

## Performance Analysis of Parallel Algorithms

Asist. Felician ALECU

Catedra de Informatică Economică, ASE București

*The main reason of parallelization a sequential program is to run the program faster. When running a parallel program on a real parallel system there is an overhead coming from processors load imbalance and from communication times needed for changing data between processors and for synchronization. This is the reason why the execution time of the program will be greater than the theoretical value.*

**Keywords:** parallel processing, execution time, parallel speedup, efficiency, costs, parallelization.

**T**impul de execuție al unui algoritm secvențial ( $T_s$ ), se exprimă ca funcție de dimensiunea datelor de intrare. Timpul de execuție paralelă ( $T_p$ ) reprezintă timpul scurs din momentul în care programul paralel este lansat în execuție până când ultimul procesor își încheie execuția. Acesta depinde însă nu numai de dimensiunea intrărilor ci și de arhitectura calculatorului paralel și de numărul de procesoare existente în sistem.

În general, timpii de execuție sunt exprimați cu ajutorul numărului de operații elementare ce sunt executate, în cazul cel mai defavorabil, pentru rezolvarea unei probleme de dimensiune dată.

Exprimarea timpului de execuție în funcție de dimensiunea problemei poate conduce la obținerea unor expresii complicate care sunt reduse la forme mai simple cu ajutorul cu ajutorul limitelor asimptotice ( $O$  - limita asimptotică superioară,  $\Omega$  - limita asimptotică inferioară,  $\Theta$  - ordinul exact).

Indicatorii de performanță asociați sistemelor paralele sunt accelerarea, eficiența și costul paralel.

**Accelerarea paralelă** se notează cu  $S$  și reprezintă câștigul de viteză ce se obține datorită execuției algoritmului pe un sistem paralel

$$T_p = T_s \cdot \alpha + T_s \cdot (1 - \alpha) / p$$

$$S = \frac{T_s}{T_p} = \frac{T_s}{T_s \cdot \alpha + T_s \cdot (1 - \alpha) / p} = \frac{1}{\alpha + (1 - \alpha) / p} = \frac{p}{\alpha \cdot (p - 1) + 1}$$

Atunci când  $p \rightarrow \infty$  avem:  $\lim_{p \rightarrow \infty} S = \frac{1}{\alpha}$

Din acest motiv, accelerarea maximă care se poate obține atunci când o fracție  $\alpha$  a pro-

în comparație cu timpul de execuție în varianta secvențială. Acest indicator se exprimă ca raport între timpul de execuție în cazul cel mai defavorabil al celui mai rapid algoritm secvențial care rezolvă o anumită problemă și timpul de execuție necesar unui algoritm paralel pentru a soluționa aceeași problemă pe un sistem cu  $p$  procesoare identice cu cel al calculatorului secvențial:  $S = \frac{T_s}{T_p}$

Teoretic, accelerarea paralelă nu poate avea o valoare mai mare ca  $p$ , numărul de procesoare. O accelerare paralelă egală cu numărul de procesoare din sistem poartă numele de accelerare liniară.

În conformitate cu legea lui Amdahl, chiar și într-un sistem paralel ideal este foarte dificil de obținut o accelerare paralelă egală cu numărul de procesoare datorită faptului că în cadrul oricărui program există o fracție  $\alpha$  care nu poate fi paralelizată și care trebuie executată secvențial. Restul de  $(1 - \alpha)$  pași de calcul se pot executa în paralel pe procesoarele disponibile în sistem. Astfel, timpul paralel de execuție și accelerarea paralelă devin:

gramului nu poate fi paralelizată este  $1/\alpha$  indiferent de numărul de procesoare din sistem. De exemplu, dacă în cazul unui sistem cu 4 procesoare o fracție de 20% dintr-un program

nu poate fi paralelizată, timpul paralel de execuție și accelerarea paralelă vor fi:

