

UNIVERSITATEA POLITEHNICA TIMIȘOARA
Facultatea de Automatică și Calculatoare

UNIVERSITATEA "POLITEHNICA"
TIMIȘOARA
BIBLIOTECA CENTRALĂ

r. Inv. 623. 720

ulap 181 Lit. U

TEZĂ DE DOCTORAT
Creșterea fiabilității la transmiterea și stocarea
informației

Conducător științific
Prof.dr.ing.Mircea Vlăduțiu

Doctorand
Ing. Cezar Morun

1998

623.746
181 c

CUPRINS

| | |
|---|-----------|
| CUPRINS..... | 1 |
| 0 INTRODUCERE..... | 5 |
| 0.1 Toleranța la defectare ca scop al proiectării..... | 5 |
| 0.2 Toleranța la defecte în subsistemul I/O..... | 7 |
| 0.3 Structura tezei de doctorat..... | 9 |
| 1 METODE DE CREȘTERE A FIABILITĂȚII ȘI DISPONIBILITĂȚII | 13 |
| 1.1 Indicatori pentru caracterizarea fiabilității și disponibilității | 13 |
| 1.1.1 Clasificarea defectiunilor | 13 |
| 1.1.2 Estimarea fiabilității..... | 15 |
| 1.1.3 Disponibilitatea..... | 16 |
| 1.1.4 Dependabilitatea | 17 |
| 1.1.5 Mentenabilitatea..... | 17 |
| 1.1.6 Cerințele fiabilității aplicate..... | 18 |
| 1.1.7 Tehnici de redundanță..... | 19 |
| 1.2 Metode de autocontrol, implementate <i>hardware</i> | 20 |
| 1.2.1 Controlul prin copii..... | 22 |
| 1.2.2 Controale de codificare..... | 22 |
| 1.2.3 Controale de temporizare..... | 23 |
| 1.2.4 Controale de excepție..... | 23 |
| 1.2.5 Restabilirea erorii..... | 24 |
| 1.2.6 Clasificarea procedurilor de restabilire | 24 |
| 1.2.6.1 Restabilirea totală..... | 24 |
| 1.2.6.2 Reconfigurarea..... | 26 |
| 1.3 Metode de autocontrol implementate <i>software</i> | 26 |
| 1.3.1 Funcția unui proces | 27 |
| 1.3.2 Secvența de comandă a unui proces..... | 27 |
| 1.3.3 Date de proces..... | 27 |
| 1.3.4 Restabilirea <i>software</i> | 27 |
| 1.3.5 Diagnosticarea erorii..... | 28 |
| 1.3.6 Autocontrolul ca mijloc de creștere a fiabilității și disponibilității sistemelor de calcul | 28 |
| 1.3.6.1 Metode de autocontrol aplicate la diferite niveluri de structură..... | 28 |
| 1.3.6.1.1 Caracteristicile metodelor de proiectare structurată la nivel de bistabil..... | 28 |
| 1.3.6.1.2 Caracteristicile metodelor de proiectare structurată la nivel de registru | 31 |
| 1.3.6.1.3 Proiectarea structurată la nivel de bloc | 32 |
| 1.3.6.1.3.1 Autocontrolul bazat pe principiul BILBO..... | 33 |
| 1.3.6.1.3.2 Modulul comutator de testare, încorporat, pentru izolarea defectelor..... | 33 |
| 1.3.6.1.4 Controlul erorilor la nivel de procesor..... | 34 |
| 1.3.6.2 Problematika construcției <i>checker</i> -elor de erori..... | 35 |
| 1.3.6.2.1 Circuite de verificare dublă cale | 35 |
| 1.3.6.2.1.1 <i>Checker</i> de egalitate..... | 35 |
| 1.3.6.2.1.2 <i>Checker</i> M din N..... | 36 |
| 1.3.6.2.2 Circuite cu autoverificare totală, încorporate..... | 36 |
| 1.4 Concluzii..... | 40 |

| | | |
|------------|--|-----------|
| 2 | ANALIZA SUBSISTEMULUI DE DISCURI DIN PUNCT DE VEDERE A FIABILITĂȚII | 41 |
| 2.1 | Coduri pentru discuri | 41 |
| 2.1.1 | Codul <i>Fire</i> pentru memoriile bazate pe disc magnetic | 42 |
| 2.1.2 | Codul întrețesut R-S SbEC-DbEC pentru memoriile cu discuri magnetice | 47 |
| 2.1.3 | Codul R-S întrețesut-încrucișat pentru discuri optice | 49 |
| 2.2 | Matrici de discuri independente și redundante | 53 |
| 2.2.1 | Organizarea datelor | 56 |
| 2.2.1.1 | Distribuția datelor | 56 |
| 2.2.1.2 | Adresarea independentă | 57 |
| 2.2.1.2.1 | Divizarea fină | 57 |
| 2.2.1.2.2 | Divizarea grosieră | 58 |
| 2.2.1.2.3 | Alegerea mărimii unității de divizare | 59 |
| 2.2.2 | Mecanismele de redundanță | 60 |
| 2.2.2.1 | Replicarea datelor | 61 |
| 2.2.2.2 | Protecția bazată pe paritate | 63 |
| 2.2.2.3 | Alte scheme | 65 |
| 2.2.3 | Performanța versus fiabilitatea | 66 |
| 2.2.3.1 | Negocierea fiabilității în favoarea performanței | 66 |
| 2.2.3.2 | Îmbunătățirea performanței și a fiabilității | 67 |
| 2.2.4 | Tratarea căderilor | 68 |
| 2.3 | Clasificarea subsistemelor de discuri redundante din punct de vedere a fiabilității.. | 69 |
| 2.4 | Concluzii | 71 |
| 3 | O ANALIZĂ A ARHITECTURILOR RAID | 73 |
| 3.1 | RAID level 0 | 73 |
| 3.1.1 | Distribuția datelor în RAID 0 | 73 |
| 3.1.2 | Comportarea RAID 0 în aplicații ce solicită rate ridicate de transfer | 73 |
| 3.1.3 | Comportarea RAID 0 în aplicații ce solicită rate ridicate de cereri de transfer | 74 |
| 3.2 | RAID level 1 | 74 |
| 3.2.1 | RAID 1 și variantele sale arhitecturale | 75 |
| 3.2.1.1 | Mirrored declustering | 75 |
| 3.2.1.2 | Chained declustering | 76 |
| 3.2.1.3 | Group-rotate declustering | 76 |
| 3.2.1.4 | Amplasarea copiilor pe discuri | 77 |
| 3.2.1.5 | Un model de disc | 78 |
| 3.2.2 | Strategii de accesare pentru RAID 1 | 82 |
| 3.2.2.1 | Strategia Random Join (RJ) | 82 |
| 3.2.2.2 | Strategia Shortest Queue (SQ) | 82 |
| 3.2.2.3 | Strategia Minimum Seek (MS) | 82 |
| 3.2.3 | Analiza comportării tipurilor arhitecturale RAID 1 cu diverse strategii | 83 |
| 3.2.3.1 | RAID1 cu strategie RJ | 84 |
| 3.2.3.2 | RAID1 cu strategie SQ | 85 |
| 3.2.3.3 | RAID1 cu strategie MS | 86 |
| 3.2.3.4 | Analizarea diverselor strategii în cazul Group-Rotate Declustering | 87 |
| 3.2.3.5 | O analiza a metodei Chained Declustering | 88 |
| 3.3 | RAID level 2 | 93 |
| 3.4 | RAID level 3 | 93 |
| 3.5 | RAID level 4 | 94 |

| | | |
|------------|--|------------|
| 3.6 | RAID level 5 | 95 |
| 3.6.1 | Strategia AR pentru RAID 5 | 95 |
| 3.6.2 | Strategia BS pentru RAID 5 | 97 |
| 3.6.3 | Tipuri semnificative de RAID 5 | 97 |
| 3.6.3.1 | RAID5/UPGD | 97 |
| 3.6.3.2 | RAID5/VSP | 100 |
| 3.6.3.2.1 | Organizarea unui RAID/VSP | 103 |
| 3.6.3.2.2 | Analiza performanței în cazul scrierii aleatoare ale unui singur bloc | 108 |
| 3.6.3.2.3 | Comparație VSP cu LFS | 111 |
| 3.6.3.2.4 | Costul recuperării datelor de pe discul căzut | 114 |
| 3.6.3.2.5 | Concluzii | 117 |
| 3.6.3.3 | RAID5/ NRP | 119 |
| 3.6.3.3.1 | Construirea distribuției bazate pe NRP | 121 |
| 3.6.3.3.2 | Algoritmul “Thorp’s shuffle” | 122 |
| 3.6.3.3.3 | Generarea secvenței de biți aleatori | 123 |
| 3.6.3.3.4 | Un model pentru RAID5/NRP | 124 |
| 3.6.3.3.5 | Analiza performanței RAID5/NRP în modul de operare normal | 126 |
| 3.6.3.3.6 | Analiza performanței RAID5/NRP în modul de operare în avarie | 127 |
| 3.6.3.3.7 | Validarea rezultatelor | 130 |
| 3.6.3.3.8 | Evaluarea performanțelor | 132 |
| 3.7 | RAID level 6 | 137 |
| 3.7.1 | Descrierea RAID level 6 | 138 |
| 3.7.1.1 | Comportare RAID level 6 în funcționare normală | 139 |
| 3.7.1.2 | Protecția la căderi în RAID level 6 | 139 |
| 3.7.1.3 | Domeniile de utilizare pentru RAID level 6 | 139 |
| 3.8 | Evaluarea performanțelor matricelor de discuri – RAID | 140 |
| 3.8.1.1 | Analiză comparativă între RAID1 și RAID5 pentru volum mare de cereri | 141 |
| 3.8.1.2 | Analiză comparativă între RAID1 și RAID5 pentru cereri de volum mare | 145 |
| 3.8.1.3 | Concluzii | 150 |
| 3.9 | Concluzii | 151 |
| 4 | PROPUNERE RAID LEVEL 7 | 152 |
| 4.1 | UPGD și rezervarea distribuită | 153 |
| 4.2 | UPGD/DS și căderile multiple | 155 |
| 4.3 | O analiză a algoritmului Ng-Mattson de creare a UPGD/DS | 156 |
| 4.4 | Optimizarea matricelor de discuri UPGD/DS compacte | 163 |
| 4.5 | Algoritm optimizat de creere UPGD/DS | 166 |
| 4.6 | Algoritm pentru stabilirea corespondenței biunivoce între elementele multimedierilor A_1 și B_1 | 170 |
| 4.7 | Exemplu de aplicare a algoritmului îmbunătățit de creare a UPGD/ DS | 174 |
| 4.8 | RAID tolerant la mai multe defecte | 201 |
| 4.8.1 | Algoritmul de creare UPGD/DS- MFT | 201 |
| 4.8.2 | Generalizare | 210 |
| 4.9 | Concluzii | 214 |

| | | |
|-----------|---|------------|
| 5 | RAIC – MATRICE DE CANALE REDUNDANTE ȘI INDEPENDENTE | 216 |
| 5.1 | Introducere..... | 216 |
| 5.2 | Construirea sistemului RAIC..... | 219 |
| 5.3 | Aranjamentul octeților de control..... | 220 |
| 5.4 | Organizarea datelor..... | 221 |
| 5.4.1 | Distribuția datelor | 222 |
| 5.4.1.1 | Adresarea independentă | 222 |
| 5.4.1.1.1 | Divizarea fină..... | 223 |
| 5.4.1.1.2 | Divizarea grosieră | 224 |
| 5.4.1.1.3 | Alegerea mărimii unității de divizare..... | 225 |
| 5.4.2 | Mecanismele de redundanță..... | 226 |
| 5.4.2.1 | Replicarea datelor | 227 |
| 5.4.2.2 | Protecția bazată pe paritate..... | 229 |
| 5.4.2.3 | Alte scheme..... | 231 |
| 5.5 | Performanța versus fiabilitatea | 232 |
| 5.5.1 | Negocierea fiabilității în favoarea performanței..... | 232 |
| 5.5.2 | Îmbunătățirea performanței și a fiabilității | 233 |
| 5.5.3 | Tratarea căderilor | 234 |
| 5.6 | Propunere de clasificare pentru RAIC..... | 235 |
| 5.7 | Aplicarea conceptului RAIC în transmisiile seriale de date..... | 236 |
| 5.8 | Concluzii | 239 |
| 6 | CONCLUZII | 240 |
| 7 | BIBLIOGRAFIE | 242 |

0 INTRODUCERE

Pentru câștigul, dar și pentru menținerea unui segment de piață, un anumit produs tehnic trebuie să posede atributul de calitate, constituind unul dintre cei mai importanți factori de concurență. Calitatea este exprimată prin intermediul mai multor parametri, dintre care se distinge, în contextul lucrării de față, fiabilitatea. Din vasta problematică a fiabilității, cercetările cuprinse în prezenta teză sunt focalizate către subsistemul I/O, cu puternice accentuări în domeniul discurilor magnetice, pornind de la ideea că se poate construi un ansamblu fiabil din componente mai puțin fiabile, prin introducerea redundanței. Prin analize de profunzime asupra celui mai vulnerabil nucleu *hardware* al unui sistem de calcul, discurile, care prin cădere pot genera cele mai costisitoare pierderi de date, în teză sunt propuse elemente de structură originale, care permit eficientizarea implementării toleranței funcționale la defectare.

Localizând aportul tezei, noile elemente redundante de structură propuse asigură garantarea corectei funcționări la nivelul subsistemului I/O, cu particularizări importante pentru structurile de discuri.

0.1 Toleranța la defectare ca scop al proiectării

Definind toleranța la defectare, în strânsă legătură cu fiabilitatea, drept capacitatea unui sistem de a-și îndeplini misiunea, în intervalul de timp stabilit, în condiții de exploatare extreme și eventual viciat de un număr limitat de subsisteme defecte.

Strategiile, conceptele și mecanismele specifice toleranței la defectarea sistemelor de calcul sunt sintetizate în Figura 0-1.

Câteva precizări sunt necesare pentru structura din Figura 0-1, unde prin mascarea defectelor se înțelege, de exemplu, aplicarea votării prin majoritate, dar și corectarea unor erori la operații de citire/scriere asupra unui disc/canal. Altă soluție pentru crearea toleranței la defectare o reprezintă reconfigurarea care reprezintă procesul de eliminare a unității defecte din sistem și restaurarea stării operaționale. Apelând la reconfigurare, trebuie avute în vedere următoarele procese: detecția defectelor reprezentând procesul de recunoaștere a anomaliei funcționale apărute; localizarea defectului prin determinarea unității cu funcționare eronată; limitarea propagării defectului prin acțiuni asupra unităților transportatoare de defect; restabilirea operațională din starea de defectare prin reconfigurare.

| | | | |
|------------------------------|--|-------------------------|--|
| Toleranța la defectare | Specificarea Defectelor | Modelarea defectelor | Specificarea claselor de defecte tolerate |
| | Mijloace utilizate – redundanța | | |
| | Acoperirea și diagnoza defectelor | | |
| | Măsurile împotriva defectelor și tratarea erorilor | | |
| | Acoperirea Erorii | Compensarea erorilor | Eliminare logică a componentelor defecte |
| | | Mascarea defectelor | Reconfigurare |
| | Corecția erorilor | | |

Figura 0-1

Pentru soluționarea complexe sarcini de proiectare a unui sistem de calcul cu toleranță la defectare, introducerea redundanței se dovedește măsura cea mai eficientă, aceasta fiind, însă, opusa metodelor convenționale de proiectare care tind spre rezultate care oferă un raport performanță/cost cât mai bun (prin eliminarea redundanței).

Prin redundanță, în contextul prezentei lucrări, vom înțelege prezența în stare funcțională pregătită a mai multor mijloace tehnice decât sunt necesare pentru executarea unor sarcini specificate. Astfel, vom putea determina patru tipuri de redundanță:

- redundanța structurală legată strict de componente *hardware*, reprezentată de extensii ale sistemului cu componente *software* și *hardware* suplimentare față de cele strict necesare execuției funcțiilor specificate;
- redundanța funcțională constă într-o extensie a funcțiilor unui sistem;
- redundanța informațională se bazează pe utilizarea codurilor detectoare și corectoare de erori;

- redundanța temporală presupune un interval de timp suplimentar față de unul necesar funcționării normale, care este necesar la un sistem redundant funcțional pentru executarea tuturor funcțiilor.

0.2 Toleranța la defecte în subsistemul I/O

Subsistemul de I/O este partea cea mai vulnerabilă a unui sistem de calcul deoarece este, în mare măsură, în contact cu mediul înconjurător și, în special, cu operatorii umani. Tot acest subsistem este cel cu evoluția cea mai lentă din punct de vedere al performanței și al fiabilității.

În contextul considerat mai sus, se remarcă discurile magnetice, principalele mijloace de stocare de date, care printr-o funcționare defectuoasă pot compromite total activitatea unui sistem de calcul.

Asigurarea protecției datelor față de situațiile anormale, inclusiv erori umane, căderi de echipamente și condiții ostile a apărut ca o necesitate stringentă chiar cu primul sistem de calcul. *Backup*-ul, copierea periodică a datelor pe un suport mobil păstrat în afara sistemului, este soluția esențială practică, dar lentă, în ciuda abilității date de creșterea numărului de discuri utilizate pentru a preveni alterarea accesului la date corecte în situații de căderi sau dezastre. Și mai puține progrese s-au realizat în domeniul minimizării erorilor umane ceea ce impune utilizarea în continuare a procedurilor de *backup*, chiar și pentru sistemele de discuri care pot supraviețui la căderi și dezastre.

Necesitatea ca sistemul de discuri să mențină un acces continuu *on-line* a apărut în anii '70, o dată cu creșterea masivă a utilizării sistemelor de calcul *on-line*. *Backup*-ul și restaurarea *off-line* erau considerate suficiente pentru a acoperi o cădere. Deoarece accesul la date corupte sau distruse nu are nici o valoare, proiectanții au fost obligați să furnizeze soluții de protejare a datelor și a accesului la acestea în cazul unor căderi. Deoarece datele sunt memorate, în principal, pe discuri, prima soluție de creștere a fiabilității, bazată pe redundanță, a constat în memorarea acelorași informații pe două suporturi diferite (*mirroring*), metoda care protejează datele și menține accesul la ele chiar și în situația căderii unui disc din perechea utilizată.

Soluția *mirroring* presupune costuri cel puțin duble pentru subsistemul de discuri, față de cel clasic, astfel că a apărut necesitatea practică de a găsi soluții mai puțin costisitoare pentru creșterea protecției datelor și asigurarea accesului continuu la acestea. Astfel, la sfârșitul anilor '80 s-au pus bazele conceptului de RAID – Redundant Array of Independent Disks [PDGG88]. RAID utilizează redundanța sub forma unor informații de control (paritate). Față de *mirroring*, unde redundanța era de 100%, la RAID ea poate scădea și sub 10%.

Pentru o tratare unitară, soluția *mirroring* a fost inclusă în familia RAID, care a fost divizată în două părți:

- *mirrored* RAID
- *parity* RAID – RAID protejat prin paritate.

Spre deosebire de *mirrored* RAID, primele *parity* RAID-uri sufereau de o foarte scăzută performanță datorată necesității generării și scrierii parității în timpul operației de scriere și a regenerării *on-line* ca răspuns la o cerere către un disc căzut. Această problemă a fost rezolvată parțial prin tehnici precum *caching*, scriere asistată pe discuri, etc. și a condus la o dezvoltare exponențială a utilizării RAID-ului după mijlocul anilor '90.

Două decade de cercetări în domeniul protecției datelor și menținerii *on-line* a accesului în cazul unei căderi sau a unui dezastru au generat o largă paletă de soluții de sisteme de discuri cu diverse grade de rezistență. Pentru evitarea confuziilor și pentru a ușura alegerea celei mai potrivite soluții pentru o anumită aplicație a fost înființată o asociație internațională a cercetătorilor din domeniul RAID la care s-au afiliat și cei mai importanți producători, asociație denumită RAID Advisory Board's (RAB).

În multe sisteme de calcul, subsistemul I/O reprezintă adesea o pondere importantă în cadrul performanțelor sistemului. De-a lungul ultimilor ani, viteza CPU a crescut în proporții de 40% până la 100% pe an, pe când timpul căutării pe disc s-a îmbunătățit cu doar 7% pe an (creșteri asemănătoare de performanță se pot observa și la celelalte componente ale subsistemului I/O).

Acest fapt a condus la o mare discrepanță între viteza CPU/memorie principală și cea a subsistemului operațiilor I/O, și este de așteptat ca aceasta să crească și mai mult în viitorul apropiat. Bazându-ne pe aceste observații, este de prevăzut ca viitoarele creșteri ale vitezei de lucru ale procesoarelor să determine doar o mică îmbunătățire a performanțelor

totale ale sistemului de calcul. Evident, datorită limitărilor tehnologice, soluțiile pot fi găsite doar la nivel structural.

O idee atractivă este cea numită matrice de discuri (*disk array*), unde subsistemul de discuri este compus din mai multe discuri, cu datele împrăștiate pe toate aceste discuri. Conceptele de „disc întrețesut“ (*disk interleaving*) și „disc secționat“ (*disk striping*) au fost pentru prima oară folosite de către Kim [Kim85], respectiv Salem și Garcia-Molina [GaMo87]. De atunci, o mare cantitate de muncă a fost depusă asupra mai multor modele, referitor la performanțele matricelor de discuri și a fiabilității acestora. Matricele de discuri au fost împărțite în șase clase denumite RAID Levels (*Redundant Arrays of Inexpensive Disks*) [PDGG88]. Printre cele șase, cele mai promițătoare două candidate la obținerea de performanțe ridicate într-un sistem de calcul par a fi RAID 1 (matrice de discuri oglindită - *mirrored disk array*) și RAID 5 (matrice cu paritate rotită - *rotated parity array*), RAID 6 nefiind încă unanim acceptată drept RAID standard, iar RAID 7 reprezentând un rezultat al autorului acestei lucrări.

Primele mele cercetări în domeniul fiabilității transmiterii și stocării datelor s-au focalizat, începând cu anul 1991, în subdomeniul RAID, în principal abordând problematica legată de „Algoritmii de distribuție a datelor, parității și spațiilor de rezervă“.

0.3 Structura tezei de doctorat

Prezenta teză de doctorat abordează problematica creșterii fiabilității și performanței transmiterii și stocării datelor prin soluții structurale îmbinate cu utilizarea unor mecanisme redundante de autocontrol și autocorecție, care să asigure o rezistență sporită a subsistemelor de I/O la căderi și dezastre. În partea finală a lucrării, deschid un nou domeniu de cercetare prin generalizarea principiilor RAID către întreg sistemul de calcul, ceea ce, în opinia mea, va conduce în viitor la un mai mare grad de unificare a algoritmilor utilizați în scopul creșterii fiabilității și performanței.

Lucrarea cuprinde șapte capitole și o listă bibliografică.

Precizez că, în vederea unei prezentări unitare, în lucrare pentru referirea elementelor, structurilor și tehnicilor uzitate am folosit termeni și abrevieri consacrate în limba engleză, acolo unde în limba română nu există o corespondență convenabilă.

După o scurtă introducere care constituie prezentul capitol, capitolul 1 cuprinde o tatonare a domeniului, precum și o trecere în revistă a metodelor generale de creștere a fiabilității transmiterii și stocării datelor.

În capitolul 2, pornind de la focalizarea inițială a cercetărilor mele în domeniul subsistemelor de discuri, se realizează o incursiune în domeniul metodelor consacrate de creștere a fiabilității acestora prin utilizarea codurilor și a redundanței, în special în cadrul structurilor de tip RAID. În final, prezint o clasificare a structurilor de discuri redundante din punct de vedere al fiabilității, în scopul clarificării problemelor existente în domeniul meu de interes.

Capitolul 3 cuprinde o analiză originală a diverselor soluții RAID, în cadrul căreia sunt prezentate câteva dintre rezultatele originale obținute de autor, cum ar fi: variante arhitecturale noi pentru RAID level 1 și 5, strategii de accesare pentru RAID 1 și 5.

În capitolul 4 propun un nou nivel RAID 7, algoritmi de distribuție a datelor și a spațiilor de rezervă pentru RAID 7, etc.

În capitolul 5, autorul propune o generalizare originală a conceptului de RAID, RAIC (*Redundant Array of Independent Channels*) generalizare care deschide practic noi direcții de cercetare, unificând abordarea din punct de vedere al performanței și fiabilității pentru toate subsistemele unui sistem de calcul. După o trecere în revistă a problemelor legate de RAIC (organizarea datelor, distribuția datelor, mecanisme de redundanță, tratarea căderilor etc.) se propune o clasificare a soluțiilor RAIC, analogă celei existente pentru RAID. Menționez că numerele de niveluri nu presupun nici un fel de ierarhizare din punct de vedere al fiabilității sau performanței.

În finalul capitolului 5 se prezintă o descriere a unor experimentări cu structuri RAIC și câteva rezultate practice obținute.

Capitolul 6 prezintă concluziile și sinteza contribuțiilor originale aduse de autor pe parcursul lucrării.

Autorul tezei mulțumește în mod deosebit conducătorului științific, domnului prof.dr.ing. Mircea Vlăduțiu care, prin îndrumarea competentă într-o direcție de un puternic interes actual, prin profesionalism în abordarea problemelor specifice și validarea soluțiilor alese, și nu în ultimul rând prin răbdarea și timpul acordat de-a lungul perioadei de pregătire, a contribuit decisiv la elaborarea tezei.

Autorul mulțumește, de asemenea, colegilor de la Catedra de Calculatoare a Facultății de Automatică și Calculatoare din cadrul Universității Politehnice din Timișoara, pentru sugestiile și recomandările primite.

1 METODE DE CREȘTERE A FIABILITĂȚII ȘI DISPONIBILITĂȚII

1.1 Indicatori pentru caracterizarea fiabilității și disponibilității

1.1.1 Clasificarea defecțiunilor

Într-un sistem există trei categorii majore de surse de erori:

- greșeli de proiectare;
- defecțiuni fizice;
- erori de interacțiune ale operatorului uman (*procedural errors*).

Aceste surse de erori contribuie la căderile sistemului. Pe de o parte, greșelile de proiectare apar atât în *hardware*, cât și în *software*. Acestea din urmă sunt predominante și se elimină greu din sistem [Toy88]. Pe de altă parte, defecțiunile fizice sunt rezultatul îmbătrânirii componentelor sau a influenței mediului, care cauzează devierea unor caracteristici ale echipamentului peste limitele câmpului specificat al toleranței. Se produc, astfel, erori determinate de așa-numitele defecțiuni parametrice *hardware*.

Acțiunile improprii ale operatorului la panourile de control și de întreținere pot determina funcționarea defectuoasă a sistemului. Ceea ce introduce operatorul are prioritate mare, sistemele fiind foarte vulnerabile la comenzi improprii și la erori de operare. Proiectarea unui sistem tolerant la defecțiuni necesită găsirea unei metode care să împiedice ca o componentă *hardware* deteriorată să determine o defecțiune logică și ca aceasta să determine, la rândul ei, o eroare.

În relativ scurta istorie a designului tolerant la defecte au existat multe confuzii între definiții și concepte fundamentale. De exemplu, se foloseau, în mod interschimbabil, termeni ca: *defecțiune*, *deteriorare* și *eroare*. Unii autori susțin că o deteriorare a avut loc atunci când calculatorul nu mai răspunde la comenzi. Dimpotrivă, alți autori afirmă că o deteriorare este o defecțiune fizică specifică, mai degrabă, componentelor electronice. Datorită diversității de păreri în privința înțelegerii acestor noțiuni și pentru a se evita echivocul în lectura acestei lucrări, consider utilă definirea unor termeni, care, deși au serioase trimiteri bibliografice - [Lala85] și [John89] -, sunt încă departe de acceptarea universală. Aceste definiții sunt:

- Deteriorarea (*failure*) apare atunci când ieșirea obținută diferă de cea proiectată. Ieșirea poate fi privită la diferite niveluri de abstractizare: la nivel de componentă electronică, la nivel de bloc sau la nivel de sistem;
- Defecțiunea sau defectul (*fault*) este o stare eronată a mașinii sau a sistemului de programe, stare generată de: deteriorări ale componentelor, interferențe fizice cu mediul, erori ale operatorului uman sau proiectare incorectă;
- Eroarea (*error*) este manifestarea unei defecțiuni într-un program sau într-o structură de date. Ea poate apărea la o anumită distanță în spațiu și timp față de locul, respectiv momentul, defectării;

Definițiile de mai sus pot fi caracterizate de următoarele adjective:

- Adjectivul „permanent“ (idem în engleză) descrie o deteriorare, o defecțiune sau o eroare cu caracter continuu și stabil. Dacă este vorba despre echipament, atunci o deteriorare permanentă reflectă o schimbare fizică ireversibilă;
- Adjectivul „intermitent“ (*intermittent*) descrie o defecțiune sau o eroare care se repetă ocazional;
- Adjectivul „tranzitoriu“ (*transient*) descrie o defecțiune sau o eroare non-repetitivă, cauzată de condiții temporare. Exemple sunt acțiunea particulelor α asupra memoriilor semiconductoare sau variația tensiunii de alimentare a echipamentelor. Astfel de erori sunt nereparabile, pentru că nu există defect fizic al echipamentului. Evitarea lor se poate face doar printr-o proiectare mai riguroasă.

Deteriorările sistemului sunt clasificate, funcțional vorbind, în defecțiuni *hardware*, erori *software* și erori de procedură. Există relații strânse între cele trei tipuri de deteriorări de sistem și categoriile surselor de defecțiuni. Defecțiunile fizice apar numai la echipament. Erorile de proiectare, chiar dacă apar în *hardware*, produc, în majoritatea lor, erori de tip *software*. Erorile de interacțiune sunt cauzate, de obicei, de greșelile făcute de operator.

În afara deteriorărilor normale ale dispozitivelor, la defecțiunile *hardware* mai contribuie condițiile mediului înconjurător (perturbațiile), erorile de proiectare (greșelile) și cronografiile improprie (*timing problems*). Greșelile de proiectare, ca în cazul erorilor *software*, sunt foarte greu de apreciat cantitativ. Calculele de fiabilitate se bazează mai mult pe defecțiunile dispozitivelor.

1.1.2 Estimarea fiabilității

Estimarea fiabilității este un proces de apreciere a fiabilității realizabile de un articol (element, subsistem sau sistem), având disponibile datele ratei de defectare, adică, prin estimarea fiabilității, se apreciază probabilitatea îndeplinirii obiectivelor articolului respectiv pentru o aplicație specifică. Aceste calcule sunt utile în stadiile de început ale unui proiect.

După îndepărtarea defecțiunilor timpurii, componentele se mențin pe o perioadă lungă la o rată de defectare aproximativ constantă. În timpul acestei perioade, rata de defectare este, de obicei, scăzută și este puțin probabil ca defecțiunile să provină de la o singură cauză. Rata constantă de defectare, reprezentată în porțiunea de viață utilă, scoate în evidență faptul că probabilitatea de defectare este independentă de vârstă. Pentru orice rată constantă de defectare, valoarea fiabilității depinde numai de timp. Funcția fiabilității $R(t)$ (*reliability*), care este caracterizată de o rată constantă de defectare este o distribuție negativă exponențială și are forma:

$R(t) = e^{-\lambda t}$, unde λ este rata de defectare, iar t este timpul.

Timpul mediu dintre defectări (*Mean Time Between Failures*; MTBF) este timpul mediu, exprimat, de obicei, în ore, în care un articol ar putea să funcționeze corect înainte de a se deteriora și reprezintă măsura cantitativă a fiabilității [Lala85], [Toy88]. În general, MTBF-ul unui sistem ar putea fi tratat ca integrala funcției de fiabilitate a întregului sistem: $\int_0^{\infty} R(t) dt$

În câteva cărți de literatură, intervalul este definit ca timp mediu de defectare (*Mean Time To Failure*; MTTF). Din punct de vedere tehnic, MTTF și MTBF nu sunt identici, MTBF-ul poate fi definit, într-o formă alternativă, astfel:

$MTBF = MTTF + MTTR$, unde MTTR este timpul mediu de reparație (*Mean Time To Repair*).

Rata de reparație μ (*repair rate*) sau inversul timpului mediu de reparație, adică

$$\mu = \frac{1}{MTTR}$$

este un alt factor care afectează substanțial fiabilitatea și mentenabilitatea unui sistem. Când o unitate dintr-un sistem duplicat este defectă, sistemul este dependent de a doua unitate pentru a-și

continua operația. Dacă unitatea defectă este reparată repede, riscul deteriorării întregului sistemului este mic deoarece unitatea a doua va opera astfel încât să păstreze integritatea funcționării sistemului. Întrucât sistemul este vulnerabil numai până la reparația unității defecte, un timp scurt de reparație poate crește foarte mult fiabilitatea sistemului.

În cazul RAID, pentru estimarea fiabilității, nu mai sunt suficienți MTBF, MTTF și MTTR, parametri care descriu cu precădere caracteristicile fizice ale dispozitivelor. Astfel, ca de curând RAB [Mass97] a introdus:

- **MTDA** – *Mean Time until (Loss of) Data Availability*, timpul mediu până ce căderea unei componente cauzează o pierdere a accesibilității într-o matrice de discuri. Pierderea accesibilității nu implică pierderea datelor (pentru multe tipuri de căderi datele rămânând intacte – ex. Căderea unui controler neredundant). Evident că MTDA se poate folosi și pentru RAIC, la fel ca și MTDL.
- **MTDL** – *Mean Time until Data Loss*, timpul mediu până ce căderea unei componente generează o pierdere permanentă de date într-o matrice de discuri. Conceptual, MTDL este similar cu MTBF, dar ține cont de posibilitatea ca redundanța RAID-ului să protejeze împotriva pierderii de pentru un număr limitat de căderi de componente.

În prezenta teză de doctorat, autorul introduce trei noi indicatori necesari pentru evaluarea fiabilității, pentru RAID 7 și RAIC 7:

MTDR – *Mean Time until Data is Recovered*, timpul mediu de la apariția unei căderi de disc până când sunt regenerate toate blocurile pierdute în spațiile de rezervă

MTRD – *Mean Time until Data is Regenerated*, timpul mediu în care RAID sau RAIC răspunde la o cerere dintr-o zonă afectată, prin regenerarea datelor utilizând redundanța.

MTAR – *Mean Time until Array is Reconstructed*, timpul mediu de reconstrucție (în starea inițială) a matricei după înlocuirea discului sau canalului căzut.

1.1.3 Disponibilitatea

Disponibilitatea $A(t)$ (*availability*) unui echipament este o funcție de timp. Ea se poate defini ca probabilitate de funcționare a echipamentului la momentul t , atât timp cât este folosit conform specificațiilor.

Conceptul de disponibilitate este utilizat în măsurarea eficacității sistemului. Atât fiabilitatea, cât și mentenabilitatea, își aduc aportul la definirea conceptului de disponibilitate. Această conexiune poate fi exprimată în felul următor [Lala85], [Toy88]:

$$A(t) = \frac{MTTF}{MTTF + MTTR}$$

unde funcția fiabilității se presupune că este distribuția exponențială $R = e^{-\lambda t}$. Deoarece MTTF-ul este o măsură a fiabilității, iar MTTR-ul este o măsură a mentenabilității, ele se pot negocia astfel încât să se obțină o disponibilitate dată.

1.1.4 Dependabilitatea

În 1982, Laprie și Coste au definit formal dependabilitatea ca un cadru de specificații de nivel înalt, care include mărimi ca fiabilitatea și disponibilitatea, ca două fațete distincte ale specificațiilor sistemului. Pe baza caracteristicilor unui sistem, dependabilitatea lui se descrie fie prin fiabilitate, fie prin disponibilitate. De exemplu, putem caracteriza mai bine un sistem nereparabil prin fiabilitatea lui, iar un sistem reparabil prin disponibilitatea lui.

În literatura mai recentă [John89], se menționează că termenul de dependabilitate include, pe lângă conceptele de fiabilitate și disponibilitate menționate anterior, și pe cele de siguranță în funcționare, mentenabilitate, performabilitate și testabilitate. Dependabilitatea reprezintă calitatea deservirii furnizate de către un anumit sistem. Fiabilitatea, disponibilitatea, siguranța, mentenabilitatea, performabilitatea și, în fine, testabilitatea sunt mărimi utilizate pentru estimarea cantitativă a dependabilității unui sistem. Ținta proiectării tolerante la defecțiuni este îmbunătățirea dependabilității prin înarmarea unui sistem pentru realizarea funcției cerute, în prezența unui număr dat de defecte. Totuși, este de notat că un sistem tolerant la defecțiuni nu este neapărat de înaltă dependabilitate și nici înalta dependabilitate nu necesită neapărat toleranța la defectare.

1.1.5 Mentenabilitate

Sistemele de calcul, ca orice alte produse, se pot clasifica în două categorii: a) sistemele cu funcție unică (simplă), la care prima defectare constituie și finalul vieții lor și b) sistemele cu funcție repetată sau sistemele cu reînnoire (restabilire), la care elementele defecte pot fi înlocuite cu altele bune, caz în care aceste produse au caracter reparabil. Prin mentenanță, se înțelege ansamblul tuturor acțiunilor tehnico-organizatorice necesare, efectuate în scopul menținerii sau restabilirii

unui produs în starea necesară îndeplinirii funcției cerute. Aceste acțiuni pot avea caracter corectiv (depistarea cauzei unei defecțiuni, repararea defectului prin înlocuirea elementelor defecte, verificarea corectitudinii operațiilor de mentenanță întreprinse) sau caracter preventiv (revizii, reglaje, verificări și reparații planificate, executate în vederea evitării unor viitoare defecțiuni).

În aceste condiții, este necesar să se exprime cantitativ aptitudinea unui produs de a fi repus în funcțiune în urma unui defect. Exprimarea se face, ca și în cazul fiabilității, printr-o probabilitate în funcție de timp. Mentenabilitatea este, așadar, aptitudinea unui produs ca, în condiții date, să fie menținut sau restabilit în starea de a-și îndeplini funcția specificată, atunci când acțiunile de mentenanță se efectuează în condiții precizate și într-un timp dat, cu procedee și remedii prescrise. Corespunzător acestei definiții, legătura dintre aspectul probabilistic și cel funcțional se exprimă astfel:

$$M(t_r) = \text{Prob}(t_r \leq T_r)$$

unde t_r este timpul de restabilire, T_r este o limită impusă duratei de restabilire, iar $M(t_r)$ este funcția de mentenabilitate.

Relația între mentenabilitatea M și rata de reparație μ sau timpul mediu de reparație MTTR este următoarea [Lala85]:

$$M(t_r) = 1 - e^{-(\mu \cdot t_r)} = 1 - e^{-\left(\frac{t_r}{\text{MTTR}}\right)}$$

La fel ca fiabilitatea, testabilitatea, dependibilitatea etc., mentenabilitatea trebuie planificată și estimată sau evaluată încă din fazele cele mai timpurii ale conceperii unui produs.

1.1.6 Cerințele fiabilității aplicate

Cerințele fiabilității variază considerabil de la o aplicație la alta. De exemplu, obiectivele fiabilității unui oficiu telefonic cu comutare digitală, proiectat pentru aplicații telefonice, sunt:

- timpul total în care sistemul rămâne neoperativ să nu depășească trei minute pe an, cu un timp mediu de reparație (MTTR) de patru ore;
- în timpul funcționării să nu se piardă sau să nu se trateze incorect mai mult de 0,01% dintre apeluri.

Funcționarea satisfăcătoare nu înseamnă, în acest caz, o fiabilitate de 100%. Sunt permise, totuși, câteva apeluri incorecte, avându-se în vedere că utilizatorul va forma din nou numărul dorit, obținând corecția erorii. Pe de altă parte, o defecțiune într-un echipament critic, cum ar fi

transponderele de telecomunicații de pe sateliții artificiali, ar determina căderea unor întregi sisteme, ceea ce este inadmisibil. În astfel de cazuri, funcționarea satisfăcătoare impune o fiabilitate de 100%.

Durata timpului de funcționare pentru un sistem de comutare telefonică este, pur și simplu, durata de viață a echipamentului care a fost folosit în sistem. Sistemul de comutare trebuie să funcționeze continuu, fără întreruperi, până când echipamentul este înlocuit la sfârșitul vieții lui sau este îndepărtat pentru orice alt motiv. Deoarece servirea trebuie să fie asigurată 24 de ore pe zi, nu poate să existe timp programat pentru reparații sau pentru întreținere, timp în care sistemul să fie neoperational. Obiectivul fiabilității de trei minute de neoperare pe an (două ore de nefuncționare în patruzeci de ani [Prad86]) și MTTR-ul de patru ore pot fi traduși într-o disponibilitate de sistem de 99,9995%.

În contrast, aplicații critice pentru controlul avioanelor, ca și în cazul lui SIFT (*Software-Implemented Fault Tolerance*) și FTMP (*Fault-Tolerant MultiProcessor*) pretind fiabilitate ultra-înaltă. Parametrii proiectării cer ca sistemul să treacă prin mai puțin de 10^{-9} deteriorări în timpul unei misiuni de 10 ore cu un MTTR de 10 ore. Factorul de disponibilitate este egal, în această situație, cu 99,9999999%.

Cerințele fiabilității determină într-un grad mai mare arhitectura sistemului, tipul și cantitatea de redundanță precum și costul sistemului. Redundanța duală se dovedește a fi adecvată aplicațiilor telefonice. Pe de altă parte, pentru controlul critic al avioanelor, triplarea cu rezervă este necesară pentru a achita un grad mare de fiabilitate.

1.1.7 Tehnici de redundanță

Defecțiunile *hardware* pot afecta secvențele de comandă sau cuvintele de date care se află în interiorul mașinii. Acest lucru duce la două tipuri de erori:

- secvența programului este neschimbată, dar defecțiunea a afectat rezultatele finale;
- secvența programului este schimbată, iar programul nu mai execută algoritmul specificat.

Erorile *software* sunt rezultatul unor traduceri (translatări) greșite sau ale unor implementări improprii ale algoritmilor originali și, de asemenea, deviază execuția instrucțiunilor de la secvența corectă. În multe situații, din păcate, defecțiunile *hardware* nu pot fi distinse de cele *software*. Din acest motiv, sistemele trebuie să fie tolerante la defecțiuni, indiferent de sursa acestora.

O modalitate bună de a produce calculatoare tolerante la defecte este aceea de a introduce redundanța prin multiplicarea părților lor. Redundanța permite calculatoarelor să ocolească erorile, în felul acesta rezultatele fiind corecte. Această metodă este cunoscută ca redundanță de protecție, fiind compusă din redundanțe *hardware*, *software* și de timp. Redundanța *hardware* este constituită din componente adiționale, care detectează și corectează erorile. Redundanța *software* conține programe adiționale, care au capacitatea de a restabili funcționalitatea sistemului, în condiții problematice. De asemenea, mai poate conține programe de detecție a erorilor, de diagnoză și de autocontrol, prin care să testeze periodic toate circuitele logice ale calculatorului. Redundanța temporală este obținută prin repetarea unei operații eronate, imediat după repetarea unei erori. Redundanța temporală se implementează adesea prin *hardware*. De exemplu, logica *hardware* poate iniția recitirea automată a unei locații de memorie în care s-a detectat o eroare de paritate.

Chiar dacă redundanța de protecție este clasificată funcțional în trei tipuri diferite, în practică ele se contopesc adesea, pentru asigurarea unei fiabilități sporite. În cazul redundanței *software*, programul de control are nevoie atât de spațiu de memorie (*hardware*), cât și de timp de execuție.

Fiecare dintre aceste tipuri de redundanță, precum și diferitele lor combinații, au fost folosite în proiectarea calculatoarelor tolerante la defecțiuni. Prioritatea acordată la implementare unuia dintre tipurile de redundanță enumerate mai sus se face în funcție de tipul aplicației și de restricțiile impuse de costuri, de cerințele de timp și de cerințele de fiabilitate.

1.2 Metode de autocontrol, implementate *hardware*

Atât structura redundanței dinamice cât și cea a redundanței statice, descrise mai sus, sunt metode care folosesc părți de rezervă pentru a face sistemul capabil să tolereze defecțiunile. Atunci când se folosește redundanța statică, unitățile de rezervă (circuite, componente sau subsisteme) sunt părți permanent active ale sistemului. Ele corectează erorile sau le maschează împiedicând, astfel, propagarea lor în sistem. Funcția mascării intervine automat, iar acțiunea de corectare este imediată și încorporată (*wired in*). Tipuri statice de redundanță au fost folosite, mai întâi, în aplicațiile militare care pretindeau o înaltă fiabilitate pentru o perioadă scurtă de timp, iar, mai recent, în aplicațiile comerciale. Redundanța dinamică, la care subsistemele adiționale

servesc ca rezerve în interiorul sistemului, s-a folosit în aplicațiile comerciale.

Componentele majore ale tolerării defecțiunilor sunt: detectarea erorii, diagnosticarea erorii (izolarea), restabilirea și repararea. În cazul redundanței dinamice, cea mai importantă componentă este detectarea erorii. Dacă toate erorile ar fi detectate și s-ar aplica tehnicile corespunzătoare de restabilire, atunci nici o defecțiune n-ar conduce sistemul la deteriorare. Din păcate, această acoperire nu se poate realiza în practică.

Viteza detectării erorii influențează procesul de localizare a defecțiunii și stăpânirea ei. De asemenea, este important să se realizeze cât mai repede următorul pas, mânguirea sau izolarea defecțiunii, astfel încât defecțiunea să nu se propage. Detectarea întârziată poate deforma date importante în tot sistemul. Viteza detectării este, de asemenea, importantă pentru localizarea sursei de erori: cu cât întârzierea este mai mare, cu atât este mai greu să se găsească sursa. Diagnosticarea incorectă împiedică rezervele funcționale de la înlocuirea corectă a unităților defecte. Mai mult decât atât, viteza de detectare afectează direct revenirea sistemului.

În general, detecția erorii este realizată prin utilizarea de *hardware*, *firmware* și *software*. Tipul schemelor de control utilizate depinde atât de structura logică a mașinii cât și de utilizarea operațională și funcțională a datelor și a semnalelor de comandă.

Schemele *hardware* de detectare a erorii încorporate într-un sistem de calcul pot avea mai multe forme. Majoritatea acestor tehnici se încadrează în clasificarea următoare:

- controale prin copii (*replication checks*);
- controale de codificare (*coding checks*);
- controale de temporizare (*timing checks*);
- controale de excepții (*exception checks*).

Detectarea erorii poate fi amplasată strategic în interiorul unei unități funcționale sau la interfață. Este mai avantajos dacă detecția erorii se face intern, la stadiul ei timpuriu, în timpul funcționării sistemului care generează rezultatele.

Controalele interne sau timpurii minimizează cantitatea activității sistemului și tranzițiile eronate cauzate de o defecțiune. Astfel, nu este timp suficient pentru a se propaga defecțiunea în interiorul sistemului,

iar acțiunile necesare pentru izolarea ei și pentru revenire la normal vor fi, probabil, mai simple. Pe de altă parte, controalele de interfață sau cele de ultimul moment sunt activate înainte de a se transmite orice fel de rezultat de la o unitate funcțională la alta. Acest fapt împiedică propagarea erorilor în exterior spre o altă unitate funcțională și simplifică cea mai dificilă problemă, și anume revenirea globală. Eroarea este conținută în interiorul nivelului în care a fost detectată.

1.2.1 Controlul prin copii

Controlul prin copii este una dintre cele mai complete metode de detectare a erorilor dintr-un sistem de calcul. Din cauza echipamentului *hardware* necesar, această metodă este, de asemenea, și cea mai costisitoare tehnică de redundanță. Numărul copiilor dintr-un sistem nu trebuie să fie limitat la doi. Atunci când trei sau mai multe copii ale unui sistem sunt implicate, controlul comparării este referit ca votare. Votarea asigură suprimarea unei ieșiri eronate prin mascarea în favoarea majorității ieșirilor. Un exemplu de controlare prin copii într-un sistem de redundanță triplă, modulară este FTMP (*Fault Tolerant Multiple Processor*) dezvoltat pentru calculatoarele avioanelor.

Controalele prin copii bazate pe copii identice ale unui sistem nu detectează consecințele greșelilor de proiectare, deoarece sunt afectate toate copiile unității. Din acest motiv, este necesară apelarea la tehnici formale de verificare a proiectării astfel încât să se asigure că nu există greșeli de proiectare în structură care ar putea ocoli protecția realizată de copiere.

În [MaMc88], Mahmood și McCluskey prezintă câteva tehnici de detecție a erorilor bazate pe utilizarea unui procesor *watchdog*. Acest procesor, care poate fi numit coprocesor de control, este mai puțin complex decât procesorul principal și detectează erorile prin monitorizarea comportamentului sistemului. Informația furnizată către coprocesorul de control poate fi legată de comportamentul acceselor la memorie, cursul comenzilor, semnalele de comandă sau chiar de justificarea rezultatelor. Asemenea duplicării, pentru detecția erorilor, această tehnică nu depinde de modelul de defecțiuni adoptat și are, în schimb, avantajul că necesită mai puțin *hardware*.

1.2.2 Controale de codificare

Codurile detectoarelor de erori se formează prin adăugarea unor biți de control la un cuvânt de date. În acest caz, capacitatea de detecție a erorilor este direct proporțională cu numărul de biți de control incluși într-un cuvânt. Există două tipuri de controale de codificare: separabile

și neseparabile. În literatura de specialitate, ele mai poartă denumirea de controale de codificare sistematice și nesistematice. Controalele separabile (*separable checks*), cum sunt controlul de paritate și codurile aritmetice, sunt caracterizate prin adăugarea biților de control la cuvintele de date. Controalele neseparabile (*nonseparable checks*), cum este codul m-din-n, sunt codificate în formate specializate.

1.2.3 Controale de temporizare

Controlul de temporizare este o formă eficientă de a controla *software*-ul pentru detecția erorilor în programele duplicate, atunci când specificația unei componente include constrângeri de timp.

Un temporizator *hardware*, numit și *watch-dog timer*, este folosit în sisteme de comutare telefonică pentru apărarea împotriva erorilor de programare din cauza cărora sistemul nu-și mai poate reveni. Conceptul este relativ simplu. Un ceas *hardware* merge continuu în procesor și este resetat periodic de către programul principal dacă nu se întâmplă nimic neobișnuit care să devieze programul din secvența lui normală de execuție.

Dacă dintr-un anumit motiv (eroare *software* sau defecțiune *hardware*) secvența execuției nu se mai întoarce la programul principal, ceasul nu se resetează. În această situație, ceasul pune în circulație o cerere de întrerupere de înaltă prioritate și întreprinde acțiunile necesare reinițializării sistemului. În cazul în care există în sistem încă un procesor *stand-by*, întreruperea ceasului poate produce comutarea automată a controlului pe mașină *stand-by*, într-o tentativă de revenire din eroarea de *time-out*.

Controalele de temporizare asupra funcționării unui circuit nu sunt controale complete pentru sistem. Ele relevă prezența defecțiunilor (erorilor), dar nu și absența lor. Ca urmare, arhitecții calculatoarelor folosesc controalele de temporizare ca să suplimenteze alte controale și ca să acopere un procent mai mare de defecțiuni ale unui sistem.

1.2.4 Controale de excepție

Programele rulează în medii protejate folosind seturi de constrângeri prestabilite. Dacă programele nu conțin erori, ele țin cont de aceste constrângeri și își realizează în mod corespunzător funcțiile specificate. Totuși, greșelile de proiectare din *software* violează, adesea, aceste constrângeri, influențând, în mod nefavorabil, întregul sistem. Unele circuite *hardware* de detecție sunt adesea proiectate în sistem pentru a recunoaște erorile de proiectare și pentru a le trata ca excepții. Deci,

tratarea excepțiilor se referă la detectarea și „mânuierea“ evenimentelor anormale și nedorite.

Constrângerile pot fi atribuite structurii *hardware* sau celei *software*. Exemple de constrângeri *hardware* sunt: alinierea improprie de adrese, locații de memorii neechipate, opcod neutilizat, depășirea stivei etc. Structura *software* pune, de asemenea, câteva constrângeri sistemului pentru a-i îmbunătăți rezistența și pentru a furniza un mediu protejat programelor de aplicații. Câteva controale *software* de excepție sunt: execuția ilegală a instrucțiunilor privilegiate, depășirea de adresă (*out of range*), violarea memoriei (*memory acces violation*), operanzi ilegali, operații aritmetice necorespunzătoare, depășirea sau împărțirea cu zero. Detectarea unei erori inițiază o excepție care este urmată de invocarea automată a subrutinei de mânăuire a excepției respective. În multe cazuri, excepția este atribuită unei greșeli de proiectare a programului.

1.2.5 Restabilirea erorii

Restabilirea este cea mai complexă și mai dificilă funcție pentru toate sistemele din cauza multitudinii stărilor posibile ale sistemului care pot să apară sub condiții problematice [Toy88]. Neajunsurile proiectării *hardware* sau *software* de a detecta erorile atunci când apar au un efect direct asupra abilității sistemului de a se restabili. Când erorile continuă să fie nedetectate, sistemul rămâne defect până când problema este recunoscută. Un alt tip de restabilire apare atunci când sistemul nu este capabil să izoleze corespunzător un subsistem defect și să reconfigureze, în jurul lui, un sistem operațional. Detectarea cât mai rapidă a unei erori ușurează determinarea componentei defecte și restabilirea sistemului. Deci, detectarea rapidă a erorii și izolarea ei sunt condițiile fundamentale ale restabilirii.

1.2.6 Clasificarea procedurilor de restabilire

Există trei clase de proceduri de restabilire: restabilirea totală (*full recovery*), restabilirea degradată (*degraded recovery*) și închiderea pentru salvare (*safe shutdown*). Procedurile de restabilire pot fi invocate automat (fără nici o interacțiune între operatorul de întreținere și sistem) sau manual. Restabilirile inițiate manual fac uz de secvențe vaste, controlate prin program.

1.2.6.1 Restabilirea totală

Într-o facilitate de timp real, sistemul care furnizează servicii în mod continuu trebuie să rămână operațional chiar și într-un mediu defect. Aceasta înseamnă că simptomele trebuie să fie recunoscute repede și

unitățile defecte trebuie să fie reparate cu puțină sau chiar fără intervenție din partea utilizatorului.

Pentru a dispune de sistem, tot timpul și la întreaga sa capacitate, subsistemele corespunzătoare de schimb trebuie să fie disponibile ca unități de înlocuire pentru cele defecte.

O procedură de restabilire totală are nevoie, în mod obișnuit, de toate cele cinci aspecte ale calculării tolerante la defecte: detectarea erorii, izolarea erorii, restabilirea sistemului, diagnosticarea erorii și repararea ei. Perioada de întrerupere a servirii, în timpul procesului de restabilire, este minimizată și, de obicei, nu este observată de către utilizator.

Diagnosticarea și repararea, care consumă cel mai mult timp în cadrul procedurii de restabilire, pot fi amânate și intercalate printre operațiunile normale ale sistemului.

Restabilirea degradată

Ca și în cazul restabilirii totale, toți pașii implicați în tolerarea defecțiunii (detectare, izolare, restabilirea sistemului, diagnoză și reparare) trebuie să fie incluși în procedură. Aceeași secvență de evenimente apare la procedura degradată, exceptând faptul că, în acest caz, nu se cuplează alt subsistem în locul celui defect. În schimb, componenta defectă este scoasă din funcție și sistemul se întoarce la o stare operațională fără defecțiuni. Funcțiile de calculare selectate sunt lăsate să opereze în sisteme care utilizează proceduri de restabilire degradată, iar din această cauză caracteristicile performanței lor în timp real scad sub un standard normal acceptat până când se realizează reparațiile corespunzătoare.

Închiderea pentru salvare

Aceasta este, adesea, denumită și operație de salvare din malfuncționare (*fail-safe operation*). Ea apare ca un caz limită al restabilirii degradate, când capacitatea de calcul a sistemului a scăzut sub limita minimă acceptabilă a operației. Scopurile închiderii pentru salvare sunt:

- evitarea distrugerii elementelor sistemului sau a modulelor *software* stocate (programe sau date) care mai pot fi compromise după ce a intervenit operația de restabilire degradată;
- încetarea interacțiunii cu orice sistem asociat (utilizator) într-un mod ordonat;
- distribuirea diagnosticului și a mesajelor de închidere la utilizatori, la anumite sisteme și la personalul de întreținere.

Categoria procedurii de închidere pentru salvare este similară unui sistem fără redundanță. Acțiunea de izolare a defecțiunii trebuie să fie inițiată pentru a se determina identitatea ei și locul în care se află. Funcționarea normală a sistemului, care a fost întreruptă la momentul detectării defecțiunii, trebuie să fie, acum, suspendată datorită diagnozei și reparării. Sistemul trebuie să fie restabilit, atât din punct de vedere *hardware*, cât și *software*, astfel încât să-și poată relua procesarea normală.

1.2.6.2 Reconfigurarea

În sistemele redundante, un ansamblu de piese de schimb sau de unități multifuncționale asigură continuitatea funcționării sistemului. Structura cea mai simplă este configurația duplex, în care fiecare unitate funcțională este duplicată. Dacă una dintre unități se defectează, unitatea duplicată se cuplează păstrând, astfel, continuitatea funcționării, în timp ce unitatea defectă se repară. În timpul perioadei de reparație poate să apară o defecțiune în unitatea duplicată, caz în care întregul sistem se va deteriora. Dacă intervalul de reparație este scurt, atunci probabilitatea de defectare, simultană, a două unități identice este foarte mică.

Configurația duplex cu o singură unitate poate fi partiționată într-o configurație duplex multi-unitate. În acest aranjament, fiecare subunitate conține un număr mic de componente care pot fi comutate într-un sistem în lucru. Sistemul va pica numai dacă o eroare apare în subunitatea redundantă, în timp ce subunitatea originală este supusă reparației. Din moment ce fiecare subunitate conține câteva componente, probabilitatea de apariție a defecțiunilor simultane într-o pereche dublă de subunități este redusă.

1.3 Metode de autocontrol implementate *software*

Necesitatea dezvoltării mecanismelor *software* pentru detectarea erorilor de programare este foarte acută atunci când se pune problema construirii unor sisteme fiabile. Mecanismele *hardware*, chiar dacă sunt foarte avansate din punct de vedere tehnologic, nu sunt specializate în analiza erorilor *software*.

Spre deosebire de ratele de deteriorare ale componentelor *hardware*, despre comportamentul și distribuția erorilor de programare nu există statistici de încredere. Detectia erorilor *software* se poate realiza numai prin recunoașterea comportamentului anormal al sistemului.

Pentru detectarea eventualelor devieri, este nevoie de un set de standarde pentru descrierea comportamentului normal. Prin controlarea

motivațiilor comportării programului în anumite stadii ale calculării, se reușește detectarea rapidă a erorilor și limitarea propagării datelor eronate.

Există mai multe tehnici generale pentru urmărirea comportamentului unui sistem de calcul și detectarea devierilor de la comportamentul normal. Printre acestea sunt: funcția unui proces, secvența de comandă a unui proces, date de proces. Un proces este definit aici ca fiind o porțiune de calcule de sine stătătoare, care, o dată inițiată, se realizează complet fără a mai avea nevoie de intrări adiționale.

1.3.1 Funcția unui proces

Motivarea ieșirilor pentru un set dat de intrări este controlată la nivelul blocului funcțional. Anumite măsuri de performanță pot fi folosite pentru a se indica funcționarea corespunzătoare a sistemului.

Atunci când încărcarea aplicată este normală, dar măsurătorile care caracterizează performanța sistemului (timpul de răspuns, debitul și timpul necesar pentru a se realiza o funcție standard) sunt în afara valorilor limită, sistemul are, probabil, una sau mai multe erori.

1.3.2 Secvența de comandă a unui proces

Secvența de calcule făcută de către un proces care se află în execuție este denumită secvență de comandă (*control sequence*). Fiecare calcul schimbă multe stări ale sistemului și detectează multe dintre neregulile acestuia.

Este evident că *hardware*-ul necesar calcului secvenței de comandă trebuie să fie cât mai fiabil, pentru a nu constitui el însuși o sursă suplimentară de erori sau de defectări.

1.3.3 Date de proces

Integritatea datelor unui sistem și structura lui pot fi observate în timp ce sistemul execută secvența programului. Deteriorarea datelor poate fi cauzată atât de o defecțiune *hardware*, cât și de o eroare *software*.

1.3.4 Restabilirea *software*

Obiectivul restabilirii după o eroare este revenirea sistemului la o stare consecventă și la o funcționare corespunzătoare.

Există trei strategii de restabilire:

- restabilirea cu întoarcere (*backward error recovery*);
- restabilirea cu salt înainte (*forward error recovery*);

- reinițializarea (*reset-ul*).

1.3.5 Diagnosticarea erorii

Ținând cont că detectarea erorii determină corectitudinea funcționării unui circuit, diagnoza erorii localizează defecțiunea la o unitate înlocuibilă. Unitatea înlocuibilă poate fi o componentă, un circuit sau un subsistem. Rutina de diagnosticare a erorii utilizează *hardware* pentru detectarea erorii și secvențe de test în scopul localizării unității defecte. Dacă sursa erorii o constituie o singură entitate, atunci defecțiunea este corectată simplu, prin înlocuirea entității respective, fără să mai fie necesară diagnoza. Dacă circuitul de detecție este mai puțin specializat, atunci rutina de diagnosticare a erorii poate fi solicitată la izolarea, în continuare, a unității afectate.

Cele mai dificile sarcini ale proiectării întreținerii sunt: restabilirea sistemului și diagnosticarea. Eficacitatea lor poate fi determinată prin simularea modului de comportare a sistemului în prezența unei anumite defecțiuni. Prin intermediul simulării deficiențele proiectării pot fi identificate și corectate înainte de a se da sistemul spre utilizare. Este necesar să se evalueze abilitatea sistemului la detectarea erorilor, la revenirea automată la funcționarea normală și la furnizarea informației de diagnoză (de exemplu: locația defecțiunii). Simularea, fizică sau digitală, a defecțiunii este un aspect important al proiectării întreținerii.

1.3.6 Autocontrolul ca mijloc de creștere a fiabilității și disponibilității sistemelor de calcul

1.3.6.1 Metode de autocontrol aplicate la diferite niveluri de structură

1.3.6.1.1 Caracteristicile metodelor de proiectare structurată la nivel de bistabil

Progresele recente din tehnologia VLSI au avut o influență mare asupra testării sistemelor digitale (numerice). Sistemele digitale de astăzi au de la o sută de mii până la un milion de porți de logică aleatoare și celule de memorie ceea ce face ca generarea testului și simularea defecțiunilor să fie extrem de dificile. Chiar dacă se pot folosi mașini *hardware* dedicate sau super-computere pentru a genera teste eficiente, costul acestora face imposibilă aplicarea lor în practică. Pentru a se depăși această problemă, cercetarea s-a îndreptat către găsirea altor metode de testare a sistemelor digitale, incluzând proiectarea pentru testabilitate și generarea de test la nivel funcțional (simularea defecțiunilor).

Dintre toate tehnicile propuse până în prezent, proiectarea pentru testabilitate (DFT, *Design For Testability*) este cea mai renumită, iar dintre toate tehnicile DFT, cea mai folosită este tehnica scanării (*scan-design*). În tehnica scanării (în [Marc93] i se spune tehnică sau metodă SCAN), bistabilele sunt conectate într-un circuit serial și sunt folosite ca terminale de tip I/O. Prin aplicarea tehnicii de scanare putem transforma un circuit secvențial în unul combinațional, făcând astfel mai simplă generarea *pattern*-urilor de test pentru circuit. De asemenea, se poate partiționa logic circuitul în câteva subcircuite și se pot genera, pentru fiecare în parte, *pattern*-urile de test.

Firma NEC a aplicat cu succes tehnica proiectării de scanare la sistemele comerciale de calcul încă din anul 1968. Tehnica numită cale de scanare (*scan-path*) a fost implementată de la nivel de capsulă la nivel de sistem și a devenit o ustensilă puternică pentru testarea și diagnoza sistemelor de calcul VLSI. Conceptul *scan in/scan out*, care stă la baza tehnicii de scanare, a fost introdus pentru prima dată în anul 1964 de către Carter s.a. pentru dezvoltarea testelor de localizare a defecțiunilor la sisteme IBM 360. NEC a dezvoltat metoda căii de scanare pentru a rezolva problemele din anii 70, cum ar fi utilizarea intensivă a MSI/LSI din sistemele comerciale de calcul, precum și numărul mare de porți aflate pe plăcile acestor sisteme.

Metoda căii de scanare pare să fie o soluție promițătoare pentru rezolvarea acestor probleme deoarece permite mai multă controlabilitate și observabilitate asupra circuitului supus testării. În plus, se pot testa atât circuitele combinaționale cât și cele secvențiale. Conceptele de controlabilitate și observabilitate se definesc în felul următor [Lala85], [Prad86], [John89]:

- **Controlabilitatea** arată cât de ușor se poate produce un semnal arbitrar, valid la intrările unei componente (subcircuit) prin excitarea intrărilor primare ale circuitului. De fapt, ea reprezintă abilitatea de a comanda un semnal doar prin intermediul intrărilor primare. O linie care poate fi poziționată pe „1” logic se numește 1-controlabilă, altă linie care poate fi poziționată pe „0” logic se numește 0-controlabilă, iar linia care poate avea ambele stări logice se numește complet controlabilă;
- **Observabilitatea** arată cât de ușor se poate determina la ieșirile primare ale circuitului ceea ce se întâmplă la ieșirile unei componente (subcircuit). Cu alte cuvinte, ea este abilitatea activării unei căi de la semnalul cercetat până la un punct măsurabil.

Prin componente se înțeleg fie circuite integrate standard (SSI și MSI) pentru circuite la nivel de placă, fie celule standard de module de bibliotecă pentru circuite LSI și VLSI. De asemenea, se presupune că mai multe componente sunt conectate cu legături unidirecționale.

Circuitul testat prin metoda căii de scanare poate fi partiționat logic în mai multe subcircuite, iar programele de test pot fi generate, independent pentru fiecare subcircuit. Costul total al întregului circuit scade la α^2 / n , unde n este numărul de subcircuite, iar α este rata de suprapunere. Tabelul următor arată o comparație între mărime și cost.

| Denumire | Mărime | Cost |
|------------------|--------------|------------------|
| Circuit original | 1 | 1 |
| Subcircuit | α / n | $(\alpha / n)^2$ |
| Circuit expandat | α | α^2 / n |

Figura 1-1

Evident că α depinde de n și, din acest motiv, proiectanții trebuie să decidă cu atenție numărul de subcircuite necesare. Pentru un circuit bine proiectat, care este pregătit pentru partiționare, valoarea lui α este sub 1,7. Calea de scanare poate facilita localizarea defecțiunii pentru că procesorul de diagnoză poate să urmărească ușor starea internă a sistemului. *Hardware*-ul adițional este foarte puțin (de la 4% la 10%) și este potrivit pentru un calculator VLSI deoarece necesită relativ puțini pini adiționali de I/O.

Bistabilul transformat cu cale de scanare, utilizat ca și circuitul de bază, include două tipuri de semnale de *clock*, C_1 pentru operația normală și C_2 pentru operația de deplasare. Astfel, bistabilele unei plăci logice sunt conectate în serie printr-o cale de scanare și operează ca un registru serial de deplasare utilizând *Clock II*, intrarea de test (*scan in*) și ieșirea de test (*scan out*). Utilizând o adresă de selecție a plăcii logice, putem selecta o anumită placă dintr-o unitate logică.

Chiar dacă tehnica scanării este o tehnică DFT puternică, pentru testarea logicii aleatoare, aplicarea ei poate pune probleme în diferite tipuri de circuite. Un astfel de tip sunt șirurile de memorie încorporate în circuite de logică aleatoare. Aplicarea directă a tehnicii de scanare la fiecare

celulă de memorie din șir costă foarte mult. O soluție posibilă este aplicarea tehnicii de scanare la un anumit cuvânt din șirul de memorie și accesarea cuvântului fixat în timpul testării circuitului. Oricum, din moment ce, în timpul testării, memoria funcționează ca un registru, ea trebuie testată utilizând altă procedură.

Tehnica scanării s-a dovedit foarte eficace pentru proiectarea circuitelor bipolare, dar este greu de aplicat la un circuit MOS și mai ales la un circuit MOS dinamic. Dificultatea apare, în primul rând, din cauza *hardware*-ului adițional. Un bistabil MOS dinamic este limitat la câteva tranzistoare și supraîncărcarea *hardware*-ului adițional este mare pentru un circuit VLSI MOS.

