





Calculatoare Numerice (2)

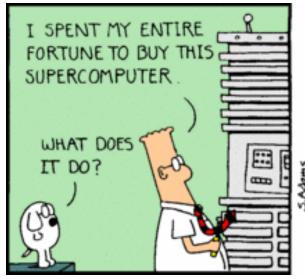
- Cursul 12 -

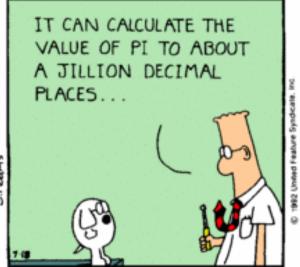
Multiprocesoare 2

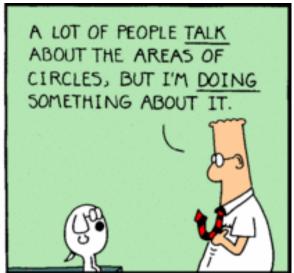
Facultatea de Automatică și Calculatoare Universitatea Politehnica București

Comic of the day









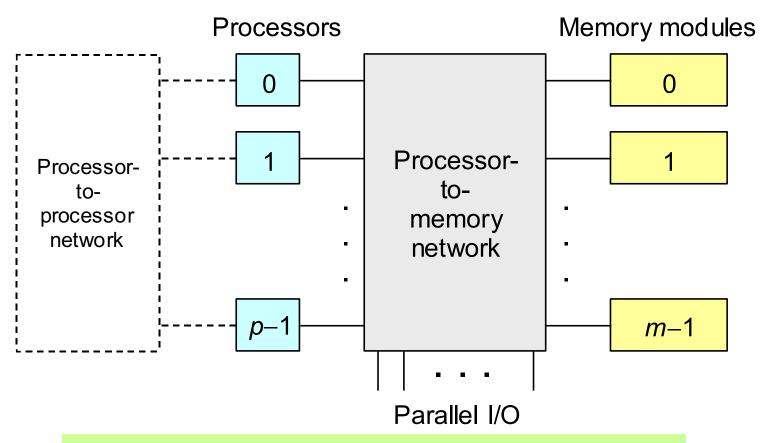
http://dilbert.com/strips/comic/1992-07-18/





Memoria partajată centralizată





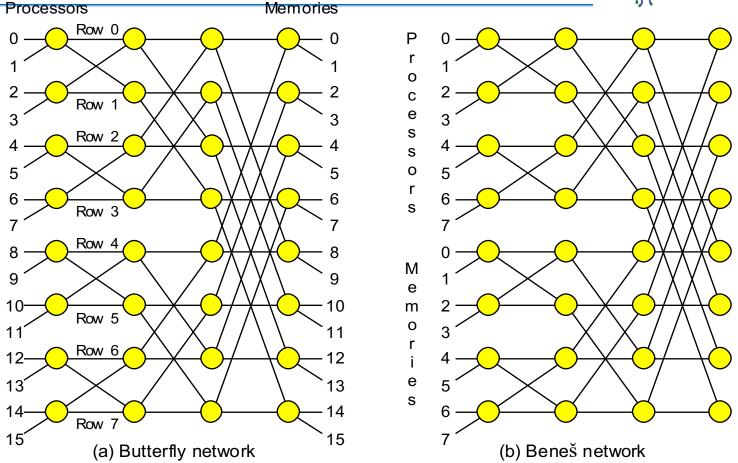
Structură multiprocesor cu memorie partajată





Rețele de interconectare Procesor-Memorie





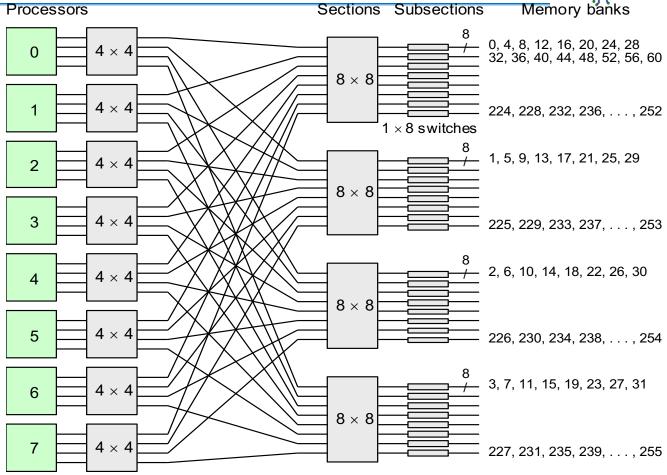
Rețelele fluture și Beneš: exemple de rețele de interconectare procesormemorie





