

Arhitecturi Avansate de Calculatoare

(Arhitecturi Paralele)

Arhitecturi Avansate de Calculatoare

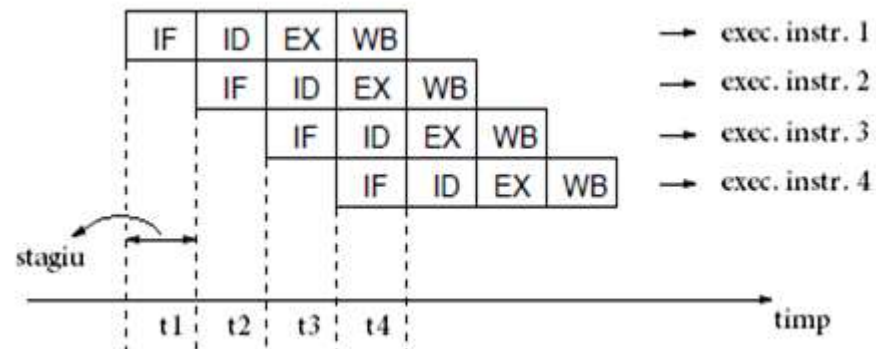
1. Notiuni introductive

- definitii; taxonomia arhitecturilor de calculatoare;
- familii de procesoare;
- metode de evaluare a performantelor; legea lui Amdahl;
- exemple de arhitecturi paralele;
- granularitatea sistemelor paralele

Arhitecturi Avansate de Calculatoare

2. Tehnologii utilizate in arhitecturile de procesoare

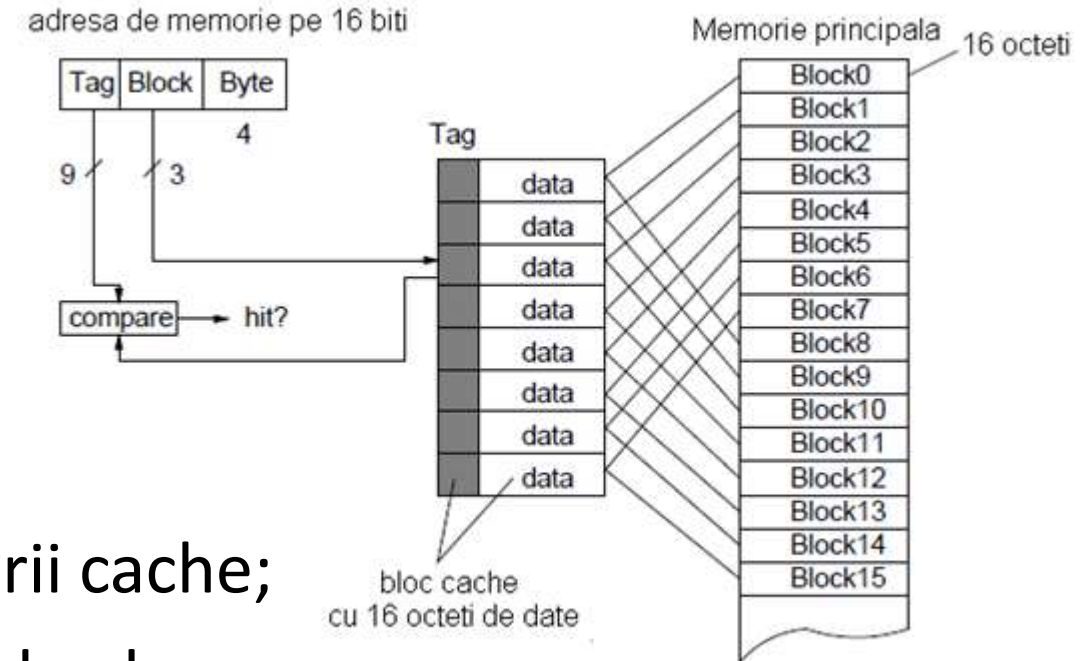
- Procesoare in sisteme paralele:
- transputerul;
- procesoare RISC;
- arhitectura pipeline; procesoare superscalare si superpipeline;
- procesoare VLIW; procesoare vectoriale



Arhitecturi Avansate de Calculatoare

3. Memorii cache

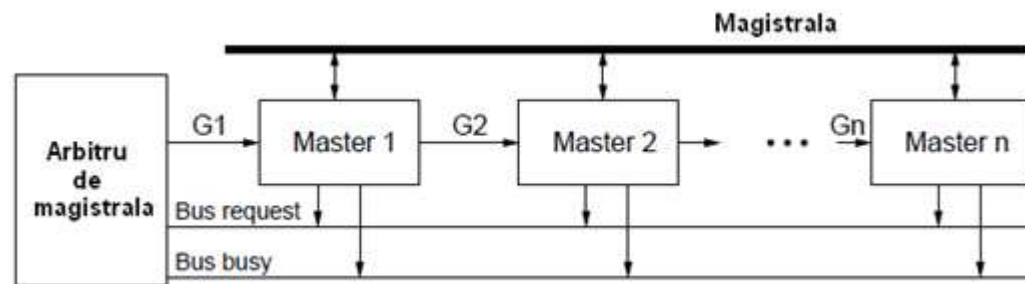
- memorii cache cu mapare directa;
- memorii cache complet asociative;
- memorii cache set-asociative;
- strategii de implementare pentru memorii cache;
- memorii write-through si memorii write-back



Arhitecturi Avansate de Calculatoare

4. Retele de comunicatie pentru sisteme paralele (I)

- comutarea si rutarea;
- clasificarea retelelor în functie de topologie
- retele directe; retele indirecte;
- retele bazate pe magistrala; magistrale sincrone si asincrone; magistrale cu tranzactii multiple



Arhitecturi Avansate de Calculatoare

5. Retele de comunicatie pentru sisteme paralele (II)

- retele multistagiu; retele blocante si nebloccante;
- retelele Omega;
- retele crossbar;
- comparatie între retelele directe si cele indirecte;
- comutarea prin pachete: metodele *Store & Forward*, *Wormhole Routing* si *Virtual Cut-through*

Arhitecturi Avansate de Calculatoare

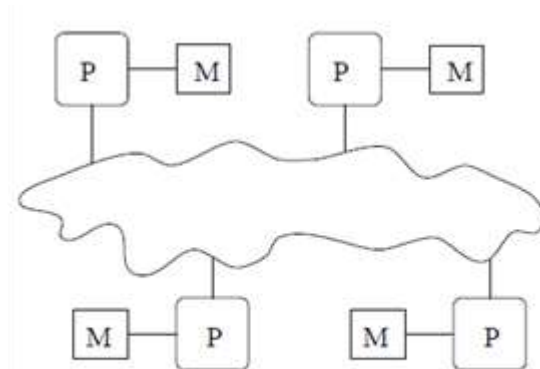
6. Tehnici de rutare

- rutare stabilita la nivelul nodului sursa;
- rutarea locala;
- rutarea adaptiva si rutarea determinista;
- eliminarea blocajelor de rutare; folosirea canalelor virtuale;
- metode pentru implementarea comunicatiei multicast

Arhitecturi Avansate de Calculatoare

7. Multicalculatoare

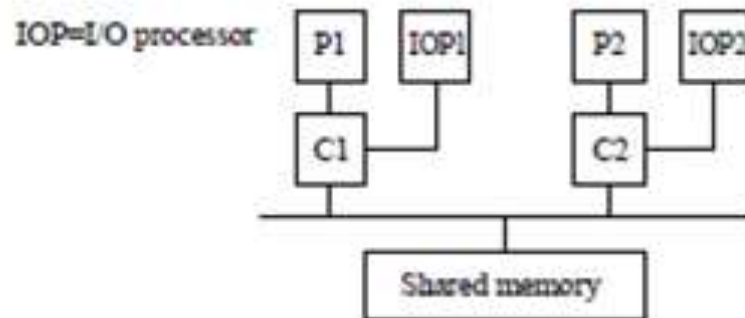
- implementările VSM și SVM;
- abordări și metode folosite în proiectarea sistemelor multicalculator;
- metoda cu server central; metoda migrației; metoda replicării datelor la citire; metoda *Full Replication*;
- exemple de multicalculatoare: IBM SP2 și Parsytec CC



Arhitecturi Avansate de Calculatoare

8. Sisteme multiprocesor (I)

- coerența memoriilor cache în sisteme cu memorie partajată;
- surse de inconsistență a datelor;
- protocoalele write-invalidat și write-update;
- arhitecturile UMA, NUMA, CC-NUMA



Arhitecturi Avansate de Calculatoare

9. Sisteme multiprocesor (II)

- ascunderea latentei la accesul memoriei de la distanta;
- modele de consistenta a memoriei;
- arhitectura COMA; comparatie între arhitectura COMA si arhitectura CC-NUMA;
- arhitecturile S-COMA, R-NUMA

Arhitecturi Avansate de Calculatoare

10. Sisteme multiprocesor (III)

Coerenta memoriilor cache in sisteme multiprocesor:

- protocoalele *Snoopy-bus* si bazate pe director;
- metode de implementare pentru memoria director: implementarile *Full Map* si *Limited Map*; metoda cu director înlantuit;
- diagramele de stari si tranzitii pentru memoriile *write-through* si *write-back*

Arhitecturi Avansate de Calculatoare

11. Sisteme multiprocesor (IV)

Sisteme de operare pentru multiprocesoare:

- exploatarea concurenței, detectarea paralelismului în programe, mecanisme de sincronizare, exemple

12. Standarde și medii de programare pentru arhitecturi paralele:

- standardul MPI, mediul PVM, limbajul OCCAM

Referinte

- Fundamentals of Parallel Multicore Architecture, Yan Solihin, Chapman and Hall/CRC June 30, 2020, ISBN-13: 978-0367575281
- O. Buza, Arhitecturi Paralele de Calculatoare, Ed. Grinta, Cluj-Napoca, 2018, ISBN 978-606-037-012-3
- J. L. Hennessy, D. A. Patterson, Computer Architecture, 6th Edition (The Morgan Kaufmann Series in Computer Architecture and Design), Elsevier, 2017, ISBN 978-0128119051
- J. L. Hennessy, D. A. Patterson, Computer Organization and Design RISC-V Edition: The Hardware Software Interface (The Morgan Kaufmann Series in Computer Architecture and Design) 1st Edition, Elsevier, 2017, ISBN-13: 978-0128122754
- G. Lerman, L. Rudolph, Parallel Evolution of Parallel Processors (Evaluation in Education and Human Services), Springer, 2013, ISBN-13: 978-1461362371
- Michel Dubois, Parallel Computer Organization and Design, 1st Edition, Cambridge University Press August 1, 2012, ISBN-13 : 978-0521886758

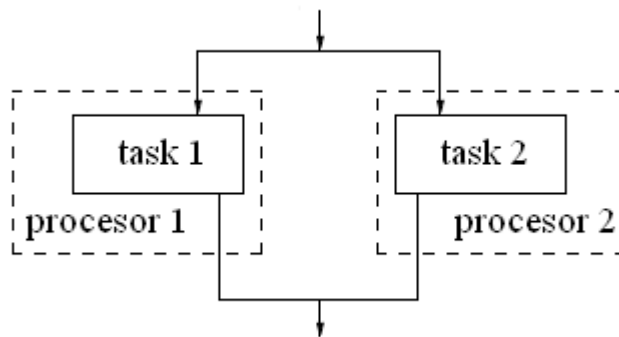
Online references

- Computer Architecture, 2020, <https://www.sciencedirect.com/topics/computer-science/computer-architecture>
- Computer Organization and Architecture, 2020, <https://www.geeksforgeeks.org/computer-organization-and-architecture-tutorials>

Curs 2 AAC

ARHITECTURI PARALELE DE CALCULATOARE NOȚIUNI INTRODUCTIVE

Prelucrarea paralelă reprezintă utilizarea mai multor procesoare pentru execuția simultană a mai multor părți dintr-un program (task-uri).



Avantaje

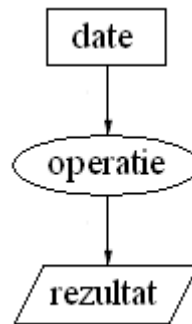
- Se elimină limitarea de viteză a procesoarelor secvențiale (timp de execuție mai scurt).
- Se elimină limita de miniaturizare a procesoarelor secvențiale (nu putem construi un procesor complex la dimensiuni oricât de mici).
- Se elimină limitarea economică (imposibilitatea de a fabrica ieftin un procesor foarte rapid => e mai bine să se folosească două sau mai multe procesoare ieftine care să ruleze în paralel pe același calculator).

Evoluția

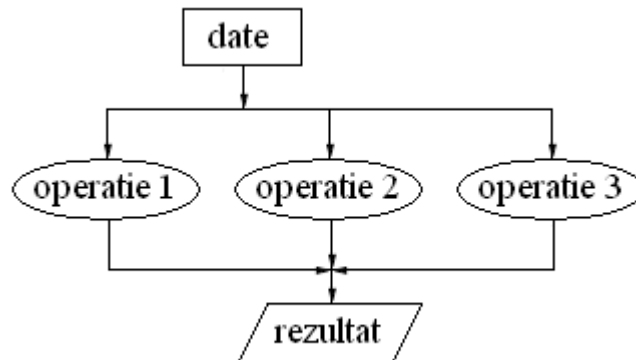
- Procesoare secvențiale obișnuite, numite procesoare scalare: execută maxim o instrucțiune într-un ciclu procesor.
- Procesoare superscalare: pot executa mai multe instrucțiuni într-un ciclu procesor.
- Procesoare vectoriale: nu mai lucrează pe operanzi scalari ci pe structuri de date multiple: vectori uni / multi dimensionali.
- Arhitecturi multiprocesor: într-un singur cip există mai multe nuclee de procesare (dual core / quad core / octa core).

Clasificarea arhitecturilor de calculatoare

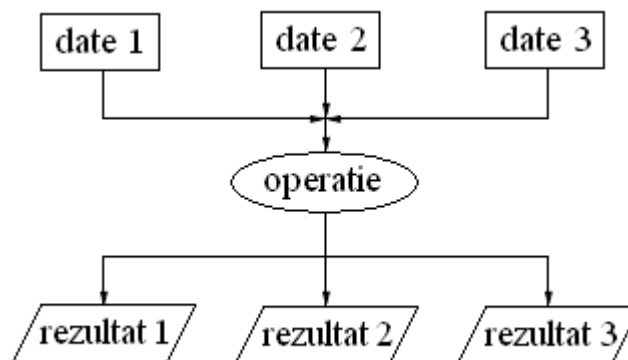
- **SISD** (*Single Instruction Single Data*): arhitectura secvențială obișnuită.



- **MISD** (*Multiple Instruction Single Data*): un program e descompus în mai multe task-uri sau secvențe de operații; operațiile se execută în paralel asupra aceluiași set de date.

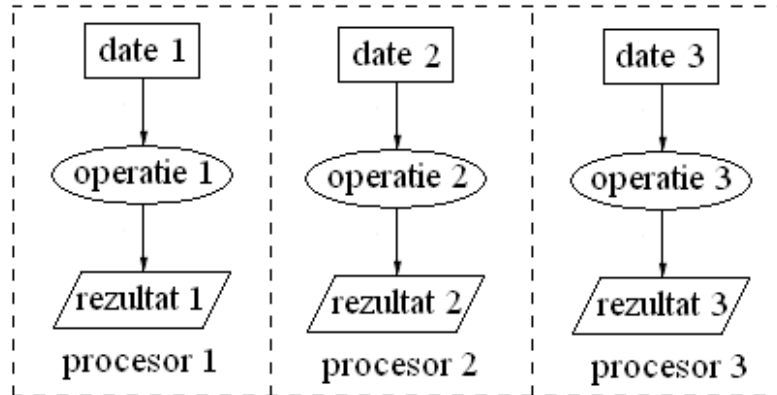


- **SIMD** (*Single Instruction Multiple Data*): un același program se execută pe mai multe seturi de date.



ex: prelucrari de masive (siruri, matrici); procesoare vectoriale

- **MIMD** (*Multiple Instruction Multiple Data*): reprezintă tipul general de arhitectură pentru calculatoarele paralele; se realizează operații diferite asupra unor seturi de date diferite.



Viteza de execuție a programelor paralele

Accelerația: expresie a vitezei de execuție: $Sp = \frac{T_{es}}{T_{ep}}$ unde:

T_{es} – timpul de execuție secvențial;

T_{ep} – timpul de execuție paralel pe mai multe procesoare al aceluiași program.

Ideal: $Sp = \frac{T_{es}}{T_{es}/P} = P$ (numărul de procesoare)

Real: $T_{ep} \neq T_{es}/P$; $T_{ep} = T_{secv} + T_{paralel}$ (o porțiune din program se poate executa doar secvențial în timpul T_{secv} și avem o porțiune care se poate executa în paralel)

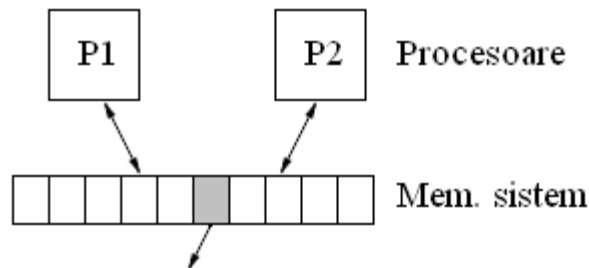
$$\Rightarrow Sp = \frac{T_{es}}{T_{secv} + T_{paralel}} \Rightarrow Sp < \frac{T_{es}}{T_{secv}} = K \quad \text{- Legea lui Amdahl}$$

Oricât am crește numărul de procesoare într-un sistem, accelerația (viteza de execuție în paralel) a unui program e întotdeauna limitată superior.

Accesul la memorie în procesarea paralelă

a) Acces partajat la memorie.

În acest caz, procesoarele împart același spațiu de memorie.

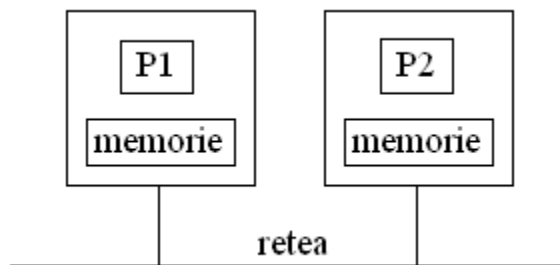


La un anumit moment, o anumită locație din memoria partajată poate fi accesată doar de către un singur procesor.

Avantaj: realizarea unui schimb rapid de date între 2 task-uri care rulează pe 2 procesoare diferite.

Exemplu: arhitectura UMA – procesoarele accesează în mod identic orice locație din memoria partajată a sistemului; arhitectura Pentium multicore; arhitectura GPU.

b) Acces distribuit la memorie.

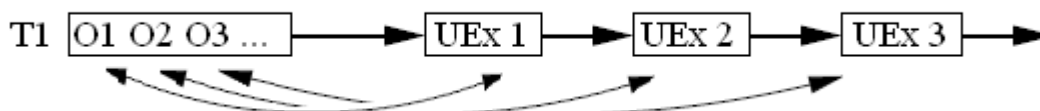


Procesoarele pot avea acces atât la o memorie locală aflată pe calculatorul din care fac parte, cât și la o memorie aflată la distanță (alt calculator).

Exemplu: arhitectura NUMA – procesoarele au acces rapid la memoria locală și acces lent la memoria aflată la distanță; arhitectura COMA – în care memoriile locale ale procesoarelor sunt organizate ca niște memorii cache.

Exemple de arhitecturi paralele

- 1) **Arhitectura Pipeline** – structură cu un singur procesor, dar care poate executa mai multe operații în paralel; astfel un anumit task se descompune într-o serie de operații care se execută independent de către unități de execuție specializate.



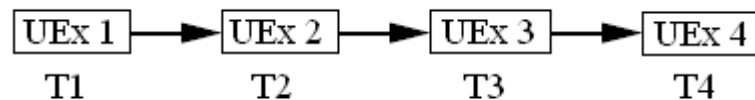
Există 2 tipuri de Pipeline: - de instrucțiuni; în care o instrucțiune se descompune într-o serie de microinstrucțiuni.

- aritmetic; în care o operație aritmetică se descompune într-o serie de microoperații.

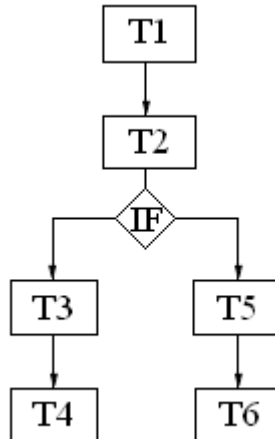
În rezolvarea unui task pot exista anumite situații speciale, în care fluxul de operații din Pipeline se oprește. Aceste situații speciale se numesc hazarduri (incidente) și sunt de 3 tipuri:

- a) Hazardul structural. El e generat de faptul că nu există destule unități de execuție pentru operațiile planificate.
- b) Hazardul de date. El e generat din cauza dependențelor de date. Un anumit task așteaptă după datele unui alt task.
- c) Hazardul de control. El e datorat ramificărilor care apar în program (instrucțiunile de salt).

De exemplu, dacă în structura Pipeline se încarcă la început patru operații:



în cazul schimbării fluxului de instrucțiuni din program, instrucțiunile încărcate în avans în structura Pipeline (T3 și T4) trebuie eliminate și înlocuite cu instrucțiunile corecte (T5 și T6).

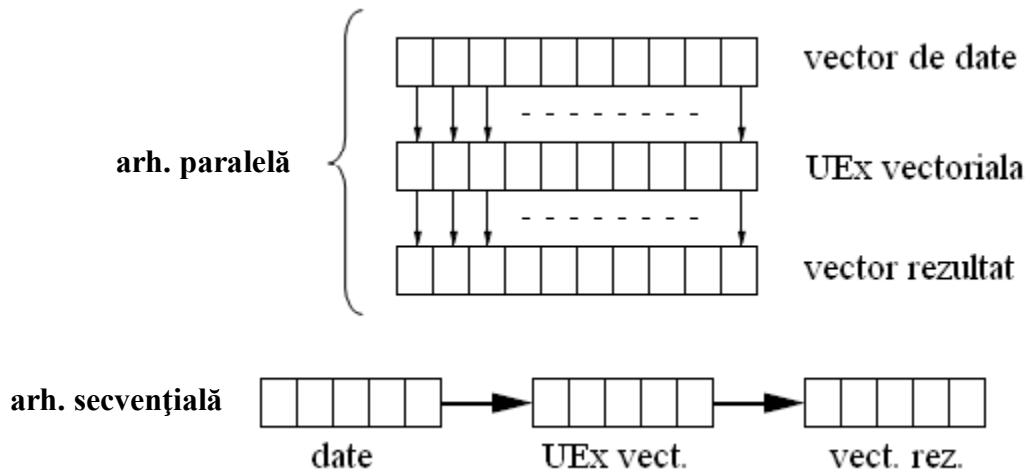


- 2) **Arhitecturi cu legături multiple** – acestea sunt structuri multiprocesor conectate în diverse configurații în care fiecare procesor realizează anumite operații simple într-un timp foarte scurt.

- a) **Arhitecturi vectoriale.**

Exemple:

- **Procesoare vectoriale**



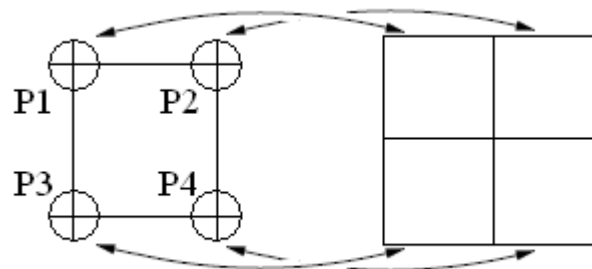
- **Procesoare sistolice:** la fiecare tact datele avansează o poziție prin structura sistolică.



Există un șir de elemente de prelucrare simple, conectate în diferite topologii; elementele de prelucrare pot avea o memorie locală și pot fi de tipul SIMD sau nu (fiecare element poate efectua aceeași sau o altă operație).

b) **Arhitecturi de tip rețea.**

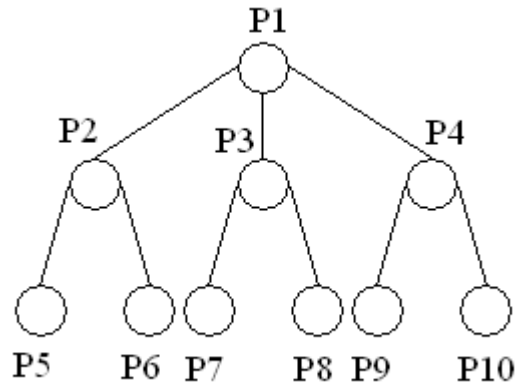
Exemplu: - procesoare matriciale



Unitățile de prelucrare sunt autonome și au o memorie locală atașată.

De obicei există o unitate de control care distribuie sarcinile la structura matricială de procesare.

- c) **Arhitectura de tip arbore.** Se folosește la calcule științifice / economice, programe de decizie.



- d) **Arhitectura hipercub.** E o arhitectură mixtă între rețea și arbore. Exemplu: fiecare nod al rețelei poate fi un arbore.

