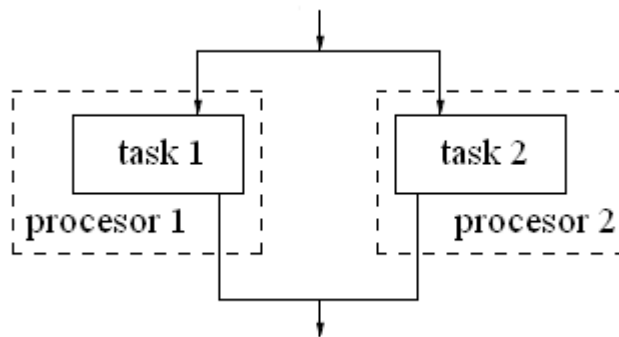


Curs 2 AAC

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

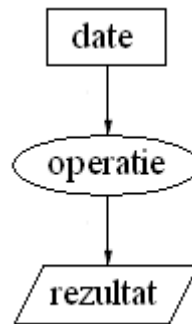
- Se elimină limitarea de viteză a procesoarelor secvențiale (timp de execuție mai scurt).
- Se elimină limita de miniaturizare a procesoarelor secvențiale (nu putem construi un procesor complex la dimensiuni oricât de mici).
- Se elimină limitarea economică (imposibilitatea de a fabrica ieftin un procesor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

Evoluția

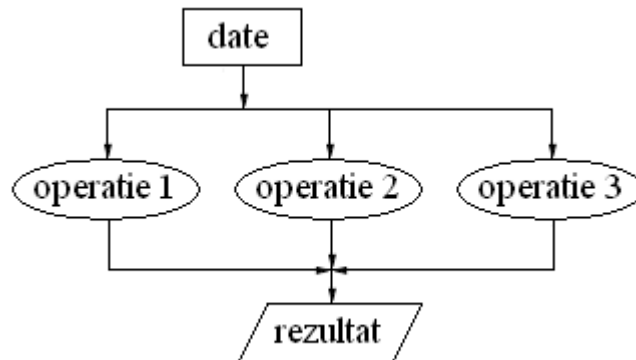
- Procesoare secvențiale obișnuite, numite procesoare scalare: execută maxim o instrucțiune într-un ciclu procesor.
- Procesoare superscalare: pot executa mai multe instrucțiuni într-un ciclu procesor.
- Procesoare vectoriale: nu mai lucrează pe operanzi scalari ci pe structuri de date multiple: vectori 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

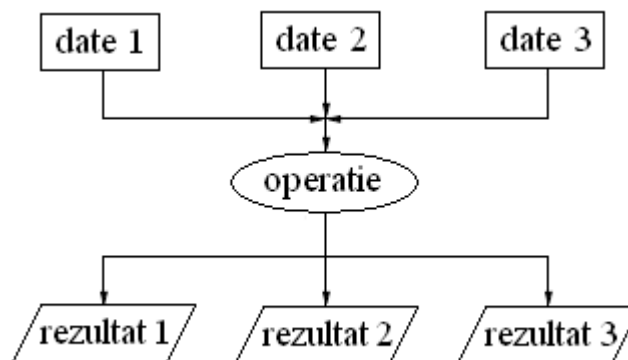
- **SISD** (*Single Instruction Single Data*): arhitectura secvențială obișnuită.



- **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.

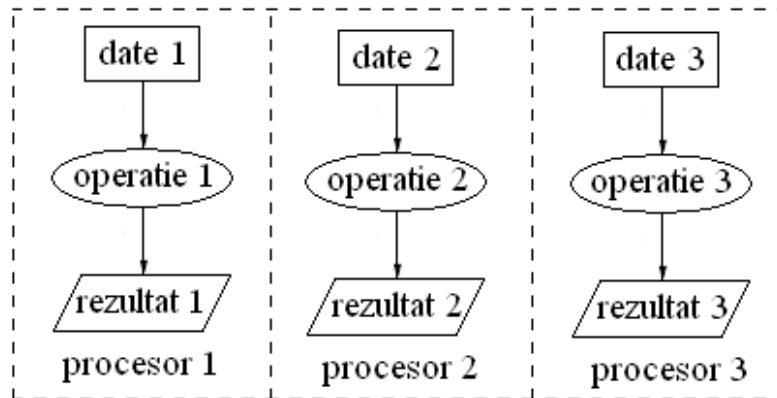


- **SIMD** (*Single Instruction Multiple Data*): un același program se execută pe mai multe seturi de date.



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)

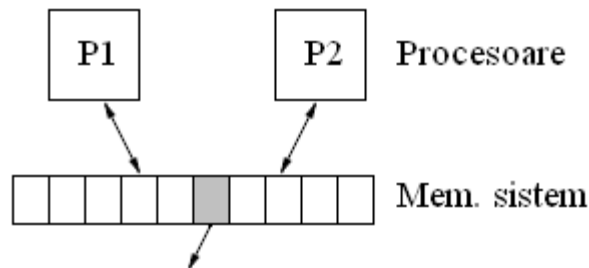
$$\Rightarrow Sp = \frac{T_{es}}{T_{secv} + T_{paralel}} \Rightarrow Sp < \frac{T_{es}}{T_{secv}} = K \quad \text{- Legea lui Amdahl}$$

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.