În 1977, IBM introduce tehnica LSSD (*Level Sensitive Scan Design*) care include tehnica scanării și, în plus, impune ca toate schimbările de stare să fie controlate de nivelul semnalului de ceas și nu de front (*Level Sensitive Design*) [Russ85], [Lala85], [Yarm90]. Această abordare reduce dependența funcționării de timpii de propagare, eliminând cursele sau hazardul. Celula de bază este elementul de memorare a informației, dependent de nivel, care asigură și tratarea semnalului de scan în modul test, element numit SRL (*Shift Register Latch*).

1.3.6.1.2 Caracteristicile metodelor de proiectare structurată la nivel de registru

Deși s-au scris foarte multe despre standardul de test *boundary-scan* (IEEE 1149.1 *Standard Test Acces Port and Boundary-Scan Architecture*), doar câtorva dispozitive li s-a implementat această arhitectură. Acordarea completă cu acest standard, ceea ce înseamnă că proiectanții ar putea să folosească orice dispozitiv, cu asigurarea că este compatibil cu *boundary-scan*, va apare numai după câțiva ani. Până atunci, plăcile cu cablaj imprimat vor fi hibride, conținând circuite (dispozitive) scanabile (după standardul IEEE) și dispozitive nescanabile. Aceste plăci necesită metodologii de testare bazate pe ambele tehnici: pat de cuie (*bed of nails*) și *boundary-scan*.

Testarea cu *boundary-scan* poate micșora drastic timpul de dezvoltare al sistemelor complexe. De asemenea, ea oferă posibilitatea testării sistemelor care utilizează componente cu montare la suprafață de mare densitate și plăci multistrat complexe. Astăzi, aceste sisteme se pot testa cu ajutorul tehnicii patului de cuie, dar există probleme din cauza geometriei micșorate a pinilor, geometrie impusă de tehnologia montării la suprafață. Și în acest caz, având unelte potrivite și dispozitive de

susținere pentru *boundary-scan*, toată placa poate fi testată complet, folosind numai calea *boundary-scan*.

Pe de altă parte, patul de cuie și, în general, metodele de testare în circuit (*in-circuit*) oferă câteva avantaje importante spre deosebire de *boundary-scan*. Un astfel de avantaj cheie al testării în circuit este abilitatea diagnozei deteriorărilor multiple la o singură trecere prin test. În plus, detectarea scurt-circuitelor poate fi realizată înainte de punerea sub tensiune a plăcii, reducând posibilitatea deteriorării catastrofale cauzate de acestea între pini și tensiunea de alimentare. De asemenea, testele în circuit permit realizarea simultană a testării parametrice cu cea digitală.

Boundary-scan, în general, nu se referă la testarea dispozitivelor analogice sau pasive. Contrar impresiei larg răspândite, *boundary-scan* este compatibil cu testarea în circuit și folosind o combinație a celor două se poate simplifica semnificativ testarea plăcilor.

Dispozitivul XC4005 FPGA este primul dintr-o familie de circuite compatibile cu standardul *boundary-scan*. Fiecare pin de I/O al acestui dispozitiv poate fi complet controlat și observat prin utilizarea unor *pattern*-uri de date introduse serial în registrul de *boundary-scan* al lui XC4005 prin pinul TDI (*Test Data Input*) și prin controlul prin pinii TMS (*Test Mode Select*) și TCK (*Test Clock*) proveniți din canalele ATE (*Automatic Test Equipment*). Acest lucru este echivalent cu stimularea complexă a intrărilor pentru a obține aceeași acoperire de test fără *boundary-scan*.

1.3.6.1.3 Proiectarea structurată la nivel de bloc

Încorporarea facilităților de test (BIT, *Built-In Test*) poate micșora considerabil timpul de întreținere și costurile pentru sistemele electronice complexe prin accelerarea detecției și izolării defecțiunilor. La nivel de placă, testul încorporat (BIT) se realizează, de obicei, prin aplicarea unui set de *pattern*-uri pseudo-aleatoare pentru a se verifica existența defecțiunilor și prin comprimarea informației răspunsului de test pentru a se obține o semnătură. Aceste tehnici sunt eficace pentru plăcile utilizate în scopuri generale atât timp cât seturile de cipuri care se testează reacționează la cele două tehnici, dând rar semnături identice pentru plăci bune sau defecte.

Notă: În [D&TR89], McCluskey precizează că nu este recomandată utilizarea expresiei *test-response compression* ca un sinonim al expresiei *test-response compaction*, deoarece prin compactare se pierde din informație, în timp ce prin comprimare se elimină doar redundanța, fără

pierderea informației. În majoritatea literaturii de specialitate pe care am consultat-o - [Lala85], [John89], [Yarm90], [Fuji90a] - nu este respectată această opinie. Din această cauză, precum și datorită faptului că în limba română, când este vorba despre reducerea volumului secvențelor binare răspuns, este consacrat termenul de comprimare [Vlad82], [Vlad89], în continuare vom folosi denumirea „comprimare“.

1.3.6.1.3.1 Autocontrolul bazat pe principiul BILBO

Tehnica de determinare a stării unui chip sau a unei plăci prin folosirea informației comprimate a răspunsului de test se numește analiza semnăturii. În general, metodele de analiză a semnăturii pentru testarea plăcilor utilizează circuite cu observatori logici, de bloc, încorporați (BILBO, *Built-In-Logic Block Observers*) pentru a genera *pattern*-uri pseudo-aleatoare și pentru a comprima răspunsurile. Deci, BILBO este o schemă încorporată de generare a testului care utilizează, în conjuncție *scan*, design-ul cu analiza de semnături. Componenta de bază a acestei tehnici de autotest este un registru de deplasare numit registru BILBO care are patru moduri de funcționare.

Această tehnică elimină nevoia generării și aplicării *pattern*-urilor de test din exterior. Totuși, este nevoie de simularea defecțiunilor pentru a se determina acoperirea de defecțiuni, realizată de secvențe de vectori pseudo-aleatorii. Mai mult decât atât, circuitul trebuie să fie simulat pentru a se obține valorile semnăturilor în cazul funcționării corecte. Cantitatea datelor de test care trebuie stocate și analizate este redusă din moment ce răspunsul circuitului la N *pattern*-uri de test este comprimat la un singur cuvânt.

1.3.6.1.3.2 Modulul comutator de testare, încorporat, pentru izolarea defectelor

În unele aplicații, în timpul operației normale, circuitele BILBO supraîncălzesc în mod nejustificat, schema din punct de vedere al consumului de energie, iar izolarea defecțiunilor durează, de obicei, mult. În locul registrului BILBO, s-a propus [KaVa89] o altă variantă de analiză a semnăturii, care utilizează un modul numit comutator de testare SW (*testing switch*). Modulele SW pot îmbunătăți, la nivel de placă, timpul de testare și supraîncărcarea de putere a schemelor de analiză a semnăturii implementate utilizând circuitele BILBO.

În multe probleme de izolare a defecțiunii este mai convenabilă izolarea defecțiunii la un grup de cipuri decât la un singur circuit. Aceste grupuri se numesc grupuri de ambiguitate (AG). Modulul SW generează

pattern-uri de test pseudo-aleatorii pentru a fi aplicate la intrările unui set AG, în timp ce, simultan, comprimă datele răspunsului de test de la ieșirile unui alt set AG. Această caracteristică poate minimiza total al testului de izolare a defecțiunii.

Utilizarea modulelor SW pentru izolarea defecțiunilor introduce întârzieri adiționale între AG-uri, în timpul modului normal de operare. Totuși, dacă proiectanții implementează cipuri cu comutatorul de testare utilizând tehnici de împachetare și tehnologii avansate, pot minimiza întârzierile adiționale. Se pare că, în cazul multor aplicații, această întârziere va avea un efect neglijabil asupra performanțelor sistemului.

1.3.6.1.4 Controlul erorilor la nivel de procesor

În anul 1972, Pradhan și Reddy au arătat că, utilizând redundanța de ordinul duplicării, codul de paritate poate realiza cu eficacitate controlarea erorilor. Afirmatia că numai duplicarea poate realiza corecția erorilor în procesare a fost pusă în practică, parțial, în calculatorul (4,2). Conceptul este deja utilizat ca un produs comercial prin sistemul de comutare Philips S2500 [Krol86].

S2500 utilizează patru perechi de procesoare-memorie, fiecare procesor fiind pe 16 biți, iar memoria pe 8 biți. De aceea, există redundanță de ordinul patru pentru procesoare și de ordinul doi pentru memorie. Ieșirea pe 16 biți a fiecărui procesor este codificată într-o versiune $(4,2)$ $GF(2^n)$ de cod R-S, producând un cod de 32 de biți. Acest cod poate corecta orice bloc singular de 8 biți. Fiecare codificator al ieșirii procesorului produce doar 8 biți din cei 32 de biți ai cuvântului de cod. Memoria asociată a celui de-al i -lea procesor stochează al i -lea octet. De aceea, codificatorul pentru procesorul i produce 8 biți care corespund celui de-al i -lea octet. La operațiunea de citire, toți cei patru octeți sunt aduși din memorie și decodificați de către decodificator, care poate corecta orice octet singular eronat sau orice eroare de doi biți în doi octeți diferiți. Este de notat faptul că deteriorarea unui singur procesor sau a unei singure memorii afectează doar un singur octet din cuvântul de cod. Din acest motiv, sistemul poate tolera deteriorarea oricărei perechi singulare procesor-memorie. În plus, codul poate corecta ștergerea unui octet singular precum și erorile singulare de bit. Astfel, dacă se demontează o pereche procesor-memorie pentru reparație, sistemul continuă să opereze utilizând decodificarea de erori-ștergeri (*error-erasure decoding*). La nivel de bit poate fi corectată orice subsecvență de eroare singulară.

Această utilizare a codificării pentru controlarea erorii de procesor este atractivă. Singurul dezavantaj real este acela că, datorită necodificării liniilor de adrese la memorie, nu are loc corecția de eroare la operația de scriere. În general, prin utilizarea căilor netradiționale descrise mai sus, codul controlului de paritate poate furniza un control eficace al erorilor procesorului. Utilizarea lor va elimina necesitatea de conversie a codului și întârzierea asociată în transfer procesor-memorie, furnizând astfel controlul uniform al erorii [Krol86], [Fuji90b].

1.3.6.2 Problematika construcției *checker*-elor de erori

Multe sisteme digitale includ detectoare de erori, adică circuite care detectează și semnalizează apariția erorilor. Există două motive principale pentru includerea detectoarelor de erori într-un sistem și aceste motive sunt: prevenirea ajungerii erorii până la ieșirea sistemului și localizarea poziției erorii. O cale pentru a preveni propagarea erorii este verificarea datelor recepționate după transferul de date și, dacă se detectează o eroare, atunci fie se cere retransmisia, fie se corectează datele utilizând un cod corector de erori (ca în memoriile RAM). O altă cale este verificarea rezultatului unei operații aritmetice și repetarea operației dacă rezultatul conține o eroare detectată. Sistemele mari includ, de obicei, detectoare de erori împreună cu circuite pentru înregistrarea locației și frecvenței evenimentelor de eroare [Mccl90]. Această informație este utilizată pentru a facilita întreținerea rapidă și exactă.

1.3.6.2.1 Circuite de verificare dublă cale

Un circuit de verificare dublă cale (*two-rail checker*), adică un circuit care controlează dacă fiecare pereche de intrări are valori complementare, este utilizat pentru a converti cele n perechi de semnale într-o singură pereche de semnale complementare dacă și numai dacă toate cele n perechi de intrare au semnale complementare.

1.3.6.2.1.1 *Checker* de egalitate

Un tip foarte important de *checker* este cel care verifică egalitatea, adică acela care compară două cuvinte de intrare pentru a determina dacă biții corespunzători din ambele cuvinte au aceeași valoare. *Checker*-ul de egalitate este componenta cheie în *checker*-ele pentru coduri separabile și pentru compararea ieșirilor circuitelor duplicate.

Un *checker* de egalitate cu autotestare poate fi obținut prin complementarea tuturor biților din unul dintre cuvintele care trebuie să fie comparate pentru a forma un cod dublă cale (*two-rail code*).

Checker-ele dublă cale discutate mai sus pot fi utilizate pentru a realiza un *checker* testabil de egalitate. O altă posibilitate este aceea de a concatena cuvântul complementat și cel necomplementat pentru a forma un cuvânt de cod care are jumătate din biții lui egali cu 1. Un checker k din $2k$ poate fi utilizat pentru a forma un *checker* testabil de egalitate.

1.3.6.2.1.2 *Checker M din N*

Checker-ele M din N formează o altă clasă de *checker*-e încorporate. Deși sunt mai puțin importante decât celelalte, totuși sunt utilizate. *Checker*-ul 1 din n este folosit pentru a monitoriza operarea corectă a rețelelor complete de decodificare. O implementare cu autotestare a acestui *checker* este formată prin conectarea tuturor ieșirilor decodificate care corespund intrărilor cu paritate pară la o poartă SAU și tuturor ieșirilor care corespund intrărilor cu paritate impară la o altă poartă SAU. Ieșirile celor două porți au fie 10, fie 01, atunci când decodificatorul operează corect. Schema este cu autotestare deoarece porțile SAU recepționează toate configurațiile de intrare (00, 10, 01) necesare testării lor pentru defecțiuni singulare de tip punere pe (*stuck-at*) în timpul operării normale. Această schemă poate fi convertită într-un circuit testabil cu o singură ieșire.

Codurile cu ponderea fixă m -din- n se utilizează atunci când se dorește atât detectarea tuturor defectelor singulare de tip *stuck-at*, cât și a celor unidireționale.

Codul zecimal 2 din 5 este cel mai cunoscut exemplu al unui cod M din N . Codurile k din $2k$ sau k din $2k+1$ sunt implementările uzuale deoarece conțin numărul maxim de cuvinte de cod pentru o lungime cunoscută a cuvântului.

1.3.6.2.2 *Circuite cu autoverificare totală, încorporate*

În 1968, Carter și Schneider au ridicat obiectivele autoverificării (*self-checking*) și au dezvoltat circuitul TSC (*Totally Self-Checking*), care are mai multe intrări și ieșiri. În timpul operării normale (fără defecțiuni) un circuit TSC nu recepționează toate combinațiile posibile la intrări, ci numai un subset al tuturor combinațiilor posibile numit spațiul codului de intrare (*input-code space*). De asemenea, el nu produce toate combinațiile posibile de ieșire, ci numai un subset numit spațiul codului de ieșire (*output-code space*).

Checker-ul TSC este un circuit combinațional care simplifică sarcina unui observator. Dacă unitatea funcțională TSC produce o ieșire care aparține spațiului codului de ieșire, *checker*-ul TSC indică faptul că nu

există nici o eroare. O ieșire care nu aparține spațiului codului de ieșire se numește ieșire ilegală (*noncode*) și, dacă *checker*-ul găsește o astfel de ieșire la unitatea funcțională, atunci indică existența unei erori.

O situație interesantă apare atunci când *checker*-ul are o defecțiune. Din moment ce se presupune că unitatea are numai o singură defecțiune la un moment dat, se poate considera că, dacă *checker*-ul are o defecțiune, unitatea funcțională nu are nici o defecțiune și, deci, nu va produce nici o ieșire ilegală. Astfel, este nevoie de un mecanism prin care fiecare defect din interiorul *checker*-ului să se releve ca o eroare la ieșirea lui. Este necesar, de asemenea, să se codifice ieșirile *checker*-ului printr-un cod detector de erori unidirecționale. *Checker*-ul trebuie să producă un număr cât mai mic de ieșiri deoarece acestea vor fi monitorizate din exterior. Cel mai mic cod nesistematic cunoscut, numit și cod dublă cale (*dual-rail code*) este $\{01, 10\}$. O eroare unidirecțională transformă acest cod în $\{00$ sau $11\}$ și astfel ieșirile *checker*-ului vor indica o eroare. Deci, se poate detecta o defecțiune/eroare prin observarea unei ieșiri ilegale la un *checker*, iar pentru a detecta o defecțiune în interior este nevoie de aplicarea unor intrări specifice (configurații de biți, vectori) numite intrări de test. Astfel, se presupune, în mod implicit, că toate intrările de test ale unui sistem TSC sunt aplicate în modul operării normale.

În literatura de specialitate se întâlnesc următoarele definiții [Gait85], [Lala85], [NaKa88], [John89], [BuJh94]:

- Un circuit se numește *fault-secure* (FS) pentru un set Φ de defecțiuni, dacă, pentru orice defecțiune singulară $\tau \in \Phi$, circuitul produce fie o ieșire de cod, așteptată pentru intrarea respectivă, fie o ieșire ilegală (*noncode*), dar nu produce un cuvânt de cod neașteptat pentru acea intrare;
- Un circuit este numit *self-testing* (ST) pentru un set Φ de defecțiuni dacă și numai dacă, pentru orice defecțiune singulară $\tau \in \Phi$, circuitul produce o ieșire ilegală (*noncode*) pentru cel puțin o intrare din cod;
- Un circuit este numit *totally-self-checking* (TSC) dacă este atât *self-testing*, cât și *fault-secure* pentru orice defecțiune singulară $\tau \in \Phi$;
- Un circuit se definește drept *code-disjoint* (CD) dacă în absența defectelor, aplicând intrărilor primare un vector binar și valid al codului, la ieșire se obține, de asemenea, un vector binar valid al respectivului cod, așa cum, dacă aplicând intrărilor primare un vector

binar ilegal pentru codul dat, la ieșire se obține, de asemenea, un vector binar ilegal pentru al respectivului cod. Un circuit TSC este numit TSC *checker* dacă este *code-disjoint*.

Notă: Pentru circuitele de verificare nu este esențială proprietatea de *fault-secure*.

Inițial, definiția circuitelor TSC a fost dată pentru circuitele combinaționale, extinzându-se, ulterior, atât asupra circuitelor secvențiale [Gait85], cât și, în general, asupra unor sisteme întregi (despre care se poate spune că se comportă ca și mașinile secvențiale).

Revenind la *checker*-ul TSC, se punctează faptul că ieșirea lui devine \$\$ sau 11, dacă și numai dacă o eroare este detectată la ieșirea unei unități funcționale sau dacă *checker*-ul are o defecțiune *hardware*. Deci, un *checker* TSC trebuie să fie *code-disjoint*.

Adeesea există situații în care, chiar dacă un *checker* este proiectat pentru *code-disjointness*, nu sunt disponibile de la unitatea funcțională TSC toate intrările de test, necesare pentru testarea lui. Pentru a face un circuit atât *self-testing*, cât și *code-disjoint* ar trebui să se suspende, din când în când, operația normală pentru a alimenta (*to flush*) *checker*-ul cu intrări de test care în mod normal nu sunt disponibile. Dar, în această situație sistemul nu mai este *on line* și există dificultăți practice în alimentarea *checker*-ului. O cale mai bună este proiectarea unor *checker*-e care să fie atât *self-testing*, cât și *code-disjoint* și care să nu aibă nevoie de acces din exterior, fiind ușor de încorporat.

În conformitate cu cele discutate anterior, un circuit TSC, ale cărui ieșiri sunt codificate într-un cod detector de erori, produce întotdeauna un cuvânt ilegal ca primă ieșire eronată din cauza unei defecțiuni. Acest comportament este foarte avantajos și este cunoscut sub denumirea de obiectivul proiectării TSC (*TSC goal*). Există, totuși, și alte clase de circuite logice care îndeplinesc acest obiectiv.

Proprietatea de *strongly fault-secure* caracterizează o clasă mai largă de circuite logice care realizează acest obiectiv sub aceleași premise privitoare la defecțiunile care se iau în considerație. Proprietatea se definește în felul următor: un circuit (sistem) este *strongly fault-secure* (SFS) pentru un set Φ de defecțiuni, dacă pentru orice defecțiune singulară $\tau \in \Phi$, circuitul poate fi:

- TSC pentru $\{\tau\}$ sau
- FS pentru $\{\tau\}$ și, dacă defectul τ apare, atunci circuitul (sistemul) rezultat este tot SFS pentru $\Phi - \{\tau\}$.

Din punct de vedere al proiectantului, este mai avantajoasă proiectarea circuitelor SFS care îndeplinesc obiectivul TSC spre deosebire de proiectarea circuitelor TSC, pentru că cele din urmă impun proiectării mai multe restricții. De asemenea, chiar dacă proiectarea sistemelor *fail-safe* a fost studiată înainte și independent de cea a circuitelor TSC și SFS, circuitele *fail-safe* sunt echivalente cu circuitele *fault-secure* cu o anumită codificare.

O altă proprietate interesantă se definește după cum urmează: un circuit este *strongly code-disjoint* (SCD) pentru un set Φ de defecțiuni, dacă pentru orice defecțiune singulară $\tau \in \Phi$, circuitul este fie:

- CD și este ST pentru $\{\tau\}$, fie
- CD și, dacă defectul τ apare, atunci circuitul rezultat este tot SCD pentru $\Phi - \{\tau\}$.

Combinăția dintre un circuit TSC (care realizează o anumită funcție) și un *checker* TSC aparține clasei rețelelor logice cu componente TSC. Aceste rețele îndeplinesc obiectivul TSC sub presupunerea defecțiunilor singulare, ceea ce înseamnă că defecțiunile apar una câte una, iar între două defecțiuni este timp suficient astfel încât toate intrările din cod să fie aplicate rețelei. Presupunând că există defecțiuni duble, rețelele de mai sus nu îndeplinesc obiectivul TSC. Prin defecțiuni duble se înțelege că între apariția a două defecțiuni care afectează circuitul care realizează funcția dorită sau între apariția a două defecțiuni care afectează *checker*-ul, există suficient timp pentru ca toate intrările din cod să fie aplicate rețelei. În [Gait88c], Gaitanis prezintă o nouă clasă de circuite de verificare TSC cu proprietăți mai puternice de autoverificare (*self-checking*) pentru a îndeplini obiectivul TSC în cazul rețelelor de *checker*-e, când se iau în considerație defecțiunile duble.

Un TSC *checker* este *strongly-self-checking* (SSC) pentru un set Φ de defecțiuni, dacă F_c , F_e și E sunt submulțimi disjuncte ale mulțimii de ieșiri ilegale [Gait88c], unde:

- F_c este mulțimea ieșirilor neașteptate, cauzate de defecțiuni interne $\tau \in \Phi$ și intrări din cod ale *checker*-ului TSC;
- F_e este mulțimea ieșirilor neașteptate, cauzate de defecțiuni interne $\tau \in \Phi$ și intrări (ilegale) din afara codului ale *checker*-ului TSC;
- E este mulțimea ieșirilor produse de *checker*-ul TSC fără defecțiuni cauzate de intrările ilegale.

1.4 Concluzii

Pregătind aplicarea unor metode care să asigure implementarea eficientă a toleranței la defectare în subsistemele I/O, în acest prim capitol de consistență al lucrării au fost supuse unei analize sistematice aspectele de defectare specifice domeniului abordat care să constituie obiective de combătut prin metodele care vor fi dezvoltate în următoarele capitole.

Rezumând, în acest capitol sunt cuprinse următoarele contribuții:

- 1) un aport în sistematizare metodelor de creștere a fiabilității și disponibilității prin autocontrol implementat atât *hardware*, cât și *software*
- 2) definirea unor noi indicatori de fiabilitate specifici matricelor de discuri, necesari pentru evaluarea structurilor de tip RAID și RAIC, în special RAID7 și RAIC 7.

2 ANALIZA SUBSISTEMULUI DE DISCURI DIN PUNCT DE VEDERE AL FIABILITĂȚII

În multe sisteme de calcul, subsistemul I/O reprezintă adesea o pondere importantă în cadrul performanțelor sistemului și din punct de vedere al fiabilității reprezintă zona cea mai vulnerabilă, defectarea acestuia putând avea efecte dezastruoase prin pierderea definitivă a datelor.

Dintre componentele subsistemului I/O, discurile magnetice prezintă cel mai ridicat grad de risc în utilizare (componente mecanice în mișcare, etc.).

De asemenea, tot ele sunt cele care au o importanță decisivă, în performanță globală a unui sistem.

Creșterea performanțelor și a fiabilității discurilor este limitată tehnologic, această limită neputând fi depășită decât prin utilizarea redundanței informaționale și structurale.

Pornind de la metodele generale de creștere a fiabilității sistemelor de calcul, în acest capitol sunt analizate metodele particulare dedicate discurilor cum ar fi: coduri specifice detectoare și corectoare de erori, precum și structuri redundante de discuri RAID. În final, se prezintă o clasificare, din punct de vedere al fiabilității, pentru structurile redundante de discuri.

2.1 Coduri pentru discuri

De mai mulți ani, discurile magnetice joacă un rol important în memorarea rapidă și de mare capacitate a fișierelor în sistemele computerizate. Încă din anii '60, tehnologia discurilor magnetice s-a îmbunătățit continuu din punct de vedere al densității, capacității de stocare, vitezei de transfer, costului și fiabilității. De exemplu, în ceea ce privește densitatea spațială de memorare, au avut loc îmbunătățiri de patru ordine de mărime — de la 10^3 la 10^7 biți pe țol pătrat.

Viteza de transfer a datelor este mai mare în cazul discurilor magnetice, decât în cazul benzilor magnetice, în special datorită vitezei mai mari de rotație. De aceea, logica sa trebuie să admită o mai mare toleranță. În etapele incipiente ale dezvoltării discurilor se foloseau doar verificările simple ale parității pentru a testa integritatea datelor. În etapa următoare s-a făcut uz de coduri polinomiale pe 16 biți pentru detectarea erorilor salvă. În anii '70 discurile de înaltă densitate și de mare viteză (de exemplu, IBM 3330) au pus bazele unei noi tehnologii. Acest tip de disc

folosea codurile *Fire* pentru a avea capacitatea detecției și corecției erorilor salvă. În proiectele recente, acest tip de cod este înlocuit cu un cod întrețesut *Reed-Solomon* pentru corecția erorilor. Pe lângă codurile detectoare și corectoare de erori se mai folosesc și alte tehnici recuperatorii pentru a mări fiabilitatea și a crește securitatea informațiilor, cum ar fi: ignorarea defectelor, blocuri alternative de date și recitirea .

În sistemele de discuri magnetice predomină erorile de tip salvă. Majoritatea acestora este legată de imperfecțiunile suprafeței discurilor sau de neregularitățile suprafeței acestora. Restul erorilor este, de regulă, datorat capetelor magnetice, susceptibile de erori aleatoare induse de zgomote [Howe84].

Discurile optice reprezintă o tehnologie relativ nouă pentru înregistrarea informațiilor. Fiecare disc este acoperit cu o peliculă reflectorizantă. Citirea și scrierea sunt realizate prin intermediul unui fascicul laser. Mai precis, citirea se reduce la a stabili reflectivitatea unei anumite zone a discului, în vreme ce la scriere un laser de putere mai mare realizează, prin topire, adâncituri în pelicula reflectorizantă. Densitatea informațiilor pe un disc optic este mult mai mare decât în cazul discurilor magnetice, de vreme ce tehnologiile optice permit o „rezoluție“ mai mare. De aceea, se estimează că acest mediu de stocare va avea un rol tot mai important în sistemele de memorare.

Erorile sunt datorate, în principal, discului optic însuși, de vreme ce nu se pot încă fabrica la un preț rezonabil pelicule reflectorizante de înaltă stabilitate și înaltă fiabilitate. Aproape toate erorile au drept cauză imperfecțiunile discului, iar restul sunt datorate variațiilor de focalizare la înregistrare sau zgomotului aleator la citire. Ca atare, sunt necesare atât facilități pentru corecția erorilor salvă, cât și pentru corecția celor aleatoare [SaTN86, IYH187]. În prezent se folosesc tehnici puternice de corecție a erorilor, cum ar fi codul dublu codificat *Reed-Solomon* sau o combinație a codului *Reed-Solomon* cu un cod ciclic, pentru a obține același nivel de fiabilitate ca în cazul discurilor magnetice.

2.1.1 Codul *Fire* pentru memoriile bazate pe disc magnetic

Codul *Fire* [Fire59], care este capabil de corecția oricărei salve de lungime l sau mai mică și de detectarea simultană a oricărei salve de lungime $d \geq l$, este generat de următorul polinom generator:

$$g(x) = (x^l + 1) p(x)$$

Aici $p(x)$ este un polinom ireductibil de grad m cu exponentul e și $c \geq l + d - 1$. Lungimea n a acestui cod este egală cu cel mai mic multiplu comun (CMMMC) dintre e și c :

$$n = \text{CMMMC}(e, c)$$

Numărul biților de verificare a acestui cod este $(c + m)$.

Fie $b(x)$ polinomul erorii salvă, al cărui grad este mai mic sau egal cu $l - 1$. Atunci $r(x)$ are expresia:

$$r(x) = v(x) + x^l \cdot b(x)$$

unde i arată poziția de începere a erorii salvă. Dacă $b(x) = 0$, atunci $r(x)$ poate fi divizat atât de $(x^c + 1)$, cât și de $p(x)$, reziduul fiind, deci, egal cu zero. Dacă $b(x) \neq 0$, atunci registrele de deplasare vor conține următoarele reziduuri:

$$x^i \cdot b(x) = s_1(x) \pmod{(x^c + 1)} \quad (1)$$

$$x^i \cdot b(x) = s_2(x) \pmod{p(x)} \quad (2)$$

Folosind aceste reziduuri, $s_1(x)$ și $s_2(x)$, trebuie să obținem $b(x)$ și valoarea lui i . După $(n - l)$ deplasări ale registrelor de deplasare, acestea vor conține următoarele reziduuri:

$$\begin{aligned} s_1(x) \cdot x^{n-l} &= x^n \cdot b(x) \\ &= b(x) \pmod{(x^c + 1)} \\ s_2(x) \cdot x^{n-l} &= x^n \cdot b(x) \\ &= b(x) \pmod{p(x)} \end{aligned}$$

Cu alte cuvinte, atunci când conținutul părții inferioare l a registrului tipului de eroare este egal cu cel al registrului localizării erorii după δ deplasări ale conținutului acestor registre, tipul erorii salvă, exprimat sub forma lui $b(x)$, poate fi introdus în aceste registre. Pe lângă acest polinom $b(x)$, poziția i a erorii se poate calcula din δ (care este egal cu $n - i$), fiind, deci, $i = n - \delta$.

Această metodă de decodificare necesită un timp mai mare atunci când i este mic în comparație cu n . De aceea, vom lua în considerație metoda decodificării de înaltă viteză [Chie69]. Să presupunem că c și e sunt prime unul față de altul.

În această metodă de decodificare, cuvântul recepționat $r(x)$ este introdus simultan în registrul tampon, registrul tipului de eroare și

registru localizării erorii. Apoi, se obțin reziduurile exprimate în ecuațiile 1 și 2 în registrul tipului erorii și, respectiv, în registrul localizării erorii. În acest moment, putem deduce cuvântul recepționat, pe baza reziduurilor astfel obținute, astfel:

$$s_1(x) = s_2(x) = 0 \rightarrow r(x) : \text{fara eroare}$$

$$\left. \begin{array}{l} s_1(x) \neq 0 \text{ and } s_2(x) = 0 \\ s_1(x) = 0 \text{ and } s_2(x) \neq 0 \end{array} \right\} \rightarrow r(x) : \text{eroare \cdot necorectabila}$$

Dacă $s_1(x) \neq 0$ și $s_2(x) \neq 0$, atunci decodificarea se poate face conform următorului algoritm:

- 1) Deplasează registrul tipului erorii până când în registrul părții superioare ($c - l$) sunt doar zerouri și memorează numărul de deplasări δ_0 . Dacă respectivul conținut este diferit de zero după ($c - l$) deplasări, atunci a apărut o eroare necorectabilă și, deci, decodificarea se oprește;
- 2) Deplasează registrul localizării erorii până când conținutul său este egal cu partea inferioară l a registrului tipului de eroare și memorează numărul de deplasări δ_l . Dacă valoarea părții inferioare a registrului tipului de eroare nu este egală cu cea a registrului localizării erorii după ($e - l$) deplasări, atunci se presupune că a apărut o eroare necorectabilă și, deci, procesul va fi oprit;
- 3) Poziția i a erorii se poate calcula din valorile δ_0 și δ_l , eroarea corectându-se, deci, prin folosirea tipului de eroare $b(x)$. Cu alte cuvinte, conținutul părții inferioare l a registrului tipului de eroare denotă polinomul $b(x)$.

În algoritmul de mai sus, poziția i a erorii se poate calcula folosind Teorema Restului Chinezesc, astfel încât

$$\begin{aligned} \delta_0 &= c - i \\ &= -i \pmod{c} \\ \delta_l &= e - i \\ &= -i \pmod{e} \end{aligned}$$

după care $-i = \delta_0 \cdot A_0 \cdot e + \delta_l \cdot A_l \cdot c \pmod{n}$, unde A_0 și A_l sunt întregi care satisfac relația:

$$A_0 \cdot e + A_l \cdot c = 1 \pmod{n}$$

Calcularea lui A_0 și A_l se poate face cu ușurință *off-line*, cu numerele $A_0 \cdot e$ și $A_l \cdot c$ memorate în calculatorul poziției erorii. În cazul acestei decodificări, numărul maxim de deplasări este $(e + c - 2)$, în timp ce la

decodicarea anterioară era $(e \cdot c - 1)$, ceea ce face acest decodificator mult mai rapid decât cel precedent.

În continuare, ne vom concentra asupra unei noi clase de coduri care este o generalizare a codurilor *Fire* [Chie69]. Acestea sunt generate de un polinom de forma

$$g(x) = (x^c - 1) \prod_{i=1}^h p_i(x) \quad (3)$$

unde $p_i(x)$ este un polinom ireductibil distinct, având exponentul e_i și gradul r_i . Se mai presupune că exponenții e_i , unde $i=1, 2, \dots, h$, nu divid c . Lungimea codului este

$$n = \text{CMMC}(c, e_1, e_2, \dots, e_h)$$

Teorema 1 [Chie69]. *Codul generat de $g(x)$, dat de ecuația 3 detectează toate erorile salvă, mai puțin cele cu lungimea mai mică sau egală cu d ; corectează toate salvele $b(x)$ de lungime mai mică sau egală cu l , care sunt prime în raport cu $\prod_{i=1}^h p_i(x)$, cu condiția ca*

$$c \geq l + d - 1 \text{ și } l \leq \sum_{i=1}^h r_i.$$

Ecuțiile care determină poziția i pot fi scrise astfel:

$$-i = \delta_0 \cdot A_0 \cdot \prod_{j=1}^h e_j + c \cdot \sum_{j=1}^h \delta_j \cdot A_j \cdot \prod_{\substack{k=1 \\ k \neq j}}^h e_k \quad (4)$$

$$A_0 \cdot \prod_{j=1}^h e_j + c \cdot \sum_{j=1}^h A_j \cdot \prod_{\substack{k=1 \\ k \neq j}}^h e_k = l \quad (5)$$

Sistemele de discuri compatibile cu IBM 3330 folosesc codul *Fire* generalizat, generat cu ajutorul polinomului următor:

$$g(x) = (x^{22} + 1) \prod_{i=1}^3 p_i(x)$$

unde

$$p_1(x) = 1 + x + x^6 + x^7 + x^{11}$$

$$p_2(x) = 1 + x + x^2 + x^3 + x^4 + x^5 + x^6 + x^7 + x^8 + x^9 + x^{10} + x^{11} + x^{12}$$

$$p_3(x) = 1 + x + x^5 + x^6 + x^7 + x^9 + x^{11}$$

sunt polinoame ireductibile distincte. Exponenții polinoamelor $p_1(x)$, $p_2(x)$ și $p_3(x)$ sunt 89, 13 și, respectiv, 23. Lungimea codului este, deci:

$$n = \text{CMMC}(22, 89, 13, 23) = 585\,442$$

El are $22 + 11 + 12 + 11 = 56$ biți de verificare și este capabil să corecteze orice salvă de lungime 11 sau mai mică.

Din ecuațiile (4) și (5) se pot deduce cele care determină poziția erorii după cum urmează:

$$i = \delta_0 A_0 e_1 e_2 e_3 + \delta_1 A_1 c e_2 e_3 + \delta_2 A_2 c e_1 e_3 + \delta_3 A_3 c e_1 e_2 \quad (6)$$

unde A_0, A_1, A_2 și A_3 satisfac relația:

$$A_0 e_1 e_2 e_3 + A_1 c e_2 e_3 + A_2 c e_1 e_3 + A_3 c e_1 e_2 = -1 \pmod{n} \quad (7)$$

Pentru acest cod putem afla că:

$$\begin{aligned} A_0 e_1 e_2 e_3 &= -1 \pmod{c} \\ A_1 c e_2 e_3 &= -1 \pmod{e_1} \\ A_2 c e_1 e_3 &= -1 \pmod{e_2} \\ A_3 c e_1 e_2 &= -1 \pmod{e_3} \end{aligned} \quad (8)$$

Cei mai mici întregi care satisfac ecuațiile (7) și (8) sunt: $A_0 = 5$, $A_1 = 78$, $A_2 = 6$ și $A_3 = 10$. Acești întregi se introduc în ecuația (6), rezultând următoarea ecuație care ne indică poziția erorii:

$$i = 133\,055 \delta_0 + 513\,048 \delta_1 + 270\,204 \delta_2 + 254\,540 \delta_3 \pmod{585\,442}$$

Operațiile de deplasare ale registrului tipului de eroare se execută până când cele mai semnificative 10 ranguri ale acestuia au toate valoarea 0. Deplasarea registrelor de localizare a erorii se execută până când întregul conținut al tuturor celor trei registre coincide cu cele mai puțin semnificative 11 ranguri ale registrului tipului de eroare. Considerând că numărul acestor deplasări este $\delta_0, \delta_1, \delta_2$ și, respectiv, δ_3 , rezultă $K_0 = 133\,055 \delta_0$, $K_1 = 513\,048 \delta_1$, $K_2 = 270\,204 \delta_2$ și, respectiv, $K_3 = 254\,540 \delta_3$. Putem, acum, obține i ca fiind $i = K_0 + K_1 + K_2 + K_3 \pmod{n}$.

O altă metodă rapidă de corecție, bazată pe codul *Fire* obișnuit, a fost propusă în [Adi84]. Aceasta este o tehnică hibridă care combină capturarea erorilor secvențiale și tehnicile de căutare în tabele. Decodificatorul necesită cel mult $(c - m - 1)$ cicluri de deplasare pentru a găsi atât eroarea de tip salvă, cât și locația ei. Implementarea decodificatorului necesită memorii programabile numai de citire

(*PROM*-uri) și matrice logice programabile (*FPGA*-uri), potrivite implementărilor pe scară largă.

2.1.2 Codul întrețesut R-S SbEC-DbEC pentru memoriile cu discuri magnetice

O schemă de corecție a erorilor de tip salvă, de înaltă performanță, care folosește un cod corector de eroare pe *byte* este dat [Hori75b]. Discurile magnetice apărute mai recent, cum ar fi cele din sistemele IBM 3370 și IBM 3380, adoptă un cod de distanță 4, *Reed-Solomon* (R-S), discutat în [Hodg80]. Acest sistem utilizează un format fix al blocului de date. Fiecare bloc are 512 *bytes*. Pentru controlul erorilor, acestuia i se mai adaugă încă nouă biți. Codul folosit în acest sistem este o prescurtată a codului *Reed-Solomon* de distanță 4, cu simboluri din câmpul $GF(2^8)$. Polinomul generator al codului este:

$$g(x) = (x + 1)(x + \alpha)(x + \alpha^2)$$

unde α este elementul primitiv al lui $GF(2^8)$ și rădăcină a lui $p(x) = 1 + x^2 + x^3 + x^4 + x^8$. Acest cod este capabil să corecteze orice eroare monosimbolică (un *byte*) și să detecteze orice eroare bisimbolică (doi *bytes*).

Pentru codificare, fiecare bloc este subîmpărțit în trei subblocuri: D_{-1} , D_0 și D_1 , fiecare având 171 de *bytes*. Cele trei blocuri se reprezintă în formă polinomială după cum urmează:

$$D_i = D_{i,170} x^{170} + D_{i,169} x^{169} + \dots + D_{i,1} x + D_{i,0}$$

$$i = \{-1, 0, 1\}$$

unde fiecare *byte* $D_{i,k}$, $i = \{-1, 0, 1\}$, $k = \{0, 1, \dots, 170\}$, este privit ca un simbol în $GF(2^8)$. Cei trei *bytes* de verificare pentru D_i , notați C_{-1} , C_0 , C_1 , unde $i = \{-1, 0, 1\}$, se obțin ca mai jos:

$$C_{i,j} = x \cdot D_i(x) \bmod (x + \alpha^j)$$

$$i = \{-1, 0, 1\} \quad j = \{-1, 0, 1\}$$

A se observa diferența dintre această metodă de codificare și cea folosită în cazul codurilor ciclice, unde toți cei trei *bytes* sunt calculați prin împărțirea lui $x^3 \cdot D_i(x)$ la $g(x)$. După ce s-au format cei trei *bytes* de control, aceștia sunt adăugați la D_i pentru a forma un cuvânt codificat W_i .

Cele trei cuvinte codificate W_1 , W_0 și W_{-1} sunt întreșute pentru a forma un bloc codificat. Ca urmare, codul general rezultat este (174,171) R-S, întreșut de gradul 3. Acest cod este capabil să corecteze orice salvă de erori de cel mult trei *bytes* și să detecteze salve eronate de cel mult șase *bytes*.

Atunci când un bloc codificat este citit de pe disc, el este descompus în trei subcuvinte W_r , W_0 și W_{-1} , unde:

$$W_r = (C_{r,1}, C_{r,0}, C_{r,-1}, D_r)$$

$$i = \{1, 0, -1\}$$

Sindromul $S_{i,j}$ ($i = \{1, 0, -1\}$, $j = \{1, 0, -1\}$) se poate obține în felul următor:

$$S_{i,j} = C'_{i,j} + x \cdot D'_i(x) = C'_{i,j} + \alpha \cdot D'_i(\alpha) \pmod{x + \alpha}$$

unde $D_i(x) = D_{i,170} x^{170} + \dots + D_{i,1} x + D_{i,0}$. Dacă $S_{i,1} = S_{i,0} = S_{i,-1} = 0$, cuvântul nu are erori. Dacă toate sindroamele $S_{i,1}$, $S_{i,0}$ și $S_{i,-1}$ sunt diferite de zero, atunci există erori de un singur *byte* în cuvântul citit. Acest lucru se datorează faptului că, dacă o eroare α^m ($0 \leq m \leq 254$) este adăugată simbolului de date $D_{i,l}$ ($l = \{0, 1, \dots, 170\}$), atunci expresia cuvântului citit este

$$D_{i,j} = D_{i,l} + \alpha^m$$

$$D_{i,k} = D_{i,k} \quad k \neq l, k = \{0, 1, \dots, 170\}$$

$$C_{i,j} = C_{i,j}$$

și, deci, sindroamele nu toate egale cu zero se pot obține din:

$$S_{i,1} = \alpha^{m \cdot 171}$$

$$S_{i,0} = \alpha^m$$

$$S_{i,-1} = \alpha^{m \cdot (171)}$$

În acest caz, tipul erorii poate fi exprimat ca $S_{i,0}$, iar poziția erorii se exprimă cu ajutorul următoarelor relații:

$$\frac{S_{i,1}}{S_{i,0}} = \frac{\alpha^{m \cdot 171}}{\alpha^m} = \alpha^{170}$$

$$\frac{S_{i,-1}}{S_{i,0}} = \frac{\alpha^{m \cdot 171}}{\alpha^m} = \alpha^{170}$$

2.1.3 Codul R-S întreșesut-încrucișat pentru discuri optice

După cum s-a arătat mai înainte, erorile din sistemele optice sunt aproape întotdeauna datorate suprafeței discului și, de aceea, în acest caz avem de a face nu numai cu erori de tip salvă, ci și cu erori aleatoare. De aceea avem nevoie de un cod pentru corecția ambelor tipuri de erori.

Codul R-S întreșesut-încrucișat pentru discuri optice (*CIRC — Cross Interleaved Reed-Solomon Code*) face parte dintr-o nouă clasă de coduri dublu codificate, în care cel de al doilea cod *Reed-Solomon* codifică ieșirea întârziată și dispersată (de exemplu, *întreșesut-încrucișat*) a primului cod codificat *Reed-Solomon*. Pentru a ilustra conceptul acestei clase de coduri este oferit un mic exemplu.

Exemplu

Fie două coduri C_1 și C_2 , capabile să corecteze, fiecare în parte, erori singulare. Să presupunem că simbolurile de verificare P și Q sunt reprezentate prin două simboluri, pentru C_1 și, respectiv, C_2 .

Simbolurile de verificare P_{4n} , sunt obținute din simbolurile consecutive: D_{4n} , D_{4n-1} , D_{4n-2} și D_{4n-3} , obținute prin folosirea unui codificator C_2 . Simbolurile codificate D_{4n} , D_{4n-1} , D_{4n-2} , D_{4n-3} și P_{4n} sunt întârziate cu lungimi inegale, devenind D_{4n} , D_{4n-3} , D_{4n-6} , D_{4n-9} și P_{4n-16} . Întârzierile sunt diferite pentru fiecare dintre simboluri. Simbolurile de verificare Q_{4n} sunt produse pe baza simbolurilor întârziate, de către un codificator C_1 . Ca rezultat al întreșeserii încrucișate, chiar dacă apar erori consecutive de patru simboluri, de exemplu, erori în D_7 , D_{10} , D_{13} și D_{16} , care nu pot fi corectate în C_1 , acestea sunt dispersate ca erori singulare în fiecare cod C_2 , putând fi, deci, corectate de către acest decodificator. Dacă apar erori în trei simboluri, de exemplu în D_{13} , D_{16} și D_{17} , acestea pot fi, de asemenea, corectate în pași secvențiali, astfel: eroarea singulară din D_{13} poate fi corectată în decodificarea C_2 , după care eroarea singulară din D_{16} este corectată în C_1 și, în final, eroarea singulară din D_{17} se corectează în C_2 . Totuși, se poate arata că erorile din patru simboluri, de exemplu în D_{13} , D_{14} , D_{16} și D_{17} , nu pot fi corectate, de vreme ce sunt câte două erori în sarcina decodoadelor C_1 și C_2 , încă din prima etapă. Pentru a corecta în mod adecvat aceste patru erori se apelează la tehnica de corecție a ștergerii, care permite corecția erorilor duble prin folosirea pointerilor de eroare pentru codurile corectoare a erorilor de simbol. De

aceea, în prima etapă, când erori de două simboluri sunt detectate în timpul decodificării C_1 , (de exemplu, în acest caz codul C_1 joacă rolul detectorului de eroare dublă), sunt inițializați pointeri de eroare pentru toate simbolurile de pe rampa crescătoare. În etapa următoare, decodificarea C_2 poate corecta aceste erori duble, folosind respectivii pointeri. Acest lucru este posibil datorită faptului că amintitele coduri corectoare C_1 și C_2 sunt coduri de distanță 3, putând, deci, să corecteze ștergeri cel mult duble.

În cazul sistemelor de înregistrare a compact-discurilor audio se folosesc următoarele două coduri corectoare de erori *Reed-Solomon*, de distanță 5, peste $GF(2^n)$ [VINH80]. Codurile C_1 și C_2 sunt după cum urmează:

C_1 : (32,28) cod corector *Reed-Solomon* de eroare dublă;

C_2 : (28,24) cod corector *Reed-Solomon* de eroare dublă.

În această schemă de codificare se pot lua în considerație multe strategii. Printre acestea, sistemele de tip compact-disc adoptă următoarea strategie, numită *super-strategie*:

1) Decodificarea C_1

- Sindromul „totul pe zero“ → nu sunt erori;
- Erori de un singur simbol → corecție a erorii singulare;
- Erori de două simboluri → corecție a erorii duble. Pointerii de localizare a erorii sunt transmiși la decodificatorul C_2 ;
- Erori de mai mult de trei simboluri → doar detectarea erorilor. Pointerii de localizare a erorilor sunt transmiși la decodificatorul C_2 ;

2) Decodificarea C_2 :

- Sindromul „totul pe zero“ → nu sunt erori;
- Erori de un singur simbol → corecție a erorii singulare;
- Erori de un singur simbol →
 - ◆ Numărul pointerilor erorilor de la decodificatorul C_1 este mai mic sau egal cu 4 și, în același timp, localizarea erorilor a doi dintre pointerii de eroare este egală cu cea calculată de către sindrom → corecția erorii duble;
 - ◆ Altele → doar detecția erorilor;

- (iv) Erori de mai mult de trei simboluri → doar detectarea erorilor.

A se observa că, atunci când erorile duble sunt corectate în timpul decodificării C_1 , acești doi pointeri de localizare a erorilor sunt transmiși decodificării C_2 corectarea greșită a erorilor triple, confundate ca erori duble în timpul decodificării C_2 . În cazul în care nici decodicatorul C_2 nu poate corecta erorile, el lasă să treacă 24 de simboluri necorectate de date, dar marcate cu pointerii de eroare proveniți de la C_1 . În acest fel chiar dacă decodicatorul C_2 nu poate decodifica, majoritatea simbolurilor sunt foarte lipsite de erori, iar valorile eșantioanelor marcate și necorectate pot fi refăcute pe baza unei interpolări lineare [VINH80] în sistemele audio digitale. În cazul unui eșantion marcat ca eronat, valoarea acestuia se poate deduce din cea a eșantioanelor corecte care îl mărginesc și care nu au fost marcate cu pointeri de eroare. Această schemă *CIRC* cu super-strategie are o capacitate maximă de corecție completă a salvelor de erori de până la 4000 de biți.

Pentru sistemele de stocare digitală, numite *CD-ROM*, interpolarea lineară nu constituie o soluție. De aceea trebuie adăugate alte metode de creștere a calității datelor. Un sistem adoptă atât codurile dublu codificate *R-S SbcEC*, cât și o verificare ciclică a redundanței (*CRC*) pentru datele deja cuprinse în *CIRC*. Cu alte cuvinte, *CD-ROM* folosește o metodă atât de puternică pentru corecția erorilor, încât biții *CRC* sunt adăugați la codificarea *CIRC*, rezultând informație dublu codificată în codurile *R-S SbcEC*. Astfel, datele sunt cvadruplu codificate, iar dacă se include și *CRC*, atunci sunt cvintuplu codificate. Un disc are o capacitate de circa 650 MB.

1) Codurile cu verificare ciclică a redundanței (*CRC*). Acest cod a fost adăugat pentru a verifica dacă nu a apărut o corectare eronată a codurilor *Reed-Solomon* sau pentru a detecta erori necorectabile. Polinomul generator al codului *CRC* este dat mai jos:

$$g(x) = (x^{16} + x^{15} + x^2 + 1)(x^{16} + x^2 + x + 1)$$

2) Codurile dublu codificate *Reed-Solomon*. Aceste coduri sunt determinate de către următorul polinom primitiv $g'(x)$, iar α este un element prim în $GF(2^8)$ (de exemplu, o rădăcină a lui $g'(x)$).

$$g'(x) = x^8 + x^4 + x^3 + x^2 + 1$$

Codul include două coduri *Reed-Solomon* (de exemplu, codul C_1 (26,24) *R-S* și codul C_2 (45,43) *R-S*). Ambele coduri sunt coduri *Reed-Solomon* cu distanța 3.

Octeții de verificare pentru C_1 sunt generați din informațiile incluse în direcția verticală. Octeții de verificare pentru C_2 sunt generați din informațiile conținute pe diagonala descrescătoare din zona datelor utile (care are dimensiunea de 24H43 de octeți, adică 1032 de octeți, în format bidimensional). Matricele H ale codurilor C_1 și C_2 sunt după cum urmează:

$$H_1 = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ \alpha^{25} & \alpha^{24} & \alpha^{23} & \dots & \alpha \end{bmatrix}$$

$$H_2 = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ \alpha^{44} & \alpha^{43} & \alpha^{42} & \dots & \alpha \end{bmatrix}$$

Întâi se generează codul C_1 , apoi codul C_2 . De observat este că decodificarea este permisă pentru oricare dintre secvențe, adică din C_1 în C_2 și invers. În plus aceste secvențe pot fi și repetitive, ca de exemplu: decodificare $C_1 \rightarrow$ decodificare $C_2 \rightarrow$ decodificare $C_1 \rightarrow$ decodificare $C_2 \dots$ Mai mult, putem folosi rezultatele decodificării *CIRC*. Fiecare dintre coduri are următoarele posibilități de corecție a erorilor:

- Ⓜ Posibilitatea corecției erorilor de un singur octet;
- Ⓞ Posibilitatea corecției ștergerilor de doi octeți.

Ca atare, există multe strategii posibile de decodificare. De exemplu, folosind capacitatea de corecție a ștergerilor, marcată cu Ⓞ, se poate aplica următorul model:

$$C_1 \textcircled{C} \rightarrow C_2 \textcircled{C} \rightarrow C_1 \textcircled{C} \rightarrow C_2 \textcircled{C} \rightarrow \dots$$

Pentru a ne proteja împotriva corecțiilor greșite din Ⓞ, codul *CRC* verifică aceste corecturi greșite și identifică erorile necorectabile. O astfel de strategie de decodificare poate atinge aceeași rată a erorilor ca unitățile comerciale de discuri magnetice sau de benzi magnetice.

Alt exemplu de *CD-ROM*, cu o capacitate de 2,6 GB, adoptă atât un cod *Reed-Solomon* întretesut, de distanță 5, cât și un cod de verificare ciclică a redundanței pentru fiecare grupare 512 octeți de date utile, numită sector.

În acest sistem este adoptat un cod C Reed-Solomon (39,35) de distanță 5 peste $GF(2^8)$, putându-se astfel corecta orice eroare de dublu-octet. Astfel, codul întretesut folosit Reed-Solomon poate corecta până la 30 de octeți greșiți într-un sector. În plus, acest sistem adoptă un cod de verificare ciclică a redundanței (unul cu patru octeți de verificare), care verifică informațiile corectate de către codurile Reed-Solomon și oferă protecție împotriva tentativei de corectare a erorilor necorectabile, care depășesc posibilitatea de corecție a codurilor Reed-Solomon. Codul de verificare ciclică a redundanței este adăugat la sfârșitul zonei datelor utile.

Alt tip de sistem de discuri optice, așa numitul „disc optic cu o scriere și multiple citiri“ (*CD-WORM - rite Once Read Many*), având un diametru de 5,25 țoli, este o propunere relativ recentă pentru memoria de masă a calculatoarelor. În cazul acestor discuri, există două tipuri de distribuții ale erorilor: una este cauzată de o scurtă salvă greșită (mai puțin de zece biți) cu probabilitate ridicată de apariție, iar cealaltă de o lungă salvă eronată (între 10 și 100 de biți) cu probabilitate scăzută de apariție [SaTN86]. Următoarele strategii au fost propuse pentru corectarea acestor tipuri de erori: un cod Reed-Solomon de mare distanță ($d \approx 17$) cu o lungime a codului de 120 până la 140 de octeți, întretesut într-un raport de la 4 până la 10, precum și un cod produs care folosește două coduri Reed-Solomon (având distanța cuprinsă între 3 și 5), ceea ce ar da naștere unui cod cu o distanță cuprinsă între 15 și 25, cu o lungime a codului de la 30 până la 50 de octeți [Yama86], [Kurt87].

Discurile optice reinscriptibile, dezvoltate recent, funcționează pe baza unui principiu magneto-optic și prezintă caracteristici similare ale erorilor apărute. De aceea, codurile Reed-Solomon de mare distanță li se aplică și lor cu aceeași eficiență [FRSH87, IYHI87, KKIS87].

2.2 Matrice de discuri independente și redundante

Asigurarea protecției datelor față de situațiile anormale, inclusiv erori umane, căderi de echipamente și condiții ostile a apărut ca o necesitate stringentă chiar cu primul sistem de calcul. Backup-ul, copierea periodică a datelor pe un suport mobil păstrat în afara sistemului, rămâne soluția esențială practică, în ciuda abilității date de creșterea numărului de discuri utilizate pentru a preveni alterarea accesului la date corecte în situații de căderi sau dezastre. Și mai puține progrese s-au realizat în domeniul minimizării erorilor umane ceea ce impune utilizarea în continuare a procedurilor de backup, chiar și pentru sistemele de discuri care pot supraviețui la căderi și dezastre.

Necesitatea ca sistemul de discuri să mențină un acces continuu *on-line* a apărut în anii '70, o dată cu creșterea masivă a utilizării sistemelor de calcul *on-line*. *Backup*-ul și restaurarea *off-line* erau considerate suficiente pentru a acoperi o cădere. Deoarece accesul la date corupte sau distruse nu are nici o valoare, proiectanții au fost obligați să furnizeze soluții de protejare a datelor și accesul la acestea în cazul unor căderi. Deoarece datele sunt memorate, în principal, pe discuri, prima soluție apărută a constat în memorarea acelorași informații pe două suporturi diferite (*mirroring*), metoda care protejează datele și menține accesul la ele chiar și în situația căderii unui disc din perechea utilizată.

Soluția „*mirroring*“ presupune costuri cel puțin duble pentru subsistemul de discuri, față de cel clasic, astfel că a apărut necesitatea practică de a găsi soluții mai puțin costisitoare pentru creșterea protecției datelor și asigurarea accesului continuu la acestea. Astfel, la sfârșitul anilor '80 s-au pus bazele conceptului de RAID – Redundant Array of Independent Disks. RAID utilizează redundanța sub forma unor informații de control (paritate). Față de *mirroring*, unde redundanța era de 100%, la RAID ea poate scădea și sub 10%.

Pentru o tratare unitară, soluția *mirroring* a fost inclusă în familia RAID care a fost divizată în două:

- *mirrored* RAID
- *parity* RAID – RAID protejat prin paritate

Spre deosebire de *mirrored* RAID, primele *parity* RAID-uri sufereau de o foarte scăzută performanță datorată necesității generării și scrierii parității în timpul operației de scriere și a regenerării *on-line* ca răspuns la o cerere către un disc căzut. Această problemă a fost rezolvată parțial prin tehnici precum *caching*, scriere asistată pe discuri etc. și a condus la o dezvoltare exponențială a utilizării RAID-ului după mijlocul anilor '90.

Două decade de cercetări în domeniul protecției datelor și menținerii *on-line* a accesului în cazul unei căderi sau a unui dezastru au generat o largă paletă de soluții de sisteme de discuri cu diverse grade de rezistență. Pentru evitarea confuziilor și pentru a ușura alegerea celei mai potrivite soluții pentru o anumită aplicație a fost înființată o asociație internațională a cercetătorilor din domeniul RAID la care sau afiliat și cei mai importanți producători, asociație denumită RAID Advisory Board's (RAB).

Primele mele cercetări în domeniul fiabilității transiterii și stocării datelor s-au focalizat, începând cu anul 1991, în subdomeniul RAID, în principal abordând problematica legată de „Algoritmii de distribuție a datelor, parității și spațiilor de rezervă“. În multe sisteme de calcul, subsistemul I/O reprezintă adesea o pondere importantă în cadrul performanțelor sistemului. De-a lungul ultimilor ani, viteza CPU a crescut în proporții de 40% până la 100% pe an, pe când timpul căutării pe disc s-a îmbunătățit cu doar 7% pe an (creșteri asemănătoare de performanță se pot observa și la celelalte componente ale subsistemului I/O).

Acest fapt a condus la o mare discrepanță între viteza CPU/memorie principală și cea a subsistemului operațiilor I/O, și este de așteptat ca aceasta să crească și mai mult în viitorul apropiat. Bazându-ne pe aceste observații, este de prevăzut ca viitoarele creșteri ale vitezei de lucru ale procesoarelor să determine doar o mică îmbunătățire a performanțelor totale ale sistemului de calcul. Se ridică următoarea problemă fundamentală: bazându-ne pe viitoarele dezvoltări predictibile ale tehnologiei în domeniul discurilor, cum se poate reduce discrepanța de viteză dintre CPU/memorie principală și subsistemul I/O?

Evident, datorită limitărilor tehnologice, soluțiile pot fi găsite doar la nivel structural.

O idee atractivă este cea numită matrice de discuri (*disk array*), unde subsistemul de discuri este compus din mai multe discuri, cu datele împrăștiate pe toate aceste discuri. Conceptele de „disc întrețesut“ (*disk interleaving*) și „disc secționat“ (*disk striping*) au fost pentru prima oară folosite de către Kim [Kim85], respectiv Salem și Garcia-Molina [GaMo87]. De atunci, o mare cantitate de muncă a fost depusă asupra mai multor modele, referitor la performanțele matricelor de discuri și a fiabilității acestora. Matricele de discuri au fost împărțite în șase clase denumite RAID (*Redundant Arrays of Inexpensive Disks*) [PDGG88]. Printre cele cinci, cele mai promițătoare două candidate la obținerea de performanțe ridicate într-un sistem de calcul par a fi RAID 1 (matrice de discuri oglindită - *mirrored disk array*) și RAID 5 (matrice cu paritate rotită - *rotated parity array*), RAID 6 nefiind încă unanim acceptată drept RAID de sine stătător, iar RAID 7 reprezentând un rezultat al autorului acestei lucrări.

Pe măsură ce performanța altor componente ale sistemului continuă să se dezvolte rapid, performanța subsistemului de stocare devine tot mai importantă. Performanța subsistemului de stocare și fiabilitatea pot fi crescute prin gruparea logică a mai multor unități de disc în matrice de

discuri. Organizările matricelor de discuri sunt definite de schemele de distribuție și de mecanismele de redundanță. Diferitele combinații ale acestor două componente fac matricele de discuri potrivite pentru o scară largă de aplicații. Multe decizii de implementare a matricelor rezultă, de asemenea, din compromisul dintre performanță și fiabilitate.

2.2.1 Organizarea datelor

Organizarea datelor pentru matricele de discuri poate fi partiționată în două componente ortogonale: schemele de distribuție a datelor și mecanismele de redundanță. Distribuția de date definește translația adreselor logice, vizibile extern, în locații de stocare din cadrul subsistemului. Sunt două căi pentru a face acest lucru: prin adresarea independentă a fiecărui disc sau prin divizarea discului. Cea din urmă oferă, adesea, grade îmbunătățite de încărcare echilibrată, transfer paralel de date, sau acces concurrent. Mecanismele de redundanță specifică tipul, scopul, și locația oricărei informații redundante din matrice (de exemplu, replicarea datelor și paritatea divizată). Matricele redundante de discuri continuă să ofere acces deplin la date și după defectarea unui disc și în timp ce datele pierdute sunt reconstruite pe un disc de înlocuire. Multe matrice de discuri combină divizarea discului cu redundanța pentru a oferi atât performanță cât și fiabilitate.

Produsele cu matrice de discuri furnizează, de obicei, *software*-ul și *firmware*-ul necesare care permit cumpărătorului să aleagă din mai multe scheme posibile de organizare. Deoarece aceasta lărgeste domeniul de folosire și îi crește performanța relativă la cost, administratorii de sistem trebuie să înțeleagă compromisurile care apar pentru a putea exploata pe deplin posibilitățile matricei.

2.2.1.1 Distribuția datelor

Distribuția de date constă în maparea adreselor logice folosite de sistemul gazdă pe discurile din subsistem. Există trei metode principale folosite în realizarea acestei translații. Metoda convențională este de a adresa, independent, fiecare disc și de a mapa numerele de bloc logic în numere de bloc de disc, direct. Distribuția este făcută manual de către administratorii de sistem, programe de aplicații sau sisteme *software*. Divizarea de disc, numită și *disk striping* sau *disk interleaving*, reunește mai multe spații de adrese de disc într-un singur spațiu, unificat, văzut de sistemul gazdă. Aceasta se face prin distribuirea consecutivă de unități logice de date (unități de divizare) din cadrul discului într-o manieră *round-robin*. Se distinge între două tipuri diferite de divizări de discuri: divizare fină, în care discurile cooperează toate pentru

satisfacerea unei cereri și divizare grosieră, în care discurile cooperează la cereri mari de date și lucrează independent la cereri mici de date.

2.2.1.2 Adresarea independentă

Dacă fiecare disc este adresat independent, utilizatorii sau sistemele gazdă trebuie să distribuie explicit datele. În general, administratorii de sistem sunt responsabili pentru răspândirea datelor pe fiecare disc. Încărcarea echilibrată sau distribuirea datelor astfel încât să se încarce echilibrat sistemul, este mai mult o artă decât o știință și necesită specialiști. Adrese de date adiacente creează *hot spots* făcând *task*-ul complex.

| DISC 0 | DISC 1 | DISC 2 | DISC 3 |
|--------|--------|--------|--------|
| A0 | B0 | C0 | D0 |
| A1 | B1 | C1 | D1 |
| A2 | B2 | C2 | D2 |
| ... | ... | ... | ... |

Figura 2-1

2.2.1.2.1 Divizarea fină

Aici, toate cele n discuri conțin o fracțiune din fiecare bloc accesibil. Numărul discurilor și mărimea unității de divizare sunt, în general, alese ținând cont de cea mai mică unitate de date accesibilă. Mărimi mai folosite sunt un bit, un octet, un sector. Mărimea folosită nu este neapărat importantă deoarece fiecare unitate accesibilă este răspândită pe toate discurile. Încărcarea este echilibrată, deoarece toate discurile primesc aceeași încărcare. Rata de transfer este de aproximativ n ori mai mare față de cea a unui disc individual, deoarece fiecare disc transferă $1/n$ din datele cerute. Totuși, componentele de poziționare ale capului ale timpului de acces - căutare și rotire - sunt ori neafectate ori crescute (dacă discurile sunt nesincronizate). De asemenea, doar o singură cerere poate fi făcută la un moment dat, deoarece toate cele n discuri lucrează la aceeași cerere. Aceste limitări restrâng, de obicei, folosirea metodei la domenii în care timpii de transfer domină timpii în care nu se lucrează cu discul.

| DISC 0 | DISC 1 | DISC 2 | DISC 3 |
|--------|--------|--------|--------|
| A0 | A0 | A0 | A0 |
| A1 | A1 | A1 | A1 |
| A2 | A2 | A2 | A2 |
| ... | ... | ... | ... |
| A7 | A7 | A7 | A7 |
| ... | ... | ... | ... |
| B0 | B0 | B0 | B0 |
| B1 | B1 | B1 | B1 |
| ... | ... | ... | ... |

Figura 2-2

2.2.1.2.2 Divizarea grosieră

Cu aceasta, nu mai este necesar ca toate discurile să coopereze la fiecare cerere. Această tehnică exploatează paralelismul transferului de date pentru cereri mari, în timp ce permite discurilor separate să satisfacă cereri mici concurrent. Câteva efecte se combină pentru a oferi o încărcare echilibrată și automată. Primul efect este acela că un fișier va tinde să aibă blocurile sale, componentele împrăștiate pe mai multe discuri. Divizarea grosieră, ca la *hashing* efectiv, randomizează primul disc accesat pentru fiecare cerere, celelalte discuri fiind identificate în maniera *round-robin*. Această încărcare, echilibrată statistic, rămâne intactă și în cazul adreselor de date adiacente și distribuții rapid schimbătoare de accese. Din motivele de mai sus, metoda are posibilitatea să ofere performanță bună cu un minim de întreținere.

| DISC 0 | DISC 1 | DISC 2 | DISC 3 |
|--------|--------|--------|--------|
| A0 | A1 | A2 | A3 |
| A4 | A5 | A6 | A7 |
| ... | ... | ... | ... |
| B0 | B1 | B2 | B3 |
| B4 | B5 | B6 | B7 |
| ... | ... | ... | ... |

Figura 2-3

2.2.1.2.3 Alegerea mărimii unității de divizare

Realizarea performanței de la metoda de mai sus necesită alegerea corectă a mărimii unității de divizare. Această alegere determină, pe de o parte, numărul de discuri peste care sunt răspândite datele fiecărei cereri. Mărimi mari ale unității permit ca discuri separate să satisfacă cereri multiple (concurrent), reducând timpul general de poziționare, rezultând adesea reducerea cozilor de așteptare și o rată de transfer mai mare. Pe de altă parte, unitățile mici determină mai multe discuri să acceseze datele în paralel, reducând timpii de transfer pentru cereri individuale. Acest compromis între paralelismul transferului și concurența accesului este guvernată de mărimea unității de divizare.

Când avem cereri de lungime variabilă, fără o secvențializare sau localizare a adreselor, parametrul cheie este cantitatea de concurență din încărcare. Mărimea unității de divizare ar trebui să crească o dată cu rata de sosire a cererilor, deci descrescând numărul mediu de discuri accesate pe cerere.

Pentru cereri constând în primul rând din cereri aliniată, de mărime fixă, mărimea unității de divizare poate fi aleasă astfel încât numărul de accese la disc, pe cerere, să fie constant. În particular, folosind orice multiplu de mărimea fixă a cererii, se permite fiecărei cereri să fie servită de către un singur disc.