Rețele de interconectare Procesor-Memorie





Interconectarea a opt procesoare la 256 bancuri de memorie la Cray Y-MP (1988), un supercomputer cu procesoare vectoriale multiple

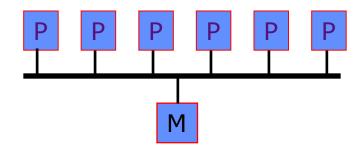




Consistența secvențială

Model de memorie





"A system is sequentially consistent if the result of any execution is the same as if the operations of all the processors were executed in some sequential order, and the operations of each individual processor appear in the order specified by the program"

Leslie Lamport

Sequential Consistency = întrețesere arbitrară cu păstrarea ordinei referințelor la memorie pentru programele secvențiale





Consistența secvențială



Task-uri secvențiale concurente: T1, T2

Variabile partajate: X, Y (inițial X = 0, Y = 10)

T1: T2: Store (X), 1
$$(X = 1)$$
 Load R₁, (Y) Store (Y), 11 $(Y = 11)$ Store (Y'), R₁ $(Y' = Y)$ Load R₂, (X) Store (X'), R₂ $(X' = X)$

Care sunt răspunsurile corecte pentru X' și Y'?

$$(X',Y') \in \{(1,11), (0,10), (1,10), (0,11)\}$$
?

Dacă y este 11 atunci x nu poate fi 0





Consistența secvențială



Consistența secvențială impune mai multe contrângeri de ordonare de memorie ca și cele impuse de dependențele de memorie ale programelor uni-procesor (\longrightarrow)

Care sunt cele din exemplele noastre?

T1:

Store (X), 1
$$(X = 1)$$
Store (Y), 11 $(Y = 11)$
Store (Y'), R_1 $(Y' = Y)$
Load R_2 , (X)

Cerințe adiționale SC

Store (X'), R_2 $(X' = X)$

Poate un sistem cu cache și out-of-order execution să pună la dispoziție o imagine consistentă secvențial a memoriei?





Excluziunea mutuala si instructiuni blocante

```
Instrucțiuni blocante atomice read-modify-write
e.g., Test&Set, Fetch&Add, Swap
vs
Instrucțiuni atomice non-blocante read-modify-write
e.g., Compare&Swap,
Load-reserve/Store-conditional
vs
```

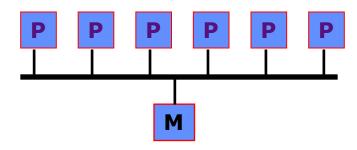
Protocoale bazate pe operații Load Store obișnuite

Performanța depinde de mai mulți factori:

degree of contention,
cache-uri,
out-of-order execution și Loads & Stores

Probleme în implementarea Consistenței Secvențiale





Implentarea CS este complicată de două probleme

• Capabilități de execuție *Out-of-order*

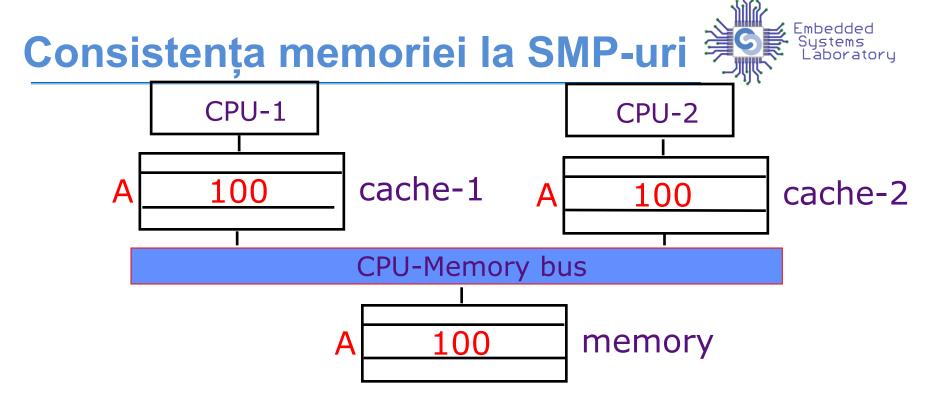
Load(a); Load(b) yesLoad(a); Store(b) yes if a \neq b Store(a); Load(b) yes if a \neq b Store(a); Store(b) yes if a \neq b

