

Planificarea programelor paralele

Prep. Felician ALECU

Catedra de Informatica Economica, A.S.E. Bucuresti

In this paper we will focus on task allocation methods used to improve performances of parallel programs executed on parallel systems. The parallel program has to be divided into several independent tasks that will be executed concurrently. These tasks have to be assigned to the system processing units using a static or dynamic fault tolerant task-scheduling algorithm. Using the static approach, the schedule is computed at compile time. Dynamic scheduling means computing at runtime either by a central master or in a distributed manner.

Keywords: *Parallel program, task allocation, task scheduling problem, static task scheduling, dynamic task scheduling, machine configuration graph, program graph, scheduling algorithms and fault tolerance.*

Introducere

În ultimii ani sistemele paralele și distribuite au devenit din ce în ce mai atractive pentru aplicații cu cerințe de calcul intensive cum ar fi cele destinate simulării unor sisteme complexe (aerodinamica, meteorologie). Principalul avantaj al unor astfel de sisteme este raportul mai mult decât atractiv dintre pret și performanțele care pot fi obținute.

Prin planificarea unui program se înțelege alocarea taskurilor componente pe procesoarele existente în sistem în funcție de relațiile de dependență dintre acestea.

Planificarea poate fi realizată static sau dinamic. În planificarea statică, se cunoaște, încă din faza de compilare, procesorul pe care va fi rulat fiecare task component al programului. Datorită faptului că problema planificării este o problemă NP-completă, soluții exacte pot fi obținute doar pentru probleme de dimensiuni reduse (vezi [Cof72], în celelalte cazuri fiind necesară aplicarea unor algoritmi euristici. În toate cazurile însă este necesară cunoașterea în avans a timpilor de execuție a taskurilor. Pentru marea majoritate a aplicațiilor practice, estimarea, la momentul compilării programului, a timpilor de execuție este aproape imposibilă.

În planificarea dinamică alocarea taskurilor pe procesoarele sistemului se face la momentul execuției. Sarcina planificării poate reveni unui singur procesor (numit master) sau poate fi distribuită la nivelul unităților de procesare prezente în sistem.

Planificarea este strâns legată de modalitățile de divizare a programelor paralele în unități de execuție. Pentru ca aceste unități de execuție să poată fi executate concurent este necesar ca ele să fie independente. Independența a două taskuri înseamnă obținerea aceluși rezultat indiferent dacă cele două taskuri sunt executate secvențial în orice ordine sau în paralel. Pe lângă restricțiile legate de ordinea de execuție, conflictele de acces la resursele partajate limitează independența între două sau mai multe taskuri. Pentru a analiza dependențele dintre taskurile unui program, uzual se folosește un graf direcționat numit graf de dependență al programului. Nodurile grafului sunt unitățile de execuție ale programului iar arcele reprezintă relațiile de dependență între acestea. Din analiza grafului de dependență reiese în mod clar posibilitățile de paralelizare a programului respectiv.

Cele mai multe limbaje paralele de programare pun la dispoziția utilizatorilor instrumente ce permit definirea taskurilor componente (procese, threaduri) și a mecanismelor de comunicație (transfer de mesaje, variabile partajate) și sincronizare (semafoare, lockuri, bariere, monitoare, evenimente) dintre acestea.

Există o serie de limbaje de programare, numite limbaje de programare cu paralelizare implicită, care, plecând de la descrierea secvențială a algoritmului, realizează în mod automat paralelizarea acestuia. Rolul principal îl are compilatorul care analizează codul sur-

sa al programului si cauta sa identifice segmente de cod care sa fie executate în paralel. Aceste segmente de cod iau forma proceselor si threadurilor. Tot compilatorul este cel care specifica modul de comunicare si sincronizare dintre elementele de executie identificate. Cu toate acestea, misiunea identificarii segmentelor independente din cadrul unui program nu este una facila si este limitata de complexitatea programului. Cu cât acesta are o structura mai neregulata, cu multiple ramificatii, bucle, apeluri, cu atât sansele de paralelizare automata scad.

Formularea problemei de planificare

Sistemul paralel poate fi descris cu ajutorul unui graf neorientat numit graful de configuratie al sistemului. Nodurile grafului reprezinta procesoarele sistemului împreuna cu memoria asociata. În dreptul unui nod se poate trece viteza procesorului si dimensiunea memoriei asociata acestuia. Arcele semnifica legaturile de comunicare (conexiunile). În dreptul unui nod se poate trece capacitatea legaturii dintre cele doua noduri.

În marea majoritate a cazurilor se considera ca un procesor poate comunica cu oricare altul din sistem. Din acest motiv graful de configuratie al sistemului este un graf complet, între oricare doua noduri existând o muchie care le uneste.

Comunicatia între procesoare se realizeaza asincron si este gestionata de catre un subsistem de comunicare care are si sarcina de a trata eventualele erori de comunicare care pot aparea. Chiar daca legatura directa dintre doua noduri a picat, acestea se considera totusi conectate daca exista cel puțin un drum în graf care sa le uneasca.