Pentru încărcări cu un oarecare grad de aliniament sau localizare a adresei, echilibrarea trebuie, de asemenea, considerată. Când cererile sunt mari, echilibrarea nu afectează semnificativ alegerea mărimii unității de divizare. Pentru încărcări mici, unitățile mari exploatează

beneficiile timpului de poziționare datorate localizării și secvențialității adreselor.

Un compromis mai complex este cel pentru încărcări generale, adesea caracterizate prin cereri mici, cereri mari scurte, și adrese puternic adiacente. În aceste cazuri, abilitatea metodei divizării grosiere de a echilibra încărcarea se îmbunătățește pe măsură ce mărimea unității de divizare scade. Mărimi mai mari ale unității tind să sufere din cauza formării unor *hot spots* de scurtă durată. S-a sugerat că unitatea de divizare trebuie să fie de 10 ori mai mare decât mărimea medie a cererii, rotunjită la un multiplu al mărimii pistei. Alegerea mărimii unității divizare poate fi rafinată considerând informație adițională, cum ar fi distribuțiile de mărime a cererii și caracteristicile localizării.

O întrebare interesantă este dacă mărimea unității de divizare trebuie să fie multiplu de mărime a pistei. Pe de o parte, o cerere pentru unitatea de divizare completă are avantajul unui disc suportând acces cu latență zero. Pe de altă parte, mărimea pistei poate corespunde, în mod dificil, mărimii și aliniamentului cererii, rezultând cereri multidisc, excesive. Multe încărcări constau din cereri de blocuri de mărime specifică, să zicem 4k sau 8k. Mărimile curente de piste sunt de multipli de sectoare (adesea 512 octeți) alese, mai degrabă, să maximizeze capacitatea decât să potrivească limitările datorate alinierii. Cele mai multe unități de disc *state-of-the-art* folosesc înregistrarea pe zone multiple pentru a crește densitatea de stocare. Folosind aceste discuri, potrivirea mărimii unității de divizare cu cea a pistei cere o mărime diferită pentru fiecare zonă de disc. Sunt necesare mai multe studii de performanță înainte ca divizarea bazată pe piste să poată fi identificată ca o alegere bună, în general.

2.2.2 Mecanismele de redundanță

Adăugând mecanisme de redundanță la un sistem de stocare, invariabil îi crește costul, i se reduce capacitatea, și/sau i se reduc performanțele. În schimb, multe sisteme fac, periodic, copii ale datelor critice, de obicei, pe banda magnetică. Orice informație coruptă sau pierdută din sistemul primar este restocată de pe copie și orice actualizări care au apărut după efectuarea copiei sunt irecuperabile. Din fericire, căderile individuale de discuri sunt foarte rare, iar căderile de sistem sunt și mai rare. Timpul mediu dintre defectări (MTBF) la un disc modern este de la 200000 până la 1 milion de ore sau, altfel spus, 20-100 de ani.

Totuși, pe măsură ce numărul de discuri continuă să crească, MTBF-ul unui sistem tipic neredundant se micșorează până la perioade de luni sau chiar săptămâni. Presupunând că MTBF-ul fiecărui disc este independent și exponențial distribuit, MTBF-ul unui grup este invers

proporțional cu numărul discurilor. Aceasta amenință fezabilitatea sistemelor nonredundante de discuri.

În plus, numărul de discuri afectate de cădere crește dramatic o dată cu folosirea divizării discurilor. Un subsistem divizat poate pierde $1/n$ din fiecare fișier, față de $1/n$ fișiere din matricea de discuri în cazul adresării independente. Un subsistem cu redundanță în organizarea datelor poate supraviețui căderii uneia sau mai multor componente. Totuși, la apariția unor noi căderi va rezulta inaccesibilitatea sau pierderea datelor. O dată cu creșterea subsistemelor de stocare, redundanța *on-line* va fi necesară pentru datele cruciale, devenind, probabil, un standard pentru stocare, în general.

2.2.2.1 Replicarea datelor

Cel mai simplu mod de a proteja datele este replicarea lor. Copii separate ale fiecărui bloc sunt menținute pe discuri separate. Datele devin indisponibile doar dacă toate discurile, conținând o copie, cad. Această formă de redundanță, folosită de mai mulți ani, se numește *disk mirroring*, *disk duplexing*, *disk shadowing* sau, mai nou, RAID 1 (Figura 2-4). Costul replicării datelor este proporțional cu $d-1$, deoarece costul discului domină prețul subsistemelor de stocare. Deoarece două copii oferă înaltă fiabilitate, alte copii, în plus, nu se folosesc. De fapt, costul mare al replicării obligă, adesea, la folosirea altor mecanisme de redundanță care nu dau însă același nivel de performanță și fiabilitate.

| DISC 0 | DISC 1 | DISC 2 | DISC 3 |
|--------|--------|--------|--------|
| 0 | 0 | 2 | 2 |
| 1 | 1 | 3 | 3 |

Figura 2-4

Menținerea unor copii multiple ale datelor afectează performanța în diferite moduri. Primul este acela că scrierile trebuie făcute pe toate copiile. Depinzând de implementare, aceasta poate crea oportunități pentru pierderea datelor sau timpi crescuți de răspuns. Performanța la citire este, pe de altă parte, mai bună în acest caz. În primul rând, cererile de citire se satisfac simultan, crescând rata de transfer și scăzând timpii de așteptare. În plus, o cerere de citire poate fi programată pe discul cu cel mai mic timp de acces. Această performanță bună la citire poate face raportul cost/performanță chiar mai mic decât în cazul

sistemelor cu discuri nonredundante, deși costul absolut este de d ori mai mare.

| DISC 0 | DISC 1 | DISC 2 | DISC 3 |
|--------|--------|--------|--------|
| 0 | 1 | 2 | 3 |
| 3 | 0 | 1 | 2 |

Figura 2-5

Metoda convențională de implementare a replicării datelor dublează fiecare disc (Figura 2-4). Fiecare set de d discuri identice pot supraviețui la $d-1$ căderi fără pierdere de date. Într-o matrice cu seturi multiple replicate, pot apărea multe căderi, fără să cauzeze vreo pierdere de informații, atâta timp cât d din ele nu sunt din același set unic, replicat. Totuși, fiecare cădere de disc reduce performanțele unui set relativ la alte seturi. Acest dezechilibru poate duce la scăderea performanțelor generale ale sistemului, dacă un set căzut devine o gâtuire.

| DISC 0 | DISC 1 | DISC 2 | DISC 3 |
|--------|--------|--------|--------|
| 0 | 1 | 2 | 3 |
| 1c | 0a | 0b | 0c |
| 2b | 2c | 1a | 1b |
| 3a | 3b | 3c | 2a |

Figura 2-6

Alternativ, seturile se pot suprapune astfel încât această degradare de performanță este împărțită în mod egal de toate discurile rămase. Termenul *declustering* se referă la o metodă pentru combinarea seturilor multiple replicate. Metoda *chained declustering* se referă la divizarea fiecărui disc în d secțiuni. Prima secțiune conține copia primară a unor informații, cu fiecare disc având date primare diferite (Figura 2-5). Cu o aranjare inteligentă, această metodă poate menține o încărcare echilibrată după o cădere a unui disc, deși probabilitatea supraviețuirii la căderi multiple este redusă.

Al doilea tip de *declustering*, cunoscut și sub numele de interpusă partiționează, de asemenea, fiecare disc, dar divizează copiile

nonprimare din toate discurile în loc să le păstreze în spații contigue. Aceasta echilibrează încărcarea suplimentară cauzată de o cădere de disc, dar reduce, în continuare, fiabilitatea (pierderea oricăror d discuri rezultă într-o cădere critică, (Figura 2-6).

2.2.2.2 Protecția bazată pe paritate

Proiectanții de subsisteme de discuri s-au folosit de rezultatele studiilor extinse ale codurilor detectoare și corectoare de erori. Din teoria codurilor se știe că pierderea unui singur bit poate fi recuperată știind bitul de paritate, ceilalți biți de protejat și faptul că acel bit este eronat. Deoarece sistemele actuale de discuri au o logică de control care, în general, detectează orice cădere de disc, poate fi folosită pentru protejarea unui set de discuri față de erorile singulare. Fiecare bit de paritate este calculat din bitul de date asociat de pe fiecare disc protejat și stocat pe un disc separat. Locațiile unor biți particulari de date dați de un bit de paritate pot varia de la schemă la schemă, așa cum poate varia și aranjarea parității. Capacitatea adițională cerută (un disc complet) pentru a menține redundanța prin intermediul parității este considerabil mai mică decât atunci când se folosește replicarea datelor. Deși costul este rezonabil, menținerea parității poate reduce în mod semnificativ performanța. Cererile de citire sunt efectuate ca în matricele nonredundante. Dacă se adaugă un disc pentru a oferi capacitate de stocare pentru paritate, el oferă servicii în plus pentru cereri de citire (în cele mai multe scheme de așezare a parității). Cererile de scriere necesită lucru suplimentar cu discul pentru actualizarea parității, adesea ducând la o performanță scăzută. Dacă toți biții de date asociați cu un bit de paritate sunt actualizați de o cerere de scriere, noul bit de paritate este calculat și inclus ca o parte a cererii totale. Dacă doar un subset al biților de date este scris, actualizarea bitului de paritate necesită lucru suplimentar.

Aceasta constă în citirea biților de pe discuri pentru construirea noului bit de paritate. Avem două metode:

- **Read-modify-write**, unde noul bit de paritate este construit din valorile vechi și cele noi ale biților care se scriu;
- **Regenerate write**, în care se generează bitul nou de paritate din valorile tuturor biților (adică noile valori ale biților care se scriu și valorile curente ale biților neschimbați).

RMW necesită, mai întâi, accese de citire la toate discurile care se actualizează (incluzând discul care conține paritatea), pe când R-W cere accese de citire la toate discurile care nu se actualizează. Pierderea

performanței poate fi redusă prin combinarea dinamică a acestor două opțiuni.

| DISC 0 | DISC 1 | DISC 2 | DISC 3 |
|--------|--------|--------|--------|
| A | A | A | a |
| B | B | B | b |
| C | C | C | c |
| ... | ... | ... | ... |

Figura 2-7

Datorită nevoii de a face citiri de pe disc înainte de actualizarea parității, performanța poate fi semnificativ mai mică față de un sistem neprotejat sau față de un sistem protejat prin replicarea datelor. Reducerea cât mai mult a acestei degradări de performanță este, încă, un subiect aflat în cercetare.

| DISC 0 | DISC 1 | DISC 2 | DISC 3 |
|--------|--------|--------|--------|
| A | A | A | a |
| B | B | b | B |
| C | c | C | C |
| ... | ... | ... | ... |

Figura 2-8

Există trei scheme majore pentru împrăștierea informației de paritate în cadrul discurilor din subsistem. Metoda directă este de a plasa toți biții de paritate pe un disc dedicat de paritate (Figura 2-7). Deși simplifică maparea adreselor logice în adrese de disc, fiecare scriere trebuie să actualizeze biții asociați de pe unicul disc de paritate. Matricele care folosesc divizarea fină sunt potrivite pentru această metodă, deoarece ele satisfac doar o cerere la un moment dat. Cu alte scheme de distribuție, totuși, discurile pot limita performanța, mai ales când vechea paritate este cerută să participe în operațiile de la RMW. Pentru evitarea acestei potențiale gâturi, informația de paritate poate fi divizată printre discuri (Figura 2-8). Folosind divizarea parității, matricea poate efectua mai multe actualizări în paralel.

| DISC 0 | DISC 1 | DISC 2 | DISC 3 |
|--------|--------|--------|--------|
| A | A | A | a |
| B | b | C | B |
| B | C | C | C |
| ... | ... | ... | ... |

Figura 2-9

O a treia opțiune s-a propus: *declustered parity*. Pentru înțelegerea acestei scheme ne putem imagina combinarea câtorva matrice mai mici, protejate de paritatea divizată. Datele și paritatea fiecărei matrice sunt distribuite în întregul set de discuri disponibile (Figura 2-9). Cu un algoritm inteligent de distribuție, impactul performanței datorat pierderii unui disc poate fi împărțit, în mod egal, de toate discurile rămase. Deși matricea poate supraviețui doar unei singure căderi, ar trebui să se refacă mai repede. Ca o alternativă, aceasta metodă se poate imagina ca o cale de a realiza fiabilitate mai bună și performanță în funcționarea la defecte prin creșterea cantității de date de paritate menținută.

2.2.2.3 Alte scheme

Menținerea parității protejează o matrice împotriva căderii unui singur disc. O matrice cu două copii ale tuturor datelor poate supraviețui până la o cădere pe perechea de copii. În ambele cazuri, datele pot fi pierdute dacă apare o a doua cădere, înainte ca prima să fie reparată. Oferirea protecției necondiționate împotriva căderilor multiple necesită redundanță suplimentară. În cele mai promițătoare metode se folosesc două coduri *Reed-Solomon*, corectoare de erori. Cu asemenea redundanță, o matrice poate reface toate datele după pierderea accesului la oricare două discuri. Compromisurile referitoare la performanță și opțiunile de amplasare a redundanței se potrivesc pentru acele matrice cu protecție bazată pe paritate, în afară de faptul că doi biți sunt mai bine menținuți decât unul.

Alte scheme, cum ar fi paritatea multidimensională și codurile mai agresive, corectoare de erori, pot proteja împotriva unui număr mai mare de căderi. Totuși, costul și implicațiile de performanță, combinate cu faptul că fiabilitatea sistemelor de stocare crește tot mai mult, pot obliga folosirea lor.

2.2.3 Performanța versus fiabilitatea

Cererile de scriere pot duce la întârzieri mari între momentul inițierii lor și momentul scrierii efective a noii informații pe suportul stabil. Până când sistemul gazdă primește mesajul de satisfacere a cererii orice resurse de sistem nu pot fi reutilizate. De asemenea, procesele pot aștepta pentru astfel de confirmări, putând rezulta cicluri inutile de procesor. Acceptând confirmarea cât mai repede posibil, gazda poate utiliza mai bine resursele sale și are, astfel, mai puține cicluri inutile de procesor, rezultând o performanță mai bună. Totuși, în cazul în care confirmarea cererii este dată înainte ca datele să fie scrise efectiv pe disc, sistemul este vulnerabil la coruperea sau pierderea informațiilor, datorită căderii tensiunii sau defectării componentelor. Când sistemul renunță la constrângerile de fiabilitate în scopul câștigării performanței, trebuie luate măsuri de precauție sau se acceptă consecințele.

2.2.3.1 Negocierea fiabilității în favoarea performanței

Un controler de matrice de discuri cu un *cache* (sau *buffer*) poate semnaliza confirmarea gazdei de îndată ce ultimul octet a fost copiat din memoria principală. Astfel de scheme îmbunătățesc mult performanța văzută de gazdă la scriere, deși datele sunt vulnerabile până la scrierea pe disc. Managerul de disc, care este responsabil de determinarea ordinii de servire pentru accesele la disc aflate în așteptare, are, acum, o problemă. Fiabilitatea cere managerului să expedieze accesele la disc care transferă datele din *cache* pe suportul fizic, în timp ce performanța obligă managerul să amâne această activitate de fundal.

În matricele care folosesc replicarea informațiilor, raportarea confirmării după ce doar o fracțiune din copii a fost actualizată poate și ea să îmbunătățească performanța. Datele sunt pierdute doar dacă apar căderi multiple înainte să se efectueze scrierile rămase. Pentru matricele protejate prin paritate, punând paritate în *cache* poate să se îmbunătățească performanța prin acumularea noii parități pentru scrieri multiple și efectuarea actualizării de paritate într-o singură scriere.

În schemele de paritate care cer cicluri de actualizare, sistemul este vulnerabil la coruperea datelor dacă setul de scrieri pe disc nu se face simultan. În cazul în care cade tensiunea, la un moment nepotrivit, doar câteva cereri vor fi satisfăcute. Aceasta face ca paritatea să fie incorectă, ducând la coruperi ulterioare dacă sunt lăsate necorectate și apoi folosită la refacerea datelor. Pe de altă parte, impunerea scrierilor simultane necesită comportare sincronă din partea mai multor discuri. Aceasta poate produce programări neoptime ale cererilor și un timp de disc inutil

crescut. Când controlerul își programează, independent, cererile individuale de scriere, utilizarea discurilor poate fi maximizată.

2.2.3.2 Îmbunătățirea performanței și a fiabilității

Cei mai mulți utilizatori nu sunt deranjați de coruperea datelor, exceptând cazurile inacceptabile. Există câteva tehnici pentru îmbunătățirea performanței și a fiabilității. Ele se realizează prin reducerea latentelor la scriere în timp ce este garantată siguranța datelor vechi și se asigură securitatea datelor noi înainte de semnalizarea confirmării. Aceste scheme pot fi implementate independent de sistemul gazda sau în cadrul matrice de discuri. O soluție simplă este folosirea memoriei RAM *cache*, non-volatilă, pentru a scrie datele sau paritatea. Aceste memorii mențin informațiile chiar și după întreruperea tensiunii de alimentare. Această organizare ar trebui să ofere o fiabilitate echivalentă cu cea a unei scheme cu redundanță.

Remaparea scrierii, o formă de paginare dublă, a fost propusă pentru multe organizări de date diferite. Această remapare constă în alterarea dinamică a mapării adreselor logice în cele fizice, astfel încât cererea de scriere de date este totdeauna plasată la locații noi. Prin menținerea copieii anterioare a blocului logic până la satisfacerea cererii, metoda îmbunătățește fiabilitatea. Desigur, mapările trebuie menținute fiabil pentru a preveni pierderea informațiilor. O metodă este stocarea informației în NVRAM. Altă metodă, mai complexă, dar mai puțin scumpă este să mărim fiecare bloc fizic de date cu adresa logică corespunzătoare și cu un indicator de timp, care pot fi folosite să reconstruiască mapările.

Alegând locații libere lângă poziția curentă a capului de citire/scriere, remaparea scrierii poate reduce întârzierea produsă de căutare și rotire. Algoritmul pentru alegerea locației disponibile trebuie să prevină consumarea tuturor locațiilor libere dintr-o anumită regiune și lăsarea altor regiuni nefolosite. De asemenea, în timp, datele secvențial logice vor fi disperate pe spațiul fizic. Ca rezultat, cereri de citire ulterioare s-ar putea să aibă de suferit din cauza întârzierilor de poziționare dacă nu se face o oarecare relocare pentru a reface spațiul fizic.

A treia opțiune - *X logging* - scrie informația despre actualizări într-o zonă de log înainte de actualizările propriu-zise. Scrierile log sunt secvențiale și au, de obicei, un timp de răspuns mai mic, mai ales când se folosește un disc log dedicat. După scrierea intrării de log, controlerul de matrice poate, în siguranță, să transmită un mesaj de terminare către gazdă și gazda accesează discul în fundal. În caz de cădere, conținutul log-ului permite repetarea oricărei actualizări incomplete. Informația

scrisă în log poate să conțină o descriere scurtă a actualizării, o copie a noii informații sau copii ale amândurora. O formă de verificare este folosită pentru a reduce timpul de refacere și permite refolosirea spațiului de log.

2.2.4 Tratarea căderilor

Prin menținerea redundanței, o matrice de discuri poate supraviețui căderii unuia sau mai multor discuri. În timp ce matricea continuă să ofere acces la toate datele, performanța ei este aproape întotdeauna afectată de căderi. În plus, fiabilitatea este redusă, deoarece căderile ulterioare pot fi critice. După cădere, datele pierdute trebuie să fie reconstruite de pe discurile rămase și scrise pe un disc de înlocuire, proces numit refacere. Până când acest proces se termină, matricea este vulnerabilă. Discurile suplimentare, numite rezerve rapide, sunt adesea incluse în matrice ca să permită refacerii să înceapă imediat.

În cele mai multe scheme, pierderea accesului la disc reduce numărul de cereri concurente de citire care pot fi deservite. Matricele care folosesc replicarea de date primesc cereri ca de obicei, dar mai puține discuri efectuează citiri. De asemenea, scrierea necesită actualizarea unei copii mai puțin. La matricele protejate prin paritate, o cerere de citire spre un disc defect cere o regenerare a datelor lipsă. Fiecare bit de date pierdut poate fi refăcut din bitul de paritate asociat, ceilalți biți contribuind la paritate. Deci, lucrul cu cereri de citire necesită accese la toate discurile care conțin acești biți. Performanța la scriere este și ea afectată în matricele protejate de paritate. Datorită faptului că este afectată și vechea informație, *Regenerate-Write* trebuie folosit pentru toate actualizările de date pe discul căzut și RMW pentru toate actualizările pe discurile rămase.

Cea mai simplă formă de refacere reconstruiește secvențial informațiile indisponibile și le scrie pe un disc de înlocuire. Într-o matrice replicată, datele pierdute sunt copiate de pe un alt disc conținând aceleași date. Într-una cu protecție bazată pe paritate datele trebuie refăcute din biții de paritate și cei de date. Dacă cerințele de fiabilitate prevăd ca timpul de reparare să fie minim, cererile gazdei pot fi anulate în timp ce refacerea are loc la viteză maximă. Altfel, refacerea și sistemul gazdă concurează pentru deservire, rezultând un alt compromis între fiabilitate și performanță. Acest compromis este, în particular, important în refacerea matricelor cu protecție bazată pe paritate, unde regenerarea datelor cere accese la mai multe discuri. Alegerea unui algoritm de refacere, mărirea unei cereri individuale de refacere, și complexitatea

circuisticii de coordonare multidisc trebuie luate în considerație când se ajustează procesul de refacere.

Piața de matrice de discuri se extinde pe măsură ce costul unităților de disc scade și capacitatea de stocare de uz general are nevoie de o creștere. Ca rezultat, matricele de discuri devin un echipament standard pe o scară tot mai mare de mașini. În curând, cerințele de stocare ale utilizatorului mediu pot fi satisfăcute fără folosirea părților mecanice. Până atunci, matricele de discuri pot oferi niveluri necesare de fiabilitate și performanță.

Dacă subsistemele de stocare continuă să se dezvolte ca și până acum, nivelul inteligenței încorporate în diferitele lor componente va continua să crească. Comunicația între aceste componente va crește mult, rezultând optimizări globale, un control mai distribuit, și noi algoritmi de distribuție care răspund dinamic la mereu schimbătorul flux de cereri.

2.3 Clasificarea subsistemelor de discuri redundante din punct de vedere a fiabilității

Pentru a marca mai precis domeniul meu de interes, voi prezenta o clasificare, din punct de vedere al fiabilității, pentru sistemele de discuri, propusa de RAB [Mass97].

Pentru a putea realiza o clasificare a subsistemelor redundante de discuri din punct de vedere a fiabilității, este necesară prezentare unor criterii [Mass97]:

- 1) protecție împotriva pierderii de date sau a accesului la date datorat căderii unor discuri
- 2) reconstrucția conținutului unui disc defect pe un înlocuitor
- 3) protecție împotriva pierderii de date în timpul întreruperii scrierii
- 4) protecție împotriva pierderii de date în timpul căderii magistralelor de I/O
- 5) protecție împotriva pierderii de date în timpul căderii oricărei componente implicate în transferul de date
- 6) monitorizare unităților rezistente la căderi precum și semnalizarea căderilor
- 7) disc activ pentru *swap*

- 8) protecție împotriva pierderii de date în timpul căderii memoriilor *cache*
- 9) protecție împotriva pierderii de date în timpul căderii tensiunii externe de alimentare
- 10) protecție împotriva pierderii de date în condiții de temperatură în afara domeniului de valori acceptate
- 11) semnalizarea căderilor de componente și a căderilor din subsistemele adiacente;

| Tip sistem de discuri | Criterii suportate |
|---------------------------------------|--------------------|
| Rezistent la căderi - FRDS | 1-6 |
| Rezistent la căderi plus – FRDS+ | 1-11 |
| Tolerant la căderi - FTDS | 1-13 |
| Tolerant la căderi plus - FTDS+ | 1-17 |
| Tolerant la căderi plus plus - FTDS++ | 1-18 |
| Tolerant la dezastre – DTDS | 1-19 |
| Tolerant la dezastre plus – DTDS+ | 1-21 |

Figura 2-10

- 12) protecție împotriva pierderii de date în timpul căderii canalului de transfer;
- 13) protecție împotriva pierderii de date în timpul căderii *contoller*-ului;
- 14) protecție împotriva pierderii accesului la date în timpul căderii memoriei *cache*;
- 15) protecție împotriva pierderii accesului la date în timpul căderii surselor de alimentare;
- 16) protecție împotriva pierderii accesului la date în timpul căderii magistralelor de I/O;
- 17) protecție împotriva pierderii accesului la date în timpul căderii surselor externe de alimentare;
- 18) protecție împotriva pierderii accesului la date în timpul înlocuirii FRU;

- 19) disc de rezervă activ;
- 20) protecție împotriva pierderii accesului la date în zona defectă;
- 21) protecție împotriva pierderii accesului de la mare distanță la date în zona căzută.

Clasificarea subsistemelor de discuri din punct de vedere al fiabilității pe care o voi folosi în continuare este data în Figura 2-10.

2.4 Concluzii

Pornind de la metodele de creștere a fiabilității prezentate în capitolul 1, în perspectiva conturării unui nou mod de abordare a proiectării sistemelor de calcul tolerante la defecte, prezentul capitol își propune o trecere în revista problematicii legate de:

- coduri detectoare și corectoare de erori pentru discuri;
- matricele discuri redundante și independente – RAID.

În această parte, teza cuprinde următoarele contribuții:

- 1) sistematizarea metodelor generale de distribuție a datelor în RAID;
- 2) sistematizarea mecanismelor de redundanță utilizate în RAID.

3 O ANALIZĂ A ARHITECTURILOR RAID

3.1 RAID level 0

RAID 0 nu este chiar un membru din familia RAID deoarece nu utilizează nici un fel de redundanță în scopul creșterii fiabilității. Totuși, există aplicații în care performanța, capacitatea și costul sunt determinante, fiabilitatea ocupând un loc secundar.

În RAID 0 datele sunt distribuite pe mai multe discuri din matrice, astfel încât pot fi deservite în paralel atât mai multe cereri de dimensiuni mici cât și cereri de dimensiuni mari.

3.1.1 Distribuția datelor în RAID 0

Datele sunt împărțite în secțiuni (blocuri) de mărime un multiplu al mărimii unui sector, și sunt distribuite pe toate discurile. Acest mod de distribuție aduce marele avantaj ca o cerere de mărime n secțiuni, unde n este numărul de discuri din RAID, poate fi deservită simultan de toate discurile, obținându-se o rată de transfer de aproape n ori mai mare decât în cazul utilizării unui singur disc (Figura 3-1, unde cu S_i s-a notat secțiunea i).

| DISC 0 | DISC 1 | DISC 2 | DISC 3 |
|--------|--------|--------|--------|
| S0 | S1 | S2 | S3 |
| S4 | S5 | S6 | S7 |
| S8 | S9 | S10 | S11 |
| ... | ... | ... | ... |

Figura 3-1

3.1.2 Comportarea RAID 0 în aplicații care solicită rate ridicate de transfer

Performanța oricărui nivel de RAID depinde, în principal, de tipul de cereri de acces și de modelul de distribuție a datelor. Acest fapt este foarte clar la RAID 0, unde mecanismele de redundanță nu există, deci nu vor influența analiza noastră.

Mai întâi vom considera că se folosește un RAID 0 pentru atingerea unei rate mari de transfer, fapt ce presupune îndeplinirea următoarelor cerințe:

- Întreaga cale de comunicație între memoria principală și discuri trebuie să poată susține nivelul cerut al ratei de transfer (magistralele memoriei și I/O, controlerul I/O etc.)
- Aplicația trebuie să lanseze cereri care să ducă la o utilizare eficientă a RAID, ceea ce, în acest caz, înseamnă că cererea tipică se va referi la mari cantități de date contigue din punct de vedere logic, care să fie mult mai mari decât mărimea secțiunii utilizate, astfel încât să poate fi valorificată posibilitatea accesului în paralel la discurile din matrice.

3.1.3 Comportarea RAID 0 în aplicații ce solicită rate ridicate de cereri de transfer

În aplicații orientate pe tranzacții, devine mult mai interesant timpul de răspuns decât rata de transfer. Pentru o cerere individuală, de volum mic, dominante sunt întârzierile datorate componentelor mecanice, în mișcare ale discurilor.

O matrice RAID 0 poate asigura o rată ridicată de execuții de cereri I/O, prin distribuirea acestora în paralel și independent, pe toate discurile.

3.2 RAID level 1

În RAID 1 se folosește redundanța prin simpla duplicare a datelor. Aceasta aduce unele avantaje dar și dezavantaje.

Principalele avantaje constau în:

- cererea de citire poate fi deservită de către oricare dintre discurile care stochează data cerută, de obicei utilizându-se cel care asigură timpii căutare (*seek*) și rotație cei mai mici;
- recuperarea datelor pierdute prin căderea unui disc, se poate executa simplu, în timp real, prin utilizarea copiilor;
- dezavantajul, legat de performanță, pentru RAID 1 constă într-o ușoară scădere la scriere, care se va executa ținând cont de timpii maximi de acces și rotație pentru discurile implicate.

Din punct de vedere economic, RAID 1 este cea mai costisitoare soluție RAID, necesitând o capacitate reală de stocare dublă față de cea utilă.

3.2.1 RAID 1 și variantele sale arhitecturale

În acest subcapitol vom descrie arhitectura RAID 1, cunoscută și sub numele de *mirrored array* (matrice de discuri „în oglindă“). Vom considera diverse variante care diferă una de cealaltă după modul în care sunt plasate datele pe discuri. Aceste variante sunt *mirrored declustering*, *chained declustering* și o variantă nouă, *goup-rotate declustering*, care combină avantajele primelor doua metode.

3.2.1.1 Mirrored declustering

În acest model cele N discuri sunt configurate ca $N/2$ perechi de discuri „în oglindă“, perechea i fiind denumită *mirror i*. Două copii ale datelor sunt plasate de-a lungul acestor $N/2$ perechi. Fiecare pereche conține aceleași informații, nefiind neapărat necesară plasarea în aceeași locație pe cele două discuri. Pentru ușurința discuției, presupunem că cele două copii sunt plasate în perechi la un deplasament de m cilindri. Când deplasamentul este zero, perechile sunt „oglundite fizic“. Blocuri logice contigue ale celor două copii sunt alocate secvențial pe secțiunile celor $N/2$ perechi de discuri. *Figura 3-2* prezintă un exemplu, în care fișierul f_i , compus din șase blocuri, este plasat pe o matrice formată din opt discuri. În figură s-a notat cu F_i prima copie, și cu f_i a doua copie a celui de-al i -lea bloc al fișierului. O cerere de citire poate fi satisfăcută de oricare copie, dar o cerere de scriere necesită ca fiecare copie să fie reactualizată. Unul dintre dezavantajele acestei arhitecturi este acela când defectându-se un disc și aria de discuri lucrează într-unul dintre modurile *failure* sau *rebuild*, tot traficul inițial orientat spre discul defectat este redirectat spre „oglanda“ lui, situație care poate deveni fără ieșire.

| Sector | <i>Mirror 1</i> | | <i>Mirror 2</i> | | <i>Mirror 3</i> | | <i>Mirror 4</i> | |
|--------|-----------------|-------|-----------------|--------|-----------------|--------|-----------------|--------|
| 0 | F1 | | F2 | | F3 | g1 | F4 | g2 |
| 1 | F5 | | F6 | | | | | |
| ... | | | | | | | | |
| I | | f1 | | f2 | | f3 | | f4 |
| I+1 | | f5 | | f6 | G1 | | G2 | |
| | Disc 0 | Disc1 | Disc 2 | Disc 3 | Disc 4 | Disc 5 | Disc 6 | Disc 7 |

Figura 3-2

3.2.1.2 Chained declustering

Pentru a depăși anteriorul dezavantaj pentru recuperarea informațiilor, Hsiao și DeWitt au propus arhitectura *chained declustering*, în care sunt păstrate două copii fizice ale datelor, denumite *primary copy* (copie primară) și *backup copy* (copie de rezervă). După cum se arată în Figura 3-3, dacă un bloc de date primar este alocat pe discul i , blocul său de rezervă este alocat pe discul $\{(i + 1) \bmod N\}$. Ca și în cazul *mirrored declustering*, o citire este satisfăcută de oricare dintre copii, în vreme ce actualizările necesită accesul la ambele copii.

| Bloc | | | | | | | | |
|-------|--------|--------|--------|--------|--------|--------|--------|--------|
| 0 | F1 | F2 | F3 | F4 | F5 | F6 | F7 | F8 |
| 1 | F9 | F10 | F11 | F12 | F13 | F14 | F15 | F16 |
| ... | | | | | | | | |
| i | f8 | f1 | f2 | f3 | f4 | f5 | f6 | f7 |
| $i+1$ | f16 | f9 | f10 | f11 | f12 | f13 | f14 | f15 |
| | Disc 0 | Disc 1 | Disc 2 | Disc 3 | Disc 4 | Disc 5 | Disc 6 | Disc 7 |

Figura 3-3

3.2.1.3 Group-rotate declustering

Profitând atât de *mirrored declustering*, cât și de *chained declustering*, se propune o a treia arhitectură, numită *group-rotate declustering*, după cum se arată în Figura 3-4. În această arhitectură, prima copie este stocată în același mod ca la *mirrored declustering*, dar a doua copie este rotită printre cele $N/2$ discuri rămase. În mod clar, *group-rotate declustering* poate efectua două citiri „mari“ în același timp precum anteriorul *mirrored declustering*, dar în timpul unei defecțiuni distribuie povara discului defectat în mod egal pe cele $N/2$ discuri. Dezavantajul acestei metode este că, în cazul defectării a doua discuri, este mai probabilă pierderea datelor decât în cazurile *mirrored/chained declustering*.

| | | | | | | | |
|---------|-----|-----|-----|---------|----|-----|-----|
| F1 | F2 | F3 | F4 | f1 | f5 | f9 | f13 |
| F5 | F6 | F7 | F8 | f2 | f6 | f10 | f14 |
| F9 | F10 | F11 | F12 | f3 | f7 | f11 | f15 |
| F13 | F14 | F15 | F16 | f4 | f8 | f12 | f16 |
| Copie 1 | | | | Copie 2 | | | |

Figura 3-4

O altă arhitectură, *interleaved declustering*, a fost propusă în literatura de specialitate. Deoarece această arhitectură nu este propice aplicațiilor în care cererile I/O necesită actualizarea unei mari cantități de date, o vom exclude din discuția prezentă.

3.2.1.4 Amplasarea copiilor pe discuri

Există câteva variante ale acestei arhitecturi depinzând de modul în care cele două copii ale datelor sunt plasate pe disc. Vom studia două dintre acestea, după cum se arată în Figura 3-5. În mod obișnuit, axa unui disc reprezintă o colecție de platan. Fiecare conține la rândul său un anumit număr de piste. Un cilindru este totalitatea pistelor aflate în aceeași poziție pe fiecare platan. Pentru discurile unei stații de lucru, de obicei, capetele de citire/scriere sunt atașate unui dispozitiv care poziționează capetele pe pista potrivită, și doar unul din capete transferă date la un moment dat.

Cele două abordări amintite sunt descrise după cum urmează. În ambele cazuri blocurile care aparțin aceleiași copii a fișierului aflate pe același disc sunt stocate continuu.

- *Horizontal Layout* (HL)

Jumătate din platanele de pe fiecare axă sunt folosite pentru a stoca prima copie, iar restul pentru a stoca cea de a doua copie. Cele două copii ocupă aceiași cilindri.

- *Vertical Layout* (VL)

Cilindrii sunt partiționați în două zone continue și de dimensiune egală. Cele două copii sunt stocate în partiții diferite pe cilindrii la aceeași poziție relativă în cadrul partițiilor.

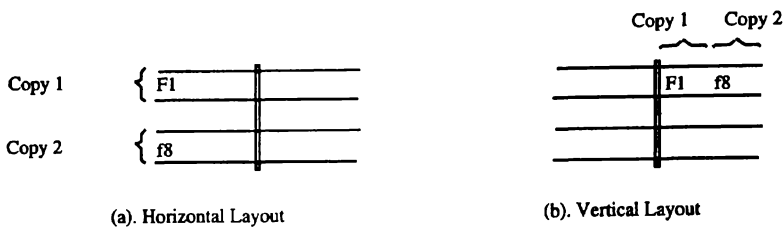


Figura 3-5

De exemplu, considerăm axa unui disc formată din patru platane numerotate de la 0 la 3 și 1000 de cilindri. Abordarea orizontală plasează copia 1 pe platanele 0 și 1, ocupând toate pistele acestor două platane, iar copia 2 pe platanele 2 și 3. Blocurile aparținând aceluiași fișier sunt stocate continuu, exemplu: în cazul Figura 3-4, blocul F9 este adiacent cu F1. Pe deasupra, cele două copii se află și pe același cilindru, exemplu: f8 este pe același cilindru cu F1 (de remarcat că f1, copia identică a lui F1, este pe un disc separat, dar la aceeași locație ca f8). Pe de altă parte, abordarea verticală plasează copia 1 pe cilindrii 0 până la 499, cuprinzând toate pistele din fiecare cilindru, iar copia 2 pe cilindrii 500 până la 999 (presupunem că fiecare pistă are aceeași capacitate). În cazul de mai sus, F1, F9, ..., încep cu cilindrul 0, și f8, f16, ..., încep cu cilindrul 500 (de remarcat că această reprezentare nu presupune ca spațiul de pe disc să fie folosit complet, deoarece fiecare grup poate avea cilindri liberi pentru a fi folosiți ulterior). Abordarea HL favorizează scrierile, în timp ce abordarea VL favorizează citirile la cantități mari de date.

Se observă că *mirrored declustering* și *chained declustering* diferă în felul de tratare a cererilor de citire pentru cantități mari de date; *mirrored declustering* poate efectua două astfel de citiri „mari” concurrent, în vreme ce *chained declustering* poate efectua doar o singură citire la un moment dat (presupunând că există un singur braț pe disc). Totuși, când un disc se defectează, *chained declustering* poate distribui fără probleme „povara discului” defectat către toate cele $N - 1$ discuri operaționale rămase.

3.2.1.5 Un model de disc

Prezentăm modelul de disc care va fi folosit în discuțiile ulterioare. Ne vom concentra asupra timpului de răspuns I/O, care este timpul scurs între sosirea cererii către subsistem și răspunsul subsistemului.

În mod obișnuit, timpul de răspuns I/O constă în patru componente:

- întârzierea „cozii“ la controler (*queuing delaz*);
- timpul de căutare (*seek time*);
- latentă de rotație (*rotational latency*);
- timpul de transfer al datelor (*data transfer time*).

Presupunem că fiecare dispozitiv are propria sa cale a datelor. Ne vom concentra doar asupra *queuing delay*, care este în funcție de încărcarea I/O, și timpul de acces la disc, care constă în timpul de căutare, de rotație și de transfer.

Se definesc:

- C : numărul de cilindrii pentru fiecare disc;
- N_b : numărul de blocuri de pe fiecare pistă;
- τ : timpul de transfer pentru un singur bloc;
- S_{max} : timpul maxim de căutare;
- R_{max} : timpul unei rotații complete;
- D : o variabilă aleatoare (*random variable - r.v.*) care definește distanța căutării pe disc;
- V : r.v.; arată dacă un acces este secvențial ($V = 0$) sau nu ($V = 1$);
- p_s : probabilitatea ca distanța de căutare să fie egală cu zero, $p_s = P\{V=0\}$;
- S : r.v.; indică timpul de căutare pe disc;
- R : r.v.; indică latentă de rotație;
- X : r.v.; indică suma timpilor de căutare și de rotație, $X=S+R$;
- B : r.v.; indică numărul de blocuri care trebuie transferați la o cerere;
- T : timpul de transfer al datelor, $T=\tau B$;
- Y_r, Y_w, Y : r.v.; indică timpii de citire, scriere și total de pe disc;
- Z_r, Z_w, Z : r.v.; indică timpii de răspuns I/O pentru citire, scriere și total.

În primul rând să considerăm modelul căutării pe disc. Deoarece s-a observat în realitate că, în majoritatea timpului, brațul discului este nemișcat, vom introduce o probabilitate a accesului secvențial, p_s , care

se definește prin probabilitatea ca distanța de căutare pe disc să fie egală cu zero. Identificăm două tipuri de accese:

- accese secvențiale;
- accese nesevențiale.

La accese secvențiale brațul discului nu se mișcă, în vreme ce la accesese nesevențiale brațul trebuie să se miște pentru a răspunde cererii. În acest ultim caz se consideră că brațul se mișcă pe oricare alt cilindru cu probabilitate egală.

Conform cu definiția lui V , avem în continuare distribuțiile condiționate pentru distanța de căutare:

$$P\{D = i | V = 0\} = \begin{cases} 1 & i = 0 \\ 0 & i = 1, 2, \dots, C-1 \end{cases}$$

$$P\{D = i | V = 1\} = \begin{cases} 0 & i = 0 \\ \frac{2(C-i)}{C(C-1)} & i = 1, 2, \dots, C-1 \end{cases}$$

unde termenul $2(C-i)/[C(C-1)]$ rezultă din presupunerea că brațul și cilindrul țintă sunt localizați aleator pe disc și că sunt distincți.

Îndepărtarea câmpurilor de condiție conduce la:

$$P\{D = i\} = \begin{cases} p_s & i = 0 \\ (1-p_s) \frac{2(C-i)}{C(C-1)} & i = 1, 2, \dots, C-1 \end{cases}$$

Se folosește, de asemenea, și următoarea expresie pentru timpul de căutare S , ca funcție de distanța de căutare D :

$$S = \begin{cases} 0 & D = 0 \\ a + b\sqrt{D} & D > 0 \end{cases}$$

unde a este timpul de accelerare a brațului și b este factorul de căutare mecanic.

Pentru ușurința analizei, aproximăm că distanța de căutare D este continuă, cu funcția de densitate a probabilității (*pdf*):

$$f_D(x) = \begin{cases} p_s u_0(x) & x = 0 \\ (1-p_s) \frac{2(C-x)}{C^2} & 0 < x \leq C \end{cases}$$

unde $u_0(x)$ este funcția impuls unitar.

Funcția de distribuție cumulativă a lui S , $F_S(x) = P\{S < x\}$:

$$F_S(x) = \begin{cases} F_D(0) & 0 \leq x \leq a \\ F_D\left(\left(\frac{x-a}{b}\right)^2\right) & a \leq x \leq S_{\max} \\ 1 & S_{\max} < x \end{cases}$$

unde timpul de căutare maxim este $S_{\max} = a + b\sqrt{C-1}$. Timpul de căutare mediu poate fi:

$$E[S] \approx (1-p_s) \left[a + b\sqrt{C-1} \frac{8}{15} \right]$$

Timpul de rotație se presupune a fi uniform distribuit în intervalul $[0, R_{\max}]$ cu pdf:

$$f_R(x) = 1/R_{\max} \quad 0 \leq x \leq R_{\max}$$

Fie $X = S + R$ suma timpilor de căutare și de rotație, atunci pdf-ul lui X este:

$$f_X(x) = \int_0^{x-R_{\max}} f_R(x-s) dF_S(s)$$

Pentru parametrii care ne interesează (a se vedea sfârșitul acestei părți),

$R_{\max} \leq b\sqrt{C-1}$ și $f_X(x)$ ia următoarea formă:

$$f_X(x) = \begin{cases} f_R(x)p_s + \int_0^x f_R(x-s) dF_S(s) & 0 \leq x \leq R_{\max} \\ \int_0^x f_R(x-s) dF_S(s) & R_{\max} < x \leq R_{\max} + a \\ \int_{x-R_{\max}}^a f_R(x-s) dF_S(s) & R_{\max} + a < x < S_{\max} \\ \int_{x-R_{\max}}^{S_{\max}} f_R(x-s) dF_S(s) & S_{\max} < x \leq S_{\max} + R_{\max} \end{cases}$$

Timpul de servire a discului pentru o cerere obișnuită este:

$$Y = X + T$$

cu media:

$$E[Y] = E[S] + E[R] + E[T] = E[X] + E[T]$$

unde T este timpul de transfer, $T = \tau B$.

În cele din urmă, parametrii discului folosiți în studiul nostru sunt rezumați în cele ce urmează: numărul de cilindri $C=12000$; rata de

transfer = 3MB/sec; timpul rotației complete $R_{\max} = 16.7\text{ms}$; dimensiunea blocului = 4096 octeți; numărul de blocuri per pistă $N_b=12$; timpul de accelerare $a = 3\text{ms}$; factorul de căutare $b=0.5$. Astfel, timpul necesar transferului unui singur bloc este $\tau = 1.3\text{ms}$.

3.2.2 Strategii de accesare pentru RAID 1

Deoarece există două copii ale datelor stocate în matrice, o cerere de citire poate fi satisfăcută de oricare dintre copii, în vreme ce o cerere de scriere trebuie executată asupra ambelor copii. În continuare, ne vom concentra asupra a trei strategii. Pentru toate aceste strategii descrise mai jos, sunt folosite de către sistem două cozi, câte una pentru fiecare copie (Figura 3-6). Fiecare coadă este servită după modelul primul sosit - primul servit. În momentul sosirii, o cerere de scriere generează două sarcini care implică ambele cozi. Diversele strategii diferă între ele prin modul în care se face alegerea discului pentru citirea datelor.

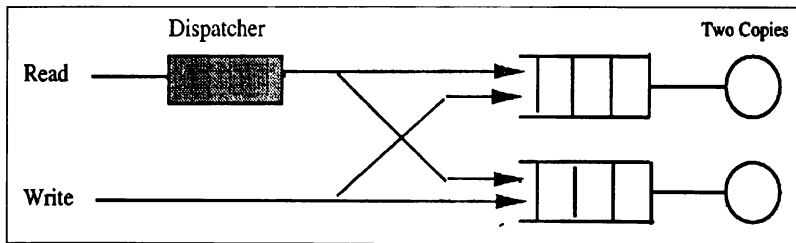


Figura 3-6

3.2.2.1 Strategia Random Join (RJ)

În cadrul acestei strategii, o cerere de citire este îndreptată în momentul sosirii în mod aleator spre coada asociată discului i având probabilitatea:

$$\alpha_i, i = 1, 2 (\alpha_1 + \alpha_2 = 1).$$

3.2.2.2 Strategia Shortest Queue (SQ)

În cadrul acestei strategii, o cerere de citire selectează discul cu cea mai scurtă coadă în momentul sosirii cererii. Legăturile sunt rupte într-o manieră arbitrară. Când cererea de citire sosește în cadrul unui sistem gol, ea este îndrumată către discul al cărui braț este mai apropiat de cilindrul care conține blocul dorit.

3.2.2.3 Strategia Minimum Seek (MS)

În momentul sosirii unei cereri de citire, strategia MS verifică adresa cilindrului corespunzător ultimei cereri și îndrumă cererea către coada

cu cea mai mică distanță de căutare. Când fiecare coadă este tratată în maniera primul sosit - primul servit, această strategie reduce timpul de căutare pentru cererile de citire.

Remarcă:

Dintre toate aceste strategii, strategia RJ este cea mai potrivită sistemelor în care cele două discuri conținând cele două copii țintă sunt localizate fizic în două poziții îndepărtate, deoarece nu necesită schimbare de informații între cele două discuri. Celelalte strategii, însă, sunt mai potrivite pentru sistemele în care toate discurile din matrice sunt localizate apropiat unele de altele deoarece acestor strategii trebuie să se cunoască starea fiecărei cozi pentru a programa cererile de citire.

3.2.3 Analiza comportării tipurilor arhitecturale RAID 1 cu diverse strategii

În acest subcapitol vom studia și compara RAID 1 și variantele sale, *mirrored declustering*, *chained declustering* și *group-rotate declustering*, în combinație cu strategiile RJ, SQ și MS. În cartea lor, Kim și Tantawi raportează existența câtorva direcții reale de referință a discurilor, care arată că procentajul acceselor secvențiale se întinde între 25% și 50%. Pentru matricele de discuri, deoarece datele sunt depuse pe mai multe discuri, probabilitatea de acces secvențial la fiecare disc este de așteptat să fie mai mică decât cea a discurilor din sistemele convenționale. De aici încolo vom alege $p_s=0.2$ în studiul de față. Vom considera și cazul $p_s=0.5$ și vom observa comportări similare.

Vom începe cu o sumară privire de ansamblu despre cum anume este simulată aria de discuri. În momentul sosirii unei cereri, simulatorul alege o coadă pentru a adăuga datele în concordanță cu strategia folosită. De exemplu, în cazul folosirii strategiei SQ împreună cu *mirrored declustering*, simulatorul selectează în primul rând aleator un disc dintre $N/2$ discuri, să spunem i , $0 \leq i < N/2$, și selectează perechea discului ca fiind $i+N/2$. Dacă cererea este de citire, ea este îndrumată către cea mai scurtă coadă dintre i sau $i+N/2$. La metoda *group-rotate declustering*, discul pereche este selectat aleator dintre cele $N/2$ discuri. În cazul *chained declustering*, primul disc i este selectat aleator dintre toate cele N discuri, iar discul pereche este $(i+1) \bmod N$.

Deoarece cele două blocuri de date duplicate pentru a fi accesate pot fi stocate în diferite locații pe discurile pereche, în loc de a păstra informații despre locația blocului de date, pentru fiecare cerere, simulatorul generează aleator o adresă a unui cilindru pentru un disc

(presupunând că orice bloc este în mod egal accesat) și apoi obține adresa cilindrului pereche prin adăugarea unui deplasament. În cazul *chained declustering* cu VL, prima adresă de cilindru este selectată cu probabilitate egală dintre cilindrul 1 până la $C/2$, iar deplasamentul este $C/2$. Pentru alte metode, prima adresă a cilindrului este aleasă cu o probabilitate egală dintre toți cilindrii. De asemenea, un deplasament de $C/2$ este utilizat pentru *mirrored* și *group-rotate declustering*. Am încercat și un deplasament egal cu zero, dar diferențele în performanțe au fost minore.

În cele din urmă, deoarece există o mică diferență între performanțele obținute prin *chained declustering* cu VL și HL (cu VL evoluând puțin mai bine), vom considera doar performanțele modelului *Vertical Layout*.

3.2.3.1 RAID1 cu strategie RJ

Performanțele sistemelor RAID 1 și ale variantelor acestuia, în cazul strategiei RJ, sunt prezentate în Figura 3-7_pentru cazurile în care probabilitatea citirii este de $p_r=0.75$, respectiv $p_r=0.25$. După cum se vede în figură, când majoritatea cererilor sunt citiri, toate cele trei strategii prezintă performanțe similare. În cazul în care majoritatea cererilor sunt scrieri, *mirrored declustering* furnizează performanțe ușor crescute în comparație cu celelalte două metode, acest lucru deoarece toate datele de pe un disc sunt copiate pe un singur disc. Timpii de căutare corespunzători unei cereri de scriere pe cele două discuri sunt identici atunci când aceasta cerere urmează unei alte cereri de scriere.

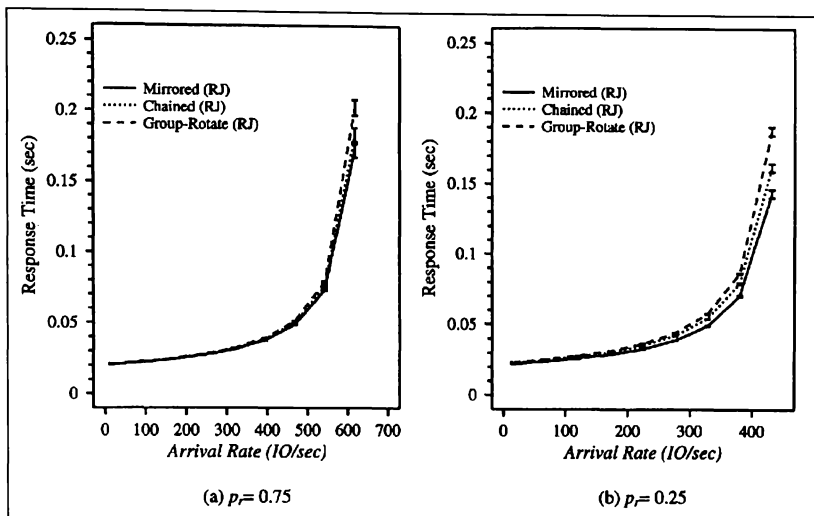


Figura 3-7

În cazul *group-rotate declustering* este, totuși, improbabil ca timpii de acces la cele două discuri să fie identici deoarece anterioara operație de scriere a fost făcută cel mai des doar pe unul din discuri, nu pe ambele discuri. Deoarece operația de scriere este considerată încheiată doar după ce se termină operațiile pe ambele discuri, timpul de răspuns este mai scăzut pentru metoda *group-rotate declustering* decât pentru *mirrored declustering*. Comportamentul în cazul *chained declustering* se află între cele două metode discutate anterior.

3.2.3.2 RAID1 cu strategie SQ

Când se folosește strategia SQ, cele trei variante prezintă comportamente foarte diferite (Figura 3-8). Când majoritatea cererilor sunt citiri, se observă că *group-rotate declustering* are rezultatele cele mai bune, iar *mirrored declustering* cele mai slabe. Motivul este că, în cadrul SQ, *mirrored declustering* furnizează o încărcare echilibrată pentru citiri între cele două cozi, iar fiecare pereche oglindită este independentă de celelalte perechi. *Group-rotate declustering* reproduce datele de pe fiecare disc din cadrul a $N/2$ discuri. Prin urmare, strategia SQ are abilitatea de a echilibra încărcările pe cele $N/2$ cozi.

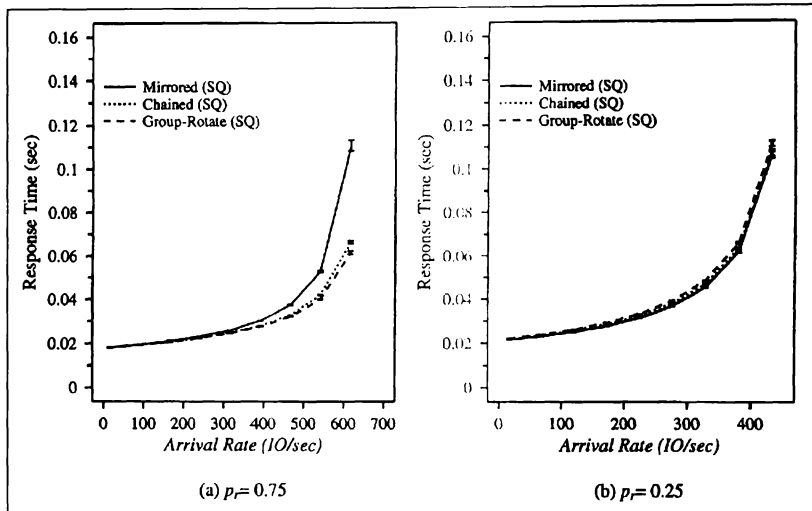


Figura 3-8

Să considerăm cazul *chained declustering*. Deși datele sunt secționate pe toate cele N discuri, conținutul fiecărui disc este duplicat doar către cei doi vecini ai săi. Prin urmare, dacă apare un dezechilibru pe unul dintre discuri, problema poate fi ușor rezolvată de către unul dintre vecini.

Deși cererile de citire beneficiază încă de o încărcare echilibrată când majoritatea cererilor sunt scrieri, timpii de răspuns la cererile de scriere sunt mai mari decât cei în cazul *mirrored* sau *chained declustering*, din motivele enumerate în anteriorul subcapitol. Prin urmare, există doar mici diferențe între performanțele totale în cadrul celor trei variante.

3.2.3.3 RAID1 cu strategie MS

În cazul strategiei MS, se observă un comportament similar pentru cele trei variante, la fel ca în cazul SQ, (metoda *group-rotate declustering* are performanțele cele mai bune când majoritatea cererilor sunt citiri, iar *mirrored declustering* are performanțele cele mai bune când majoritatea cererilor sunt scrieri).

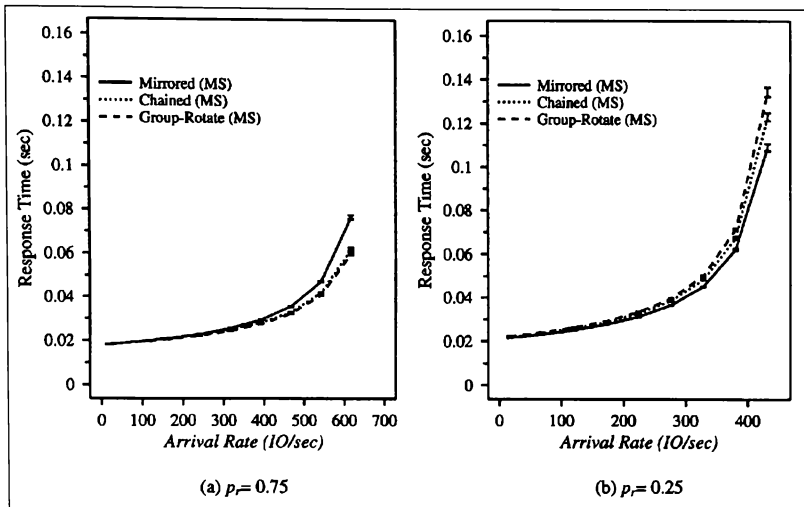


Figura 3-9

3.2.3.4 Analizarea diverselor strategii în cazul Group-Rotate Declustering

Încheiem acest subcapitol prin compararea diverselor strategii în cazul arhitecturii *group-rotate declustering*. Pentru rezultatele prezentate în Figura 3-10 utilizarea discului variază aproximativ de la 0.02 la 0.93. Din figură se observă ca RJ este în mod semnificativ mai neperformant decât celelalte două. Când majoritatea cererilor sunt citiri, sunt doar mici diferențe între strategiile MS și SQ. Când majoritatea cererilor sunt scrieri, SQ depășește clar MS. Totuși, deoarece timpii de servire a discului pentru cererile de citire în cazul MS sunt mai scăzuți decât în cazul SQ, rezultatul direct va fi mai ridicat în cazul MS.

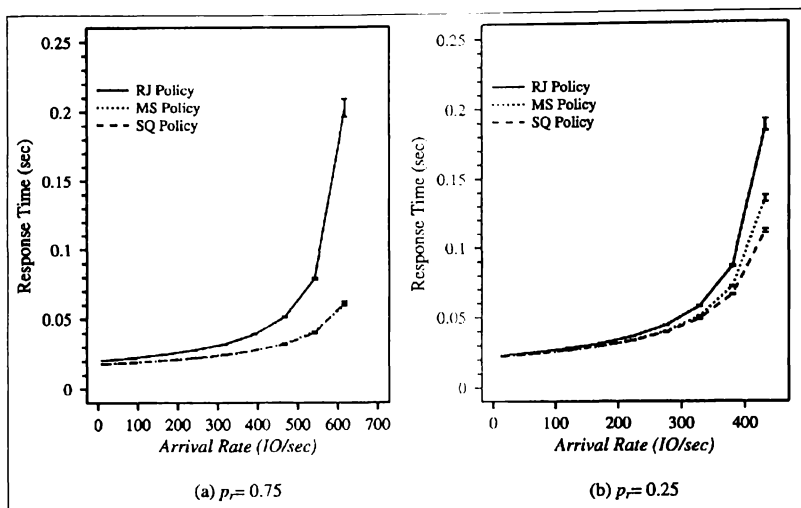


Figura 3-10

3.2.3.5 O analiză a metodei Chained Declustering

După cum s-a descris anterior, există două metode de a plasa datele în cadrul *chained declustering*: *horizontal layout (HL)* și *vertical layout (VL)*.

HL favorizează operațiile de scriere deoarece datele din aceeași secțiune aparținând celor două copii sunt localizate în același cilindru. Deci, o scriere masivă poate actualiza cele două copii prin mișcarea brațului către cilindrul destinație, actualizând prima copie, urmată apoi de o întârziere de rotație, iar apoi actualizând a doua copie (de reamintit că blocurile aparținând celor două copii aflate pe un disc nu sunt aceleași). Prin urmare, este necesară doar o singură căutare. Timpii de servire pentru citiri și scrieri sunt:

$$Y_r = X + T$$

$$Y_w = X + T + R + T$$

unde T este timpul de transfer și depinde de numărul blocurilor care trebuie transferate de pe fiecare disc.

Acest sistem poate fi modelat ca o coadă $M/G/1$. Timpul de servire este:

$$Y = \theta Y_r + (1 - \theta) Y_w$$

cu valorile momentane:

$$\bar{Y} = \bar{X} + \bar{T} + p_w(\bar{R} + \bar{T})$$

$$\bar{Y}^2 = \bar{X}^2 + 2\bar{X}\bar{T} + \bar{T}^2 + p_w(3\bar{T}^2 + \bar{R}^2 + 2\bar{X}\bar{R} + 2\bar{X}\bar{T} + 2\bar{R}\bar{T})$$

După formula *Pollaczek-Khinchin*, timpul mediu de așteptare în coadă este:

$$\bar{Q} = \frac{\lambda \bar{Y}^2}{2(1 - \rho)}$$

unde $\rho = \lambda \bar{Y}$ arată modul de folosire a dispozitivului.

În final, timpul mediu de răspuns I/O, \bar{Z} , este dat de

$$\bar{Z}_r = \bar{Q} + \bar{Y}_r$$

$$\bar{Z}_w = \bar{Q} + \bar{Y}_w$$

$$\bar{Z} = p_r \bar{Z}_r + p_w \bar{Z}_w$$

unde \bar{Z}_r și \bar{Z}_w sunt timpii medii de răspuns pentru cereri de citire, respectiv de scriere.

VL favorizează cererile de citire, deoarece cele două copii ale datelor sunt localizate pe cilindri diferiți, la o distanță de $C/2$ piste între copii, unde C este numărul total de cilindri ai fiecărui dispozitiv. O citire poate fi făcută din copia cea mai apropiată de poziția brațului curent. În orice caz, o citire nu caută niciodată pe mai mult de $C/2$ piste. În timp ce strategia VL reduce timpul de servire a citirilor, ea crește timpul total al operațiilor de scriere. Când începe o operație de scriere, brațul se mută pe copia cea mai apropiată de poziția sa curentă, actualizează copia și se mută apoi la $C/2$ piste pentru a actualiza și cealaltă copie. Astfel, timpii de servire a citirilor și scrierilor $Y_r = X + T$ și $Y_w = X + T + S_c + R + T$, unde $X = S + R$ este suma timpului de căutare și a latentei de rotație, iar S_c este timpul de căutare a $C/2$ piste.

Analog, putem modela sistemul printr-o coadă *M/G/1* și să calculăm valorile momentane ale timpului de servire total prin

$$\bar{Y} = \bar{X} + \bar{T} + p_w(\bar{R} + \bar{T} + S_c)$$

$$\bar{Y}^2 = \bar{X}^2 + 2\bar{X}\bar{T} + \bar{T}^2 + p_w(3\bar{T}^2 + \bar{R}^2 + 2\bar{X}\bar{R} + 2\bar{X}\bar{T} + 4\bar{R}\bar{T} + S_c(S_c + 2\bar{X} + 2\bar{R} + 4\bar{T}))$$

Performanțele celor trei variante de RAID 1 sunt prezentate aici făcând uz de modelul anterior prezentat *group-rotate declustering* și de

simulările pentru modelele *mirrored* și *chained declustering*. Numărul de discuri din matrice a fost ales $N=64$. Dacă un disc are o capacitate de la 1G la 3G *bytes*, o matrice de discuri cu 64 de astfel de discuri poate atinge o capacitate de la 64G la 192G *bytes*.

Fiecare cerere se presupune că accesează N_s secțiuni plus o secțiune parțială, dimensiunea cerută este $B - N_s W + T_s$ blocuri, unde W este lățimea secțiunii. Presupunem că N_s este o distribuție geometrică cu parametrul q , iar T_s are următoarea distribuție:

$$P\{T_s = i\} = \begin{cases} p_f & i = 0 \\ (1 - p_f) \frac{1}{W - i} & i = 1, \dots, W - 1 \end{cases}$$

unde p_f este probabilitatea ca o cerere să execute o operație I/O pe întreaga secțiune (*full stripe I/O*). Cu alte cuvinte, fiind dată o cerere care conține și o secțiune parțială, dimensiunea acestei secțiuni parțiale este uniform distribuită în secțiune. Dimensiunea medie a cererii B este 256 blocuri (1M *bytes*). Atunci valorile momentane pentru N_s sunt:

$$\overline{N_s} = \frac{B - T_s}{W}$$

$$N_s^2 = \frac{2 - q}{q}$$

unde $\overline{T_s} = (1 - p_f)W/2$ și $q = 1/\overline{N_s}$.

Prin urmare, valorile momentane pentru timpul de transfer pe fiecare disc T , sunt:

$$\overline{T} = (\overline{N_s} + 1 - p_f)\tau$$

$$\overline{T^2} = \left[\overline{N_s^2} + (1 - p_f)(2\overline{N_s} + 1) \right] \tau^2$$

Studiind efectul pe care îl are p_f asupra performanțelor se observă că timpul mediu de răspuns este insensibil la modificările lui p_f în cazul RAID 1, acest lucru nefiind valabil și pentru RAID 5. Rezultatele prezentate în Figura 3-11 sunt obținute pentru $p_f=0.5$.

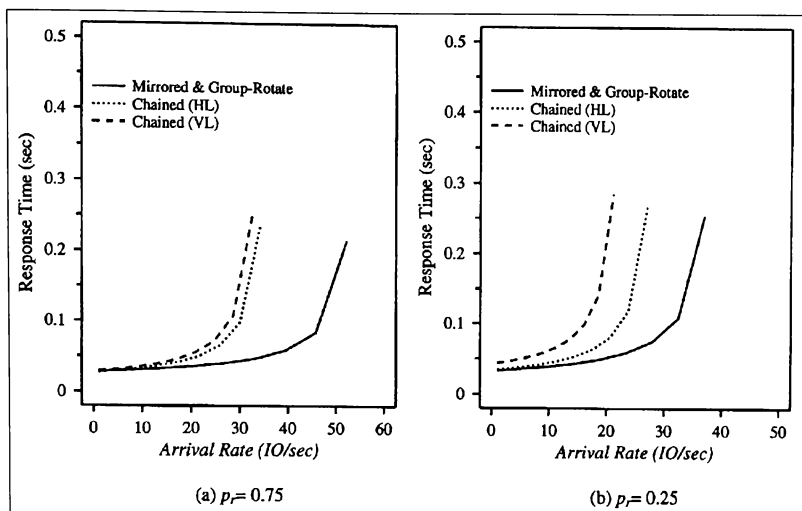


Figura 3-11

Figura 3-11 prezintă timpul mediu de răspuns ca o funcție a probabilităților de citire $p_r=0.75$ și $p_r=0.25$. În acest caz *mirrored* și *group-rotate declustering* se comportă la fel și mai bine decât *chained declustering*. Dintre cele două strategii de plasare a datelor, în cazul *chained declustering*, s-a observat ca HL are rezultate mai bune decât VL din cauza costului ridicat al scrierilor în cazul VL. O dată cu creșterea probabilității de citire, performanțele VL se apropie de cele ale HL, și s-a observat că la $p_r > 0.85$, VL depășește performanțele HL.

În Figura 3-12 și în Figura 3-13 sunt ilustrate procentajele maxime suportate de cele trei arhitecturi, drept funcție de dimensiunea medie a cererilor și de numărul de discuri din matrice. În aceste figuri se observă că *mirrored* și *group-rotate declustering* sunt mai performante decât *chained declustering*.