• Cache-uri

Cache-urile pot preveni ca efectul unui store să fie văzut de alte procesoare







Presupunem că CPU-1 actualizează A la 200. write-back: memoria si cache-2 au valori vechi write-through: cache-2 are valoarea veche

Contează aceste valori neactualizate? Cum este văzută memoria partajată de software?





Write-back Caches & SC



T1 is executed

prog T1 ST X, 1 ST Y,11

cache-1

X = 1Y = 11 memory

X = 0Y = 10X' =

cache-2

prog T2 LD Y, R1 ST Y', R1 LD X, R2 ST X',R2

cache-1 writes back Y

X=1

X = 0Y = 11

T2 executed

X = 1Y = 11 X = 0Y = 11X'=

cache-1 writes back X

X = 1Y = 11

X = 1Y = 11X'=

Y = 11Y' = 11X = 0

Y = 11Y' = 11X = 0

 cache-2 writes back X' & Y'

X = 1Y = 11

X = 1Y = 11X' = 0

Write-through Caches & SC



prog T1					
ST	X, 1				
ST	Y,11				

• T1 executed

• T2 executed

Nici cache-urile write-through nu mențin consistența secvențială

Menținerea consistenței secvențiale (CS)



CS este suficientă pentru programe tip producer-consumer și cu excluziune mutuală (e.g., Dekker)

Copiile multiple ale unei locații în diferite Cache-uri pot cauza degradarea CS.

Este nevoie de suport hardware pentru

- doar un singur procesor la un moment dat are permisiune de write la o locație de memorie
- nici un procesor nu poate să încarce o copie a vechii locații după o scriere
 - ⇒ Protocoale de coerență a cache-ului





Protocoale de menţinere a coerenţei cache pentru CS



write request:

adresa este *invalidată* (*actualizată*) în toate cache-urile înainte (după) o operație de write

read request:

dacă o este gasită o copie "murdară" într-un cache, se efectuează un write-back înainte de citirea memoriei

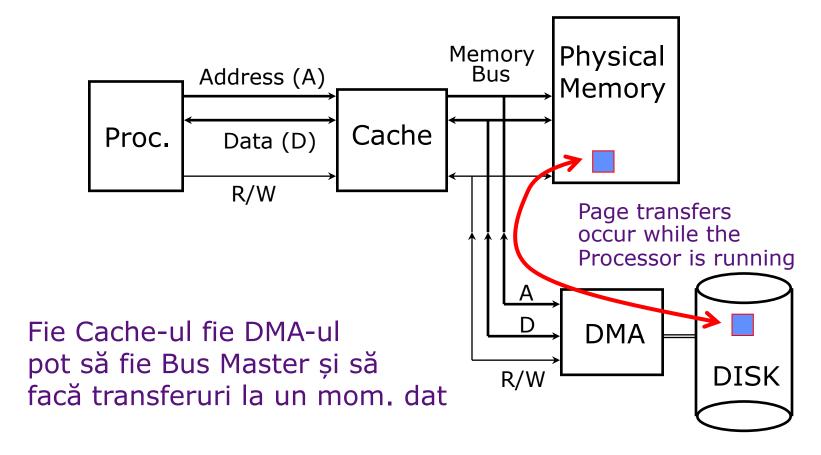
Ne vom concentra pe protocoale de Invalidare și nu pe protocoale Update