$$T_p = T_s \cdot 0.2 + T_s \cdot 0.8 / 4 = 0.4 \cdot T_s$$

$$S = \frac{T_s}{0.4 \cdot T_s} = \frac{1}{0.4} = 2.5$$

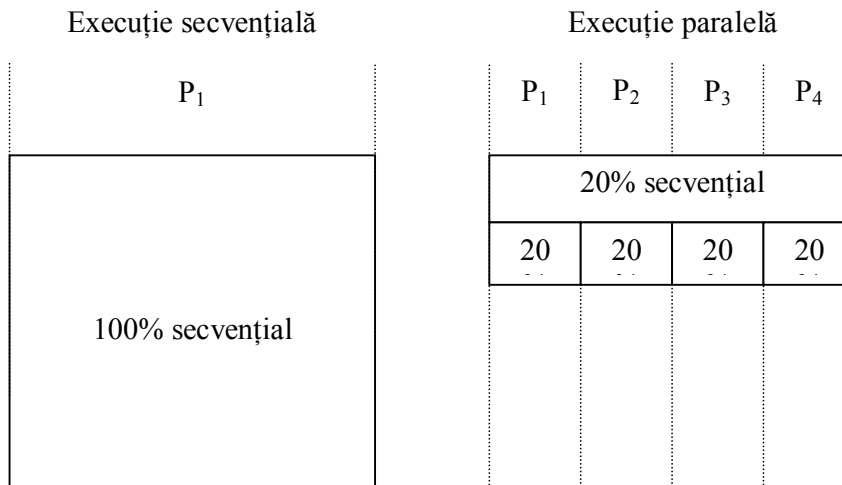


Fig. 1. Execuția unui program în cadrul unui sistem ideal

Accelerarea maximă care poate fi obținută este  $1/0.2 = 5$  ceea ce semnifică faptul că timpul paralel de execuție nu va fi niciodată mai mic de 20% din cel secvențial, indiferent de numărul de procesoare din sistem.

Legea lui Amdahl exprimă în mod clar necesitatea minimizării fracției  $\alpha$  ce nu poate fi paralelizată prin stabilirea unei limite superioare a accelerării paralele

Atunci când programul paralel este executat pe un sistem real, apar o serie de costuri suplimentare generate de comunicația dintre procesoare (schimb de date, sincronizare) cât și de încărcarea neechilibrată a procesoarelor. Din acest motiv timpul paralel de execuție poate depăși valoarea teoretică calculată (Fig. 2).

**Eficiența paralelă** se definește ca fiind egală cu raportul dintre accelerarea paralelă și numărul de procesoare utilizate:  $E = \frac{S}{p}$

Datorită faptului că accelerarea are o valoare limitată la numărul de procesoare, eficiența paralelă poate lua valori numai între 0 și 1.

Eficiența paralelă cuantifică măsura în care procesoarele din sistem au efectuat operații utile de calcul de-a lungul execuției programului paralel.

Timpul paralel de execuție va fi 40% din cel secvențial ceea ce înseamnă că programul paralel va fi doar de 2.5 ori mai rapid decât cel secvențial (figura 1).

**Costul paralel** înglobează timpii consumați de toate procesoarele sistemului paralel pentru obținerea soluției problemei:  $C_p = p \cdot T_p$

În general prin cost înțelegem suma tuturor timpilor consumați pentru rezolvarea unei probleme. Din acest motiv costul secvențial va fi întotdeauna egal cu timpul secvențial de execuție deoarece acesta reprezintă timpul total consumat pentru rularea programului:  $C_s = T_s$

Eficiența paralelă mai poate fi exprimată cu ajutorul raportului dintre costul secvențial și cel paralel:

$$E = \frac{S}{p} = \frac{T_s / T_p}{p} = \frac{T_s}{p \cdot T_p} = \frac{C_s}{C_p}$$

Un algoritm paralel se consideră a fi optimal din punct de vedere al costului dacă costul paralel al acestuia este proporțional cu timpul de execuție al celui mai rapid algoritm serial care rezolvă respectiva problemă:  $C_p = k \cdot T_s$

Folosind relația anterioară, eficiența paralelă devine:

$$E = \frac{C_s}{C_p} = \frac{T_s}{p \cdot T_p} = \frac{T_s}{k \cdot T_s} = \frac{1}{k} = O(1)$$

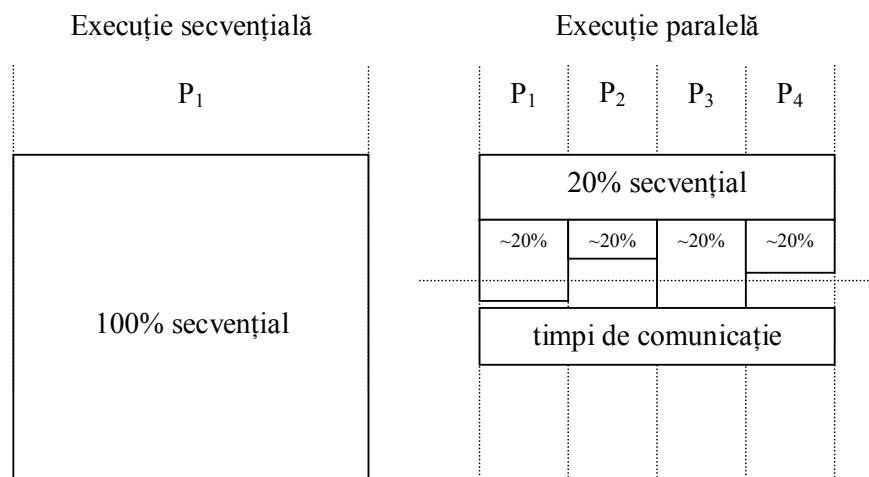
Rezultă că un algoritm cost-optimal are o eficiență de ordinul  $O(1)$ , constantă în raport cu

dimensiunea problemei și cu numărul de procesoare utilizate pentru rezolvarea acesteia.