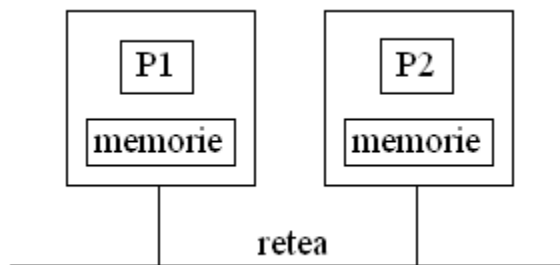


La un anumit moment, o anumită locație din memoria partajată poate fi accesată doar de către un singur procesor.

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.

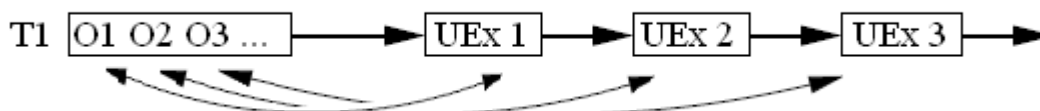


Procesoarele pot avea acces atât la o memorie locală aflată pe calculatorul din care fac parte, cât și la o memorie aflată la distanță (alt calculator).

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.



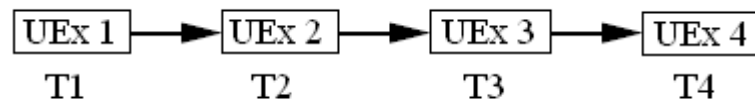
Există 2 tipuri de Pipeline: - de instrucțiuni; în care o instrucțiune se descompune într-o serie de microinstrucțiuni.

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

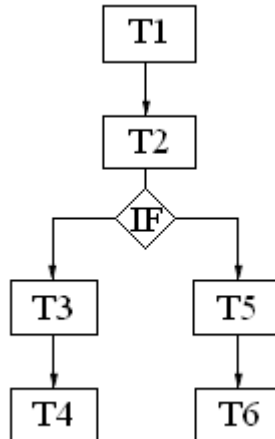
În rezolvarea unui task pot exista anumite situații speciale, în care fluxul de operații din Pipeline se oprește. Aceste situații speciale se numesc hazarduri (incidente) și sunt de 3 tipuri:

- Hazardul structural. El e generat de faptul că nu există destule unități de execuție pentru operațiile planificate.
- Hazardul de date. El e generat din cauza dependențelor de date. Un anumit task așteaptă după datele unui alt task.
- 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).

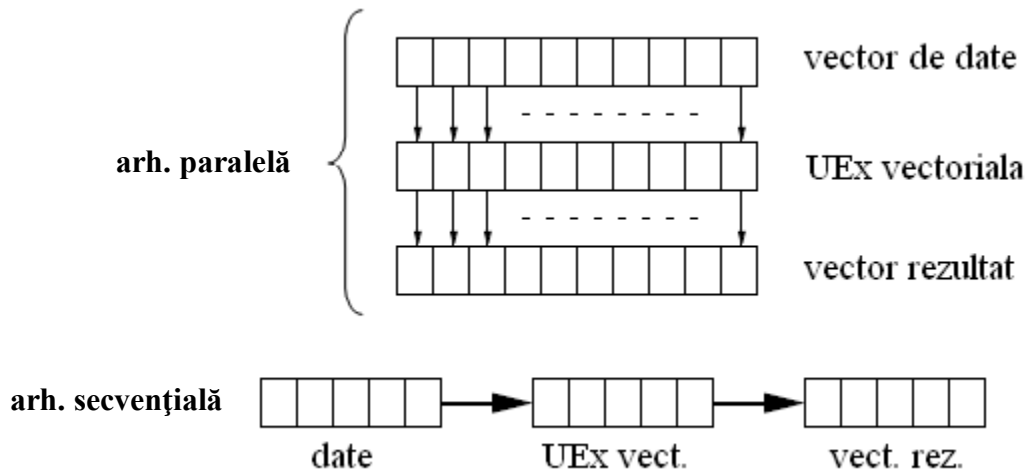


- 2) **Arhitecturi cu legături multiple** – acestea sunt structuri multiprocesor conectate în diverse configurații în care fiecare procesor realizează anumite operații simple într-un timp foarte scurt.

- a) **Arhitecturi vectoriale.**

Exemple:

- **Procesoare vectoriale**



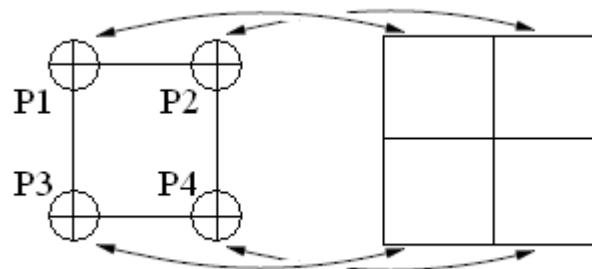
- **Procesoare sistolice:** la fiecare tact datele avansează o poziție prin structura sistolică.



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.**

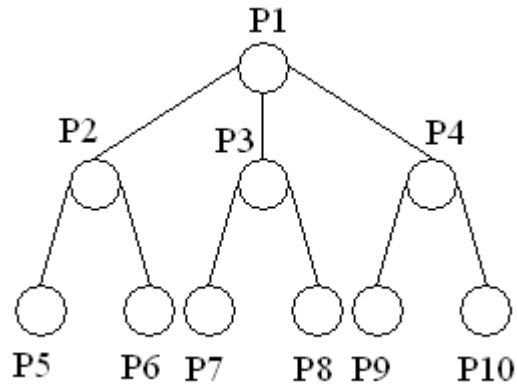
Exemplu: - procesoare matriciale



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.