Warmup: Parallel I/O





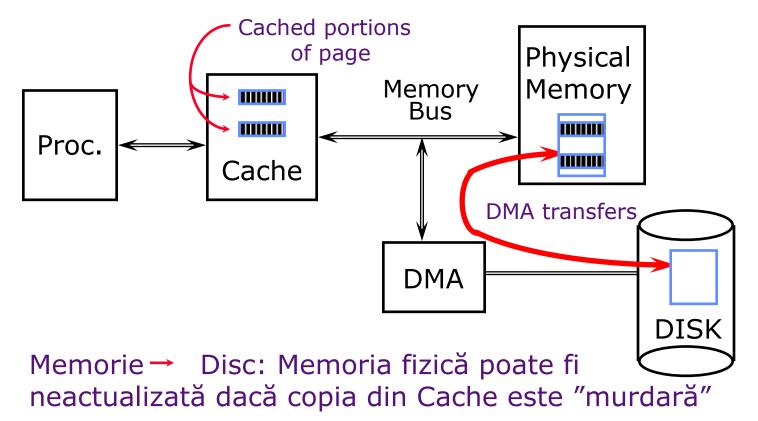
(DMA vine de la Direct Memory Access)





Probleme cu Parallel I/O





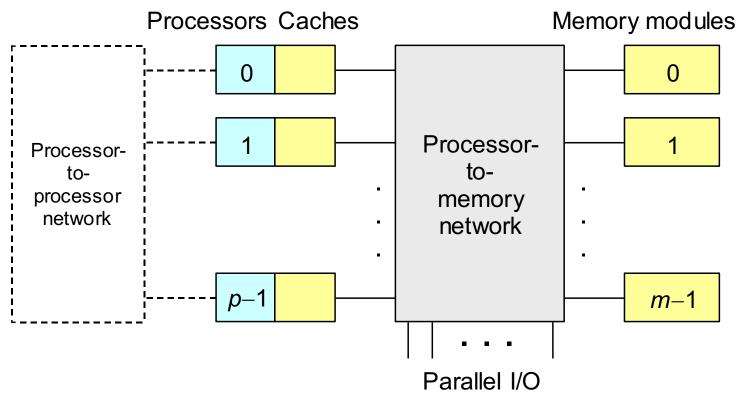
Disc → Memorie: Cache-ul poate să conțină date vechi și să nu vadă write-urile la memorie





Cache-uri multiple și coerența cache-urilor





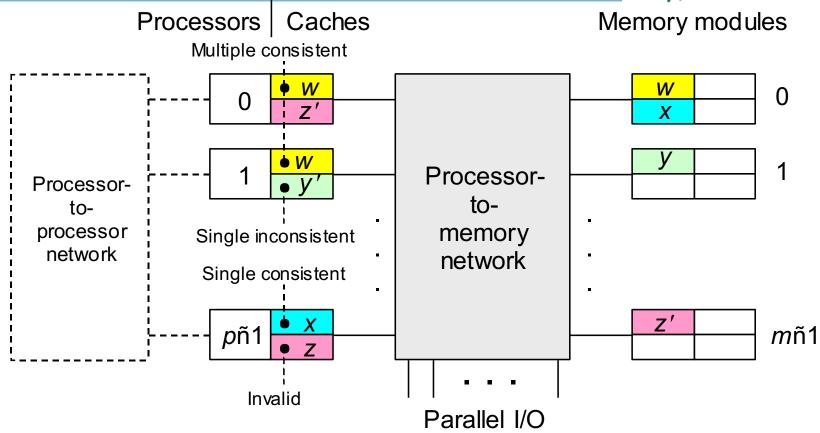
Memoriile cache dedicate fiecărui procesor reduc traficul la memoria procipală (prin rețeaua de interconectare) dar introduc o serie de probleme de consistență.





Starea copiilor de date





Diferite tipuri de blocuri de date din cache pentru un procesor paralel cu memorie principală centralizată și cache-uri locale pentru fiecare procesor

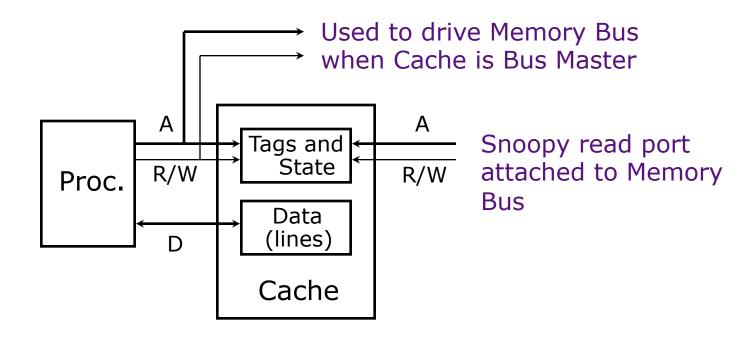




Snoopy Cache Goodman 1983



- Idee: Să avem un cache care spionează (snoop upon) transferurile DMA, și atunci "do the right thing"
- Etichetele snoopy cache sunt dual-port