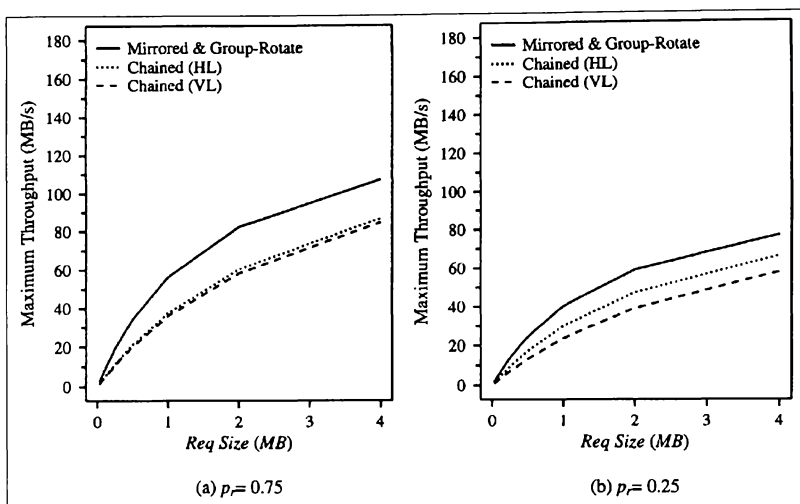


Figura 3-12

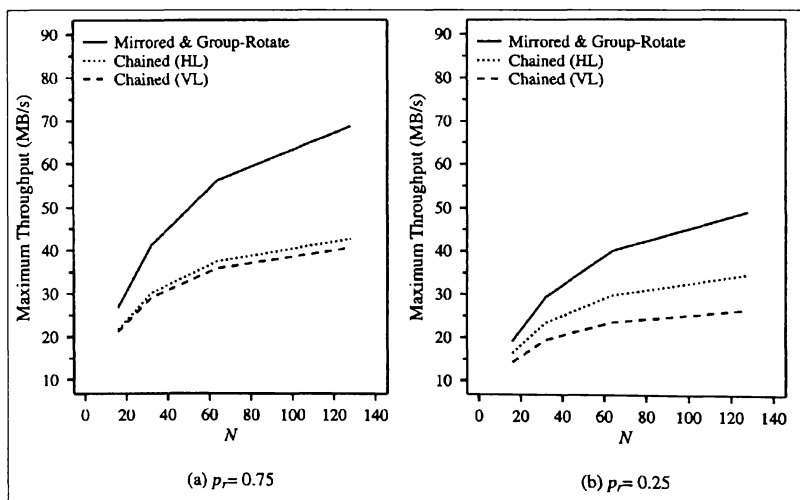


Figura 3-13

Se remarcă faptul că motivul principal pentru care performanțele *mirrored* și *group-rotate declustering* sunt superioare celor ale metodei *chained declustering* este că, pentru date secționare de-a lungul celor N discuri, *chained declustering* suportă în mod obișnuit o citire la un moment dat, în vreme ce celelalte suportă două citiri simultan. Însă,

chained declustering poate fi modificată pentru a putea efectua două citiri la un moment dat. Să considerăm o matrice cu 8 discuri. Dacă împărțim discurile 0, 2, 4, și 6 într-un grup, și discurile 1, 3, 5, 7 în alt grup, atunci fiecare dintre cele două grupuri conțin blocurile de date necesare pentru o operație „mare” I/O (de dimensiune mare). Prin urmare, cele două grupuri pot suporta două citiri de dimensiuni mari în mod concurrent. Această abordare crește, însă, complexitatea *software*-ului și a sistemelor de fișiere, deoarece fiecare bloc de date citit de pe fiecare disc în parte nu este în ordinea cerută și trebuie reordonat în *buffer*-ul memoriei principale. Se pot găsi și alte abordări pentru a îmbunătăți performanțele *chained declustering*, dar lasăm explorarea acestora ca un domeniu pentru viitoarele cercetări.

3.3 RAID level 2

RAID 2 utilizează o tehnică paralelă de acces la discurile din matrice, fiecare dintre acestea participând la satisfacerea fiecărei cereri. Datele sunt distribuite pe discurile din matrice în secțiuni foarte mici, de obicei de mărimea unui octet sau cuvânt. Pentru fiecare cuvânt stocat se calculează un cod ECC (uzual *Hamming*), care se depune într-o poziție omoloagă pe discuri dedicate.

RAID 2 este mai mult o soluție teoretică decât una practică.

3.4 RAID level 3

RAID 3 este organizată într-un mod similar cu RAID 2, cu deosebirea că se utilizează un singur disc suplimentar pentru memorarea informațiilor de control (biți de paritate).

În cazul căderii unui disc, datele de pe acesta se pot reconstrui utilizând datele omoloage existente pe discurile rămase valide. Imediat ce discul defect este înlocuit, datele pierdute se pot reconstrui.

Operația de reconstrucție este simplă. Pentru exemplificare voi considera un RAID cu patru discuri notate D_i , unde pentru $i=3$ considerăm discul cu informația de paritate calculată astfel:

$$D_3(k) = D_2(k) \text{ XOR } D_1(k) \text{ XOR } D_0(k)$$

Presupunând că discul D_1 a căzut, și adăugând $D_3(k) \text{ XOR } D_1(k)$ la ambii membri de mai sus se obține:

$$D_1(k) = D_2(k) \text{ XOR } D_1(k) \text{ XOR } D_3(k)$$

În concluzie, datele de pe orice secțiune de pe orice disc, în cazul unei căderi singulare, pot fi refăcute utilizând conținuturile secțiunilor de pe discurile rămase valide.

În cazul unei căderi de disc, toate datele rămân disponibile, matricea operând în modul redus, cu performanțe diminuate, în special la citire, când se regenerează datele pierdute. Revenirea la performanța maximă presupune înlocuirea discului căzut și reconstruirea datelor corespunzătoare pe acesta.

Din punct de vedere al performanței, datorită faptului că datele sunt distribuite în secțiuni foarte mici, RAID 3 asigură rate foarte bune de transfer. În cazul aplicațiilor orientate pe rate mari ale cererilor de transfer, performanța suferă datorită faptului că RAID 3 poate satisface o singură cerere la un moment dat.

3.5 RAID level 4

La RAID 4 se permit accese independente la discuri, astfel încât cereri separate pot fi satisfăcute în paralel. Din acest motiv, RAID 4 este mai potrivită pentru aplicații orientate pe număr mare de cereri de acces, decât pentru cele orientate pe transferuri masive.

Și la RAID 4 se utilizează distribuția datelor pe discuri diferite, dar secțiunile sunt de dimensiuni relativ mari. Protecția este asigurată prin paritate, stocată pe un disc dedicat.

Din punct de vedere al performanței, RAID 4 prezintă o scădere în cazul scrierilor de volum mic, deoarece controlerul de matrice trebuie să actualizeze atât datele, cât și biții de paritate. Considerând o matrice de patru discuri, notate cu D_i , unde D_3 este de paritate, rezultă:

$$D_3(k) = D_2(k) \text{ XOR } D_1(k) \text{ XOR } D_0(k)$$

După actualizarea unei secțiuni, de exemplu pe discul D_1 , și notată ', avem:

$$\begin{aligned} D'_3(k) &= D_2(k) \text{ XOR } D'_1(k) \text{ XOR } D_0(k) = \\ &= D_2(k) \text{ XOR } D'_1(k) \text{ XOR } D_0(k) \text{ XOR } D_1(k) \text{ XOR } D_1(k) = \\ &= D_3(k) \text{ XOR } D'_1(k) \text{ XOR } D_1(k) \end{aligned}$$

Rezultă că pentru calcularea noii parități, controlerul de matrice trebuie să citească vechea paritate și vechea dată, după care să scrie noua dată și paritate. Deci, fiecare scriere logică implică două citiri și două scrieri fizice.

În cazul scrierilor de volum mare, citirile nu mai sunt necesare, noua paritate calculându-se din noile date, astfel că discul de paritate poate fi accesat în paralel cu cele de date și deci nu mai sunt necesare citiri, respectiv scrieri, suplimentare.

În orice tip de scriere, însă, discul de paritate este implicat, ceea ce uneori poate să genereze întârzieri.

3.6 RAID level 5

În general, RAID 5 este organizată într-un mod asemănător cu RAID 4, cu deosebirea că și paritatea este distribuită pe toate discurile, eliminându-se, astfel, posibilele gâtuiuri întâlnite la RAID 4.

RAID 5, fiind una dintre matricele de discuri cele mai folosite (alături de RAID 1), cercetările mele au fost orientate în special către acestea.

3.6.1 Strategia AR pentru RAID 5

În acest subcapitol, se prezintă o strategie de programare sincronizată a I/O, strategia *after read-out* (AR), pentru RAID 5 care se potrivește aplicațiilor în care cererile sunt făcute pentru cantități mici de date. Această strategie susține problema sincronizării scrierilor, descrisă anterior. În toate cazurile, o cerere de scriere se consideră încheiată doar când, atât datele cât și paritatea, au fost actualizate.

Descriem mai întâi strategia AR pentru cazul când o cerere de scriere necesită actualizarea unui singur bloc. În cazul acestei strategii sunt păstrate două cozi pentru fiecare disc din matrice, una pentru sosirile de cereri de citire/scrieri de date (coada D) și alta pentru cererile de actualizare a parității (coada P), după cum se ilustrează în Figura 3-14. Când sosește o cerere de citire sau de scriere, dispecerul o trimite în coada D a discului țintă.

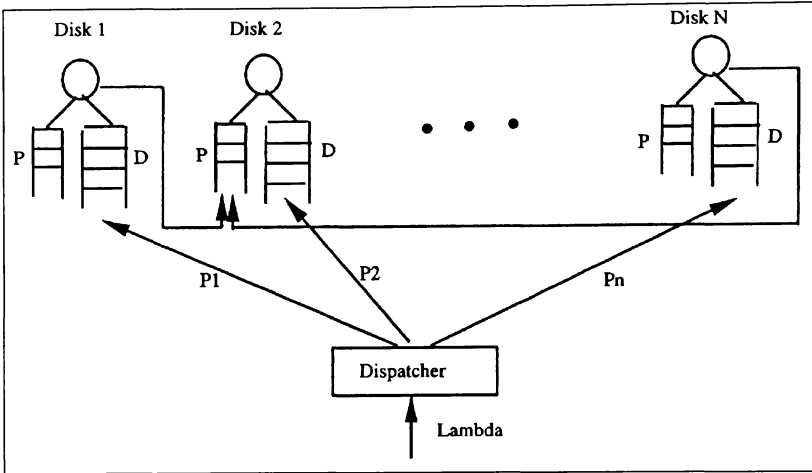


Figura 3-14

Pentru a servi o cerere de scriere, cererile de actualizare a datelor și a parității vor urma următoarele stadii: căutare, rotație, citire (pentru calculul noii parități), rotație completă, și scris. Figura 3-15 prezintă scenariul pentru servirea unei cereri de citire. Când se atinge capul cozii și este programat pentru servire, se generează o cerere de actualizare a parității către coada P a discului corespunzător care conține blocul (blocurile) de paritate, după ce vechile date au fost citite. De aceea, această strategie se numește *after read-out* („după citire“). Cererile din coada P au o prioritate mai ridicată decât cele din coada D pentru a fi siguri că o cerere de paritate mai deosebită este servită cel mai curând posibil, deoarece corespunzătoarea operație de actualizarea datelor este deja începută. Operațiile de disc se presupun a fi non-preemptive, deci o cerere de paritate nu poate fi servită până când operația I/O curentă nu ia sfârșit.

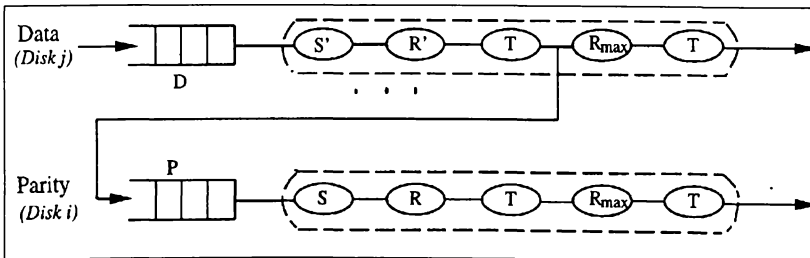


Figura 3-15

Dacă o cerere de scriere dorește să modifice mai mult decât un bloc dintr-o secțiune, atunci cererea de paritate este generată când ultimul bloc vechi de date a fost citit.

3.6.2 Strategia BS pentru RAID 5

O altă strategie, numită BS, a fost propusă. Aceasta generează o cerere de actualizare a parității imediat după ce o scriere își începe execuția. Am observat diferențe minore între performanțele acestor două strategii (strategia BS se comportă mai bine la încărcări scăzute). Totuși, strategia AR este mai simplă de implementat, în special în sisteme distribuite în care discurile care conțin blocurile de date și de paritate pot fi plasate în diferite locuri. Într-un astfel de mediu, strategia AR necesită trimiterea unui singur mesaj unui disc de paritate îndepărtat care conține datele vechi și comenzile, în vreme ce strategia BS necesită transmiterea a doua mesaje separate la momente diferite, unul conținând comenzile, iar celalalt conținând vechile date.

3.6.3 Tipuri semnificative de RAID 5

3.6.3.1 RAID5/UPGD

Câteva tipuri noi de matrice de discuri au fost recent propuse în care grupurile de paritate au o distribuție uniformă, astfel încât încărcarea suplimentară creată de căderea unui disc să poată fi împărțită în mod egal de toate discurile supraviețuitoare, rezultând performanța cea mai bună în modul degradat. În prezent, se utilizează soluții cu mai multe discuri de rezervă astfel încât să poată fi servite continuu cereri de acces foarte importante. Mai mult, în noile scheme de distribuție numite distribuție rezervată, spațiile de rezervă sunt, de fapt, distribuite în toată aria. Aceasta înseamnă că, după o refacere, noua matrice va fi logic diferită de cea originală.

Arhitecturile RAID 4 și RAID 5 oferă o toleranță la defecte prin utilizarea protecției prin paritate pentru blocuri de date din diferite discuri. Totuși, chiar dacă datele unui disc căzut pot fi încă refăcute cu ajutorul parității, acesta este un proces încet și poate duce la creșterea substanțială a încărcării pe discurile supraviețuitoare. Dacă o matrice este proiectată fără luarea în considerație a consecințelor căderii, atunci performanța sistemului se va degrada la un nivel inacceptabil.

Fie o matrice constând din N discuri în care fiecare grup de paritate implica K discuri. Presupunem că această matrice primește o încărcare de R citiri și W scrieri, distribuite uniform pe matrice. Este evident că în timpul funcționării normale a matricei încărcarea totală este $R+4*W$, deoarece o scriere implică citirea vechii date și a vechii parități, apoi scrierea noilor date și a noii parități. Când unul dintre discuri a căzut, cererile către acesta sunt deservite de discurile rămase valide în grupurile de paritate corespunzătoare. Se poate arăta că încărcarea totală în modul degradat este $R+4*W+[(K-2)*R+(K-8)*W]/N$, care este suportată de noua matrice, având un disc mai puțin, ca efect al căderii lui. Câteva articole recente au arătat că dacă N este mai mare decât K și dacă grupurilor de paritate pot fi distribuite uniform în matrice, atunci încărcarea adițională cauzată de o cădere poate fi împărțită egal de toate discurile rămase. Aceasta este în contradicție cu ceea ce se întâmplă la RAID 5 standard, cu grupuri de paritate multiple, unde discurile din același grup de paritate cu discul căzut trebuie să suporte toată încărcarea adițională, pe când alte discuri nu sunt afectate deloc. Mai mult, se recomandă ca valoarea lui K să fie mult mai mică față de N , după cum se vede și din ecuația de mai sus, pentru a rezulta o încărcare generală, ce rezultă în urma căderii unui disc, cât mai mică. Rezultatul net este acela că fiecare disc rămas „vede” doar o mică încărcare suplimentară și, deci, suferă mai puțin performanța.

Pasul următor este de a proiecta matrice astfel încât asocierea grupurilor de paritate este distribuită într-o manieră uniformă în matrice (**clustered RAID**). Se definește un *BIBD* (b, v, r, k, λ) (*Balanced Incomplete Block Design*) prin gruparea a v diferite varietăți de obiecte în b blocuri astfel încât:

- fiecare bloc conține k varietăți distincte de obiecte;
- fiecare varietate apare în r blocuri diferite;
- oricare două varietăți apar împreună în λ blocuri.

Numărând toate replicările de varietăți, rezultă că $v*r=b*k$.

Similar, considerând numărul de replicări ale varietăților dintr-un set de blocuri ce conține o varietate dată, avem $r*(k-1)=\lambda*(v-1)$.

O matrice generală de paritate cu N discuri este definită după cum urmează. Se partiționează logic fiecare disc în M domenii de mărime egală. Fiecare grup de paritate implică un domeniu de pe fiecare din cele K discuri, K fiind **dimensiunea grupului de paritate**. Pentru a distribui egal locațiile de paritate ca în RAID 5, fiecare domeniu este subdivizat

logic în K subdomenii. Un subdomeniu poate conține mai multe blocuri care nu necesită să fie neapărat contigue fizic ca domeniile, iar subdomeniile sunt considerate entități logice (un mod alternativ de a privi lucrurile este că întâi se formează un grup de paritate de bază și apoi se multiplică în matrice până când toate blocurile sunt atribuite. Toate blocurile dintr-un disc aparținând aceluiași grup de paritate pot fi considerate ca un subdomeniu).

| | | | |
|---|---|---|---|
| A | A | A | A |
|---|---|---|---|

Figura 3-16

| | | | |
|----|----|----|----|
| a1 | a1 | a1 | p1 |
| a2 | a2 | p2 | a2 |
| a3 | p3 | a3 | a3 |
| p4 | a4 | a4 | a4 |

Figura 3-17

În continuare, vom considera un grup A de paritate, format din patru domenii care va fi notat conform Figura 3-16, înțelegând că subdomeniile sunt organizate conform Figura 3-17 cu locațiile de paritate rotite. RAID 5 standard este un caz special de matrice generală de paritate, unde $M=1$ domenii pe disc, iar N este un întreg multiplu de K .

O distribuție uniformă a grupurilor de paritate (UPGD - Uniform Parity Group Distribution) este definită ca o matrice generală cu paritate, în care fiecare pereche de discuri are în comun exact L grupuri de paritate. Astfel, dacă oricare disc din aria UPGD cade, toate discurile sunt afectate în mod egal. O matrice de discuri poate fi construită dintr-un BIBD prin mapările:

- varietate->disc
- bloc->grup de paritate,

rezultatul este o matrice UPGD cu $N=v$, numărul de grupuri de paritate $=b$, mărimea grupului de paritate $K=k$, numărul de domenii pe disc $M=r$, și numărul de grupuri de paritate în comun oricărei perechi de discuri $L=\lambda$. Ca exemplu, vom considera BIBD-ul de mai jos, cu $b=15$, $v=10$, $r=6$, $k=4$ și $\lambda=2$. Maparea sugerată va duce la o matrice UPGD

de 10 discuri cu 15 grupuri de paritate (de la A la Q), fiecare disc având șase domenii, fiecare grup de paritate având patru domenii și fiecare pereche de discuri având două grupuri de paritate în comun:

| A | B | C | D | E | F | G | H | J | K | L | M | N | P | Q |
|---|---|----|---|---|---|----|---|---|---|---|----|----|----|---|
| 1 | 1 | 1 | 1 | 3 | 2 | 2 | 1 | 2 | 3 | 4 | 1 | 2 | 3 | 4 |
| 2 | 3 | 4 | 2 | 4 | 4 | 3 | 5 | 6 | 5 | 5 | 8 | 5 | 6 | 7 |
| 5 | 6 | 7 | 3 | 5 | 6 | 7 | 6 | 7 | 7 | 6 | 9 | 9 | 8 | 8 |
| 8 | 9 | 10 | 4 | 8 | 9 | 10 | 7 | 7 | 8 | 9 | 10 | 10 | 10 | 9 |

Figura 3-18

| D1 | D2 | D3 | D4 | D5 | D6 | D7 | D8 | D9 | D10 |
|----|----|----|----|----|----|----|----|----|-----|
| A | A | B | C | A | B | C | A | B | C |
| B | D | D | D | E | F | G | E | F | G |
| C | F | E | E | H | H | H | J | K | L |
| D | G | G | F | K | J | J | M | M | M |
| H | J | K | L | L | L | K | P | N | N |
| M | N | P | Q | N | P | Q | Q | Q | P |

Figura 3-19

3.6.3.2 RAID 5/VSP

Matricea de discuri oferă avantaje deosebite în privința costului, performanței și a rentabilității. O matrice de N discuri mici poate fi mai puțin costisitoare decât un singur disc de capacitate mare, chiar dacă este adăugat un disc pentru informația de paritate. Prin adăugarea informației de paritate aria de discuri devine și mai fiabilă: datele pot fi recuperate la defectarea unui singur disc al oricărui din matricea de discuri.

Considerațiile asupra încărcării echilibrate și asupra ușurinței recuperării informațiilor au condus la dezvoltarea unor scheme în care informația de paritate este distribuită pe cele $N+1$ discuri. Organizarea clasei de

matrice de discuri în care paritatea este distribuită în acest mod a devenit cunoscută sub numele de RAID 5. În varianta cea mai simplă de organizare a RAID 5 datele sunt organizate în secțiuni cuprinzând N blocuri de date și un bloc de paritate (în care blocul de paritate este, de exemplu, un XOR la nivel de bit a celor N blocuri de date). Fiecare bloc din secțiune este plasat pe discuri diferite din matrice, cu poziția blocurilor de date în secțiune astfel încât fiecare disc să conțină o cantitate egală a blocurilor de paritate.

Deși această structură oferă câteva avantaje distincte, menținerea informației de paritate poate introduce o întârziere semnificativă la scriere în cazul unei fragmentări mai mici decât o secțiune. De exemplu, considerând un singur bloc scris într-o secțiune dată: pentru a înnoi informația din blocul de paritate sunt necesare informații asupra vechiului bloc de date, vechiului bloc de paritate și noului bloc de date. Dintre acestea toate, cu excepția noului bloc de date, trebuie citite de pe disc și noile blocuri de date și paritate trebuie scrise pe disc, toate acestea reprezentând până la patru operații I/O.

Întârzierea celor patru operații I/O pentru scrierea unui singur bloc nu trebuie privită ca un dezavantaj.

Dacă se compară timpul celor $N+1$ accese la disc per *byte* cu timpul de acces în cazul unui disc de capacitate mai mare. Există cazuri când performanțele pot fi îmbunătățite reducând acest dezavantaj. Un mod de a realiza acest lucru este acela de a înseria scrierile: dacă toate cele N blocuri de date dintr-o secțiune sunt scrise în același timp, blocul de paritate poate fi calculat în paralel necesitând $N+1$ scrieri pe disc (și nici o citire de pe disc).

O metodă a organizării datelor astfel încât să se beneficieze de proprietatea descrisă anterior este un sistem de fișiere *log-structured* (SFLS).

Oricum, această metodă prezintă un dezavantaj numit „curățirea segmentului“, care necesită periodic citiri și scrieri de zone relativ mari de pe discuri pentru a se întrepătrunde spațiul liber pentru următoarea utilizare.

Amintim aici două abordări ale acestei probleme: *floating parity* și *parity logging*:

- *floating parity* nu reduce numărul total de operații I/O; uneori unele perechi de operații pot fi realizate într-o singură rotație a discului;

- *parity logging* mărește, de fapt, numărul total de operații de I/O, reducând, însă, raportul cost/operație ca rezultat al inserării informațiilor de paritate, dar aceste operații necesită un *buffer* de memorie nonvolatil și un spațiu adițional pe disc.

În continuare vom considera o abordare diferită în scopul reducerii costului operației de scriere bazată pe noțiunea VSP a blocurilor de paritate.

Considerând dată o secțiune cuprinzând N blocuri de date D_1, D_2, \dots, D_N și un bloc de paritate P , P conținând paritatea pentru un subset arbitrar de blocuri de date $D_{i_1}, D_{i_2}, \dots, D_{i_q}$, unde subsetul protejat este identificat în P (plasat, de exemplu, ca un vector de biți). Informația de paritate nu este menținută pentru blocurile neprotejate de date din secțiune și astfel de blocuri pot fi scrise și adăugate setului protejat fără a avea cunoștințe asupra conținutului blocului. În plus, dacă blocurile de paritate excedentare sunt plasate într-un *cache* pentru secțiuni conținând blocuri neprotejate, scrierea sau modificarea unui bloc, dintr-una dintre aceste secțiuni necesită doar două operații I/O. În final, recuperarea datelor în urma defectării unui disc este, de asemenea, facilitată deoarece recuperarea blocurilor neprotejate nu este necesară.

Presupunem că doar blocurile de date libere (nealocate) pot fi neprotejate. Acestea pot atinge o stare în care devin protejate și este de dorit a se trece aceste blocuri în stare neprotejată. Acest proces îl vom numi *cleaning* (curățire). Acest proces se poate încadra ușor în majoritatea sistemelor de stocare existente implementate prin sisteme de fișiere, sisteme de gestiune a bazelor de date etc.

Pentru a reduce fragmentarea informației stocate, recuperarea blocurilor pierdute, îmbunătățirea performanțelor discurilor s.a.m.d., aceste sisteme efectuează în mod obișnuit o reorganizare periodică a discului.

Scopul nostru este de a investiga potențialul îmbunătățirii performanței în sisteme utilizând *variable scope protection* (VSP). VSP prezintă grade de libertate în plus, cu mențiunea că se păstrează blocurile de paritate pentru a crește fiabilitatea. Considerăm două măsuri de evaluare a performanței: media costului total pentru scrierea unui singur bloc și costul pentru refacerea informației. Aceste informații preliminare se folosesc pentru a obține mai rapid rezultatele dorite (care se obțin prin simulare sau pe cale analitică). Deși studiul teoretic arată la prima vedere rezultate nesatisfăcătoare, sistemele care utilizează VSP dau rezultate bune pe cale practică.

3.6.3.2.1 Organizarea unui RAID/VSP

Presupunem un sistem format dintr-o matrice de $N+1$ discuri. Acesta conține un controler cu un *cache*, format dintr-un număr de blocuri (de date și/sau de paritate) recent accesate. Controlerul gestionează *cache*-ul și aria de discuri și procesează cererile provenite de la unul sau de la mai multe sisteme. În general, o colecție de blocuri de date și blocuri de paritate asociate este denumită „grup de paritate”. O varietate largă de metode propun organizarea grupurilor de paritate în matrice de discuri.

Aici se folosește metoda în care fiecare grup de paritate constă dintr-o secțiune, de exemplu un bloc ales din fiecare disc: N blocuri de date și un bloc de paritate. Aceasta nu presupune o reducere a generalității problemei (în prezentarea de față operațiile controlerului se aplică la fel de bine pentru orice grup de paritate), dar simplifică prezentarea noastră. Controlerul de disc gestionează secțiunile (sau, mai general, grupurile de paritate): fiind dată o cerere de citire a unui bloc oarecare specificat printr-un identificator de bloc de disc, controlerul determină secțiunea de care aparține blocul și discul pe care se află blocul considerat în acea secțiune (desigur, se verifică în primul rând existența blocului în *cache*). După cum s-a menționat în introducere, este de dorit ca, în secțiunile discurilor, blocurile de paritate să fie distribuite în mod egal în cele $N+1$ discuri.

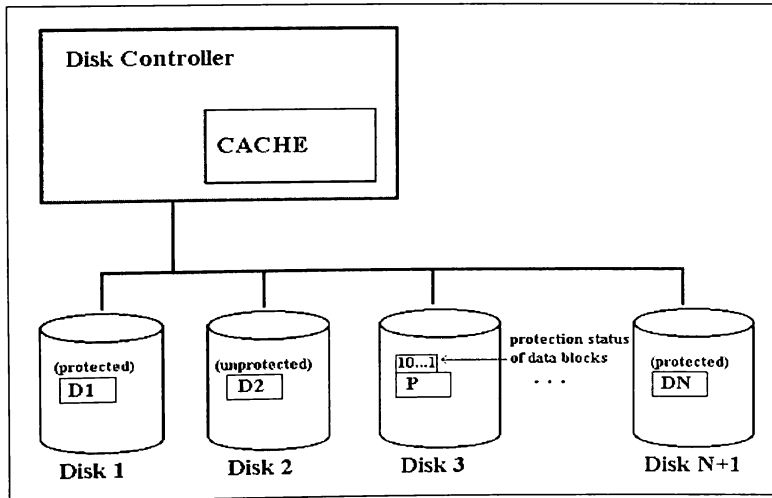


Figura 3-20

Fiecare bloc de date din subsistemul matricei de discuri se poate afla în una din următoarele două stări :

- alocat în acest moment de un sistem de un nivel superior (exemplu: un sistem de fișiere sau un sistem de baze de date) de stocare a datelor;
- nealocat, deoarece fie nu a fost alocat încă, fie a fost anterior alocat, dar la următoarea alocare datele au fost șterse sau mutate în altă parte.

Ne referim la blocurile de date alocate *in-use* (în folosință) și blocurile nealocate *free* (libere). Spre deosebire de arhitecturile RAID anterioare, nu sunt necesare blocuri de paritate pentru toate blocurile de date din secțiune. De asemenea, se presupune că este necesar ca blocul de paritate să conțină informații de paritate numai pentru blocurile aflate *in-use*. Poate, de asemenea, să conțină informații de paritate pentru un subset oarecare (incluzând și un set vid) de blocuri goale într-o secțiune. Dacă blocul de paritate conține informații de paritate pentru un bloc de date dat, spunem că blocul de date este protejat; altfel este neprotejat.

În cele ce urmează, fiind dată o secțiune, fiecare bloc de date se poate afla în una din următoarele trei stări:

- *in-use* (obligatoriu protejat);
- *free* și protejat;
- *free* și neprotejat.

Dacă toate blocurile de date din secțiune sunt *free*, vom numi aceasta o secțiune *free*. Un subset de blocuri de date aflate într-o secțiune *free* poate fi protejat. Dacă toate blocurile de date dintr-o secțiune sunt *free* și neprotejate, vom numi aceasta *clean free*. De remarcat că un bloc de paritate a unei secțiuni *clean free* nu conține nici o informație utilă (el conține informații de paritate pentru un set de blocuri vide). Putem considera și că blocurile de paritate se pot afla în două stări: *in-use* (unul sau mai multe blocuri de date sunt protejate) sau *free* (toate blocurile de date sunt neprotejate).

Există mai multe moduri de a identifica blocurile de date într-o secțiune. De exemplu: starea fiecărui bloc de date (protejat sau neprotejat) poate fi reținută într-un vector de biți aflat în blocul de paritate. Acesta prezintă avantajul pe care îl au blocurile de paritate aflate în *cache* conținând stări protejate/neprotejate asociate blocurilor de date. Ca un

exemplu pentru un model alternativ, controlerul de disc poate păstra informații pentru toate secțiunile într-un RAM nonvolatil.

Este necesară, de asemenea, păstrarea informațiilor despre stările *in-use* sau *free* ale fiecărui bloc de date. Acesta se asigură la nivelul unui sistem *software* de un strat din sistemul de fișiere sau din sistemul de baze de date, de exemplu. Pentru a fi mai clari, ne vom referi la acesta drept *storage manager*. *Storage manager* implementează următoarele funcții:

- harta identificatorilor datelor din blocurile de disc (unde aceste blocuri de disc sunt, în contextul de față, la rândul lor specificate de identificatorii blocurilor de disc);
- procese care cer citirea sau scrierea datelor;
- procese care cer alocare și dealocare de date.

De asemenea, controlerul are disponibilă harta identificatorilor blocurilor de disc dintr-o secțiune și poate, de exemplu, să determine când o intrare dintr-o secțiune este *free*.

În condițiile actuale este de presupus că *storage manager*-ul este implementat ca un strat în sistemul analizat. De asemenea, se poate implementa *storage manager*-ul ca o parte a controlerului de disc. În acest caz, câteva optimizări devin fezabile. Exemple posibile ar fi:

- integrarea hărților de identificatori logici de date ca și ID-uri a blocurilor de disc și ID-urile blocurilor de disc ca adrese fizice ale blocurilor.
- integrarea controlului informațiilor *in-use/free* și protejate/neprotejate.

Discutarea problemelor atinse anterior nu constituie obiectul lucrării de față și pentru următoarele descrieri presupunem că *storage manager*-ul este implementat independent de controlerul de disc.

Operațiile obișnuite ale controlerului de disc (în absența defecțiunilor) diferă de un controler RAID fără VSP doar în privința scrierilor și a convertirii blocurilor protejate în blocuri neprotejate.

În primul rând vom discuta despre scrieri. Pentru simplitate, vom descrie doar scrierea unui singur bloc (scrierea multibloc este considerată o extensie a scrierii unui singur bloc). Presupunem că paritatea este calculată printr-un XOR pe biți a blocurilor de date. O scriere într-o secțiune *free* poate fi tratată ca un caz special:

privitor la actuala stare de protecție a blocurilor de date *free*, după scriere doar un singur bloc *in-use* ar trebui să fie protejat.

Pentru aceasta, scrierea unui bloc D într-o secțiune liberă S este procesată după cum urmează (unde P este un bloc de paritate).

- se setează $P=D$;
- se marchează D ca protejat în S ;
- se scrie P, D pe disc (de remarcat că scrierile pe disc pot avea loc în mod concurrent);
- dacă oricare din blocurile de date din S , altul decât D , este prezent în *cache*, se eliberează spațiul din *cache* pentru fiecare astfel de bloc.

În continuare, scrierea unui bloc D într-o secțiune S (*non free*) cu blocul de paritate P este tratată după cum urmează (unde D_{old} se referă la anterioara versiune a lui D).

Dacă D_{old} nu este protejat, atunci:

- citește P de pe disc, dacă nu este în *cache*;
- setează $P=P \text{ XOR } D$;
- marchează D ca protejat în S ;
- scrie P, D pe disc (de remarcat că scrierile pe disc pot fi concurente);
- ieșire.

Dacă D_{old} este protejat, se execută următoarele:

- Dacă P și/sau D_{old} nu sunt în *cache*, atunci citește de pe disc (de remarcat că dacă nici unul nu se află în *cache*, citirile de pe disc pot fi concurente);
- Setează $P=P \text{ XOR } D_{old} \text{ XOR } D$;
- scrie P și D pe disc (scrierile pot fi concurente);
- Dacă D_{old} se află în *cache*, eliberează spațiul ocupat de D_{old} .

În cele din urmă se descriu operațiile pentru „curățirea“ secțiunilor, prin care se înțelege convertirea blocurilor de date protejate *free* către o stare neprotejată. Avantajul utilizării unei secțiuni curate (în care toate blocurile *free* sunt neprotejate) este la scrierea inserată în secțiune, care este mai puțin costisitoare în termeni de accese la disc. În cazul unei

scrieri a unei secțiuni *free*, curățirea unei secțiuni în care toate blocurile de date sunt *free* poate fi tratată ca un caz special, după cum urmează:

- marchează toate blocurile de date din S ca neprotejate
- pentru fiecare dintre aceste blocuri de date, dacă se află în *cache* eliberează spațiul ocupat din *cache*;
- scrie P pe disc (presupunând că informația de protecție este depusă în P);

Dacă secțiunea de curățat are blocuri *in-use*, atunci operația este puțin mai complicată. Următoarea implementare minimizează numărul de accese la disc în acest caz (unde S este secțiunea, P este bloc de paritate, blocurile de date *in-use* sunt D_i și blocurile de date protejate *free* sunt B_j).

Fie d numărul celor D_i care nu se află în *cache* și fie b numărul B_j care nu se află în *cache*:

- dacă $d \leq b$ atunci citește fiecare D_i citește fiecare D_i care nu se află în *cache* (citirile pot fi concurente) și calculează P prin XOR pe toate D_i ;
 - altfel, citește fiecare B_j care nu se află în *cache* (citirile pot fi concurente) și recalculează P prin XOR între P și toți B_j ;
- marchează fiecare B_j ca neprotejat;
- pentru fiecare B_j din *cache* eliberează spațiul din *cache*-ul blocului;
- scrie P pe disc;

Aceasta încheie descrierea operațiilor controlerului de disc, fără nici un detaliu privitor la algoritmi sau operațiile la nivelul *storage manager*.

Din cele anterior arătate, se poate observa că pentru anumite operații sunt îmbunătățiri semnificative ale unui sistem fără VSP. De exemplu, o scriere a unei secțiuni *free* îi este necesar doar două accese la disc fără a sonda conținutul *cache*-ului și o scriere a unui bloc neprotejat *free* îi sunt necesare trei accese la disc. În mod ideal, *storage manager*-ul ar trebui proiectat în așa mod încât să beneficieze de avantajele acestor proprietăți.

Există un număr de idei în modelul unui astfel de *storage manager*, unele dintre ele fiind prezentate în continuare.

Pentru a reduce costurile de scriere, *storage manager*-ul poate să realoce scrierile blocurilor neprotejate *free*. De asemenea, doar dacă secțiunea

este *free*, costul scrierii poate fi redus și mai mult dacă blocul de paritate se află în *cache*. În consecință, pentru dispozitivul de stocare este indicat să se mențină în *cache* un număr de blocuri de paritate pentru secțiunile cu blocuri neprotejate *free*.

După scrierea blocurilor neprotejate *free*, acestea devin blocuri protejate, iar dacă scrierile se datorează realocării, fără alte operații, locațiile anterioare ale blocurilor vor rămâne protejate. Astfel, curățirea periodică este necesară pentru a converti blocuri protejate în blocuri neprotejate.

Similar, dacă numărul blocurilor de paritate aflate în *cache* pentru secțiuni cu blocuri neprotejate este epuizat, reverificarea acestora este necesară pentru îmbunătățirea performanței.

În următoarele două secțiuni, se vor prezenta rezultatele obținute în legătură cu câteva dintre problemele principale din modelul *storage manager* cu VSP. În primul rând, considerăm un model minimal al *storage manager*, proiectat spre a reduce costul scrierilor prin curățirea periodică a celor mai puțin utilizate secțiuni, împreună cu scrierea de blocuri neprotejate în aceste secțiuni. Apoi, vom compara această clasă de metode cu cea a sistemului de fișiere structurat (*LFS*) în care scrierile sunt plasate într-un *buffer* (*bufferate*) și relocate astfel încât secțiuni întregi (inclusiv blocul de paritate) sunt scrise în același moment (de remarcat că reorganizarea periodică a discului, operată de *LFS* pentru a întrepătrunde spațiul liber, este denumită tot „curățire“. Aici, pentru a evita ambiguitatea, vom folosi termenul „curățire“ când ne vom referi la conversia blocurilor protejate *free* în blocuri neprotejate, iar „reorganizarea *LFS*“ pentru procesul de întrepătrundere a spațiului liber de pe disc).

3.6.3.2.2 Analiza performanței în cazul scrierii aleatoare ale unui singur bloc

Vom descrie performanțele obținute în simularea unui model simplificat de sistem care folosește VSP. Fiecare operație se presupune că este o scriere într-un bloc logic ales la întâmplare (uniform distribuit). *Storage manager*-ul conține o hartă a identificatorilor blocurilor logice corespunzătoare identificatorilor blocurilor de pe disc (identificatorii blocurilor de pe disc sunt transformați în adrese fizice de către controlerul de disc). Folosind această hartă, *storage manager*-ul este liber să realoce blocuri, astfel încât fiecare scriere logică poate elibera un bloc *in-use* (precedenta locație a blocului logic) și alocă un bloc *free* (noua locație a blocului). Un model mai apropiat de realitate ar putea

include operații constând în citiri multiple și/sau scrieri în aceeași secțiune, în care costul modificării blocului de paritate poate fi amortizat prin mai mult decât o scriere.

Presupunem că:

- există i blocuri de date *in-use* în orice moment, fiecare specificat prin identificatorul blocului logic $1, 2, \dots, i$ care indică spre unul din numărul d de blocuri de date din aria de discuri (astfel că raportul de utilizare este de i/d);
- fiecare operație ulterioară este o scriere a unui singur bloc într-unul din cele i blocuri logice, orice bloc fiind la fel de prioritar.

Storage manager-ul execută această secvență a scrierii unui singur bloc după cum urmează:

- găsește cea mai puțin utilizată secțiune, adică secțiunea cu numărul cel mai mare de blocuri *free* (alege orice astfel de secțiune în cazul că sunt mai multe);
- golește secțiunea (în modul descris în secțiunea anterioară), convertind toate blocurile protejate *free* în blocuri neprotejate;
- realocă fiecare bloc scris în continuare într-un bloc *free*, până când toate blocurile sunt *in-use*;
- repetă de la pasul 1.

Deoarece blocul de paritate a fost găsit în pasul 2 și toate scrierile din pasul 3 se fac în blocuri neprotejate, fiecare operație de scriere nu necesită nici o operație de citire, ci doar două operații de scrieri pe disc (una pentru date și una pentru paritate) care pot avea loc concurrent.

Această tactică a fost simulată cu scopul de a găsi costul mediu al curățării blocurilor (măsurat în accese la disc) pentru scrierea unui singur bloc. S-au folosit arhitecturile matricelor de discuri 3+1 și 7+1 (unde $N+1$ înseamnă N blocuri de date și un bloc de paritate pentru o secțiune). Pentru a determina acest cost, s-a presupus folosirea metodelor de curățire prezentate anterior, cu presupunerea cazului cel mai negativ al nici unui *cache-hit* al blocurilor de date. Se deduce că pentru o secțiune *free* costul este nul (un bloc de paritate nou este scris la prima scriere în secțiune) și pentru o secțiune cu q blocuri *in-use* și $N-q$ blocuri protejate *free* costul este mai mic decât q sau $N-q+1$ citiri de disc (în primul caz sunt citite blocurile *in-use* pentru a calcula paritatea, iar în al doilea caz blocurile neprotejate *free* împreună cu actualul bloc de paritate sunt citite pentru a recalcula paritatea). Nu luăm în calculul

costului scrierea pe disc a blocului de paritate deoarece se presupune că *storage manager*-ul face ca scrierea blocului de paritate să fie întârziată până la scrierea primului bloc de date într-o secțiune (astfel că scrierea blocului de paritate este parte a costului scrierii unui singur bloc).

Simulările au fost făcute pentru sisteme conținând 1000 și 10000 de secțiuni. Interesant, rezultatele au fost aproape identice în fiecare caz (totuși, au fost obținute rezultate diferite pentru arhitecturile 3+1 și 7+1). Concluzionăm că, pentru o arhitectură cu un număr dat de discuri, costul mediu al curățirii la scriere nu depinde de numărul total de secțiuni (presupunând că acest număr este suficient de ridicat). Figura 3-21 rezumă aceste rezultate în care fiecare număr este medie a trei rânduri a câte 25000 de scrieri (a fost observată o foarte mică variație între ture).

De remarcat este faptul că la o încărcare de 95% și peste, costul mediu este același pentru ambele arhitecturi: 3+1, 7+1. Aceasta se datorează următoarelor: pentru încărcări apropiate de 100%, majoritatea secțiunilor selectate pentru curățire vor avea un bloc *free*, rezultând în costul de curățire două citiri de disc (una pentru blocuri protejate *free* și una pentru blocul de paritate) pentru scrierea ulterioară a unui bloc de date. Astfel, costul mediu devine independent de numărul de discuri la o încărcare aproape de limite și, pe cât se apropie de o încărcare de 100%, acest cost are o valoare limitată la doi.

| | 3 + 1 RAID | | 7 + 1 RAID | |
|-------|------------|---------|------------|---------|
| Stor. | 1,000 | 10K | 1,000 | 10K |
| Util | Setiune | Setiune | Setiune | Setiune |
| 50% | 0 | 0 | 0.13 | 0.13 |
| 60% | 0.14 | 0.14 | 0.30 | 0.30 |
| 70% | 0.34 | 0.34 | 0.58 | 0.58 |
| 80% | 0.52 | 0.51 | 1.08 | 1.08 |
| 90% | 1.33 | 1.33 | 1.45 | 1.45 |
| 95% | 1.68 | 1.68 | 1.68 | 1.68 |
| 98% | 1.88 | 1.88 | 1.88 | 1.88 |

Figura 3-21

Costul mediu la scriere, măsurat în număr de accese la disc, nu depășește valoarea doi (pentru scrierea blocului de date și a blocului de paritate) plus costul mediu de curățire per scriere (de remarcat că pentru RAID 3+1 la o încărcare de 50% costul a fost nul: a fost cazul în care s-a găsit o secțiune cu trei blocuri *free*). Comparând acestea cu un sistem fără *VSP*, în care fiecare scriere (de asemenea a unui singur bloc) necesită patru accese la disc, există întotdeauna o îmbunătățire a performanțelor. Totuși, rezultate de un acces la disc per scriere au fost obținute în cazul 3+1 la o încărcare de până la 80%, iar în cazul 7+1 la o încărcare de până la 70%.

3.6.3.2.3 Comparație VSP cu LFS

După cum s-a menționat anterior, un *LFS* este o abordare cunoscută pentru reducerea întârzierii întâlnită la operația de scriere pe discuri. Deși metodele *LFS* nu au fost inițial proiectate pentru matrice de discuri, ele se aplică la fel de bine în acest context și permit ca datele să fie scrise în unități de secțiuni: folosind o matrice de discuri cu secțiuni conținând N blocuri de date și unul de paritate, scrierile pe disc sunt păstrate într-un *buffer* până când se adună minim N blocuri și, apoi, se scriu într-unul sau în mai multe secțiuni *free*. Pentru simplitate vom neglija problemele care pot apărea cu privire la pierderea de date într-un *buffer* parțial plin (în practică, aceste *buffer*-e trebuie scrise pe disc după un anumit timp) și vom presupune că datele sunt întotdeauna scrise în unități de secțiuni.

Folosind această abordare, scrierea a N blocuri necesită $N+1$ accese la disc (unul pentru fiecare bloc de date și unul pentru blocul de paritate). Totuși, acesta poate fi atins prin reorganizarea periodică a segmentelor matricei de discuri pentru a genera matrice conținând secțiuni *free*, astfel că pentru a compara costul scrierilor, costul reorganizării trebuie luat în calcul.

Segmentele care vor fi reorganizate sunt în mod tipic alese drept cele cu cea mai mică încărcare sau după o funcție care combină nivelul de încărcare cu vârsta presupusă a datelor în fiecare segment. Indiferent de metoda folosită la selecția segmentului, va exista o încărcare medie pentru astfel de segmente, să o notăm u , și setul de segmente care trebuie reorganizat va conține un număr total de secțiuni, să spunem M . Folosind metoda de copiere și compactare, reorganizarea decurge după cum urmează: blocurile *in-use* din toate cele M secțiuni sunt citite pentru un total de uMN citiri de disc și, apoi, scrise cu blocurile *in-use* împachetate în secțiuni. Folosind o aproximare continuă, aceasta necesită un total de $uM(N+1)$ scrieri (de reamintit că blocul de paritate

este scris pentru fiecare secțiune în faza de scriere). Aceasta conduce la $(1-u)M$ secțiuni *free* cu un cost de $uM(2N+1)$ accese la disc. După reorganizare, $(1-u)MN$ blocuri pot fi scrise în secțiuni *free*, ceea ce necesită $(1-u)M(N+1)$ accese la disc ținând cont și de scrierea blocului de paritate. Fie C_{LFS} costul scrierii folosind această metodă, definit ca media numărului total de operații I/O pentru un bloc scris. Din cele prezentate anterior deducem:

$$C_{LFS} = \frac{uM(2N+1) + (1-u)M(N+1)}{(1-u)MN} = \frac{N + uN + 1}{(1-u)N}$$

Acum, să presupunem aceleași segmente cu M secțiuni și același factor de încărcare u selectate pentru curățire și, ulterior, realocate scrierii într-un sistem VSP. Spre deosebire de abordarea LFS, aici nu există un proces de reorganizare a datelor. În schimb, blocurile protejate *free* din toate astfel de segmente sunt convertite în blocuri neprotejate, iar ulterioarele scrieri se referă la aceste locații. Fiind dată o secțiune particulară, să presupunem că încărcarea ei este de u' . Deci, există $u'N$ blocuri *in-use* și $(1-u')N$ blocuri protejate *free* în aceasta secțiune. Presupunem că după curățire noul bloc de paritate nu va fi scris pe disc până la prima scriere în secțiune. Prin urmare, costul curățirii este numărul de citiri de pe disc necesare pentru a recalcula paritatea, care este $\min(u'N, 1+(1-u')N)$ pentru $u' < 1$ (pentru $u'=1$ costul este zero, deoarece nu există blocuri *free* pentru curățat). Din păcate, această funcție nu este liniară în u' , ceea ce înseamnă că pentru costul curățirii sunt necesare informații suplimentare asupra distribuției lui u' . Totuși, putem obține marginea superioară a costului prin:

- neținând cont de cazul special $u'=1$ și, în schimb, folosind $\min(u'N, 1+(1-u')N)$ drept costul curățirii pentru toate secțiunile;
- folosind banalele inegalități $a \geq \min(a,b)$ și $b \geq \min(a,b)$.
- presupunem că, în cazul în care costul curățirii pentru fiecare secțiune având încărcarea u' ar fi $u'N$, costul total ar fi uMN . Similar, când costul ar fi întotdeauna $1+(1-u')N$, costul total ar fi $M+(1-u)MN$. De vreme ce acestea două sunt marginile superioare ale costului total, acest cost este mai mic sau egal cu $M\min(uN, 1+(1-u)N)$ citiri de pe disc.

Apoi, după curățire, să considerăm scrierea în blocurile neprotejate dintr-o secțiune cu încărcarea u' a unei secțiuni cu $(1-u')N$ blocuri neprotejate. Bufferizând scrierile precum în cazul LFS și incluzând scrierea blocului de paritate, aceste operații necesită $1+(1-u')N$ scrieri pe disc pentru $u' < 1$ (dacă $u'=1$, atunci nici o scriere nu poate avea loc în

această secțiune). Scrierea tuturor celor $(1-u)MN$ blocuri neprotejate rezultate din curățire necesită cel mult $M+(1-u)MN$ scrieri pe disc. Rezultă, deci, costul mediu pentru scriere C_{vsp} trebuie să satisfacă relațiile următoare:

$$C_{vsp} \leq \frac{M \min(uN, 1 + (1-u)N) + M + (1-u)MN}{(1-u)MN} = \frac{\min(N+1, 2+2(1-u)N)}{(1-u)N}$$

Pentru a compara aceste două metode, să considerăm raportul C_{vsp}/C_{LFS} . Avem următorul rezultat:

$$\frac{C_{vsp}}{C_{LFS}} \leq \frac{\min(N+1, 2+2(1-u)N)}{N+uN+1} \leq \frac{N+1}{N+uN+1} \leq 1$$

Se observă că acest cost la scriere, măsurat în număr de operații I/O per bloc scris, este întotdeauna mai mic folosind metoda VSP (neținând cont de cazul limita $u=0$ în care toate secțiunile sunt *free*, în care costurile de scriere sunt evident egale). Marginea superioară a raportului C_{vsp}/C_{LFS} dat anterior este arătată în Figura 3-22 în funcție de u , pentru matricele de discuri 3+1 și 7+1. De remarcat că îmbunătățirile cresc o dată cu creșterea lui u . Pentru mărirea factorului de încărcare se poate deduce din relațiile anterioare următoarea relație:

$$\lim_{u \rightarrow 1} \frac{C_{vsp}}{C_{LFS}} \leq \frac{2}{2N+1}$$

Acest rezultat este interesant deoarece metoda LFS este cunoscută a da rezultate slabe la un factor de încărcare ridicat. În practică este de așteptat ca performanțele relative ale metodelor VSP, în comparație cu metodele LFS, să fie mai bune decât cele sugerate anterior, ținând cont de faptul că expresiile derivate din raportul costului mediu la scriere între metodele VSP și LFS sunt referitoare la marginile superioare, în practică, însă, raportul real fiind mai mic.

În general, presupunerile simpliste dau un avantaj îndoielnic metodelor LFS (exemplu: se presupune că metodele LFS întotdeauna scriu datele la un anumit nivel de ocupare, deși, după cum s-a remarcat anterior, în practică *buffer*-ele parțial încărcate pot fi forțate să scrie datele pe disc chiar la un nivel de ocupare mai mic decât cel considerat).

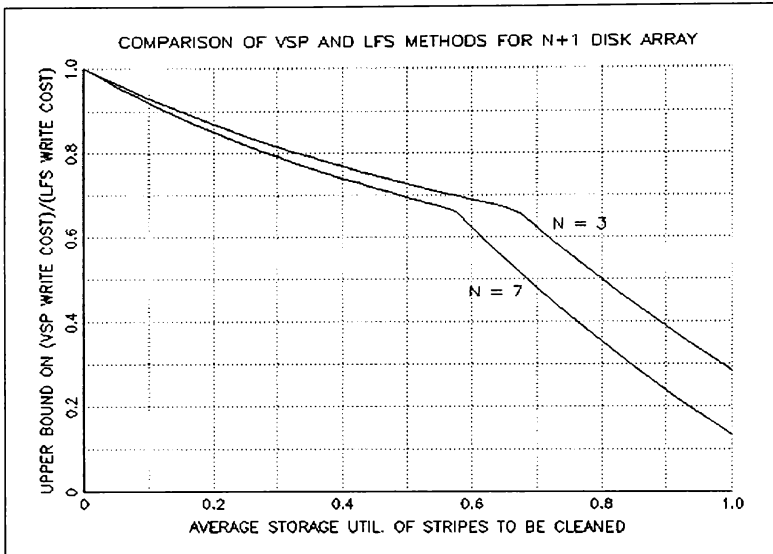


Figura 3-22

3.6.3.2.4 Costul recuperării datelor de pe discul căzut

În acest capitol, vom considera costul recuperării blocurilor de date și de paritate în cazul defectării unui singur disc din aria de $N+1$ discuri. O analiză a costului recuperării discului defect este îngreunată de faptul că, în practică, pentru a grăbi recuperarea, mari cantități de date (chiar cilindri întregi) trebuie citite de pe fiecare disc nedefectat în același timp. Totuși, vom estima aici munca totală necesară recuperării doar prin numărarea blocurilor care trebuie citite de pe fiecare disc nedefectat pentru a recupera un bloc ales la întâmplare de pe discul defectat. Fără VSP și cu grupuri de paritate implementate în secțiuni, acest număr este întotdeauna N , deoarece un bloc trebuie citit de pe fiecare celelalte discuri.

Să presupunem acum că folosim o schemă cu VSP. Atunci, o parte p din toate blocurile de date este protejată, iar restul blocurilor sunt neprotejate. Pentru a simplifica următoarele analize, să presupunem că fiecare bloc de date este protejat cu probabilitatea p , independent de starea de protecție a celorlalte blocuri de date.

Există cel puțin două moduri în care poate avea loc recuperarea, depinzând de detaliile modelului sistemului, precum și de modul de defectare al discului. Într-unul din cazuri, nici o informație de protecție

alta decât cea stocată în aria de discuri nu este disponibilă și, de asemenea, informația *in-use* păstrată de *storage manager* nu este disponibilă. Totuși, s-ar putea ca informațiile asupra stării de protecție a blocurilor de disc să fie disponibile altundeva, de exemplu într-o hartă păstrată în *storage manager*. În acest caz, informațiile *in-use* ar fi probabil disponibile, dar pentru simplitate vom examina cazul recuperării folosind doar informația de protecție. Ambele cazuri vor fi tratate mai jos, în primul rând analizând cazul în care nu este disponibilă informație de protecție (sau *in-use*).

În primul rând, fiind dat un bloc ales arbitrar de pe discul defectat, deoarece blocurile de paritate sunt egal distribuite pe toate discurile, acest bloc are probabilitatea de a fi un bloc de paritate de $1/(N+1)$ și de date de $N/(N+1)$.

Dacă blocul de pe discul defect este un bloc de paritate, atunci, dacă informația asupra stării blocurilor de date din secțiune nu se găsește disponibilă în alt loc, este de presupus că fiecare bloc de date poate fi *in-use*, necesitând citirea tuturor celor N blocuri de date de pe celelalte discuri pentru a calcula noul bloc de paritate (alternativ, dacă informația de protecție este disponibilă, numărul de blocuri de date necesar pentru a calcula noul bloc de paritate ar fi pN).

În caz contrar, blocul de pe discul defectat este un bloc de date. Presupunând că informația de protecție este disponibilă doar din blocul de paritate al secțiunii, cel puțin un bloc (de exemplu, blocul de paritate) trebuie citit pentru a recupera informația de protecție (dacă informația de protecție este disponibilă în alt loc, această citire nu este necesară). Recuperând informația, blocul va fi neprotejat cu o probabilitate de $1-p$, caz în care nici o acțiune ulterioară nu e necesară. Blocul este protejat, cu probabilitatea p , caz în care numărul cerut de blocuri de date pentru a efectua recuperarea informațiilor este $p(N-1)$.

Combinând rezultatele anterioare, numărul de blocuri necesare pentru a recupera blocul considerat de pe discul defectat este:

$$N_b = \frac{N}{N+1} + \frac{N(1-p)}{N+1} + \frac{Np}{N+1}(1+p(N-1)) = N \left(\frac{2+p^2(N-1)}{N+1} \right)$$

Alternativ, să presupunem că informația de protecție este disponibilă pe o sursă, alta decât blocurile de paritate din aria de discuri (dar nu informația *in-use*). Repetând analizele anterioare, considerate pentru acest caz, numărul de blocuri necesare pentru a recupera blocul considerat de pe discul defectat este:

$$N_b = \frac{pN}{N+1} + \frac{N(1-p)0}{N+1} + \frac{Np}{N+1}(1+p(N-1)) = N \left(\frac{2p + p^2(N-1)}{N+1} \right)$$

Să presupunem, de exemplu, că $N=7$ și că jumătate din blocurile de date din aria de discuri sunt protejate. În cazul în care informația de protecție este disponibilă doar în blocurile de paritate din aria de discuri, vom găsi că $(2+0.5^2 \times 6)/8=43.75\%$ blocuri trebuie citite pentru a recupera fiecare bloc defectat în comparație cu recuperarea în cazul RAID 5. În mod similar, dacă informația de protecție este disponibilă în alt loc, doar $(2 \times 0.5 + 0.5^2 \times 6)/8 = 31.25\%$ blocuri trebuie citite. Reducerea costului recuperării, comparat cu RAID 5 folosind anterioarele două metode este arătată în Figura 3, ca o funcție de raportul blocurilor de date protejate, pentru ariile de discuri 3+1 și 7+1.

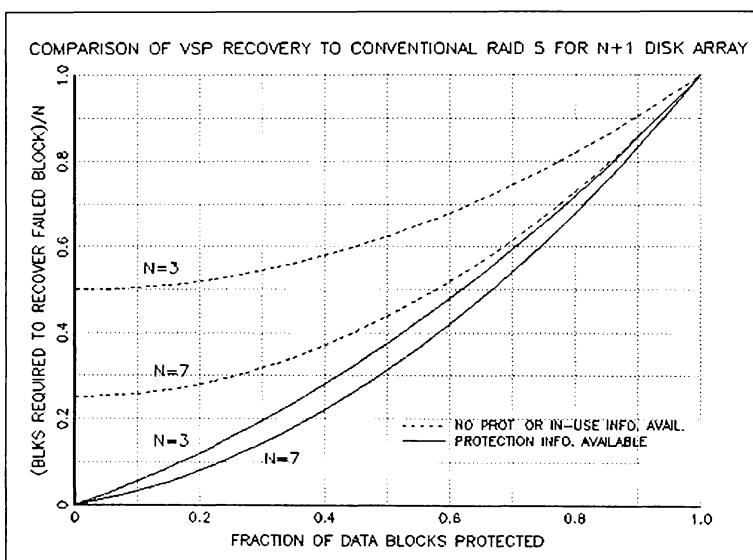


Figura 3-23

Dacă atât informația de protecție, cât și cea *in-use* sunt disponibile (de exemplu, din *storage manager*) și sunt folosite în timpul recuperării, costul recuperării devine și mai mic (doar dacă nu cumva aria de discuri se află în starea în care toate secțiunile sunt curate). O astfel de analiză implică adăugarea unui număr de presupuneri cu privire la numărul și distribuția blocurilor protejate *free*. Deoarece aceste presupuneri sunt dependente de felul în care lucrează *storage manager*-ul și de algoritmi utilizați de acesta, aici vom considera o situație mai simplificată (de remarcat că rezultatele anterioare ale costului recuperării când

informația de protecție este disponibilă în alt loc este marginea superioară a costului recuperării în cazul unor metode care folosesc atât informația de protecție, cât și cea *in-use*). Prin urmare, în practică, este de așteptat ca acest cost să fie mai mic în acest caz decât cel indicat de rezultatele anterioare. De asemenea, este de remarcat faptul că asemenea metode nu doar recuperează blocurile de pe un disc defectat, dar (cu nici un cost suplimentar, în timpul procesului de recuperare) au și un efect lateral, de curățire a secțiunilor.

3.6.3.2.5 Concluzii

Am prezentat o extensie a arhitecturii RAID 5 în care proporția informației de protecție pentru fiecare secțiune poate varia. Folosind VSP, informația de paritate este menținută pentru subseturi de blocuri de date în fiecare secțiune. Dacă paritatea este păstrată pentru un bloc de date dat, se spune că este protejat; altfel este neprotejat. Presupunem că, în practică, toate blocurile de date *in-use* trebuie să fie protejate. Blocurile de date *free* (nealocate) sunt inițial neprotejate. Cum scrierile au loc în blocuri de date neprotejate, blocurile devin protejate. Dacă data stocată într-un bloc *in-use* (protejat) este ștersă sau mutată, starea blocului devine protejat *free*.

Scrierile pe disc la o valoare mai mică decât granularitatea secțiunii sunt probleme de performanță cunoscute pentru arhitecturile RAID 5: scrierea unui singur bloc poate necesita până la patru operații I/O. Prin scrierea în blocuri neprotejate, această problemă se reduce, deoarece scrierea unui singur bloc necesită cel mult trei operații I/O în acest caz. Dacă scrierea se face într-un bloc deja alocat, este necesară realocarea scrisului. Când se realocă o scriere, noua locație devine protejată *in-use* și anterioara locație este protejată *free*. Totuși, pentru a menține o cantitate de blocuri neprotejate, este necesară o curățare periodică prin care blocuri protejate *free* sunt convertite în blocuri neprotejate.

Simulările au fost făcute de către un *storage manager* „minimal“ care efectuează următoarele:

- găsește și curăță cea mai puțin utilizată secțiune;
- realocă toate scrierile ulterioare a câte unui singur bloc aceasta secțiune până când aceasta este plină.

S-a descoperit că există întotdeauna o îmbunătățire în costul scrierii (măsurat în media numărului total de operații I/O per scriere) comparat cu modelul convențional de RAID 5. Totuși, îmbunătățiri substanțiale au apărut doar pentru încărcări mici, până la medii. De exemplu, o

economie de mai mult de un acces la disc per scriere a fost obținută pentru matricele de discuri 3+1 pentru încărcări de până la 80%, iar pentru matricele de discuri 7+1 pentru încărcări de până la 70%. În practică, este de așteptat ca performanțele să fie simțitor mai bune, deoarece scrierile nu sunt aleatoare și uniform distribuite (exemplu: seturi de blocuri dintr-un fișier dat pot fi citite și scrise împreună).

Un mod alternativ de a reduce costul scrierii este de a forța scrierea la o anumită granularitate a secțiunii. Un mod de a face acest lucru este utilizarea LFS. Însă, aceasta necesită reorganizarea periodică a discului (cunoscută în literatura LFS drept curățire) pentru a genera noi zone libere. Folosind câteva presupuneri simple, rezultatele analitice au fost obținute comparând costul mediu de scriere pentru metodele VSP și LFS. S-a arătat că metodele VSP au un cost mediu de scriere mai scăzut în toate cazurile, cu îmbunătățiri mărite la creșterea încărcării. Acest rezultat este interesant deoarece este știut că metodele LFS prezintă deficiențe ale performanței la o încărcare ridicată.

Costurile recuperării informației au fost analizate pentru două cazuri. Într-unul dintre acestea, nici o informație alta decât cea aflată pe discurile funcționale din aria de discuri nu este disponibilă. În al doilea caz, toată informația de protecție (dar nu informația *in-use*) este disponibilă în alt loc (de exemplu, din sistemul *storage manager*-ului). În ambele cazuri, costul recuperării informației a fost mai mic decât cel al sistemelor convenționale RAID 5, cu o reducere și mai mare o dată cu scăderea numărului de blocuri protejate. De exemplu, cu 50% blocuri de date protejate, costul recuperării a fost de 44% în raport cu sistemele convenționale RAID 5 în primul caz, și cu 31% în cel de-al doilea. Dacă și informația *in-use* se folosește pentru recuperarea informațiilor, costul va fi și mai mic și, pe deasupra, ca un efect lateral benefic al recuperării, un anumit număr de secțiuni vor fi curățate prin procesul de recuperare.

Rezumând, permițând blocului de paritate să protejeze un subset de blocuri de date dintr-o secțiune este o idee simplă care poate duce la micșorarea problemelor întâlnite la sistemele RAID 5 în cazul scrierilor mici. În general, aceste modele necesită relocarea scrierilor și curățirea periodică. În multe sisteme aceasta se potrivește foarte bine cu sistemul de fișiere existent și/sau programele de aplicații. De exemplu, în unele sisteme, din motive care țin de recuperarea informațiilor, modificările unui fișier existent sunt efectuate întâi prin scrierea noii versiuni a fișierului (în spațiul liber de pe disc), și abia apoi ștergând versiunea anterioară a fișierului (eliberând spațiul ocupat de vechea versiune). Similar, reorganizarea periodică a discului este făcută astăzi în mod

obișnuit pentru a reduce fragmentarea, a recupera blocuri pierdute din cadrul blocului de fișiere s.a.m.d.

3.6.3.3 RAID 5/ NRP

Matricele de discuri au fost propuse pentru a asigura un spațiu mare de stocare de date, asigurând o înaltă fiabilitate la un cost redus. În organizarea RAID (matrice redundante de discuri independente) sunt $N+1$ discuri, acestea conținând un bloc de paritate pentru N blocuri de date, fiecare fiind stocat pe un disc diferit. (În acest subcapitol vom lua în considerație specificațiile RAID 5). Vom numi setul de blocuri de date, împreună cu blocul de paritate, grup de paritate.

Dacă un disc se defectează, fiecare bloc de date de pe acest disc se poate reface din celelalte N blocuri aflate în grupul de paritate.

În orice fel, aceasta poate duce la o substanțială degradare de performanță în timpul refacerii, deoarece discul defect poate dubla încărcarea pe discurile rămase, presupunând că fiecare disc este egal încărcat.

Acest inconvenient al RAID-ului tradițional poate fi rezolvat dacă se utilizează conceptul *clustered* RAID (RAID grupat) propus de Muntz și Lui.

Clustered RAID este un sistem RAID avansat în care grupurile de paritate de dimensiune G se întind peste *cluster*-ele discurilor de dimensiune C ($C > G$).

Într-un astfel de sistem încărcarea datorată discului defect este distribuită eventual peste toate cele $C-1$ discuri rămase, rezultând o degradare de performanță mai mică. Un *cluster* de dimensiune mare poate fi utilizat pentru a reduce timpul de refacere și reconstrucție a informațiilor de pe discul defect, ceea ce conduce la o performanță mai mare.

Cu toate acestea mărirea dimensiunii *cluster*-ului are un efect negativ prin mărirea ratei la al doilea eșec. Analizele arată că acest efect poate fi compensat de un timp de refacere mic, conducând la o însemnată îmbunătățire a timpului sistemului defect.

Deși sunt multe metode de refacere a informației de pe discul defect, nici una nu asigură o corespondență între blocurile de date și adresele de disc sau adresele blocurilor de date și adresele grupurilor de paritate pentru C și G date.

O schemă simplă care folosește toate combinațiile de G discuri luate câte C este greu de implementat, pentru că sunt (C/G) alegeri diferite și o tabelă care să rețină adresele blocurilor poate ocupa un spațiu foarte mare. În plus, spațiul se poate mări dacă numărul grupurilor de paritate din aria de discuri este substanțial mai mic decât (C/G) această schemă de distribuție nu asigură o încărcare aleatoare.

Considerăm $C=6$ și $G=3$. Presupunem că avem numai cinci grupuri de paritate; în acest caz o harta a grupurilor de discuri poate fi:

(1, 2, 3), (1, 2, 4), (1, 2, 5), (1, 2, 6), (2, 3, 4). Este clar că nu se asigură o încărcare aleatoare.

În continuare, prezint un algoritm rapid de distribuire a grupurilor de paritate de dimensiune G pe cele C discuri. Muntz și Lui notează că aceasta poate fi tratată ca o problemă de proiectare a unui bloc incomplet balansat. Holland și Gibson sugerează scheme bazate pe descrierea incompleta a blocurilor când acestea sunt disponibile. Convenim să egalizăm încărcarea de pe fiecare disc prin maparea aleatoare a fiecărui grup. Maparea este creată pentru fiecare grup de paritate bazându-ne pe permutări aleatoare. Această schemă este aplicată oricărei combinații de G și C ($G \leq C$). Cu toate acestea atât teoria cât și simularea practică arată că este o metodă rapidă și se ocupă foarte puțin spațiu.

După ce am stabilit caracteristicile conceptului *clustered* RAID, ne vom axa pe prezentarea modelului analitic, care este tratat în continuare și pe evaluarea performanțelor.

Procesele de citire și de scriere din aria de discuri sunt destul de complicate, astfel fiecare acces logic necesită un număr variabil de citiri și scrieri fizice, depinzând de faptul că se face o citire logică sau o scriere logică și de faptul ca operația s-a efectuat cu succes sau nu. Dacă accesul se face la un disc defect depinde de configurația datelor de pe acel disc și de amplasarea blocului de paritate. În plus, în timpul operației de refacere, procesele de reconstrucție a informației modifică nivelul de încărcare și, dacă se folosește o metodă complexă ca „reconstrucția cu redirectarea citirilor“, încărcarea sistemului poate fi foarte mare.