Granularitatea sistemelor multiprocesor

Granularitatea este rezoluția (gradul de finețe) cu care e rezolvată o anumită sarcină. În cazul de față, sarcina e execuția paralelă a unui program \Rightarrow granularitatea indică gradul de paralelizare pentru acel program.

Există 3 tipuri:

- Granularitate mică (grosieră) – realizarea paralelismului se face la nivel de proces.
- Granularitate medie – paralelismul se realizează la nivel de *thread* (fir de execuție din interiorul unui proces alcătuit din mai multe instrucțiuni care utilizează aceleași resurse).
- Granularitate mare (fină) – paralelismul se realizează la nivel de instrucțiune.

PROCESOARE ÎN SISTEME PARALELE

În trecut, calculatoarele paralele conțineau procesoare proiectate special pentru a deservi sarcina paralelizării. Astăzi calculatoarele paralele utilizează procesoare obișnuite de uz general, soluție mult mai convenabilă din punct de vedere economic dar și tehnologic computațional.

Un exemplu de procesor special dezvoltat pentru calcul paralel este **transputerul**, care apărut în anii 1980. O unitate de calcul bazată pe transputer conține (Fig. 1):

- un procesor capabil să execute operații aritmetice de bază precum și operații de intrare/ieșire;
- o memorie locală de tip SRAM de dimensiuni mici dar foarte rapidă;
- 4 canale de comunicație cu exteriorul;
- o interfață cu memoria externă.

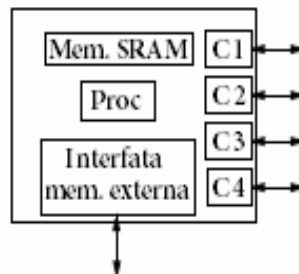


Figura 1. Structura de bază a unui transputer

Pe cele 4 canale de comunicație se pot lega alte transputere, formându-se o rețea multiprocesor cu topologie (configurație) diversă. Un exemplu de celulă cu 5 transputere este ilustrat în figura 2:

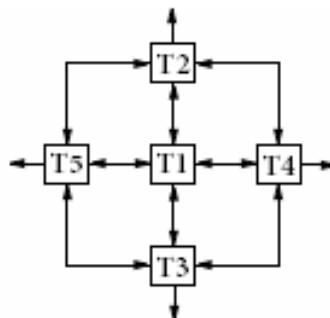


Figura 2. O celulă formată din 5 transputere

Dezavantajul folosirii transputerelor era dat de o putere computațională redusă, insuficientă pentru aplicațiile tot mai complexe care au fost dezvoltate.

În anii 1990 au apărut procesoarele **RISC** (*Reduced Instruction Set Computer* – Calculatoare cu set redus de instrucțiuni), spre deosebire de procesoarele **CISC** (*Complex Instruction Set Computer*) care existau până atunci. Aceste procesoare RISC dădeau puterea computațională de care era nevoie în aplicațiile în care ele au devenit din ce în ce mai ieftine, fiind produse pe scară largă.

Câteva exemple de procesoare RISC: procesorul MIPS de la SGI, procesorul POWER (IBM), SPARC (SUN), procesorul ALPHA (Cray).

Caracteristici ale procesoarelor RISC:

- au pu ine moduri de adresare;
- au un format fix al instruc iunilor (32b sau 64b);
- au un num r mare de registre care ajut la comutarea rapid a contextului între programe;
- folosesc instruc iuni de tip LOAD / STORE, instruc iuni dedicate pentru înc rcarea / salvarea datelor din memoria calculatorului în registrele procesorului (Fig.3):

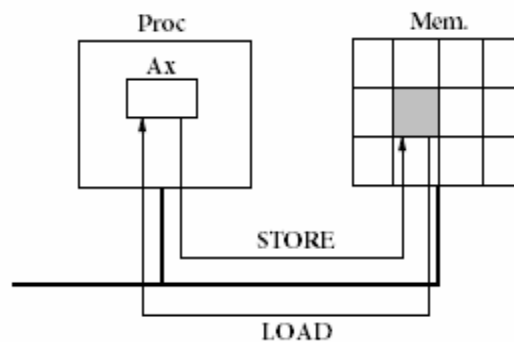


Figura 3. Instruc iunile LOAD / STORE

- utilizeaz memoria cache pentru înc rcarea datelor în procesor. Memoria cache aduce în avans instruc iunile procesorului sau datele din memoria sistem DRAM, conferind avantajul unei viteze m rite de acces fa de accesarea direct a datelor din memoria extern . (Fig.4):

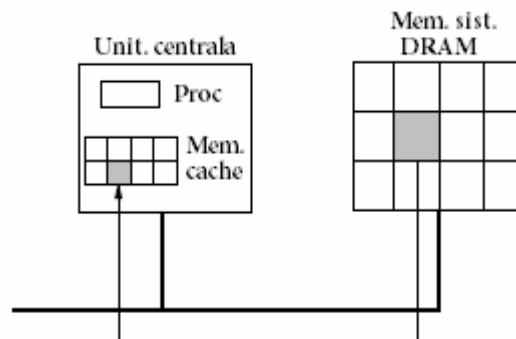


Figura 4. Memoria cache

Procesoarele RISC utilizeaz **arhitectura pipeline** (band de asamblare), în care o instruc iune sau un task se descompune în mai multe opera ii care sunt executate independent de c tre unit i de execu ie specializate (Fig. 5):

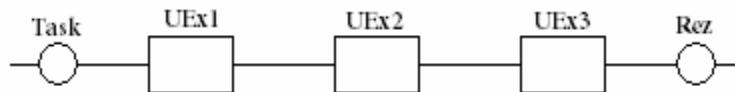


Figura 5. Arhitectura *pipeline*

În cazul procesoarelor RISC, arhitectura *pipeline* definește mai multe stadii (etape) de execuție ale unei instrucțiuni. De exemplu, putem avea următoarele 4 stadii (Fig. 6):

Stagiul 1 – IF (*instruction fetch*): încărcarea instrucțiunii din memorie

Stagiul 2 – ID (*instruction decoding*): decodificarea instrucțiunii

Stagiul 3 – EX (*execute*): execuția propriu-zisă a instrucțiunii

Stagiul 4 – WB (*write back*): scrierea rezultatului în registru sau memorie

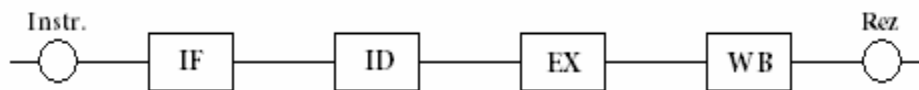


Figura 6. Stagiile de execuție ale unei instrucțiuni

Diagrama de funcționare în timp a acestei structuri este ilustrată în figura 7:

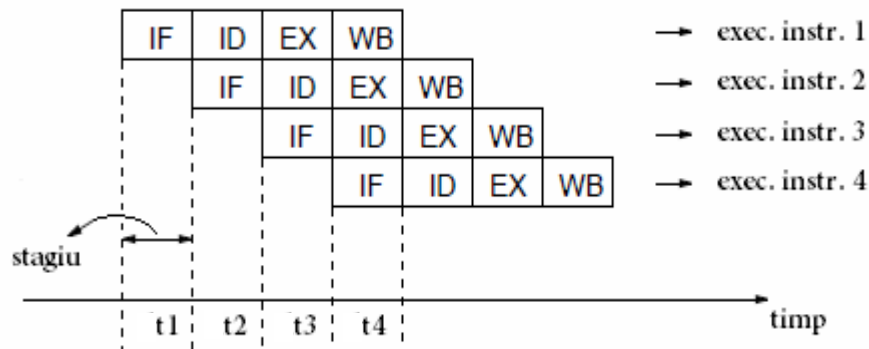


Figura 7. Diagrama execuției în timp pentru structura *pipeline*

Din această diagramă se poate observa că în orice moment de timp (de exemplu la momentul t4) se execută 4 stadii în paralel pentru 4 instrucțiuni diferite, rezultând astfel un paralelism temporal. La fiecare ciclu temporal, structura *pipeline* obține un rezultat, astfel că putem executa câte o instrucțiune la fiecare ciclu procesor (și nu în 4 cicluri cum ar fi fost fără structura *pipeline*). Acest exemplu se aplică pentru procesoarele scalare. În acest caz există o singură linie *pipeline* și se poate executa maxim o instrucțiune pe ciclu.

În prezent, majoritatea procesoarelor au o structură superscalară. La această structură există mai multe linii *pipeline* care operează în paralel. Ca rezultat, se pot executa mai multe instrucțiuni pe un ciclu procesor.

În figura 8 se prezintă o structură superscalară cu 3 linii *pipeline*:

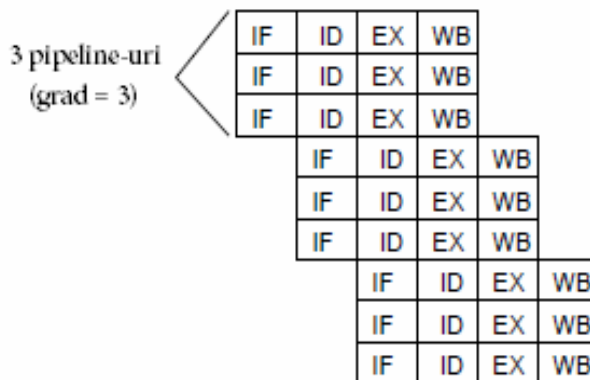


Figura 8. Diagrama de timp pentru structura superscalară

Câteva caracteristici ale arhitecturii superscalare:

- arhitectura permite executarea în paralel a mai multor instrucțiuni;
- se pot executa mai multe instrucțiuni pe un ciclu procesor;
- permite execuția *out of order* (în afara ordinii) a instrucțiunilor: instrucțiuni independente pot fi executate în afara ordinii normale din program; condiția pentru execuția *out of order* este să nu existe dependențe de date între instrucțiuni.

O altă metodă pentru creșterea eficienței *pipeline* este folosirea arhitecturii *superpipeline*. În acest caz se reduce numărul de operații care se efectuează într-un stadiu, dar se mărește numărul de stagii prin împărțirea fiecărui stadiu în mai multe substagii. Astfel putem crește frecvența de tact a procesorului, pentru că se execută mai puține operații într-un ciclu procesor decât la un stadiu obișnuit.

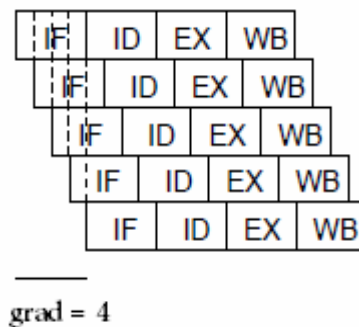


Figura 9. Diagrama de timp pentru structura *superpipeline*

În exemplul de mai sus, datorită faptului că un stadiu a fost împărțit în 4 substagii, vom putea crește frecvența de tact a procesorului de 4 ori.

Avantajele arhitecturii *superpipeline*:

- crește viteza de execuție a instrucțiunilor datorită creșterii frecvenței de ceas a procesorului;
- crește numărul de operații procesate în paralel datorită creșterii numărului de stagii;
- crește eficiența operațiilor executate într-un stadiu. De exemplu, dacă într-un *pipeline* obișnuit avem o operație care se execută într-un timp mai mic decât un ciclu procesor, rezultatul va fi disponibil tot la sfârșitul ciclului (Fig. 10 a). În schimb, la arhitectura *superpipeline*, rezultatul unei operații este disponibil cu frecvența tactului, care acum este mai rapid (Fig. 10 b):

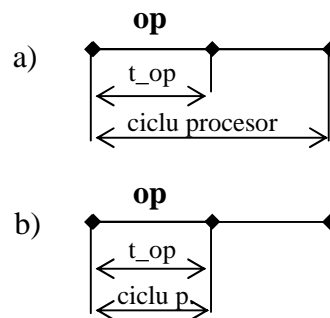


Figura 10. Execuția unei operații corespunde toare unui substadiu: a) într-o structură *pipeline*; b) într-o structură *superpipeline*

Dezavantaje ale arhitecturii *super pipeline*:

Principalul dezavantaj este reprezentat de faptul că dacă *pipeline*-ul trebuie golit de instrucțiunile încărcate în avans (de exemplu când apar instrucțiuni de salt în program), atunci există un număr mai mare de stagii care trebuie reinițializate.

Așa cum am văzut, procesoarele RISC sunt procesoare de uz general, dar există și alte tipuri de procesoare care se aplică unui domeniu mai restrâns, obținând performanțe mai bune pe acel domeniu specific (de exemplu procesoare DSP, procesoare VLIW și procesoare vectoriale). În continuare vor fi prezentate pe scurt procesoarele VLIW și procesoarele vectoriale.

Procesoarele VLIW (Very Large Instruction Word)

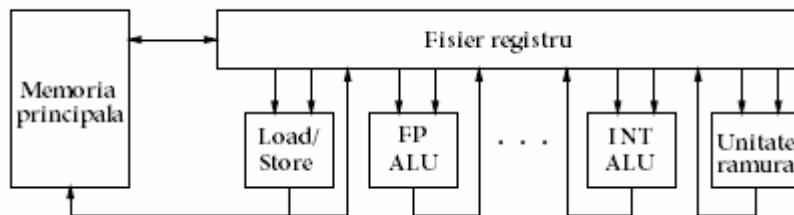


Figura 10. Structura internă a unui procesor VLIW

Procesorul VLIW este alcătuit din mai multe unități funcționale echivalente unităților de execuție din procesoarele superscalare, și un fișier de registre care conține un număr mare de registre procesor (pot exista mai mult de 128 registre). O caracteristică a acestor procesoare este că ele dispun de un format mare al instrucțiunilor. Fiecare instrucțiune este compusă din mai multe sloturi de operații în care se plasează diferite operații de tip RISC, în așa fel încât toate operațiile dintr-o instrucțiune să poată fi executate în paralel.

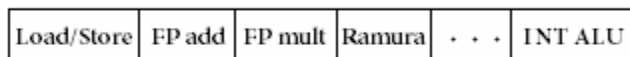


Figura 11. Formatul instrucțiunii la procesorul VLIW

Programul compilator este cel care umple instrucțiunile cu operații care se pot executa în paralel. De aici rezultă caracteristica de *static scheduling* (planificare statică), adică determinarea paralelismului se face în timpul compilării și nu al execuției. Rezultă astfel un avantaj al procesoarelor VLIW: ele nu necesită un hardware complex, cum era cazul la procesoarele superscalare (unde determinarea paralelismului se făcea prin hardware în timpul execuției).

Ca și dezavantaj al procesoarelor VLIW: codul obiect rezultat în urma compilării este mai puțin compact decât la un procesor obișnuit, datorită faptului că nu toate sloturile dintr-o instrucțiune pot fi încărcate cu operații care se execută în paralel.

Diagrama execuției în timp a unei instrucțiuni este ilustrată în figura 12 (pentru un procesor având 3 sloturi pe instrucțiune).

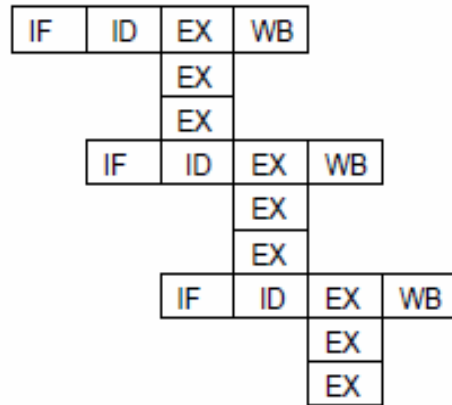


Figura 12. Diagrama execuției unei instrucțiuni la procesorul VLIW

Câteva exemple de procesoare VLIW:

- Philips Trimedia (utilizat pentru aplicații multimedia);
- arhitectura Intel pe 64 biți (IA64) - o combinație între RISC și VLIW.

Procesoare vectoriale

În cele mai multe cazuri, procesoarele vectoriale nu funcționează independent, ci sunt folosite pe post de coprocesoare alături de un procesor principal. Spre deosebire de procesoarele obișnuite, care operează cu scalari, procesoarele vectoriale pot opera pe vectori, adică pe grupe de scalari. Ele pot fi de tipul registru-registru, în care vectorii se încarcă din registre vectoriale, iar rezultatele se salvează tot în registre vectoriale, sau pot să fie de tipul memorie-memorie, în care vectorii și rezultatele se iau / se pun în memoria sistemului. Unitățile de execuție vectoriale au o funcționare de tip *pipeline* (fig.13).

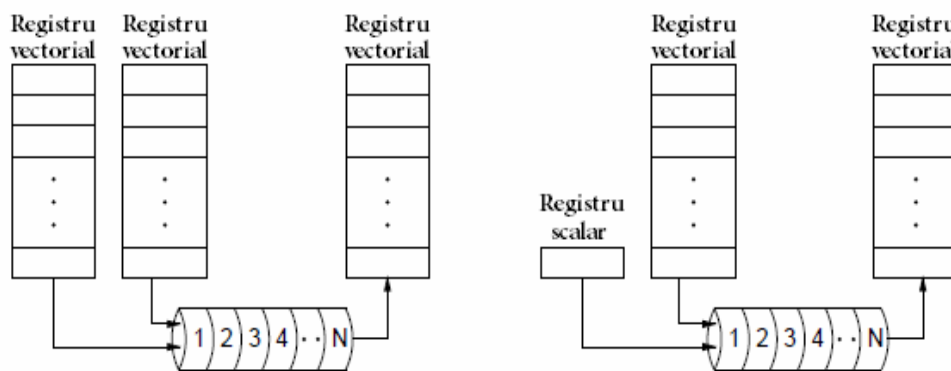


Figura 13. Unitățile de execuție într-un procesor vectorial

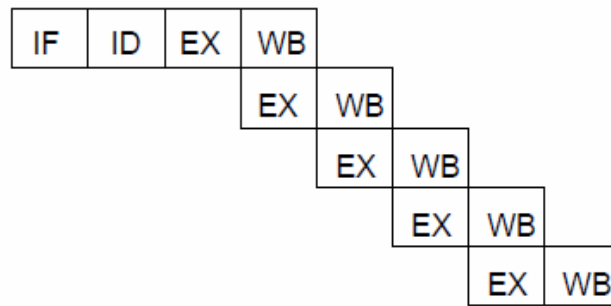


Figura 14. Diagrama de execuție a unei instrucțiuni vectoriale

Domeniul de aplicabilitate al procesoarelor vectoriale este reprezentat de calculele științifice și aplicațiile multimedia, unde este necesară realizarea unui număr mare de operații pe structuri vectoriale.

MEMORIA CACHE

Memoria cache a apărut din necesitatea de a elimina diferența de viteză dintre procesorul rapid și memoria sistem, mult mai lentă. Este o memorie statică (SRAM) de dimensiuni mici dar foarte rapidă, folosită pentru a aduce datele mai aproape de unitățile de execuție din procesor. Ca urmare a performanțelor de funcționare a memoriei cache, actualmente doar un procent de aproximativ 10% din accesările la memorie ale procesorului se fac din memoria sistem, iar restul de 90% se fac din memoria cache.

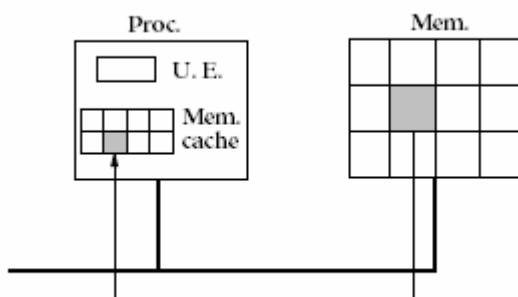


Figura 1. Memoria cache

Memoria cache se bazează pe principiul localității. Avem două tipuri de localitate:

- localitate temporală – locațiile de memorie accesate la un moment dat tind să fie accesate din nou în viitor;
- localitate spațială – datele aflate în vecinătatea unor date accesate tind să fie și ele accesate.

Exemple de memorii cache:

- memoria TLB (*Translation Look-aside Buffer*) este un cache special pentru traducerea între adresele fizice și adresele virtuale din calculator. Datorită folosirii memoriei TLB, calculul adreselor pentru memoria virtuală nu se face la fiecare acces al memoriei, ci se ia din tabela TLB.

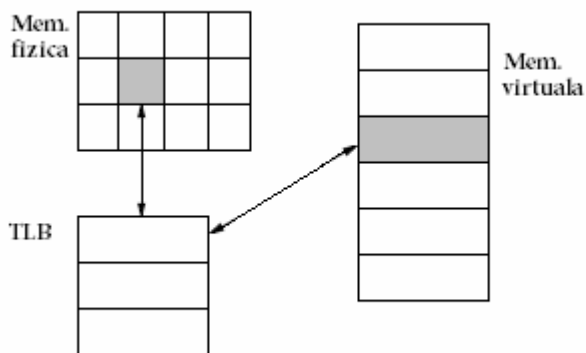


Figura 2. Memoria TLB

- memoria cache pentru instrucțiuni - este folosită pentru încărcarea în avans a instrucțiunilor în procesor ;
- memoria cache de date - folosită pentru încărcarea în avans în datelor necesare din memoria sistem.

Există 3 moduri de mapare a datelor din memoria fizică în memoria cache:

- 1) mapare directă
- 2) mapare cu asociere completă
- 3) mapare cu asociere pe blocuri (set asociativ)

1) Memoria cache cu mapare directă

În acest caz există o mapare directă, secvențială între adresele din memoria principală și locațiile memoriei cache. Aici memoria fizică este împărțită în pagini, fiecare pagină fiind de dimensiunea memoriei cache. Atâtea paginile din memoria sistem cât și memoria cache sunt împărțite în blocuri de aceeași dimensiune.

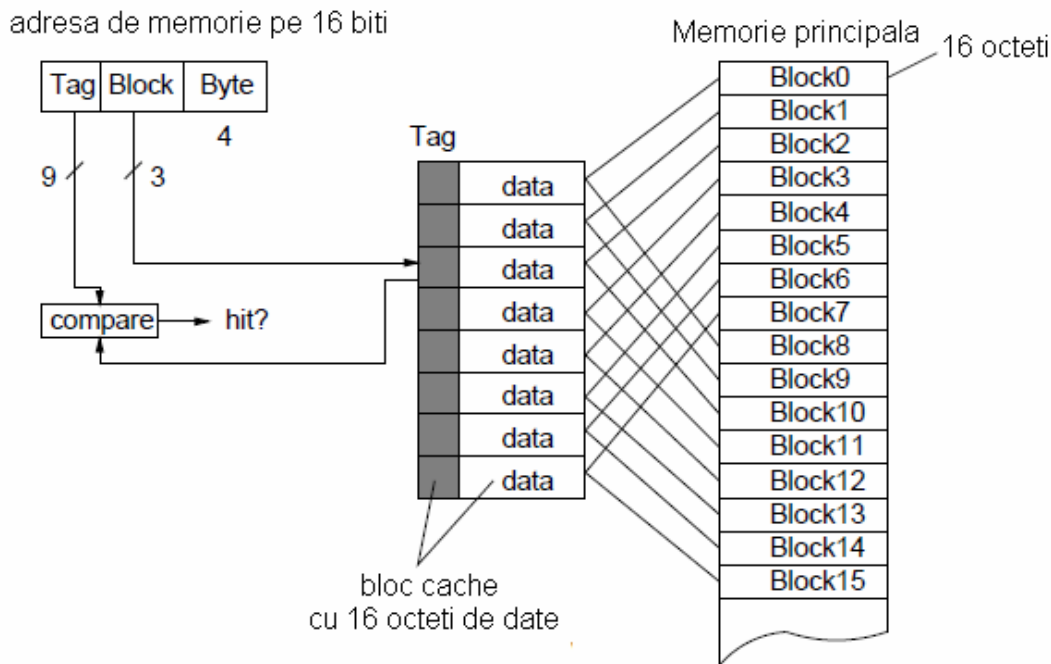


Figura 3. Memoria cache cu mapare directă

În exemplul de mai sus, o pagină conține 8 blocuri; astfel, blocul 0 din memoria cache corespunde cu blocurile 0, 8, 16, etc. din memoria sistem.

Pentru a ști exact cu ce bloc din memorie corespunde un anumit bloc din memoria cache, se folosește un tag (indicator) la fiecare intrare din memoria cache. Tag-ul conține cei mai semnificativi biți ai adresei blocului din memoria principală, reprezentând de fapt adresa de pagină a blocului.

Dacă procesorul vrea să acceseze un bloc de date din memorie situat la o anumită adresă, mai întâi trebuie să verifice dacă blocul cerut se află sau nu în memoria cache. Pentru aceasta, se va proceda astfel:

- a) se selectează o intrare din memoria cache pe baza informației de bloc din adresa cerută;

Tag	Block	Byte
-----	-------	------

- b) se compară câmpul tag din adresa cerută de procesor cu cel din intrarea cache;
- c) dacă cele două sunt egale, avem așa-numitul cache-hit, adică blocul cerut din memoria principală se regăsește în memoria cache;
- d) dacă cele două tag-uri sunt diferite, atunci blocul cerut nu există în memoria cache și el va fi încărcat din memoria principală.