Acțiunile Snoopy Cache pentru DMA



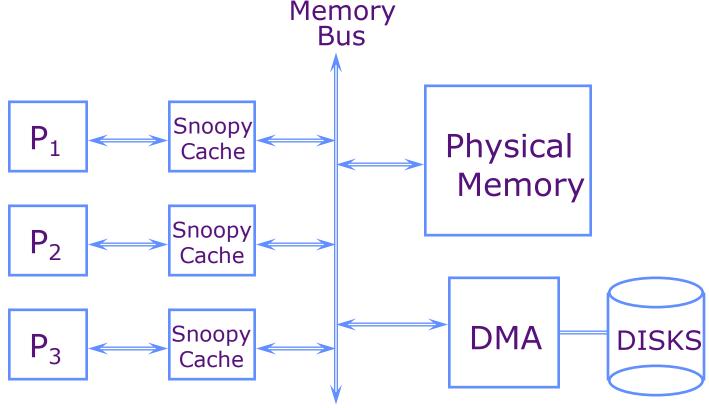
Observed Bus Cycle	Cache State	Cache Action	
	Address not cached	No action	
DMA Read	Cached, unmodified	No action	
Memory → Disk	Cached, modified	Cache intervenes	
	Address not cached	No action	
DMA Write	Cached, unmodified	Cache purges its copy	
Disk →Memory	Cached, modified	???	





Shared Memory Multiprocessor



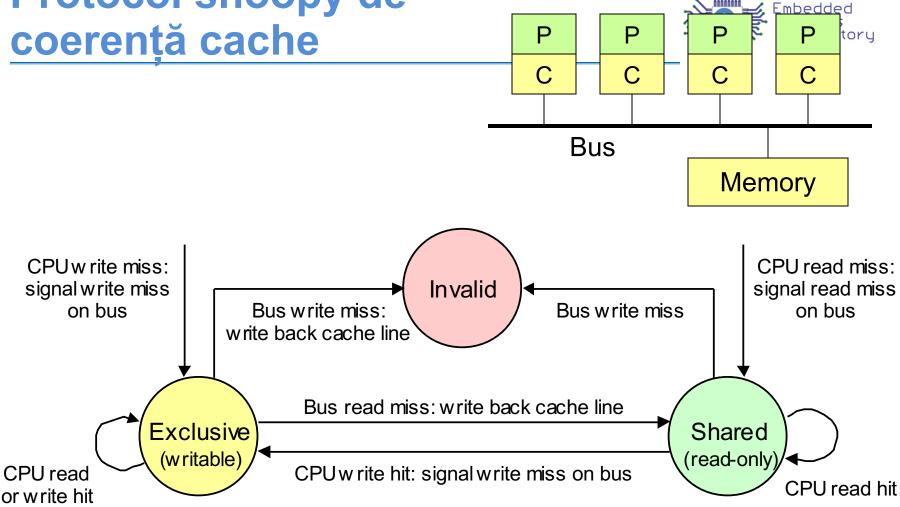


Folosește mecanismul snoopy pentru a păstra consistența memoriei pentru toate procesoarele