Modelul nostru analizează timpul de răspuns la citire și scriere într-un caz normal și în cazul când se defectează un disc, precum și timpul de refacere a informațiilor, urmărind efectele prin variația raportului C/G și încărcarea externă. Modul normal de funcționare este modelat utilizând o coadă $M/G/1$, în timp ce efectul unui singur disc defect și procesul de

reconstrucție a informației este captat de o coadă de tip *permanent customer*. Acesta este un model analitic larg utilizat, fiind general și aplicabil și altor probleme.

Rezultatele sunt validate prin comparare cu datele obținute din simulare. Analiza noastră arată că modelul *clustered RAID* este semnificativ mai tolerant la defectarea unui disc decât modelul RAID clasic. Atât timpul de refacere cât și degradarea de performanță din timpul refacerii datelor sunt substanțial reduse în modelul *clustered RAID*; mai mult decât atât, acest avantaj poate fi obținut prin utilizarea unui raport C/G mic.

3.6.3.3.1 Construirea distribuției bazate pe NRP

Considerăm o matrice de C discuri gestionate de un controler, cu un grup de paritate format din G blocuri ($G \leq C$, $G-1$ blocuri de date și un bloc de paritate). Fiecare bloc este plasat pe un disc diferit, astfel oricare bloc aflat pe discul defect poate fi refăcut din celelalte blocuri aflate în grupul de paritate. În continuare, presupunem că un bloc este o pistă de disc.

Presupunem că sunt B blocuri disponibile pe disc. Din punctul de vedere al sistemului de operare, întreaga matrice de discuri apare ca un singur disc logic, cu blocurile numerotate $0, 1, \dots, BC-1$; aceste blocuri includ și blocul de paritate. Subseturile $\{0, 1, \dots, G-1\}$, $\{G, G+1, \dots, 2G-1\}$, ... formează grupuri de paritate, iar ultimul bloc din fiecare grup este blocul de paritate.

Această configurație logică este oglindită într-o matrice fizica adresată astfel: $\text{Address}(i) = (\text{Disc}(i), \text{LocalBloc}(i))$, unde indicele i indică numărul blocului $\text{LocalBloc}(i)$ aflat pe discul $\text{Disc}(i)$.

Această schemă de adresare trebuie să îndeplinească unele condiții:

- trebuie ca maparea numerelor dintr-un grup de paritate să conțină discuri diferite;
- blocurile de paritate trebuie să fie reactualizate ori de câte ori un bloc aflat în grupul de paritate este modificat, acestea trebuie să fie distribuite uniform pe toate discurile;
- pagina de date cerută, aflată pe un disc defect cauzează citiri de pe celelalte $G-1$ discuri rămase. Pentru a echilibra încărcarea din timpul defectului trebuie ca blocurile de pe fiecare disc să fie uniform distribuite;
- adresele trebuie să fie ușor de calculat.

În continuare se prezintă o schemă de adresare simplă bazată pe permutări aleatoare. Pentru simplitate, presupunem că C este par și că G divide C , dar mai târziu se va înlătura această restricție.

Fie P_0, P_1, \dots permutări aleatoare ale $(0, 1, \dots, C-1)$ unde P_n mută j în $P_{n,j}$. Fie $n = \lfloor i/C \rfloor$ și $j = i \bmod C$, atunci schema de mapare asociază blocului i adresa $\text{Address}(i) \equiv (\text{Disk}(i), \text{LocalBlock}(i)) = (P_{n,j}, n)$. Este ușor de observat că următoarea schemă satisface condițiile 1), 2) și 3). Dacă G divide C nici un grup de paritate nu este obținut din permutări distincte și astfel nici care două blocuri din același grup de paritate, nu se pot afla pe același disc.

Cum fiecare permutare este aleatoare, blocul de paritate din fiecare grup, este distribuit aleator pe discuri și fiecare grup de paritate este mapat pe un set de discuri care sunt neuniform selectate din cele (C/G) subseturi de discuri.

Problema rămasă ar fi aceea dacă $P_{n,j}$ se poate calcula eficient. O variantă ar fi de a calcula P_n de fiecare dată când este nevoie de $P_{n,j}$. Aceasta este ineficientă atunci când C este mare și, prin urmare, se propune calcularea lui $P_{n,j}$ în $O(\log_2 C)$ pași utilizând metoda *Thorp's shuffle*.

3.6.3.3.2 Algoritmul „Thorp's shuffle“

Vom descrie acest algoritm exemplificând pe un teanc de cărți de joc:

- Se taie cărțile în două jumătăți egale A și B
- Se extrage la întâmplare o carte din teancul A și una din teancul B , apoi se extrage cea mai de jos carte din primul teanc și peste ea se pune cea mai de jos carte din al doilea teanc.
- Se repetă pasul al doilea până se termină cele două teancuri.

Exista $C=2N$ cărți numerotate $0, 1, \dots, 2N-1$, „amestecarea lui Thorp“ mută cartea j în poziția $T(j)$, unde U_j este un bit aleator și

$$T(j) = 2j + U_j \quad \text{pentru } j < N \quad (1)$$

$$T(j+N) = 2j + 1 - U_j \quad \text{pentru } j < N \quad (2)$$

Această amestecare se poate repeta până când cărțile sunt permutate aleator, astfel se generează a n -a permutare aleatoare P_n , începând cu permutarea inițială și aplicându-se K amestecări, utilizând o secvență de N biți aleatori $U_{m1}^{(n)}, U_{m2}^{(n)}, \dots, U_{mN}^{(n)}$ pentru m amestecări. Pentru a găsi $P_{n,j}$ vom urmări poziția cărții j după fiecare amestecare; poziția

după K amestecări este $P_{n,j}$. Experimentul se face pentru $C \leq 32$ și $U_{mj}^{(n)}$, d variind valorile lui G indicând $K=2\log_2(C)$ amestecări aleatoare ale permutărilor suficient pentru a evalua încărcarea datorată discului defect. În continuare, vom pune problema generării eficiente a biților aleatori.

3.6.3.3.3 Generarea secvenței de biți aleatori

Fiecare permutare din metoda de mai sus necesită K amestecări și fiecare amestecare folosește N biți aleatori. Este ușor de văzut că $U_{mj}^{(n)}$ este un număr de biți aleatori $NKn+N(m-1)+j$.

Atribuirem $U_{mj}^{(n)} = [2X_{NKn+N(m-1)+j} / M]$, unde X_i este un număr pseudo-aleator uniform distribuit în $[0, M)$.

Dacă X_i ar fi generat printr-o recurentă linear-congruenta $X_{i+1} = (aX_i + c) \bmod M$, atunci:

$$X_{i+k} = (a_k X_i + c_k) \bmod N \quad (3)$$

unde $a_k \equiv a^k \bmod M$ și $c_k \equiv (C(a^k - 1)/(a - 1)) \bmod M$

Scriem k în baza b : $k = k_0 + k_1 b + \dots + k_i b^i$. Acum X_k se poate calcula în $r = \lceil \log_b k \rceil$ pași utilizând (3),

$$X_{k_0+k_1b+\dots+k_i b^i} = (a_{k_i b^i} X_{k_0+k_1b+\dots+k_{i-1} b^{i-1}} + c_{k_i b^i}) \bmod M \quad i = 0, \dots, r \quad (4)$$

Atunci acestea sunt blocurile BC în aria de discuri și se arată ușor că

$k \leq BCk/2$ și, prin urmare, este nevoie să se stocheze numai valorile a_i și c_i pentru $i = pb^r$, $p = 1, 2, \dots, b - 1$ și $n = 0, 1, 2, \dots, \lceil \log_b (BCK/2) \rceil$

Stocarea necesită pentru aceasta $O(b \log_b (BCK/2))$. Alegând o

valoare pentru b potrivita (ex: $b=100$), X_k poate fi calculat rapid, fără cerințe mari de stocare.

Algoritmul *Thorp's shuffle* caută adresele de disc $Disk(nC+j) = P_{n,j}$ utilizând $K=2 \log_2 C$ numere pseudo-aleatoare $X_{q_1}, X_{q_2}, \dots, X_{q_K}$ pentru subsecvența $X_{NKn}, X_{NKn+1}, \dots, X_{NKn+1}$. Folosim această metodă o singură dată pentru a calcula X_{q_1} , iar X_{q_2}, \dots, X_{q_K} se calculează direct din (3):

$$X_{q_{i+1}} = X_{q_i+(q_{i+1}-q_i)} = (a_{(q_{i+1}-q_i)} X_{q_i} + b_{(q_{i+1}-q_i)}) \bmod M$$

Atunci $q_{i+1} - q_i \leq 2N=C$, X_{q_i} ; se poate afla într-un singur pas dacă se antecalculază a_k și b_k pentru $k \leq C$.

Timpul total pentru calcularea adreselor de disc pentru orice bloc folosind această metodă este $O(\log_b B + \log_b C)$ și costul de memorare a informației este $O(\log_b B + C \log_b C)$.

În prezentarea anterioară am impus ca G să fie un multiplu de C și C să fie par, pentru simplitate. Fiecare dintre restricții se poate înlătura.

Impunem ca G să dividă C ca să se respecte condiția (1) ca nici un grup de paritate să nu fie spart peste permutări distincte.

Pentru un G generalizat nu folosim o permutare numai o dată ci de $R = LCM(C, G)/C$ ori consecutiv. Astfel, fiecare permutare distribuie un set de blocuri RC consecutive, astfel nici un grup de paritate nu este spart peste permutări distincte, satisfăcând condiția (1).

Observație: C este necesar să fie par pentru a se putea aplica *Thorpe's shuffle*, iar dacă C nu este par se modifică *Thorpe's shuffle* alegând $N = \lfloor C/2 \rfloor$

și $U_{N,1} = 0$ totdeauna. Simularea arată că aceste modificări creează permutări aleatoare satisfăcătoare.

3.6.3.3.4 Un model pentru RAID5/NRP

Considerând aria de discuri formata din C discuri, cu un grup de paritate de dimensiune G . Traficul este dat de distribuția neuniformă a discurilor, sosirile cererilor de citire și scriere de pe fiecare disc trebuie sa fie în conformitate cu procesele date de metoda *Poisson*, procese care au ratele de sosire λ_r și λ_w ; $\lambda = \lambda_r + \lambda_w$. Unitatea de date cerută se presupune uniform distribuită pe discuri.

Fiecare disc are o coada FIFO pentru cereri. Timpul de servire pentru o cerere fizică este compus din timpul necesar căutării, latența și timpul necesar transferului de date. Timpul de căutare se poate aproxima astfel:

$Seek = a + b\sqrt{\Delta}$, unde Δ este fracțiunea pistei traversate de capul de disc. Latența se presupune uniform distribuită în $(0, t_r)$, unde t_r este timpul de rotație al discului (valorile considerate se pot vedea în Figura 3-24). Timpul de transfer pentru o pagină citită este constant t_p .

| Parameter | Denoted | Value |
|---------------------|---------|---------|
| arm movement time | a | 3ms |
| seek factor | b | 27ms |
| disk rotation time | t_r | 16.8ms |
| page transfer time | t_p | 2ms |
| page size | | 4KB |
| block size | | 1 track |
| block transfer time | t_b | 16.8ms |
| disks in cluster | C | 16 |
| blocks in group | G | 8 |
| pages per track | | 8 |
| tracks per disk | | 16,000 |

Figura 3-24

La scriere procesul este mai complex. Presupunem că dorim să înnoim informația din pagina V_0 , grupul de paritate conține paginile V_1, \dots, V_{G-2} și P (din blocul de paritate). Dacă noua valoare a paginii V_0 este V_0' , atunci noua valoare a lui P este, $P' = P \text{ XOR } V_0 \text{ XOR } V_0'$. Operația de înlocuire constă în două operații de citire/scriere pe discuri diferite. Prima citește V_0 și scrie V_0' , iar a doua citește P , calculează P' și îl scrie.

Timpul de calcul/transfer pentru fiecare din aceste operații este $t_r + t_p$, exceptând cazul când V_0 se află pe discul defect, când noua valoare de paritate se construiește după formula

$$P' = V_0' \text{ XOR } V_1 \text{ XOR } \dots \text{ XOR } V_{G-2}$$

Întârzierea în acest caz este maximum dintre întârzierile citirilor și timpul de scriere P' .

Când un disc se defectează, datele cerute de pe el sunt refăcute prin citirea datelor din celelalte blocuri aflate în grupul de paritate.

Presupunem că un proces de refacere reconstruiește datele de pe discul defect accesând datele de pe discurile valide. Presupunem că acest proces se desfășoară asincron. Procesul (inițiat de unitatea de control) citește de pe fiecare disc valid un bloc care face parte din grupul de paritate al discului defect și inserează într-un *buffer* aflat în unitatea de control, pentru a furniza la următoarea cerere direct din coadă. Pe discul disponibil, procesul reface blocul prin aducerea grupului de paritate din *bufferul* unității de control.

Vom urmări metoda de reconstrucție cu redirectarea citirilor.

Înnoirea informațiilor necesită două scrieri: una în pagina țintă și una în pagina corespunzătoare din grupul de paritate. Dacă se încearcă o scriere pe discul defect aceasta se redirecțiază spre o zonă deja refăcută.

3.6.3.3.5 Analiza performanței RAID5/NRP în modul de operare normal

Într-o operație efectuată cu succes ne interesează întârzierile operațiilor de citire și scriere și rezultatul maxim care se poate obține. Acesta se poate obține prin modelarea fiecărui disc printr-o coadă *M/G/1*, procesul dorit se obține prin aplicarea transformatei *Laplace* și *Poisson* distribuției timpului de căutare, latenței și timpului de transfer astfel:

$$\begin{aligned} Seek^*(s) &= E[e^{-sSeek}] = E[e^{-s(a+b\sqrt{|d|})}] = \\ &= \frac{(24 + 24bs + 8b^2s^2)e^{-(a+b)s} - (24 - 4b^2s^2)e^{-as}}{b^4s^4} \end{aligned}$$

$$Latency^*(s) = E[e^{-sLatency}] = (1 - e^{-st_r})/(t_r s)$$

$$Transfer_r^*(s) = E[e^{-sTransfer_r}] = e^{-st_r}$$

$$Transfer_w^*(s) = E[e^{-sTransfer_w}] = e^{-s(t_r + t_f)}$$

Prin urmare LT aplicată distribuției timpului total este:

$$\begin{aligned} B^*(s) &= E[e^{-sService}] = \\ &= Seek^*(s)Latency^*(s)(\lambda_r Transfer_r^*(s) + \lambda_w Transfer_w^*(s))/(\lambda_r + \lambda_w) \end{aligned}$$

LT aplicată distribuției timpului de așteptare în coadă este data de:

$$W^*(s) = E[e^{-sW}] = (s(1 - \lambda E[Service]))/(s - \lambda(1 - B^*(s)))$$

Întârzierea totală pentru citirile de pe discul i este suma dintre timpul de așteptare și timpul necesar citirii și, analog, întârzierea la scriere.

Cum traficul este neuniform distribuit pe discuri, întârzierea la scriere este egală cu cea la citire, astfel întârzierea locală la citire este egală cu cea la citirea din matrice și întârzierea la scriere în matrice este dublă față de întârzierea scrierii locale. Pentru a obține rezultatul maxim pentru fiecare cerere externă de modificare de date se fac două scrieri locale. Dacă ratele de sosire a citirilor și a scrierilor în matrice sunt Λ_r și Λ_w , respectiv pentru discuri Λ_r/C și Λ_w/C . Aplicând restricția $\lambda E[Service] \leq 1$ se obține rezultatul maxim în matrice.

$$\Lambda_r(a + 8b/15 + t_r/2 + t_p) + 2\Lambda_w(a + 8b/15 + 3t_r/2 + t_p) \leq C.$$

3.6.3.3.6 Analiza performanței RAID5/NRP în modul de operare în avarie

Analiza matricei de discuri în caz de defecțiune, este complicată de interacțiunea dintre accesul externe și procesul de reconstrucție *on-line* descris anterior.

Considerăm fiecare disc în parte. Fie T_m timpul de sosire al celei de-a m -a cereri de reconstrucție, $Y_m = T_m - T_{m-1}$ întârzierea pentru a $m-1$ cerere de reconstrucție și N_m numărul cererilor externe dinaintea sosirii cererii m . Dacă $N_m > 0$, atunci distribuțiile timpilor de căutare și de latență pentru a m -a cerere de reconstrucție sunt date de (6) și (7). Dacă $N_m = 0$ atunci cererea de reconstrucție urmează anterioarei cereri de reconstrucție imediat și atunci devine ușor de urmărit timpii de căutare și latență împreună. Deoarece fiecare bloc este o pistă și capul trebuie să schimbe pistele și apoi să aștepte sosirea începutului fiecărei piste, timpii combinați de căutare și latență trebuie să fie egali cu cel puțin timpul unei rotații de disc. Totuși, în timpul unei rotații, capul poate căuta de-a lungul unui număr mare de piste pentru presupusele valori ale parametrilor. Prin urmare, timpii de căutare și latență egalează t_r , cu LT:

$$SeekLatency^*(s) = E[e^{-sSeekLatency}] = e^{-st_r}$$

Rezumând, aceasta este o coadă cu sosiri de tip *Poisson*, cu rata λ și timpul de servire corespunzător $LT B^*(s)$ și un client permanent (procesul de reconstrucție) care accesează coada din nou îndată ce își termină serviciul, cu timpul de servire:

$$B_q = \begin{cases} SeekLatency + BlockTransfer & N_m = 0 \\ Seek + Latency + BlockTransfer & N_m > 0 \end{cases}$$

Fie Y întârzierea stării uniforme pentru cererile de reconstrucție și W timpul de așteptare al cozii pentru citiri și scrieri locale. Rezultă următoarele ecuații:

$$\begin{aligned}
 Y^*(s) &\equiv E[e^{-sY}] = Y^*(\lambda)a(s) + b(s)Y^*(\Phi(s)) \\
 E[Y] &= -(pa'(0) + b'(0))/(1 - \Phi'(0)) \\
 W^*(s) &\equiv E[e^{-sW}] = -(Y^*(s) - Y^*(\lambda(1 - B^*(s))))/((s - \lambda(1 - B^*(s))))B^*(s) \\
 E[W] &\equiv W^*(0) = (1 - \lambda B^*(0))E[Y^2]/(2E[Y])
 \end{aligned}$$

unde:

$$a(s) \equiv \text{BlockTransfer}^*(s)(\text{SeekLatency}^*(s) - \text{Seek}^*(s)\text{Latency}^*(s))$$

$$b(s) \equiv \text{Latency}^*(s)\text{BlockTransfer}^*(s)\text{Seek}^*(s)$$

$$p \equiv Y^*(\lambda) = \frac{\prod_{i=0}^{\infty} b(\Phi^i(\lambda))}{1 - \sum_{i=0}^{\infty} a(\Phi^i(\lambda))\prod_{j=0}^{i-1} b(\Phi^j(\lambda))}$$

$$\Phi(s) \equiv \lambda(1 - B^*(s)); \Phi'(s) \equiv \Phi(\Phi^{-1}(s)); \Phi^0(s) \equiv s$$

Valorile momentelor Y se pot obține prin diferențiere. Rezultatul acestui disc pentru procesul de refacere este $1/E[Y]$. Întârzierile pentru citiri și scrieri locale de pe acest disc sunt date prin suma timpilor corespunzători de servire. Totuși, citirile externe făcute de pe discul defect și toate actualizările externe pot implica discuri multiple.

Deoarece o parte a citirilor de pe discul defect pot fi redirectate către spațiul disponibil dacă datele au fost deja refăcute, restul trebuie reconstruite *on-line* prin citirea celor $G-1$ grupuri corespunzătoare. Întârzierea efectivă pentru o astfel de citire reprezintă maximum întârzierilor pentru citiri locale de pe cele $G-1$ discuri implicate. Deoarece traficul este uniform pe aceste discuri, întârzierile citirilor locale sunt distribuite în mod identic și media lor μ și variația σ^2 pot fi găsite ca și în cazurile de mai sus. Pentru simplitate vom aproxima aceste întârzieri ale citirilor, ca valori extreme independente ale variabilelor aleatoare cu media μ și variația σ^2 . Deci, maximum acestor $G - 1$ variabile aleatoare este o valoare extremă cu media $\mu + (\sqrt{\delta}\sigma/\pi)\ln(G-1)$ și variația σ^2 .

Acestea ne conduc la o aproximare pentru întârzierea citirilor externe echivalente unor citiri de pe discul defect.

Deoarece o actualizare externă depune date la fel în blocul de date precum și în corespunzătorul bloc de paritate, întârzierea totală depinde de cazul direcțării către discul defect. Rezultă trei cazuri:

- ambele scrieri locale sunt fie pe un disc valid, fie pe deja refăcuta porțiune a discului liber. Întârzierea totală este suma întârzierilor unor scrieri locale în modul normal;
- scrierea parității se face către discul defect. Nu este necesară paritatea scrisă, iar scrierea blocurilor de date nu necesită citire înainte de scriere. Întârzierea totală este egală cu întârzierea pentru citire a blocului de date;
- scrierea datelor către discul defect. Nu sunt necesare datele scrise, dar pagina de paritate trebuie reconstruită, necesitând citirea a $G-2$ pagini.

Aproximând întârzierile citirilor locale, din nou, ca valori extreme independente a variabilelor aleatoare cu media μ și variația σ^2 , întârzierea medie pentru actualizare externă este $t_m = 2\mu + (\sqrt{\delta}\sigma/\pi)\ln(G-2)$

După cum s-a menționat anterior, toate scrierile îndreptate către porțiunea discului defect, deja refăcută trebuie puse în coadă pentru a păstra integritatea blocurilor refăcute. O parte din citirile îndreptate către discul defect pot fi îndreptate către spațiul disponibil dacă porțiunea de date a fost reconstruită. Fie procentul citirilor și scrierilor externe către discul defect care sunt redirectate către discul disponibil la timpul t : $f_r(t)$ și $f_w(t)$. Fie discul defect numit disc 0, iar discul liber discul C. Atunci ratele de sosire efective la timpul t ale citirilor de pe discul valid i și de pe discul liber sunt:

$$\lambda_r(i, t) = \Lambda_r(1 + (G-1)(1 - f_r(t))/(C-1))/C + (G-2)(1 - f_w(t))\Lambda_w/(C(C-1))$$

$$\lambda_r(C, t) = \Lambda_r f_r(t)/C$$

Ratele de sosire efective la timpul t ale scrierilor de pe discul valid i și de pe discul liber sunt:

$$\lambda_w(i, t) = 2\Lambda_w/C - \delta$$

pentru $0 < i < C$

$$\lambda_w(C, t) = 2\Lambda_w f_w(t)/C$$

Termenul $\delta = 2A_w(1-f_w(t))/(C(C-1))$ ajustează încărcările operațiilor de citire și scriere deoarece anumite scrieri care implică discul defect prezintă o încărcare similară cu cea a citirilor. Folosind aceste rate de sosire externe, întârzierile instantanee pentru fiecare disc și rezultatul maxim instantaneu pot fi calculate. Întârzierile medii totale pentru citiri și scrieri externe se pot găsi și prin considerarea mediilor ponderate:

$$E[D_r^{(overall)}(t)] = ((C-1)E[D_r(1)] + f_r(t)E[D_r(0)] + (1-f_r(t))E[D_r(C)]) / C$$

$$E[D_w^{(overall)}(t)] = ((C-2)D_w(1,2) + f_w(t)D_w(1,C) + D_w(C,1) + (1-f_w(t))D_w(1,2))$$

În analizele anterioare am presupus că procesul de refacere de pe discurile valide și de pe cel liber operează independent. În realitate, acest lucru nu este adevărat; într-o situație stabilă, rata în care procesul de reconstrucție al discului liber consuma blocuri trebuie să egaleze rata la care procesul de refacere ale discurilor valide (considerate împreună) produce blocuri.

Aproximând, alegem această rată ca fiind minimum dintre cele două rate găsite. Fiecare bloc reconstruit pe discul liber consuma $G-1$ blocuri produse de discurile valide. Prin urmare, rata instantanee a reconstrucției este, aproximativ:

$$\lambda_{rec} \approx \min(1/E[Y_m^{(spare)}], (C-1)/(G-1)E[Y_m^{(surviving)}])$$

Procentul de disc defect deja recuperat la timpul t , $f_w(t)$ depinde de rata de reconstrucție $f_w(t) = B^{-1} \int_0^t \lambda_{rec}(t) dt$

De remarcat că la timpul t rata instantanee de refacere $\lambda_{rec}(t)$ depinde de $f_w(t)$.

Deși nu se poate obține o formă a expresiei mai clară pentru

$$f_w(0) = 0; f_w(t + \delta t) \approx f_w(t) + \lambda_{rec}(t)\delta t / B$$

considerând drept 0 timpul de start al refacerii și δt mic.

Recuperarea este completă când $f_w(t) = 1$. Procentul de citiri de pe discul defect care sunt redirectate către discul liber la timpul t este $f_r(t)$; în mod evident acesta nu poate depăși $f_w(t)$ deoarece doar acele blocuri deja refăcute pot fi citite de pe discul liber. În analizele noastre vom presupune că $f_r(t)/f_w(t) = \beta$ (constant).

3.6.3.3.7 Validarea rezultatelor

S-a construit un simulator *software* pentru a verifica soluțiile aproximative ale modelului descris. Acest simulator modelează

comportamentul discurilor în detaliu, de fapt urmărind mișcările capetelor de citire/scriere, iar timpii de servire necesari pentru fiecare cerere se bazează pe aceste date. Nici o aproximare nu a fost făcută cu privire la independența discurilor. Sunt testate atât performanțele în modul normal cât și de recuperare prin forțarea repetată a căderii discurilor, urmărindu-se procesul de refacere; sunt măsurați timpii de refacere și întârzierile de citire/actualizare. Aceste operații se repetă de un număr de ori suficient de mare. Se dă în Figura 3-25 și Figura 3-26 o comparație între procentul datelor recuperate și întârzierilor medii pentru accesele externe la diferite momente de timp în timpul procesului de refacere.

Toate comparațiile sunt făcute pentru grupuri de paritate de dimensiune $G=8$, încărcările corespunzând unei utilizări de 54% în mod normal, cu 80 % dintre accese fiind citiri și redirectarea totală a citirilor ($\beta=1$).

După cum se poate vedea din grafice există o asemănare mare între aproximările analitice și rezultatele experimentale. În majoritatea cazurilor diferența este mai mică decât 6%. Alte comparații, pentru grade diferite de utilizare, arată de asemenea o potrivire.

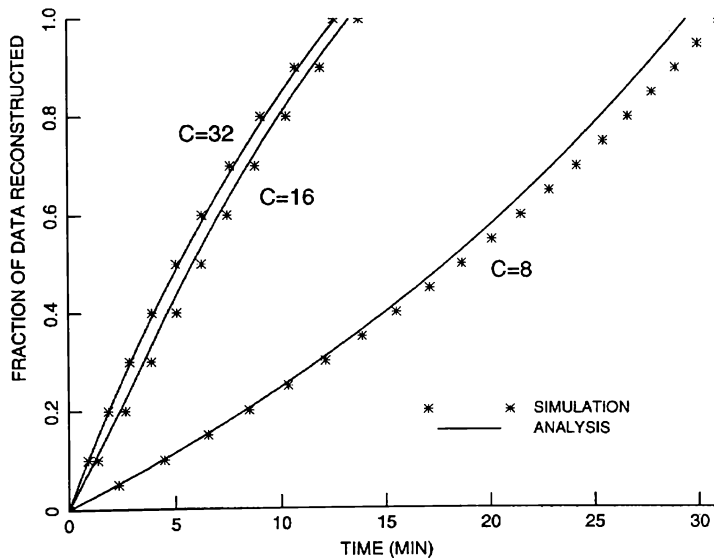


Figura 3-25

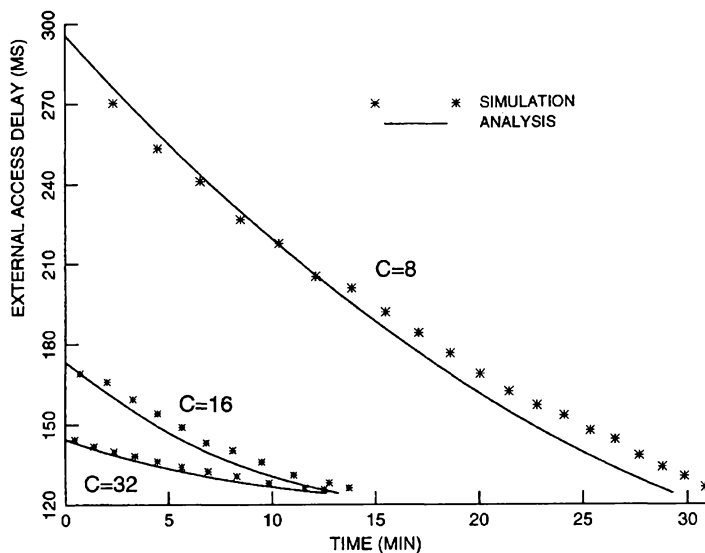


Figura 3-26

De remarcat că procesul măsurat este unul tranzitoriu și experimentele sunt prin urmare mari consumatoare de timp. În majoritatea cazurilor unei valori generate în câteva secunde de modelul analitic îi sunt necesare câteva ore pentru a fi produsă prin simulare.

3.6.3.3.8 Evaluarea performanțelor

În acest subcapitol vom prezenta evaluarea performanțelor modelului *clustered* RAID bazat pe modelul analitic prezentat anterior. Am calculat timpii de recuperare (reconstrucție) și întârzierile la citire/actualizare în modul normal și în timpul recuperării informațiilor. Obiectul acestui studiu este acela de a înțelege:

- impactul raportului C/G asupra performanțelor obținute în mod normal și în cel al recuperării;
- efectul redirectării strategiei de citire.

După cum vom vedea pentru *unclustered* RAID, discurile „supraviețuitoare“ reprezintă o gâtuire a performanțelor, iar regândirea strategiei operației de citire poate îmbunătăți performanțele. Pentru *clustered* RAID chiar și cu un raport C/G mic, de valoare 2, spațiul liber de pe disc devine o povară, iar redirectarea strategiei de citire prelungeste timpul de reconstrucție.

În studiul de față am presupus o matrice cu $G=8$ și C variind între 8 și 64. Ratele de sosire pentru citire/actualizare sunt considerate pentru o utilizare de 45% în modul normal și un procent al citirilor de 0.8 (80% dintre cereri sunt citiri). În cadrul redirectării strategiei de citire, valoarea inițială pentru raportul $\beta \equiv f_r(t)/f_w(t)$ este 1, toate cererile de citire care pot fi servite de spațiul liber sunt așezate în coadă. Toate valorile prezentate se bazează pe o soluție analitică aproximativă.

Primul grafic, Figura 3-27, arată schimbarea în timpul de recuperare în momentul în care crește raportul C/G , pentru rate de sosire a cererilor corespunzând unor utilizări de 25%, 45% și 65% în mod normal. Pentru o utilizare de 45%, timpul de recuperare cu $C/G=2$ este jumătate din cel al *unclustered* RAID ($C/G=1$); la o utilizare de 65%, *unclustered* RAID se saturează — nu este posibilă o recuperare *on-line* (și nu este punctat acest lucru în grafic). Cu $C/G \geq 2$, totuși, este posibilă recuperarea.

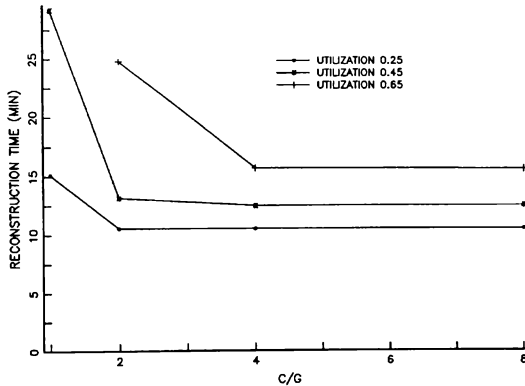


Figura 3-27

Figura 3-28 și Figura 3-29 arată întârzierile de citire/scriere în timpul recuperării pentru $C/G=1, 2, 4$ și 8 . În ambele cazuri, întârzierile în sistemele *clustered* sunt mai scăzute decât cele din sistemele *unclustered*. Câștiguri mari se obțin prin mărirea raportului C/G de la 1 la 2; o mărire în continuare a lui C reduce destul de puțin întârzierile medii. (Această diferență este ridicată pentru încărcări superioare, dar principala tendință se păstrează — o reducere mare a întârzierilor când C/G crește de la 1 la 2 și avantaje mici dacă C crește în continuare). Trebuie remarcat că procesul de reconstrucție impune o încărcare ridicată a sistemului, după cum arată marea diferență dintre întârzierile normale și întârzierile din timpul recuperării.

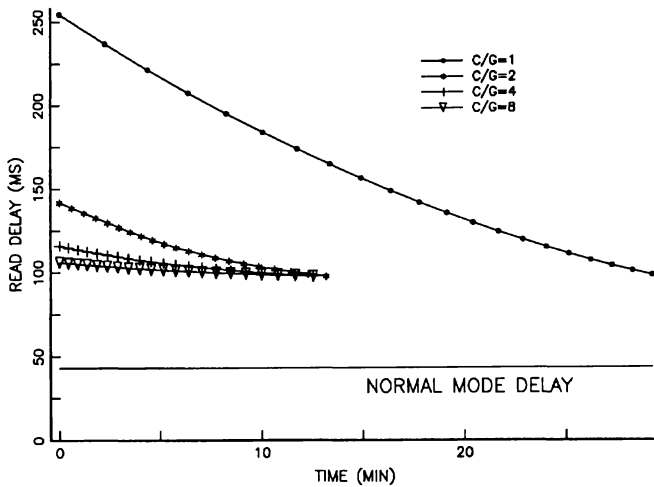


Figura 3-28

În Figura 3-30 vedem cum funcția $f_w(t)$ se modifică cu timpul. O trăsătură interesantă a acestui grafic este aceea că panta curbei pentru $C/G=1$ crește cu timpul, în vreme ce panta curbei pentru *clustered arrays* ($C/G \geq 2$) scadește cu timpul. Aceasta se explică prin faptul că în primul caz, discurile „supraviețuitoare“ reprezintă o gâtuire a procesului de reconstrucție; în vreme ce spațiul liber este realocat, unele dintre citirile de pe discul defectat sunt transferate în spațiul liber, iar rata de reconstrucție a discurilor funcționale crește. Pentru *clustered arrays* spațiul liber reprezintă o gâtuire a procesului de reconstrucție și, în vreme ce crește $f_w(t)$, încărcarea spațiului liber crește și ea, iar rata de reconstrucție se află într-un ușor declin. Acest fapt explică și comportamentul din Figura 3-31, care arată că pentru $C/G=1$, timpul de reconstrucție scade dacă $f_r(t)/f_w(t)$ crește și, prin urmare, scade încărcarea discurilor nedefectate; timpul de reconstrucție pentru $C/G \geq 2$ arată un comportament opus. Prin urmare, strategia de redirectare a citirilor este folositoare pentru *unclustered RAID*, dar nu și pentru *clustered RAID*, chiar și pentru un raport scăzut al C/G de valoare 2. Cum citirile aleatoare deranjează procesul normal de reconstrucție pistă cu pistă de pe discul liber în care apar și întârzierile căutărilor în plus, strategia este valabilă doar când discurile nedefectate reprezintă o gâtuire a performanțelor în cazul *unclustered RAID*.

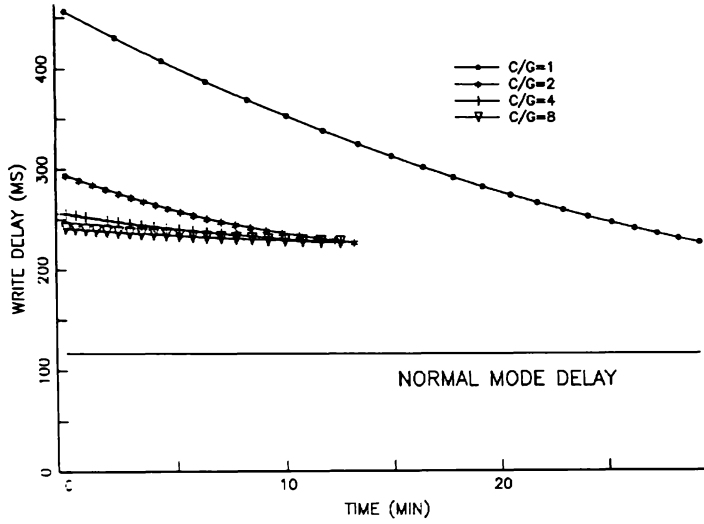


Figura 3-29

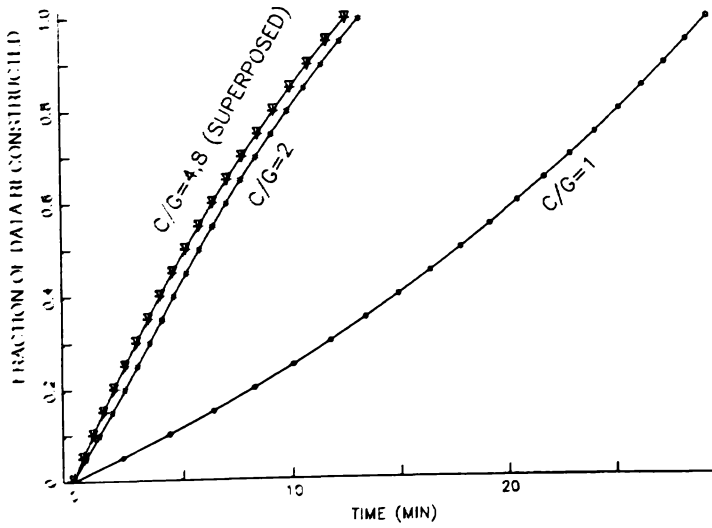


Figura 3-30

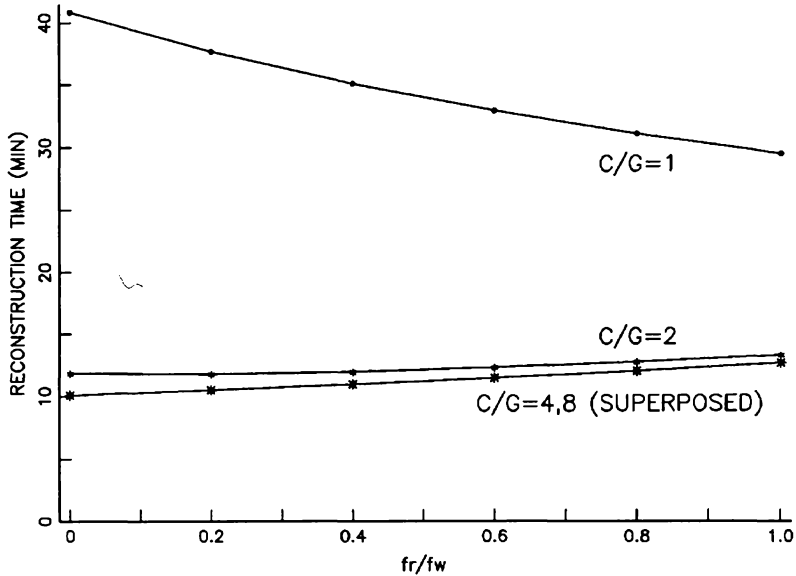


Figura 3-31

Sistemul *unclustered* este, de asemenea, sensibil și la raportul operațiilor de citire (procentul de cereri care sunt citiri), după cum arată Figura 3-32. În timpul recuperării, o citire de pe discul defectat care trebuie reconstruit implică $G-1$ citiri de pe discul funcțional; scrierea pe acest disc implică doar $G-2$ citiri de pe discurile funcționale. Prin urmare, măbind procentul citirilor se mărește încărcarea discurilor funcționale în timpul recuperării; aria de discuri *unclustered* prezintă o creștere semnificativă a timpului de reconstrucție o dată cu creșterea procentului de citiri. Matricele *clustered* sunt puțin afectate în acest caz.

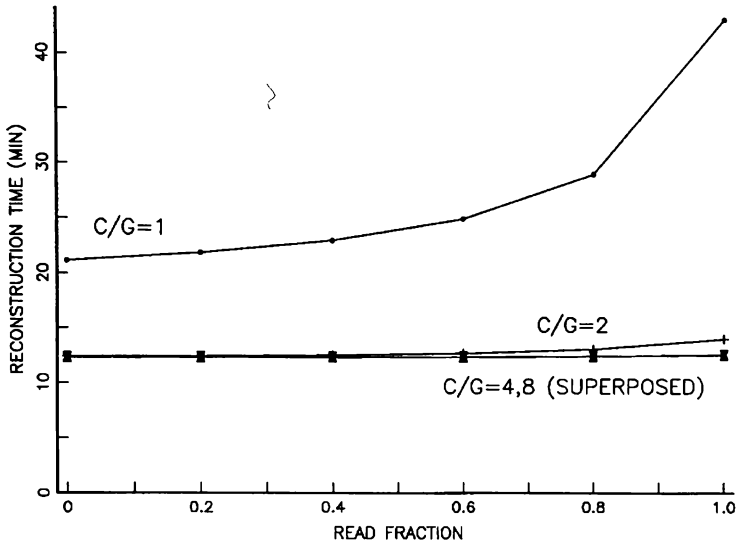


Figura 3-32

Reiese în mod clar din aceste rezultate că folosind *clustered* RAID se îmbunătățește toleranța sistemului la defectarea vreunui disc, luând în considerație fie timpul de recuperare, fie performanțele obținute în timpul recuperării. Prin urmare, majoritatea acestor câștiguri pot fi obținute prin folosirea unor rapoarte C/G de valoare mică.

3.7 RAID level 6

RAID level 6, definit pentru prima dată în [Katz89] este o dezvoltare a RAID level 4 și 5 care asigură o înaltă disponibilitate a datelor. RAID level 6 utilizează două sume de control calculate independent care protejează toate blocurile utile de date. În cazul scrierilor două blocuri de paritate independente trebuie actualizate astfel ca rezultă o puternică scădere de performanță, ceea ce impune obligatoriu utilizarea unor metode auxiliare, cum ar fi *logging* sau *write-back cache*, pentru a reduce aceste influențe negative.

Principalul avantaj adus de RAID level 6 constă în înalta disponibilitate a datelor. Pentru ca datele să ajungă inaccesibile datorită unor căderi de discuri, trei discuri trebuie să cadă în intervalul de timp MTTR. Această posibilitate este de mii de ori mai puțin probabilă decât cădere a două discuri într-un RAID 1 sau RAID 5.

3.7.1 Descrierea RAID level 6

După cum am mai amintit, este necesară calcularea independentă a două date de control. Aceasta se poate realiza prin utilizarea a doi algoritmi diferiți care să acopere aceleași date sau utilizând aceiași algoritmi, dar cu acoperiri diferite ale datelor.

Utilizare a doi algoritmi diferiți de calculare a datelor de control peste aceleași date este cunoscută sub denumirea de *redundanță unidimensională* sau *P+Q parity*. În Figura 3-33 este prezentat un exemplu în care cu B_i am notat blocul de date utile și cu $P-i,j$, respectiv $Q-i,j$ am notat datele de control corespunzătoare blocurilor i și j .

| DISC 0 | DISC 1 | DISC 2 | DISC 3 | DISC 4 | DISC 5 |
|--------|---------|---------|--------|--------|--------|
| B0 | B1 | B2 | B3 | P-0,3 | Q-0,3 |
| B4 | B5 | B6 | P-4,7 | Q-4,7 | B7 |
| B8 | B9 | P-8,11 | Q-8,11 | B10 | B11 |
| B12 | P-12,15 | Q-12,15 | B13 | B14 | B15 |
| ... | ... | ... | ... | ... | ... |

Figura 3-33

| DISC 0 | DISC 1 | DISC 2 | DISC 3 | DISC 4 | DISC 5 | DISC 6 | DISC 7 | DISC 8 | DISC 9 | DISC 10 |
|---------|---------|---------|--------|--------|--------|------------|--------|--------|--------|------------|
| P-0,6 | P-2,8 | P-4,10 | B0 | B2 | B4 | P-0,2,4 | B6 | B8 | B10 | P-6,8,10 |
| P-1,7 | P-3,9 | P-5,11 | B1 | B5 | B5 | P-1,3,5 | B7 | B9 | B11 | P-7,9,11 |
| P-12,18 | P-14,20 | P-16,22 | B12 | B14 | B16 | P-12,14,16 | B18 | B20 | B22 | P-18,20,22 |
| P-13,19 | P-15,21 | P-17,23 | B13 | B15 | B17 | P-13,16,19 | B19 | B21 | B23 | P-19,21,23 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| | | | DATE | DATE | DATE | | DATE | DATE | DATE | |

Figura 3-34

Utilizarea unui singur algoritm de calculare a datelor de control care acoperă diferite subseturi de date este cunoscută sub denumirea de *redundanța bidimensională*. În Figura 3-34 este prezentat un exemplu în

care cu B_i am notat blocul i de date utile și cu $P_{-i,j,k}$ am notat datele de control corespunzătoare blocurilor i,j,k .

Se poate observa că discurile 3, 4, 5 și 6 formează un RAID 4, la fel ca și discurile 7, 8, 9 și 10. De asemenea, discurile 0, 3 și 7, la fel ca și 1, 4 și 8 și 2, 5 și 9.

Referitor la costurile implicate de RAID level 6 cu redundanța bidimensională se poate afirma că $M \times N$ discuri de date utile presupun $M+N$ discuri de paritate. Pentru matrice de dimensiuni mici este scumpă din punct de vedere al capacității redundante (în exemplul considerat fiind de aproximativ 82%). În cazul matricelor de dimensiuni mai mari, de exemplu 100 de discuri, rezultă o redundanță de 20%.

Deoarece funcția XOR nu mai prezintă (față de altele) avantajul timpului de evaluare (datorită progreselor tehnologice – ex. ASICs), tendința în cadrul RAID 6 este către redundanța unidimensională cu utilizarea unor coduri *Reed-Solomon*.

3.7.1.1 Comportare RAID level 6 în funcționare normală

În funcționare normală, procesele de scriere și citire sunt asemănătoare cu cele întâlnite la RAID level 4 și 5. La citire, cererile sunt direcționate către două discuri independente obținându-se performanțe foarte bune.

Performanța la scriere este, însă, extrem de slabă, datorită faptului că trebuie actualizate amândouă blocurile de control implicate. Din acest motiv, utilizare unor tehnici de accelerare a scrierii, precum *logging* sau *write-back cache* se impun cu necesitate.

3.7.1.2 Protecția la căderi în RAID level 6

Când un disc cade, fie de date sau paritate, regenerare se realizează identic cu cea în cazul RAID 5 sau 4.

O a doua cădere de disc, înainte de înlocuirea primului, va necesita pentru regenerare rezolvarea unor sisteme de două ecuații cu două necunoscute, întotdeauna compatibile.

3.7.1.3 Domeniile de utilizare pentru RAID level 6

Datorită prețului foarte ridicat, dar și a fiabilității deosebite, această soluție se pretează doar în acele aplicații care necesită asigurarea unei siguranțe foarte bune a datelor, cu acceptare unei performanțe slabe în cazul unor operații intensive de scriere.

3.8 Evaluarea performanțelor matricelor de discuri – RAID

În continuare, se vor compara performanțele obținute în cazul în care toate discurile sunt operaționale (modul normal) pentru RAID 1 și RAID 5. Toate aceste evaluări de performanță sunt considerate pentru două tipuri de aplicații:

- aplicații în care cererile I/O sunt pentru cantități mici de date (exemplu: stații de lucru);
- aplicații în care cererile I/O sunt pentru cantități mari de date (exemplu: supercalculatoare, procesarea de imagini).

Considerăm aria de discuri fiind formată din N discuri identice. Presupunem că rotațiile discurilor sunt sincronizate, dar nu și cele ale brațelor. În arhitecturile considerate, data este reprezentată printr-un bloc întretesut de-a lungul a W discuri, unde W este lățimea secțiunii, $W \leq N$. De exemplu, un fișier este divizat logic în blocuri f_1, f_2, \dots, f_m care sunt stocate pe discuri consecutive în cadrul secțiunii (n poate fi mai mare decât lățimea secțiunii W). Dacă o cerere accesează exact un întreg multiplu de W blocuri de date dintr-o secțiune (aliniat la marginile secțiunii), atunci se numește o cerere I/O *full stripe* (pe toată secțiunea). Altfel, ea se numește cerere I/O *partial stripe* (pe o parte a secțiunii). Principalul avantaj al întreteserii blocurilor este acela că astfel este posibilă execuția concomitentă a mai multor cereri I/O mici. Alte alternative sunt întreteserile la nivel de bit sau la nivel de *byte*, exemplificate prin RAID 2 și RAID 3. Totuși, ambele necesită ca toate discurile din matrice să fie implicate în servirea fiecărei cereri I/O și, astfel, cel mult o cerere I/O poate fi executată la un moment dat. Cel mai semnificativ avantaj obținut din aceste arhitecturi este o rată de transfer ridicată. În vreme ce aceste arhitecturi sunt potrivite pentru aplicații în care orice cerere transferă cantități mari de date, ele nu sunt potrivite pentru prelucrarea informațiilor de afaceri sau a mediului stațiilor de lucru. În acest caz, deoarece fiecare cerere transferă o cantitate scăzută de date, avantajul ratei de transfer ridicate scade. Pe de altă parte, deoarece toate discurile trebuie să deservească o singură cerere, aria de discuri nu poate suporta operații I/O multiple în mod concurrent. Un avantaj evident al întreteserii datelor de-a lungul secțiunilor într-o matrice de discuri este încărcarea echilibrată a fișierelor des folosite pe discuri multiple.

3.8.1.1 Analiză comparativă între RAID 1 și RAID 5 pentru volum mare de cereri

În continuare, se studiază implicarea matricelor de discuri în aplicații cum ar fi procesarea informațiilor din domeniul afacerilor, precum și al mediului stațiilor de lucru. Se vor prezenta rezultate analitice și experimentale pentru arhitecturi de matrice de discuri împreună cu diverse strategii de proiectare. Deoarece s-a observat în literatura de specialitate că pentru aceste tipuri de aplicații cererile I/O accesează în mod obișnuit o cantitate mică de date, facem presupunerea că fiecare cerere accesează un singur bloc de date (4096 bytes) și că fiecare bloc distinct este în mod egal accesat. În acest caz, discurile pot servi cererile în mod independent una de cealaltă. În cele din urmă presupunem că cererile sosesc în concordanță cu procesul de sosire *Poisson*.

Cel mai important avantaj al RAID 5 în comparație cu RAID 1 este prețul mai scăzut al acestuia, deoarece RAID 1 necesită un număr de discuri aproape dublu pentru a atinge aceeași capacitate, sau, folosind același număr de discuri, RAID 1 ajunge doar la jumătate din capacitatea RAID 5. Dat fiind costul mare plătit de RAID 1, ne interesează câștigul de performanță atins de RAID 1 față de RAID 5.

În acest subcapitol vom analiza RAID 5 în conjuncție cu o strategie AR. Apoi, folosind același model, vom compara performanțele RAID 5 cu varianta *group-rotate declustering* a RAID 1 în conjuncție cu o strategie SQ. De aici înainte, pentru a nu apărea vreo confuzie, vom folosi RAID 1 pentru a ne referi la varianta *group-rotate declustering* în combinație cu SQ.

Considerăm o matrice de discuri cu N discuri. Bazat pe modelul prezentat anterior, timpul de servire al discului pentru o cerere de citire este:

$$Y_r = X + T \quad (1)$$

iar pentru o cerere de scriere este

$$Y_w = X + T + R_{\max} \quad (2)$$

unde X este suma latențelor de căutare și de rotație, T este timpul de transfer (deoarece presupunem un acces la un sigur bloc, $T = \tau$) și R_{\max} este timpul pentru o rotație completă. (De remarcă că după timpul de transfer la o citire T , discul trebuie să se rotească până la începutul

blocului de date și apoi să scrie din nou. Suma timpilor de rotație și de transfer la scriere este egală cu timpul pentru o rotație completă R_{max}).

Pentru ușurința studiului, se introduc următoarele notații în analizele noastre:

- λ : rata de sosire a cererilor către subsistemului I/O;
- $\lambda_p(i), \lambda_d(i)$: ratele de sosire pentru cererile prioritare de actualizare a parității și cele de prioritate scăzută a citirilor/scrierilor pe discul i , $i=1, 2, \dots, N$;
- $Q_p(i), Q_d(i)$: r.v.; indică timpii necesari operațiilor efectuate asupra cozilor în urma unor cereri de prioritate ridicată, respectiv scăzută asupra discului i ;
- Y_p, Y_d : r.v.; pentru timpul de servire a cererilor cu prioritate ridicată și scăzută;

În modelul considerat, cererile I/O se presupun că sosesc în aria de discuri după algoritmul *Poisson* cu rata λ . Sunt păstrate două cozi în fața fiecărui disc, una pentru cererile de citire/scriere a datelor (coada D), iar cealaltă pentru cererile de paritate (coada P). Când sosește o cerere, ea este îndreptată către discul având probabilitatea p_i , $i=1, 2, \dots, N$. Intrările în coada D a discului i sunt descrise în algoritmul *Poisson* prin parametrul $\lambda_d(i) = p_i \lambda$, unde cererea poate fi de citire, având probabilitatea p_r , iar cea de scriere cu probabilitatea $p_w = 1 - p_r$. Intrările în coada P a discului i sunt descrise printr-o suprapunere de $N-1$ cereri de paritate. Când N este mare, el poate fi aproximat în algoritmul *Poisson* prin parametrul $\lambda_p(i) = (1-p_i)p_w \mathcal{N}(N-1)$. Cele $N-1$ procese de generare a parității nu sunt independente unul față de celălalt. Totuși, când N este mare, corelațiile între procese se micșorează și, astfel, se ajunge la concluzia că pentru $N \rightarrow \infty$ ele devin independente. De fapt, când $N=16$, rezultatele experimentale arată o mare apropiere față de rezultatele analitice. În sfârșit, strategia primul sosit-primul servit este folosită în cazul ambelor cozi.

În acest fel, discurile i ($i=1, \dots, N$) prin două cozi bazate pe priorități, în care timpul de servire a cererilor cu prioritate scăzută (citiri/scrieri în coada D) sunt obținute din relațiile (1) și (2). Se definește θ ca fiind o variabilă asociată unei cereri și care ia valoarea 1 dacă este o citire sau 0, în caz contrar. Atunci vom avea: $Y_p = Y_w$, și $Y_d = \theta Y_r + (1-\theta) Y_w$, unde

$$P\{\theta = 1\} = p_r \text{ și } P\{\theta = 0\} = p_w.$$

Primele două valori pentru Y_p și Y_d sunt :

$$\overline{Y}_p = \overline{X} + \tau + R_{\max}$$

$$\overline{Y}_p^2 = \overline{X}^2 + 2(\tau + R_{\max})\overline{X} + (\tau + R_{\max})^2$$

$$\overline{Y}_d = \overline{X} + \tau + p_w R_{\max}$$

$$\overline{Y}_d^2 = p_r(\overline{X}^2 + 2\tau\overline{X} + \tau^2) + p_w(\overline{X}^2 + 2(\tau + R_{\max})\overline{X} + (\tau + R_{\max})^2)$$

unde

$$\overline{X} = \overline{S} + \overline{R}, \quad \overline{X}^2 = \overline{S}^2 + 2\overline{S}\overline{R} + \overline{R}^2 \quad (2)$$

Bazându-ne pe rezultatele cozilor cu priorități non-preemptive, întârzierile medii a cozilor P și D corespunzătoare discurilor i sunt:

$$\overline{Q}_p(i) = \frac{\lambda_d(i)\overline{Y}_d^2 + \lambda_p(i)\overline{Y}_p^2}{2[1 - \lambda_p(i)\overline{Y}_p]} \quad (3)$$

$$\overline{Q}_d(i) = \frac{\lambda_d(i)\overline{Y}_d^2 + \lambda_p(i)\overline{Y}_p^2}{2[1 - \lambda_p(i)\overline{Y}_p][1 - \lambda_d(i)\overline{Y}_d - \lambda_p(i)\overline{Y}_p]} \quad (4)$$

Timpu mediu de răspuns la o citire de pe discul i este:

$$\overline{Z}_r(i) = \overline{Q}_d(i) + \overline{X} + \tau \quad (5)$$

iar timpul mediu de răspuns la o scriere este:

$$\overline{Z}_w(i) = \overline{Q}_d(i) + \frac{1}{N-1} \sum_{j=1, j \neq i}^N \overline{Q}_p(j) + 2\overline{X} + 2\tau + R_{\max} \quad (6)$$

În (6), primii doi termeni corespund întârzierilor medii în timp ce se așteaptă citirea blocurilor vechi de date și de paritate; al treilea termen corespunde latențelor medii de căutare și de rotație pentru operațiile de citire; al patrulea termen corespunde timpilor de transfer asociați celor două operații de citire; al cincilea termen corespunde unei necesare rotații și scrieri a următorului bloc de paritate după citirea vechiului bloc de paritate.

În final, timpii medii de răspuns sunt calculați ca o medie a tuturor discurilor:

$$\overline{Z}_r = \sum_{i=1}^N p_i \overline{Z}_r(i)$$

$$\overline{Z}_w = \sum_{i=1}^N p_i \overline{Z}_w(i)$$

$$\bar{Z} = p_r \bar{Z}_r + p_w \bar{Z}_w \quad (7)$$

Rezultatele prezentate aici pentru RAID 5 sunt obținute prin alegerea lui $p_i = 1/N$, $i = 1, 2, \dots, N$. Deoarece în studiul de față RAID 5 folosește o unitate mică a secțiunii (4K bytes) și păstrează datele pe discuri multiple, este de înțeles presupunerea făcută cum că vin cereri cu o probabilitate egală pe fiecare disc.

Compararea performanțelor este orientată în două direcții. În primul rând, considerăm că ambele arhitecturi (RAID 1 și RAID 5) au aceeași capacitate (numărul de discuri la RAID 1 este $2n$, iar la RAID 5 este $n+1$), unde n a fost ales egal cu 8. În al doilea vom compara cele două arhitecturi având același număr de discuri, caz în care RAID 1 pierde jumătate din capacitatea sa. Rezultatele sunt prezentate în Figura 3-35, din care se observa că, în primul caz, RAID 1 poate suporta mai mult decât de trei ori rata de sosire a cererilor I/O decât RAID 5. În al doilea caz, RAID 1 încă poate suporta un număr de două ori mai mare de cereri I/O într-o secundă decât RAID 5. De aici deducem că RAID 1 se comportă în mod semnificativ mai bine decât RAID 5 în medii în care cererile I/O sunt scăzute, dar cu costul sporit prin dublarea numărului de discuri sau de pierdere a jumătate din capacitate.

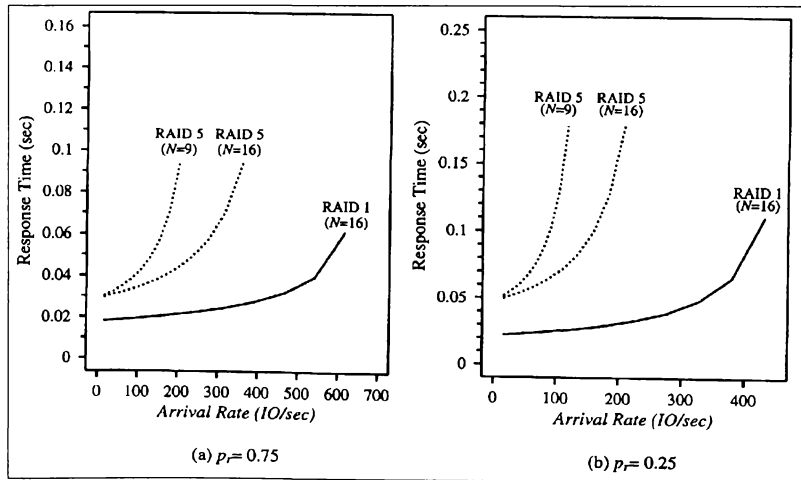


Figura 3-35

3.8.1.2 Analiză comparativă între RAID 1 și RAID 5 pentru cereri de volum mare

Aplicații precum calcule superioare (*supercomputing*) și procesare a imaginilor necesită o programare la un nivel mai scăzut și, prin urmare, sunt necesare mai puține cereri I/O per unitatea de timp. Totuși, fiecare cerere I/O accesează de obicei o cantitate ridicată de date. În general, parametrii de calcul sunt mutați în cantitate mare de pe disc în structurile de date rezidente în memorie, iar rezultatele sunt scrise periodic înapoi pe disc. În acest caz mai multe discuri lucrează împreună ca un dispozitiv logic unic furnizând o rată de transfer mai ridicată. Cum s-a arătat anterior, se presupune că rotațiile tuturor discurilor din matrice sunt sincronizate. Datele sunt depuse în secțiuni de-a lungul tuturor discurilor în același cilindru. Prin urmare, deoarece fiecare cerere I/O accesează una sau mai multe secțiuni, timpii de căutare se presupun a fi identici pentru fiecare disc. Uneori poate exista o diferență de o pistă între discuri, dar acest lucru este neglijabil. În cazul arhitecturilor *mirrored declustering* și *group-rotate declustering*, discurile sunt împărțite în două secțiuni, unul pentru fiecare copie a datelor. Căutările se consideră că sunt sincronizate într-un grup, dar cele două grupuri pot rezolva în mod independent cereri și pot fi modelate ca servere cu două cozi de cereri separate. Deoarece toate discurile din matrice pot fi implicate în servirea unei cereri în cazul *chained declustering* și RAID 5, se presupune că toate sunt sincronizate și sunt modelate ca un singur server cu o singură coadă de cereri.

În cele ce urmează vom prezenta, pe scurt, un model analitic pentru RAID 5 sincronizat și îl vom compara cu RAID 1 pentru cazurile în care cele două matrice au aceeași capacitate și același număr de discuri, la fel cum am făcut în contextul operațiilor reduse de I/O.

Presupunând că toate discurile din matrice sunt sincronizate, timpul de servire pentru o cerere de citire este $Y_r = X + T$. Pentru cereri de scriere, deoarece operații mari de I/O pot implica mai multe secțiuni de date, se disting două cazuri. În cazul *full stripe write*, în care blocurile de date care trebuie actualizate încep și se sfârșesc la limitele secțiunii, timpul de servire a scrierilor este același cu al citirilor deoarece blocul de paritate al fiecărei secțiuni poate fi calculat din noile blocuri de date. În acest caz nu sunt necesare alte operații de citire în plus. Totuși, dacă blocul de început sau de sfârșit nu este aliniat la marginile secțiunii, apare cazul parțial *stripe write* și trebuie urmată o procedură *read-update*. Pentru simplitate, presupunem că scrierea unei secțiuni parțiale este necesară doar pentru ultima secțiune a fiecărei operații I/O,

considerând că blocul de început este întotdeauna aliniat la marginea secțiunii. Aceste fapte pot conduce la rezultate bune ale RAID 5.

Fie p_f probabilitatea că o cerere de scriere este o scriere în toată secțiunea. Pentru o scriere parțială în secțiune, deoarece accesează mai multe blocuri de pe fiecare disc, după citirea unui bloc din ultima secțiune (pentru calculul noului bloc de paritate), discul trebuie să se rotească la începutul celor N_s blocuri de date. Prin urmare, în locul unei rotații complete, este necesar ca rotația să fie făcută doar pe distanța rămasă până la scrierea noilor date. Să reamintim că N_b arată numărul de blocuri de pe fiecare pistă. Timpul de rotație este:

$$R = \left(1 - \frac{N_s \bmod N_b}{N_b}\right) R_{\max}$$

cu valorile momentane:

$$\bar{R} = (1 - \eta) R_{\max}$$

$$\bar{R}^2 = \left(1 - 2\eta + \frac{1}{N_b^2} E[(N_s \bmod N_b)^2]\right) R_{\max}^2$$

unde:

$$\eta = E\left[\frac{N_s \bmod N_b}{N_b}\right] =$$

$$= \frac{1}{N} \sum_{k=0}^{N_b-1} k \sum_{m=0}^{\infty} P\{N_s = mN_b + k\}$$

și

$$E[(N_s \bmod N_b)^2] = \sum_{k=0}^{N_b-1} k^2 \sum_{m=0}^{\infty} P\{N_s = mN_b + k\}$$

Atunci, timpul de servire pentru o scriere parțială în secțiune este:

$$Y_w = X + \tau + R + T$$

Valorile momentane pentru timpul de servire total sunt:

$$\bar{Y} = p_r \bar{Y}_r + p_w p_f \bar{Y}_r + p_w (1 - p_f) \bar{Y}_w =$$

$$= \bar{X} + \bar{T} + p_w (1 - p_f) (\tau + \bar{R})$$

$$\bar{Y}^2 = \bar{X}^2 + 2\bar{X}\bar{T} + \bar{T}^2 + p_w (1 - p_f) (\bar{R}^2 + \tau^2 + 2\bar{R}\tau + 2\bar{X}\tau + 2\bar{R}\tau + 2\bar{T}\tau)$$

Având aceste rezultate RAID 5 poate fi modelat ca o coadă $M/G/1$.

În continuare, vom compara performanțele RAID 1 și RAID 5 în cazul unor operații I/O de dimensiuni mari, cu o medie de 1M *bytes*.

În primul rând, vom compara RAID 1 și RAID 5 cu aceeași capacitate, RAID 1 având $2n$ discuri, în vreme ce RAID 5 conține doar $n+1$ discuri (cu $n=32$). Rezultatele sunt prezentate în Figura 3-36 pentru diferite valori ale lui p_f și se observă că RAID 1 se comportă mai bine în toate cazurile decât RAID 5. Figura 3-37 ilustrează situația în care cele două arhitecturi au același număr de discuri. Se observă din Figura 3-37a că RAID 1 depășește RAID 5 când majoritatea cererilor sunt citiri. Acest lucru se întâmplă deoarece RAID 1 servește două cereri de citiri simultan, în vreme ce RAID 5 secvențial. Deși timpul de transfer la RAID 5 este doar jumătate din cel al RAID 1, întârzierea cauzată de servirea concurentă a mai multor operații de citire la RAID 1 echilibrează dezavantajul unui timp de transfer mai mare. Când majoritatea cererilor sunt scrieri, RAID 1 își pierde avantajul de a efectua în paralel citiri. Prin urmare RAID 5 are performanțe mai bune datorită timpului de transfer mai scăzut. După cum se arată în Figura 3-37b, performanțele RAID 5 scad rapid o dată cu scăderea probabilității p_f datorită timpului de servire adițional necesar procedurii *read-update* pentru scrierea parțială a unei secțiuni. Când 30% din scrieri sunt parțiale într-o secțiune, RAID 1 este echivalent cu RAID 5 la încărcări mari; iar când jumătate din scrieri sunt parțiale RAID 1 depășește RAID 5. În Figura 3-38 și Figura 3-39 este punctată sarcina maximă suportată de RAID 1 și RAID 5 în funcție de dimensiunea transferului sau de numărul de discuri. În general, RAID 1 se comportă mai bine decât RAID 5. Excepția apare când se folosește același număr de discuri, majoritatea cererilor sunt scrieri și dimensiunea transferului per cerere este mare (mai mare decât 2MB). În acest caz, rezultatele RAID 5 sunt mai bune decât RAID 1.