Avantajul modului de mapare direct :

Memoriile cu mapare directă sunt memorii simple și de mare viteză. Ele se pot utiliza ca memorii cache de nivel 2, putând avea dimensiuni mai mari.

2) Memorii cache cu asociere completă

Aici, un bloc din memoria sistem poate fi plasat oriunde în memoria cache.

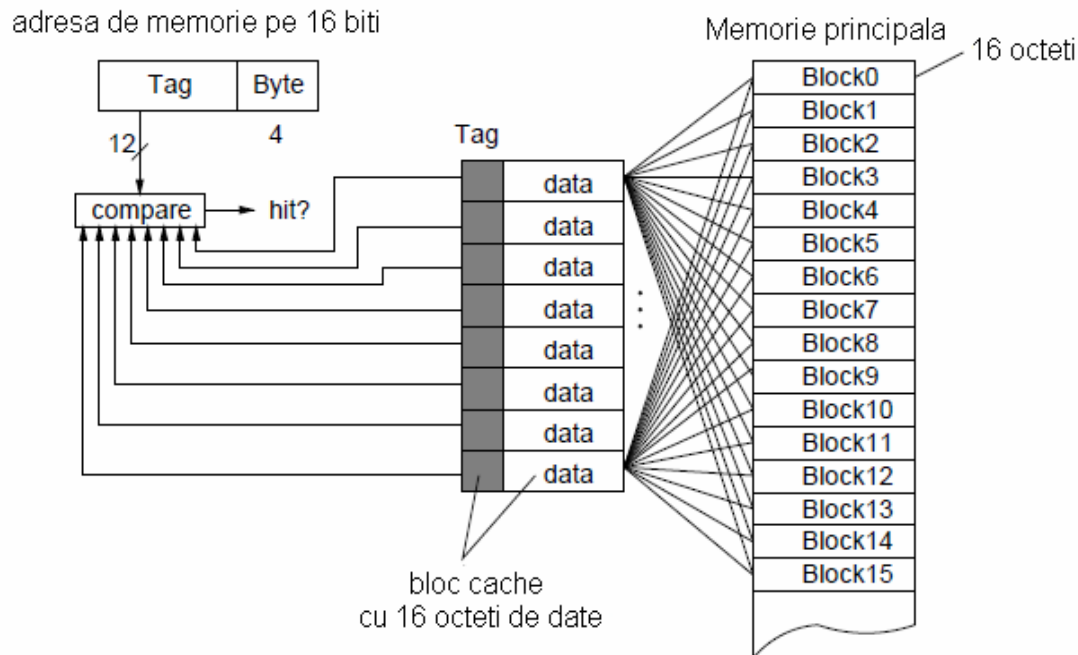


Figura 4. Memoria cache cu asociere directă

Pentru a detecta dacă un bloc din memoria sistem există în memoria cache, se face o comparație în paralel a tag-urilor din toate blocurile memoriei cache. Dacă se identifică un tag, înseamnă că s-a găsit blocul căutat.

Avantajul acestui tip de memorie: folosește un mod de mapare mai flexibil.

Dezavantaj: este mai costisitor de implementat datorită comparațiilor care se fac în paralel.

3) Memorii cache set asociative

Acestea folosesc o combinație între maparea directă și maparea cu asociere completă. Aici, memoria cache este împărțită în mai multe seturi de înțiriri, existând o mapare directă între blocurile din memoria principală și seturile din memoria cache. În interiorul seturilor se face o mapare cu asociere completă.

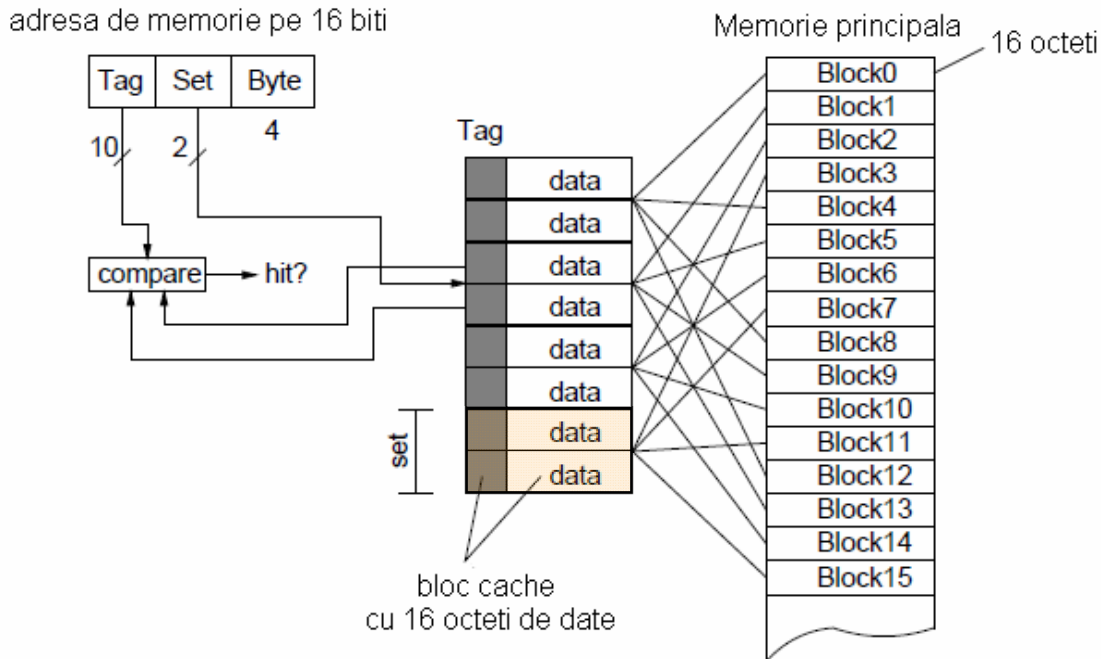


Figura 5. Memoria cache set-asociativ

Pentru a vedea dacă blocul dorit se găsește în memoria cache, se face o comparație în paralel a tag-urilor doar pentru înțirile dintr-un singur set.

Avantajul memoriilor set-asociative: se reduce numărul de comparații efectuate în paralel; sunt mai puțin costisitoare decât memoriile complet asociative; se folosesc ca memorii cache de nivelul 1, având performanțe ridicate.

Observa ie: În sistemele actuale exist o ierarhie de memorii cache organizat pe mai multe nivele; în func ie de pozi ia memoriei cache fa de procesor, putem avea: memorie cache de nivel 1 (care se g se te în interiorul capsulei procesorului), memorie cache de nivel 2, de nivel 3, etc.

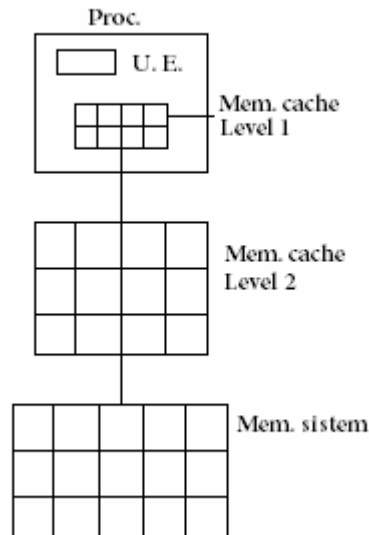


Figura 6. Ierarhia de memorii cache într-un sistem de calcul

Strategii utilizate pentru memoriile cache:

În leg tur cu anumite evenimente ce pot sa apar la accesarea memoriei cache, se stabilesc urm toarele strategii:

S1) Strategia *cache-miss*. Aceast strategie se refer la memorii cache complet asociative. Astfel dac e nevoie de un bloc de memorie care nu se reg se te în cache (*cache-miss*), trebuie s se stabileasc care anume dintre blocurile existente în memoria cache va fi înlocuit de noul bloc care va fi adus din memoria sistem.

Strategia de înlocuire poate fi:

- aleatorie
- *first in, first out* - se înlocuie te primul bloc introdus în cache
- se înlocuie te cel mai pu in recent utilizat. Aici se înlocuie te blocul care nu a fost accesat de cel mai mult timp (cel mai vechi).

S2) Strategii de scriere a memoriilor cache. S presupunem c procesorul efectueaz o opera ie STORE (scriere în memorie). Exist 2 posibilit i:

a) valoarea se scrie în memoria cache și în memoria sistem. Aceste memorii se numesc *write-through*; valoarea respectiv se scrie prin cache în memoria sistemului.

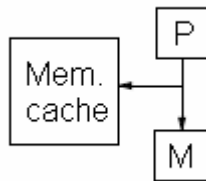


Figura 7. Memoria *write-through*

b) scrierea se face doar în memoria cache. Scrierea în memoria sistem se face ulterior doar când este necesar acest lucru. Aceste memorii se numesc *write-back*.

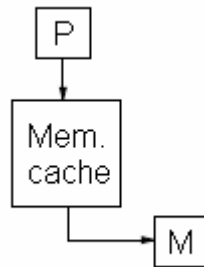


Figura 8. Memoria *write-back*

Avantajele memoriei *write-through*:

- memoria principală este întotdeauna consistentă cu memoria cache.

Dezavantajele memoriei *write-through*:

- această memorie este mai încetă pentru că se așteaptă scrierea în memoria sistem. La o operație STORE se așteaptă scrierea în memoria sistem, care este mai lentă.
- se mărește traficul între magistrală și memoria principală datorită faptului că fiecare acțiune STORE a procesorului se propagă în memoria principală.

Acest lucru nu se petrece la memoria *write-back*. În schimb, la acest tip de memorie este necesar un bit suplimentar pentru fiecare bloc din memoria cache, care să specifice dacă blocul este consistent cu memoria principală. Acest bit se numește bit *dirty*.

Dacă bitul *dirty* este 0, înseamnă că blocul din memoria cache este consistent (identic) cu blocul din memoria sistem. Dacă bitul *dirty* este 1, înseamnă că blocul nu este consistent (nu este actualizat în memoria sistem). Când un astfel de bloc *dirty* trebuie să fie eliminat din memoria cache, el va fi scris mai întâi în memoria principală.

Dezavantaj al memoriei *write-back*:

La schimbarea contextului program de către un sistem de operare multi-tasking se folosește întotdeauna un context diferit de memorie, de aici rezultă că un mare număr de blocuri din memoria cache vor trebui să fie înlocuite și astfel timpul de comutare a contextului program crește foarte mult.

S3) Strategii *write-miss*. Evenimentul *write-miss* apare când exist un *cache-miss* la o instruc iune STORE, adic blocul dorit nu se afl în memoria cache atunci când procesorul vrea s scrie în acesta.

Memoriile *write-back* folosesc 2 tipuri de strategii:

- a) Strategia “aloc la scriere” (*allocate on write*). Aici memoria cache aloc un bloc cache, dup care realizeaz ac iunea de a scrie în acel bloc.
- b) Strategia “încarc la scriere” (*fetch on write*). Aici memoria cache mai întâi cite te blocul din memoria principal , apoi scrie data în blocul citit.

Memoriile *write-through* folosesc strategia “nu se aloc la scriere” (*no allocate on write*). Aici, dac exist un *cache-miss* la scriere, atunci nu se mai aloc un bloc în memoria cache, ci data trimis de procesor va fi scris direct în memoria principal a sistemului.

RE ELE DE INTERCONECTARE

Scopul unei re ele de interconectare este de a permite schimbul de date între procesoarele unui sistem paralel. Exist două no iuni fundamentale legate de re ele: comutarea i rutarea.

1. Comutarea (*network switching*)

Se refer la modul de transmitere a datelor între procesoarele din re ea. Exist două tipuri principale:

- a) Comutarea prin circuite
- b) Comutarea prin pachete

Comutarea prin circuite se refer la transmiterea datelor prin circuite dedicate. Aici se stabile te o leg tur fizic între un procesor surs i un procesor destina ie, leg tur ce trebuie s r mân stabil pe întreaga durat a transmisiei datelor.

Comutarea prin pachete - aici datele de transmis se împart în blocuri de dimensiuni mici numite pachete, iar apoi se aloc un canal de comunica ie doar pentru transmiterea unui singur pachet.

2. Rutarea (*network routing*)

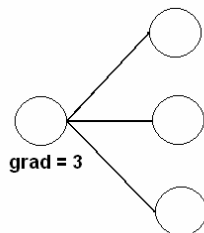
Se refer la modalitatea de a conduce datele prin re ea. Rutarea define te astfel o rut pe care datele o vor urma pentru a ajunge la destina ie. Rutarea este în strâns leg tur cu topologia (structura) re elei.

Clasificarea re elelor în func ie de topologie:

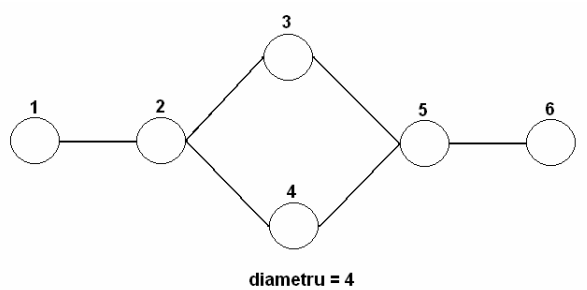
- a) re ele directe; aici exist o conectivitate punct la punct între procesoarele vecine
- b) re ele indirecte; acestea utilizeaz canale de comunica ie care sunt comune mai multor procesoare

No iuni fundamentale legate de re ele:

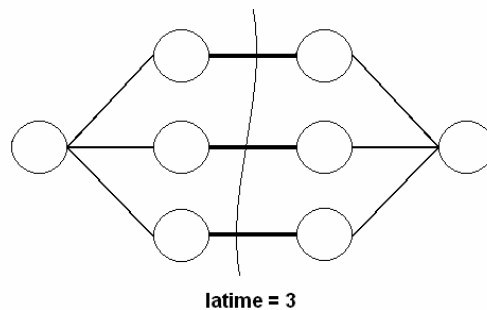
- gradul unui nod: reprezint num rul de canale de comunica ie care sunt legate de acel nod.



- diametrul re elei: reprezint distan a maxim între două noduri ale re elei.



- lăimea secțiunii: reprezintă numărul de canale de comunicație între două jumătăți ale rețelei.

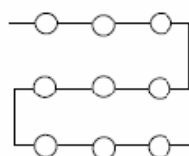


- toleranța la erori: reprezintă numărul de rute alternative între două noduri.
- scalabilitatea rețelei: reprezintă posibilitatea de expansiune a rețelei.
- rata de transfer: reprezintă cantitatea de date transferate în unitatea de timp.
- latența: reprezintă durata maximă de transfer a unei date prin rețea.

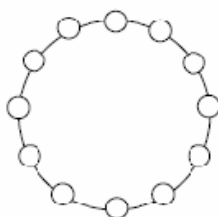
A. Rețele directe

În acest caz există legături punct la punct între nodurile învecinate ale rețelei. Ele sunt numite rețele statice pentru că aceste legături punct la punct definite sunt fixe (nu se schimb pe întreaga durată a transmisiei datelor).

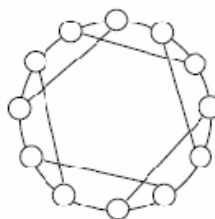
Exemple:



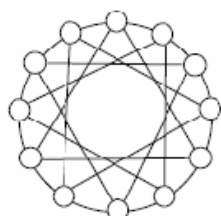
Linear array



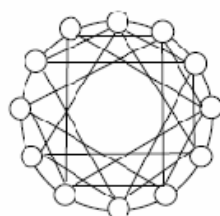
Ring



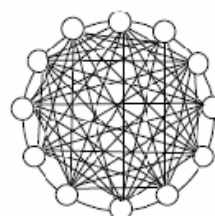
Chordal ring of degree 3



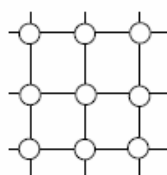
Chordal ring of degree 4
(same as Illiac mesh)



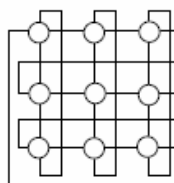
Barrel shifter



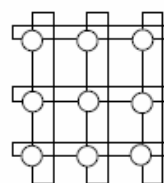
Completely connected



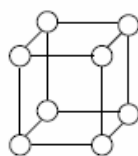
Mesh



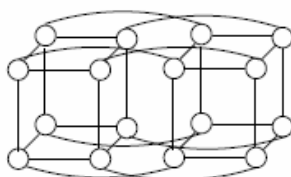
Illiach mesh



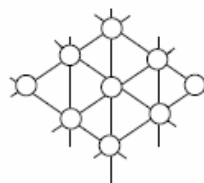
Torus



3-Hypercube



4-Hypercube



Systolic array

Tip de retea	Gradul	Diametrul	Latimea Sectiunii
Linear array	2	$N - 1$	1
Ring	2	$\lfloor \frac{N}{2} \rfloor$	2
Completely conn.	$N - 1$	1	$(\frac{N}{2})^2$
Binary tree	3	$2(\log_2 N - 1)$	1
2D-mesh	4	$2(\sqrt{N} - 1)$	\sqrt{N}
2D-torus	4	$2\lfloor \frac{\sqrt{N}}{2} \rfloor$	$2\sqrt{N}$
Hypercube	$\log_2 N$	$\log_2 N$	$\frac{N}{2}$

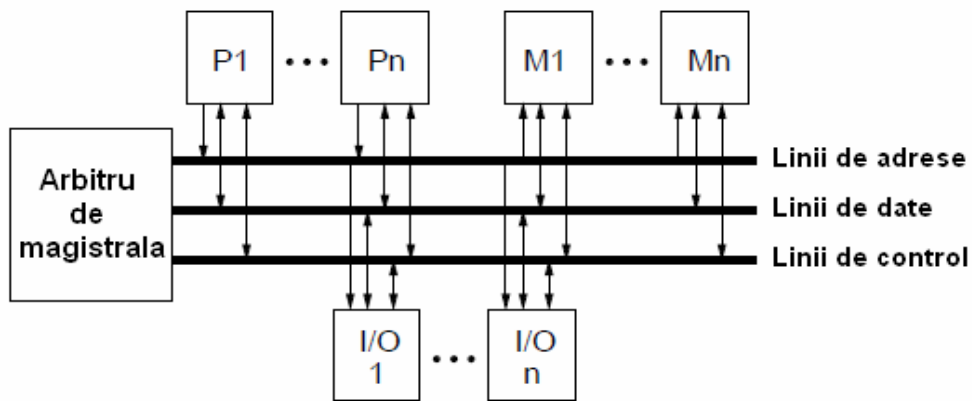
B. Re ele indirecte

Re elele indirecte mai sunt numite i re ele dinamice. Spre deosebire de re elele statice, aici un nod nu are vecini fixa i. Topologia re elei poate fi schimbat dinamic în func ie de cerin ele aplica iei. Exist trei tipuri de re ele dinamice:

- a) Re ele cu magistral
- b) Re ele multi-stagiu
- c) Re ele *crossbar* (cu comutatoare)

a) Re ele bazate pe magistral :

Ele se bazeaz pe o magistral alc tuit din mai multe linii de bit, la care se conecteaz diferite resurse (procesoare, memorii, dispozitive de intrare/ie ire etc). Magistrala poate con ine linii de adres , de date i de control.



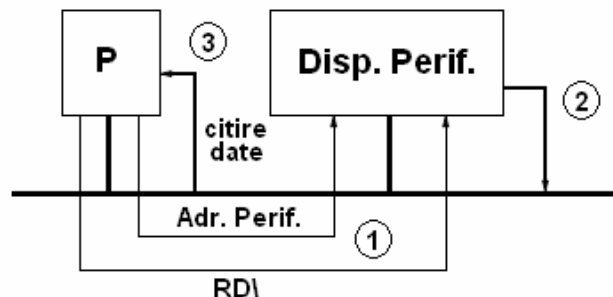
La o magistral se pot conecta atât dispozitive de tip *master* (care ini iaz un transfer de date pe magistral), cât i dispozitive de tip *slave*. Dacă la magistral sunt conectate mai multe dispozitive *master*, atunci va fi necesar un dispozitiv de arbitrare a magistralei.

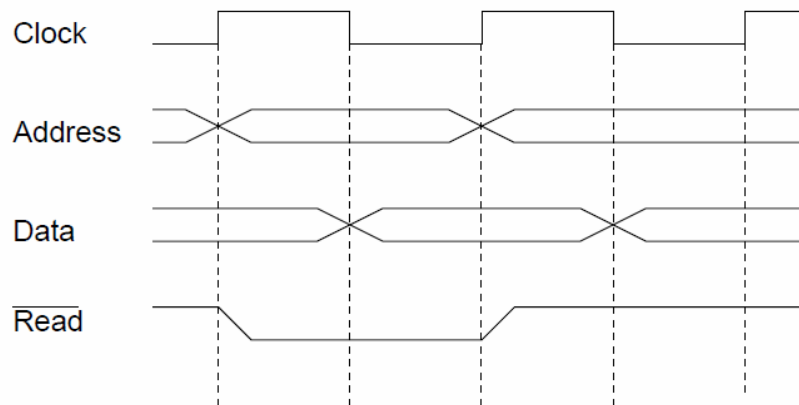
Exist două tipuri de magistrale:

- 1) Magistrale sincrone
- 2) Magistrale asincrone

1) Magistrale sincrone

Aici, toate tranzac iile pe magistral sunt sincronizate folosind un semnal de ceas. Exemplu: citirea unei date de la un dispozitiv periferic.





2) Magistrale asincrone

Aici, acțiunile nu sunt la fel de predictibile ca la o magistrală sincronă. La magistrala asincronă se utilizează un protocol de tip *handshaking*, pentru a se indica apariția unui eveniment pe magistrală, de exemplu validitatea datelor pe linia de date a magistralei poate fi indicată de un semnal de comandă numit *Strobe*.

În general, magistralele asincrone sunt mai complexe, mai scumpe și mai puțin eficiente decât magistralele sincrone datorită protocolului care trebuie realizat cu semnale fizice. Totuși, magistralele asincrone au avantajul că sunt mai flexibile și prin intermediul lor se pot ușor conecta resurse cu viteze diferite.

Exemplu de magistrală sincronă: magistrala procesor-memorie.

Exemplu de magistrală asincronă: magistrala SCSI, care asigură comunicația cu dispozitivele periferice.

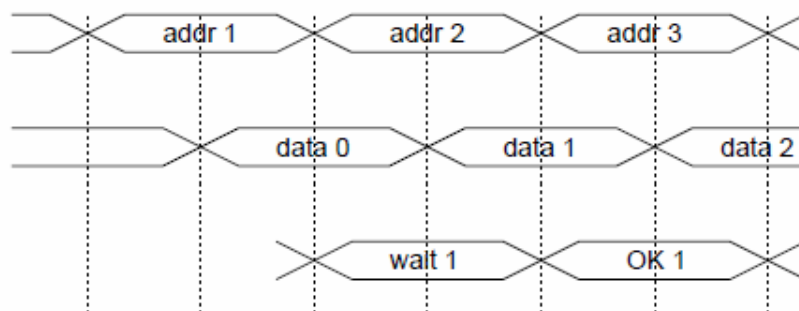
Magistrale *split transaction*:

În cazul unei magistrale obișnuite, magistrala rămâne ocupată pe întreaga durată a unei tranzacții, chiar și atunci când magistrala așteaptă date de la un periferic, lucru care este inefficient. Pentru a mări eficiența magistralei, s-au creat așa numitele magistrale cu tranzacții multiple (*split transaction bus*). În acest tip de magistrală, pot fi active la un moment dat mai multe tranzacții diferite.

Exemplu:

Considerăm o tranzacție de citire de date de la un dispozitiv periferic pe o magistrală sincronă cu tranzacții multiple. Această tranzacție funcționează astfel:

- se pune adresa perifericului de la care se citesc datele pe magistrala de adrese, și semnalul de citire READ pe magistrala de comenzi.
- după ce adresa este citită de dispozitivul periferic, liniile de adresă sunt pregătite să primească o altă adresă de la un alt periferic.
- când dispozitivul periferic este pregătit să furnizeze datele cerute, el va plasa aceste date pe magistrala de date și va specifica destinatarul datelor prin intermediul unei etichete (*tag*) furnizate pe linii suplimentare de comandă.



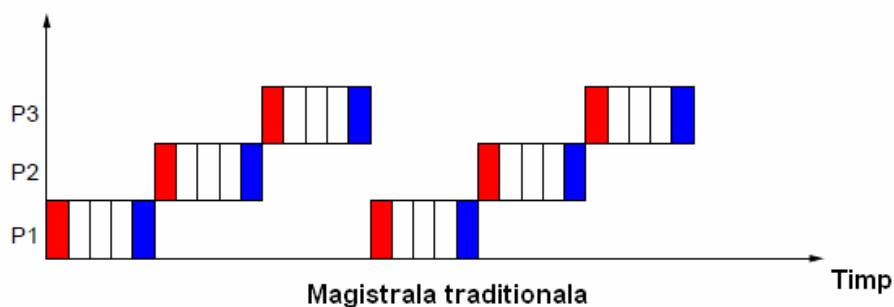
Magistrala cu tranzactii multiple

Avantaj: magistrala *split transaction* m re te rata de transfer a datelor.