Protocol snoopy de coerentă cache



Automat FSM pentru un protocol de coerență cache ce folosește cache write-back





Diagrama de tranziții pentru Cache

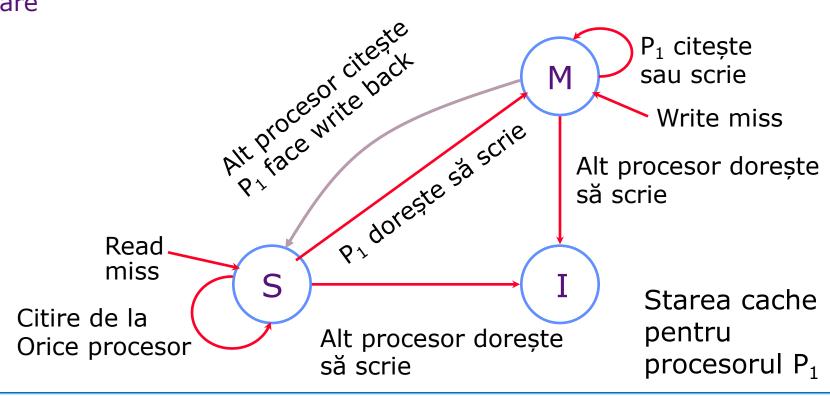
Protocolul MSI



S: Shared I: Invalid

Address tag

Biţi stare



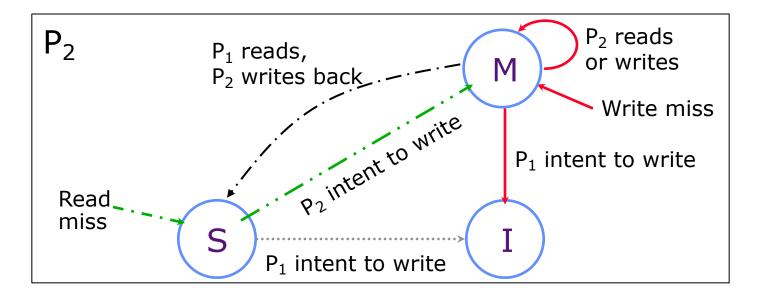




Exemplu cu două procesoare

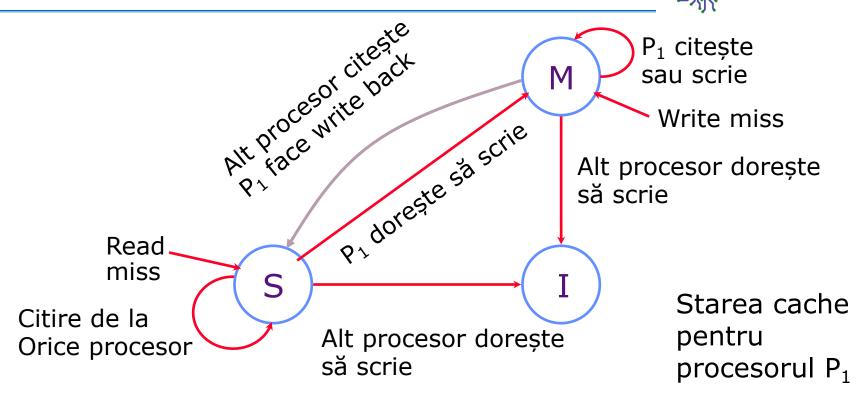
(Citesc și scriu aceeași linie din cache)

P₁ reads P_1 P₁ reads P₂ reads, or writes M P₁ writes P₁ writes back Write miss P₂ reads P₂ writes P₂ intent to write P₁ reads Read P₁ writes miss P₂ writes P₂ intent to write P₁ writes



Observație





- Dacă o linie este în starea M atunci nici un alt cache nu poate să aibă o copie a linei!
 - Memoria rămâne coerentă, nu pot exista copii diferite