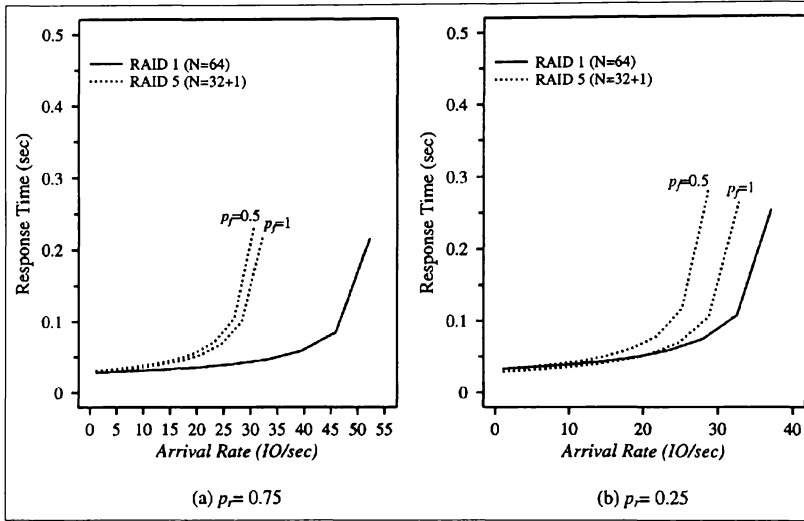


Figura 3-36

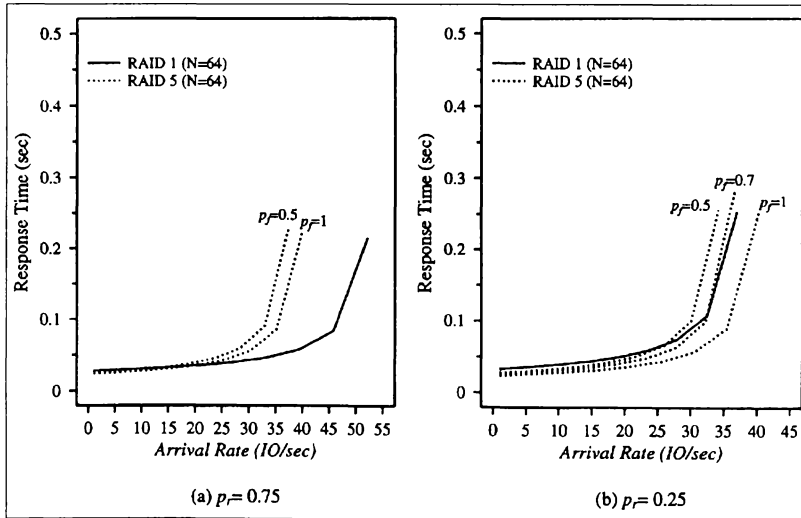


Figura 3-37

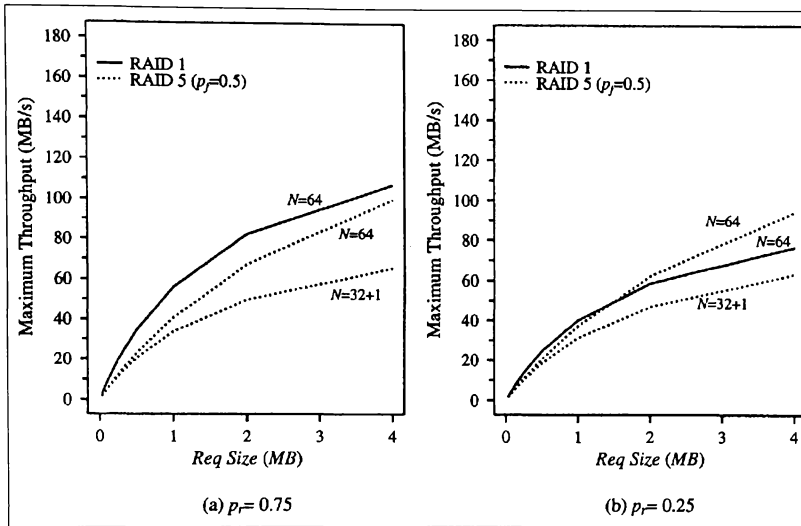


Figura 3-38

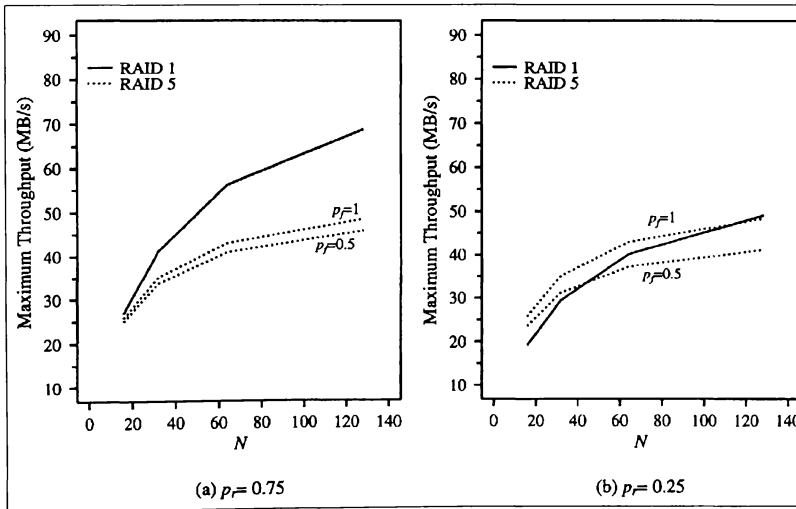


Figura 3-39

În discuțiile de mai sus am presupus ca majoritatea cererilor I/O ale RAID 5 efectuează operații I/O *full stripe* ($p_f > 0.5$). Această condiție necesită ca sistemul de fișiere sau managerul de *cache* să cunoască lățimea secțiunii matricei de discuri. Dacă numărul de discuri crește sau

scade, sistemul de fișiere sau managerul de *cache* trebuie să se adapteze acestei schimbări pentru a păstra un nivel ridicat al performanțelor.

3.8.1.3 Concluzii

În modul obișnuit de operare, o tactică ce repartizează o cerere de citire de pe disc unei cozi de lungime minimă, pentru orice tip posibil de date, furnizează cel mai scăzut timp mediu de răspuns dintre toate variantele RAID 1 în medii în care aplicațiile necesită cantități scăzute de date. În aplicații în care se folosesc cantități mari de date, politica cozii de lungime minimă trebuie corelată cu operații cum ar fi *mirrored declustering* sau mai recenta operație *group-rotate declustering* („deciorchinizare prin rotație de grup“), care permite ca citiri independente să se execute independent.

Mai sus amintitele arhitecturi RAID 1 depășesc semnificativ RAID 5 când aplicațiile generează cereri I/O pentru cantități scăzute de date. Acest lucru este valabil pentru cazul în care ambele arhitecturi au același număr de discuri. În cazul aplicațiilor care generează cereri I/O pentru cantități mari de date, rezultatele nu sunt la fel de evidente. RAID 5 prezintă performanțe mai ridicate când majoritatea cererilor sunt pentru cantități foarte mari de date, când majoritatea cererilor sunt scrieri și când majoritatea scrierilor se efectuează pe o secțiune întregă.

Studiile s-au efectuat prin combinarea de experimente și analize. De exemplu, arhitectura RAID 5 a fost analizată aproximativ, dar corect, prin descompunere, în care fiecare disc este modelat după sistemul de cozi bazate pe priorități.

Într-o lucrare recentă, Lee și Katz prezintă un model analitic al matricelor de discuri. În aceste lucrări, operațiile de citire și scriere sunt tratate în același mod la evaluarea performanțelor sistemelor RAID 5. Prin urmare, costul ridicat obținut prin scrieri parțiale în secțiuni sub RAID 5 este ignorat.

Performanțele I/O pot fi, de asemenea, îmbunătățite prin introducerea unui „*cache* de disc“ (*disk cache*). Eficacitatea *cache*-ului depinde de modul acceselor I/O, precum și de dimensiunea *cache*-ului. Pentru aplicații în care accesele la disc denotă un factor de localizare ridicat, *disk caching*-ul poate satisface majoritatea cererilor de citire și reduce traficul I/O cu discul. Pentru operațiile de transfer, maparea este mai puțin eficace deoarece cererile I/O accesează în mod aleator discul și nu este deloc practic să se mapeze întreaga baza de date. În orice caz, *cache*-urile de disc pot fi combinate cu matricele de discuri, iar studiul de față rămâne valabil și în cazul acestor sisteme.

Principalele aspecte abordate în prezentul capitol sunt studiile asupra problemelor de proiectare, dezvoltarea unor modele matematice și examinarea performanțelor diverselor arhitecturi ale matricelor de discuri, atât pentru operații I/O în cantitate redusă, cât și ridicată.

În particular, am studiat arhitecturile de matrice de discuri existente și am propus, în următorul capitol, o nouă arhitectură, RAID level 7. Se prezintă, de asemenea, strategii de proiectare potrivite pentru diverse arhitecturi de matrice de discuri. În cele din urmă, din studiile făcute s-a observat că principalul impact al performanțelor RAID 5 este costul ridicat al *partial stripe write*. Cercetătorii au căutat diverse modalități de a reduce acest cost. Cel mai atractiv mod este LSFS (*Log-Structured File System*), în care multe scrieri „mici” sunt cumulate în *cache* și apoi se face o scriere mare completă. Un alt mod este cel numit *floating parity*, în care blocurile de paritate sunt împărțite în cilindri, fiecare conținând o pistă de rezervă. Un nou bloc de paritate, în loc să fie scris la locul său, este scris în locul celui mai apropiat bloc nealocat următor vechiului bloc de paritate. Costul plătit este următorul: controlerul trebuie să păstreze o hartă a locațiilor tuturor blocurilor de paritate.

3.9 Concluzii

Dintre soluțiile de creștere a fiabilității unităților de discuri, prezentate anterior, în prezentul capitol, abordez, în detaliu, matricele de discuri.

Analiza performanțelor tuturor nivelurilor acceptate de RAID, împreună cu diverse strategii de accesare a discurilor, au în vedere propunerea unei noi structuri RAID, pe parcursul următorului capitol:

În această parte, teza cuprinde următoarele contribuții:

1. o originală analiză a performanțelor structurilor RAID 1 și RAID 5;
2. prezentarea unor modele utile în evaluarea performanței unor structuri RAID.

4 PROPUNERE RAID LEVEL 7

Se cunosc propuneri de utilizare a unor discuri suplimentare active, dar nefolosite, în scopul oferirii unor spații de reconstrucție pentru datele stocate pe discurile căzute.

Ținând cont de viitorul acestor soluții, propun un nou nivel de RAID, level 7, care presupune existența unor spații de rezervă, distribuite pe toate discurile din matrice, spații care să fie utilizate în cazul căderii oricărui disc. În funcție de dimensiunea acestor spații, RAID 7 poate suporta căderi multiple succesive, fără o degradare de performanță excesivă. Evident, soluția nu este ieftină, dar din punct de vedere a fiabilității, este optimă, iar din punct de vedere al performanței este echivalentă cu RAID 5.

RAID 7 este singurul tip de RAID care poate reveni, fără intervenție umană, după o cădere de disc, în scurt timp, la operare în regim protejat, cu performanțe apropiate de cele anterioare căderii.

Utilizare unor discuri de rezerva în matrice de tip 1–4, nu se justifică nici din punct de vedere tehnic, dar nici economic. Din acest motiv în prezenta lucrare mă voi ocupa doar de un RAID 7 provenit din dezvoltarea unei matrice RAID 5.

Este evident faptul că cu cât timpul de reparare este mai mic cu atât crește timpul până la pierderea datelor, deoarece aria revine la o stare protejată mai repede. Un mod de a „repara“ un disc căzut este de a avea rezerve de discuri, astfel încât ele pot înlocui pe cele căzute în mod automat. Timpul de reparare este practic egal cu cel de refacere a matricei.

S-a arătat că menținerea unei rezerve active neutilizată în matrice este inefficientă. O soluție mai bună, numită „rezervare distribuită“, este de a lăsa spațiu suficient pe fiecare disc din matrice astfel încât datele refăcute de pe un disc căzut pot fi stocate pe acel spațiu. Aceasta permite discurilor să fie folosite în permanență, mai puține date să se piardă la o cădere de disc și, în timpul refacerii unui disc, scrierea datelor refăcute să se facă pe mai multe discuri din matrice. Un exemplu simplu este:

| D1 | D2 | D3 | D4 | D5 |
|----|----|----|----|----|
| A | A | A | A | S |
| B | B | B | S | B |
| C | C | S | C | C |
| D | S | D | D | D |
| S | E | E | E | E |

Figura 4-1

unde fiecare domeniu notat cu **S** reprezintă domeniul de rezervă. Dacă, să zicem, cade discul 1, atunci domeniul A este refăcut pe domeniul de rezervă al discului 5, domeniul B pe discul 4 s.a.m.d. Aria rezultată, având patru discuri este, încă o matrice cu toleranță la defecte singulare, deși fără spațiu de rezervă.

4.1 UPGD și rezervarea distribuită

Dacă rezervarea distribuită se combină cu o matrice UPGD, atunci rezultă o matrice cu distribuție uniformă a grupurilor de paritate cu rezerve distribuite (UPGD/DS) (*Uniform Parity Group Distribution with Distributed Sparing*), ca în Figura 4-2. De notat că spațiile goale dintre domenii sunt puse doar cu scopul unei înțelegeri mai exacte, aceste spații nefiind fizic pe disc. Deoarece avem o UPGD, vor fi necesare exact $\lambda=2$ accese de citire pe fiecare din cele $v-1=6$ discuri din matrice pentru a reface cele $r=6$ domenii de pe discul căzut și domeniile refăcute trebuie să fie scrise câte unul pe fiecare disc rămas, pentru un total de exact $\lambda+1=3$ accese fiecare. Astfel, aceste UPGD/DS-uri dau o distribuție uniformă a încărcării suplimentare la căderea unui disc, atât în funcționarea în mod degradat cât și în timpul refacerii datele pe spațiul de rezervă.

| D1 | D2 | D3 | D4 | D5 | D6 | D7 | |
|----|----|----|----|----|----|----|--------------------|
| A | A | A | | | | | Grup de paritate A |
| B | | | B | B | | | Grup de paritate B |
| C | | | | | C | C | Grup de paritate C |
| D | D | | D | | | | Grup de paritate C |
| E | | E | | | E | | Grup de paritate D |
| F | | | | F | | F | Grup de paritate E |
| | G | | | G | G | | Grup de paritate F |
| | | H | H | H | | | Grup de paritate G |
| | I | I | | | | I | Grup de paritate I |
| | | | J | | J | J | Grup de paritate J |
| | K | | K | | K | | Grup de paritate K |
| | | L | L | | | L | Grup de paritate L |
| | M | | | M | | M | Grup de paritate M |
| | | N | | N | N | | Grup de paritate N |
| □ | □ | □ | □ | □ | □ | □ | Blocuri de rezerva |

Figura 4-2

Referindu-ne la Figura 4-2, după orice cădere de disc, domeniile de pe discul căzut sunt refăcute pe spațiul de rezervă, vor fi $v'=v-1=6$ discuri cu $r'=r+1=7$ domenii în fiecare grup. În scopul obținerii unui BIBD, astfel încât cea de-a doua cădere să genereze din nou o încărcare echilibrată, atunci $v'*r'=42$ trebuie să fie egal cu $b'*k'=42$, ceea ce se și întâmplă, și $\lambda'=(r'*(k'-1))/(v'-1)=2.8$ trebuie să fie întreg, și nu este. Astfel, orice disc ar cădea, sau oricum s-ar reface datele de pe discul căzut, rezultatul nu poate fi o matrice UPGD. O a doua cădere va cauza o încărcare suplimentară, ce nu va fi distribuită în mod egal pe discurile rămase.

4.2 UPGD/DS și căderile multiple

În continuare, se va descrie o metodă pentru matrice de discuri, astfel încât, plecându-se de la o matrice UPGD/DS, după o succesiune de căderi ale unor discuri (fiecare urmată de refacere) să rezulte tot o secvență de matrice UPGD/DS, fiecare bazată pe discurile rămase valide în matrice. Scopul este de a plasa grupurile de paritate și spațiul liber în matrice astfel încât de fiecare dată când cade un disc, aria să fie UPGD/DS înainte de cădere, dar și după ce datele sunt refăcute de pe discul căzut și scrise pe spațiul liber distribuit (sau doar o UPGD dacă nu mai este spațiu liber). Fiind realizat acest obiectiv, sarcina suplimentară în modul degradat sau/și refacerea domeniilor de date pe spațiul de rezervă vor fi totdeauna uniform distribuite pe discurile ramase funcționale.

În continuare se vor folosi notațiile:

- N -numărul de discuri din matrice;
- K -mărimea grupului de paritate;
- b -numărul total de grupuri de paritate;
- (x,y) -BIBD (*Balanced Complete Block Design*) format prin enumerarea tuturor combinațiilor de x luate câte y ;
- $E(x,y)$ -numărul tuturor acestor combinații, care este:

$$x!/((x-y)!*y!)$$

- f -numărul de căderi la care pot fi refăcute și remapate datele înainte ca să se înlocuiască vreun disc;
- aria se află în starea i dacă datele au fost refăcute și remapate în spațiul liber după i căderi. Cea de-a $i+1$ -a cădere poate sau nu să fi apărut încă: ($0 \leq i \leq f$);
- v_i -numărul de discuri ce conțin date/paritate din starea i . Unul din aceste discuri poate să fi căzut, dar datele/paritatea nu au fost încă remapate. Avem: $v_0=N$; $v_{i+1}=v_i-1$; $v_i=N-i$;
- r_i -numărul de domenii de date/paritate de pe fiecare disc din starea i , și are valoarea: $r_i=(b*K)/v_i$;
- s_i -numărul de domenii de rezervă din fiecare disc, când aria este în starea i .

4.3 O analiză a algoritmului Ng-Mattson de creare a UPGD/DS

Numărul total de grupuri de paritate într-o matrice b nu se schimbă pe măsură ce domeniile dintr-un disc căzut se mapează în domeniile de rezervă.

Un mod simplu de aranjare pentru matrice în scopul obținerii unei UPGD în starea i este de a avea grupurile de paritate posibile de mărime K din v_i discuri, replicate de un număr de ori. Pentru a obține acest lucru, trebuie ca $E(v_i, K)$ să-l dividă pe b . Mai mult, trecând de la o matrice cu v_i discuri la una v_{i-1} , condiția de a menține proprietatea de simetrie a distribuției uniforme este să existe v_i-K copii ale fiecărui grup de paritate și să mapăm câte un domeniu de pe discul căzut, ce aparține grupului, pe fiecare din cele v_i-K discuri neimplicate în grupul de paritate. Acest lucru se va demonstra mai târziu. Deci, trebuie ca v_i-K să dividă și el pe b . Prin urmare, punctul de plecare este selectarea numărului b , astfel încât să fie cel mai mic multiplu comun al numerelor $(v_i-K)*E(v_i, K)$, pentru toți $i=0..f-1$.

Pașii ce urmează și dovada că funcționează se bazează pe presupunerea ca pentru toți $i=0..f-1$, înainte de cea de-a $i+1$ -a cădere, aria este aranjată ca o UPGD/DS bazată pe o structură (v_i, K) -BIBD replicată de un număr de $(v_i-K)*Q_i$ ori, unde Q_i este un întreg. Această presupunere este adevărată pentru $i=0$, prin definiție datorită pasului 3 de mai jos și se va arăta că va fi adevărat, prin inducție, și pentru $i=1..f-1$, conform pasului 4. Algoritmul pentru crearea și funcționarea unei matrice UPGD/DS are următorii pași:

- selectarea numărului $f \leq N-K$, care indică numărul de căderi după care se pot reface datele pe spațiile de rezervă. Acesta este numărul de căderi până la care este nu este necesară intervenția umană, pentru schimbarea discului sau a discurilor;
- determinarea numărului de grupuri de paritate din matrice, ca fiind c.m.m.m.c. al numerelor $(v_i-K)*E(v_i, K)$, pentru $i=0..f-1$. Aceasta înseamnă că, pentru fiecare valoare a lui i , există un întreg Q_i astfel încât $b=Q_i*(v_i-K)*E(v_i, K)$
- Formăm o matrice inițială UPGD/DS prin:
 - ◆ enumerarea tuturor combinațiilor de N discuri luate câte K și fiecare combinație formează un grup de paritate;

- ◆ replicarea fiecărui grup de paritate format mai sus de Q_0 ori să obținem în total b grupuri de paritate, cu $r=(b*K)/N$ domenii atribuite fiecăruia din cele N discuri;
- ◆ adăugând $s_0=(b*K*f)/(N*(N-F))$ domenii de rezervă pentru rezervarea distribuită;

Când aria se afla în starea i , presupunem că este aranjată ca o UPGD/DS bazată pe o structură (v_i, K) -BCBD replicată de $Q_i * (v_i - K)$ ori și cade cel de-al $i+1$ -lea disc. Aria va funcționa în mod degradat și se va reconstrui informația pe spațiul de rezervă disponibil conform astfel:

- imediat ce a căzut al $i+1$ -lea disc, există $v_i - 1$ discuri bune. Pentru ușurință, vom nota discurile: $D_1, \dots, D_{k, \dots}, D_{v_i - 1}$, și D_{v_i} , ultimul fiind cel care s-a defectat; pentru o structură (v_i, K) -BCBD grupăm conceptual toate grupurile de paritate care nu folosesc discul D_{v_i} într-un set numit A_i . Vor fi exact $Q_i * (v_i - K) * E(v_i - 1, K)$, adică de $Q_i * (v_i - K)$ ori numărul de moduri în care se pot selecta K discuri din $D_1, D_2, \dots, D_{v_i - 1}$ -similar, fie B_i mulțimea tuturor grupurilor de paritate care folosesc discul D_{v_i} . Vor fi exact $r_i = (b * K) / v_i = Q_i * (v_i - K) * E(v_i - 1, K - 1)$ adică $Q_i * (v_i - K)$ copii ale aceluiași grup de paritate unic care folosește discul D_{v_i} sau, altfel spus, de $Q_i * (v_i - K)$ ori numărul de moduri în care se pot selecta $K - 1$ discuri din D_1, \dots, D_{v_i} .
- pentru fiecare grup de paritate unic din B_i , se iau toate cele $Q_i * (v_i - K)$ copii ale sale și se reconstruiesc Q_i domenii de pe discul D_{v_i} și se remapează pe unul din cele $v_i - K$ discuri nefolosite de acest grup, apoi se reconstruiesc alte Q_i domenii de pe discul D_{v_i} pe un alt disc nefolosit de acest grup de paritate s.a.m.d., până când toate cele $Q_i * (v_i - K)$ domenii de pe discul D_{v_i} care aparțin copiilor grupului de paritate unicat care au fost reconstruite și remapate pe cele $v_i - K$ discuri nefolosite de acest grup de paritate, punând Q_i domenii pe fiecare disc. Remaparea domeniilor de pe discul căzut doar pe cele ce nu folosesc acest grup de paritate oferă toleranță la defecte deoarece nici un disc nu va conține două domenii din același grup de paritate.

Deoarece există un număr de $E(v_i - 1, K - 1)$ grupuri unice de paritate în mulțimea B_i , pasul anterior se repetă de tot atâtea ori. Numărul de domenii care sunt remapate pe oricare disc D_j este de Q_i ori numărul grupurilor unicate care conțin pe D_{v_i} , dar nu conțin pe D_j , adică $Q_i * E(v_i - 2, K - 1)$, care se poate arăta că este totuna cu $r_i / (v_i - 1)$, numărul

de domenii r_i din D_{v_i} distribuit în mod egal printre cele v_i-1 discuri rămase.

Propoziție

Când aria se afla în starea $i < f$, dacă este o UPGD/DS bazată pe o structura (v_i, K) -BCBD replicat de $Q_i^*(v_i-K)$ ori, atunci pasul va produce o nouă structură constând din $Q_{i+1}^*(v_{i+1}-K)$ replicări ale structurii (v_i, K) -BCBD cu $v_{i+1} = v_i - 1$ discuri.

Demonstrație

Pentru orice K discuri G_1, \dots, G_K , avem:

- există un singur grup de paritate unicat care conține domenii de pe discurile G_2, \dots, G_K și D_{v_i} , și care nu conține pe G_1 . Acest grup unicat este replicat de $Q_i^*(v_i-K)$ ori. Fie P_1 această mulțime de grupuri de paritate și ele aparțin mulțimii B_i . Q_i din domeniile ei vor fi remapate de pe discul D_{v_i} pe G_1 ;
- fie P_2 mulțimea de $Q_i^*(v_i-K)$ grupuri de paritate din mulțimea B_i care nu folosesc discul G_2 , dar folosesc discurile G_1, G_3, \dots, G_K , și D_{v_i} . Q_i din domeniile sale vor fi remapate de pe discul D_{v_i} pe G_2 ;
-; și, în final fie P_k mulțimea de $Q_i^*(v_i-K)$ grupuri de paritate din B_i care nu folosesc discul G_K , dar le folosesc G_1, \dots, G_{K-1} , și D_{v_i} . Q_i din domeniile sale au fost remapate de pe D_{v_i} pe G_k .

Astfel, Q_i^*K grupuri de paritate din setul B_i au fost remapate pentru a forma Q_i^*K grupuri care folosesc discurile G_1, \dots, G_K . În plus, există $Q_i^*(v_i-K)$ grupuri din mulțimea A_i care, de asemenea, folosesc aceste discuri. Deci, există un total de Q_i^*K grupuri care folosesc discurile G_1, \dots, G_K în noua matrice după remapare. Deoarece acest lucru este adevărat pentru orice K discuri, este adevărat pentru oricare dintre cele $E(v_i-1, K)$ mulțimi unice de K discuri.

Deci, rezultatul remapării îl reprezintă cele $Q_i^*v_i$ copii ale structurii (v_i-1, K) -BCBD. Mai mult, deoarece avem:

$$\begin{aligned} Q_i^*v_i &= b / [E(v_i, K)^*(v_i-K)] = (\text{din definiția lui } Q_i) \\ &= b / \{(v_i-K)^* v_i! / [(v_i-K)! * K!]\} = b / \{(v_i-1)! / [(v_i-1-K)! * K!]\} = \\ &= b / E(v_{i+1}, K) = Q_{i+1}^*(v_{i+1}-K) \quad (\text{tot din definiție}), \end{aligned}$$

noua matrice este tot una UPGD/DS bazată pe o structură (v_{i+1}, K) -BCBD replicată de $Q_i^*(v_{i+1}, K)$ ori.

Deoarece algoritmul începe cu o matrice UPGD/DS în starea 0, demonstrația de mai sus a arătat, prin inducție, că în urma unui număr de până la f căderi succesive discuri se produce o serie de noi matrice UPGD/DS. Pentru toți $i=0, 1, \dots, f-2$, pasul 4 produce o matrice după cea de-a $(i+1)$ -a cădere. După f căderi, nemaifiind spațiul liber pe discuri, va rezulta doar o matrice UPGD. Deci, algoritmul conduce la o matrice cu o funcționare foarte bună în mod degradat și o performanță bună a reconstruirii.

După cea de-a $(i+1)$ -a cădere și după ce datele au fost remapate, cantitatea de spațiu liber rămas în fiecare disc este redus de domenii remapate pe fiecare disc, adică r_i/v_{i-1} . Se poate arăta că după cea de-a f -a cădere și după f remapări, cantitatea totală de spațiu liber folosit de fiecare disc este s_0 , adică spațiul liber cu care începe fiecare disc.

Exemplu

Algoritmul este ilustrat în Figura 4-3 pentru o matrice UPGD/DS cu $N=5$, mărimea grupului de paritate $K=3$, și rezervare pentru $f=2$ căderi. Pentru această matrice $b=20$, c.m.m.m.c al numerelor

$(v_0-K)*E(v_0, K)=20$ și $(v_1-K)*E(v_1, K)=4$. Aceasta duce la $Q_0=1$ și $Q_1=5$.

În Figura 4-4, grupurile de paritate și domeniile de rezervă sunt plasate în matrice astfel încât este formată o matrice UPGD/DS din toate combinațiile de $N=5$ discuri luate câte $K=3$: 1-2-3 (grupul A), 1-2-4 (grupul B) s.a.m.d. până ce $K=3$ domenii din grupul de paritate J sunt atribuite la al zecelea triplet 3-4-5.

Fiecare grup de paritate format mai sus este replicat de $Q_0^*(v_0-K)=2$ ori, prima replicare fiind grupurile A, B, ..., J, iar cea de-a doua fiind K, L, ..., T. Rezultă astfel 12 domenii atribuite la fiecare din cele 5 discuri, $s_0=(b*K*f)/(N*(N-f))=8$ domenii de rezervă care se adaugă la fiecare disc.

| D1 | D2 | D3 | D4 | D5 |
|----|----|----|----|----|
| A | A | A | | |
| B | | B | B | |
| C | C | | | C |
| D | D | D | | |
| E | | E | | E |
| F | | | F | F |
| | G | G | G | |
| | H | H | | H |
| | I | | I | I |
| | | J | J | J |
| K | K | K | | |
| L | L | | L | |
| M | M | | | M |
| N | | N | N | |
| O | | O | | O |
| P | | | P | P |
| | Q | Q | Q | |
| | R | R | | R |
| | S | | S | S |
| | | T | T | T |
| □ | □ | □ | □ | □ |
| □ | □ | □ | □ | □ |
| □ | □ | □ | □ | □ |
| □ | □ | □ | □ | □ |
| □ | □ | □ | □ | □ |
| □ | □ | □ | □ | □ |
| □ | □ | □ | □ | □ |
| □ | □ | □ | □ | □ |
| □ | □ | □ | □ | □ |

Figura 4-3

| D1 | D2 | D3 | D4 | D5 |
|----|----|----|----|----|
| A | A | A | | |
| B | | B | B | |
| D | | D | D | |
| | G | G | G | |
| K | K | K | | |
| L | L | | L | |
| N | | N | N | |
| | Q | Q | Q | |
| C | C | C | | C |
| M | M | | M | M |
| E | E | E | | E |
| O | | O | O | O |
| F | F | | F | F |
| P | | P | P | P |
| H | H | H | | H |
| | R | R | R | R |
| I | I | | I | I |
| | S | S | S | S |
| J | | J | J | J |
| | T | T | T | T |
| □ | □ | □ | □ | |
| □ | □ | □ | □ | |
| □ | □ | □ | □ | |
| □ | □ | □ | □ | |
| □ | □ | □ | □ | |

Figura 4-4

Notăm discurile D1, D2, D3, D4, D5 (D5 fiind cel care se defectează).

În Figura 4-4 toate grupurile de paritate care nu folosesc discul D5 sunt puse în mulțimea A0.

Vor fi exact $Q_0^*(v_0-K)*E(v_0-1,K)=8$, adică de $Q_0^*(v_0-K)=2$ ori numărul de moduri în care se pot selecta $K=3$ discuri din $v_0-1=4$ discuri. Aceste grupuri sunt A, B, D, G și K, L, N, Q puse în topul figurii 4b. Toate grupurile care folosesc discul D5 sunt puse în mulțimea B0, unde vor fi exact r_0 elemente, adică $Q_0^*(v_0-K)=2$ copii ale celor $E(v_0-1,K-1)$ grupuri unice care folosesc discul D5. Aceste grupuri sunt C, M, E, O, F, P, H, R, L, S, J și T puse în josul figurii (grupurile C și M reprezintă 2 copii ale aceluiași unicat). Pentru un unicat (de exemplu C) în mulțimea B0 se iau toate cele $Q_0^*(v_0-K)=2$ copii (C și M) și se reconstruiesc $Q_0=1$ domeniu de pe D5 pe fiecare din cele $v_0-K=2$ discuri nefolosite de aceste grupuri (de exemplu se reconstruiește domeniul C de pe D5 pe D3 și domeniul M de pe D5 pe D4). Această procedură se repetă până ce toate domeniile de pe discul D5 au fost remapate pe discurile D1..4. Amplasarea rezultantă a grupurilor de paritate este o matrice UPGD/DS cu $v_1=4$, $K=3$, $b=20$, $r_1=15$ și $s_1=5$. Ea este bazată pe o matrice cu o structură (3, 4)-BCBD replicată de $Q_1^*(v_1-K)=5$ ori.

La o a doua cădere, procedura se repetă cu aceleași valori ale lui b și K , dar cu $v_1=4$ și $Q_1=5$. Amplasarea grupurilor de paritate este dată în Figura 4-5 și este o matrice UPGD/DS cu $v_2=3$, $K=3$, $b=20$, $r_2=20$ și $s_2=0$. Următoarea cădere va forța sistemul să funcționeze în mod degradat până ce unul din discurile căzute este înlocuit și îi sunt restaurate datele.

| | | | |
|----|----|----|----|
| D1 | D2 | D3 | D4 |
| A | A | A | |
| K | K | K | |
| C | C | C | |
| E | E | E | |
| H | H | H | |
| B | B | B | B |
| M | M | M | M |
| F | F | F | F |
| I | I | I | I |
| D | D | D | D |
| N | N | N | N |
| O | O | O | O |
| P | P | P | P |
| J | J | J | J |
| G | G | G | G |
| Q | Q | Q | Q |
| R | R | R | R |
| S | S | S | S |
| T | T | T | T |

Figura 4-5

4.4 Optimizarea matricelor de discuri UPGD/DS compacte

Algoritmul descris anterior va funcționa pentru orice triplet (n, k, f)

unde:

- n reprezintă numărul inițial de discuri în matrice;
- k reprezintă dimensiunea grupului de paritate;
- f reprezintă numărul de căderi de disc pentru care s-a prevăzut spațiu de rezervă pe fiecare dintre discuri.

Tabelul care urmează conține numărul grupurilor de paritate b pentru diferite structuri UPGD/DS, caracterizate fiecare prin tripleta (n, k, f) , numărul de domenii de date r_0 atribuit fiecărui disc din matrice, precum și numărul domeniilor de rezervă s_0 cu care a fost înzestrat fiecare disc pentru ca aria să poată susține cele f căderi pentru care a fost proiectată. Menționăm că cei trei parametri care interesează, adică b , r_0 și s_0 au fost generați de algoritmul prezentat anterior:

| n | k | f | b | r_0 | s_0 |
|----|----|---|--------|--------|-------|
| 6 | 3 | 2 | 60 | 30 | 15 |
| 8 | 4 | 2 | 840 | 420 | 140 |
| 12 | 5 | 2 | 5544 | 2310 | 462 |
| 16 | 5 | 3 | 720720 | 225225 | 51975 |
| 16 | 10 | 3 | 240240 | 150150 | 34650 |

Figura 4-6

Se observă că în toate cazurile, în faza inițială, fiecare matrice este aranjată într-o structura de tip UPGD/DS bazată pe enumerarea tuturor combinațiilor de n discuri luate câte k , replicate de un anumit număr de ori (fiecare combinație de discuri astfel enumerată formând un grup de paritate):

- $E(6,3)$ replicate de 3 ori;
- $E(8,4)$ replicate de 12 ori;
- $E(12,5)$ replicate de 7 ori;
- $E(16,6)$ replicate de 165 ori;
- $E(16,10)$ replicate de 30 ori.

De asemenea, domeniile de rezervă care trebuie alocate fiecărui disc din matrice este multiplicat de atâtea ori de câte ori sunt replicate cele $E(n, k)$ grupuri de paritate din matrice.

Acest număr mare de grupuri de paritate cerut pentru ca algoritmul *Ng-Mattson* să funcționeze implică o dimensiune mare a tabelii care va administra aria, deci, capacitatea memoriei *cache* din controler trebuie să fie mare, ceea ce va ridica în mod substanțial costul sistemului.

O cale de a reduce acest cost este aceea de a încerca o reducere a numărului de grupuri de paritate din matrice. Mai clar spus, având n discuri inițial în matrice, iar dimensiunea grupului de paritate este egală cu k , vom căuta o soluție care să implice în starea de plecare (înaintea apariției primului disc defect) **un număr mai mic de replici** ale tuturor combinațiilor de n discuri luate câte k , numărul acestor combinații dând numărul de **modele** de grupuri de paritate din matrice. Practic, se caută minimizarea numărului b , ceea ce automat implică micșorarea lui r_0 și a lui s_0 .

În continuare, este prezentat un exemplu în care valoarea inițială a lui b este mai mică decât cea cerută de pasul 2 al algoritmului de la capitolul precedent.

Astfel, pentru o matrice UPGD/DS cu $n=6$ discuri, dimensiunea grupului de paritate $k=3$ și spațiu de rezervă pentru $f=2$ defecțiuni, valoarea b ce s-a ales, $b=20$, care reprezintă numărul grupurilor de paritate din matrice, este mai mic decât cel dat de condiția:

$$(v_0-k) * E (v_0-k) = 60.$$

Totuși, și în acest caz se poate obține o matrice UPGD/DS repetată la care îi vom atribui atributul de **compactă**. În subcapitolul următor, prezentăm structura matricei înaintea defectării primului disc.

După ce a avut loc defectarea primului disc (fie acesta D_6), informația care se găsea pe acesta este reconstruită în spațiile de rezervă ale discurilor funcționale rămase în matrice într-o asemenea manieră încât aria rezultată este o matrice UPGD/DS, având $v_1=5$ discuri, $b=20$, $r_1=12$ domenii de date per disc, $k=3$ și $s_1=3$ domenii de rezervă pe fiecare dintre cele 5 discuri rămase.

După ce a avut loc defecțiunea celui de-al doilea disc (fie acesta D_5), informația ce se afla pe acest disc este reconstruită în spațiul de rezervă al discurilor rămase în matrice (în număr de 4), devenind doar o matrice UPGD (nu mai există domenii de rezervă pe discuri) cu $v_2=4$ discuri, $b=20$, $r_2=15$, $k=3$.

Acest exemplu demonstrează că matrice UPGD/DS replicate mai mici decât cele construite cu algoritmul *Ng-Mattson* există, dar găsirea unei căi de a reconstrui domeniile de pe discul defect pe celelalte discuri din matrice în scopul de a produce o nouă matrice UPGD este dificilă.

Găsind un b mai mic, din algoritmul *Ng-Mattson* putem folosi fără nici o modificare pașii 1 și 3. Pasul 2, care decide numărul de replici se va

modifica. În fine, pasul 4 și cel mai important, care realizează practic remaparea domeniilor discului defect în domeniile de rezervă ale discurilor funcționale ramase în matrice, se va adapta în așa manieră încât după remaparea despre care am vorbit, noua matrice să-și păstreze structura de tip UPGD/DS (sau UPGD în situația în care nu mai există spațiu de rezervă pe discuri).

Este necesară observația că numărul de domenii de rezervă atribuite inițial fiecărui disc din matrice, s_0 , va fi și el corespunzător mai mic:

$s_0 = (b * k * f) / (n * (n - f))$, (b fiind mai mic, s_0 va rezulta de asemenea mai mic).

4.5 Algoritm optimizat de creare UPGD/DS

Cercetările întreprinse de autor, având ca scop găsirea unei soluții de a micșora numărul b al grupurilor de paritate din matrice, au fost încununuate de succes: am reușit găsirea unei metode prin care, pentru cazul: $\{n=2*k, f=2\}$ obținem un număr semnificativ mai mic al grupurilor de paritate.

În acest caz, dimensiunea grupului de paritate trebuie să fie jumătate din numărul inițial de discuri ($k=n/2$), și sunt admise două defecțiuni succesive de discuri în matrice ($f=2$).

Lucrurile stau în felul următor: înainte de defectarea primului disc, aria este aranjată într-o structură UPGD/DS bazată pe enumerarea tuturor combinațiilor de $n=2*k$ discuri luate de k ori (fiecare din această combinație reprezentând un model unic de grup de paritate), replicate de un număr de x ori. Despre acest număr x vom vorbi și îl vom calcula puțin mai târziu, rezumându-ne acum doar la a spune că rezultă mai mic decât acel număr $Q_0 * (v_0 - k)$, care reprezenta numărul inițial de replici ale fiecărui model unic de grup de paritate, număr cerut de pasul 2 al algoritmului descris în cadrul capitolului anterior.

Revenind la ideea începută mai sus, după ce primul disc din matrice se defectează, din numărul total de grupuri de paritate, care este egal cu $b = x * E(n, k)$, vom avea o parte din ele care utilizează discul defect, iar cealaltă parte (restul până la numărul b) care nu utilizează discul defect:

- $x * E(n-1, k-1)$ grupuri de paritate utilizează discul defect;
- $x * E(n-1, k)$ grupuri de paritate nu utilizează discul defect.

Deoarece în cazul nostru avem $n=2*k$, obținem $E(n-1, k-1) =$

$$= E(2*k-1, k-1) = E(2*k-1, 2*k-1-(k-1)) = E(2*k-1, k) = E(n-1, k), \text{ deci:}$$

$$E(n-1, k-1) = E(n-1, k) \quad | \text{înmulțim cu } x$$

și rezultă în final:

$$x * E(n-1, k-1) = x * E(n-1, k) = b/2,$$

ceea ce înseamnă că **numărul grupurilor de paritate care nu utilizează discul defect este egal cu numărul grupurilor de paritate care utilizează discul defect**, lucru foarte important, care va sta la baza algoritmului de construire a matricei UPGD/DS compacte.

Să notăm cu **A** mulțimea tuturor grupurilor de paritate care utilizează discul defect, iar cu **B** mulțimea tuturor grupurilor de paritate care nu utilizează discul defect. Pe baza demonstrației de mai sus, rezultă:

$$\text{card } A = \text{card } B = b/2.$$

Pentru ca, după remaparea domeniilor discului defect, într-o manieră uniformă, în domeniile de rezervă ale discurilor funcționale rămase în matrice, aceasta să-și păstreze proprietatea de distribuție uniformă a grupurilor de paritate, este necesar ca domeniile discului defect să fie în așa fel remapate, astfel încât să rezulte o matrice aranjată într-o structură UPGD/DS bazată pe enumerarea tuturor combinațiilor de $n-1=2*k-l$ discuri luate de k ori, replicate de $y=2*x$ ori.

Notăm cu N numărul de discuri rămase în matrice după defectarea primului disc: $N=n-1$, iar cu p numărul de domenii ale discului defect (este vorba de primul disc ce s-a defectat în matrice) ce vor fi remapate în fiecare din cele N discuri funcționale rămase în matrice:

$$p = (\text{card } A)/N = (\text{card } B)/N = b/(2*N)$$

Din punct de vedere matematic, problema se traduce prin aceea de a stabili o **corespondență bijectivă** între mulțimile **A** și **B**, ținând cont în stabilirea acestei corespondențe ca, în final, exact p domenii aparținând discului defect să fie remapate în fiecare din cele N discuri funcționale rămase în matrice, în domeniile de rezervă pe care fiecare disc le deține.

Observăm că, atât în mulțimea **A** cât și în mulțimea **B**, fiecare element apare de *exact* x ori (fiecare combinație de $n-1$ discuri luate câte $k-1$ ori apare replicată de x ori în mulțimea **A**, iar în mod similar, fiecare combinație de $n-1$ discuri luate câte k ori apare replicată de asemenea de x ori în mulțimea **B**). Vom descompune în continuare mulțimea **A** în x submulțimi identice din toate punctele de vedere (atât ca număr de elemente, cât și din punct de vedere al conținutului): A_1, A_2, \dots, A_x .

Fiecare submulțime A_i , unde $i=1, 2, \dots, x$ este compusă din toate grupurile de paritate care utilizează discul defect, **luate o singură dată**:

$$\text{card } A_i = (\text{card } A)/x = b/(2^*x) = E(N, k-1), i = 1, 2, \dots, x.$$

Într-o manieră asemănătoare, vom descompune mulțimea B în x submulțimi identice B_1, B_2, \dots, B_x . Fiecare submulțime B_i , unde $i=1, 2, \dots, x$, este compusă din toate grupurile de paritate care utilizează discul defect, luate o singură dată:

$$\text{card } B_i = (\text{card } B)/x = b/(2^*x) = E(N, k), i = 1, 2, \dots, x.$$

Stabilind o corespondență biunivocă între submulțimile A_1 și B_1 , de exemplu, aceste două mulțimi având același cardinal, respectând condiția ca fiecare disc funcțional rămas în matrice să primească exact $p'=p/x$ domenii aparținând discului defect, urmând apoi să aplicăm această corespondență asupra celorlalte perechi de submulțimi (A_i, B_i), $i=2, 3, \dots, x$ și, în final, putem spune că am stabilit o corespondență bijectivă între mulțimile A și B , ceea ce ne-am și propus, iar numărul de domenii ale discului defect care vor fi remapate pe fiecare din cele $N=n-1$ discuri funcționale rămase în matrice va fi: $p^*x=p$, adică atât cât ne-am propus.

Așadar, problema se reduce la a stabili corespondența biunivocă între elementele mulțimilor A_1 și B_1 , având grijă să îndeplinim și condiția de remapare uniformă enunțată mai sus.

Pentru ca explicațiile care urmează să fie mai „aerisite“, vom folosi următoarele notații:

- $N=n-1$ este numărul de discuri care rămân în matrice după defectarea primului disc (acesta putând fi oricare disc din matrice);
- $M = \{ D_1, D_2, D_3, \dots, D_N \}$ reprezintă o mulțime având ca elemente discurile rămase în matrice după defectarea primului disc. Etichetarea acestor discuri s-a făcut doar conceptual. Considerăm aceasta mulțime ordonată după indici (în ordine crescătoare);

Din submulțimi ale mulțimii M se construiesc cele două mulțimi A_1 și B_1 , astfel:

- $A_1 = \{\text{toate submulțimile ordonate cu } k-1 \text{ elemente din mulțimea } M\}$;
- $B_1 = \{\text{toate submulțimile ordonate cu } k \text{ elemente din mulțimea } M\}$.

Cu aceste definiții, avem:

$$A_1 = \{\{D_1, D_2, \dots, D_{k-2}, D_{k-1}\}, \{D_1, D_2, \dots, D_{k-2}, D_k\}, \dots, \{D_i, D_{i+1}, \dots, D_{k+i-3}, D_{k+i-2}\}, \{D_i, D_{i+1}, \dots, D_{k+i-3}, D_{k+i-1}\}, \dots, \{D_{N-k+2}, D_{N-k+3}, \dots, D_{N-1}, D_N\}\};$$

$$B_1 = \{\{D_1, D_2, \dots, D_{k-1}, D_k\}, \{D_1, D_2, \dots, D_{k-1}, D_{k+1}\}, \dots, \{D_i, D_{i+1}, \dots, D_{k+i-2}, D_{k+i-1}\}, \{D_i, D_{i+1}, \dots, D_{k+i-2}, D_{k+i}\}, \dots, \{D_{N-k+1}, D_{N-k+2}, \dots, D_{N-1}, D_N\}\}.$$

Mulțimile A_1 și B_1 au același cardinal: $\text{card } A_1 = \text{card } B_1 = b/(2*x)$.

A stabili o corespondență bijectivă între elementele mulțimilor A_1 și B_1 înseamnă a asocia fiecărui element al mulțimii A_1 exact un element al mulțimii B_1 .

Reamintim că elementele mulțimilor A_1 și B_1 sunt la rândul lor mulțimi (pe parcursul acestui capitol, vom face deseori referire la ele numindu-le **elemente-mulțime**), mai exact, sunt submulțimi ale mulțimii M , iar un element al mulțimii A_1 are cardinalul $=k-1$, în timp ce un element al mulțimii B_1 are cardinalul $=k$.

Să notăm cu l numărul de elemente ale mulțimilor A_1 și B_1 . Cu această notație putem scrie cele două mulțimi astfel:

$$A_1 = \{A_{11}, A_{12}, \dots, A_{1l}\}$$

$$B_1 = \{B_{11}, B_{12}, \dots, B_{1l}\}$$

Numărul $p' = l/N$ ne indică l/x din numărul de domenii ale discului defect reconstruite în fiecare disc rămas în matrice, în spațiul de rezervă al acestora.

Acum, presupunând că perechea de elemente (A_{1i}, B_{1j}) se află în corespondență, A_{1i} aparținând mulțimii A_1 și B_{1j} aparținând mulțimii B_1 , făcând diferența $B_{1j} \setminus A_{1i} = \{D_t\}$, trebuie să rezulte o mulțime formată dintr-un singur element, ceea ce reprezintă o primă restricție în stabilirea legii de corespondență. Elementul D_t rezultat este un element al mulțimii M (evident).

Făcând reuniunea tuturor mulțimilor (în număr total de l) de forma:

$$\{B_{1j} \setminus A_{1i}\},$$

unde i, j aparțin mulțimii $\{1, 2, 3, \dots, l\}$, iar A_{1i} și B_{1j} sunt două elemente aparținând mulțimii A_1 , respectiv B_1 , între care s-a stabilit corespondența, trebuie ca mulțimea creată ca rezultat al acestei reuniuni

să conțină toate elementele mulțimii M , fiecare fiind **replacat de p' ori**. Aceasta reprezintă a doua restricție impusă (această restricție este cel mai greu de realizat), care este chiar condiția cerută încă din start și dacă ea se respectă, aria de discuri rezultată după remaparea domeniilor discului defect în spațiul de rezervă al fiecărui disc rămas în ea (remapare care se va face în acest caz **uniform**) va păstra proprietatea de UPGD/DS.

4.6 Algoritm pentru stabilirea corespondentei biunivoce între elementele mulțimilor A_1 și B_1

După cum am mai spus, elementele mulțimilor A_1 și B_1 sunt ordonate, în ordinea crescătoare a indicilor. Fiecărui element A_{1i} aparținând mulțimii A_1 îi atașăm mulțimea complement a elementului A_{1i} (care la rândul său este o submulțime a mulțimii M) față de mulțimea $M: A_{1i} \leftarrow M \setminus A_{1i}$

Ca rezultat al acestor asocieri, alcătuim o mulțime C_1 care are tot l elemente, iar fiecare element din C_1 este o mulțime, în fapt fiind mulțimea complement a unui element din A_1 față de mulțimea M , pentru toate elementele din A_1 .

Pentru situația noastră, în care $N=2*k-1$ este identică cu mulțimea B_1 , singura diferență fiind aceea că elementele mulțimii C_1 sunt ordonate în ordinea descrescătoare a indicilor.

Vom explica principial, în cele ce urmează, modul în care se stabilește o corespondență între un element al mulțimii A_1 și un element al mulțimii B_1 :

Dintr-un element al mulțimii C_1 , care la rândul ei este o mulțime de forma $C_{1z}=\{D_{zk}, D_{zk+1}, \dots, D_s, \dots, D_{zN}\}$ extragem elementul D_s . Această alegere a lui D_s se face după o regulă ce va fi definită mai târziu. Elementul asociat lui C_{1z} din mulțimea A_1 este de forma: $A_{1z}=\{D_{z1}, D_{z2}, \dots, D_{zk-1}\}$

Adăugăm elementul D_s lui A_{1z} , rezultând:

$$A_{1z}=\{D_s, D_{z1}, D_{z2}, \dots, D_{zk-1}\},$$

mulțime care, după ordonarea în ordine crescătoare a indicilor, este identică cu un element din mulțimea B_1 , $A_{1z}=B_{1z}$.

În acest fel, am stabilit o corespondență între elementul A_{1z} aparținând mulțimii A_1 și elementul B_{1z} care aparține mulțimii B_1 .

Se elimină elementul A_{1z} din mulțimea A_1 , elementul asociat C_{1z} din mulțimea C_1 și elementul B_{1z} din mulțimea B_1 . De asemenea, se vor căuta și se vor elimina din mulțimile C_{1r} (care sunt elemente ale lui C_1) toate elementele D_r care, adăugate elementului A_{1r} asociat lui C_{1r} , formează, după ordonare, o mulțime identică cu B_{1z} . Aceste suprimări vor ajuta la direcționarea căutarilor ulterioare pentru a stabili o corespondență între alte două elemente din A_1 , respectiv B_1 , pentru ca funcția de corespondență să fie bijectivă.

Revenind la aria de discuri, această corespondență semnifică faptul că un domeniu aparținând discului defect a fost reconstituit, pe baza parității, utilizând domeniile grupului de paritate asociat elementului A_{1z} și a fost remapat într-unul din domeniile de rezervă ale discului numărul s (acesta fiind unul din cele N discuri funcționale rămase în matrice. Din cele k posibilități în care acest domeniu putea fi remapat (însemnând cele k discuri care nu utilizează grupul de paritate din care face parte domeniul luat în discuție) a fost aleasă, pe baza unei reguli precise, una din ele. Modul de remapare este următorul:

Primul disc rămas în matrice în care vom reconstrui un domeniu al discului defect este D_i , unde $i=1$ (rezultând așadar discul D_1).

Se repetă de un număr de $N \cdot p'$ ori următorii 5 pași:

Se caută, consultând lista elementelor-mulțimi ale mulțimii C_1 începând cu primul din lista, până se găsește primul element-mulțime în a cărui componentă se află D_i . Fie C_{iz} acest element-mulțime. Acesta are asociat elementul-mulțime A_{iz} , aparținând mulțimii A_1 .

Adăugând elementul D_i la A_{iz} rezultă, după ordonarea noului element-mulțime A_{iz} , că acesta este identic cu un element-mulțime aparținând mulțimii B_1 , notat B_{iz} . Se consideră că s-a stabilit o corespondență între elementul A_{iz} din A_1 și elementul B_{iz} din B_1 , corespondență care se memorează.

Se consultă din nou lista elementelor-mulțimi ale mulțimii C_1 și, dacă se găsesc în C_1 elemente-mulțimi care au în componența lor pe D_i , iar dacă prin adăugarea lui D_i la elementul-mulțime asociat, aparținând lui A_1 , se obține, după ordonare, un element identic cu B_{iz} , se suprimă D_i din elementul-mulțime căruia îi aparținea înaintea executării acestui pas.

De asemenea, se elimină D_i și din elementul-mulțime asociat, aparținând lui A_1 (în acest element-mulțime, D_i a fost adăugat la începutul acestui pas).

Se elimina urmatoarele elemente multime: $A_{iz} \subset A_1$; $B_{iz} \subset B_1$; $C_{iz} \subset C_1$.

Se atribuie variabilei i valoarea $i = imod n$ și se trece la pasul 1).

În final, mulțimile A_1 , B_1 și C_1 vor fi vide.

Ca urmare a executării acestui algoritim, deoarece fiecare element D_i , i luând valorile de la 1 la N face obiectul prelucrării mulțimilor A_1 , B_1 și C_1 de p' ori, înseamnă că fiecare din cele N discuri rămase în matrice va remapa p' domenii ale discului defect, deci, această remapare se realizează uniform printre cele N discuri rămase în matrice.

Aplicând analog această lege de corespondență pentru toate celelalte perechi de mulțimi (A_i , B_i), $i=2,3,\dots,x$ care sunt simple replici ale mulțimilor A_1 , respectiv B_1 , vom avea situația dorită, adică aceea că fiecare din cele N discuri rămase în matrice va remapa exact $p=p'*x$ domenii ale discului defect. Mai mult decât atât, noua matrice va fi aranjată ca o structură UPGD/DS bazată pe enumerarea tuturor combinațiilor de N discuri ($N=n-1$) luate câte k (fiecare din aceste combinații reprezentând un model unic de grup de paritate), **replicate de $2*x$ ori**.

A mai rămas să explicăm cum se alege numărul x (numărul inițial de replici a fiecărui model unic de grup de paritate). Acest număr este legat doar de cea de-a doua defecțiune de disc din matrice și este ales astfel încât după reconstruirea domeniilor celui de-al doilea disc defect în spațiul de rezervă al celor $n-2$ discuri rămase în matrice, aceasta să fie aranjată într-o structură UPGD/DS.

Am reușit să demonstrăm faptul că, alegând corespunzător numărul x , remaparea domeniilor în urma defecțiunii celui de-al doilea disc se poate realiza utilizând algoritmul prezentat la capitolul 4.

Pentru aceasta, ținând cont de faptul că după prima defecțiune, avem în matrice un număr de $E(n-1,k)$ modele unice de grupuri de paritate replicate de $2*x$ ori, din acestea, la apariția celui de-al doilea defect, situația se prezintă astfel:

- $2*x*E(n-2,k-1)$ grupuri de paritate utilizează discul defect;
- $2*x*E(n-2,k)$ grupuri de paritate nu utilizează discul defect.

Deoarece ne rezumăm doar la două defecte, condiția ce se impune pentru ca aria rezultată după remapare să fie UPGD, este ca din cele $2^x \cdot E(n-2, k-1)$ grupuri de paritate ce utilizează discul defect să formăm un număr întreg de replici ale fiecărui model unic de grup de paritate ce nu utilizează discul defect. Aceasta se traduce matematic prin faptul că trebuie ales cel mai mic număr x pentru care 2^x divide pe $n-2-(k-1)$, adică $2^x / (2^k - k - 1) = 2^x / (k-1) = q$ și q este cel mai mic întreg pozitiv.

Cu acest x calculat, după remaparea domeniilor celui de-al doilea disc defect, aria va fi aranjată într-o structură de tip UPGD bazată pe enumerarea tuturor combinațiilor de $n-2$ discuri luate câte k , replicate de $(2^x + q)$ ori.

Exemple:

$$n = 6, k = 3 \Rightarrow k-1 = 2, x = 1, q = 1;$$

$$n = 8, k = 4 \Rightarrow k-1 = 3, x = 3, q = 2;$$

$$n = 10, k = 5 \Rightarrow k-1 = 4, x = 2, q = 1.$$

Referitor la numărul domeniilor de rezervă care trebuie adăugate inițial fiecărui disc din matrice, acesta este calculat cu expresia:

$$s_0 = b \cdot k \cdot f / (n \cdot (n-f)) = b \cdot k \cdot f / ((2^k) \cdot (2^k - f)) = b \cdot f / (2^k \cdot (2^k - f)),$$

deoarece $k = n/2$.

Algoritmul descris în acest capitol a fost verificat experimental, pentru mai multe triplete $(2^k, k, 2)$, ultimul element al tripletului reprezentând numărul de defecte admise, care în descrierea noastră trebuie să fie maxim 2.

Următorul tabel este oferit spre a face o comparație între parametrii cei mai importanți necesari celor 2 algoritmi:

| n | k | F | Algoritmul inițial | | | Algoritmul îmbunătățit | | |
|----|---|---|--------------------|----------------|----------------|------------------------|----------------|----------------|
| | | | b | r ₀ | s ₀ | b | r ₀ | s ₀ |
| 6 | 3 | 2 | 60 | 30 | 15 | 20 | 10 | 5 |
| 8 | 4 | 2 | 840 | 420 | 140 | 210 | 105 | 35 |
| 10 | 5 | 2 | 2520 | 1260 | 315 | 504 | 252 | 63 |
| 12 | 6 | 2 | 27720 | 13860 | 2772 | 4620 | 2310 | 462 |
| 14 | 7 | 2 | 72072 | 36036 | 6006 | 10296 | 6148 | 858 |
| 16 | 8 | 2 | 720720 | 360360 | 51480 | 90090 | 45045 | 6435 |

4.7 Exemplu de aplicare a algoritmului îmbunătățit de creare a UPGD/ DS

Algoritmul pe care tocmai l-am descris îl vom exemplifica în cele ce urmează.

Vom considera cazul în care numărul de discuri în matrice este inițial egal cu 6, dimensiunea grupului de paritate este 3, iar spațiul de rezervă alocat fiecărui disc pentru 2 defecțiuni ($n=6, k=3, f=2$).

Alegerea numărului inițial de replici x :

Se caută cel mai mic număr întreg pozitiv q astfel ca $2*x/(k-1)=q$.
Rezultând $x=1, q=1$.

Determinarea lui b :

$$b = x * E(n, k) = 1 * E(6, 3) = 20$$

Formarea matricei UPGD/DS inițiale:

| D1 | D2 | D3 | D4 | D5 | D6 |
|----|----|----|----|----|----|
| A | A | A | | | |
| B | B | | B | | |
| C | C | | | C | |
| D | D | | | | D |
| E | | E | E | | |
| F | | F | | F | |
| G | | G | | | G |
| H | | | H | H | |
| I | | | I | | I |
| J | | | | J | J |
| | K | K | K | | |
| | L | L | | L | |
| | M | M | | | M |
| | N | | N | N | |
| | O | | O | | O |
| | P | | | P | P |
| | | Q | Q | Q | |
| | | R | R | | R |
| | | S | | S | S |
| | | | T | T | T |
| □ | □ | □ | □ | □ | □ |
| □ | □ | □ | □ | □ | □ |
| □ | □ | □ | □ | □ | □ |
| □ | □ | □ | □ | □ | □ |
| □ | □ | □ | □ | □ | □ |

Figura 4-7

| D1 | D2 | D3 | D4 | D5 | D6 |
|----|----|----|----|----|----|
| A | A | A | | | |
| B | B | | B | | |
| C | C | | | C | |
| E | | E | E | | |
| F | | F | | F | |
| H | | | H | H | |
| | K | K | K | | |
| | L | L | | L | |
| | N | | N | N | |
| | | Q | Q | Q | |
| D | D | | | | D |
| G | | G | | | G |
| I | | | I | | I |
| J | | | | J | J |
| | M | M | | | M |
| | O | | O | | O |
| | P | | | P | P |
| | | R | R | | R |
| | | S | | S | S |
| | | | T | T | T |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |

Figura 4-8

Enumerăm toate combinațiile de 6 discuri luate câte 3 și fiecare combinație de discuri astfel enumerată formează un grup de paritate.

Deoarece $x=1$, inițial în matrice fiecare model de grup de paritate apare o singură dată.

Numărul de domenii asignate fiecărui disc va fi $r_0=(b*k)/n=10$

Numărul de domenii de rezervă pentru fiecare disc este $s_0=b*f/((2*(2*k-f))=5$

În Figura 4-7 structura matricei UPGD/DS înaintea defectării primului disc.

Presupunem că discul etichetat cu D₆ se defectează, aria fiind formată acum din discurile D₁, D₂, D₃, D₄, D₅. În Figura 4-8, toate grupurile de paritate care nu utilizează discul defect D₆ au fost grupate conceptual (să nu se înțeleagă greșit că au fost mutate fizic), formând mulțimea B, care are $x*E(5,3)=10$ elemente (B coincide cu B₁). Elementele acestei mulțimi sunt grupate în prima parte a figurii.

Analog, toate grupurile de paritate care utilizează discul defect D₆ au fost grupate formând mulțimea A, care are $x*E(5,2)=10$ elemente (și, aici, A coincide cu A₁). Elementele acestei mulțimi au fost grupate în continuarea mulțimii B.

Deoarece în acest caz mulțimile A și B sunt totuna cu mulțimile A₁ și B₁, vom aplica algoritmul folosind notațiile A, respectiv B.

Se observă că mulțimea A este formată din toate combinațiile de $N=5$ discuri luate câte $k-1=2$, iar mulțimea B este formată din toate combinațiile de $N=5$ discuri, luate câte $k=3$. În ambele mulțimi, fiecare combinație astfel enumerată formează un grup de paritate.

După remapare, fiecare din cele $n=5$ discuri din matrice va avea: $p=p'=r_0/N=2$ domenii reconstruite ale discului defect D₆, iar mulțimea A trebuie să se „transforme” într-una de tip B, adică să fie formată din enumerarea tuturor combinațiilor de $N=5$ discuri luate câte $k=3$ ori. Această „transformare” se traduce de fapt prin stabilirea unei corespondențe bijective între elementele mulțimii A și elementele mulțimii B.

În cele ce urmează, vom parcurge ciclul cu ciclul algoritmul descris în acest capitol și vom prezenta după fiecare ciclu starea matricei de discuri (modul cum se remapează domeniile discului defect în spațiul de rezervă al fiecărui disc rămas în matrice).

Generarea mulțimilor A, B și C

Figura următoare ilustrează cele trei mulțimi (structura lor în faza inițială). Pentru a ușura discuția, etichetăm elementele celor trei mulțimi cu A_i , B_i , respectiv C_i , unde $i=1,2,\dots,10$, cu observația că există o asociere între elementul A_i al mulțimii A și elementul C_i al mulțimii C, pentru fiecare valoare a lui i .

| Mulțimea A | Mulțimea C | Mulțimea B |
|----------------------|--|--------------------------|
| $A^1=\{D_1,D_2\}$ | $\langle\text{asociere}\rangle C^1=\{D_3,D_4,D_5\}$ | $B^1=\{D_1,D_2,D_3\}$ |
| $A^2=\{D_1,D_3\}$ | $\langle\text{asociere}\rangle C^2=\{D_2,D_4,D_5\}$ | $B^2=\{D_1,D_2,D_4\}$ |
| $A^3=\{D_1,D_4\}$ | $\langle\text{asociere}\rangle C^3=\{D_2,D_3,D_5\}$ | $B^3=\{D_1,D_2,D_5\}$ |
| $A^4=\{D_1,D_5\}$ | $\langle\text{asociere}\rangle C^4=\{D_2,D_3,D_4\}$ | $B^4=\{D_1,D_3,D_4\}$ |
| $A^5=\{D_2,D_3\}$ | $\langle\text{asociere}\rangle C^5=\{D_1,D_4,D_5\}$ | $B^5=\{D_1,D_3,D_5\}$ |
| $A^6=\{D_2,D_4\}$ | $\langle\text{asociere}\rangle C^6=\{D_1,D_3,D_5\}$ | $B^6=\{D_1,D_4,D_5\}$ |
| $A^7=\{D_2,D_5\}$ | $\langle\text{asociere}\rangle C^7=\{D_1,D_3,D_4\}$ | $B^7=\{D_2,D_3,D_4\}$ |
| $A^8=\{D_3,D_4\}$ | $\langle\text{asociere}\rangle C^8=\{D_1,D_2,D_5\}$ | $B^8=\{D_2,D_3,D_5\}$ |
| $A^9=\{D_3,D_5\}$ | $\langle\text{asociere}\rangle C^9=\{D_1,D_2,D_4\}$ | $B^9=\{D_2,D_4,D_5\}$ |
| $A^{10}=\{D_4,D_5\}$ | $\langle\text{asociere}\rangle C^{10}=\{D_1,D_2,D_3\}$ | $B^{10}=\{D_3,D_4,D_5\}$ |

Fiecare element al mulțimii A are atașat elementul complementar în raport cu mulțimea $M=\{D_1,D_2,D_3,D_4,D_5\}$

Avem în total $N \cdot p' = 10$ treceri prin algoritm.

Prima trecere

Parcurgând lista elementelor mulțimii C de sus în jos, căutam primul element C_i care-l conține pe D_1 . Acesta este $C^5=\{D_1,D_4,D_5\}$.

Adăugăm elementului asociat lui C^5 , și anume A^5 , pe D_1 . După ordonarea elementului A^5 extins, obținem: $\{D_1,D_2,D_3\}$ care este chiar B^1 . Astfel, am stabilit o corespondență între elementul A^5 aparținând mulțimii A și elementul B^1 aparținând mulțimii B.

Spunem că am remapat un domeniu al discului defect D_6 și anume acela care intră în componența modelului de grup de paritate asociat lui A^5 neextins într-unul din domeniile de rezervă ale discului D_1 , formându-se cu acest prilej un model de grup de paritate reprezentat de B^1 .

Concret, domeniul M din D_6 se remapează într-unul din domeniile de rezervă ale discului D_1 , după această remapare aria arătând ca în Figura 4-9. Se observă că acum discul D_1 are cu un spațiu de rezervă mai puțin decât celelalte discuri funcționale rămase în matrice. În continuare, se inspectează din nou lista elementelor mulțimii C , căutând acele elemente C_i care conțin un D_k care adăugat lui A_i , face ca după ordonarea acestuia, să fie identic cu B^l . În acest caz, elementul D_k se elimină din C_i . Aplicând această regulă pe exemplul dat, vom suprima următoarele elemente: D_3 din mulțimea C^1 și D_2 din mulțimea C^2 .

În finalul acestei treceri, vom memoram corespondența stabilită și eliminăm elementele-mulțime între care s-a stabilit corespondența, precum și pe cel corespunzător din mulțimea C . Eliminăm pe: A^5 din mulțimea A , C^5 asociat lui A^5 din mulțimea C , B^1 din mulțimea B .

Noua structură a mulțimilor A , B și C :

| Mulțimea A | Mulțimea C | Mulțimea B |
|-------------------------|--|------------------------------|
| $A^1 = \{D_1, D_2\}$ | $\langle \text{asociere} \rangle C^1 = \{D_4, D_5\}$ | $B^2 = \{D_1, D_2, D_4\}$ |
| $A^2 = \{D_1, D_3\}$ | $\langle \text{asociere} \rangle C^2 = \{D_4, D_5\}$ | $B^3 = \{D_1, D_2, D_5\}$ |
| $A^3 = \{D_1, D_4\}$ | $\langle \text{asociere} \rangle C^3 = \{D_2, D_3, D_5\}$ | $B^4 = \{D_1, D_3, D_4\}$ |
| $A^4 = \{D_1, D_5\}$ | $\langle \text{asociere} \rangle C^4 = \{D_2, D_3, D_4\}$ | $B^5 = \{D_1, D_3, D_5\}$ |
| $A^6 = \{D_2, D_4\}$ | $\langle \text{asociere} \rangle C^6 = \{D_1, D_3, D_5\}$ | $B^6 = \{D_1, D_4, D_5\}$ |
| $A^7 = \{D_2, D_5\}$ | $\langle \text{asociere} \rangle C^7 = \{D_1, D_3, D_4\}$ | $B^7 = \{D_2, D_3, D_4\}$ |
| $A^8 = \{D_3, D_4\}$ | $\langle \text{asociere} \rangle C^8 = \{D_1, D_2, D_5\}$ | $B^8 = \{D_2, D_3, D_5\}$ |
| $A^9 = \{D_3, D_5\}$ | $\langle \text{asociere} \rangle C^9 = \{D_1, D_2, D_4\}$ | $B^9 = \{D_2, D_4, D_5\}$ |
| $A^{10} = \{D_4, D_5\}$ | $\langle \text{asociere} \rangle C^{10} = \{D_1, D_2, D_3\}$ | $B^{10} = \{D_3, D_4, D_5\}$ |

A doua trecere

Parcurgând lista elementelor mulțimii C de sus în jos, căutam primul element C_i care-l conține pe D_2 . Acesta este $C^3 = \{D_2, D_3, D_5\}$. Adăugăm elementului asociat lui C^3 și anume A^3 , pe D_2 . După ordonarea elementului A^3 extins, obținem $\{D_1, D_2, D_4\}$ care este chiar B^2 . Astfel, am stabilit o corespondență între elementul A^3 aparținând mulțimii A și elementul B^2 aparținând mulțimii B .