Dezavantaj: poate s m reasca laten a tranzac iilor, pentru c acum magistrala multiplexeaz mai multe tranzac ii n acela i timp.

■ = se foloseste magistrala de adrese
 = nu e folosita magistrala
 ■ = e folosita magistrala de date

Procesoare



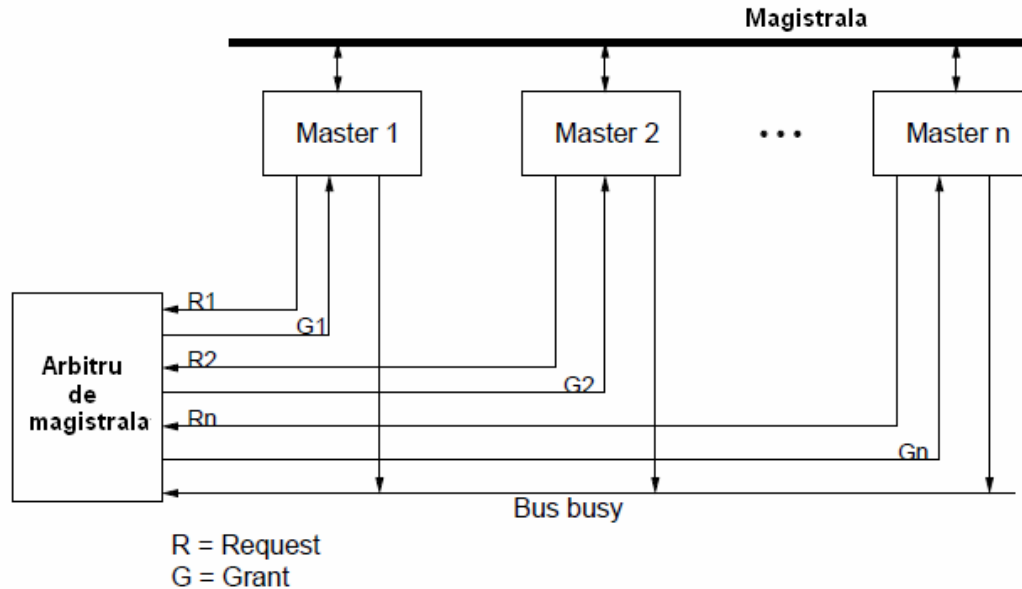
Procesoare



Arbitrarea magistralei:

Arbitrarea este necesar atunci cnd exist mai multe dispozitive *master* conectate pe magistral . Dac dou sau mai multe dispozitive *master* vor s acceseze magistrala n acela i timp, controlerul de magistral va decide c rui dispozitiv i se va permite mai ntm accesul. Exist dou modalit i de arbitrare:

- a) Schema de arbitrare bazat pe un semnal de *Request* (cerere) și *Grant* (acceptare).
 Când un *master* vrea să acceseze magistrala, trimite o cerere către arbitru (controler). În momentul în care arbitrul permite accesul, el va trimite dispozitivului *master* semnalul de grant. Dacă există mai multe dispozitive care vor să acceseze simultan magistrala, atunci va fi aplicat o politică de selecție bazată pe prioritate.

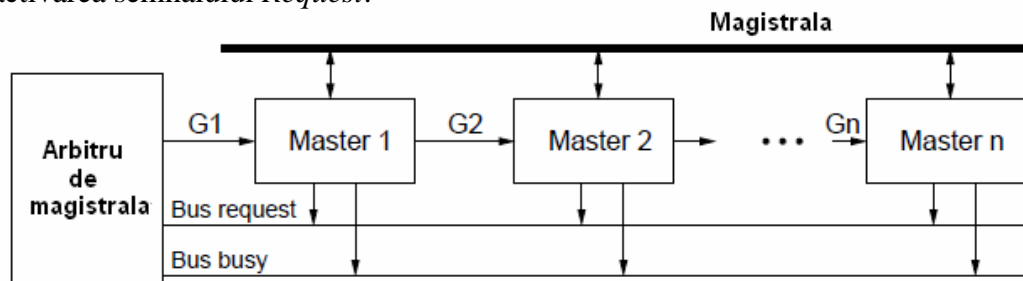


Când un *master* a câștigat accesul la magistrală, el va activa semnalul *Bus Busy*. Cât timp acest semnal este activ, arbitrul nu va mai acorda acces la magistrală unui alt dispozitiv *master*.

Avantaj: eficiență în funcționare.

Dezavantaj: reprezintă o implementare scumpă datorită unui număr mare de linii de control.

- b) Schema de arbitrare *daisy-chaining* (arbitrare cu înlanțuire). Aici există un singur semnal de grant care se propagă de-a lungul unui lanț de *master*e. Când un *master* primește un semnal de grant, el poate face o cerere de acces la magistrală prin activarea semnalului *Request*.



Dacă un *master* nu vrea să acceseze magistrala, el nu mai departe semnalul de grant către următorul *master*.

Avantaj: este mai puțin costisitoare decât schema precedentă.

Dezavantaj: datorită faptului că semnalul de grant circulează mai lent pe lanțul de *master*e, va dura mai mult până când o cerere de acces la magistrală va fi acceptată.

Rețele Multistagiu

O rețea multistagiu constă din mai multe nivele sau stagii de *switch*-uri (comutatoare) conectate între ele. Un *switch* asigură legături interne între mai multe intrări și mai multe ieșiri.

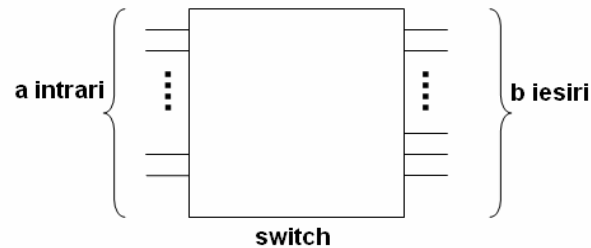


Figura 1. Un *switch* asigură conexiunea dintre a intrări și b ieșiri.

Într-o rețea multistagiu, cu cât numărul de stagii este mai mare, cu atât întârzierea rețelei va fi mai mare. Alegând un anumit tipar de interconectare a *switch*-urilor, se pot realiza diferite topologii de rețele.

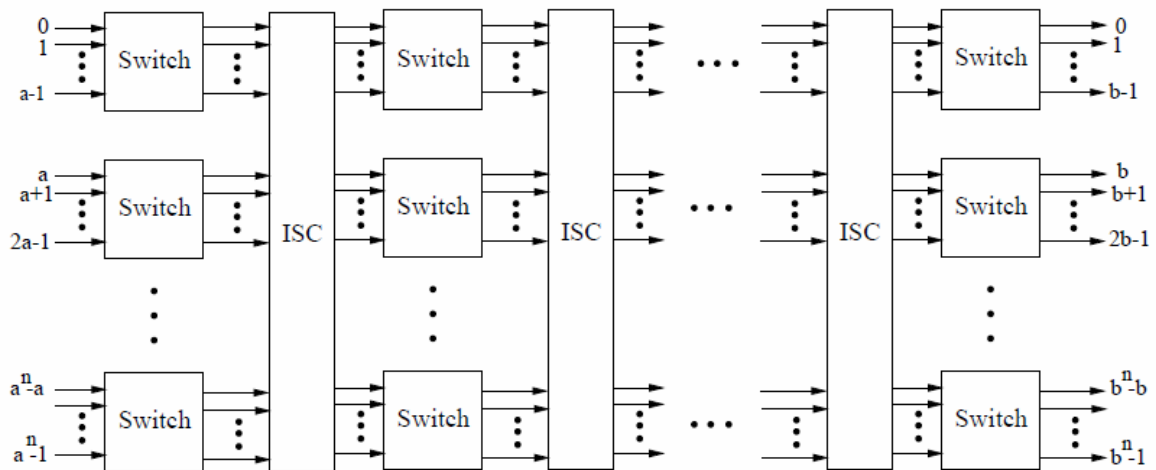


Figura 2. Generalizare a unei rețele multistagiu interconectate (MIN), construită cu *switch*-uri $a \times b$ și un tipar pentru conexiunile interstagiului (ISC).

Un exemplu de rețele multistagiu sunt rețelele de Omega (fig. 3). Aici, conectarea *switch*-urilor realizează o distribuție completă a semnalelor de intrare către ieșiri (de la orice intrare se poate ajunge la orice ieșire).

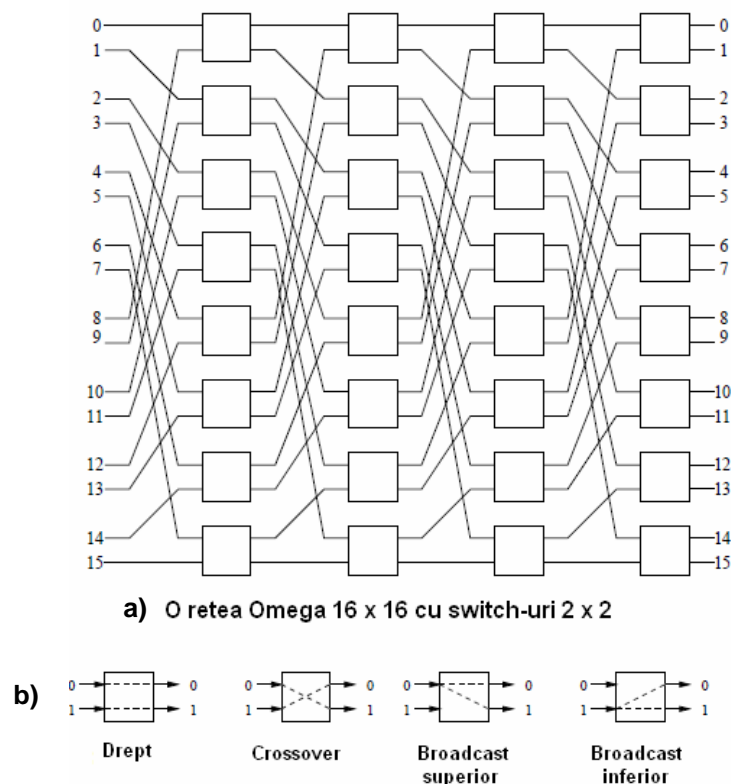


Figura 3. a) Exemplu de retea Omega. b) Conexiuni acceptate ale *switch*-urilor

Rețelele multistagii pot să fie **blocante** sau **neblocante**. Într-o rețea blocantă nu se pot realiza simultan toate legăturile de la intrare la ieșire din cauza unor conflicte de comutare internă a *switch*-urilor. La *switch*-urile unde apar astfel de conflicte de comutare, o parte din semnalele de la intrare sunt blocate până când canalele de ieșire necesare devin disponibile. Rețelele neblocante suportă orice conexiune fără blocare. De exemplu, rețelele Omega sunt rețele blocante (fig. 4):

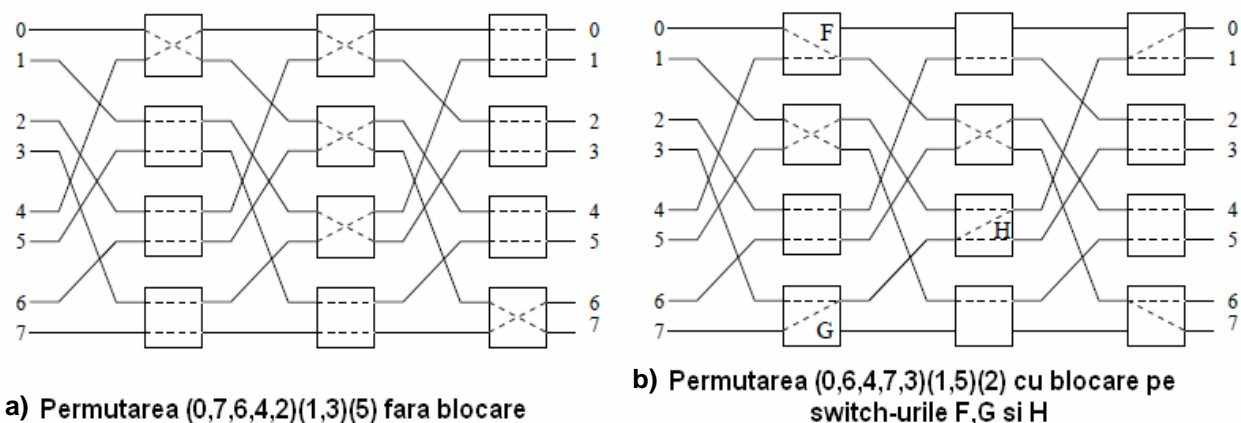


Figura 4. Exemplu de conexiuni într-o rețea Omega.
a) fără blocare; b) cu blocare

În figura 4 a), permutarea (0,7,6,4,2)(1,3)(5) semnific următoarele conexiuni (neblocante):

0->7; 7->6; 6->4; 4->2; 2->0; 1->3; 3->1; 5->5.

În figura 4 b), conexiunile dorite sunt blocante:

0->6; 6->4; 4->7; 7->3; 3->0; 1->5; 5->1; 2->2.

Rețea Crossbar

O rețea de tip *crossbar* constă dintr-o matrice de comutatoare simple de tip *on/off* care pot realiza sau întrerupe o legătură de la o intrare către o ieșire:

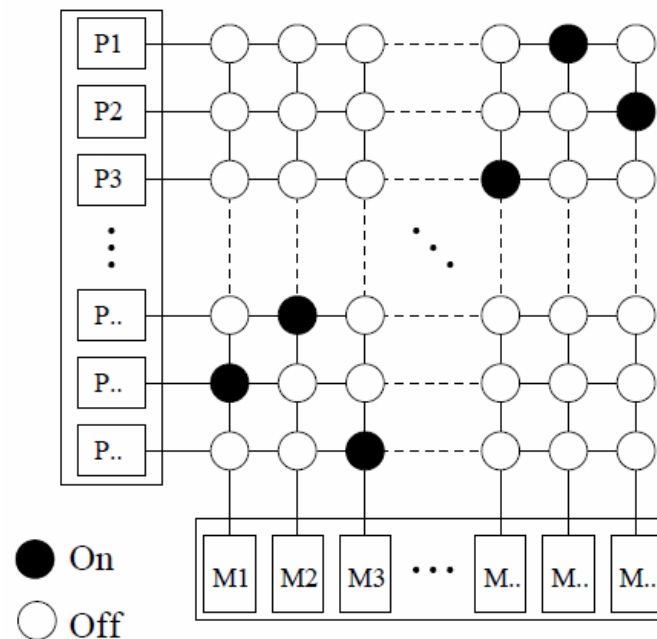


Figura 5. Rețea Crossbar

Într-o rețea de tip *crossbar*, pe o linie verticală nu putem avea mai multe *switch*-uri setate pe *on*, adică nu se poate ca mai multe procesoare să acceseze simultan același bloc de memorie. Pe o linie orizontală putem avea mai multe *switch*-uri setate, rezultând în acest caz o emisie *multicast* de la un procesor la mai multe blocuri de memorie.

Avantaj: comunicarea prin *switch*-uri *crossbar* este neblocaant, adică toate combinațiile de legături între procesoare și memorii sunt posibile fără conflicte (respectând condițiile anterioare).

Sunt posibile 3 implementări hardware pentru o rețea de tip *crossbar*(fig.6):

- folosind un lanț de multiplexoare;
- folosind două rețele de semnale: una de intrări, alta de ieșiri, conectate prin circuite *three-state*;
- folosind un buffer de memorie în care intrările corespund cu procesoarele; fiecare element din buffer indică modulul de memorie RAM care poate fi accesat de procesorul respectiv.

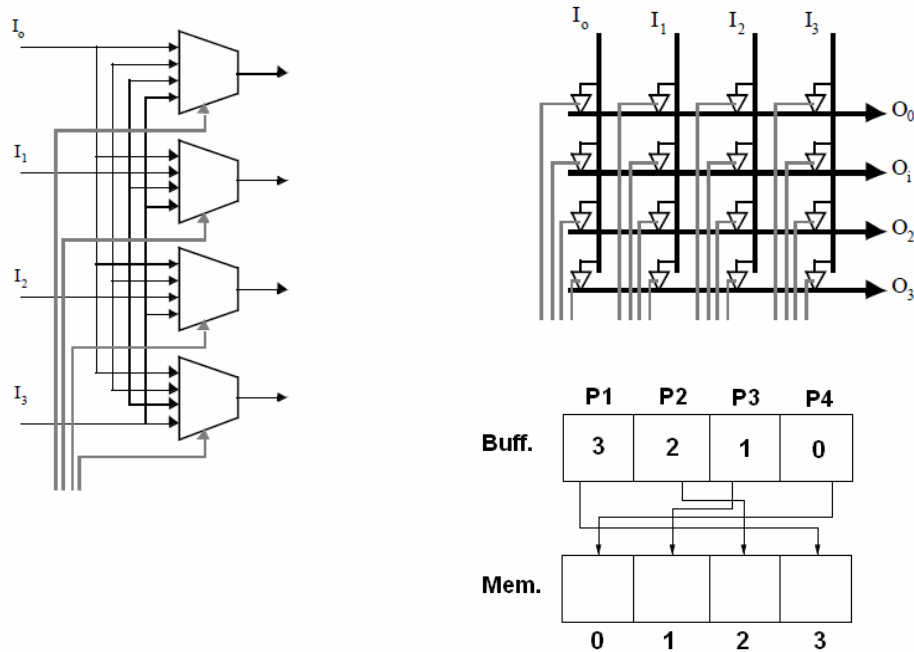


Figura 6. Implementări hardware pentru o rețea de tip *crossbar*

Comparație între rețelele directe și cele indirecte:

- La rețelele directe există canale de comunicație dedicate între procesoarele vecine. La rețelele indirecte nu există legături punct la punct între procesoare.
- La rețelele directe este necesar un algoritm de rutare pentru ca informațiile să ajungă de la sursă la destinație. La rețelele indirecte doar rețelele multistagiu necesită rutare. Aici *switch*-urile comutate într-un anumit mod vor asigura legătura între sursă și destinație. Magistralele nu necesită rutare iar rețelele *crossbar* necesită comutarea unui singur *switch* pentru efectuarea conexiunii.
- Rețelele directe pot fi ușor scalabile, adică se pot adăuga noduri în rețea cu modificări minime. Rețelele indirecte sunt mai greu scalabile, astfel magistralele pot ajunge la saturație când se depășește un anumit număr de procesoare. La rețelele multistagiu și *crossbar* creșterea numărului de procesoare devine scump datorită conexiunilor și *switch*-urilor care trebuie adăugate în plus.
- Rețelele directe folosesc de obicei comutarea prin pachete; rețelele *crossbar* și magistralele folosesc comutarea prin circuite. La rețelele multistagiu, dacă

switch-urile sunt setate înainte începerii comunicărilor, atunci va fi realizat comutarea prin circuite. Dacă *switch*-urile pot comuta dinamic, în funcție de ruta unui mesaj, atunci se poate implementa comutarea prin pachete.

Comutarea prin pachete (*packet switching*)

Rețelele care utilizează comutarea prin pachete folosesc *switch*-uri dedicate pentru transmiterea pachetelor. Aceste *switch*-uri transmit pachetele în rețea până când ele ajung la nodul destinație.

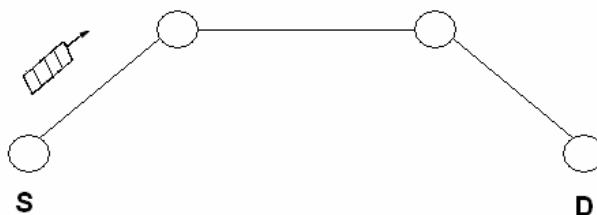


Figura 7. Transmiterea pachetelor de la surs la destinație

La rețele directe, *switch*-urile fac parte din nodurile rețelei. Aici, *switch*-urile transmit pachetul mai departe dacă nodul curent este diferit de nodul destinație, sau copiază pachetul în memoria locală dacă nodul curent este același cu nodul destinație.

La rețele indirecte, *switch*-urile nu fac parte din nodurile de procesare.

Există 3 modalități de comutare a pachetelor:

- 1) *Store & Forward*
- 2) *Wormhole routing*
- 3) *Virtual cut-through*

1) *Store & Forward*

Aici, pachetul este cea mai mică unitate de transfer a datelor. *Switch*-urile din rețea dispun de buffere pentru memorarea locală a pachetelor.

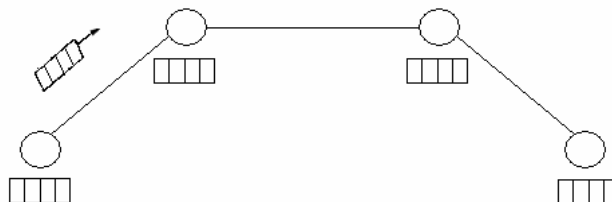


Figura 8. Transmisia *Store & Forward*

În drumul de la surs la destinație, un pachet ce tranzitează un nod va fi mai întâi memorat în întregime în nodul respectiv (etapa ***Store***), apoi va fi transmis mai departe către nodul destinație (etapa ***Forward***). Tehnica este simplă de implementat, dar toleranța la erori nu e foarte bună.

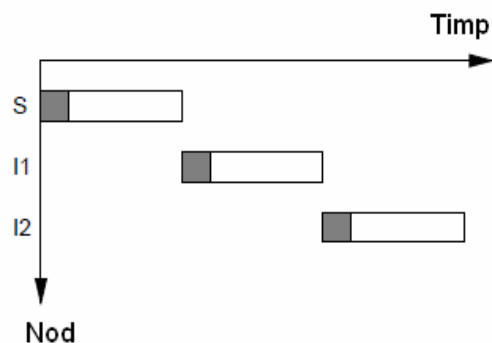


Figura 9. Diagrama de timp a transmisiei unui pachet - *Store & Forward*

2) Wormhole routing

În cadrul acestei metode, pachetele sunt împărțite în blocuri de dimensiuni mici (de la 1 la câțiva octeți). Blocul de date minimal care este transmis pe liniile de comunicație se numește **flit** (*Flow Control Digit*). Datele sunt conduse în rețea în manieră *pipeline* flit cu flit, toate flit-urile urmând același traseu. Când un nod recepționează un flit, îl transmite imediat către nodul următor și apoi va recepționa următorul flit.

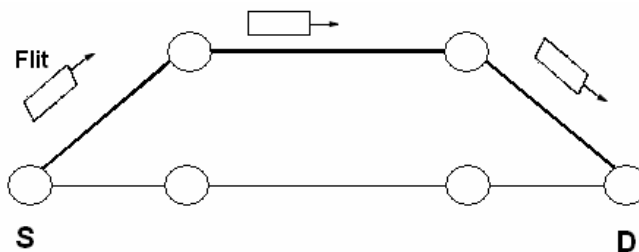


Figura 10. Transmisia *Wormhole routing*

Informația de rutare este stocată în primele flit-uri din pachet. Pe toată durata transmiterii unui pachet, canalele de comunicație alocate la început rămân ocupate până la transmiterea ultimului flit din pachet, care le eliberează.

Avantajul metodei: timpul de tranziție al datelor în noduri scade, fliturile având dimensiuni mai mici decât dimensiunea unui pachet.

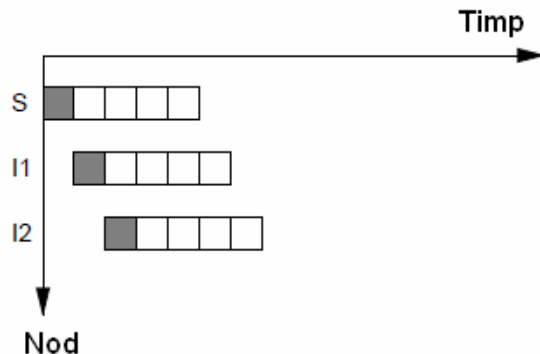


Figura 11. Diagrama de timp a transmisiei unui pachet - *Wormhole routing*

Dezavantaj: metoda poate duce la saturarea rezelei. Astfel, atunci când *header*-ul unui pachet este blocat, rămân blocate toate canalele de comunicație alocate pentru transmiterea pachetului, fiecare flit așteptând într-un nod intermediar. Aceasta situație poate duce de asemenea și la blocarea altor pachete care circulă în rezele.