MESI: Protocol MSI îmbunătățit

performanțe mărite pentru date locale



Fiecare linie are un tag

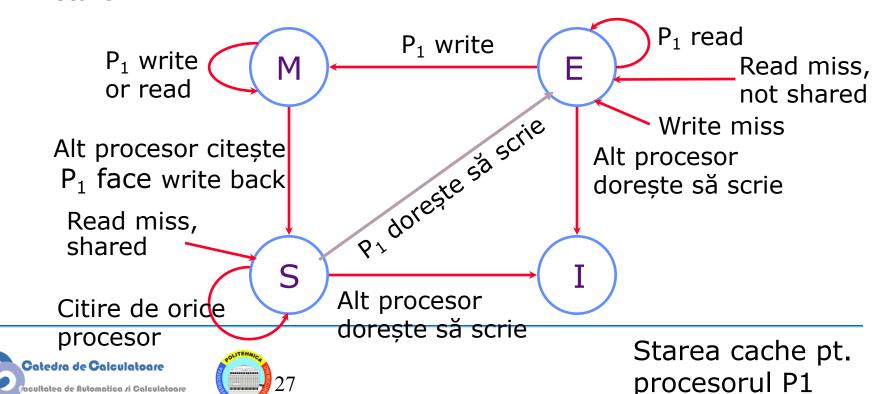
Address tag
Biţi
stare

M: Modified Exclusive

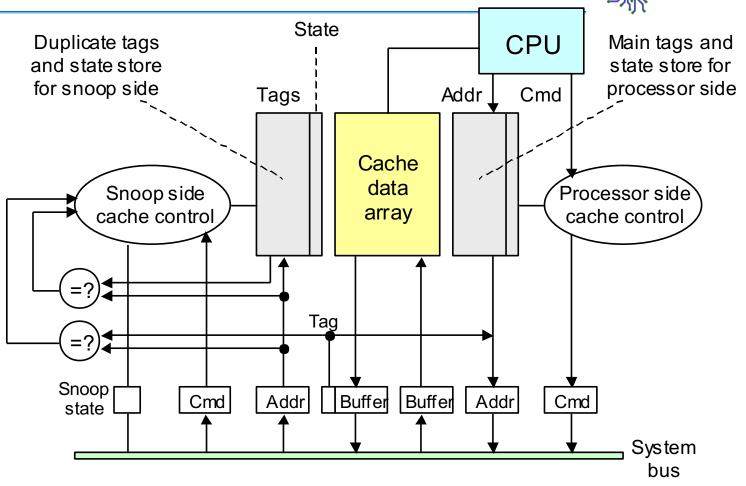
E: Exclusive, unmodified

S: Shared

I: Invalid



Implementarea Algoritmului Snoopy Cache



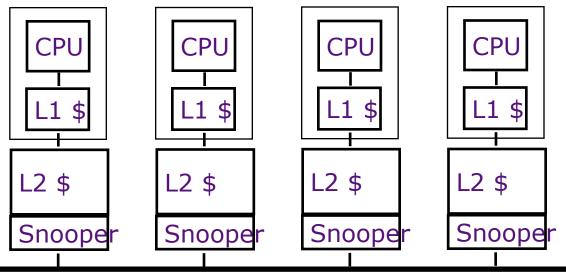
Structura principală a unui algoritm snoopy de coerență cache





Snoop optimizat cu cache-uri Level-2





- Procesoarele au de obice cache pe două niveluri
 - mic L1, mare L2 (de obicei ambele on chip acum)
- Proprietatea de incluziune: intrările din L1 trebuie să fie în L2

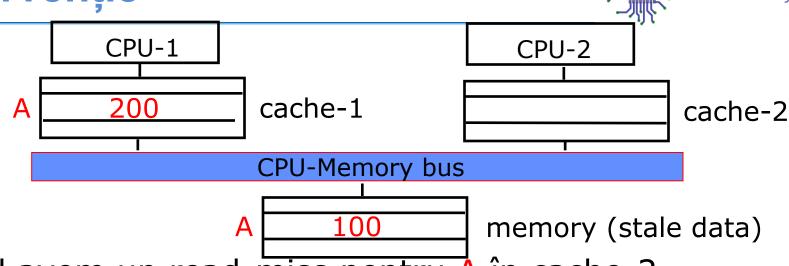
invalidare în L2 \Rightarrow invalidare în L1

Datedra de Calculatoare

ultatea de Automatica si Calculatoare

• Snooping în L2 nu afectează lățimea de bandă CPU-L1
Ce probleme pot să apară?

Intervenție



Când avem un read-miss pentru A în cache-2, se emite pe bus un read request pentru A

- Cache-1 trebuie să facă vizibilă și să își schimbe starea în "shared"
- Memoria poate să răspundă și ea la cerere!

Știe memoria că are date vechi?

Cache-1 trebuie să intervină prin controllerul de memorie pentru a da datele corecte pentru cache-2





False Sharing



state	blk addr	data0	data1	 dataN
				 01 01 0 01 1

Un bloc cache conține mai mult de un cuvânt de date

Coerența cache-ului este făcută la nivel e bloc și nu la nivel de cuvânt

Presupunem că P_1 scrie word_i și P_2 scrie word_k și ambele cuvinte au aceeași adresă de bloc.

Ce se poate întâmpla?