Spunem că am remapat un domeniu al discului defect D_6 și anume acela care intră în componența modelului de grup de paritate asociat lui A^3 neextins, într-unul din domeniile de rezervă ale discului D_2 , formându-se cu acest prilej un model de grup de paritate reprezentat de B^2 .

Concret, domeniul I din D_6 se remapează într-unul din domeniile de rezervă ale discului D_2 , după această remapare aria arătând ca în Figura 4-10. Se observă că acum discul D_2 are, de asemenea, cu un spațiu de rezervă mai puțin. Aplicând în continuare pasul 3 al algoritmului, vom suprima următoarele elemente: D_4 din mulțimea C^1 și D_1 din mulțimea C^6 .

În finalul acestei treceri, vom memora corespondența stabilită și eliminăm elementele-mulțime între care s-a stabilit corespondența, precum și pe cel corespunzător din mulțimea C. Eliminăm pe A^3 din mulțimea A, C^3 asociat lui A^3 din mulțimea C, B^2 din mulțimea B.

Noua structură a mulțimilor A, B și C:

| MULȚIMEA A | MULȚIMEA C | MULȚIMEA B |
|---|------------|------------------------------|
| $A^1 = \{D_1, D_2\} \leftarrow \text{asociere} \rightarrow C^1 = \{D_5\}$ | | $B^3 = \{D_1, D_2, D_5\}$ |
| $A^2 = \{D_1, D_3\} \leftarrow \text{asociere} \rightarrow C^2 = \{D_4, D_5\}$ | | $B^4 = \{D_1, D_3, D_4\}$ |
| $A^4 = \{D_1, D_5\} \leftarrow \text{asociere} \rightarrow C^4 = \{D_2, D_3, D_4\}$ | | $B^5 = \{D_1, D_3, D_5\}$ |
| $A^6 = \{D_2, D_4\} \leftarrow \text{asociere} \rightarrow C^6 = \{D_3, D_5\}$ | | $B^6 = \{D_1, D_4, D_5\}$ |
| $A^7 = \{D_2, D_5\} \leftarrow \text{asociere} \rightarrow C^7 = \{D_1, D_3, D_4\}$ | | $B^7 = \{D_2, D_3, D_4\}$ |
| $A^8 = \{D_3, D_4\} \leftarrow \text{asociere} \rightarrow C^8 = \{D_1, D_2, D_5\}$ | | $B^8 = \{D_2, D_3, D_5\}$ |
| $A^9 = \{D_3, D_5\} \leftarrow \text{asociere} \rightarrow C^9 = \{D_1, D_2, D_4\}$ | | $B^9 = \{D_2, D_4, D_5\}$ |
| $A^{10} = \{D_4, D_5\} \leftarrow \text{asociere} \rightarrow C^{10} = \{D_1, D_2, D_3\}$ | | $B^{10} = \{D_3, D_4, D_5\}$ |

A treia trecere

Parcurgând lista elementelor mulțimii C de sus în jos, căutăm primul element C_i care-l conține pe D_3 . Acesta este $C^4 = \{D_2, D_3, D_4\}$. Adăugăm elementului asociat lui C^4 și anume A^4 , pe D_3 . După ordonarea elementului A^4 extins obținem: $\{D_1, D_3, D_5\}$ care este chiar

B^5 . Astfel, am stabilit o corespondență între elementul A^4 aparținând mulțimii A și elementul B^5 aparținând mulțimii B .

Spunem că am remapat un domeniu al discului defect D_6 și anume acela care intră în componența modelului de grup de paritate asociat lui A^4 neextins, într-unul din domeniile de rezervă ale discului D_3 , formându-se cu acest prilej un model de grup de paritate reprezentat de B^5 .

Concret, domeniul J din D_6 se remapează într-unul din domeniile de rezervă ale discului D_3 , după această remapare aria arătând ca în Figura 4-11. Discul D_3 are acum cu un spațiu de rezervă mai puțin.

| D1 | D2 | D3 | D4 | D5 | D6 |
|----|----|----|----|----|----|
| A | A | A | | | |
| B | B | | B | | |
| C | C | | | C | |
| E | | E | E | | |
| F | | F | | F | |
| H | | H | H | | |
| | K | K | K | | |
| | L | L | | L | |
| | N | | N | N | |
| | | Q | Q | Q | |
| D | D | | | | D |
| G | | G | | | G |
| I | | | I | | I |
| J | | | | J | J |
| M | M | M | | | |
| | O | | O | | O |
| | P | | | P | P |
| | | R | R | | R |
| | | S | | S | S |
| | | | T | T | T |
| | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |

Figura 4-9

| | | | | | |
|----|----|----|----|----|----|
| D1 | D2 | D3 | D4 | D5 | D6 |
| A | A | A | | | |
| B | B | | B | | |
| C | C | | | C | |
| E | | E | E | | |
| F | | F | | F | |
| H | | H | H | | |
| | K | K | K | | |
| | L | L | | L | |
| | N | | N | N | |
| | | Q | Q | Q | |
| D | D | | | | D |
| G | | G | | | G |
| I | | I | I | | |
| J | | | | J | J |
| M | M | M | | | |
| | O | | O | | O |
| | P | | | P | P |
| | | R | R | | R |
| | | S | | S | S |
| | | | T | T | T |
| | | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |

Figura 4-10

| | | | | | |
|----|----|----|----|----|----|
| D1 | D2 | D3 | D4 | D5 | D6 |
| A | A | A | | | |
| B | B | | B | | |
| C | C | | | C | |
| E | | E | E | | |
| F | | F | | F | |
| H | | H | H | | |
| | K | K | K | | |
| | L | L | | L | |
| | N | | N | N | |
| | | Q | Q | Q | |
| D | D | | | | D |
| G | | G | | | G |
| I | | I | I | | |
| J | | J | | J | |
| M | M | M | | | |
| | O | | O | | O |
| | P | | | P | P |
| | | R | R | | R |
| | | S | | S | S |
| | | | T | T | T |
| | | | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |

Figura 4-11

Aplicând în continuare pasul 3 al algoritmului, vom suprima următoarele elemente D5 din mulțimea C², D1 din mulțimea C⁹. În finalul acestei treceri, memorăm corespondența stabilită și eliminăm elementele-mulțime între care s-a stabilit corespondența, precum și pe

cel corespunzător din mulțimea C. Eliminăm pe A^4 din mulțimea A, C^4 asociat lui A^4 din mulțimea C, B^5 din mulțimea B.

Noua structură a mulțimilor A, B și C:

| Mulțimea A | Mulțimea C | Mulțimea B |
|---|------------|------------------------------|
| $A^1 = \{D_1, D_2\}$ <-asociere-> $C^1 = \{D_5\}$ | | $B^3 = \{D_1, D_2, D_5\}$ |
| $A^2 = \{D_1, D_3\}$ <-asociere-> $C^2 = \{D_4\}$ | | $B^4 = \{D_1, D_3, D_4\}$ |
| $A^6 = \{D_2, D_4\}$ <-asociere-> $C^6 = \{D_3, D_5\}$ | | $B^6 = \{D_1, D_4, D_5\}$ |
| $A^7 = \{D_2, D_5\}$ <-asociere-> $C^7 = \{D_1, D_3, D_4\}$ | | $B^7 = \{D_2, D_3, D_4\}$ |
| $A^8 = \{D_3, D_4\}$ <-asociere-> $C^8 = \{D_1, D_2, D_5\}$ | | $B^8 = \{D_2, D_3, D_5\}$ |
| $A^9 = \{D_3, D_5\}$ <-asociere-> $C^9 = \{D_2, D_4\}$ | | $B^9 = \{D_2, D_4, D_5\}$ |
| $A^{10} = \{D_4, D_5\}$ <-asociere-> $C^{10} = \{D_1, D_2, D_3\}$ | | $B^{10} = \{D_3, D_4, D_5\}$ |

A patra trecere

Parcurgând lista elementelor mulțimii C de sus în jos, căutam primul element C^i care-l conține pe D_4 . Acesta este $C^2 = \{D_4\}$. Adăugăm elementului asociat lui C^2 și anume A^2 , pe D_4 . După ordonarea elementului A^2 extins, obținem $\{D_1, D_3, D_4\}$ care este chiar B^4 . Astfel, am stabilit o corespondență între elementul A^2 aparținând mulțimii A și elementul B^4 aparținând mulțimii B.

| | | | | | |
|----|----|----|----|----|----|
| D1 | D2 | D3 | D4 | D5 | D6 |
| A | A | A | | | |
| B | B | | B | | |
| C | C | | | C | |
| E | | E | E | | |
| F | | F | | F | |
| H | | H | H | | |
| | K | K | K | | |
| | L | L | | L | |
| | N | | N | N | |
| | | Q | Q | Q | |
| D | D | | | | D |
| G | | G | G | | |
| I | I | | I | | |
| J | | J | | J | |
| M | M | M | | | |
| | O | | O | | O |
| | P | | | P | P |
| | | R | R | | R |
| | | S | | S | S |
| | | | T | T | T |
| | | | | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |

Figura 4-12

| D1 | D2 | D3 | D4 | D5 | D6 |
|----|----|----|----|----|----|
| A | A | A | | | |
| B | B | | B | | |
| C | C | | | C | |
| E | | E | E | | |
| F | | F | | F | |
| H | | H | H | | |
| | K | K | K | | |
| | L | L | | L | |
| | N | | N | N | |
| | | Q | Q | Q | |
| D | D | | | D | |
| G | | G | G | | |
| I | I | | I | | |
| J | | J | | J | |
| M | M | M | | | |
| | O | | O | | O |
| | P | | | P | P |
| | | R | R | | R |
| | | S | | S | S |
| | | | T | T | T |
| | | | | | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |

Figura 4-13

| D1 | D2 | D3 | D4 | D5 | D6 |
|----|----|----|----|----|----|
| A | A | A | | | |
| B | B | | B | | |
| C | C | | | C | |
| E | | E | E | | |
| F | | F | | F | |
| H | | H | H | | |
| | K | K | K | | |
| | L | L | | L | |
| | N | | N | N | |
| | | Q | Q | Q | |
| D | D | | | D | |
| G | | G | G | | |
| I | I | | I | | |
| J | | J | | J | |
| M | M | M | | | |
| | O | | O | | O |
| | P | | | P | P |
| | | R | R | | R |
| | | S | | S | S |
| T | | | T | T | |
| | | | | | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |

Figura 4-14

Spunem că am remapat un domeniu al discului defect D_6 și anume acela care intră în componenta modelului de grup de paritate asociat lui A_2 neextins, într-unul din domeniile de rezervă ale discului D_4 ,

formându-se cu acest prilej un model de grup de paritate reprezentat de B^4 . Concret, domeniul G din D_6 se remapează într-unul din domeniile de rezervă ale discului D_4 , după această remapare aria arătând ca în Figura 4-12. Discul D_4 are acum cu un spațiu de rezervă mai puțin. Aplicând în continuare pasul 3 al algoritmului, vom suprima următorul element: D_1 din mulțimea C^8 . În finalul acestei treceri, vom memora corespondența stabilită și eliminăm elementele-mulțime între care s-a stabilit corespondența, precum și pe cel corespunzător din mulțimea C . Eliminăm pe A^2 din mulțimea A , C^2 asociat lui A^2 din mulțimea C , B^4 din mulțimea B .

Noua structură a mulțimilor A , B și C :

| Mulțimea A | Mulțimea C | Mulțimea B |
|-------------------------|------------------------------|------------------------------|
| $A^1 = \{D_1, D_2\}$ | $C^1 = \{D_5\}$ | $B^3 = \{D_1, D_2, D_5\}$ |
| $A^6 = \{D_2, D_4\}$ | $C^6 = \{D_3, D_5\}$ | $B^6 = \{D_1, D_4, D_5\}$ |
| $A^7 = \{D_2, D_5\}$ | $C^7 = \{D_1, D_3, D_4\}$ | $B^7 = \{D_2, D_3, D_4\}$ |
| $A^8 = \{D_3, D_4\}$ | $C^8 = \{D_2, D_5\}$ | $B^8 = \{D_2, D_3, D_5\}$ |
| $A^9 = \{D_3, D_5\}$ | $C^9 = \{D_2, D_4\}$ | $B^9 = \{D_2, D_4, D_5\}$ |
| $A^{10} = \{D_4, D_5\}$ | $C^{10} = \{D_1, D_2, D_3\}$ | $B^{10} = \{D_3, D_4, D_5\}$ |

A cincea trecere

Parcurgând lista elementelor mulțimii C de sus în jos, căutam primul element C^i care-l conține pe D_5 . Acesta este $C^1 = \{D_5\}$. Adăugăm elementului asociat lui C^1 și anume A^1 , pe D_5 . După ordonarea elementului A^1 extins, obținem $\{D_1, D_2, D_5\}$ care este chiar B^3 . Astfel, am stabilit o corespondență între elementul A^1 aparținând mulțimii A și elementul B^3 aparținând mulțimii B .

Spunem că am remapat un domeniu al discului defect D_6 și anume acela care intră în componența modelului de grup de paritate asociat lui A^1 neextins, într-unul din domeniile de rezervă ale discului D_5 , formându-se cu acest prilej un model de grup de paritate reprezentat de B^3 .

Concret, domeniul D din D_6 se remapează într-unul din domeniile de rezervă ale discului D_5 , după această remapare aria arătând ca în Figura 4-13. Discul D_5 are acum cu un spațiu de rezervă mai puțin. Aplicând în continuare pasul 3 al algoritmului, vom suprima următorul element D_1 din mulțimea C^7 . În finalul acestei treceri, vom memora corespondența stabilită și eliminăm elementele-mulțime între care s-a stabilit corespondența, precum și pe cel corespunzător din mulțimea C . Eliminăm pe A^1 din mulțimea A , C^1 asociat lui A^1 din mulțimea C , B^3 din mulțimea B .

Noua structură a mulțimilor A , B și C :

| Mulțimea A | Mulțimea C | Mulțimea B |
|---|------------|------------------------------|
| $A^6 = \{D_2, D_4\} \leftarrow \text{asociere} \rightarrow C^6 = \{D_3, D_5\}$ | | $B^6 = \{D_1, D_4, D_5\}$ |
| $A^7 = \{D_2, D_5\} \leftarrow \text{asociere} \rightarrow C^7 = \{D_3, D_4\}$ | | $B^7 = \{D_2, D_3, D_4\}$ |
| $A^8 = \{D_3, D_4\} \leftarrow \text{asociere} \rightarrow C^8 = \{D_2, D_5\}$ | | $B^8 = \{D_2, D_3, D_5\}$ |
| $A^9 = \{D_3, D_5\} \leftarrow \text{asociere} \rightarrow C^9 = \{D_2, D_4\}$ | | $B^9 = \{D_2, D_4, D_5\}$ |
| $A^{10} = \{D_4, D_5\} \leftarrow \text{asociere} \rightarrow C^{10} = \{D_1, D_2, D_3\}$ | | $B^{10} = \{D_3, D_4, D_5\}$ |

Deci, până acum am remapat câte un domeniu din discul defect D_6 în fiecare din discurile rămase în matrice. Începând cu următoarea trecere, vom lua iarăși pe rând câte un disc din matrice în care vom remapa domeniul din D_6 .

A șasea trecere

Parcurgând lista elementelor mulțimii C de sus în jos, căutam primul element C_i care-l conține pe D_1 . Acesta este $C^{10} = \{D_1, D_2, D_3\}$. Adăugăm elementului asociat lui C^{10} și anume A^{10} , pe D_1 . După ordonarea elementului A^{10} extins, obținem $\{D_1, D_4, D_5\}$ care este chiar B^6 . Astfel, am stabilit o corespondență între elementul A^{10} aparținând mulțimii A și elementul B^6 aparținând mulțimii B .

Spunem că am remapat un domeniu al discului defect D_6 și anume acela care intră în componența modelului de grup de paritate asociat lui A^{10} neextins, într-unul din domeniile de rezervă ale discului D_1 , formându-se cu acest prilej un model de grup de paritate reprezentat de B^6 .

Concret, domeniul T din D_6 se remapează într-unul din domeniile de rezervă ale discului D_1 , după această remapare aria arătând ca în Figura 4-14. Aplicând în continuare pasul 3 al algoritmului, nu vom suprima nimic.

În finalul acestei treceri, vom memora corespondența stabilită și eliminăm elementele-mulțime între care s-a stabilit corespondența, precum și pe cel corespunzător din mulțimea C. Eliminăm pe A^{10} din mulțimea A, C^{10} asociat lui A^{10} din mulțimea C, B^6 din mulțimea B.

Noua structură a mulțimilor A, B și C:

| MULȚIMEA A | MULȚIMEA C | MULȚIMEA B |
|--|------------|------------------------------|
| $A^6 = \{D_2, D_4\}$ <-asociere-> $C^6 = \{D_3, D_5\}$ | | $B^7 = \{D_2, D_3, D_4\}$ |
| $A^7 = \{D_2, D_5\}$ <-asociere-> $C^7 = \{D_3, D_4\}$ | | $B^8 = \{D_2, D_3, D_5\}$ |
| $A^8 = \{D_3, D_4\}$ <-asociere-> $C^8 = \{D_2, D_5\}$ | | $B^9 = \{D_2, D_4, D_5\}$ |
| $A^9 = \{D_3, D_5\}$ <-asociere-> $C^9 = \{D_2, D_4\}$ | | $B^{10} = \{D_3, D_4, D_5\}$ |

A șaptea trecere

Parcurgând lista elementelor mulțimii C de sus în jos, căutam primul element C_i care-l conține pe D_2 . Acesta este $C^8 = \{D_2, D_5\}$. Adăugăm elementului asociat lui C^8 și anume A^8 , pe D_2 . După ordonarea elementului A^8 extins obținem $\{D_2, D_3, D_4\}$ care este chiar B^7 . Astfel, am stabilit o corespondență între elementul A^8 aparținând mulțimii A și elementul B^7 aparținând mulțimii B.

Spunem că am remapat un domeniu al discului defect D_6 și anume acela care intră în componența modelului de grup de paritate asociat lui A^8 neextins, într-unul din domeniile de rezervă ale discului D_2 , formându-se cu acest prilej un model de grup de paritate reprezentat de B^7 .

Concret, domeniul R din D_6 se remapează într-unul din domeniile de rezervă ale discului D_2 , după această remapare aria arătând ca în Figura 4-15. Aplicând în continuare pasul 3 al algoritmului, nu vom suprima nimic.

În finalul acestei treceri, vom memora corespondența stabilită și eliminăm elementele-mulțime între care s-a stabilit corespondența,

precum și pe cel corespunzător din mulțimea C. Eliminăm pe A^8 din mulțimea A, C^8 asociat lui A^8 din mulțimea C, B^7 din mulțimea B.

Noua structură a mulțimilor A, B și C:

| MULȚIMEA A | MULȚIMEA C | MULȚIMEA B |
|------------|------------|------------|
|------------|------------|------------|

| | | |
|----------------------|-----------------|---------------------------|
| $A^6 = \{D_2, D_4\}$ | $C^6 = \{D_5\}$ | $B^8 = \{D_2, D_3, D_5\}$ |
|----------------------|-----------------|---------------------------|

| | | |
|----------------------|----------------------|---------------------------|
| $A^7 = \{D_2, D_5\}$ | $C^7 = \{D_3, D_4\}$ | $B^9 = \{D_2, D_4, D_5\}$ |
|----------------------|----------------------|---------------------------|

| | | |
|----------------------|----------------------|------------------------------|
| $A^9 = \{D_3, D_5\}$ | $C^9 = \{D_2, D_4\}$ | $B^{10} = \{D_3, D_4, D_5\}$ |
|----------------------|----------------------|------------------------------|

| D1 | D2 | D3 | D4 | D5 | D6 |
|----|----|----|----|----|----|
| A | A | A | | | |
| B | B | | B | | |
| C | C | | | C | |
| E | | E | E | | |
| F | | F | | F | |
| H | | H | H | | |
| | K | K | K | | |
| | L | L | | L | |
| | N | | N | N | |
| | | Q | Q | Q | |
| D | D | | | D | |
| G | | G | G | | |
| I | I | | I | | |
| J | | J | | J | |
| M | M | M | | | |
| | O | | O | | O |
| | P | | | P | P |
| | R | R | R | | |
| | | S | | S | S |
| T | | | T | T | |
| | | | | | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |

Figura 4-15

| D1 | D2 | D3 | D4 | D5 | D6 |
|----|----|----|----|----|----|
| A | A | A | | | |
| B | B | | B | | |
| C | C | | | C | |
| E | | E | E | | |
| F | | F | | F | |
| H | | H | H | | |
| | K | K | K | | |
| | L | L | | L | |
| | N | | N | N | |
| | | Q | Q | Q | |
| D | D | | | D | |
| G | | G | G | | |
| I | I | | I | | |
| J | | J | | J | |
| M | M | M | | | |
| | O | | O | | O |
| | P | P | | P | |
| | R | R | R | | |
| | | S | | S | S |
| T | | | T | T | |
| | | | | | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |

Figura 4-16

| D1 | D2 | D3 | D4 | D5 | D6 |
|----|----|----|----|----|----|
| A | A | A | | | |
| B | B | | B | | |
| C | C | | | C | |
| E | | E | E | | |
| F | | F | | F | |
| H | | H | H | | |
| | K | K | K | | |
| | L | L | | L | |
| | N | | N | N | |
| | | Q | Q | Q | |
| D | D | | | D | |
| G | | G | G | | |
| I | I | | I | | |
| J | | J | | J | |
| M | M | M | | | |
| | O | | O | | O |
| | P | P | | P | |
| | R | R | R | | |
| | | S | S | S | |
| T | | | T | T | |
| | | | | | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |

Figura 4-17

A opta trecere

Parcurgând lista elementelor mulțimii C de sus în jos, căutam primul element C^i care-l conține pe D3. Acesta este $C^7 = \{D3, D4\}$. Adăugăm elementului asociat lui C^7 și anume A^7 , pe D3. După ordonarea elementului A^7 extins obținem $\{D2, D3, D5\}$ care este chiar B^8 . Astfel,

am stabilit o corespondență între elementul A^7 aparținând mulțimii A și elementul B^8 aparținând mulțimii B.

Spunem că am remapat un domeniu al discului defect D_6 și anume acela care intră în componența modelului de grup de paritate asociat lui A^7 neextins, într-unul din domeniile de rezervă ale discului D_3 , formându-se astfel un model de grup de paritate reprezentat de B^8 .

Concret, domeniul P din D_6 se remapează într-unul din domeniile de rezervă ale discului D_3 , după această remapare aria arătând ca în Figura 4-16. Aplicând în continuare pasul 3 al algoritmului, vom suprima elementul D_2 din mulțimea C^9 .

În finalul acestei treceri, vom memora corespondența stabilită și eliminăm elementele-mulțime între care s-a stabilit corespondența, precum și pe cel corespunzător din mulțimea C. Eliminăm pe A^7 din mulțimea A, C^7 asociat lui A^7 din mulțimea C, B^8 din mulțimea B.

Noua structură a mulțimilor A, B și C:

| MULȚIMEA A | MULȚIMEA C | MULȚIMEA B |
|--|------------------------------|------------|
| $A^6 = \{D_2, D_4\} < \text{asociere} \rightarrow C^6 = \{D_5\}$ | $B^9 = \{D_2, D_4, D_5\}$ | |
| $A^9 = \{D_3, D_5\} < \text{asociere} \rightarrow C^9 = \{D_4\}$ | $B^{10} = \{D_3, D_4, D_5\}$ | |

A noua trecere

Parcurgând lista elementelor mulțimii C de sus în jos, căutam primul element C_i care-l conține pe D_4 . Acesta este $C^9 = \{D_4\}$. Adăugăm elementului asociat lui C^9 , și anume A^9 , pe D_4 . După ordonarea elementului A^9 extins, obținem $\{D_3, D_4, D_5\}$ care este chiar B^{10} . Astfel, am stabilit o corespondență între elementul A^9 aparținând mulțimii A și elementul B^{10} aparținând mulțimii B.

Spunem că am remapat un domeniu al discului defect D_6 și anume acela care intră în componența modelului de grup de paritate asociat lui A^9 neextins, într-unul din domeniile de rezervă ale discului D_4 , formându-se astfel un model de grup de paritate reprezentat de B^{10} .

Concret, domeniul S din D₆ se remapează într-unul din domeniile de rezervă ale discului D₄, după această remapare aria arătând ca în Figura 4-17.

Aplicând în continuare pasul 3 al algoritmului, nu se modifică nimic.

În finalul acestei treceri, vom memora corespondența stabilită și eliminăm elementele-mulțime între care s-a stabilit corespondența, precum și pe cel corespunzător din mulțimea C. Eliminăm pe A⁹ din mulțimea A, C⁹ asociat lui A⁹ din mulțimea C, B¹⁰ din mulțimea B.

Ilustrăm noua structură a mulțimilor A, B și C:

MULȚIMEA A

MULȚIMEA C

MULȚIMEA B

A⁶={D₂,D₄} <-asociere->C⁶={D₅}

B⁹={D₂,D₄,D₅}

A zecea trecere (ultima)

Parcurgând lista elementelor mulțimii C de sus în jos, căutam primul element C_i care-l conține pe D₅. Acesta este C⁶={D₅}. Adăugăm elementului asociat lui C⁶, și anume A⁶, pe D₅. După ordonarea elementului A⁶ extins, obținem {D₂, D₄, D₅} care este chiar B⁹. Astfel, am stabilit o corespondență între elementul A⁶ aparținând mulțimii A și elementul B⁹ aparținând mulțimii B.

Spunem că am remapat un domeniu al discului defect D₆ și anume acela care intră în componența modelului de grup de paritate asociat lui A⁶ neextins, într-unul din domeniile de rezervă ale discului D₅, formându-se cu acest prilej un model de grup de paritate reprezentat de B⁹.

Concret, domeniul O din D₆ se remapează într-unul din domeniile de rezervă ale discului D₅, după această remapare aria arătând ca în Figura 4-18, ceea ce reprezintă chiar starea matricei după remaparea tuturor domeniilor discului defect (după defectarea primului disc). Aplicând în continuare pasul 3 al algoritmului, nu se modifica nimic.

În finalul acestei treceri, vom memora corespondența stabilită și eliminăm elementele-mulțime între care s-a stabilit corespondența, precum și pe cel corespunzător din mulțimea C. Eliminăm pe A⁶ din mulțimea A, C⁶ asociat lui A⁶ din mulțimea C, B⁹ din mulțimea B.

Noua structură a mulțimilor A, B și C:

Mulțimea A Mulțimea C Mulțimea B

Figura 4-18

Se observă că sfârșitul algoritmului găsește cele trei mulțimi A, C și B fără nici un element, ceea ce înseamnă că toate domeniile discului defect D₆ au fost remapate. De asemenea, constatăm că noua matrice este aranjată tot într-o structură UPGD/DS, bazată pe enumerarea tuturor combinațiilor de N=5 discuri luate câte k=3, replicate de 2*x=2 ori:

- grupul DDD este o replică a grupului CCC,
- grupul GGG este o replică a grupului EEE,
- grupul III este o replică a grupului BBB,

și așa mai departe,

- grupul TTT fiind o replică a grupului HHH.

Cum numărul de replici este divizibil cu N-k=2, în urma defectării celui de-al doilea disc, reconstrucția domeniilor acestuia în spațiul de rezervă al celorlaltor discuri rămase în matrice se face mai simplu, pentru că se îndeplinesc condițiile cerute de algoritmul de la capitolul 4 și, deci, se va lucra cu acest algoritm.

Starea matricei după cea de-a doua remapare se prezintă în Figura 4-19.

| D1 | D2 | D3 | D4 | D5 | D6 |
|----|----|----|----|----|----|
| A | A | A | | | |
| B | B | | B | | |
| C | C | | | C | |
| E | | E | E | | |
| F | | F | | F | |
| H | | H | H | | |
| | K | K | K | | |
| | L | L | | L | |
| | N | | N | N | |
| | | Q | Q | Q | |
| D | D | | | D | |
| G | | G | G | | |
| I | I | | I | | |
| J | | J | | J | |
| M | M | M | | | |
| | O | | O | O | |
| | P | P | | P | |
| | R | R | R | | |
| | | S | S | S | |
| T | | | T | T | |
| | | | | | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |
| □ | □ | □ | □ | □ | |

Figura 4-19

| D1 | D2 | D3 | D4 |
|----|----|----|----|
| A | A | A | |
| B | B | | B |
| C | C | C | |
| E | | E | E |
| F | F | F | |
| H | | H | H |
| | K | K | K |
| L | L | L | |
| N | N | | N |
| Q | | Q | Q |
| D | D | | D |
| G | | G | G |
| I | I | | I |
| J | | J | J |
| M | M | M | |
| | O | O | O |
| | P | P | P |
| | R | R | R |
| | S | S | S |
| T | | T | T |
| | | | |
| | | | |
| | | | |
| | | | |
| | | | |

Figura 4-20

Se observă că ea este aranjată într-o structură UPGD bazată pe enumerarea tuturor combinațiilor de $n-2=4$ discuri luate câte $k=3$ (este doar UPGD pentru că nu mai există spațiu de rezervă), replicate de cinci ori.

Am realizat astfel refacerea informației pe cele 4 discuri rămase.

4.8 RAID tolerant la mai multe defecte

Voi prezenta, în continuare o metodă de proiectare a matricelor de discuri tolerante la mai multe defecte – RAID – MFT.

Modelele RAID prezentate anterior sunt bazate pe multiplicarea grupurilor de paritate.

Mai mult, în algoritmul prezentat anterior se impun restricțiile $\{n=2*k; f=2\}$, unde n este numărul de discuri și f numărul de căderi tolerate.

Matricea de discuri va exista ca sistem UPGD/DS cu $N+1$ discuri, și va putea suporta I căderi, fără intervenția omului. De fiecare dată informația de pe discul căzut va fi reconstruită pe celelalte tot într-o matrice UPGD/DS cu excepția ultimei căderi, când aria va fi UPGD fără spațiu liber. Totuși și-n această situație aria noastră va fi tolerantă la defect singular, însă imediat după aceea discul va trebui înlocuit.

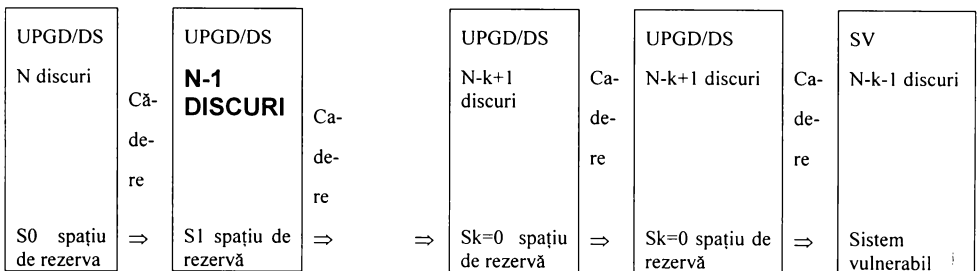


Figura 4-21

4.8.1 Algoritmul de creare UPGD/DS- MFT

Fie o UPGD/DS caracterizată prin parametrii:

- k - dimensiunea grupului de paritate;
- n - numărul de discuri;

- b - numărul de grupuri de paritate din matrice;
- c - numărul de blocuri în care se divide discul;
- r - numărul de blocuri folosite pe un disc;
- l - numărul de grupuri comune pe oricare 2 discuri;
- s - numărul de blocuri goale pe un disc (pentru DS).

Pentru a determina configurația matricei inițiale, vom pleca de la ultima, cea la care va ajunge și sistemul după I căderi (fără spațiu liber), toleranța la defect singular. Voi considera drept exemplu un RAID cu 3 discuri cu protecția bazată pe 4 grupuri de paritate uniform distribuite,

| D1 | D2 | D3 |
|----|----|----|
| A | A | A |
| B | B | B |
| C | C | C |
| D | D | D |

Figura 4-22

| D1 | D2 | D3 | D4 |
|----|----|----|----|
| A | | A | A |
| B | B | | B |
| C | C | C | |
| | D | D | D |

Figura 4-23

| D1 | D2 | D3 | D4 |
|----|----|----|----|
| A | B | C | D |
| B | C | D | A |
| C | D | A | B |
| | | | |

Figura 4-24

unde $\{A, A, A\}$ grup de paritate etc.

Adăugând un disc în sistem și împărțind informația uniform pe cele 4 discuri obținem o matrice cu spațiu liber, având parametrii (Figura 4-23): $n=4, k=3, b=4, r=3, l=2$

Noua configurație este o UPGD/DS pentru ca: $r * n = b * k = 12$, $l * (n-1) = (k-1) * r = 6$ și are spațiu liber $s = 1$.

Am extins aria prin adăugarea unui disc.

Spațiul liber de pe discuri se va ocupa în cazul defectării. Presupunând că discul 1 cade, reconstituirea are loc cum se arată în Figura 4-24. Nu are importanță unde este plasat pe disc spațiul liber. Se mai adaugă un disc.

$n = 5$

| D1 | D2 | D3 | D4 | D5 |
|----|----|----|----|----|
| A | B | C | D | |
| B | C | D | A | |
| C | D | A | B | |
| | | | | |

Figura 4-25

Dimensiunea grupului de paritate rămâne aceeași. Se va modifica însă b , numărul de grupuri. Deoarece aceste grupuri vor trebui împărțite uniform pe 5 discuri, considerăm că fiecare bloc este împărțit în 5 blocuri mai mici. Rezultă o matrice UPGD/DS cu 5 discuri și 20 de grupuri de paritate (Figura 4-26).

| | | | | |
|----|----|----|----|----|
| D1 | D2 | D3 | D4 | D5 |
| | | A1 | A1 | A1 |
| A2 | | A2 | A2 | |
| A3 | | | A3 | A3 |
| A4 | | A4 | | A4 |
| A5 | | A5 | A5 | |
| | B1 | | B1 | B1 |
| B2 | | | B2 | B2 |
| B3 | B3 | | B3 | |
| B4 | B4 | | | B4 |
| B5 | B5 | | B5 | |
| | C1 | C1 | | C1 |
| C2 | | C2 | | C2 |
| C3 | C3 | | | C3 |
| C4 | C4 | C4 | | |
| C5 | C5 | C5 | | |
| | D1 | D1 | D1 | |
| | | D2 | D2 | D2 |
| | D3 | | D3 | D3 |
| | D4 | D4 | | D4 |
| | D5 | D5 | D5 | |

Figura 4-26

Se observă că s-au pus pe discul 5 blocurile cu indicele 1 de pe D1, cele cu indicele 2 de pe D2, cele cu indicele 3 de pe D3 și cele cu 4 de pe D4.

O rearanjare duce la evidențierea spațiului liber, la o cădere de disc refacerea necesitând doar 3 locații pe fiecare disc, restul de $s' = 5$ rămânând disponibile pentru următoarea cădere:

| D1 | D2 | D3 | D4 | D5 |
|----|----|----|----|----|
| A2 | B1 | C1 | D1 | A1 |
| A3 | B3 | C2 | D2 | B1 |
| A4 | B4 | C4 | D3 | C1 |
| A5 | B5 | C5 | D5 | B2 |
| B2 | C1 | D1 | A1 | C2 |
| B3 | C3 | D2 | A2 | D2 |
| B4 | C4 | D4 | A3 | C3 |
| B5 | C5 | D5 | A5 | D3 |
| C2 | D1 | A1 | B1 | A3 |
| C3 | D3 | A2 | B2 | D4 |
| C4 | D4 | A4 | B3 | A4 |
| C5 | D5 | A5 | B5 | B4 |
| | | | | |
| | | | | |
| | | | | |
| S' | S' | S' | S' | S' |

Figura 4-27

Sistemul este UPGD/DS cu parametrii $n = 5k = 3$, $b = 20r = 12$, $l = 6s = 8$ (Figura 4-27).

Extindem numărul de discuri la 6. Se remarcă faptul că nu este necesară mărirea lui b (Figura 4-28).

| | | | | | |
|----|----|----|----|----|----|
| D1 | D2 | D3 | D4 | D5 | D6 |
| | | A1 | A1 | A1 | |
| A2 | | A2 | A2 | | |
| A3 | | | A3 | A3 | |
| A4 | | A4 | | A4 | |
| A5 | | A5 | A5 | | |
| | B1 | | B1 | B1 | |
| B2 | | | B2 | B2 | |
| B3 | B3 | | B3 | | |
| B4 | B4 | | | B4 | |
| B5 | B5 | | B5 | | |
| | C1 | C1 | | C1 | |
| C2 | | C2 | | C2 | |
| C3 | C3 | | | C3 | |
| C4 | C4 | C4 | | | |
| C5 | C5 | C5 | | | |
| | D1 | D1 | D1 | | |
| | | D2 | D2 | D2 | |
| | D3 | | D3 | D3 | |
| | D4 | D4 | | D4 | |
| | D5 | D5 | D5 | | |

Figura 4-28

Notăm altfel grupurile, de la 1 la 20 și prezentăm sistemul extins la 6 discuri. Sistemul tolerează 3+1 căderi, ultima fiind fatală (Figura 4-29).

| D1 | D2 | D3 | D4 | D5 | D6 |
|----|----|----|----|----|----|
| 1 | 1 | | 1 | | |
| 2 | 2 | | 2 | | |
| 3 | 3 | | | 3 | |
| 4 | 4 | | | 4 | |
| 5 | | 5 | | 5 | |
| 6 | | 6 | | 6 | |
| 7 | | 7 | | | 7 |
| 8 | | 8 | | | 8 |
| 9 | | | 9 | | 9 |
| 10 | | | 10 | | 10 |
| | 11 | 11 | | | 11 |
| | 12 | 12 | | | 12 |
| | 13 | 13 | 13 | | |
| | 14 | 14 | 14 | | |
| | 15 | | | 15 | 15 |
| | 16 | | | 16 | 16 |
| | | 17 | 17 | 17 | |
| | | 18 | 18 | 18 | |
| | | | 19 | 19 | 19 |
| | | | 20 | 20 | 20 |

Figura 4-29

| BAD | D2 | D3 | D4 | D5 | D6 |
|-----|----|----|----|----|----|
| | 1 | | 1 | | X |
| | 2 | | 2 | | X |
| | 3 | | X | 3 | |
| | 4 | | X | 4 | |
| | X | 5 | | 5 | |
| | X | 6 | | 6 | |
| | | 7 | | X | 7 |
| | | 8 | | X | 8 |
| | | X | 9 | | 9 |
| | | X | 10 | | 10 |
| | 11 | 11 | | | 11 |
| | 12 | 12 | | | 12 |
| | 13 | 13 | 13 | | |
| | 14 | 14 | 14 | | |
| | 15 | | | 15 | 15 |
| | 16 | | | 16 | 16 |
| | | 17 | 17 | 17 | |
| | | 18 | 18 | 18 | |
| | | | 19 | 19 | 19 |
| | | | 20 | 20 | 20 |

Figura 4-30

| BAD | D2 | D3 | D4 | D5 | BAD |
|-----|----|----|----|----|-----|
| | 1 | | 1 | X | |
| | 2 | | 2 | X | |
| | 3 | | 3 | 3 | |
| | 4 | | 4 | 4 | |
| | 5 | 5 | | 5 | |
| | 6 | 6 | | 6 | |
| | X | 7 | | 7 | |
| | X | 8 | | 8 | |
| | X | 9 | 9 | | |
| | | 10 | 10 | X | |
| | 11 | 11 | X | | |
| | 12 | 12 | X | | |
| | 13 | 13 | 13 | | |
| | 14 | 14 | 14 | | |
| | 15 | X | | 15 | |
| | 16 | | X | 16 | |
| | | 17 | 17 | 17 | |
| | | 18 | 18 | 18 | |
| | | X | 19 | 19 | |
| | | X | 20 | 20 | |

Figura 4-31

| BAD | BAD | D3 | D4 | D5 | BAD |
|-----|-----|----|----|----|-----|
| | | X | 1 | 1 | |
| | | X | 2 | 2 | |
| | | X | 3 | 3 | |
| | | X | 4 | 4 | |
| | | 5 | X | 5 | |
| | | 6 | X | 6 | |
| | | 7 | X | 7 | |
| | | 8 | X | 8 | |
| | | 9 | 9 | X | |
| | | 10 | 10 | 10 | |
| | | 11 | 11 | X | |
| | | 12 | 12 | X | |
| | | 13 | 13 | X | |
| | | 14 | 14 | X | |
| | | 15 | X | 15 | |
| | | X | 16 | 16 | |
| | | 17 | 17 | 17 | |
| | | 18 | 18 | 18 | |
| | | 19 | 19 | 19 | |
| | | 20 | 20 | 20 | |

Figura 4-32

În Figura 4-30, Figura 4-31 și Figura 4-32 se prezintă o modalitate de supraviețuire a matricei după căderea consecutivă a discurilor 1, 6 și 2.

A căzut discul 1, apoi 6, apoi 2, refacerea căsuțelor marcate cu X făcându-se la fiecare etapă.

4.8.2 Generalizare

O matrice pentru a fi UPGD/DS trebuie să îndeplinească cerințele următoare:

- $r * n = b * k = 12$

$$\bullet \quad l * (n-1) = (k-1) * r = 6$$

Există și cerințele k mic și b minim. Evident $k \leq n$.

Numărul total de blocuri este $c * n$ (numărul de blocuri/disc * numărul de discuri).

Fie N inițial. Adăugăm un disc. Datele se distribuie uniform pe $n=N+1$ discuri în noua configurație. Se procedează la împărțirea fiecărui disc în $N+1=n$ părți egale. Pe fiecare disc va fi informație $r=c * ((n-1)/n)$, iar spațiul liber va fi $s=c * (1/n)$.

$$\text{Cum } c.m.m.d.c.(n-1, n) = 1 \quad \Rightarrow \quad n \mid c$$

Extindem aria la $N+2$ discuri, distribuind informația. Rezultă pe fiecare disc avem $r=c * ((n-2)/n)$, spațiul liber fiind $s=c * (2/n)$, de unde $n \mid c * 2$ și $n \mid c * (n-2)$.

Prin inducție după i , extinzând aria la $N+i$ discuri, vom avea $c * N$ informație pe $N+i$ discuri, $r=c * (N/(N+i))$, $s=c * (i/(N+i))$, rezultând condițiile: $n \mid c * (n-i)$ și $n \mid c * i$

De unde se deduc:

- b numărul de grupuri în matrice
- r numărul de blocuri ocupate cu informație
- s numărul de blocuri goale

Acest sistem RAID este tolerant la $i+1$ defecte, dar nu la 2 simultane.

Suportă i defecte deoarece are spațiu liber, iar al $i+1$ -lea este tolerat datorită parității. Abia după aceasta este necesară intervenția omului pentru înlocuirea unui disc defect.

Atunci când cade un disc (sistemul trece de la $N+i$ la $N+i-1$ discuri), toată informația de pe acel disc este împărțită egal între cele rămase. Este ca și cum discul căzut ar fi copiat în spațiul distribuit pe celelalte, însă el nu este copiat ci este refăcut cu ajutorul parității.

În tabelul de mai jos se dau câteva valori obținute prin rularea unui program care calculează numărul de grupuri de paritate b , numărul de blocuri c pe un disc, spațiul folosit pe discuri r , spațiul liber s . Are ca parametri de intrare dimensiunea grupului de paritate k , numărul de discuri după căderi și numărul de căderi f la care poate rezista.

Ultimul algoritm prezentat dă rezultate bune când k este mic (3-5) sau apropiat de n și se obțin valori rezonabile pentru serii de până la 5-6 căderi.

| Valori de intrare | | | Valori obținute | | | | | | |
|-------------------|---|---|-----------------|----|------|------|------|-----|-----|
| K | N | F | f | n | b | c | r | l | s |
| 3 | 3 | 7 | 1 | 3 | 1 | 1 | 1 | 1 | 0 |
| | | | 2 | 4 | 4 | 4 | 3 | 2 | 1 |
| | | | 3 | 5 | 20 | 20 | 12 | 6 | 8 |
| | | | 4 | 6 | 20 | 20 | 10 | 4 | 10 |
| | | | 5 | 7 | 140 | 140 | 60 | 20 | 80 |
| | | | 6 | 8 | 280 | 280 | 105 | 30 | 175 |
| | | | 7 | 9 | 840 | 840 | 280 | 70 | 560 |
| 3 | 4 | 5 | 1 | 4 | 4 | 3 | 6 | 3 | 0 |
| | | | 2 | 5 | 20 | 15 | 5 | 2 | 1 |
| | | | 3 | 6 | 20 | 15 | 30 | 10 | 12 |
| | | | 4 | 7 | 140 | 105 | 105 | 30 | 63 |
| | | | 5 | 8 | 280 | 210 | 280 | 70 | 224 |
| 3 | 6 | 5 | 1 | 6 | 20 | 10 | 10 | 4 | 0 |
| | | | 2 | 7 | 140 | 70 | 60 | 20 | 10 |
| | | | 3 | 8 | 280 | 140 | 105 | 30 | 35 |
| | | | 4 | 9 | 840 | 420 | 280 | 70 | 140 |
| | | | 5 | 10 | 840 | 420 | 252 | 56 | 168 |
| 3 | 8 | 4 | 1 | 8 | 280 | 105 | 105 | 30 | 0 |
| | | | 2 | 9 | 840 | 315 | 280 | 70 | 35 |
| | | | 3 | 10 | 840 | 315 | 252 | 56 | 63 |
| | | | 4 | 11 | 9240 | 3465 | 2520 | 504 | 945 |
| 4 | 5 | 4 | 1 | 5 | 5 | 4 | 4 | 3 | 0 |
| | | | 2 | 6 | 15 | 12 | 10 | 6 | 2 |
| | | | 3 | 7 | 105 | 84 | 60 | 30 | 24 |
| | | | 4 | 8 | 210 | 168 | 105 | 45 | 63 |
| 4 | 6 | 4 | 1 | 6 | 15 | 10 | 10 | 3 | 0 |
| | | | 2 | 7 | 105 | 70 | 60 | 6 | 10 |
| | | | 3 | 8 | 210 | 140 | 105 | 30 | 35 |
| | | | 4 | 9 | 630 | 420 | 280 | 45 | 140 |

| K | N | F | f | n | b | c | r | l | s |
|----------|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| 6 | 6 | 4 | 1 | 6 | 1 | 1 | 1 | 6 | 0 |
| | | | 2 | 7 | 7 | 7 | 6 | 30 | 1 |
| | | | 3 | 8 | 28 | 28 | 21 | 45 | 7 |
| | | | 4 | 9 | 84 | 84 | 56 | 105 | 28 |

Figura 4-33

Valorile obținute pentru parametrii matricelor extinse sunt foarte bune, însă nu toate soluțiile se pot aplica. De altfel, pentru unele date de intrare, nu se pot obține soluții, deci k și N au anumite restricții.

4.9 Concluzii

Principalele contribuții aduse în acest capitol se referă la:

1. definire unei noi clase de structuri RAID (RAID 7) cu performanțe deosebite referitoare la timpii de regenerare a datelor și reconstrucție a matricei;
2. elaborarea a trei algoritmi performanți pentru construcția matricelor RAID 7.

5 RAIC – MATRICE DE CANALE REDUNDANTE ȘI INDEPENDENTE

5.1 Introducere

Principala activitate (din punct de vedere al consumului de timp) desfășurată de un sistem de calcul o reprezintă transferul de date, astfel în acest domeniu vom regăsi principalele surse de erori precum și posibilitățile cele mai mari de creștere a performanțelor. Chiar dacă sistemul este proiectat pentru înalta fiabilitate, el este expus la erori generate de evenimente atmosferice (umiditate, temperatură etc.), zgomot electromagnetic, căderi *hardware* și *software*. În aceste condiții, este foarte important ca sistemul să poată detecta erorile și să fie capabil să execute un set de acțiuni în vederea eliminării efectelor acestora.

În continuare, vom presupune că în orice transfer de informații sunt implicate patru entități, și anume:

- Emițătorul sau sursa de informații;
- Canalul sau mediul de comunicații;
- Receptorul sau destinația datelor;
- Sursa de erori.

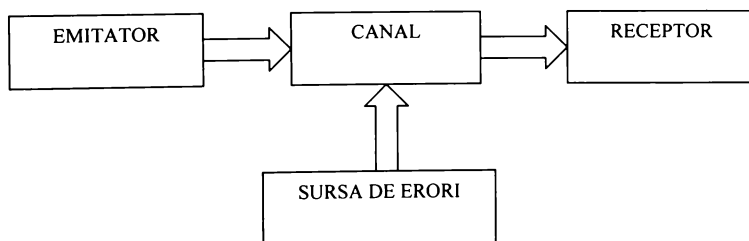


Figura 5-1

Structura ideală a suportului pentru un transfer de date

Deși pot să apară erori la oricare din cele trei niveluri, în continuare ne vom concentra atenția doar asupra celor referitoare la canale.

În prezenta lucrare voi considera următoarea clasificare a canalelor din punct de vedere al transformărilor logice suferite de date în timpul trecerii prin canale:

Canale cu transfer identic furnizează receptorului datele exact în aceeași formă ca la intrare, instantaneu (fără întârzieri).

Canalele procesor modifică datele în timpul transferului și le furnizează instantaneu receptorului.

Canalele cu memorie furnizează receptorului datele exact în aceeași formă ca la intrare, dar cu o întârziere controlată.

În realitate, nu există canale care să se încadreze strict într-una din cele trei categorii. De obicei, ele sunt o combinație ponderată a celor trei, un exemplu fiind prezentat în Figura 5-2.

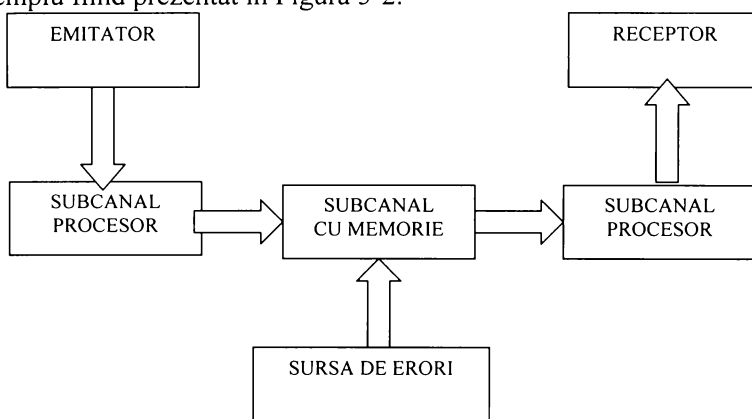


Figura 5-2

Fiecare subcanal, pe lângă funcția sa de transfer, va prezenta un timp de propagare. De exemplu, subcanalul procesor Δt_{pi} , iar subcanalul cu memorie Δt_m . Particularizând, funcțiile de transfer și timpii de propagare se pot practic obține orice tip de canale.

Exemplu

Pentru un canal de tip RS 232C putem identifica:

Subcanalul procesor 1 care se ocupă de calcularea bitului de paritate și serializarea datelor (Δt_{p1})

Subcanalul cu memorie constituit din emițătoare, receptoare și linia pe care se execută comunicația propriu-zisă, având toate împreună timpul de propagare Δt_m .

Subcanalul procesor 2 care restabilește data transmisă, recalculează bitul de paritate și verifică integritatea datei recepționate (Δt_{pi}).

Bazat pe cele prezentate anterior, devine evident că un disc magnetic poate fi considerat un canal de comunicații, conectat la un receptor și emițător, grupate în aceeași unitate logică.

Generalizarea poate merge mai departe, considerând sistemele de calcul drept un set de emițătoare, receptoare și canale de comunicații, fiecare activitate fiind de fapt un transfer de date. În concluzie, problemele legate de creșterea fiabilității și a performanțelor pot fi abordate într-un mod similar cu cele legate de transferurile de date.

Rata de transfer este cel mai important parametru de performanță pentru un canal de comunicații, dar ea este limitată drastic de tehnologia utilizată. Rezultă că, pentru a depăși aceste limite, singurele soluții sunt cele structurale, care presupun utilizarea în paralel a mai multor seturi de canale, inclusiv emițătoarele și receptoarele aferente. Creșterea numărului de componente va conduce, însă, la o scădere atât a fiabilității cât și a disponibilității. Pentru a depăși aceste probleme suplimentare propun în continuare utilizarea unor soluții bazate pe redundanță și codificare.

Pentru transferuri în care datele care pot fi divizate în pachete sau blocuri protejate prin coduri detectoare de erori voi prezenta în continuare o soluție originală de creștere a fiabilității și performanțelor transferurilor de date. Această soluție am denumit-o **RAIC (Redundant Array of Independent Channels)**.

Pentru a crește rata de transfer, datele pot fi distribuite simultan prin mai multe canale independente. Pentru un număr n rezonabil de astfel de canale se poate obține o rată globală de transfer de n ori mai mare decât în cazul utilizării unui singur canal. Desigur un număr de canale pot fi rezervate pentru redundanța în scopul acoperirii unor eventuale erori.

În Figura 5-3 prezint o sugestie de utilizare a unor coduri pentru câteva tipuri particulare de canale:

| TIP DE CANAL | COD RECOMANDAT |
|---|--|
| MEMORII RAPIDE (la nivel de bit) | HSIO (SEC-DED) HAMMING MODIFICAT (SEC-DED) IMAI- KAMIYANAGI HORIGUCHI- MORITA PARITATE INCRUCISATA |
| MEMORII RAPIDE (la nivel de cuvânt) | HAMMING BURTON FUJIWARA HONG-PATEL REED-SOLOMON KANEDA-FUJIWARA CHEN |
| DISCURI | FIRE REED-SOLOMON INTRETESUT (SI INCRUCISAT) |
| MEMORII CU BULE | HAMMING FIRE |
| PROCESOR | PARITATE BERGER |

Figura 5-3

Din păcate, RAIC nu reduce costul transferurilor de date, dar asigură o creștere spectaculoasă a performanței și a fiabilității acestora.

În continuare, voi analiza principalele probleme legate de **RAIC**, cum ar fi: organizarea și distribuția datelor, adresarea independentă, divizarea fină, divizarea grosieră, alegerea mărimii unității de divizare, mecanisme de redundanță, replicarea datelor, protecția bazată pe paritate.

5.2 Construirea sistemului RAIC

În cele ce urmează, vom explica cum putem genera octeții de control, precum și cum îi vom folosi pentru a implementa operația de siguranță la apariția unui defect în sistemul de intrare/ieșire. Să presupunem că avem trei canale și pe fiecare dintre ele transferăm un octet (adică opt biți) de date. Două vor transfera informație utilă, în timp ce al treilea va transfera un octet de control, generat cu ajutorul funcției SAU exclusiv (XOR).

În cazul în care canalul numărul 1 are transferat octetul 00110000, iar canalul numărul 2 are transferat octetul 00000011, atunci octetul de control transferat în canalul numărul 3 este generat realizând operația

XOR asupra celor doi octeți transferați în canalele 1 și 2. Operația executându-se la nivel de bit, rezultă octetul 00110011. Transferând acest octet de control prin cel de-al treilea canal, putem să reconstituim data transferată prin oricare dintre cele trei canale, folosindu-ne de informația transferată în celelalte două. Astfel, presupunând că eșuează canalul numărul 1, data de pe acesta se reconstruiește realizând un SAU exclusiv între octeții dispozitivelor 2 și 3:

```

00000011  CANAL 2 (CANAL PENTRU DATE)
XOR 00110011  CANAL 3 (CANAL PENTRU PARITATE)
00110000  CANAL 1 (DATĂ RECONSTRUITĂ)

```

Această tehnică se aplică pentru orice dimensiune a blocurilor, în mod identic. De asemenea, numărul de canale care intră în componența sistemului nu este limitat din punct de vedere teoretic.

Această caracteristică de siguranță la defect se pretează sistemelor cu misiuni critice. Chiar dacă un canal din set se defectează, sistemul continuă să fie operațional. Oricum, canalul defect va putea fi înlocuit imediat.

În cele ce urmează voi folosi următoarea convenție:

- scriere – operație de transfer la nivelul emițătorului;
- citire – operație de transfer la nivelul receptorului.

5.3 Aranjamentul octeților de control

Până acum am arătat cum oferă un sistem RAIC o performanță ridicată împărțind datele, modul în care disponibilitatea datelor este îmbunătățită utilizând funcția XOR pentru generarea octeților de control, precum și modul în care se reconstruiește informația unui canal care a eșuat. Pentru a completa această analiză, iată cum se poate aranja informația utilă pe canale, pentru obținerea unei funcționări mai eficiente.

Pentru a avea tot timpul actualizată informația de control, octeții de control trebuie să fie modificați de fiecare dată când se transferă ceva pe oricare dintre canale. Dacă proiectăm sistemul în așa fel încât toți octeții de control să fie transmiși pe același canal, se poate crea o gâtuire a canalului, rezultând o scădere a performanței sistemului RAIC.

Pentru a evita această problemă este mult mai eficient să transmitem pe fiecare canal atât date utile, cât și octeți de control. Împrăștiind octeții de control peste toate canalele din sistem, se permit transferuri multiple simultane. Într-un sistem RAIC cu patru canale, acesta oferă o triplare a performanțelor, în ceea ce privește operațiile de transfer de blocuri mici de date, în comparație cu un sistem în care toată informația de control este transmisă pe un singur canal. Cu cât numărul de canale crește, performanța sistemului crește proporțional.

În continuare, prezint un exemplu în care am un sistem RAIC cu patru canale, pe fiecare canal transmițându-se patru blocuri (nu interesează mărimea blocului, dar ea trebuie să fie aceeași pentru toate cele patru canale). În total avem 16 blocuri, dintre care în 12 blocuri se află informație utilă, iar în 4 blocuri avem informație de control (notată $Cb(a-b-c)$, unde a , b și c sunt numerele blocurilor de date asociate acestui bloc de control):

| Canal 1 | Canal 2 | Canal 3 | Canal 4 |
|----------------|-------------|-------------|-------------|
| Bloc 1 | Bloc 2 | Bloc 3 | $Cb(1-2-3)$ |
| Bloc 5 | Bloc 6 | $Cb(4-5-6)$ | Bloc 4 |
| Bloc 9 | $Cb(7-8-9)$ | Bloc 7 | Bloc 8 |
| $Cb(10-11-12)$ | Bloc 10 | Bloc 11 | Bloc 12 |

Figura 5-4

5.4 Organizarea datelor

Organizarea datelor pentru matricele de canale poate fi partiționată în două componente ortogonale: schemele de distribuție a datelor și mecanismele de redundanță. Distribuția de date definește translația adreselor logice, vizibile extern, în locații de transfer din cadrul subsistemului. Sunt două căi pentru a face acest lucru: prin adresarea independentă a fiecărui canal sau prin divizarea canalului. Cea din urmă oferă, adesea, grade îmbunătățite de încărcare echilibrată, transfer paralel de date sau acces concurrent. Mecanismele de redundanță specifică tipul, scopul și locația oricărei informații redundante din matrice (de exemplu, replicarea datelor și paritatea divizată). Matricele redundante de canale continuă să ofere acces deplin la date și după

defectarea unui canal și în timp ce datele pierdute sunt reconstruite pe un canal de înlocuire sau din informațiile recepționate corect. Matricele de canale pot combina divizarea canalului cu redundanța pentru a oferi atât performanță cât și fiabilitate.

5.4.1 Distribuția datelor

Distribuția de date constă în maparea adreselor logice folosite de sistemul gazdă pe canalele din subsistem. Există, în principiu, trei metode principale folosite în realizarea acestei translații. Metoda convențională este de a adresa, independent, fiecare canal și de a mapa numerele de bloc logic în numere de bloc de canal, direct. Distribuția este făcută manual de către administratorii de sistem, programe de aplicații sau sisteme *software*. Divizarea de canal, numită și *channel striping* sau *channel interleaving*, reunește mai multe spații de adrese de canal într-un singur spațiu, unificat, văzut de sistemul gazdă. Aceasta se face prin distribuirea consecutivă de unități logice de date (unități de divizare) din cadrul canalului într-o manieră *round-robin*. Se disting între două tipuri diferite de divizări de canale: divizare fină, în care canalele cooperează toate pentru satisfacerea unei cereri și divizare grosieră, în care canalele cooperează la cereri mari de date și lucrează independent la cereri mici de date.

5.4.1.1 Adresarea independentă

Dacă fiecare canal este adresat independent, utilizatorii sau sistemele gazdă trebuie să distribuie explicit datele. În general, administratorii de sistem sunt responsabili pentru răspândirea datelor pe fiecare canal. Încărcarea echilibrată sau distribuirea datelor astfel încât să se încarce echilibrat sistemul, este mai mult o artă decât o știință și necesită specialişti. Adrese de date adiacente creează *hot spots* făcând *task*-ul complex.

| Canal 0 | Canal 0 | Canal 3 | Canal 4 |
|---------|---------|---------|---------|
| A0 | B0 | C0 | D0 |
| A1 | B1 | C1 | D1 |
| A2 | B2 | C2 | D2 |
| ... | ... | ... | ... |

Figura 5-5

5.4.1.1.1 Divizarea fină

Aici, toate cele n canale conțin o fracțiune din fiecare bloc accesibil (Figura 5-6). Numărul canalelor și mărimea unității de divizare sunt, în general, alese astfel încât produsul lor să dividă egal cea mai mică unitate de date accesibilă. Mărimi mai folosite sunt un bit, un octet, secțiune. Mărimea folosită nu este neapărat importantă deoarece fiecare unitate accesibilă este răspândită pe toate canalele. Încărcarea este echilibrată deoarece toate canalele primesc aceeași încărcare. Rata de transfer este de aproximativ n ori mai mare față de cea a unui canal individual deoarece fiecare canal transferă $1/n$ din datele cerute. De asemenea, doar o singură cerere poate fi făcută la un moment dat, deoarece toate cele n canale lucrează la aceeași cerere. Aceste limitări, de obicei, restrâng folosirea metodei la domenii în care timpii de transfer domină timpii în care nu se lucrează cu canalul.

| Canal 0 | Canal 0 | Canal 3 | Canal 4 |
|---------|---------|---------|---------|
| A0 | A0 | A0 | A0 |
| A1 | A1 | A1 | A1 |
| ... | ... | ... | ... |
| A7 | A7 | A7 | A7 |
| ... | ... | ... | ... |
| B0 | B0 | B0 | B0 |
| B1 | B1 | B1 | B1 |
| ... | ... | ... | ... |

Figura 5-6

5.4.1.1.2 Divizarea grosieră

Cu aceasta nu mai este necesar ca toate canalele să coopereze la fiecare cerere. Această tehnică exploatează paralelismul transferului de date pentru cereri mari, în timp ce permite canalelor separate să satisfacă cereri mici concurrent. Câteva efecte se combină pentru a oferi o încărcare echilibrată și automată. Primul efect este acela că un fișier va tinde să aibă blocurile sale componente împrăștiate pe mai multe canale. Divizarea grosieră, ca la *hashing* efectiv, randomizează primul canal accesat pentru fiecare cerere, celelalte canale fiind identificate în maniera *round-robin* (Figura 5-7). Această încărcare, echilibrată statistic, rămâne intactă și în cazul adreselor de date adiacente și distribuții rapid schimbătoare de accese. Din motivele de mai sus, metoda are posibilitatea să ofere performanță bună cu un minim de întreținere.

| Canal 0 | Canal 0 | Canal 3 | Canal 4 |
|---------|---------|---------|---------|
| A0 | A1 | A2 | A3 |
| A4 | A5 | A6 | A7 |
| ... | ... | ... | ... |
| B0 | B1 | B2 | B3 |
| ... | ... | ... | ... |

Figura 5-7

5.4.1.1.3 Alegerea mărimii unității de divizare

Realizarea performanței de la metoda de mai sus necesită alegerea corectă a mărimii unității de divizare. Această alegere determină, pe de o parte, numărul de canale peste care sunt răspândite datele fiecărei cereri. Mărimi mari ale unității permit canalelor separate să satisfacă cereri multiple (concurrent), rezultând adesea reducerea cozilor de așteptare și o rată de transfer mai mare. Pe de altă parte, unitățile mici determină mai multe canale să acceseze datele în paralel, reducând timpii de transfer pentru cereri individuale. Acest compromis între paralelismul transferului și concurența accesului este guvernată de mărimea unității de divizare.

Când avem cereri de lungime variabilă, fără o secvențializare sau localizare a adreselor, parametrul cheie este cantitatea de concurență din încărcare. Mărimea unității de divizare ar trebui să crească o dată cu rata de sosire a cererilor, deci descrescând numărul mediu de canale accesate pe cerere.

Pentru cereri constând în primul rând din cereri aliniate, de mărime fixă, mărimea unității de divizare poate fi aleasă astfel încât numărul de accese la canal, pe cerere, să fie constant. În particular, folosind orice multiplu de mărimea fixă a cererii, se permite fiecărei cereri să fie servită de către un singur canal.

Pentru încărcări cu un oarecare grad de aliniament sau localizare a adresei, echilibrarea trebuie, de asemenea, considerată. Când cererile sunt mari, echilibrarea nu afectează semnificativ alegerea mărimii unității de divizare. Pentru încărcări mici, unitățile mari exploatează

beneficiile timpului de acces datorate localizării și secvențialității adreselor.

Un compromis mai complex este cel pentru încărcări generale, adesea caracterizate prin cereri mici, cereri mari scurte, și adrese puternic adiacente. În aceste cazuri, abilitatea metodei divizării grosiere de a echilibra încărcarea se îmbunătățește pe măsură ce mărimea unității de divizare scade. Mărimi mai mari ale unității tind să sufere din cauza formării unor *hot spots* de scurtă durată. Sugerez ca unitatea de divizare trebuie să fie de 10 ori mai mare decât mărimea medie a cererii, rotunjită la un multiplu al mărimii blocului. Alegerea mărimii unității de divizare poate fi rafinată considerând informație adițională, cum ar fi distribuțiile de mărime a cererii și caracteristicile localizării.

O întrebare interesantă este dacă mărimea unității de divizare trebuie să fie multiplu de mărime a blocului. Pe de o parte, o cerere pentru unitatea de divizare completă are avantajul unui canal suportând acces cu latență zero. Pe de altă parte, mărimea blocului poate corespunde, în mod dificil, mărimii și aliniamentului cererii, rezultând cereri multicanal, excesive. Multe încărcări constau din cereri de blocuri de mărime specifică, să zicem 4k sau 8k. Cele mai multe unități de canal *state-of-the-art* folosesc transferuri pe zone multiple pentru a crește densitatea de transfer. Folosind aceste canale, potrivirea mărimii unității de divizare cu cea a blocului cere o mărime diferită pentru fiecare zonă de canal.

5.4.2 Mecanismele de redundanță

Adăugând mecanisme de redundanță la un sistem de transfer, invariabil îi crește costul, i se reduce capacitatea și/sau i se reduc performanțele. Din fericire, căderile individuale de canale sunt foarte rare, iar căderile de sistem sunt și mai rare. Timpul mediu până la defectare (MTBF) la un canal modern este de sute de mii de ore, altfel spus, zeci de ani.

Totuși, pe măsură ce numărul de canale crește, MTBF-ul unui sistem tipic neredundant se micșorează până la perioade de luni sau chiar săptămâni. Presupunând că MTBF-ul fiecărui canal este independent și exponențial distribuit, MTBF-ul unui grup este invers proporțional cu numărul canalelor. Aceasta amenință fezabilitatea sistemelor nonredundante de canale.

În plus, numărul de canale afectate de cădere crește dramatic o dată cu folosirea divizării canalelor. Un subsistem divizat poate pierde $1/n$ din fiecare fișier, față de $1/n$ fișiere din matricea de canale în cazul adresării independente. Un subsistem cu redundanță în organizarea datelor poate

supraviețuii căderii uneia sau mai multor componente. Totuși, la apariția unor noi căderi va rezulta inaccesibilitatea sau pierderea datelor. O dată cu creșterea subsistemelor de transfer, redundanța *on-line* va fi necesară pentru datele cruciale, devenind, probabil, un standard pentru transfer, în general.