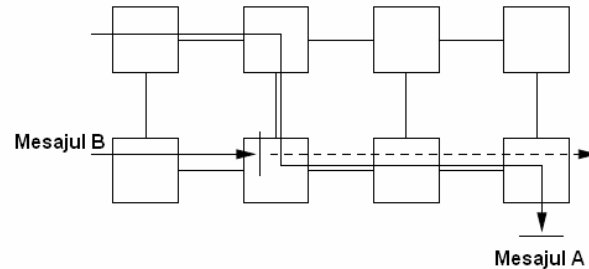


Figura 12. Exemplu de blocare a mesajelor la metoda *Wormhole routing*

Timpii de transmisie a pachetelor pe liniile de comunicație se calculează astfel:

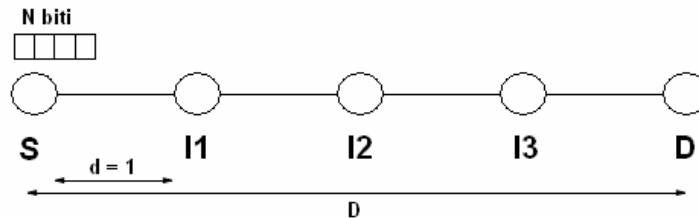


Figura 13. Transmiterea unui pachet – calculul timpilor de transmisie

- *Store and forward:* - N = numărul de biți dintr-un pachet
- W = lățimea de bandă a canalului

$$T1 = \frac{N}{W} \times D$$

- *Wormhole routing:* - F = numărul de biți dintr-un flit

$$T2 = \frac{N}{W} + \frac{F}{W} \times (D - 1)$$

unde:

$\frac{N}{W}$ - timpul în care ultimul flit din pachet ajunge la nodul vecin I1

$\frac{F}{W} \times (D - 1)$ - timpul după care ultimul flit (de dimensiune F) ajunge de la nodul I1 la destinație

3) *Virtual cut-through*

La această metodă, bufferele din fiecare nod au mărimea egală cu lungimea unui pachet. Transmiterea pachetului se face *pipeline* flit cu flit ca și la **Wormhole Routing**. Dacă se blochează un *header* al unui pachet, transmiterea celorlalte flit-uri continuă până când toate fliturile ajung în nodul care s-a blocat. Astfel canalele alocate se pot elibera, reducând gradul de saturare al rețelei.

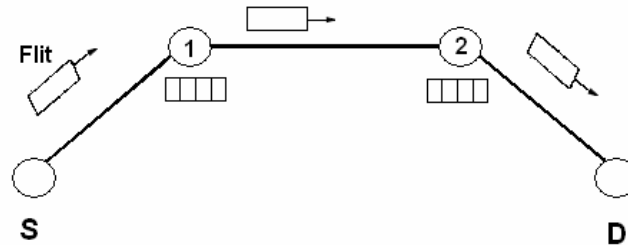


Figura 14. Metoda *Virtual cut-through*

Pentru controlul fluxului de date (asigurarea faptului că datele nu se pierd în rețea), se poate folosi un protocol de tip *handshaking*, care este bazat pe o pereche de semnale **Request/Acknowledge**:

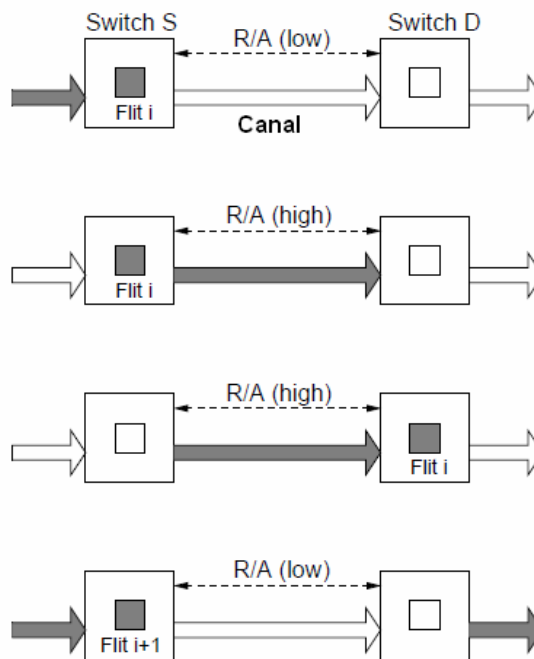


Figura 15. Protocol de transmisie bazat pe semnale **Request/Acknowledge**

TEHNICI DE RUTARE

Rutarea determină calea sau ruta pe care datele o urmează de-a lungul rețelei pentru a ajunge la destinație. Există 2 metode de rutare:

- Rutare stabilită de la nivelul nodului surs (*Source-based Routing*);
- Rutare locală (*Local Routing*), în care determinarea rutei se face în cadrul fiecărei switch din rețea.

a) Rutare stabilită la nivelul nodului surs

În acest caz, nodul surs stabilește calea de urmat prin rețea. Calea se scrie în header-ul pachetului și stabilește direcția pe care pachetul trebuie să o urmeze pentru fiecare nod.

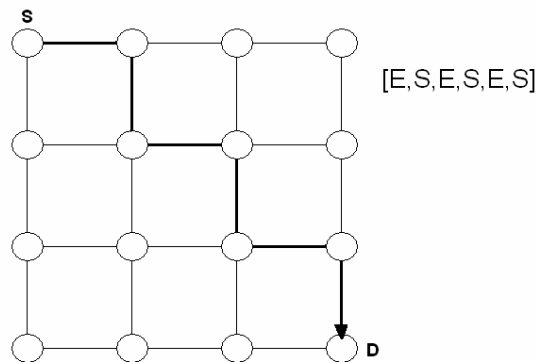


Figura 1. Rutare stabilită la nivelul nodului surs

Switch-ul din fiecare nod al rețelei va trimite pachetul prin rețea în conformitate cu informația de direcție curentă (care apoi este eliminată din header).

Dezavantaje:

- Deoarece pachetul include calea completă de parcurs prin rețea, aceasta va duce la creșterea dimensiunii fiecărui pachet.
- Pachetele nu pot anticipa defecțiunile switch-urilor din rețea și nu pot alege o altă rută dacă un canal de comunicație este blocat.

b) Rutarea locală

Decizia rutei se face la nivelul fiecărei switch din rețea. Ca urmare, nodul surs va specifica în headerul pachetului doar identificatorul nodului destinație, reducând astfel dimensiunea pachetului, dar pentru fiecare switch trebuie să decidă asupra celui de urmat, aceasta va duce la creșterea complexității constructive a fiecărei switch.

Rutarea poate duce la blocaje (*deadlocks*) datorită timpurilor pentru eliberarea canalelor de comunicație.

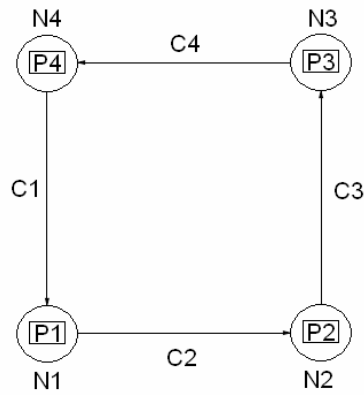


Figura 2. Blocaj datorat a tept rilor
pentru eliberarea canalelor de comunica ie

În exemplul din figura 2, pachetul P1 a teapt eliberarea canalului C2, dar în acela i timp ine ocupat canalul C1 pentru a nu fi transmis alt pachet care s îl suprascrie în bufferul din nodul N1. La fel, P2 a teapt eliberarea lui C3 i ine ocupat canalul C2. La fel P3 asteapt eliberarea lui C4 i ine ocupat C3. P4 a teapt eliberarea lui C1 i ine ocupat canalul C4. Astfel apare o dependen ciclic între canale (fig. 3):

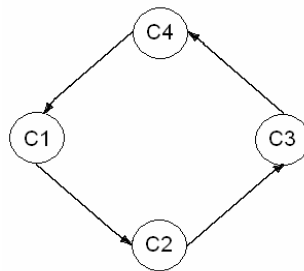


Figura 3. Graful dependen ei dintre canale

Rutarea local poate fi minimal sau non-minimal . La rutarea minimal , algoritmul de rutare alege întotdeauna calea cea mai scurt între surs i destina ie. La rutarea non-minimal se poate alege orice cale.

Exist dou tehnici de rutare local :

- 1) rutare determinist ;
- 2) rutare adaptiv .

1) Rutarea determinist

În acest caz, pachetele trimise de la un nod surs la un nod destina ie urmeaz întotdeauna aceea i rut. Ruta este aleas astfel încât s fie minimal i s nu duc la blocaje.

Dou exemple pentru rutarea determinist sunt: rutarea în coordonate XY i etichetarea intervalelor:

R1) Rutarea în coordonate XY, care se folosește pentru rețele de tip *mesh* (gril).

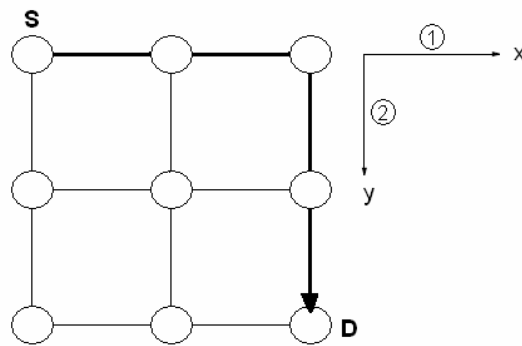


Figura 4. Rutarea în coordonate XY

Un pachet care merge de la sursă la destinație este condus mai întâi pe axa X pentru a ajunge la coordonata X a destinației, apoi pe axa Y pentru a ajunge la coordonata Y a destinației. Aici nu există blocaje datorate dependențelor ciclice. De exemplu, un pachet care circulă pe coordonata Y nu va mai fi trimis pe coordonata X (fig.5):

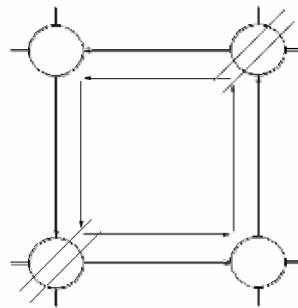


Figura 5. Eliminarea dependențelor ciclice la rutarea în coordonate XY

R2) Etichetarea intervalelor. Aici, porturile de ieșire ale fiecărei *switch* au asociat un anumit interval care indică nodurile destinație posibile pentru acel port:

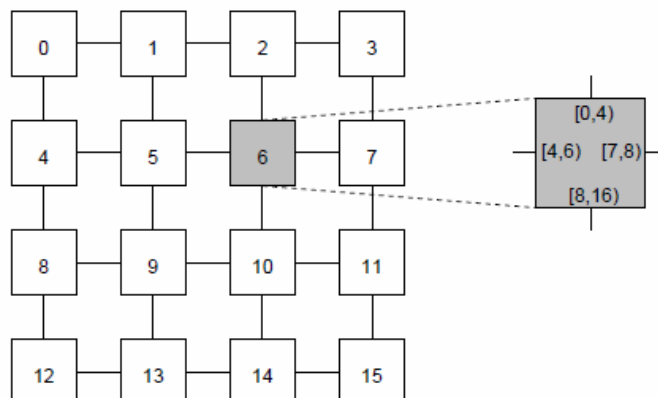


Figura 6. Etichetarea intervalelor

Un pachet care sosește într-un anumit nod va fi trimis pe portul de ieșire care specifică intervalul ce conține identificatorul nodului destinației. În acest caz se pot alege intervalele fiecărei *switch* astfel încât să nu apară blocaje în rutare.

2) Rutarea adaptivă

În acest caz, ruta se stabilește în mod dinamic la fiecare *switch* în funcție de integritatea rețelei și de canalele alocate deja pentru alte transmisii.

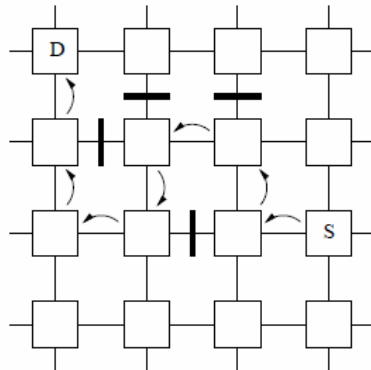


Figura 7. Rutarea adaptivă

Avantaje:

- rutarea adaptivă este mai tolerantă la erori;
- poate să evite blocajele din rețea;
- încărcarea rețelei este mai uniformă decât la rutarea deterministă :

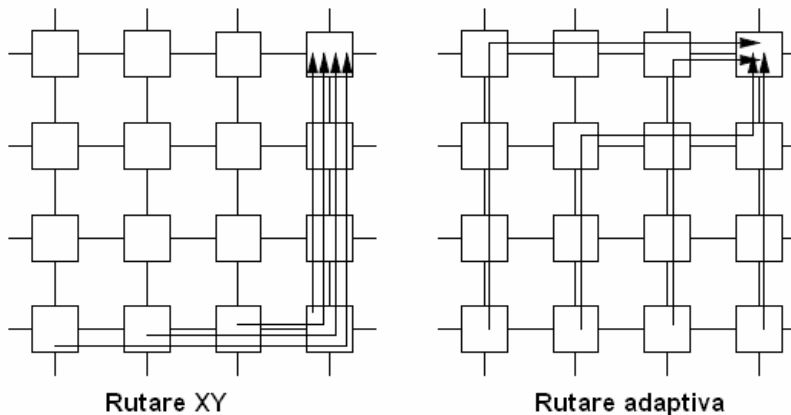


Figura 8. Uniformizarea încărcării rețelei la rutarea adaptivă față de rutarea XY

Dezavantaje:

- rutarea adaptivă necesită *switch*-uri mai inteligente, deci mai scumpe;
- algoritmul de rutare devine mai complicat, datorită necesității de a evita blocajele din rețea;

- pachetele care apar în același mesaj pot să ajungă la destinație într-o ordine aleatoare, datorită faptului că ele pot fi rutate pe căi diferite. Astfel, reasamblarea pachetelor la destinație va mări durata transmisiei.

Rutarea adaptivă poate fi de 2 feluri:

- 1) Complet adaptiv, în care nu există constrângeri asupra direcțiilor pe care le poate urma un pachet.
- 2) Parțial adaptiv, în care pot exista astfel de constrângeri.

Eliminarea blocajelor de rutare:

Există mai multe metode pentru eliminarea blocajelor datorate dependențelor ciclice. O primă metodă este folosirea rutării deterministe, de exemplu rutarea XY.

O altă metodă este folosirea unei metode parțial adaptivă care derivă din rutarea XY, numită **West First Routing**, care se folosește pentru rețele de tip *mesh*. Aici se transmite prima dată pachetul către vest (dacă e nevoie) apoi se transmite adaptiv la nord, sud sau est.

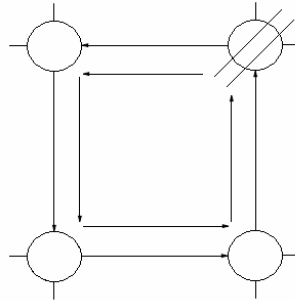


Figura 9. Eliminarea dependențelor ciclice la rutarea **West First Routing**

Rutarea **West First Routing** este non-minimală pentru a fi capabilă să elimine blocajele.

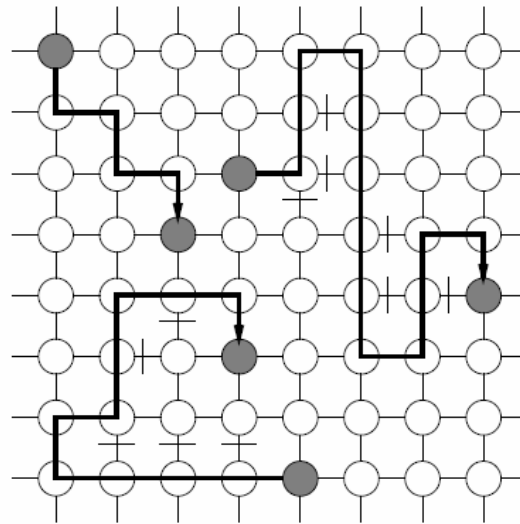


Figura 10. Exemple de rutare **West First Routing**

O alt metod utilizat pentru eliminarea dependen elor ciclice este folosirea canalelor de comunica ie virtuale. Un canal virtual este o leg tur logic între 2 noduri, mai multe canale virtuale putând fi multiplexate pe un acela i canal de comunica ii.

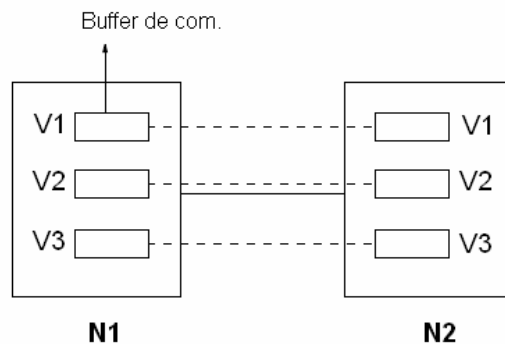


Figura 11. Folosirea canalelor de comunica ie virtuale

Multiplexarea pe canalul de comunica ie se face prin utilizarea unui mecanism rotativ prin care se selecteaz pe rând canalele virtuale pentru transmisie.

Introducerea de canale virtuale poate s evite dependen ele ciclice (fig. 12):

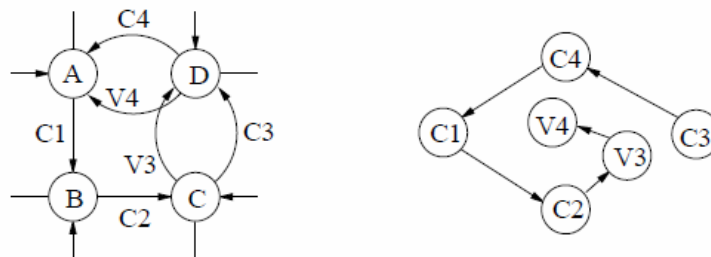


Figura 12. Eliminarea dependen elor ciclice prin introducerea de canale virtuale

Un alt exemplu de folosire a canalelor virtuale se refer la o re ea de tip *mesh* cu 2 leg turi virtuale bidirec ionale pe axa Y i o singur leg tur bidirec ional pe axa X. Pe această re ea se construie te subre ea +X” pentru transmiterea unui pachet la est de surs , i subre ea -X” pentru a transmite un pachet la vest de surs .

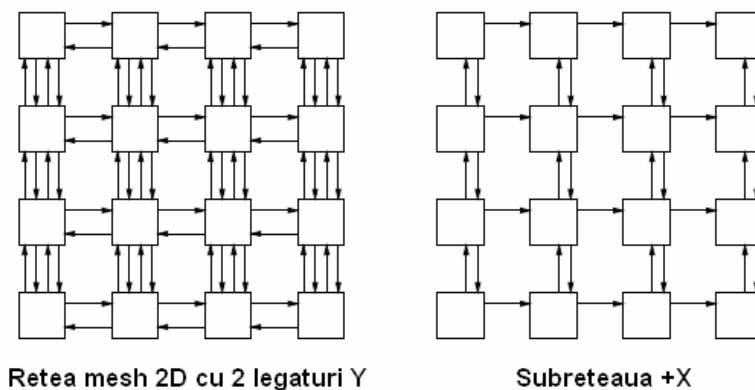


Figura 13. Re ea *mesh* cu 2 leg turi virtuale pe axa Y i o leg tur virtual pe axa X

Avantajele folosirii canalelor virtuale:

- elimină dependențele ciclice;
- se mărește capacitatea de utilizare a rețelei; astfel, dacă un pachet se blochează pe un canal virtual așteptând eliberarea unui alt canal de comunicație, alte pachete mai pot fi totuși transmise de-a lungul legăturii fizice respective folosind un alt canal virtual de comunicație:

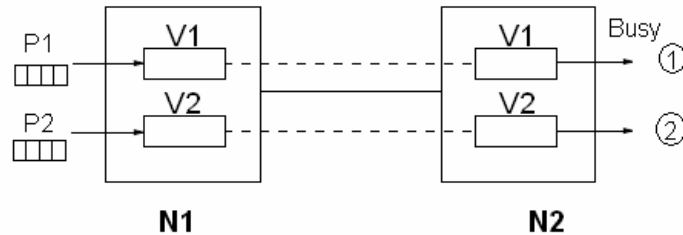


Figura 14. Alegerea unui canal virtual liber pe un canal fizic având două legături virtuale

- canalele virtuale facilitează maparea unei topologii logice peste topologia fizică a rețelei.

Dezavantaje:

- necesită resurse hardware mai complexe, de exemplu: buffere mai mari în noduri, implementarea fizică a unei politici de planificare a canalelor virtuale pe legătura fizică etc. Toate acestea vor duce la creșterea costurilor de implementare;
- canalele virtuale pot crea latență a transmisiei pentru că pachetele folosesc partajat limea de bandă a canalelor de comunicație.

Metode pentru implementarea comunicației multicast

Într-o rețea de comunicație există următoarele tipuri de transmisii:

- transmisie **unicast** - reprezintă transmiterea punct la punct între un nod surs și un nod destinație;
- transmisie **multicast** - de la un nod surs la mai multe noduri destinație;
- transmisie **broadcast** - de la un nod surs la toate celelalte noduri.

Implementarea comunicației multicast se poate face în 3 moduri:

- 1) Fără suport hardware special, caz în care se utilizează un software special pentru rutare. *Switch*-urile nu trebuie să fie inteligente, iar pentru transmiterea către nodurile destinație se folosește comunicația unicast.

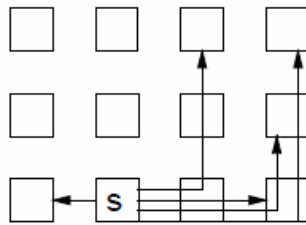


Figura 15. Implementare multicast folosind comunicația unicast

- 2) Se folosesc *switch-uri* care permit replicarea datelor. Astfel *switch*-urile copiază informația din pachete în bufferele proprii și apoi o trimit mai departe la nodurile vecine în funcție de destinație.

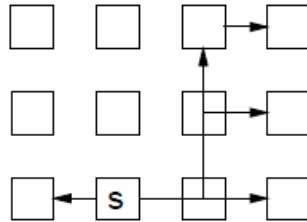


Figura 16. Implementare multicast prin replicarea datelor

- 3) Se pot folosi *switch-uri inteligente* care permit alegerea unei căi optime prin care ea, astfel încât pachetul să ajungă la toate nodurile destinație.

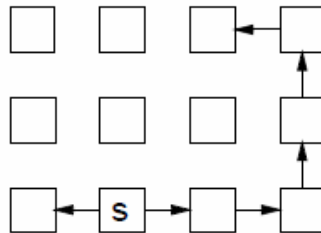


Figura 17. Implementare multicast prin folosirea *switch*-urilor inteligente

MULTICALCULATOARE

Multicalculatoarele sunt arhitecturi MIMD cu memorie distribuită. Ele folosesc metoda de comutare prin pachete pentru schimbul de date și o rețea directă (punct la punct) pentru conectarea procesoarelor între ele.

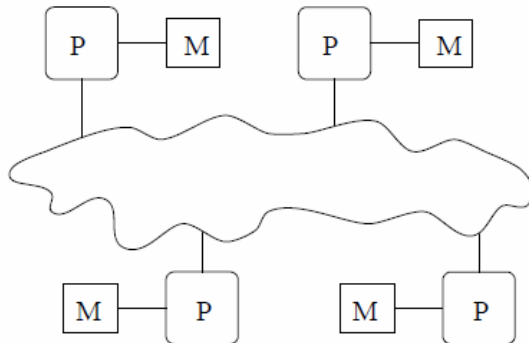


Figura 1. Structura unui multicalculator

Într-un multicalculator nu există un spațiu global de adrese, un procesor putând accesa doar memoria privată (locală). Comunicarea și sincronizarea între procesoare se realizează prin intermediul transmisiei de mesaje; utilizatorii sunt cei care trebuie să gestioneze comunicația prin codul aplicației.

La modul global, multicalculatoarele sunt arhitecturi fără partajare în care toate resursele sistem (memorie, discuri, etc.) sunt distribuite și accesibile doar de către procesorul local. În ceea ce privește rutarea, nodurile intermediare din rețea trebuie să găsească o rută pentru mesajele transmise la un nod aflat la distanță. La început, această rută se stabilea prin software de către procesorul nodului respectiv, printr-un mecanism bazat pe întreruperi. În sistemele de astăzi, există un circuit integrat specializat în fiecare nod de comunicație, care se ocupă de rutarea, comutarea și transferul pachetelor la nodul destinație.

Dezavantajul multicalculatoarelor derivă din nonexistența memoriei globale. Cu toate acestea, se poate simula funcționarea unei memorii globale prin 2 metode:

- 1) VSM – *Virtual Shared Memory*
- 2) SVM – *Shared Virtual Memory*

În ambele abordări, programatorul vede o arhitectură cu o memorie globală partajată, accesibilă de către toate procesoarele. Referințele la această memorie globală sunt rezolvate prin intermediul transmiterii unor mesaje care realizează schimbul de date între procesor și memoria globală.

1) VSM – Virtual Shared Memory

VSM presupune o implementare hardware; memoria partajat global va face parte din memoria virtual a sistemului, fiind dispus în vârful acesteia. SO nu va face distincie între un multicalculator de tip VSM și o mașină cu memorie partajat.

În această metodă, unitatea folosită pentru partajarea memoriei este blocul cache. Astfel, la accesarea unui bloc de memorie, în cazul unui *cache miss*, controlerul de memorie cache determină dacă acel bloc se găsește în memoria locală sau nu. Dacă nu, controlerul va trimite un mesaj la nodul *remote* (de la distanță) care deține blocul de memorie cerut, pentru a-l aduce prin rețea în nodul local.