**Costul suplimentar de calcul paralel** reprezintă suma timpilor de inactivitate și a celor consumați de către procesoarele sistemului paralel pentru operațiile de comunicație și sincronizare, operații fără de care obținerea rezultatelor finale ale programului nu ar fi posibilă. Acest cost nu există pentru varianta secvențială a algoritmului deoarece în sistem

există un singur procesor.

În figura 2 se prezintă modul în care se va executa un program în cadrul unui sistem paralel real format din patru procesoare, în condițiile în care 20% din program nu se poate paraleliza, încărcarea procesoarelor nu este echilibrată iar timpii de comunicație nu mai sunt neglijați. Se observă faptul că, față de cazul ideal, timpul paralel de execuție este mai mare datorită influenței celor doi factori anterior menționați.



**Fig. 2.** Execuția unui program în cadrul unui sistem real

Costul suplimentar de calcul paralel se poate exprima cu ajutorul următoarei formule:

$$C_{supl} = C_p - C_s = p \cdot T_p - T_s$$

Pentru ducerea la bun sfârșit a unei sarcini de calcul este necesar ca procesoarele sistemului să comunice între ele. Transferul de date între două procesoare ale sistemului este o operațiune mare consumatoare de timp care poate genera inactivitate la nivelul elementelor de procesare. Din acest motiv, algoritmi paraleli ar trebui să fie proiectați astfel încât majoritatea acceselor la memorie să nu necesite transportul datelor solicitate prin rețeaua de interconectare.

Existența în cadrul unui program a unui punct de sincronizare între mai multe procesoare va genera timp de așteptare deoarece execuția programului va continua numai după ce toate procesoarele au atins punctul respectiv.

Principală cauză pentru care la nivelul elementelor de procesare se înregistrează timp de inactivitate o constituie încărcarea neechilibrată a procesoarelor care survine atunci

când procesoarele din sistem au asociate sarcini de calcul de dimensiuni inegale. Într-un astfel de caz, unele procesoare își vor încheia execuția mai devreme și vor rămâne în așteptare în timp ce altele vor fi încă în faza de execuție.

O altă sursă care generează timp de inactivitate și implicit o încărcare neechilibrată a procesoarelor o reprezintă variația gradului de paralelism pe parcursul execuției programului, unde prin gradul de paralelism al unui program paralel înțelegem numărul de procesoare care sunt utilizate pentru execuția acestuia.

Chiar și prezența în cadrul unui program a unei părți ce nu poate fi paralelizată atrage după sine o astfel de încărcare neechilibrată a procesoarelor deoarece fracția respectivă va trebui să fie executată secvențial, de către un singur procesor. O astfel de situație este exprimată de către legea lui Amdahl.

Pentru ca un program paralel să fie executat cât mai rapid este necesar să se reducă cât mai mult posibil fracția de program ce nu se

poate paraleliza și în plus să se asigure o distribuție cât mai uniformă a sarcinii de calcul pe elementele de procesare în condițiile în care se are în vedere minimizarea timpilor de comunicație

Pentru obținerea unor algoritmi paraleli cât mai eficienți este necesar ca aceștia să îndeplinească o serie de condiții:

- ▶ să fie proiectați astfel încât să folosească un număr de procesoare mai mic de decât dimensiunea problemei, mai ales în situațiile în care este necesar să se prelucreze un volum considerabil de date; mai mult, algoritmul nu trebuie să fie dependent de o anumită configurație a calculatorului paralel pe care va fi executat;

- ▶ să se obțină un timp paralel de execuție cât mai mic și considerabil mai bun decât cel corespunzător celui mai rapid algoritm secvențial care rezolvă aceeași problemă; în general ne așteptăm ca timpul paralel de execuție să scadă în situația în care numărul de procesoare disponibile în sistem crește;

- ▶ costul paralel trebuie să fie optimal.

### Bibliografie

[Lad04] S. Ladd, *Guide to Parallel Programming*, Springer-Verlag, 2004

[Kar03] G. Karniadakis, R. M. Kirby, *Parallel Scientific Computing in C++ and MPI : A Seamless Approach to Parallel Algorithms and their Implementation*, Cambridge University Press, 2003

[Gra03] A. Grama et al, *An Introduction to Parallel Computing: Design and Analysis of Algorithms*, Addison Wesley, 2<sup>nd</sup> edition, 2003

[Dod02] Gh. Dodescu, B. Oancea, M. Raceanu, *Procesare paralelă*, Editura Economica, București, 2002

[Gro02] W. Gropp et al, *The Sourcebook of Parallel Computing*, Morgan Kaufmann, 2002