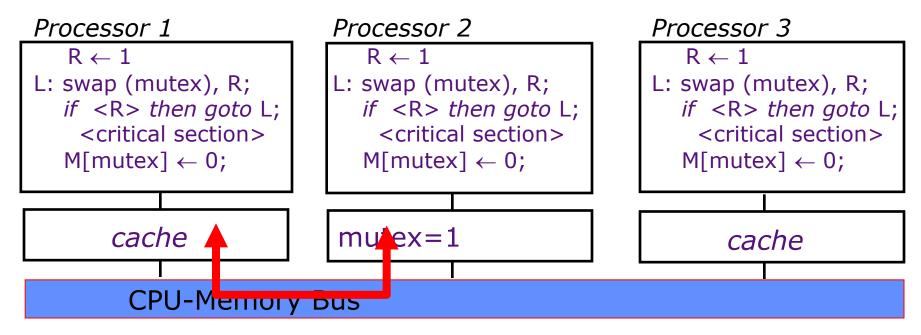




Sincronizarea și cache-urile:

Probleme de performanță





Protocoalele de coerență a cache-ului vor face mutex-ul să facă pingpong între cache-urile lui P1 și P2.

Acest fenomen poate fi redus prin citirea locației inițiale a mutex-ului (non-atomic) și execuția unui schimb doar dacă acesta are valoarea zero.





Performanța și traficul pe magistrale



În general, o instrucțiune read-modify-write necesită două operații pe magistrală, dacă nu avem alte operații la memorie de către alte procesoare

Într-un scenariu multiprocesor, accesul tuturor celorlalte procesoare la magistrală trebuie să fie blocat pe durata execuției operației atomice de read-modify-write

⇒ costisitor pentru magistrale simple ISA-urile moderne folosesc

load-reserve store-conditional





Load-reserve & Store-conditional Embedded Systems Laboratory

Registre speciale pentru stocarea flag-urilor de reserve, adresa și rezultatul store-conditional

```
Load-reserve R, (a):

<flag, adr> \leftarrow <1, a>;

R \leftarrow M[a];
```

```
Store-conditional (a), R:

if <flag, adr> == <1, a>

then cancel other procs'

reservation on a;

M[a] \leftarrow <R>;

status \leftarrow succeed;

else status \leftarrow fail;
```

Dacă un alt procesor vede o tranzacție de store la adresa din registrul de rezervare, bitul de rezervare este setat la 0

- Mai multe procesoare pot rezerva 'a' simultan
- Aceste instrucțiuni sunt similare cu load și store obișnuite, din pct de vedere al traficului pe bus





Performanță:

Load-reserve & Store-conditional



Numărul total de tranzacții pe bus nu este neapărat redus, dar spargerea unei instrucțiuni atomice în load-reserve & store-conditional:

- crește utilizarea magistralei, mai ales pentru magistralele care efectuează tranzacții în mai multe etape
- reduce efectul ping-pong pentru cache deoarece procesoarele care încearcă să obțină un semafor nu trebuie să facă un store de fiecare dată





Lățimea de bandă limitează performanțele



Avem un sistem multiprocesor cu memorie partajată construit în jurul unui singur bus cu lățimea de bandă de *x* GB/s. Cuvintele de instrucțiuni și date au fiecare 4B lățime, fiecare instrucțiune necesită accesul în medie la 1.4 cuvinte din memorie (inclusiv la instrucțiunea însăși). Rata combinată de hit a cache-urilor este de 98%. Calculați limita sperioară a performanței sistemului multiprocesor în GIPS. Liniile de adresă sunt separate și nu afectează lățimea de bandă de date.

Soluție

Execuția unei instrucțiuni implică un transfer de $1.4 \times 0.02 \times 4 = 0.112$ B. Astfel, limita absolută a perfomanței este de x/0.112 = 8.93x GIPS. Dacă presupunem o lățime a busului de 32 de biți, că nici un ciclu pe bus nu se irosește și o frecvență de ceas pe bus de y GHz, limita superioară a perfomanței devine 286y GIPS. Magistralele operează la frecvențe de 0.1 - 1 GHz. Prin urmare, o performanță apropiată de 1 TIPS (chiar și $\frac{1}{4}$ TIPS) este peste capabilitățile acestui tip de arhitectură.





Acknowledgements



- These slides contain material developed and copyright by:
 - Arvind (MIT)
 - Krste Asanovic (MIT/UCB)
 - Joel Emer (Intel/MIT)
 - James Hoe (CMU)
 - John Kubiatowicz (UCB)
 - David Patterson (UCB)
- MIT material derived from course 6.823
- UCB material derived from course CS252