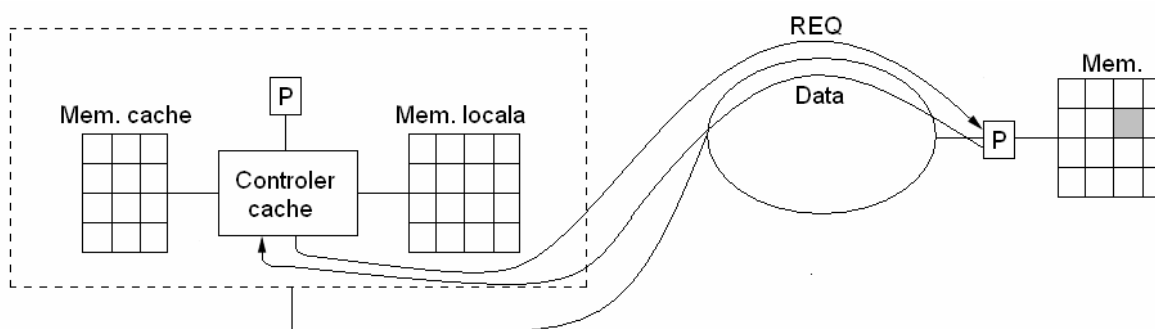


Figura 2. Implementarea VSM

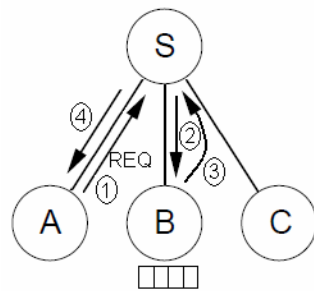
2) SVM - Shared Virtual Memory

SVM este o implementare software la nivelul SO, având ca suport hardware unitatea de management a memoriei (MMU) a procesorului. Aici, unitatea de partajare este pagina de memorie a SO. La adresarea unei locații de memorie, MMU determină dacă pagina ce conține acea locație se află în memoria locală sau nu. Într-un calculator obișnuit, dacă pagina de memorie nu se află în memoria locală, atunci va fi încărcată din memoria virtuală stocată pe hard-disk. În cazul de față, SO va încălca acea pagină din nodul *remote* care o deține.

Există 4 abordări pentru sistemele VSM, respectiv SVM:

- a) **Metoda utilizând un server central:** există un singur nod care poate accesa un anumit bloc de memorie la un moment dat; toate cererile de acces la memorie sunt trimise către un server central. Serverul este cel care cunoaște locația exactă în rețea a oricărui bloc de memorie și răspunde la cererile venite din partea nodurilor din sistem.

În figura 3 este ilustrată o situație în care nodul A formulează o cerere pentru un bloc de memorie aflat pe nodul B:



- 1) A trimite o cerere c tre S
- 2) S trimite cererea la B
- 3) B r spunde si trimite blocul de memorie c tre S
- 4) S trimite r spuns la A

Figura 3. Metoda utilizând un server central

Avantaje: algoritmul este u or de implementat; se men ine o puternic coeren si consisten a sistemului.

Dezavantaj: se pot produce gâtuiiri în traficul de date al re elei, datorit supraînc rc rii cu cereri a serverului central. Pentru rezolvarea acestei probleme, datele partajate se pot distribui între mai multe servere.

- b) **Metoda migra iei (*Full Migration*)**: nu exist un server central; fiecare nod tie unde anume se afl fiecare bloc de memorie, pe baza unei distribu ii statice a adreselor în re ea. La un anumit moment de timp, fiecare bloc de memorie se g se te într-un anumit nod al re elei sub controlul unui anumit procesor. Dacă se face o cerere prin re ea pentru o dat con inut într-un bloc de memorie, procesorul care de ine acel bloc îl transmite mai întâi c tre nodul care a f cut cererea i apoi îi invalideaz copia blocului din nodul local.

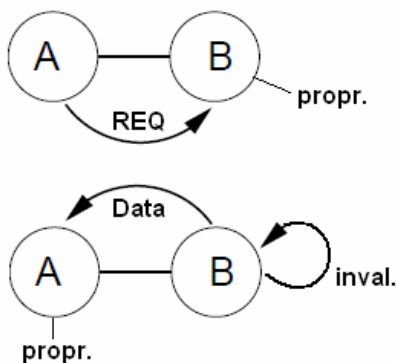


Figura 4. Metoda migra iei

În întreaga re ea, în orice moment de timp, va exista o singur copie valid pentru fiecare bloc de memorie. Astfel, fiecare bloc de memorie va avea la un moment dat un singur proprietar, adic nodul care de ine copia valid a blocului.

Dacă nodul c ruia îi este adresat o cerere pentru un bloc de memorie a trimis anterior acel bloc c tre alt procesor (nod), atunci el va trimite cererea mai departe c tre acel nod care a devenit noul proprietar (fig. 5):

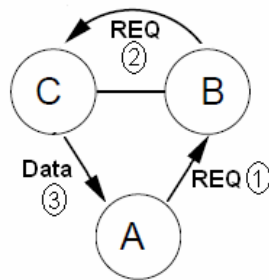
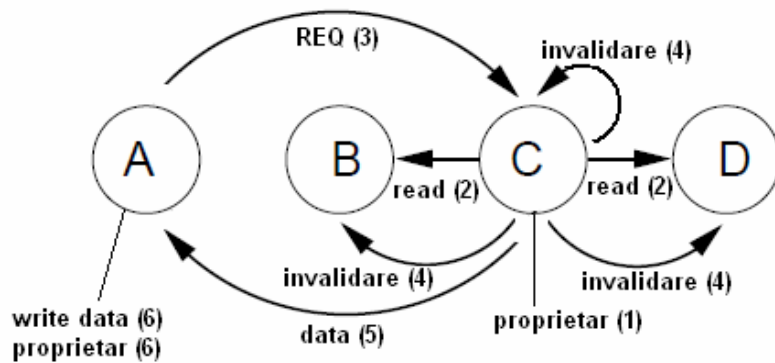


Figura 5. Migrarea cererilor în rețea

Ca urmare, apare dezavantajul că o cerere poate migra mai mult timp în rețea până să ajungă la destinatar.

- c) **Metoda Read Replication:** această metodă încearcă să reducă latența de comunicație permițând mai multor noduri să aibă copii *read-only* ale blocurilor de memorie; aici fiecare nod al rețelei ține pe baza unei distribuții statice a adreselor unde se află un anumit bloc de memorie. Se face distincție între citirea și scrierea unei date:

- 1) Când un procesor efectuează o citire non-locală, el preia prin copiere blocul de la nodul proprietar.
- 2) Dacă se efectuează o scriere, atunci putem avea de exemplu următoarea situație:



- 1) C este proprietarul datelor
- 2) B,D citesc aceleași date
- 3) A face o cerere pentru scriere
- 4) C invalidează datele trimise anterior
- 5) C trimite date la A
- 6) A scrie datele și devine noul proprietar

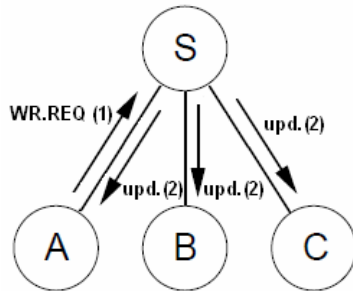
Figura 6. Metoda *Read Replication*

Diferența față de metoda anterioară (pct. b) este că citirea datelor nu schimbă proprietarul, deci nu mărește latența de comunicație. Proprietarul se schimbă doar la scrierea datelor.

- d) **Metoda Full Replication:** se permite mai multor noduri să aibă copii ale blocurilor de memorie atât la citire cât și la scriere. Toate nodurile cunosc

distribu ia blocurilor de memorie în re ea; citirea unui bloc se face de la orice nod care îl de ine; scrierea unui bloc se face astfel:

- 1) Se preia de la un proces secven iator un num r de secven unic în re ea;
- 2) Valoarea care se dore te a fi scris în blocul de date, împreun cu num rul de secven , se trimite la toate nodurile care de in o copie a blocului respectiv (fig. 7):



- 1) A trimite o cerere de scriere c tre S
- 2) S adaug un num r de secven la aceast cerere i trimite un mesaj de actualizare c tre nodurile A,B i C

Figura 7. Metoda *Full Replication*

Num rul de secven este necesar pentru cazul în care în re ea circul mai multe mesaje de actualizare pentru aceea i dat .

Câteva exemple de multicalculatoare: - IBM SP2
- Parsytec CC

IBM SP2

IBM SP2 folose te procesoare de tip Power 2, utilizând pentru conectarea procesoarelor pl ci de conectare ce con in fiecare o re ea multistagiu cu 2 nivele. Re eaua este alc tuit din 8 *switch*-uri, fiecare *switch* având 8 leg turi bidirec ionale:

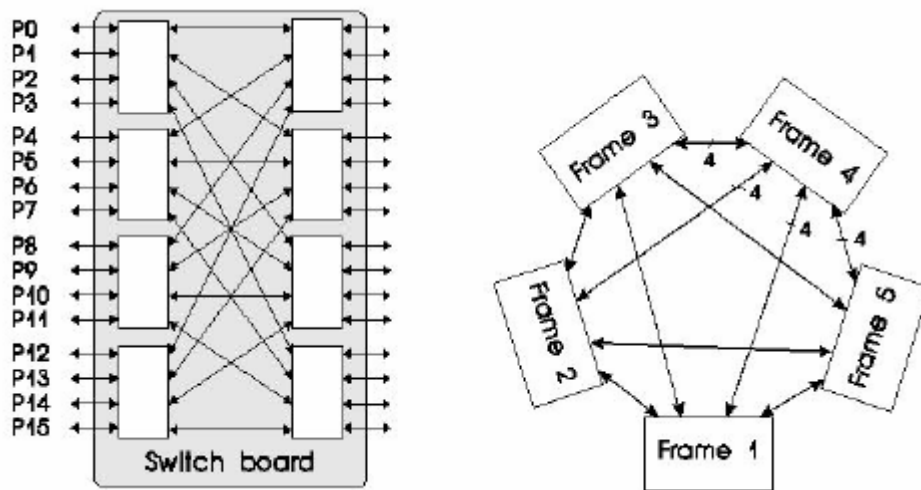


Figura 8. Structura de conectare la calculatorul IBM SP2

În figura următoare este ilustrat un exemplu de conectare cu 12 *switch*-uri:

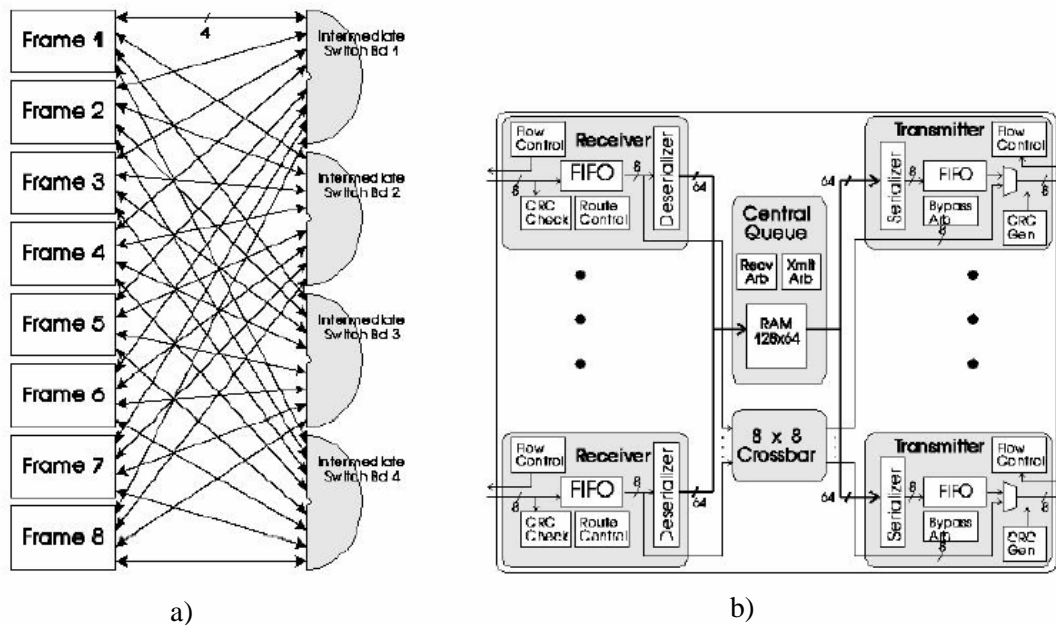


Figura 9. a) Exemplu de conectare cu 12 *switch*-uri; b) Arhitectura unui *switch*

Arhitectura unui *switch* cuprinde 8 receptoare (*receiver*-e) și 8 transmi-toare (*transmitter*-e), fiecare *receiver* este conectat la fiecare *transmitter* printr-o re-ea *crossbar* (ce conectează orice intrare cu orice ie-ire).

Se implementează rutarea de tip *wormhole routing*, bazat pe flituri. Când un flit (o parte dintr-un pachet) ajunge la intrarea unui *receiver*, sunt posibile două situa-ții:

- Dacă *transmitter*-ul dorit este liber, flitul este transmis direct la transmi-tor prin re-ea *crossbar*;
- Dacă *transmitter*-ul este ocupat, atunci *switch*-ul poate stoca până la 128x8 flituri în bufferul central; apoi, dacă *transmitter*-ul rămâne ocupat, se blochează linia de intrare a *receiver*-ului.

Parsytec CC

Acest multicalculator realizează un echilibru între eficiență și costuri prin combinarea unei re-ele multistagiu cu o re-ea de tip *mesh* (grilă). Sistemul constă din conectarea unor procesoare de tip Power PC prin intermediul unei re-ele numite ***Mesh of Clos*** bazat pe re-ea multistagiu ***Clos***.

Re-ea ***Clos*** se bazează pe *switch*-uri având un număr de $2k$ canale bidirec-ționale.

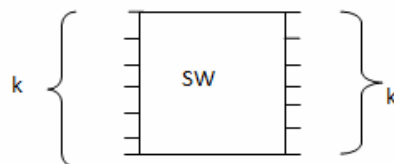


Figura 10. *Switch* având un număr de $2k$ canale bidirec-ționale

O rețea **Clos** de dimensiune 1 este alcătuită dintr-un singur *switch* la care se conectează k procesoare:

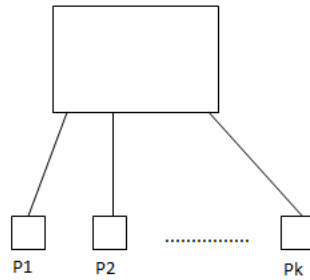


Figura 11. Rețeaua **Clos** de dimensiune 1

Rețeaua **Clos** de dimensiune h se realizează prin conectarea a k rețele de dimensiune $h-1$ printr-un număr de k^{h-1} *switch-uri*:

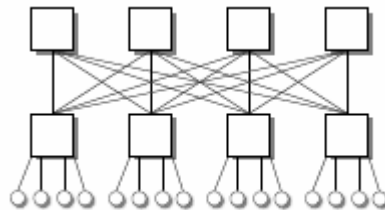


Figura 12. Rețeaua **Clos** de dimensiune $h=2$;
 $k=4$, fiecare *switch* e conectat la 4 procesoare

Rețeaua **Mesh of Clos** se realizează prin înlocuirea unui număr r de stadii din rețeaua **Clos** printr-o rețea de tip *mesh*. Pentru $r=1$ se înlocuiesc *switch-urile* din rețea cu rețele *mesh*.

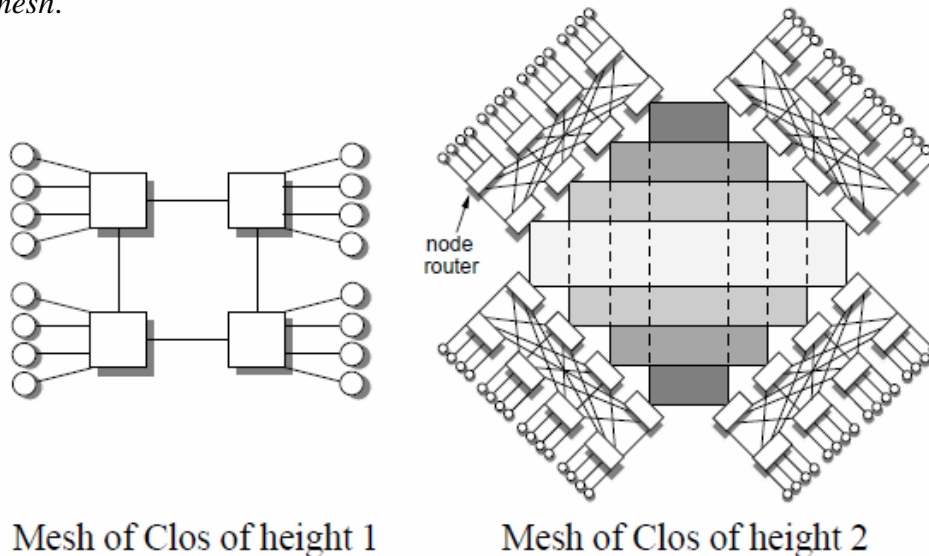


Figura 13. Rețele **Mesh of Clos** de dimensiuni $h=1$ și $h=2$

S-a demonstrat că acest tip de conectare este mai bun decât o conectare simplă de tip *mesh* (**Mesh of Clos** are un număr de *switch-uri* mai mic decât o rețea de tip *mesh*).

MULTIPROCESOARE

Multiprocesoarele sunt calculatoare paralele care dispun de o singură memorie globală partajată. Procesoarele din aceste sisteme sunt conectate printr-o rețea indirectă sau printr-o combinație de rețea indirectă și rețea punct la punct.

Comunicarea și sincronizarea dintre procesoare este realizată prin intermediul unor locații partajate din memoria globală, dar pe acestea mai inițial se poate emula și interfața pe bază de transmisie de mesaje MPI (*Message Passing Interface*) caracteristic multicalculatoarelor.

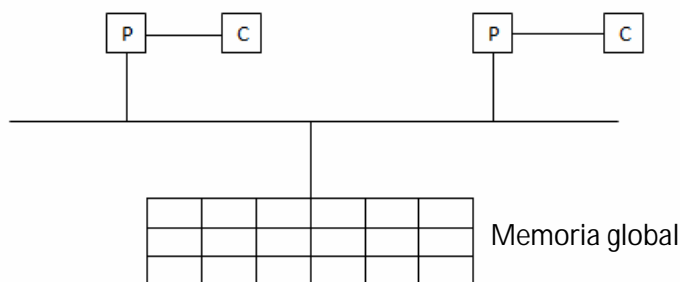


Figura 1. Arhitectura unui sistem multiprocesor

O problemă a acestor calculatoare este scalabilitatea redusă. De exemplu, calculatoarele multiprocesor bazate pe magistrală nu se pot scala la un număr mare de procesoare. Astfel, la un număr de 8-10 procesoare, magistrala intră în saturație. O altă problemă care trebuie rezolvată la calculatoarele multiprocesor este menținerea coerenței memoriilor cache.

Coerența memoriilor cache în calculatoarele multiprocesor cu memorie partajată

Cum fiecare procesor dispune de o memorie cache și fiecare procesor poate accesa o aceeași locație din memoria partajată, pot apărea inconsistențe între diferitele memorii cache care au încărcat aceeași variabilă.

Sursele inconsistențelor pot fi 3 la număr:

- 1) Scrierea datelor în zone partajate de memorie.
- 2) Migrarea proceselor între procesoare.
- 3) Activitățile de intrare-ieșire.

Exemple:

Situația 1:

Scrierea datelor într-o zonă partajată :

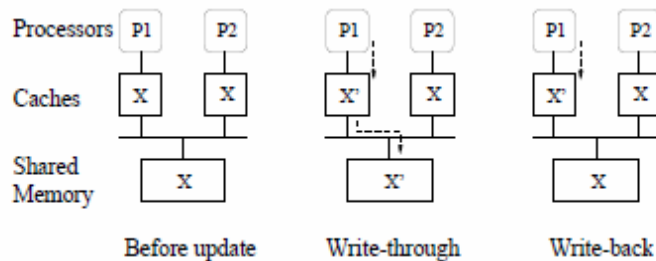


Figura 2. Inconsisten la scrierea datelor într-o zonă partajată

Situația inițială : P1 și P2 citesc aceeași memorie X.

Write-through: Memoria cache permite scrierea, valoarea este scrisă imediat în memoria globală. P1 scrie valoarea X' în variabila X după care P2 vrea să citească această valoare, dar P2 va citi X, nu valoarea actualizată (X') => inconsistență.

Write-back: P1 scrie valoarea X' în X dar nu este trecut în memoria globală. P2 citește valoarea X => inconsistență.

Situația 2:

Migrarea proceselor.

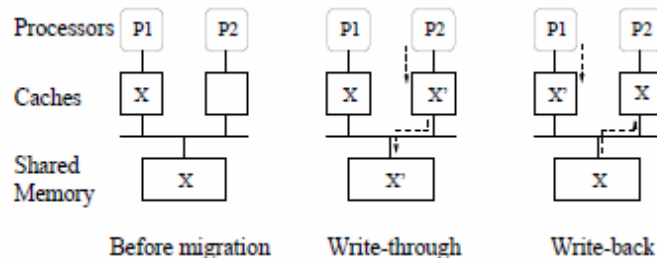


Figura 3. Inconsistență în cazul migrației proceselor

Situația inițială : P1 citește valoarea X din memoria globală.

Write-through: Un proces din P2 scrie valoarea X' în variabila X după care procesul se mută în P1 și va citi din nou variabila X => inconsistență.

Write-back: Un proces din P1 scrie valoarea X' în variabila X apoi același proces migrează în procesorul P2 care citește valoarea X din memoria globală => inconsistență.

Situația 3:

Activitățile de intrare-ieșire.

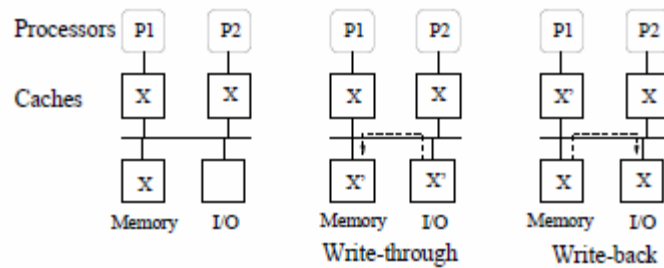


Figura 4. Inconsisten a datelor în cazul activităților de intrare-ieșire

Situația inițială: S-a adăugat un dispozitiv de intrare-iesire pe magistrala sistem, P1 și P2 citesc aceeași valoare X din memoria globală.

Write-through: Presupunem că perifericul vrea să scrie valoarea X' în memoria principală în locul variabilei X. Scrierea se face fără intervenția procesoarelor, de exemplu prin DMA. Apoi P1 sau P2 vor să citească variabila X din memorie => inconsistență.

Write-back: P1 scrie X' în variabila X după care perifericul vrea să citească valoarea X din memorie => inconsistență.

Pentru rezolvarea ultimei probleme (la memoriile *write-back*) se poate muta interfața de intrare-iesire lângă procesor astfel încât perifericul să poată accesa memoria globală doar prin memoria cache:

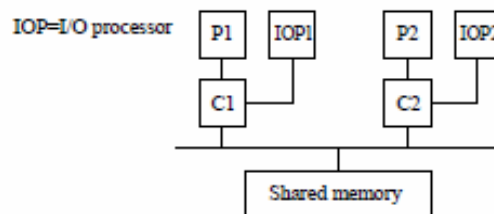


Figura 5. Rezolvarea inconsistenței prin reamplasarea interfețelor de intrare-iesire

O altă soluție tot pentru această problemă este ca procesorul să nu citească locațiile de memorie din spațiul de intrare-iesire (alocat perifericelor) din memoria cache, ci direct din memoria globală.

Pentru menținerea coerenței memoriilor cache se folosesc anumite protocoale de coerență. Există **două tipuri de protocoale pentru menținerea coerenței**:

- 1) *Write invalidate* (scriere cu invalidare)
- 2) *Write update* (actualizare la scriere)

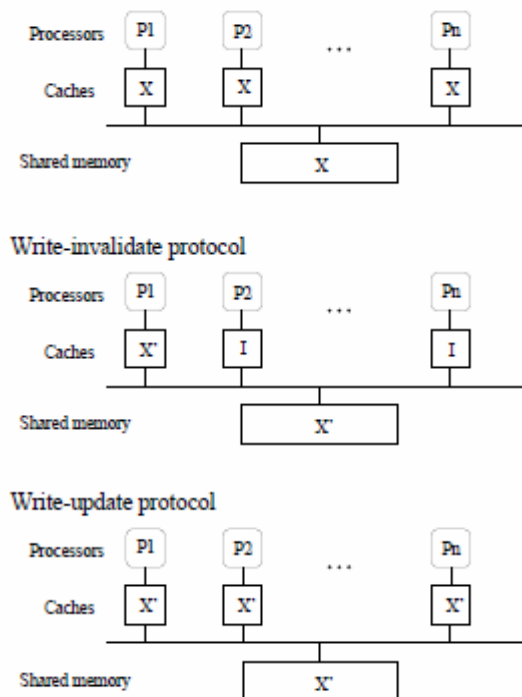


Figura 6. Protocoale pentru menținerea coerenței datelor

Protocolul *Write-Invalidate*: dacă P_1 scrie valoarea X' în variabila X din memorie atunci se va trimite un mesaj de invalidare tuturor celorlalte memorii cache pentru blocul care conține variabila X .

