

## IV.9. Arhitectura RISC

# Structura clasică a CPU

CISC (*Complex Instruction Set Computer*)

- număr mare de instrucțiuni
- complexitate mare a instrucțiunilor → timp mare de execuție
- număr mic de regiștri → acces intensiv la memorie

## Observații practice

- multe instrucțiuni sunt rar folosite
- 20% din instrucțiuni sunt executate 80% din timp
  - sau 10% sunt executate 90% din timp
  - depinde de sursa de documentare...
- instrucțiunile complexe pot fi simulate prin instrucțiuni simple

# Structura alternativă

## RISC (*Reduced Instruction Set Computer*)

- set de instrucțiuni simplificat
  - instrucțiuni mai puține (relativ)
  - și mai simple funcțional (elementare)
- număr mare de regiștri (zeci)
- mai puține moduri de adresare a memoriei

# Accesele la memorie

- format fix al instrucțiunilor
  - același număr de octeți, chiar dacă nu toți sunt necesari în toate cazurile
  - mai simplu de decodificat
- arhitectură de tip *load/store*
  - accesul la memorie - doar prin instrucțiuni de transfer între memorie și regiștri
  - restul instrucțiunilor lucrează doar cu regiștri

# Avantajele structurii RISC

- instrucțiuni mai rapide
- reduce numărul de accese la memorie
  - depinde de capacitatea compilatoarelor de a folosi regiștrii
- accesele la memorie - mai simple
  - mai puține blocaje în pipeline
- necesar de siliciu mai mic - se pot integra circuite suplimentare (ex. cache)

## **IV.10. Arhitecturi paralele de calcul**

# Calcul paralel - utilizare

- comunicare între aplicații
  - poate fi folosită și procesarea concurentă
- performanțe superioare
  - calcule științifice
  - volume foarte mari de date
  - modelare/simulare
    - meteorologie, astronomie etc.



# Cum se obține paralelismul?

- structuri pipeline
  - secvențial/paralel
- sisteme multiprocesor
  - unitățile de calcul - procesoare
- sisteme distribuite
  - unitățile de calcul - calculatoare

# Performanța

- ideal - viteza variază liniar cu numărul de procesoare
- real - un program nu poate fi paralelizat în întregime
- exemple
  - operații de I/O
  - sortare-interclasare

## Scalabilitatea (1)

- creșterea performanței o dată cu numărul procesoarelor
- probleme - sisteme cu foarte multe procesoare
- factori de limitare
  - complexitatea conexiunilor
  - timpul pierdut pentru comunicare
  - natura secvențială a aplicațiilor

## Scalabilitatea (2)

- funcțională pentru un număr relativ mic de procesoare
- număr relativ mare
  - creșterea de performanță nu urmează creșterea numărului de procesoare
- număr foarte mare
  - performanța se plafonează sau poate scădea

# Sisteme de memorie

- după organizarea fizică a memoriei
  - centralizată
  - distribuită
- după tipul de acces la memorie
  - comună (partajată)
  - locală

# Tipuri de sisteme multiprocesor

- sisteme cu memorie partajată centralizată
- sisteme cu memorie partajată distribuită
- sisteme cu schimb de mesaje

# Memorie partajată centralizată

- denumiri
  - UMA (*Uniform Memory Access*)
  - SMP (*Symmetrical Multiprocessing*)
- memorie comună
- accesibilă tuturor procesoarelor
- timpul de acces la memorie
  - același pentru orice procesor și orice locație

# Memorie partajată distribuită

- denumiri
  - DSM (*Distributed Shared Memory*)
  - NUMA (*Non-Uniform Memory Access*)
- memoria este distribuită fizic
- spațiu de adrese unic
  - văzut de toate procesoarele
- accesul la memorie - neuniform



# Sisteme cu schimb de mesaje

- multicalculatoare
- fiecare procesor are propria memorie locală
  - nu este accesibilă celorlalte procesoare
- comunicarea între procesoare
  - transmiterea de mesaje explicite
  - similar rețelelor de calculatoare