Erorile ce pot aparea la nivelul elementelor de procesare sunt gestionate de catre algoritmul de planificare. Atunci când un procesor devine izolat (nu mai exista conexiune cu nici un alt nod), se transmite un mesaj de informare tuturor nodurilor conectate din sistem. Erorile la nivelul elementelor de procesare pot fi tranziente sau permanente. În cazul unei erori tranziente, elementul de procesare poate participa din nou la executia pro-

gramului dupa ce eroarea înregistrata a fost remediata.

Programul paralel ce se doreste a fi executat poate fi si el descris cu ajutorul unui graf orientat numit graful programului. Nodurile grafului pot fi de doua feluri: taskuri (reprezentate sub forma de cercuri) si date (reprezentate sub forma de patrate). Un nod în care se gaseste un task poate fi etichetat cu durata de executie a acestuia (costul de calcul). Datele de iesire ale unui task ce sunt utilizate ca date de intrare de catre un alt task apar ca nod al grafului. Nodurile semnifica fluxurile de date din sistem generate de executia programului si introduc constrângeri de precedenta în executia taskurilor. Astfel, un task nu poate fi executat pâna când toate datele de care are nevoie sunt disponibile. Datele necesare taskului devin disponibile atunci când nodurile parinte au încheiat executia si pot comunica datele rezultate. Exemplu de graf al programului este prezentat în figura 1.

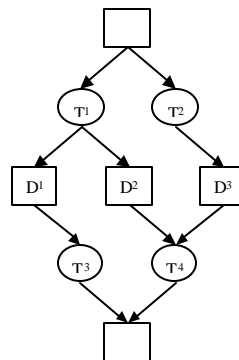


Fig. 1. Exemplu de graf al programului

Metode de planificare

În cazul în care graful programului este complet cunoscut înainte de momentul executiei programului (taskuri componente, timpi de executie, flux de date), se poate utiliza o planificare statica a taskurilor. În schimb, daca forma grafului programului nu poate fi complet determinata decât la momentul executiei, atunci se va folosi planificarea dinamica.

Datorita faptului ca se realizeaza la momentul compilarii programului (deci înainte de executia efectiva a acestuia), planificarea statica nu este consumatoare de timpi de executie. Planificarea dinamica implica executia

algoritmului de planificare concomitent cu rularea programului paralel.

Scopul principal al planificării este reducerea timpului de execuție paralela a programului. Din acest motiv, timpul consumat pentru fundamentarea deciziei de planificare trebuie să fie mai mic decât timpul câștigat prin aplicarea respectivei planificări. Nu se poate afirma că un anumit tip de planificare este cel mai bun deoarece alegerea unei metode de planificare pentru un anumit program se face în funcție de caracteristicile acestuia. Planificarea statică este atractivă datorită faptului că se execută o singură dată și nu este consumatoare de timp de execuție. Se poate afirma însă cu certitudine faptul că planificarea dinamică oferă posibilitatea tratării într-o manieră eficientă a eventualelor erori permanente sau temporare aparute la nivelul unuia sau mai multor elemente de procesare din sistem. Planificarea se poate reprezenta grafic cu ajutorul diagramei Gantt. Dintr-o astfel de diagramă reiese în mod clar procesorul pe care este executat un anumit task precum și timpul de execuție al acestuia.

Clusterizarea este o metodă frecvent folosită în cazul planificării statice. Ideea ce stă la baza acestei metode este aceea conform căreia costul de comunicație a două taskuri este diferit de zero doar dacă taskurile sunt planificate pe procesoare diferite. Prin această metodă se încearcă reducerea numărului taskurilor prin gruparea acestora. Grupele astfel obținute poartă numele de cluster și vor fi planificate pe același procesor. Astfel, dacă pentru două taskuri avem o valoare mare a costului de comunicație, atunci acestea vor fi incluse în cadrul aceleiași grupe fiind planificate pe același procesor. Planificarea pe același procesor a celor două taskuri anulează costul de comunicație. După definirea clusterelor, acestea vor fi planificate pe procesoarele sistemului în mod static, la momentul compilării programului.

În cazul planificării dinamice, sarcina planificării poate reveni unui singur procesor sau poate fi distribuită elementelor de procesare din sistem.

Prin folosirea unui singur procesor (numit procesor *master*), acesta va fi responsabil

pentru alocarea taskurilor pe procesoarele sistemului datorită faptului că pe respectivul procesor va fi rulat algoritmul de planificare.

