan, "Studiul posibilităților de teleinvestigare cu ajutorul calculatorului a agregatului aeroelectric de 300 KVA de pe muntele Semenic", *IPA*, mai 1985

- [CMP86] V. Coifan, I. C. Moș și M. Popa, “Structura sistemelor de prelucrare a datelor numerice”, *Îndrumător de lucrări*, IPTVT 1986
- [CP88] V. Coifan și O. Proștean, “Calculator personal cu facilități de conducere a proceselor realizat în jurul lui μ P HD 64180”, *Simpoziomul național de calculatoare și conducere automată a proceselor*, Timișoara 1988, pp. 19 – 22

B. Rapoarte la contracte de cercetare

- [CB87] V. Coifan și R. Boraci, “Cercetări privind realizarea unui regulator de tensiune destinat hidrogenatoarelor reversibile”, *Contract de cercetare nr. 126/1987 cu CCSITEH Reșița*
- [CB88] V. Coifan și R. Boraci, “Experimentarea și analiza regulatorului de tensiune într-o centrală hidroelectrică”, *Contract de cercetare nr. 162/1988 cu CCSITEH Reșița*
- [CP84] V. Coifan și I. Păpușoi, “Comanda cu microcalculatorul ECAROM 881 a regimului de propulsie la vehiculul MAGNIBUS 01”, *Contract de cercetare nr. 128/1984 cu CCSIT Electroputere Craiova*
- [Coif87] V. Coifan, “Posibilități de utilizare într-un mediu multiprocesor a modulelor din familia MULTIPROM”, *Contract de cercetare nr. 7405/1987 cu IPA București*

C. Invenții și inovații

- *Procedeu de extindere a posibilităților de reglare automată bipozițională a echipamentelor de calcul numerice* – brevet de invenție
- *Metodă și dispozitiv electronic pentru testarea poziției rotorului și conducerea unei mașini sincrone cu ajutorul microprocesorului* – brevet de invenție
- *Metodă de testare a întreruptoarelor din microcalculatorul ECAROM 881* – certificat de inovator, nr. 133/1985
- *Schemă logică pentru protecția memoriei partajate* – certificat de inovator nr. 200/1985

- *Modul de memorie pentru sisteme de calcul multiprocesor* – certificat de inovator nr. 489/1986
- *Comenzi de monitor specifice arhitecturilor de calcul multiprocesor* – certificat de inovator nr. 419/1986
- *Panou de testare a modulelor calculator monoplacă pe 8 bit*, certificat de inovator nr. 420/1986
- *Automat de acces la magistrala comună a unui multiprocesor pe 8 bit*, certificat de inovator nr. 423/1986
- *Logică pentru evitarea blocării magistralei comune în sistemele de calcul multiprocesor*, certificat de inovator nr. 487/1996
- *Metodă și schemă de configurare dinamică a sensului de transfer a datelor la porturile paralele programabile*, certificat de inovator nr. 1072/1989