# Comunicare

- conectarea
  - între procesoare
  - între procesoare și memorie
- variante
  - sisteme cu memorie partajată centralizată - magistrală
  - sisteme cu memorie partajată distribuită - rețele de interconectare

# Comunicare pe magistrală

- economic
- proiectare simplă
- performanță mai redusă
- nu scalează bine cu numărul de procesoare
- rol important - memoriile cache
  - probleme de coerență

# Rețele de interconectare (1)

- între
  - procesoare
  - procesoare și memorie
- scop
  - flexibilitate
  - performanță
    - mai multe accese în paralel

## Rețele de interconectare (2)

- tipuri de conectare
  - totală - fiecare cu fiecare
  - parțială - unele perechi de componente nu sunt conectate direct
- conectare procesoare-memorii
  - mai multe circuite de memorie
  - pot fi accesate în paralel de procesoare diferite

# Coerența memoriei

- trebuie ca toate procesoarele să folosească ultima valoare scrisă pentru o variabilă partajată
- problema - memoriile cache
- scopul - orice variabilă partajată să aibă aceeași valoare
  - în toate cache-urile
  - în memoria principală

# Coerența - cache *write-back*

x - variabilă partajată

Procesor	Acțiune	Valoare cache A	Valoare cache B	Valoare memorie
				9
A	i=x;	9		9
B	j=x;	9	9	9
A	x=5;	5	9	9
B	k=x;	5	9	9

# Coerența - *cache write-through*

x - variabilă partajată

Procesor	Acțiune	Valoare cache A	Valoare cache B	Valoare memorie
				9
A	i=x;	9		9
B	j=x;	9	9	9
A	x=5;	5	9	5
B	k=x;	5	9	5



# Ce înseamnă coerența? (1)

## 1. ordinea execuției

- a) procesorul P scrie în variabila X
- b) apoi procesorul P citește X
- între a) și b) nu există alte scrieri în X

→ citirea b) returnează valoarea scrisă de a)

## Ce înseamnă coerența? (2)

### 2. viziune coerentă a memoriei

- a) procesorul P scrie în variabila X
  - b) apoi procesorul Q ( $Q \neq P$ ) citește X
    - între a) și b) nu există alte scrieri în X
    - între a) și b) a trecut suficient timp
- citirea b) returnează valoarea scrisă de a)

## Ce înseamnă coerența? (3)

### 3. serializarea scrierilor

- a) procesorul P scrie în variabila X
- b) procesorul Q ( $Q=P$  sau  $Q \neq P$ ) scrie în variabila X

- toate procesoarele văd cele două scrieri în aceeași ordine
- nu neapărat a) înaintea b)

# Menținerea coerenței cache-urilor

- protocoale de menținere a coerenței
- se bazează pe informațiile despre liniile de cache
  - *invalid* - datele nu sunt valide
  - *dirty* - doar cache-ul curent deține valoarea actualizată
  - *shared* - cache-ul curent deține valoarea actualizată, la fel memoria principală și eventual alte cache-uri

# Tipuri de protocoale

- *directory based*
  - informațiile despre fiecare linie de cache - ținute într-un singur loc
- *snooping*
  - fiecare cache are o copie a liniei partajate
  - fără centralizarea informației
  - cache-urile monitorizează magistrala
    - detectează schimbările produse în liniile de cache

# Actualizarea cache-urilor

- fiecare cache anunță modificările făcute
- celelalte cache-uri reacționează
- contează doar operațiile de scriere
- variante
  - scriere cu invalidare (*write invalidate*)
  - scriere cu actualizare (*write update*)

## Scriere cu invalidare (1)

- un procesor modifică o dată
- modificarea se face în cache-ul propriu
  - toate celelalte cache-uri sunt notificate
- celelalte cache-uri
  - nu au o copie a datei modificate - nici o acțiune
  - au o copie a datei modificate - își invalidează linia corespunzătoare
  - valoarea corectă va fi preluată când va fi nevoie