Protocolul *Write-Update*: dacă P_1 scrie valoarea X' în variabila X atunci se vor trimite mesaje de actualizare cu valoarea X' la toate celelalte memorii cache. Protocolul *Write-Update* necesită mai mult lărgime de bandă din partea rețelei, pentru că această metodă este echivalentă cu o operație de broadcast pentru actualizarea datelor.

Metoda *Write-Invalidate* suferă în schimb de un fenomen numit *false-sharing*; în acest caz se pot face anumite invalidări ale datelor care nu sunt necesare în program:

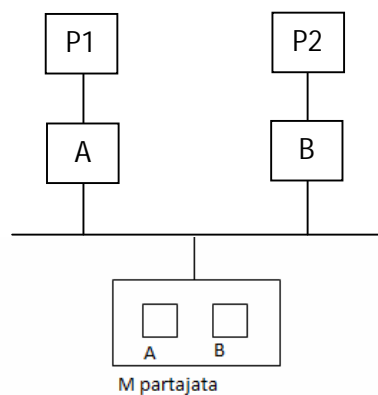


Figura 7. Situație în care poate apărea fenomenul de *false-sharing*

În exemplul de mai sus, A și B sunt două variabile ce există în același bloc de memorie. Dacă P1 modifică variabila A, atunci se invalidează și variabila B (invalidarea realizându-se la nivel de bloc), de aici din urmă invalidare nu este necesară.

1. Arhitectura UMA (*Uniform Memory Access*)

Această arhitectură se utilizează în sisteme multiprocesor tradiționale care folosesc o rețea de comunicație indirectă, de obicei o magistrală. În sistemele UMA, modalitatea de conectare și timpul de acces la memoria globală partajată sunt aceleași pentru toate procesoarele.

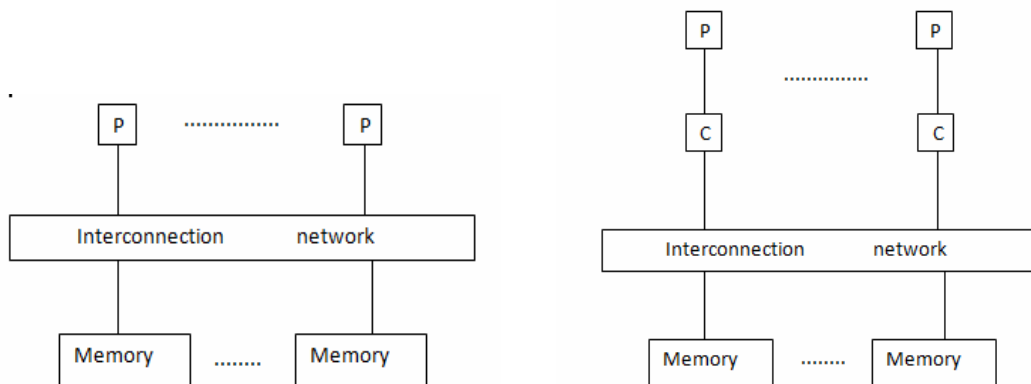


Figura 8. Arhitectura UMA

O clasă specială de calculatoare care folosește arhitectura UMA o reprezintă multiprocesoarele simetrice SMP (*Symmetric Multiprocessor*).

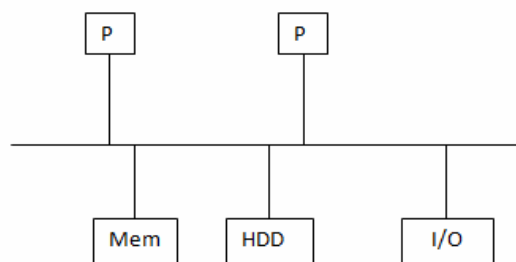


Figura 9. Arhitectura unui sistem SMP

În cazul unui sistem SMP, toate resursele sistem (memorii, hard disk-uri, dispozitive de I/O) sunt accesibile tuturor procesoarelor într-o manieră uniformă.

Câteva probleme ale calculatoarelor cu arhitectura UMA:

1) Aceste arhitecturi sunt foarte greu scalabile.

Arhitectura bazat pe magistrală suferă fenomenul de saturare dacă se conectează prea multe procesoare. Folosirea unei rețele *cross-bar* limitează scalabilitatea datorită costurilor. Folosirea rețelelor multistagiu este scumpă, adăugând o latență mare a accesului la memoriile partajate atunci când se mărește numărul de stadii.

Soluții:

a) Reducerea traficului pe magistrală prin folosirea memoriilor cache, dar și aici crescând numărul de procesoare crește numărul de evenimente *cache-miss* (crește numărul de blocuri cache invalide), evenimente care rezolvare va conduce din nou la saturarea magistralei.

b) Gruparea unui număr mic de procesoare împreună cu o singură memorie partajată într-un cluster și conectarea mai multor cluster printr-o rețea directă (arhitectura NUMA).

2) Poate apărea blocarea memoriei atunci când mai multe procesoare vor să acceseze memoria simultan.

Soluție:

Pentru rezolvarea problemei putem împărți memoria în mai multe blocuri (sau bancuri de memorie) astfel încât fiecare banc să se acceseze într-un anumit interval de adresă. De exemplu putem avea o memorie cu două bancuri, primul banc fiind accesat la adrese pare iar al doilea la adrese impare.

Dar și în cazul împărțirii memoriei în bancuri apar probleme dacă mai multe procesoare vor să acceseze simultan același banc de memorie.

Exemplu: Presupunem că avem mai multe procesoare care sunt conectate printr-o rețea multistagiu la mai multe bancuri de memorie (fig.10). În momentul în care mai multe procesoare accesează simultan un banc de memorie (B1) apare congestia sau saturarea rețelei:

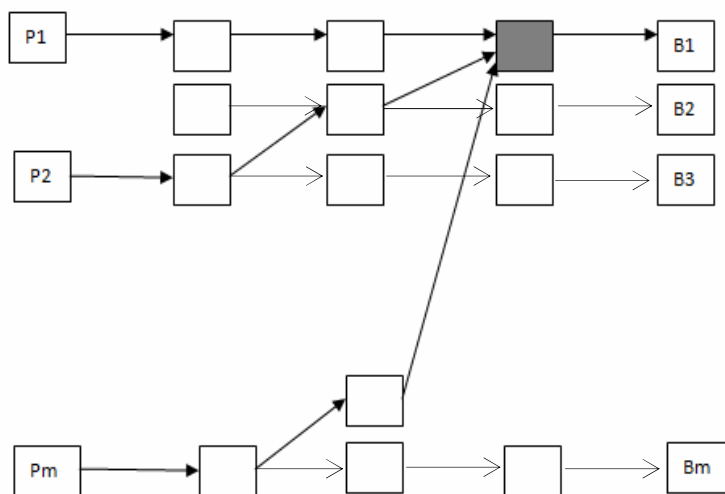


Figura 10. Congestia unei rețele multistagiu

În exemplul de mai sus, *switch*-ul conectat la bancul B1 reprezintă un punct de congestie.

Ca o soluție pentru acest problemă putem folosi combinarea mesajelor:

De exemplu se folosește o operație atomică numită *Fetch and Add* (o operație atomică nu poate fi divizată sau întreruptă de nici un alt proces). Operația *Fetch and Add* (x, e) este echivalentă cu înlocuirea variabilei x din memorie prin valoarea $x+e$.

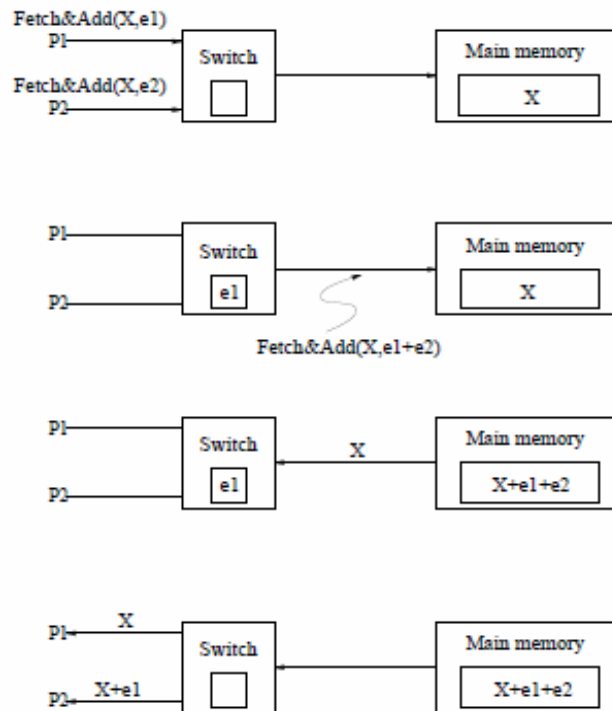


Figura 11. Accesul simultan a două procesoare la o locație de memorie prin operații *Fetch and Add*

Pasul 1: P1 și P2 vor să acceseze simultan variabila X prin operații atomice *Fetch and Add*.

Pasul 2: *Switch*-ul memorează valoarea de adunare $e1$ a primei operații, apoi combină cele două operații *Fetch and Add* într-una singură având ca valoare de adunare $e1+e2$, și trimite această operație ca un singur mesaj memoriei principale.

Pasul 3: Valoarea din memorie se actualizează și se trimite înapoi către *switch* un mesaj cu valoarea inițială a lui X .

Pasul 4: *Switch*-ul trimite valoarea inițială către cele două procesoare; astfel P1 va primi valoarea X în timp ce P2 va primi valoarea $X+e1$ (presupunem că P1 și P2 nu au memorii cache).

Acest exemplu ilustrează avantajul metodei de combinare a mesajelor: memoria este citită sau scrisă o singură dată la accesul simultan a mai multor procesoare, astfel încât nu va mai apărea fenomenul de saturare a rețelei.

2. Arhitectura NUMA (Non-Uniform Memory Access)

Un dezavantaj al arhitecturii UMA este acela că ea duce la realizarea de sisteme greu scalabile; pentru a înlătura acest dezavantaj s-au dezvoltat sistemele NUMA. În cadrul acestei arhitecturi, mai multe procesoare cu acces uniform la o memorie partajată sunt conectate într-un cluster, mai multe cluster fiind conectate printr-o rețea punct la punct scalabil, ce folosește pentru comunicare transmiterea de mesaje.

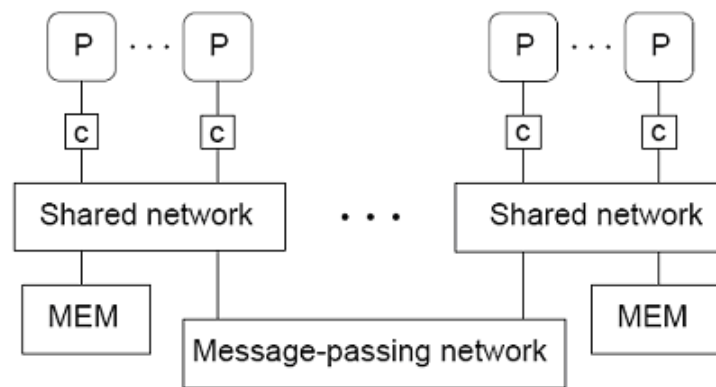


Figura 12. Arhitectura NUMA

Din punct de vedere al programării există un singur spațiu de memorie global, astfel că NUMA devine o arhitectură partajată din punct de vedere logic, deși din punct de vedere fizic ea este o arhitectură distribuită.

Când un procesor accesează memoria sistemului, controlerul de memorie cache (sau în alte implementări unitatea de management a memoriei MMU) verifică dacă locația adresată se găsește în memoria locală sau într-o memorie aflată la distanță. În cazul în care se accesează un bloc de memorie aflat la distanță (*remote*), se trimite un mesaj în rețeaua de cluster pentru accesarea blocului *remote* prin rețea. Acest mod de funcționare are ca consecință faptul că accesul la memoria locală este rapid, pe când accesul la memoria *remote* este lent (de aici vine și numele arhitecturii NUMA).

Pentru a reduce numărul de accesări ale memoriilor aflate la distanță, procesoarele aduc în memoria cache proprie blocurile *remote* dorite:

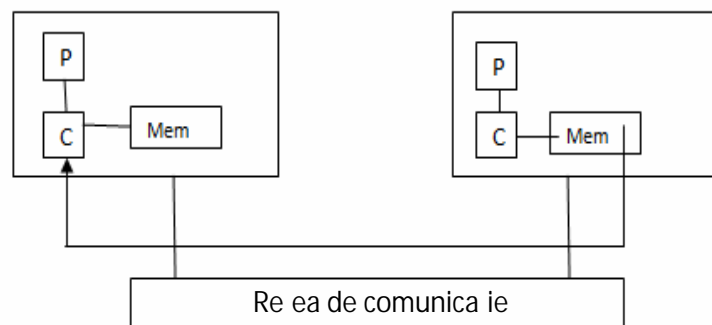


Figura 13. Aducerea în memoria locală a blocurilor *remote*

Dar în această situație va trebui rezolvată problema coerenței memoriilor cache care au citit aceeași variabilă. În acest scop, au fost create sistemele numite CC-NUMA (*Cache Coherent NUMA*), sisteme care folosesc un protocol de tip *Write-Invalidate* pentru menținerea coerenței memoriilor cache.

ARHITECTURA CC-NUMA (continuare)

Ascunderea laten ei la accesul memoriei de la distan în sistemele CC-NUMA

Arhitectura CC-NUMA bazat pe memorii cache reduce laten a acces rii memoriilor de la distan , dar este important s existe metode pentru acoperirea timpilor de a teptare datorat i laten ei acestor accese *remote*.

Exist 3 metode pentru acoperirea laten ei:

- Preînc rcarea datelor – prin care se încearc predic ia datelor necesare i înc rcarea lor în memoria cache înainte ca ele s fie propriu-zis accesate;
- Metoda *multithreading* – în care se folose te un *thread* separat pentru accesul la un bloc de memorie *remote*. Astfel, în timpul a tept rii accesului la memorie se poate comuta pe un alt *thread* care poate fi rulat:

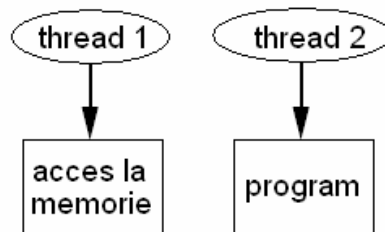


Figura 1. Metoda *multithreading*

- Folosirea unor modele relaxate de consisten a memoriei.

Modele de consisten a memoriei:

M1. Modelul nerelaxat de consisten secven ial (SC)

În acest model, toate accesele la memorie sunt atomice (nu pot fi întrerupte) i au o ordine precis :

- un procesor scrie într-un bloc partajat de memorie;
- toate celelalte procesoare primesc mesaje de invalidare pentru blocul respectiv;
- un alt procesor scrie, .a.m.d.

Conceptul de consisten secven ial se poate implementa printr-o memorie cu un singur port de intrare-ie ire. Procesoarele sunt conectate la memorie prin intermediul unui *switch* (fig.2). Acest *switch* alege la un moment dat un singur procesor care s aib acces la memorie. Un alt procesor care vrea s acceseze memoria va trebui s a tepte pân când procesorul anterior a terminat accesul.

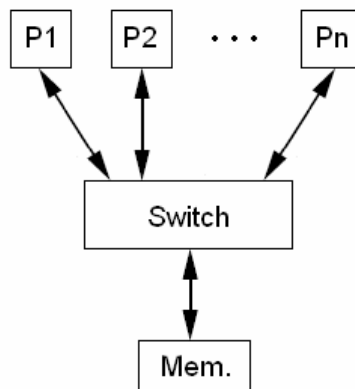


Figura 2. Implementarea modelului de consistență secvențial

Modelul de consistență secvențial nu este optimal datorită stricteții sale. Procesoarele superscalare nu folosesc modelul de consistență secvențial. Pentru ascunderea latenței accesului la memorie ele folosesc execuția *out of order* a instrucțiunilor și de asemenea scrierea datelor de ieșire în *buffere* de scriere, și nu scrierea direct în memorie.

M2. Modele de consistență relaxată a memoriei (RC)

Există mai multe modele RC:

- Modelul de consistență a procesoarelor – în care operații de tip LOAD pot devansa operații WRITE în cadrul unui program.
- Modelul cu ordine de memorare parțial – în care operațiile LOAD urmate de WRITE pot devansa operații WRITE precedente dintr-un program.
- Modelul de consistență slab – accesul la memorie se poate face fără a respecta ordinea instrucțiunilor din program (pentru aceasta nu trebuie să existe dependențe de date între instrucțiuni).

În toate modelele de consistență relaxată s-au introdus așa numite bariere de sincronizare. O barieră reprezintă o operație sincronizată care determină ca toate operațiile de acces la memorie care au fost lansate anterior să se finalizeze înainte de execuția operației sincronizate. Plasând barierele în locuri bine determinate din program se obțin avantajele specifice modelului de consistență secvențial, dar cu o performanță mai bună, pentru că accesul la memorie nu este restricționat între bariere.

Dezavantajele arhitecturii CC-NUMA:

Principalul dezavantaj derivă din faptul că datele *remote* sunt aduse în memoriile cache ale procesoarelor, memorii care au dimensiuni reduse. Astfel, dacă este necesar accesul la blocuri de date *remote* de dimensiuni mari, aceste blocuri nu vor putea fi memorate în întregime în memoriile cache, și ca urmare ele vor trebui încărcate secvențial de la distanță, timpul de acces crescând foarte mult.

Ca soluții posibile la această situație putem avea:

- 1) Creșterea dimensiunii memoriilor cache. În acest caz va crește însă și latența acceselor la memoriile locale (pentru că memoriile cache devin mai lente), și cresc de asemenea costurile de implementare.
- 2) Implementarea la nivelul sistemului de operare a unui mecanism de migrare a paginilor de memorie. Astfel, atunci când o dată *remote* este accesat frecvent, sistemul de operare poate decide replicarea la citire sau migrarea la scriere a paginii în care se găsește datele cerute. Astfel pagina va fi adusă în memoria locală a procesorului care a cerut datele. Această soluție are ca dezavantaj creșterea complexității metodei și funcționarea la o granularitate mare (la nivel de pagină).

ARHITECTURA COMA (CACHE ONLY MEMORY ARCHITECTURE)

Arhitectura COMA seamănă cu arhitectura NUMA, diferența fiind că memoriile locale DRAM funcționează precum memoriile cache (realizându-se o mapare de tip direct sau asociativ).

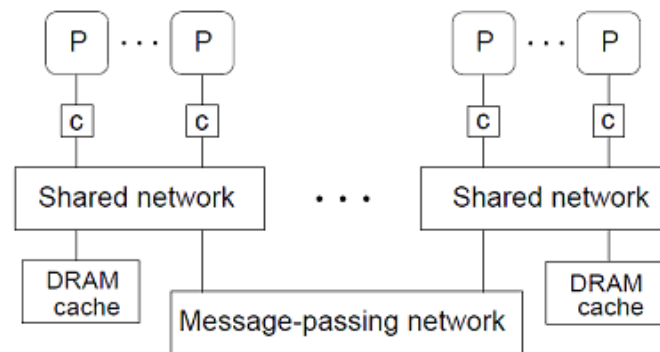


Figura 3. Arhitectura COMA

Datele încercate de la distanță (*remote*) sunt aduse în memoria locală și sunt păstrate în corespondență cu adresele lor din sistem. Blocurile de date nu au o locație fixă într-o anumită memorie locală, ci ele pot să migreze liber în sistem.

Arhitectura COMA utilizează o rețea cu transmitere de mesaje având o structură arborescentă. Fiecare nod sau *switch* din această rețea conține un director care indică datele conținute în arborele subiacent (al cărui părinte este nodul respectiv) (figura 4). Cum datele nu au o locație fixă în sistem, ele trebuie căutate în memoriile locale distribuite folosind structura de directoare a rețelei.

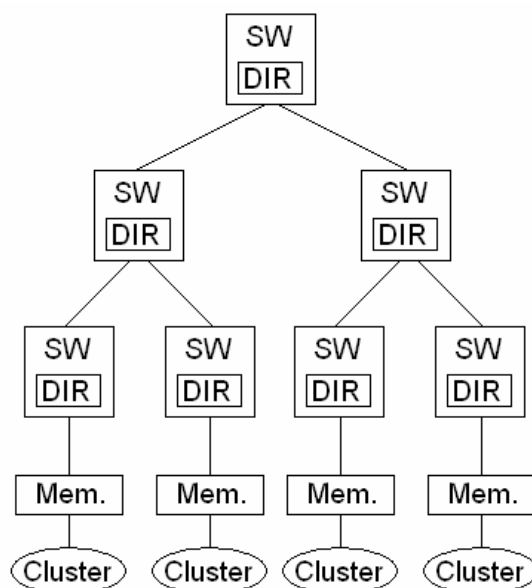


Figura 4. Structura arborescentă a rețelei de clustere (COMA)

Există și posibilitatea combinării mesajelor: dacă un *switch* primește două cereri pentru o aceeași dată de la un subarboare atunci cele două cereri se vor combina într-una singură, care va fi înaintată *switch*-ului părinte. Când datele cerute sunt disponibile, ele vor fi trimise înapoi către cele două *switch*-uri care au trimis cererile inițiale.

Pentru menținerea coerenței memoriilor cache, arhitectura COMA folosește protocolul *write-invalidate*, în care coerența se realizează la granularitatea fină a blocurilor din memoria locală. Acest protocol de coerență a arhitecturii COMA este însă mai greu de implementat decât la CC-NUMA.

De exemplu, o situație de care trebuie să se țină seamă este următoarea: dacă într-o memorie locală există o ultimă copie a unui bloc de date din sistem și dacă acel bloc urmează să fie eliberat din memorie ca urmare a algoritmului de înlocuire de blocuri din memoria cache, atunci acel bloc va trebui mutat în altă memorie locală a sistemului.

Avantajul arhitecturii COMA constă în memoriile locale cache de dimensiuni mari capabile să memoreze cantități mari de date de la distanță. Ca și cerințe constructive, la memoriile locale cache este necesar să se adauge memorii suplimentare atât pentru *tag*-urile cache cât și pentru informațiile de stare ale blocurilor cache.

Compararea între arhitectura COMA și arhitectura CC-NUMA:

Avantajul arhitecturii COMA față de CC-NUMA: arhitectura este mai flexibilă și suportă migrarea blocurilor de date. Aceasta se face la nivel hardware la o granularitate fină, reducându-se fenomenul de *false sharing*.

Dezavantajele arhitecturii COMA sunt următoarele:

- Arhitectura este mai costisitoare și mai greu de implementat;
- Protocolul de coerență este și el dificil de implementat (trebuie să se țină seama că ultimele copii a blocurilor de date să nu fie terse din sistem);
- Accesele la distanță sunt mai lente decât la CC-NUMA din cauza structurii arborescente a rețelei de comunicații.

Performanțele celor două arhitecturi depind de aplicații: dacă este necesar accesul la blocuri mari de date de la distanță, atunci arhitectura COMA va avea performanțe mai bune, în schimb dacă există o rată mare a invalidărilor blocurilor de date cache (datorită scrierilor repetate) atunci arhitectura COMA va avea o latență mai mare a acceselor *remote*, deci performanțe mai slabe, arhitectura CC-NUMA fiind mai potrivită.

ARHITECTURA S-COMA (SIMPLE COMA)

Pentru a reduce complexitatea și costurile arhitecturii COMA, a fost propus o variantă a acestei arhitecturi numită SIMPLE COMA. În această arhitectură, alocarea blocurilor în memoria locală de tip cache se face la granularitatea sistemului de operare, adică la nivelul paginilor.

Asemănător cu arhitectura SVM (*Shared Virtual Memory*), și aici sistemul de operare, prin unitatea de management a memoriei (MMU), determină dacă pagina în care se află data cerută se găsește sau nu în memoria locală. Dacă da, avem un eveniment *cache hit*, dacă nu, avem un eveniment *cache miss*. Cum MMU este cea care determină prezența unei pagini în memoria cache, nu mai este nevoie de memorii suplimentare pentru tagurile asociate blocurilor de memorie.

Spre deosebire de SVM, arhitectura S-COMA gestionează coerent blocurile de date prin hardware la o granularitate fină la dimensiunea blocurilor cache. Astfel este redus și aici fenomenul de *false sharing*. O altă consecință este că în procesul de migrare a datelor se transferă blocuri cache de dimensiuni mici, și nu pagini (fig. 5).