Folosirea unei metode descentralizate de planificare presupune delegarea sarcinii planificării către nodurile care au terminat execuția taskurilor alocate [Liy95] [Yue97]. Să considerăm cazul în care taskul T_1 (figura 1) și-a încheiat execuția. Procesorul care a executat taskul respectiv devine liber și verifică dacă sunt taskuri a căror execuție depinde de rezultatele generate de T_1 . Astfel de taskuri vor fi numite succesori. Un task succesori este pregătit pentru execuție atunci când îi sunt disponibile toate datele de intrare. În cazul grafului din figura 1, taskurile T_3 și T_4 pot deveni succesori gata de execuție. După identificarea succesorilor gata de execuție, procesorul încearcă să găsească procesoare libere în sistem cărora să le fie alocate taskurile respective (unul dintre succesori poate fi executat chiar pe procesorul respectiv). După executarea succesorilor, procesoarele cărora le-au fost alocate taskurile respective preiau sarcina planificării noilor succesori. Acest procedeu se continuă până în momentul în care toate taskurile componente ale programului au fost executate și s-a obținut setul final de date. În cazul în care un task are mai mulți predecesori (cazul taskului T_4 din figura 1), trebuie să existe implementată o politică cu ajutorul căreia să se poată stabili în mod clar care dintre predecesori este responsabil pentru alocarea respectivului task.

Tratarea erorilor într-o manieră eficientă poate fi realizată cu succes în planificarea dinamică deoarece algoritmul de planificare, fiind rulat la momentul execuției programului paralel, se poate adapta la evenimentele care pot apărea în sistem.

În cazul în care un element de procesare devine indisponibil, algoritmul de planificare va trebui să se limiteze la folosirea exclusivă a procesoarelor valide [Alv98]. Dacă planificarea se realizează la nivel centralizat, atunci procesorul master nu va mai alocă taskuri procesorului indisponibil, până în momentul în care acesta va fi raportat ca fiind din nou valid de către rutinele software specializate în diagnosticarea elementelor hardware din sistem. În cazul în care

sistem. În cazul în care planificarea se face la nivel local, atunci este suficient ca toate procesoarele valide din sistem să fie informate, prin intermediul unui mesaj distribuit în cadrul rețelei de comunicație, asupra faptului că un procesor a devenit indisponibil. În acest fel el nu va mai fi folosit în cadrul planificării până în momentul în care va deveni din nou disponibil, moment la care procesoarele sistemului vor fi informate de acest lucru tot prin intermediul unui mesaj.

O situație aparte este aceea în care taskul este lansat din nou în execuție pe procesorul în memoria căruia este stocată copia datelor de intrare. Într-un astfel de caz se realizează o nouă copie a datelor respective iar această copie va fi stocată în memoria locală a unui alt element de procesare. Această nouă copie va putea relansa taskul în cazul în care execuția acestuia ar esua din nou. După ce taskul a fost executat cu succes, datele lui de intrare nu mai sunt necesare. Din acest motiv datele originale cât și eventualele copii ale acestora sunt eliminate din memoriile locale ale elementelor de procesare din sistem. Datele de ieșire reprezintă date de intrare pentru succesorii taskului și din acest motiv acestea vor fi la rândul lor duplicate.

Problema planificării, în forma ei cea mai generală, pleacă de la existența unui graf al programului cu n noduri și a unui sistem de calcul paralel cu p procesoare. Soluția optimă se află într-un spațiu de p^n planificări posibile.

Pentru a măsura calitatea și eficiența diferitelor algoritmi și metode de planificare ar trebui să se compare performanțele obținute cu cele ale planificării optime. Planificarea optimă este aceea care asigură un timp minim de execuție a programului paralel.

Atunci când soluția optimă nu poate fi atinsă, se încearcă rezolvarea planificării într-o manieră eficientă prin utilizarea unor algoritmi euristici de planificare care să conducă la găsirea unor soluții sub-optimale în timp de execuție rezonabili.

Pentru marea majoritate a problemelor de planificare însă, soluția optimă nu poate fi obținută. Din acest motiv, concluzii asupra eficienței unui algoritm de planificare se pot trage doar comparând performanțele acestuia cu performanțele altor algoritmi de planificare cunoscuți în condițiile în care testele se efectuează asupra aceluși seturi de date [Kwo98].

Bibliografie

[Alv98] G.A. Alverson, W.G. Griswold, C. Lin, D. Notkin, D.P. Agrawal, *Optimal Scheduling Algorithm for Distributed-Memory Machines*, IEEE Transactions on Parallel and Distributed Systems, Vol. 9, Nr. 1, January 1998

[Cof72] E. Coffman, R. Graham, *Optimal Scheduling on Two Processor System*, Acta Informatica, Vol. 1, 1972

[Kwo98] Y.K. Kwok, I. Ahmad, *Benchmarking the Task Graph Scheduling Algorithm*, IEEE Transactions on Parallel and Distributed Systems, Vol. 9, Nr. 3, March 1998

[Liy95] Y. Li, F.J. Markus, J. Rost, *Agent Scheduling – A Model for Dynamic Task Scheduling*, EURO-PAR, Lectures Notes in Computer Science, Springer-Verlag, 1995

[Yue97] K.K. Yue, D.J. Lilja, *An Efficient Processor Allocation Strategy for Multiprogrammed Shared-Memory Multiprocessors*, IEEE Transactions on Parallel and Distributed Systems, Vol. 8, Nr. 12, December 1997