## Scriere cu invalidare (2)

x - variabilă partajată

Proc.	Acțiune	Reacție cache	Cache A	Cache B	Memorie
					9
A	i=x;	read miss	9		9
B	j=x;	read miss	9	9	9
A	x=5;	invalidation	5	inv.	9
B	k=x;	read miss	5	5	5



## Scriere cu actualizare (1)

- un procesor modifică o dată
- modificarea se face în cache-ul propriu
  - toate celelalte cache-uri sunt notificate
  - se transmite noua valoare
- celelalte cache-uri
  - nu au o copie a datei modificate - nici o acțiune
  - au o copie a datei modificate - preiau noua valoare

## Scriere cu actualizare (2)

x - variabilă partajată

Proc.	Acțiune	Reacție cache	Cache A	Cache B	Memorie
					9
A	i=x;	read miss	9		9
B	j=x;	read miss	9	9	9
A	x=5;	invalidation	5	5	5
B	k=x;	read hit	5	5	5

# Invalidare vs. actualizare (1)

- mai multe scrieri succesive în aceeași locație
  - *write invalidation* - o singură invalidare (prima dată)
  - *write update* - câte o actualizare pentru fiecare scriere
  - mai avantajos - invalidare

## Invalidare vs. actualizare (2)

- mai multe scrieri în aceeași linie de cache
  - modificarea unei locații necesită invalidarea/actualizarea întregii linii
  - *write invalidation* - o singură invalidare (prima dată)
  - *write update* - câte o actualizare pentru fiecare scriere
  - mai avantajos - invalidare

## Invalidare vs. actualizare (3)

- "timpul de răspuns"
  - timpul între modificarea unei valori la un procesor și citirea noii valori la alt procesor
  - *write invalidation* - întâi invalidare, apoi citire (când este necesar)
  - *write update* - actualizare imediată
  - mai avantajos - actualizare

## Invalidare vs. actualizare (4)

- ambele variante au avantaje și dezavantaje
- scrierea cu invalidare - ocupare (mult) mai redusă a memoriei și magistralelor
- scrierea cu actualizare - rata de succes a cache-urilor mai mare
- mai des folosită - invalidarea

# V. Dispozitivele periferice

# Magistrale (1)

- căi de comunicație a informației
- o magistrală leagă între ele mai mult de 2 componente printr-o cale unică
- descrierea unei magistrale
  - semnalele electrice folosite
  - reguli de comunicare - trebuie respectate de toate părțile implicate
  - mod de conectare



## Magistrale (2)

### Accesul la magistrală

- mai multe entități pot solicita simultan accesul
- este necesară o procedură de arbitraj
  - decide cine va primi accesul
  - celelalte trebuie să aștepte eliberarea magistralei

# Arbitrarea magistrelor

## Tipuri de arbitrare a magistrelor

- centralizat
  - decizia o ia un circuit specializat (arbitru)
- descentralizat
  - componentele se înțeleg între ele
  - pe baza regulilor care definesc funcționarea magistralei

# Conectarea la magistrală

- probleme de natură electrică
- mai multe circuite conectate împreună
  - în intrare și ieșire
- nu se pot conecta mai multe ieșiri
  - nivelele diferite de tensiune ar distruge circuitele
- o soluție - multiplexarea
  - toate ieșirile sunt conectate la un multiplexor

## Circuite *tri-state*

- ieșirea are 3 stări posibile
  - 0
  - 1
  - impedanță infinită (*High-Z*)
- primele două corespund valorilor obișnuite
- a treia implică decuplarea de pe magistrală
  - ca și cum ieșirea circuitului nici nu ar fi conectată la magistrală

## Circuite *open-collector*

- în unele cazuri se numesc *open-drain*
  - în funcție de tehnologia utilizată
- este posibilă conectarea mai multor ieșiri simultan
- valoarea rezultată - funcția AND între ieșirile conectate

# Magistrale - privire generală

## Avantaje

- activitatea pe magistrală - ușor de controlat
- economic - structură relativ simplă

## Dezavantaj

- performanțe mai scăzute
  - doar 2 componente pot comunica la un moment dat