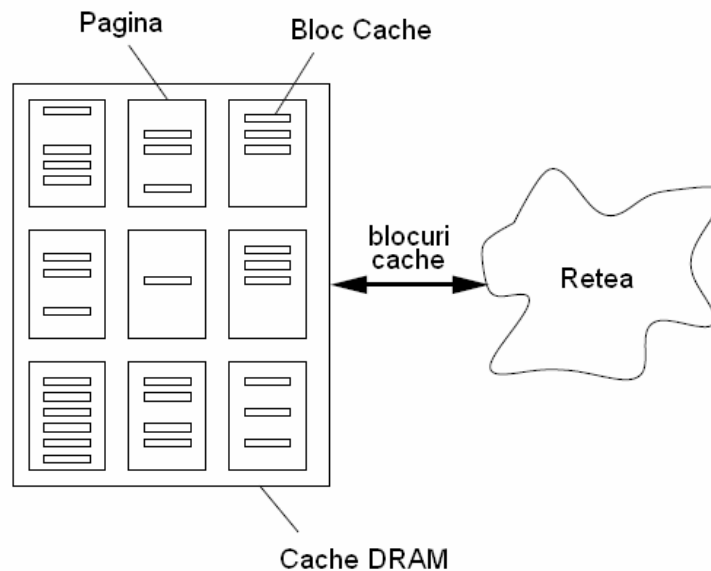


Figura 5. Gestionarea blocurilor de memorie la S-COMA

Dacă la un moment dat există un eveniment *cache miss* (pagina corespunzătoare nu se găsește în memorie) atunci:

- 1) se alocă o nouă pagină în memoria locală ;
- 2) blocul de date cerut este încărcat dintr-o altă memorie (*remote*).

O consecin ă a acestui mod de func ioneare este c ă o pagin ă din memorie poate fi la un moment dat umplut ă doar par țial cu blocuri de date. Astfel, la un eveniment *cache hit* (pagina c ăutată se afl ă în memorie), trebuie s ă se verifice și dac ă blocul de date cerut exist ă și este valid în acea pagin ă.

Arhitectura S-COMA poate implementa memorii cache locale complet asociative în care sistemul de operare poate înc ăca orice pagin ă la orice adres ă din memorie.

Dezavantaje ale arhitecturii S-COMA:

- 1) Evenimentele *cache miss* sunt tratate mai lent decăt la arhitectura COMA din cauza trat ării la nivelul sistemului de operare și nu prin hardware.
- 2) Fenomenul de *false replacement*: mecanismul de înlocuire a paginilor în memorie poate duce la înlocuirea unei pagini complet umplute cu blocuri de memorie printr-o pagin ă ce conține un singur bloc ce este necesar la un moment dat, reducându-se astfel eficiența lucrului cu date *remote*.

ARHITECTURA R-NUMA (REACTIVE NUMA)

Această arhitectură poate comuta dinamic între arhitectura CC-NUMA (când blocurile *remote* sunt aduse în memoriile cache ale procesoarelor) și arhitectura S-COMA (când blocurile de date *remote* sunt aduse în memoriile locale de dimensiuni mari).

COERENȚA MEMORIILOR CACHE ÎN SISTEME MULTIPROCESOR

Toate sistemele cu arhitectură multiprocesor folosesc protocoale pentru menținerea coerenței memoriilor cache. Protocoalele de coerență se bazează pe existența unor stări asociate blocurilor de memorie cache și a unor tranziții între stări. Aceste tranziții între stările unui bloc de memorie depind atât de acțiunile procesorului local, cât și de acțiunile altor procesoare ce dețin copii ale acelui bloc de memorie.

De exemplu, un bloc de memorie (**b**) situat în memoria cache a unui procesor (P1) va trece în starea invalid dacă un alt procesor (P2) scrie în copia acelui bloc de memorie pe care a făcut-o în propria memorie cache (fig. 1):

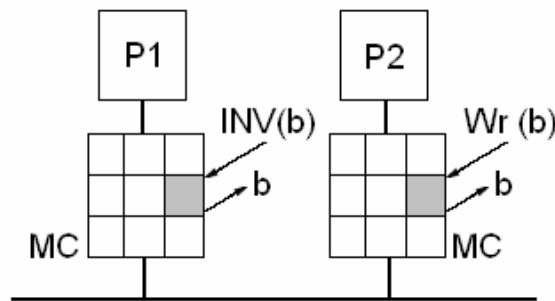


Figura 1. Invalidarea unui bloc de memorie a cărui copie a fost modificat

Din cele de mai sus rezultă necesitatea de a ține evidența tuturor copiilor care există pentru fiecare bloc de memorie cache din sistem.

Protocoalele de coerență în sistemele multiprocesor se împart în două clase:

- 1) Protocol cu urmărirea tranzacțiilor pe magistrală (*Snoopy Bus*)
- 2) Protocol bazat pe director

1) Protocolul *Snoopy Bus*

Acesta se aplică mai ales în cazul sistemelor UMA bazate pe magistrală. În cadrul acestui protocol, controlerul memoriei cache urmărește tranzacțiile pe magistrală și verifică dacă un alt procesor accesează un bloc de memorie care este prezent și în memoria sa locală. Dacă există un astfel de procesor și dacă accesul este unul de scriere în memorie, atunci controlerul cache va invalida blocul din memoria locală. Protocolul *Snoopy Bus* poate să fie atât de tip *write-invalidate* cât și de tip *write-update*.

2) Protocolul bazat pe director

Aici există o memorie director care poartă strează localizarea tuturor copiilor din sistem pentru toate blocurile din memoriile cache. Acest protocol se folosește pentru menținerea coerenței memoriilor cache în sistemele CC-NUMA și respectiv COMA. Aceste arhitecturi se bazează pe rețele cu transmitere de mesaje în mod *non-broadcast*. În aceste arhitecturi, protocolul *Snoopy Bus* (care necesită urmărirea întregului trafic din rețea) nu ar fi potrivit. De obicei protocoalele bazate pe director sunt de tipul *write-invalidate*.

Există mai multe metode de implementare pentru memoria director:

- a) **Implementarea Full-Map.** Aici pentru fiecare bloc dintr-o memorie cache se presupune localizarea copiilor blocului respectiv din toate memoriile cache ale sistemului. Astfel, pentru fiecare bloc **X** din memoria cache se presupune un vector de N biți, unde N este egal cu numărul de procesoare (și de memorii cache) din sistem; în acest vector, bitul b_i indică dacă o copie a blocului respectiv se găsește în memoria cache de indice i (fig. 2).

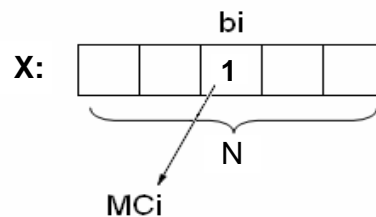


Figura 2. Vectorul din memoria director corespunzător blocului **X**

În acest caz, dacă un procesor scrie în blocul **X**, atunci se vor invalida toate copiiile blocului **X** din sistem (fig. 3a-3c):

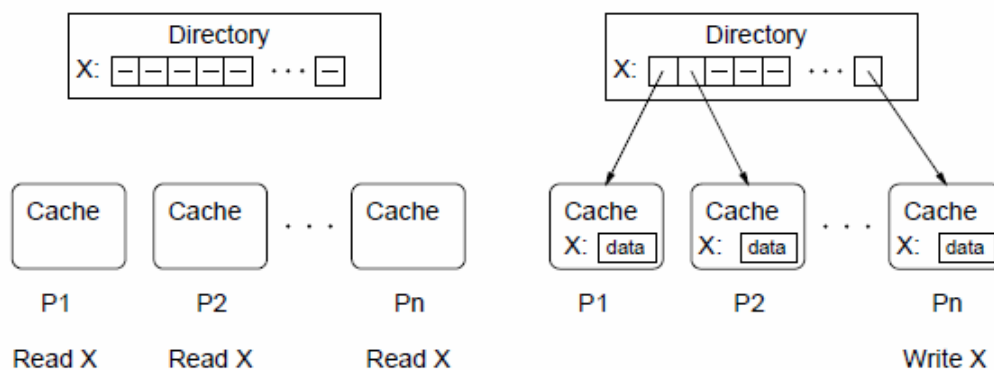


Figura 3a. Procesoarele **P1**, **P2**, **Pn** citesc blocul **X** în memoria cache proprie

Figura 3b. Se setează biții corespunzători din vectorul blocului **X** din memoria director; apoi procesorul **Pn** scrie în blocul **X**

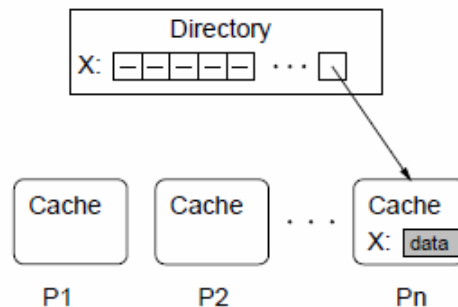


Figura 3c. Copiile blocului **X** din **P1** și **P2** se invalidează prin resetarea biților din memoria director

Un dezavantaj al acestei metode este dat de m rimea memoriei director auxiliare, dimensiune care are ordinul $O(N^2)$.

De exemplu, dac avem N memorii cache în sistem, i presupunem c fiecare memorie cache are un num r de P blocuri, rezult un num r total de $N \times P$ blocuri de memorie cache; acestea corespund cu $N \times P$ vectori în memoria director, rezultând dimensiunea total a memoriei director: $(N \times P) \times N = (N^2) \times P$.

- b) O implementare mai restrictiv a memoriei director este cea numit **Limited Map** (director cu dimensiune limitat). În acest caz poate exista doar un num r limitat de copii ale unui bloc de memorie în întregul sistem:

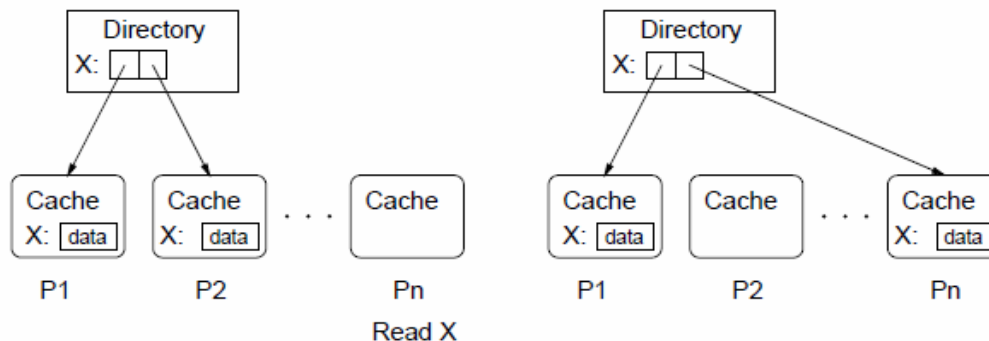


Figura 4. Implementarea **Limited Map** a memoriei director

În figura 4 se arat un exemplu în care se permite existen a doar dou copii ale fiec rui bloc de memorie în sistem. Avantajul acestei metode este c scade dimensiunea memoriei director; dezavantajul este c se limiteaz paralelismul.

- c) **Metoda cu director înl n uit (Chained Directory)**

În cadrul acestei metode, o intrare în memoria director corespunz toare unui bloc cache con ine un pointer la o list de copii ale blocului respectiv. Când un procesor recep ioneaz un bloc de memorie pentru a-l citi, se pune o referin la acel procesor în capul listei memoriei director, iar vechea referin din director e stocat în memoria cache a procesorului care a citit blocul. Astfel se creeaz o list înl n uit a tuturor copiilor unui bloc din sistem (fig. 5).

În acest caz, ac iunea de a invalida copiiile în bloc de memorie se reduce la o opera ie de parcurgere a listei înl n uite. Dezavantajul metodei este c la înlocuirea unui bloc dintr-o memorie cache, referin a la acest bloc trebuie îndep rtat din lista în care este inclus , crescând astfel complexitatea opera iei.

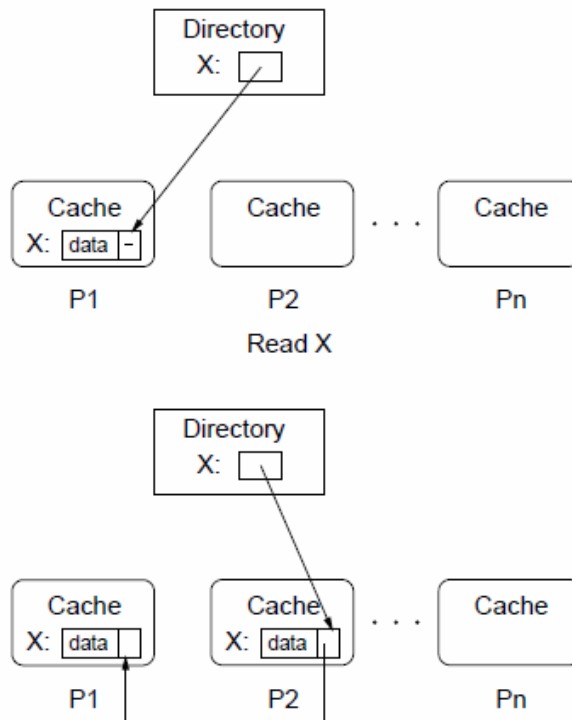
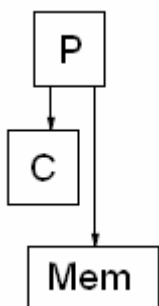


Figura 5. Metoda cu director în l n uit

ST RI I TRANZI II ÎN PROTOCOALELE DE COEREN

Protocoloalele de coeren pentru memoriile cache se bazeaz pe st ri i tranzi ii între st ri. De exemplu, procesul de invalidare a unui bloc de memorie cache presupune tranzi ia între dou st ri: valid i invalid. Dacă un bloc se afl în starea invalid, atunci procesorul nu va mai putea utiliza datele din acel bloc.

În continuare se prezint dou exemple de protocoale de coeren : pentru memoriile **Write Through**, respectiv **Write Back**, împreun cu diagramele de st ri.



1) Protocol de coeren pentru memoriile de tip **Write Through** Cache.

Fiecare bloc din memoria cache se poate afla în doar dou st ri: valid sau invalid. Un bloc se consider invalid i în momentul în care el este înlocuit din memoria cache. Se definesc trei ac iuni pentru blocul respectiv: citire bloc (**R** -Read), scriere bloc (**W** -Write) i înlocuire (**Z** -Replace). Aceste trei ac iuni se pot efectua de c tre memoria cache curent , de indice i, fiind notate: R(i), W(i), Z(i), sau aceste ac iuni se pot efectua în memoria cache de indice j, notate prin R(j), W(j), Z(j).

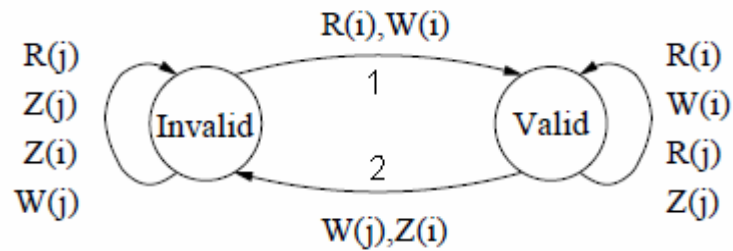
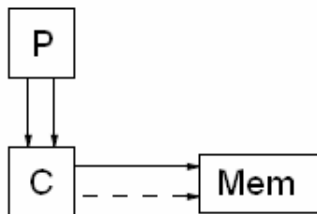


Figura 6. Diagrama de stări și tranziții pentru memoria *Write Through*

În diagrama de mai sus, tranzițiile între stările valid-invalid se realizează în următoarele situații:

1) Dacă blocul este invalid, blocul devine valid prin citire $R(i)$ - deoarece el este adus în memoria locală din locația sursă, sau prin scriere $W(i)$ - când el este actualizat în propria memorie cache.

2) Dacă blocul este valid, la scrierea blocului către alt procesor în altă memorie cache $W(j)$, blocul curent se invalidează. De asemenea, blocul se invalidează și la înlocuirea sa din memoria cache ($Z(i)$).



2) Protocol pentru memoria de tip **Write Back** cache.

Aici există trei stări în care se poate afla un bloc dintr-o memorie cache: invalid (**INV**), partajat (**S** - *Shared*) sau modificat (**M** - *Modified*).

Diagrama de stări și tranziții pentru memoria *Write Back* este ilustrată în figura 7:

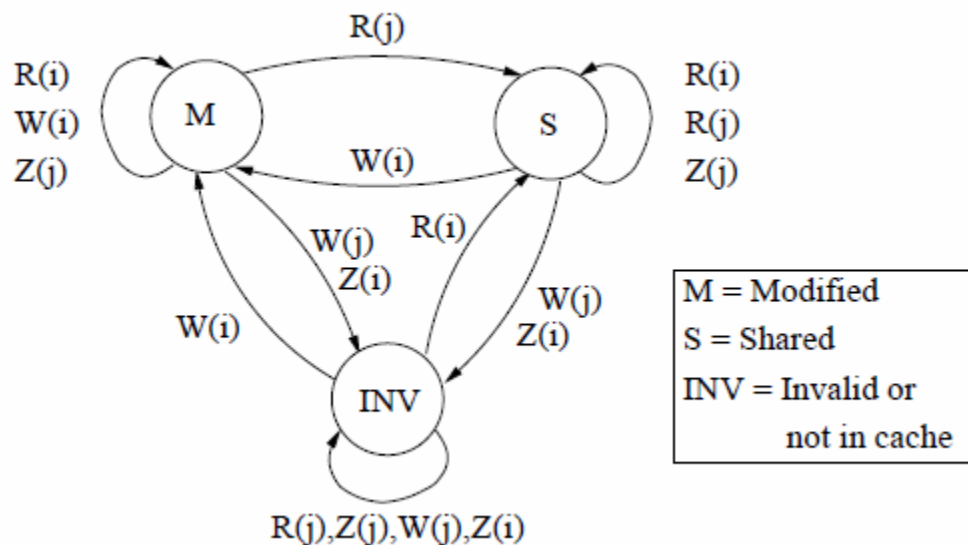


Figura 7. Diagrama de stări și tranziții pentru memoria *Write Back*

În figura 8 sunt ilustrate câteva situa ii posibile în care se poate găsi un bloc de memorie cache (de tipul **Write Back**).

În figura 8 (a) este prezentată situa ia inițială, în care avem două module de memorie cache C_i și C_j (în stânga, respectiv dreapta), și memoria principală (în centru). Să considerăm un bloc X din memoria centrală. La început, el nu este încărcat în memoriile cache; el apare ca *Invalid* în ambele memorii C_i și C_j .

În figura 8 (b) memoriile cache citesc blocul X din memoria principală. În ambele memorii cache, blocul X devine *Shared*, adică partajat; copiile multiple ale aceluiași bloc de memorie vor fi întotdeauna în starea *Shared*.

În situa ia din figura 8 (c), procesorul asociat cu memoria cache C_i scrie în blocul X din C_i ; ca urmare, blocul devine *Modified* (modificat) în C_i , iar în toate celelalte memorii C_j care dețineau copia partajată a blocului X , starea acestui bloc devine *Invalid*. Astfel, la orice moment de timp, doar o singură copie a unui bloc de memorie din sistem se poate găsi în starea *Modified*. Acesta va fi scris în memoria principală în momentul în care va fi necesară eliminarea lui din memoria cache.

Dacă apoi o altă memorie C_j citește blocul X (situa ie ilustrată în figura 8 (d) și detaliat în paragraful următor), blocul din memoria C_i aflat în starea *Modified* devine *Shared*, fiind totodată scris și în memoria principală.

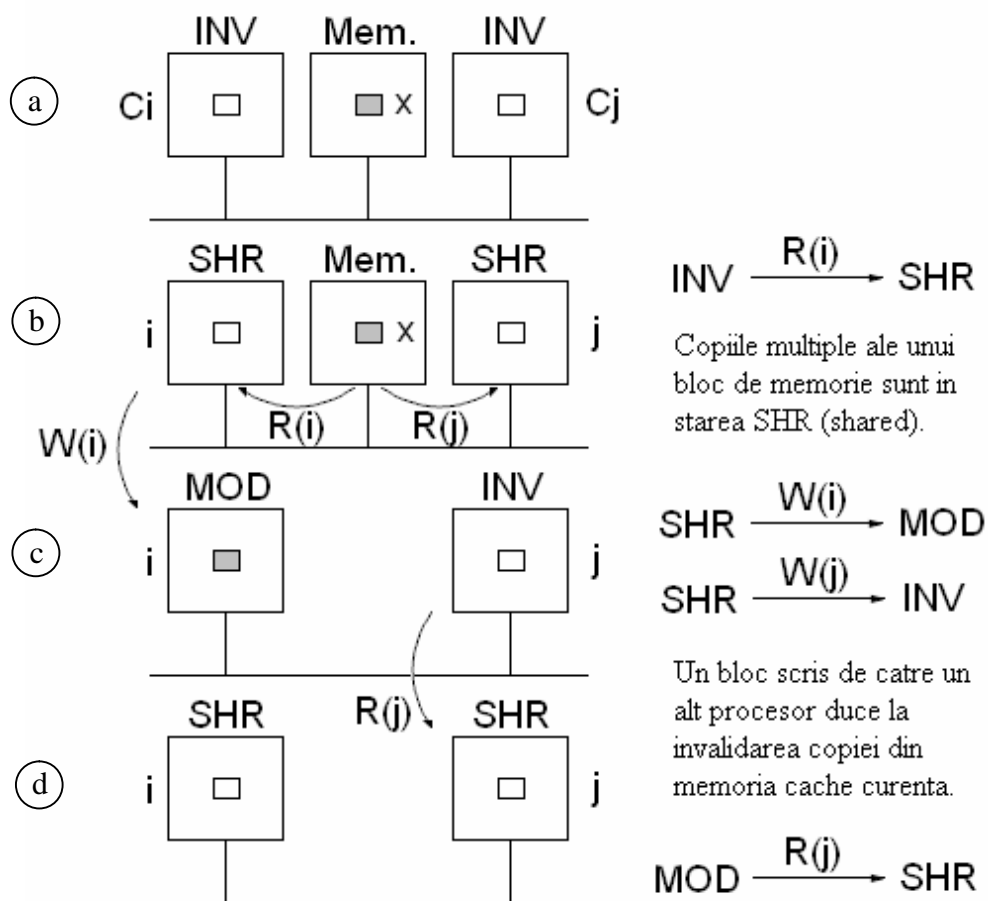


Figura 8. Stări și tranziții pentru un bloc de memorie din categoria **Write Back**

Compara ie între memoriile *Write Back* i *Write Through*

Configura ia *Write Through* cache are avantajul c memoria principal a sistemului este întotdeauna consistent . Astfel, dac apare un eveniment *cache miss* (blocul c utat nu se reg se te în memoria cache), memoria cache poate s încarce imediat blocul dorit din memoria principal .

În cazul memoriei *Write Back*, memoria cache nu poate înc rca blocul direct din memoria sistem pentru c el poate exista în starea *Modified* într-o alt memorie cache. Aceast situa ie se poate rezolva în 2 moduri:

- Modul 1.** a) Tranzac ia în curs de derulare cu memoria principal se opre te;
b) Memoria cache care de ine copia modificat scrie blocul actualizat în memoria principal i modific starea blocului în *Shared*;
c) Tranzac ia cu memoria principal se reia, ob inându-se copia actualizat a blocului cerut.

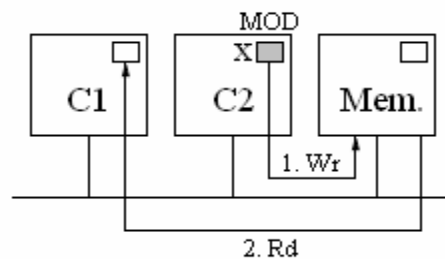


Figura 9. Prima modalitate de citire a unui bloc aflat în starea *Modified*. Memoria cache C2 scrie mai întâi blocul în memoria principal

- Modul 2.** a) Tranzac ia cu memoria principal se redirec ioneaz c tre memoria cache care de ine copia modificat a blocului;
b) Memoria cache de la care s-a citit scrie blocul în memoria principal i îi modific starea în *Shared*.

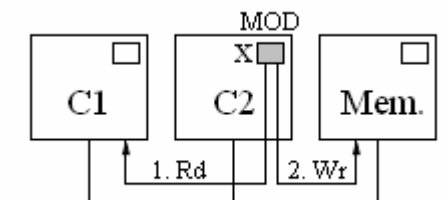


Figura 10. A doua modalitate de citire a unui bloc aflat în starea *Modified*. Memoria cache C1 cite te direct din memoria în care se g se te blocul modificat

Prima metod cre te laten a evenimentelor *cache miss*. A doua metod are o complexitate mai mare la implementare.

Dezavantajul memoriilor *Write Through* îl constituie faptul c ele cauzeaz un trafic de re ea suplimentar datorit scrierilor frecvente din memoria cache în memoria principal . În plus, un protocol de coeren pentru memoriile *Write Through* ce utilizeaz dou st ri: valid / invalid, mai are dezavantajul c pentru fiecare scriere a unui bloc de c tre o memorie cache trebuie invalidate toate copiile acelui bloc din sistem; în schimb, la memoriile *Write Back*, dac blocul se g se te în starea *Modified*, atunci va exista o singur copie actualizat în sistem care trebuie invalidat .