5.4.2.1 Replicarea datelor

Cel mai simplu mod de a proteja datele este replicarea lor. Copii separate ale fiecărui bloc sunt transferate pe canale separate. Datele devin indisponibile doar dacă toate canalele, conținând o copie, cad. Această formă de redundanță, folosită de mai mulți ani, se numește *channel mirroring*, *channel duplexing*, *channel shadowing* sau, mai nou, RAIC 1. Costul replicării datelor este proporțional cu $c-1$ (c fiind numărul de canale din matrice), deoarece costul canalului domină prețul subsistemelor de transfer. Deoarece două copii oferă înaltă fiabilitate, alte copii, în plus, nu se folosesc. De fapt, costul mare al replicării obligă, adesea, la folosirea altor mecanisme de redundanță care nu dau același nivel de performanță și fiabilitate.

Utilizarea unor copii multiple ale datelor afectează performanța în diferite moduri. Primul este acela că scrierile trebuie făcute pe toate copiile. Depinzând de implementare, aceasta poate crea oportunități pentru pierderea datelor sau timpi crescuți de răspuns. Performanța la citire este, pe de altă parte, mai bună în acest caz. În primul rând, cererile de citire se satisfac simultan, crescând rata de transfer și scăzând timpii de așteptare. În plus, o cerere de citire poate fi programată pe canalul cu cel mai mic timp de acces. Această performanță bună la citire poate face raportul cost/performanță chiar mai mic decât în cazul sistemelor cu canale nonredundante, deși costul absolut este de c ori mai mare.

Metoda convențională de implementare a replicării datelor dublează fiecare canal (Figura 5-8). Fiecare set de c canale identice pot supraviețui la $c-1$ căderi fără pierdere de date. Într-o matrice cu seturi multiple replicate, pot apărea multe căderi, fără să cauzeze vreo pierdere de informații, atâta timp cât c din ele nu sunt din același set unic, replicat. Totuși, fiecare cădere de canal reduce performanțele unui set relativ la alte seturi. Acest dezechilibru poate duce la scăderea performanțelor generale ale sistemului, dacă un set căzut devine o gâtuire.

Alternativ, seturile se pot suprapune astfel încât această degradare de performanță este împărțită în mod egal de toate canalele rămase. Termenul *declustering* se referă la o metodă pentru combinarea seturilor multiple replicate. Metoda *chained declustering* se referă la divizarea

fiecărui canal în c secțiuni. Prima secțiune conține copia primară a unor informații, cu fiecare canal având date primare diferite (a se vedea Figura 5-9). Cu o aranjare inteligentă, această metodă poate menține o încărcare echilibrată după o cădere a unui canal, deși probabilitatea supraviețuirii la căderi multiple este redusă.

Al doilea tip de *declustering*, cunoscut și sub numele de *double declustering* (vezi Figura 5-10) partiționează, de asemenea, fiecare canal, dar divizează copiile nonprimare din toate canalele în loc să le păstreze în spații contigue. Aceasta echilibrează încărcarea suplimentară cauzată de o cădere de canal, dar reduce, în continuare, fiabilitatea (pierderea oricăror c canale rezultă într-o cădere critică).

| Canal 0 | Canal 1 | Canal 2 | Canal 3 |
|---------|---------|---------|---------|
| 0 | 0 | 2 | 2 |
| 1 | 1 | 3 | 3 |

Figura 5-8

| Canal 0 | Canal 1 | Canal 2 | Canal 3 |
|---------|---------|---------|---------|
| 0 | 1 | 2 | 3 |
| 3 | 0 | 1 | 2 |

Figura 5-9

| Canal 0 | Canal 1 | Canal 2 | Canal 3 |
|---------|---------|---------|---------|
| 0 | 1 | 2 | 3 |
| 1c | 0a | 0b | 0c |
| 2b | 2c | 1a | 1b |
| 3a | 3b | 3c | 2a |

Figura 5-10

5.4.2.2 Protecția bazată pe paritate

Proiectanții de subsisteme de canale se pot folosi de rezultatele studiilor extinse ale codurilor detectoare și corectoare de erori. Din teoria codurilor se știe că pierderea unui singur bit poate fi recuperată știind bitul de paritate, ceilalți biți de protejat și faptul că acel bit este eronat. Deoarece sistemele actuale de canale au o logică de control care, în general, detectează un procent semnificativ din căderi, aceasta poate fi folosită pentru protejarea unui set de canale față de erorile singulare. Fiecare bit de paritate este calculat din bitul de date asociat de pe fiecare canal protejat și transferat pe un canal separat. Distribuirea unor biți particulari de date dați de un bit de paritate pot varia de la schemă la schemă, așa cum poate varia și aranjarea parității. Capacitatea adițională cerută (un canal complet) pentru a menține redundanța prin intermediul parității este considerabil mai mică decât atunci când se folosește replicarea datelor. Deși costul este rezonabil, menținerea parității poate reduce în mod semnificativ performanța. Cererile de citire sunt efectuate ca în matrice nonredundante. Dacă se adaugă un canal pentru a oferi capacitate de transfer pentru paritate, el va genera încărcare în plus pentru cereri de scriere (în cele mai multe scheme de așezare a parității). Cererile de scriere necesită lucru suplimentar cu canalul pentru actualizarea parității, adesea ducând la o performanță scăzută. Dacă toți biții de date asociați cu un bit de paritate sunt actualizați de o cerere de scriere, noul bit de paritate este calculat și inclus ca o parte a cererii totale. Dacă doar un subset al biților de date este scris, actualizarea bitului de paritate necesită lucru suplimentar. Aceasta constă în citirea biților de pe canale pentru construirea noului bit de paritate. Avem două metode:

- **Read-modify-write**, unde noul bit de paritate este construit din valorile vechi și cele noi ale biților care se scriu;
- **Regenerate write**, în care se generează bitul nou de paritate din valorile tuturor biților (adică noile valori ale biților care se scriu și valorile curente ale biților neschimbați).

RMW necesită, mai întâi, accese de citire la toate canalele care se actualizează (incluzând canalul care conține paritatea), pe când R-W cere accese de citire la toate canalele care nu se actualizează. Pierderea performanței poate fi redusă prin combinarea dinamică a acestor două opțiuni. Datorită nevoii de a face citiri de pe canal înainte de actualizarea parității, performanța poate fi semnificativ mai mică față de un sistem neprotejat sau față de un sistem protejat prin replicarea datelor.

Reducerea cât mai mult a acestei degradări de performanță este, încă, un subiect aflat în cercetare.

| Canal 0 | Canal 1 | Canal 2 | Canal 3 |
|---------|---------|---------|---------|
| A | A | A | a |
| B | B | B | b |
| C | C | C | c |
| ... | ... | ... | ... |

Figura 5-11

| Canal 0 | Canal 1 | Canal 2 | Canal 3 |
|---------|---------|---------|---------|
| A | A | A | a |
| B | B | b | B |
| C | c | C | C |
| ... | ... | ... | ... |

Figura 5-12

| Canal 0 | Canal 1 | Canal 2 | Canal 3 |
|---------|---------|---------|---------|
| A | A | A | a |
| B | b | C | B |
| B | C | c | C |
| ... | ... | ... | ... |

Figura 5-13

Există trei scheme majore pentru împrăștierea informației de paritate în cadrul canalelor din subsistem. Metoda directă este de a plasa toți biții de paritate pe un canal dedicat de paritate (Figura 5-11). Deși simplifică

maparea adreselor logice în adrese de canal, fiecare scriere trebuie să actualizeze biții asociați de pe unicul canal de paritate. Matricele care folosesc divizarea fină sunt potrivite pentru această metodă, deoarece ele satisfac doar o cerere la un moment dat. Cu alte scheme de distribuție, totuși, canalele pot limita performanța, mai ales când vechea paritate este cerută să participe în operațiile de la RMW. Pentru evitarea acestei potențiale gâturi, informația de paritate poate fi divizată printre canale (vezi Figura 5-12). Folosind divizarea parității, matricea poate efectua mai multe actualizări în paralel.

O a treia opțiune propusă: *declustered parity*. Pentru înțelegerea acestei scheme ne putem imagina combinarea câtorva matrice mai mici, protejate de paritatea divizată. Datele și paritatea fiecărei matrice sunt distribuite în întregul set de canale disponibile (a se vedea Figura 5-13). Cu un algoritm inteligent de distribuție, impactul performanței datorat pierderii unui canal poate fi împărțit, în mod egal, de toate canalele rămase. Deși matricea poate supraviețui doar unei singure căderi, ar trebui să se refacă mai repede. Ca o alternativă, aceasta metodă se poate imagina ca o cale de a realiza fiabilitate mai bună și performanță în funcționarea la defecte prin creșterea cantității de date de paritate menținută.

5.4.2.3 Alte scheme

Menținerea parității protejează o matrice împotriva căderii unui singur canal. O matrice cu două copii ale tuturor datelor poate supraviețui până la o cădere pe perechea de copii. În ambele cazuri, datele pot fi pierdute dacă apare o a doua cădere, înainte ca prima să fie reparată. Oferirea protecției necondiționate împotriva căderilor multiple necesită redundanță suplimentară. În cele mai promițătoare metode se folosesc două coduri *Reed-Solomon*, corectoare de erori. Cu asemenea redundanță, o matrice poate reface toate datele după pierderea accesului la oricare două canale. Compromisurile referitoare la performanța și opțiunile de amplasare a redundanței se potrivesc pentru acele matrice cu protecție bazată pe paritate, în afară de faptul că doi biți sunt mai bine menținuți decât unul.

Alte scheme, cum ar fi paritatea multidimensională și codurile mai agresive, corectoare de erori, pot proteja împotriva unui număr mai mare de căderi. Totuși, costul și implicațiile de performanță, combinate cu faptul că fiabilitatea sistemelor de transfer crește tot mai mult, pot obliga folosirea lor.

5.5 Performanța versus fiabilitatea

Cererile de scriere pot duce la întârzieri mari între momentul inițierii lor și momentul scrierii efective a noii informații pe canal. Până când sistemul gazdă primește mesajul de satisfacere a cererii orice resurse de sistem nu pot fi reutilizate. De asemenea, procesele de aplicație pot aștepta pentru astfel de confirmări, putând rezulta cicluri inutile de procesor. Acceptând confirmarea cât mai repede posibil, gazda poate utiliza mai bine resursele sale și are, astfel, mai puține cicluri inutile de procesor, rezultând o performanță mai bună. Totuși, în cazul în care confirmarea cererii este dată înainte ca datele să fie scrise efectiv pe canal, sistemul este vulnerabil la coruperea sau pierderea informațiilor, datorită căderii tensiunii sau defectării componentelor. Când sistemul renunță la constrângerile de fiabilitate în scopul câștigării performanței, trebuie luate măsuri de precauție sau se acceptă consecințele.

5.5.1 Negocierea fiabilității în favoarea performanței

Un controler de matrice de canale cu un *cache* (sau *buffer*) poate semnaliza confirmarea gazdei de îndată ce ultimul octet a fost copiat din memoria principală. Astfel de scheme îmbunătățesc mult performanța văzută de gazdă la scriere, deși datele sunt vulnerabile până la scrierea pe canal. Managerul de canal, care este responsabil de determinarea ordinii de servire pentru accesele la canal aflate în așteptare, are, acum, o problemă suplimentară. Fiabilitatea cere managerului să expedieze accesele la canal care transferă datele din *cache* pe suportul fizic, în timp ce performanța obligă managerul să amâne această activitate de fundal.

În matricele care folosesc replicarea informațiilor, raportarea confirmării după ce doar o fracțiune din copii a fost actualizată poate și ea să îmbunătățească performanța. Datele sunt pierdute doar dacă apar căderi multiple înainte să se efectueze scrierile rămase. Pentru matricele protejate prin paritate, punând paritate în *cache* poate să se îmbunătățească performanța prin acumularea noii parități pentru scrieri multiple și efectuarea actualizării de paritate într-o singură scriere.

În schemele de paritate ce cer cicluri de actualizare, sistemul fiind vulnerabil la coruperea datelor dacă setul de scrieri pe canal nu se face simultan. În cazul în care cade tensiunea, la un moment nepotrivit, doar câteva cereri vor fi satisfăcute. Aceasta face ca paritatea să fie incorectă, ducând la coruperi ulterioare dacă sunt lăsate necorectate și apoi folosită la refacerea datelor. Pe de altă parte, impunerea scrierilor simultane necesită comportare sincronă din partea mai multor canale. Aceasta

poate produce programări neoptime ale cererilor și un timp de canal inutil crescut. Când controlerul își programează, independent, cererile individuale de scriere utilizarea canalelor poate fi maximizată.

5.5.2 Îmbunătățirea performanței și a fiabilității

Cei mai mulți utilizatori nu sunt deranjați de coruperea datelor, exceptând cazurile inacceptabile. Există câteva tehnici pentru îmbunătățirea performanței și a fiabilității. Ele se realizează prin reducerea latentelor la scriere în timp ce este garantată siguranța datelor vechi și se asigură securitatea datelor noi înainte de semnalizarea confirmării. Aceste scheme pot fi implementate independent de sistemul gazdă sau în cadrul matricei de canale. O soluție simplă este folosirea memoriei RAM *cache*, non-volatilă, pentru a scrie datele sau paritatea. Aceste memorii mențin datele chiar și după întreruperea tensiunii de alimentare. Această organizare ar trebui să ofere o fiabilitate echivalentă cu cea a unei scheme cu redundanță. Remaparea scrierii, o formă de paginare dublă, poate fi utilizată pentru multe organizări de date diferite. Această remapare constă în alterarea dinamică a mapării adreselor logice în cele fizice, astfel încât cererea de scriere de date este totdeauna plasată la locații noi. Prin menținerea copiei anterioare a blocului logic până la satisfacerea cererii, metoda îmbunătățește fiabilitatea. Desigur, mapările trebuie menținute fiabil pentru a preveni pierderea informațiilor. O metodă este transferarea informației în NVRAM. Altă metodă, mai complexă, dar mai puțin scumpă este să mărim fiecare bloc fizic de date cu adresa logică corespunzătoare și cu un indicator de timp, care pot fi folosite să reconstruiască mapările.

Algoritmul pentru alegerea locației disponibile trebuie să prevină consumarea tuturor locațiilor libere dintr-o anumită regiune și lăsarea altor regiuni nefolosite. De asemenea, în timp, datele secvențial logice vor fi disperate pe spațiul fizic.

A treia opțiune - *logging* - scrie informația despre actualizări într-o zonă de log înainte de actualizările propriu-zise. Scrierile log sunt secvențiale și au, de obicei, un timp de răspuns mai mic, mai ales când se folosește un canal log dedicat. După scrierea intrării de log, controlerul de matrice poate, în siguranță, să transmită un mesaj de terminare către gazdă și gazda accesează canalul în fundal. În caz de cădere, conținutul log-ului permite repetarea oricărei actualizări incomplete. Informația scrisă în log poate să conțină o descriere scurtă a actualizării, o copie a noii informații sau copii ale amândurora. O formă de verificare este folosită pentru a reduce timpul de refacere și permite refolosirea spațiului de log.

5.5.3 Tratarea căderilor

Prin menținerea redundanței, o matrice de canale poate supraviețui căderii unuia sau mai multor canale. În timp ce matricea continuă să ofere acces la toate datele, performanța ei este aproape întotdeauna afectată de căderi. În plus, fiabilitatea este redusă, deoarece căderile ulterioare pot fi critice. După cădere, datele pierdute trebuie să fie reconstruite de pe canalele rămase și eventual scrise pe un canal de înlocuire, proces numit refacere. Până când acest proces se termină, matricea este vulnerabilă. Canalele suplimentare, numite rezerve rapide, sunt adesea incluse în matrice ca să permită refacerii să înceapă imediat.

În cele mai multe scheme, pierderea accesului la canal reduce numărul de cereri concurente de citire care pot fi deservite. Matricele care folosesc replicarea de date primesc cereri ca de obicei, dar mai puține canale efectuează citiri. De asemenea, scrierea necesită actualizarea unei copii mai puțin. La matricele protejate prin paritate, o cerere de citire spre un canal defect cere o regenerare a datelor lipsă. Fiecare bit de date pierdut poate fi refăcut din bitul de paritate asociat și din ceilalți biți de date. Deci, lucrul cu cereri de citire necesită acces la toate canalele care conțin acești biți. Performanța la scriere este și ea afectată în matricele protejate de paritate. Datorită faptului că este afectată și vechea informație, *Regenerate-Write* trebuie folosit pentru toate actualizările de date pe canalul căzut și RMW pentru toate actualizările pe canalele rămase.

Cea mai simplă formă de refacere reconstruiește secvențial informațiile indisponibile și le scrie pe un canal de înlocuire. Într-o matrice replicată, datele pierdute sunt copiate de pe un alt canal conținând aceleași date. Într-una cu protejare bazată pe paritate datele trebuie refăcute din biții de paritate și cei de date. Dacă cerințele de fiabilitate prevăd ca timpul de reparare să fie minim, cererile gazdei pot fi anulate în timp ce refacerea are loc la viteză maximă. Altfel, refacerea și sistemul gazdă concurează pentru deservire, rezultând un alt compromis între fiabilitate și performanță. Acest compromis este, în particular, important în refacerea matricelor cu protecție bazată pe paritate, unde regenerarea datelor cere accese la mai multe canale. Alegerea unui algoritm de refacere, mărimea unei cereri individuale de refacere, și complexitatea circumsticienii de coordonare multicanal trebuie luate în considerație când se ajustează procesul de refacere.

Dacă subsistemele de transfer continuă să se dezvolte ca și până acum, nivelul inteligenței încorporate în diferitele lor componente va continua să crească. Comunicația între aceste componente va crește mult,

rezultând optimizări globale, un control mai distribuit, și noi algoritmi de distribuție care răspund dinamic la mereu schimbătorul flux de cereri.

5.6 Propunere de clasificarea pentru RAIC

În continuare, voi propune o clasificare a soluțiilor RAIC, pornind de la criteriile funcționale și de la ideea considerării RAID ca o subclasă RAIC.

Diferitele niveluri de RAIC vor oferi diferențe în performanță, cost, asigurarea integrității și a disponibilității datelor.

RAIC de nivel 0 - date distribuite fără utilizarea parității

RAID 0 presupune distribuirea datelor pe mai multe canale fără a utiliza nici un fel de redundanță.

RAIC de nivel 1 – date duplicate

În RAIC 1 se utilizează mai multe seturi de date duplicate transferate pe canale separate. Dintre nivelurile de RAIC, acesta asigură cea mai mare disponibilitate a datelor, precum și cea mai mare performanță în transferul de date.

RAIC de nivel 0+1 – date distribuite și replicate

Acest nivel intermediar combină avantajele RAID 0 și RAID 1

RAIC de nivel 2 – protecție prin coduri Hamming

În RAIC 2 datele sunt distribuite pe un număr de canale separate, la care se adaugă un număr de alte canale care să asigure corecția unei erori de bloc, precum și detecția unei erori duble.

RAIC de nivel 3 – distribuția datelor și utilizarea unui canal separat pentru paritate

Datele sunt distribuite peste un număr de canale separate, la care se adaugă un canal pe care se transferă informația de paritate calculată cu ajutorul funcției XOR.

RAIC de nivel 4 – canale independente pentru distribuția datelor și paritate

Idem RAID 3 cu observația că se pot utiliza și independent canalele.

RAIC de nivel 5 – date și paritate distribuite

Idem RAID 4 cu observația că și paritatea este distribuită pe un set dintre canalele matricei.

RAIC de nivel 6 – date și paritate dublă distribuite

RAIC de nivel 7– RAID 5 cu canale de rezervă active

5.7 Aplicarea conceptului RAIC în transmisiile seriale de date

Pentru a testa parțial conceptul de RAIC, am folosit un stand de încercare format din:

- PC AT PENTIUM II 450 Mhz dotate cu cate 8 interfețe seriale RS232c – 2 bucăți.
- un cablu pentru transfer date – 8*UTP de lungime 15m

Informația transferată a fost divizată în blocuri de câte 512 cuvinte, blocurile fiind protejate prin paritate încrucișată. După fiecare bloc transferat emițătorul confirmă, cu un cuvânt de control, recepționarea corectă, corectată sau eronată.

Am folosit sisteme mult mai puternice decât ar fi fost necesare pentru aplicația dată în scopul reducerii la minimum a influențelor timpilor de calcul pentru paritate, respectiv pentru regenerarea datelor pierdute. Tot în același scop am utilizat paritatea încrucișată la nivel de blocuri.

Au fost testate, cu acest stand, RAIC 0, RAIC 1 (*chained declustering*), RAIC 4, RAIC 5, RAIC 6, RAIC 7 (MFT 3).

RAIC 2 și RAIC 3 le-am considerat nepotrivite pentru aplicația luată în considerație.

Interfețele seriale au fost setate pentru o rată de transfer de 33.6 Kbd.

Rezultatele obținute, în ceea ce privește rata globală de transfer, în condiții optime de funcționare, le prezint în Figura 5-14

| Nivel RAIC | Rata obținută [Kbd] |
|---------------|---------------------|
| RAIC 0 | 268.8 |
| RAIC 1 | 134.4 |
| RAIC 4 | 235.2 |
| RAIC 5 | 235.2 |
| RAIC 6 | 201.6 |
| RAIC 7 | 168 |

Figura 5-14

Rezultatele obținute, în ceea ce privește rata globală de transfer, în condițiile în care am provocat o cădere fizică (tăiere o legătură), le prezint în Figura 5-15

| Nivel RAIC | Rata obținută [Kbd] |
|-------------------|----------------------------|
| RAIC 0 | 0 |
| RAIC 1 | 134.4 |
| RAIC 4 | 201.6 |
| RAIC 5 | 201.6 |
| RAIC 6 | 168 |
| RAIC 7 | 134.4 |

Figura 5-15

Rezultatele obținute, în ceea ce privește rata globală de transfer, în condițiile în care am provocat o a doua cădere fizică (tăiere o legătură), le prezint în Figura 5-16

| Nivel RAIC | Rata obținută [Kbd] |
|-------------------|----------------------------|
| RAIC 0 | 0 |
| RAIC 1 | 0 – 134.4 |
| RAIC 4 | 0 |
| RAIC 5 | 0 |
| RAIC 6 | 134.4 |
| RAIC 7 | 100.8 |

Figura 5-16

Rezultatele obținute, în ceea ce privește rata globală de transfer, în condițiile în care am provocat o a treia cădere fizică (tăiere o legătură), le prezint în Figura 5-17

| Nivel RAIC | Rata obținută [Kbd] |
|------------|---------------------|
| RAIC 0 | 0 |
| RAIC 1 | 0 – 134.4 |
| RAIC 4 | 0 |
| RAIC 5 | 0 |
| RAIC 6 | 0 |
| RAIC 7 | 67.2 |

Figura 5-17

Rezultatele obținute, în ceea ce privește rata globală de transfer, în condițiile în care am supus unul dintre canale la influențele unui puternic câmp variabil electromagnetic, le prezint în Figura 5-18

| Nivel RAIC | Rata obținută [Kbd] |
|------------|---------------------|
| RAIC 0 | 0 |
| RAIC 1 | 134.4 |
| RAIC 4 | 201.6 |
| RAIC 5 | 201.6 |
| RAIC 6 | 168 |
| RAIC 7 | 134.4 |

Figura 5-18

Concluzii:

- cea mai bună performanță a oferit-o RAIC 0, în condiții normale de exploatare, dar și cea mai slabă fiabilitate în condiții defavorabile (stres fizic și electromagnetic);
- cea mai sigură soluție s-a dovedit RAIC 7.

5.8 Concluzii

Pe baza conceptelor prezentate, pe larg, în capitolul anterior, în partea a cincea a lucrării se prezintă o nouă soluție structurală de creștere a fiabilității sistemelor de calcul, soluție care cuprinde ca pe un caz particular structurile RAID. Principalele contribuții cuprinse în capitolul 5 sunt:

- 1) introducerea conceptului original – RAIC, care cuprinde complet, ca un subset, conceptul RAID;
- 2) definirea și sistematizarea metodelor fundamentale de distribuție a datelor într-un RAIC;
- 3) definirea și sistematizarea principalelor mecanisme de redundanță într-un RAIC;
- 4) o clasificare a structurilor RAIC;
- 5) realizarea practică a unei structuri RAIC.

6 CONCLUZII

Urmărind implementarea eficientă a toleranței la defectare în subsistemele I/O, precum și creșterea performanțelor acestora prin noile structuri RAID level 7 și RAIC, prezenta teză de doctorat aparține, în principal, domeniului fiabilității sistemelor de calcul.

Contribuțiile din cadrul tezei de doctorat au fost expuse, în extenso, prin cate un paragraf distinct, la sfârșitul fiecărui capitol. Aceste contribuții nu vor mai fi înșiruite, preferându-se o sinteză a celor mai importante, în acest capitol, jalonată de potențiale direcții de cercetare deschise. Teza cuprinde:

- 1) Contribuții la sistematizarea metodelor generale de creștere a fiabilității sistemelor de calcul;
- 2) Contribuții la definirea unor noi parametri cantitativi de evaluare a fiabilității subsistemului I/O;
- 3) Contribuții la sistematizare diferitelor structuri de discuri tolerante la defectare;
- 4) Dezvoltare domeniului matricelor de discuri RAID prin dezvoltarea unor algoritmi originali, precum și prin definirea unei noi clase: RAID level 7;
- 5) Utilizarea generalizării noțiunilor de canal, transfer de date, emițători și receptori, pentru reducerea formală a numărului de componente tipice ale unui sistem de calcul;
- 6) Fundamentarea teoretică a unui nou concept – RAIC – utilizabil în proiectarea sistemelor tolerante la defectare;
- 7) Realizarea practică a unei structuri RAIC.

Viitoarele posibile cercetări cu problematica cuprinsă în această lucrare ar putea fi orientate spre dezvoltări ale RAID level 7 și, mai ales, în noul domeniu al RAIC.

7 BIBLIOGRAFIE

1. **[Adap92]** Adaptec Corporation: Selecting An Array Interface, Computer Technology Review, March 24, 1992, p. 27.
2. **[BBB95]** M.Blaum, J.Brady, J. Bruck, J. Menon: EVENODD: An Efficient Scheme for Tolerating Double Disk Failures in RAID Architectures. IEEE Transaction on Computers 44(2): 192–202, 1995.
3. **[Beal91]** D.G. Beal: Data Storage System for Providing Redundant Copies of Data in Different Disk Drives, International Patent Application, Publication Number WO 91/20034, December 26, 1991.
4. **[Bell84]** C.G. Bell: The Mini and Micro Industries Computer, vol. 17, no. 10, pp. 14-30, Oct. 1984.
5. **[Bell93]** Bellcore Applied Research Area: Profile of RAID — An Emerging Robust On-Line Storage Technology , Bell Communications Research publication, copyright 1992.
6. **[Bels91]** S.J. Belsan, et al: Deleted Data File Release Mechanism for a Dynamically Mapped Virtual Data Storage Subsystem, International Patent Application, Publication Number WO 91/20025, December 26, 1991.
7. **[BhDi92]** A. Bhide and D. Dias: RAID Architectures for OLTP, IBM Research Report RC 17879, March 24, 1992.
8. **[BiGr88]** D. Bitton and J. Gray: Disk Shadowing, Proc. 14th VLDB Conf., pp. 331–338, Los Angeles, Aug. 1988.
9. **[Bitt89]** D. Bitton: Arm Scheduling in Shadowed Disks, Proc. IEEE Spring ComCon, pp. 132–136, Feb. 1989.
10. **[BoCM91]** M.F. Bond, B.E. Clark and R.S. McRoberts: Method and Apparatus for Recovering Parity Protected Data, European Patent Application, Publication Number 0 462 917 A2, May 22, 1991.
11. **[Brad90a]** Brady, et al: Method and Means for Accessing DASD Arrays with Tuned data Transfer Rate and Concurrency, European Patent Application EU: 91304503.5, May 24, 1990.

12. **[Brad90b]** J.Brady: Organisation of XOR Hardware in Disk Arrays, IBM Invention Disclosure SA90-91-004, January 1990.
13. **[BuCS91]** W.A. Burkhard, K.C. Claffy, and T.J.E. Schwarz: Performance of Balanced Disk Array Schemes, Digest of Papers, 11th IEEE Symposium on Mass Storage Systems, October 1991.
14. **[BuJh94]** S.W. Burns, N.K. Jha: A Totally Sel-Checking Checker for a Parallel Unordered Coding Scheme. IEEE Transaction on Computers, April 1994, vol. 43, pp. 490–495.
15. **[Bult88]** D.L. Bultman: High Performance SCSI Using Parallel Drive Technology, Proc. BUSCON Conf., pp. 40–44, Anaheim CA, Feb. 1988.
16. **[CaLo91]** L-F. Cabrera and D.D.E. Long: SWIFT: Using Distributed Disk Striping to Provide High I/O Data Rates, Computing Systems, Fall 1991.
17. **[CARK91]** Chervenak, Ann L., Randy H. Katz: Performance of a RAID Prototype, 1991 ACM SIGMETRICS Conference on Measurement and Modeling of Computer Systems, San Diego CA, May 1991.
18. **[CGGH96]** W.V. Courtright II, G.A. Gibson, M. Holland, J. Zelenka: RAIDframe: Rapid Prototyping for Disk Arrays. SIGMETRICS 1996: pp. 268–269.
19. **[CGKP90]** P. Chen, G. Gibson, R.H. Katz, and D.A. Patterson: An Evaluation of Redundant Arrays of Disks Using an Amdahl 5890, Performance Evaluation Rev., vol. 18, no. 1, pp. 74–85, May 1990.
20. **[Chan92]** Chantal, A Division of BusLogic: Why Disk Arrays? Chantal Systems Publication, copyright 1992.
21. **[Chen91a]** C.H.Chen: Performance Models for Disk Arrays, IBM Tehnical Report TR-L41-C-3, 1991.
22. **[Chen91b]** C.H.Chen: Reliability of Some Reduntant Systems, IBM Technical Report TR-L41-C-1, 1991.
23. **[Chen91c]** C.H. Chen: M/G/1 Models of Disk Arrays with Small Data Transfer Size, IBM May 1991.
24. **[Chen92]** S.-Z. Chen: Design, Modeling, and Evaluation of High Performance I/O Subsystems, PhD thesis, Univ. of Massachusetts at Amherst, Sept. 1992.

25. **[Cher90]** A.L. Chervenak: Performance Measurement of the First RAID Prototype, Technical Report UCB/CSD 90/674, UC Berkeley, May 1990.
26. **[Chie69]** R.T. Chien: Burst-Correcting Codes with High-Speed Decoding. IEEE Transaction on Information Theory IT-15 (January 1969): pp. 109–113.
27. **[ChKa91]** A.L. Chervenak and R.H. Katz: Performance of a Disk Array Prototype, CMG Transactions, Fall 1991.
28. **[ChKa91]** A.L. Chervenak, R.H. Katz: Performance of a RAID Prototyp. SIGMETRICS 1991: pp. 188–197.
29. **[ChLe95]** P.M. Chen, E.K. Lee: Striping in a RAID Level 5 Disk Array. SIGMETRICS 1995: pp. 136–145.
30. **[ChPa91]** P.M. Chen and D.A. Patterson: Maximizing Performance in a Striped Array, CMG Transactions, Fall 1991.
31. **[ChRe92]** J. Chandy and A.L.N. Reddy: Failure Evaluation of Disk Array Organizations, IBM Research Report RC 8706, March 30, 1992.
32. **[ChTo92]** S.-Z. Chen and D. Towsley: Raid 5 vs. Parity Striping: Their Design and Evaluation, J. Parallel and Distributed Computing, vol. 17, pp. 58–74, 1992.
33. **[ChTo96]** S. Chen, D.F. Towsley: A Performance Evaluation of RAID Architectures. IEEE Transaction on Computers 45(10): pp. 1116–1130, 1996.
34. **[ChWH95]** A.J. Chaney, I.D. Wilson, A. Hopper: The Design and Implementation of a RAID 3 Multimedia File Server. NOSSDAV 1995: pp. 306–317.
35. **[Cipr90]** Ciprico, Inc.: Redundant Disk Arrays Enhance Data Safety To Support Network Servers, Computer Technology Review, Winter 1990, p. 20.
36. **[Clar88]** B.E. Clark, et al: Parity Spreading to Enhance Storage Acces, US Patent 4 761 785, August 1988.
37. **[CLDL94]** P.M. Chen, E.K. Lee, A.L. Drapeau, K. Lutz, E.L. Miller, S. Seshan, K. Shirriff, D.A. Patterson, R.H. Katz: Performance and Design Evaluation of the RAID 2 Storage Server. Distributed and Parallel Databases 2(3): pp. 243–260, 1994.

38. **[CLGK94]** P.M. Chen, E.L. Lee, G.A. Gibson, R.H. Katz, D.A. Patterson: RAID High-Performance, Reliable Secondary Storage. *Computing Surveys* 26(2): pp. 145–185, 1994.
39. **[CLVW93]** P. Cao, S.B. Lim, S. Venkataraman, J. Wilkes: The TickerTAIP Parallel RAID Architecture. *ISCA 1993*: pp. 52–63.
40. **[CoKe89]** G. Copeland and T. Keller: A Comparison of High-Availability Median Recovery Techniques, *ACM SIGMOD Conference*, 1989.
41. **[CPDP90]** Chen, Peter M., David A. Patterson: Maximizing Performance in a Striped Disk Array, *Proceedings of the 1990 ACM SIGARCH 17th Annual International Symposium of Computer Architecture*, Seattle WA, May 1990.
42. **[CPGG90]** Chen, Peter M., Garth A. Gibson, Randy H. Katz, David A. Patterson: An Evaluation of Redundant Arrays of Disks Using an Amdahl 5890, *Proc. of the 1990 ACM SIGMETRICS Conference on Measurement and Modeling of Computer Systems*, Boulder CO, May 1990.
43. **[Croc89]** T.W. Crockett: File Concepts for Parallel I/O, *ACM* 1989.
44. **[D&TR89]** D&T Roundtable: Built-in Self Tests: Are Expectations Too High? *IEEE Design & Test of Computers*, June 1989, pp. 66–74.
45. **[Data91]** DataQuest Disk Array Forecast, January, 1991.
46. **[Digi91]** DigiData: RAID Slow Getting Off The Ground, *Computer Tehnology Review*, Winter 1991, p. 44.
47. **[DoOu89]** F. Douglass and J. Ousterhout: Log-Structured File Systems, *Proc. IEEE Spring ComCon*, pp. 124–129, Feb. 1989.
48. **[DSHM94]** A.L. Drapeau, K. Shirriff, J.H. Hartman, E.L. Miller, S. Seshan, R.H. Katz, K. Lutz, D.A. Patterson, E.K. Lee, P.M. Chen, G.A. Gibson: RAID-II: A High-Bandwidth Network File Server. *ISCA 1994*: pp. 234–244.
49. **[Dunp90]** R.H. Dunphy Jr., et al: Disk Drive Memory, *US Patent 4 914 656*, April 1990.
50. **[Dyna92]** Dynatek Corporation: Implementing Raid in Cross-Platform Environments, *Computer Technology Review*, Winter 1992, p. 43.

51. **[Endl89]** ENDL Consulting: RAID Arrays Making a Slow Start, Computer Technology Review, September 1989, p. 22.
52. **[Fire59]** P.Fire: A Class of Multiple-Error-Correcting Binary Codes for Non-Independent Errors. Sylvania Report no. RSL-E-2, Mountain View CA: Sylvania Electronic Defense Laboratory, Reconnaissance System Division, 1959.
53. **[Form93]** Formation: Choosing Hardware/Software RAID, Computer Technology Review, July, 1993, p.39.
54. **[FRSH87]** A.W. Funkenbusch, T.A. Rinehart, D.W. Siitari, Y.S. Hwang, R.N. Gardner: Magneto-Optics Technology for Mass Storage Systems, May 1987: pp. 101–106.
55. **[Fuji90a]** E. Fujiwara: Logic Testing and Design for Testability. England: MIT Press, 1990.
56. **[Fuji90b]** E. Fujiwara, K. Pradhan: Error-Control Coding in Computers. IEEE Computer, July 1990, pp. 63–71.
57. **[Gait85]** N. Gaitanis: A Totally Self-Checking Error Indicator. IEEE Transaction on Computers, August 1985, vol. C-34, pp. 758–761.
58. **[Gait88]** N. Gaitanis: Totally Self-Checking Checkers with Separate Internal Fault Indication. IEEE Transaction on Computers, October 1988, vol. 37, pp. 1206–1213.
59. **[GaKo98]** E. Gabber, H.F. Korth: Data Logging: A Method for Efficient Data Updates in Constantly Active RAIDs. ICDE1998: pp. 144–153.
60. **[GaMS88]** H. Garcia-Molina and K. Salem: The Impact of Disk Striping on Realibility, IEEE Data Base Engineering Bulletin, March 1988.
61. **[GGHK89]** Gibson, Garth A., Lisa Hellerstein, Richard M. Karp, Randy H. Katz, David A. Patterson: Coding Techniques for Handling Failures in Large Disk Arrays Third International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS III), Boston MA, April 1989, pp. 123–132.
62. **[GHKK89]** G.A. Gibson, L. Hellerstein, R.M. Karp, R.H. Katz and D.A. Patterson: Failure Correction Techniques for Large Disk Arrays, Proceedings the 3rd International Conference on Arhitectural Support for Programming Languages and Operating Systems, 1989.

63. **[Gibs89]** G.A. Gibson: Performance and Reliability in Redundant Arrays of Inexpensive Disks, CompCon, February 1989.
64. **[Gibs91]** G.A. Gibson: Redundant Disk Arrays: Reliable Parallel Secondary Storage, UC Berkeley Technical Report UCB/CSD 90/613, March 1991.
65. **[Gibs92]** Garth A. Gibson: Redundant Disk Arrays/Reliable, Parallel Secondary Storage, MIT Press, 1992.
66. **[Gibs95]** G.A. Gibson: Storage Technology: RAID and Beyond. SIGMOD Conference 1995: p. 471.
67. **[GiPa93]** G. Gibson and D. Patterson: Designing Disk Arrays for High Data Reliability, Journal of Parallel and Distributed Computing, January, 1993, pp. 4–27.
68. **[Gold91]** S. Goldstein: Array Storage Devices: Teaching Old DASD New Tricks, CMG Transaction, Fall 1991.
69. **[GrHW90]** J. Gray, B. Host, and M. Walker: Parity Striping of Disk Arrays: Low-Cost Reliable Storage with Acceptable Throughput, Proc. 16th VLDB Conf., pp. 148–161, Brisbane, Australia, 1990.
70. **[Hark81]** J.M. Harker et al.: A Quarter Century of Disk File Innovation, IBM J. Research and Development, vol. 25, pp. 677–689, Sept. 1981.
71. **[Hart88]** C.B. Hartness: Data Error Correction System, US Patent Number 4 775 978, October 4, 1988.
72. **[HCTS81]** M.Y. Hsiao, W.C. Carter, J.W. Thomas, W.R. Stringfellow: Reliability, Availability and Serviceability of IBM Computers Systems: A Quarter Century of Progress: IBM Journal of Research and Development 25 (September 1981): pp. 453–461.
73. **[HePa90]** J.L. Hennessy and D.L. Patterson: Computer Architecture — A Quantitative Approach. Morgan Kaufman, 1990.
74. **[Hida92]** HI-Data Corporation: One RAID Solution Does Not Fit All, Computer Technology Review, November 1992, p. 37.
75. **[HoGi92]** M. Holland and G. Gibson: Parity Declustering for Continuous Operation in Redundant Disk Arrays, Proc. ASPLOS-V, pp. 23–35, Sept. 1992.

76. **[HoPa93]** R.Y. Hou, Y.N. Patt: Comparing Rebuild Algorithms for Mirrored and RAID 5 Disk Arrays. SIGMOD Conference 1993: pp. 317–326.
77. **[Howe84]** T.D. Howell: Analysis of Correctable Errors in the IBM3380 Disk File: IBM Journal of Research and Development, March 28, 1984: pp. 206–211.
78. **[HsDe]** H. Hsiao and D. DeWitt: Chained Declustering: A New Availability Strategy for Multiprocessor Database Machines, Proc. Sixth Int’l Conf. Data Eng., pp. 456–465, Los Angeles, Feb. 1990.
79. **[Ibis89]** IBIS Systems: Is a Disk Array a Cost-Effective Option? Computer Technology Review, May 1989, p. 22.
80. **[Inte92]** Interphase Corporation: Planing For Disk Array Integration, Computer Technology Review, November, 1992, p. 38.
81. **[IYHI87]** K. Itao, A. Yamaji, S. Hara, N. Izawa: Magneto-Optical Mass Storage System with 130 mm Write-Once Disk Compatibility: 8th IEEE Symposium on Mass Storage Syst. (May 1987): pp. 92–97.
82. **[John89]** B.W. Johnson. Design and Analysis of Fault Tolerant Digital System. Addison-Wesley Publishing Company, 1989.
83. **[John93]** B.W. Johnson. Boundary Scan Eases Test of New Technologies. Test&Measurement Europe, Autumn 1993.
84. **[KaGP89]** R.H. Katz, G. Gibson, and D.A. Patterson: Disk System Architectures for High Performance Computing, Proc. IEEE, vol. 77, No. 12, pp. 1842–1858, Dec. 1989.
85. **[KaVa89]** N. Kanopoulos, N. Vasanthavada: A Built-in Test Module for Fault Isolation. IEEE Design & Test of Computers, June 1989, pp. 58–65.
86. **[Kim85]** M.Y. Kim: Parallel Operation of Magnetic Disk Storage Devices: Synchronized Disk Interleaving ,Proceedings of the Fourth International Workshop on Data-Base Machines, 1985.
87. **[KiM86]** M.Y. Kim: Synchronized Disk Interleaving, IEEE Trans. on Computers, vol. 35, no. 11, pp. 978–988, Nov. 1986.
88. **[KiPa85]** M.Y. Kim and A. Patel: Error-Correcting Codes for Intervelnd Disks With Minimum Redundancy, IBM Research Report RC 11185, May 1985.
89. **[KiTa87]** M.Y. Kim and A. Tantawi: Asynchronous Disk Intervelning, IBM Reasearch Report RC 12497, January 1987.

90. **[KiTa91]** M.Y. Kim and A.N. Tantawi: Asynchronous Disk Interleaving: Approximating Acces Delays, IEEE Transaction on Computers, July 1991.
91. **[KiTa91]** M.Y. Kim and A.N. Tantawi: Asynchronized Disk Interleaving: Approximating Access Delays, IEEE Trans. on Computers, vol. 40, No. 7, pp. 801–810, July 1991.
92. **[KKIS87]** O. Kato, K. Kurosawa, K. Iwakuni, M. Shimbo: Development of Error Correction Method and LSI for Optical Disk Data Storage System: Paper of Technical Group on Computer Architecture, IPS Japan (September 1987): pp. 67–73.
93. **[KRGP89]** Katz, Randy H., G.A. Gibson, D.A. Patterson: Disk System Architectures for High Performance Computing, Proceedings of the IEEE, 77, December 12, 1989, pp. 1842–1858.
94. **[KRJO89]** Katz, Randy H., John K. Ousterhout, David A. Patterson, Peter M. Chen, Ann L. Chervenak, Rich Drewes, Garth A. Gibson, Ed K. Lee, Ken C. Lutz, Ethan L. Miller, Mendel Rosenblum: A Project on High Performance I/O Subsystems, ACM Computer Architecture News, 17, 5, Invited paper, September 1989, pp. 24–31.
95. **[Krol86]** T. Krol: (N,K) Concept Fault Tolerance. IEEE Transaction on Computers, April 1986, vol. C-35, No. 4, pp. 339–349.
96. **[Kurt87]** C. Kurtz : Development of a High-Capacity Performance Optical Storage System. Dig. 8th IEEE Symposium on Mass Storerger System (May 1987): pp. 107–111.
97. **[Kurz88]** F. Kurzweil: Small Disk Arrays-The Emerging Approach to High Performance, CompCon, Spring 1988.
98. **[Lala85]** P.K. Lala: Fault Tolerant and Fault Testable Hardware Design. Prentice-Hall International, London, 1985.
99. **[Lawl81]** F.D. Lawlor: Efficient Mass Storage Parity Recovery Mechanism, IBM Technical Dislocure Bulletin 24, July 2, 1978.
100. **[Lee90]** E.K.Lee: Software and Performance Issues in the Implementation of a RAID Prototype, Technical Report UCB/CSD 90/573, UC Berkeley, May 1990.
101. **[Lee91]** E. Lee: Hardware Overview of RAID 2 , UC Berkeley RAID Retreat, Lake Tahoe, January 1991.

102. **[Leg92]** Legacy Storage Systems: Back to Basics — A Disk Array Technology Primer, Computer Technology Review, October 1992, p. 42.
103. **[LeKa91]** E.K. Lee and R.H. Katz: Performance Consequences of Parity Placement in Disk Arrays, Proc. Fourth Int'l Conf. on Architectural Support for Programming Langs. and OS, pp. 190–199, April 1991.
104. **[LeKa93]** E.K. Lee and R.H. Katz: An Analytic Performance Model of Disk Arrays, Proc. SIGMETRICS '93, pp. 98–109, 1993.
105. **[LERK91]** Lee, Edward K., Randy H. Katz: Performance Consequence of Parity Placement in Disk Arrays, Fourth International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS IV), Palo Alto CA, April 1991.
106. **[LiKB87]** M. Livny, S. Khoshafian, and H. Boral: Multi-Disk Management Algorithm, Proc. SIGMETRICS Conf., pp. 69–77, May 1987.
107. **[LiMP93]** J.J. Lim, J. Menon, D. Palmer: A Distributed Development Environment for Embedded Software. SP&E 23(11): pp. 1235–1248.
108. **[LoMa92]** R. W-M. Lo and N.S. Matloff: A Probabilistic Limit on the Virtual Size of Replicated Disk Systems, IEEE Transactions on Knowledge and Data Engineering, February 1992.
109. **[LoMo94]** D.D.E. Long, B.R. Montague: Swift/RAID: A Distributed RAID System. Computing Systems 7(3): pp. 333–359, 1994.
110. **[Lync72]** W.C. Lynch: Do Disk Arms Move? Performance Evaluation Rev., vol. 1, pp. 3–16, December 1972.
111. **[MaMc88]** A. Mahmood, E.J. McCluskey: Concurrent Error Detection Using Watchdog Processors — A Survey. IEEE Transaction on Computers, February 1988, vol. 37, pp. 160–174.
112. **[Mank91]** P.S. Manka: Direct Acces Storage Device with Independently Stored Parity, US Patent Number 0 469 924 A2, August 2, 1991.
113. **[Marc93]** R. Mărculescu: Tehnici de bază în proiectarea pentru testabilitate. Tempus Postgraduate School of Computer Aided Electrical Engineering. Editor Daniel Iona, București, 1993.

114. **[Mass97]** P. Massglia: The RAIDbook, A Storage Technology Handbook, Sixth Edition, RAID Advisory Board, February 1997.
115. **[Matt91]** R.L. Mattson: Method for Balancing the Frequency of DASD Array Access When Operating in Degraded Mode, European Patent Application, Publication Number 0 469 924 a2, August 2, 1991.
116. **[Mccl90]** E. McCluskey: Design Techniques for Testable Embedded Error Checkers. IEEE Computer, July 1990, pp. 84–88.
117. **[Mead89]** W.E. Meador: Disk Arrays, CompCon, February 1989.
118. **[MeCo93]** J. Menon, J. Cortney: The Architecture of a Fault-Tolerant Cached RAID Controller. ISCA 1993: pp. 76-86.
119. **[MeKa89]** J. Menon and J. Kasson: Methods for Improved Update Performance of Disk Arrays, Proceedings of the Twenty-Fifth Hawaii International Conference on System Sciences, 1991 (see also IBM Research Report RJ 6928, July 1989).
120. **[MeKa90]** J. Menon and J. Kasson: Methods for Improved Update Performance of Disk Arrays, Technical Report RJ 6928 (66034), IBM Almaden Research Center, Nov. 1990.
121. **[MeMa91]** J. Menon and D. Mattson: Performance of Disk Arrays in Transaction Processing Environments, CMG Transaction, Fall 1991.
122. **[MeMa92a]** J. Menon, D. Mattson: Comparison of Sparing Alternatives for Disk Arrays. ISCA 1992: pp. 318–329.
123. **[MeMa92b]** Performance of Disk Arrays in Transaction Processing Environments. ICDS 1992: pp. 302–309.
124. **[MeMa92c]** J. Menon and R. Mattson: Distributed Sparing in Disk Arrays, CompCon, February 1992.
125. **[MeMa92d]** J. Menon and D. Mattson: Comparison of Sparing Alternative for Disk Arrays, Proceedings of the International Symposium on Computer Architecture, May 1992.
126. **[MeMN91]** J. Menon, D. Mattson, and S. Ng, Performance of Disk Arrays in Transaction Processing Environments, Technical Report RJ 8230 (75424), IBM Almaden Research Center, July 1991.

127. **[Meno94]** J. Menon: Performance of RAID5 Disk Arrays with Read and Write Caching. Distributed and Parallel Database 2(3): pp. 261–293, 1994.
128. **[MeRK93]** J. Menon, J. Roche, and O.J. Kasson, Floating Parity and Data Disk Arrays, J. Parallel and Distributed Computing, vol. 17, pp. 129–139, 1993.
129. **[MeRW95]** J. Menon, Jeff Riegel, James C. Wyllie: Algorithms for Software and Low-Cost Hardware RAIDs. COMPCON 1995: pp. 411–418.
130. **[MeTr97]** J. Menon, K. Treiber: Daisy: Virtual-Disk Hierarchical Storage Manager. SIGMETRICS Performance Evaluation Review 25(3): pp. 37–44, 1994.
131. **[MeYu91]** A. Merchant and P.S. Yu: Modeling and Comparisons of Striping Strategies in Data Replication Architectures, IBM Research Report RC 19960, 1991.
132. **[MeYu91]** A. Merchant and P.S. Yu: Modeling and Comparisons of Striping Strategies in Data Replication Architectures, IBM Technical Report, T.J. Watson Research Center, Sept. 1991.
133. **[MeYu92]** A. Merchant and P.S. Yu, Design and Modelling of Clustered RAID, Proc. 22nd FTCS, pp. 140–149, Boston, July 1992.
134. **[MeYu95]** A. Merchant and P.S. Yu, Analytic Modeling and Comparisons of Striping Strategies for Replicated Disk Arrays, IEEE Trans. on Computers, vol. 44, No. 3, pp. 419–433, March 1995.
135. **[MeYu96]** A. Merchant, P.S. Yu: Analytic Modeling of Clustered RAID with Mapping Based on Nearly Random Permutation. IEEE Transaction on Computers 45(3): pp. 367–373, 1996.
136. **[Micr91]** Micropolis Corp.: Drive Crush Terror Boosts RAID Market, Computer Tehnology Review, May 1991.
137. **[MiGr91]** C.A. Milligan and J.M. Graziano: Logical Track Write Scheduling System for a Parallel Disk Drive Array Data Storage Subsystem, International Patent Application, Publication Number WO 91/16711, October 31, 1991.
138. **[MiMo96]** S.K. Mishra, P. Mohapatra: Performance Study of RAID 5 Disk Arrays with Data and Parity Cache. ICPP, vol. 1, 1996: pp. 222–229.

139. **[MoFS91]** A.N. Mourad, W.K. Fuchs and D.G. Saab: Performance Evaluation of Redundant Disk Array Support for Transaction Recovery, Coordinated Science Laboratory, University of Illinois, Urbana, Nov. 11, 1991.
140. **[MoKi94]** K. Mogi, M. Kitsuregawa: Dynamic Parity Stripe Reorganization for RAID 5 Disk Arrays. PDIS 1994: pp. 17–26.
141. **[MoKi96]** K. Mogi, M. Kitsuregawa: Hot Mirroring: A Study to Hide Parity Upgrade Penalty and Degradations During Rebuilds or RAID 5. SIGMOD Conference 1996: pp. 183–194.
142. **[Moru94]** C. Morun: Stadiul actual al implementării metodelor de fiabilizare la transmiterea și stocarea datelor. Referat de doctorat, ianuarie 1994.
143. **[Moru94]** C. Morun: Metode de protecție la transmiterea informației. Referat de doctorat, iunie 1994.
144. **[Moru94]** C. Morun: Implementarea în rețele de calculatoare a protecției informației. Referat de doctorat, iunie 1994.
145. **[Moru98a]** C. Morun: RAIC, Conferința CERF Banat '98, noiembrie 1998.
146. **[Moru98b]** C. Morun: RAID Level 7, SAFECOMP '99, 18th Int'l. Conf. on Computer Safety, Reliability & Security, Toulouse, France — lucrare trimisă spre recenzie.
147. **[Moru98c]** C. Morun : RAIC and Fault Tolerance, SAFECOMP '99, 18th Int'l. Conf. on Computer Safety, Reliability & Security, Toulouse, France — lucrare trimisă spre recenzie.
148. **[Moru98d]** C. Morun: RAID 7 vs RAID 5 with Spare Disks, Seventh Ann. Working Conf. on Information Security Management & Small Systems Security, Amsterdam, the Netherlands — lucrare trimisă spre recenzie.
149. **[Moru98e]** C. Morun: Striping Strategies for RAIC, SAFECOMP '99, 18th Int'l. Conf. on Computer Safety, Reliability & Security, Toulouse, France — lucrare trimisă spre recenzie.
150. **[Moru98f]** C. Morun: An algorithm for “nFT RAIC”, COMPSAC '99, 23rd Int'l Computer Software & Applications Conf., Phoenix AZ – lucrare trimisă spre recenzie.

151. **[MoSt97a]** C. Morun, A. Stanciu. RAID with Multiple Fault Tolerance. Buletinul Științific al Universității „Politehnica“ din Timișoara, România, tomul 42 (56), 1997, pp. 131–136.
152. **[MoSt97b]** C. Morun, A. Stanciu. Better Results with RAID. Buletinul Științific al Universității „Politehnica“ din Timișoara, România, tomul 42 (56), 1997, pp. 137–142.
153. **[MuLu90]** R.R. Muntz and J.C. Lui: Performance Analysis of Disk Arrays Under Failure, Proc. 16th VLDB Conf., pp. 162–173, Brisbane, Australia, 1990.
154. **[NaKa88]** T. Nanya, T. Kawamura: Error Secure/Propagating Concept and its Application to the Design of Strongly Fault Secure Processors. IEEE Transaction on Computers, January 1988, vol. 37, pp. 14–24.
155. **[NCR91a]** NCR Corporation: NCR Disk Arrays: Solutions For Your Business Environment NCR publication, copyright 1991.
156. **[NCR91b]** NCR corporation: What Are Disk Arrays? NCR Publication, copyright 1990, 1991.
157. **[NCR91c]** NCR Corp.: Disk Arrays Diving Into The Computer Mainstream, Computer Tehnology Review, October, 1991, p. 12.
158. **[NCR93a]** NCR Corporation: Can Fault Resilience Protect RAIDs?, Computer Technology Review, February, 1993, p. 28.
159. **[NeWO88]** M. Nelson, B. Welch, and J. Ousterhout: Caching in the Sprite Network File System, ACM Trans. on Computer Systems, vol. 6, No. 1, pp. 134–154, Feb. 1988.
160. **[Ng87]** S. Ng: Design Alternatives for Disk Duplexing , IBM Research Report RJ 5481, January 30, 1987.
161. **[Ng89]** S. Ng: Pitfalls in Deigning Disk Arrays, CompCon, February 1989.
162. **[Ng89]** S. Ng: Some Design Issues of Disk Array, Proc. IEEE Spring ComCon, pp. 137–142, 1989.
163. **[NgLS88]** S. Ng, D. Lang, and R. Selinger: Trade-Offs Between Devices and Paths in Achieving Disk Interleaving, Proc. 15th Int’l Symp. Computer Architecture, pp. 196–201, Hawaii, May 1988.
164. **[NgMa94]** S. Ng and R.L. Mattson: Uniform Parity Group Distribution in Disk Arrays with Multiple Failures, IEEE Trans. on Computers, vol. 43, No. 4, pp. 501–506, April 1994.

165. **[OgFl90]** M. Ogata and M. Flynn: A Queueing Analysis for Disk Array Systems, Technical Report CSL-TR-90-443, Stanford Univ., August 1990.
166. **[Ohre90]** E. Ohrenstein: Supercomputers Seek High Throughput and Expandable Storage, Computer Technology Review, Spring 1990.
167. **[Olso89]** T.M. Olson: Disk Array Performance in a Random 10 Environment, Computer Architecture News, September 1989.
168. **[OsDo91]** J. Osterhout and F. Dougliis: Beating the I/O Bottleneck: A Case for Log-Structured File Systems, CMG Transactions, Fall 1991.
169. **[Ouch78]** N.K. Ouchi: System for Recovering Data in Failed Memory Unit, US Patent Number 4 092 732, May 30, 1978.
170. **[OuDo89]** J.K. Ousterhout and F. Dougliis: Beating the I/O Bottleneck: A Case for Log-Structured File Systems, Operating System Rev., vol. 23, No. 1, pp. 11–28, Jan. 1989.
171. **[PaBa86]** A. Park and K. Balasubramanian: Providing Fault Tolerance in Parallel Secondary Storage Systems, Technical Report CS-TR-057-86, Princeton University, November 1986.
172. **[PAGK88]** D.A. Patterson, G. Gibson, and R.H. Katz: A Case for Redundant Arrays of Inexpensive Disks (RAID), Proc. ACM SIGMOD Conf., pp. 109–116, July 1988.
173. **[PaGK89]** D.A. Patterson, G. Gibson and R. H. Katz: Introduction to Redundant Arrays of Inexpensive Disks (RAID), CompCON, February 1989.
174. **[PCGK89]** D.A. Patterson, P. Chen, G. Gibson, and R.H. Katz: Introduction to Redundant Arrays of Inexpensive Disks (RAID), Proc. IEEE Spring ComCon, pp. 112–117, 1989.
175. **[PDGG88]** Patterson, David A., Garth A. Gibson, Randy H. Katz: A Case for Redundant Arrays of Inexpensive Disks (RAID) Proceedings of the 1988 ACM SIGMOD Conference on Management of Data, Chicago IL, June 1988, pp. 109–116.
176. **[Plan97]** J.S. Plank: A Tutorial on Reed-Solomon Coding for Fault-Tolerance in RAID-Like Systems. SP&E 27(9): pp. 995–1012, 1997.

177. **[Prad86]** D.K. Pradhan. Fault Tolerant Computing: Theory and Techniques. Vol. I & II, Prentice-Hall, 1986.
178. **[Quan93]** Quantum Corporation: SCSI Array Features Can Cost Integrators, Technology Review, July 1993, p. 14.
179. **[ReBa89]** A.L.N. Reddy and P. Banerjee: An Evaluation of Multiple-Disk I/O Systems, IEEE Trans. Computers, vol. 38, No. 12, pp. 1680–1690, Dec. 1989.
180. **[ReBa91a]** A.L.N. Reddy and P. Banerjee: Gracefully Degradable Disk Arrays, Proc. 21st Int'l Symp. on Fault-Tolerant Computing, pp. 401–408, Montreal, June 1991.
181. **[Redd91a]** A.L.N. Reddy: Reads and Writes: When I/O's aren't quite the same, IBM Research Report RJ 8033, March 1991.
182. **[Redd91b]** A.L.N. Reddy: Reads and Writes: When I/O's Aren't Quite the Same, Proceedings of the Twenty-Fifth Hawaii International Conference on System Sciences, 1991.
183. **[RiMe96]** J. Riegel, J. Menon: Performance of Recovery Time Improvement Algorithms for Software RAIDs. PDIS 1996.
184. **[RoFr94]** J.T. Robinson and P.A. Franaszek: Analysis of Reorganization Overhead in Log-Structured File Systems, Proc. 10th Int'l Conf. Data Eng., pp. 102–110, 1994.
185. **[RoOs91a]** M. Rosenblum and J.K. Osterhout: The LFS Storage Manager, CMG Transactions, Fall 1991.
186. **[RoOs91b]** M. Rosenblum and J. Osterout: The Design and Implementation of a Log Structured File System, Proceedings of the Thirteenth ACM Symposium on Operating System Principles, October 1991.
187. **[RoOu90]** M. Rosenblum and J.K. Ousterhout: The LFS Storage Manager, Proc. Summer '90 USENIX Technical Conf., pp. 315–324, Anaheim, Calif., June 1990.
188. **[RoOu91]** M. Rosenblum and J.K. Ousterhout, The Design and Implementation of a Log-Structured File System, Draft, March 1991.
189. **[RoOu92]** M. Rosenblum and J.K. Ousterhout, The Design and Implementation of a Log-Structured File System, ACM Trans. on Computer Systems, vol. 10, No. 1, pp. 26–52, Feb. 1992.

190. **[RuOk94]** T. Ruwart and M. O’Keefe: Performance Characteristics of a 100 Megabyte/Second Disk Array, presented in Storage & Interface ’94, Jan. 1994.
191. **[SaGM86]** K. Salem and H. Garcia-Molina: Disk Striping, Proceedings of the Second International Conference on Data Engineering, 1986.
192. **[SaMo86]** K. Salem and H. Garcia-Molina: Disk Striping, Proc. IEEE Data Eng. Conf., pp. 336–342, Los Angeles, Feb. 1986.
193. **[SaTN86]** M. Saito, T. Takeda, M. Nunotani: An Evaluation of Error Correcting Codes for Optical Disks: Paper of Technical Group IT86-48, IECE Japan, 1986.
194. **[ScBu92]** T.J.E. Swarcz, W.A. Burkhard: RAID Organization and Performance. ICDCS 1992: pp. 318– 325.
195. **[Schu88]** M.E. Schultze: Consideration in the Design of a RAID Prototype, Tehnical Report UCB/CSD 88/448, UC Berkeley, August 1988.
196. **[SGKP89]** M. Schulze, G. Gibson, R. Katz, and D. Patterson: How Reliable is a RAID? Proc. IEEE Spring ComCon Conf., pp. 118–123, San Francisco, Feb. 1989.
197. **[SHCG94]** D. Stodolsky, M. Holland, W.V. Courtright, and G.A. Gibson: Parity-Logging Disk Arrays, ACM Trans. on Computer Systems, vol. 12, No. 3, pp. 206–235, Aug. 1994.
198. **[Staf91]** Staff: Arrays More For Bandwidth Then Capacity, Computer Tehnology Review, October, 1991, p. 12.
199. **[Staf92a]** Staff: RAID Breaks Out Of Level 5 Straightjacket, Computer Tehnology Review, March, 1992, p. 1.
200. **[Staf92b]** Staff: RAID Board To Play The Number Game, Computer Tehnology Review, November, 1992, p. 41.
201. **[Staf92c]** Staff: Raid Becoming The People’s Choice, Computer Tehnology Review, August, 1992, p. 36.
202. **[Staf93a]** Staff: No Cop-Outs on Backup Whith RAID, Computer Tehnology Review, February, 1993, p. 1.
203. **[Staf93b]** Staff: Is SCSI Out of Sync With RAID?, Computer Tehnology Review, March, 1993, p. 1.

204. **[Staf93c]** Staff: Controllers Basic to RAID Success , Computer Technology Review, May 1993, p. 35.
205. **[Staf93d]** Staff: Arrays: Hardware of Software RAID, Computer Technology Review, July 1993, p. 26.
206. **[Staf93e]** Staff: How to Select Between RAID and Solid State, Computer Technology Review, September, 1993, p. 41.
207. **[Staf93f]** Staff: RAID Fights Fog for Channel , Computer Technology Review, October, 1993, p. 1.
208. **[Staf93g]** Staff: Raid 6 Adding Another Level of Confusion, Computer Technology Review, November, 1993, p. 1.
209. **[Staf94a]** Staff: RAID Has To Get Its Numbers Straight And Fast, Technology Review, January 1994, p. 1.
210. **[Staf94b]** Staff: Certification Coming, Technology Review, January 1994, p. 1.
211. **[Ston90]** M. Stonebreaker and G.A. Schloss: Distributed RAID.
212. **[Stor91]** Storage Concepts, Is RAID The Disk Strategy For The Masses?, Computer Tehnology Review, May 1991, p. 23.
213. **[Stor92]** Storage Concepts: RAID5 Is Less Than Expected, Computer Tehnology Review, March, 1992, p. 12.
214. **[Stor94]** Storage Dimensions: Which RAID Meets the Application?, Technology Review, February 1994, p. 37.
215. **[StSc90]** M. Stonebreaker and G.A. Schloss: Distributed RAID — a New Multiple Copy Algorithm, Proceedings of the 6th International Conference on Data Engineering, February 1990.
216. **[StSc90]** M. Stonebraker and G.A. Schloss: Distributed RAID — A New Multiple Copy Algorithm, Proc. Sixth Int'l Conf. Data Eng., pp. 430–437, Feb. 1990.
217. **[Tech93]** Technology Forums, Ltd.: Is RAID Five Levels of Confusion?, Technology Review, February 1993, p. 37.
218. **[ThMe94]** A. Thomasian, J. Menon: Performance Analysis of RAID 5 Disk Arrays with a Vacationing Server Model for Rebuild Mode Operation. ICDE 1994: pp. 111–119.
219. **[ThMe97]** A. Thomasian, J. Menon: RAID 5 Performance with Distributing Sparring. IEEE Transaction on Parralel and Distributed Systems 8 (6): pp. 640–657, 1997.

220. **[Thor73]** E. Thorp: Nonrandom Shuffling with Applications to the Game of Faro, *JASA*, vol 68, pp. 842-847, 1973.
221. **[Toy88]** W.N. Toy: Fault Tolerant Computing, 1988.
222. **[TPBG93]** F.A.Tobagi, J. Pang, R. Baird, M. Gang: Streaming RAID: A Disk Array Management System for Video Files. *ACM Multimedia*, 1993: pp. 393–400.
223. **[TsLe97]** W.J. Tsai, S.Y. Lee: Multi-Partition RAID: A New Method for Improving Performance of Disk Arrays under Failure. *The Computer Journal* 40 (1): pp. 30–42, 1997.
224. **[Vasu88]** A. Vasudeva: A Case for Disk Array Storage System, *Proceedings of Reliability Conference, Santa Clara CA*, 1988.
225. **[VINH80]** L.B. Vries, K.A. Imink, J.G. Nibor, S. Abe: The Compact Disc Digital Audio System — Modulation and Error Correction: *Proceedings of the Sixty-Seventh AES Convention (October 1980), No. 1674 (H-8)*.
226. **[Vlad82]** M. Vlăduțiu: Contribuții la creșterea coeficientului de disponibilitate al echipamentelor automate de prelucrare a informației prin testare neconvențională. Teză de doctorat, București, 1982.
227. **[VICr89]** M. Vlăduțiu, M. Crișan: Tehnica testării echipamentelor automate de prelucrare a datelor. Editura Facla, Timișoara, 1989.
228. **[VIPe94a]** M. Vlăduțiu, N.S. Petrakis.: About New Signature Generation Techniques for Binary Sequences. *International Conference on Technical Informatics. Proceedings*, vol. 5, Timișoara, 1994, pp. 11–15.
229. **[VIPe94b]** M. Vlăduțiu, N.S. Petrakis: Adapted Combinational Array for Exact Binary Division with Signed Operands. *International Conference on Technical Informatics. Proceedings*, vol. 5, Timișoara, 1994, pp. 1–10.
230. **[VTPB94]** M. Vlăduțiu, D. Todincă, N.S. Petrakis, I. Brînzaș: The Principle of a Fuzzy Processor with Parity Check. *Buletinul Științific al U.T. Timișoara*, 1994, tomul 39 (53), seria Automatizări și Calculatoare.
231. **[WeZS91]** G. Weikum, P. Zabback and P. Scheuermann: Dynamic File Allocation in Disk Arrays, *Proceedings of the 1991 ACM SIGMOD*, May 1991.

232. **[WGSS95]** J. Wilkes, R.A. Golding, C. Staelin, T. Sullivan: The HP AutoRAID Hierarchical Storage Syst. SOSP 1995: pp. 96–108.
233. **[WGSS96]** J. Wilkes, R.A. Golding, C. Staelin, T. Sullivan: The HP AutoRAID Hierarchical Storage Syst. TOCS 14(1): pp. 108–136, 1996.
234. **[Yama86]** A. Yamagishi: Encoder and Decoder LSIs for BCH Code and RS Code: Proceedings of 1986 Workshop on Coding Theory and Its Applications WCTA86-2 (August 1986).
235. **[Yarm90]** V.N. Yarmolik. Fault Diagnostics of Digital Circuits. New York Wiley, 1990.