

ALGORITMI DE CALCUL PARALEL

BOGDAN DUMITRESCU

septembrie 2001

Cuprins

1	Arhitecturi paralele	7
1.1	O clasificare generală	8
1.1.1	Arhitecturi SIMD	9
1.1.2	Arhitecturi MIMD	11
1.2	Arhitecturi MIMD cu memorie distribuită	15
1.2.1	Caracteristici ale nodurilor	15
1.2.2	Câteva noțiuni din teoria grafurilor	16
1.2.3	Topologii curente	17
1.2.4	Scufundări	22
1.3	Criterii de performanță	25
2	Modele de programare	29
2.1	Descrierea algoritmilor SPMD	29
2.2	Modelul de comunicație prin mesaje	33
2.3	Standardul MPI	39
2.4	Modelul PRAM	44
3	Algoritmi de comunicație	47
3.1	Modele de comunicație	48
3.2	Operații de comunicație globală	52
3.3	Comunicații globale	54
3.3.1	Comunicație de la un procesor la altul	55
3.3.2	Difuzare	63
3.3.3	Difuzare generală	70
3.3.4	Difuzare personalizată (distribuție)	74
3.3.5	Schimb complet	78
3.3.6	Câteva concluzii	80
3.4	Comunicații în modelul comutare de circuit	82
3.5	Dirijare (routing)	87

4	Algoritmi paraleli: concepție și apreciere	91
4.1	Performanțele unui algoritm paralel	92
4.1.1	Criterii de bază	92
4.1.2	Despre overhead	94
4.1.3	Legea lui Amdahl	95
4.1.4	Alte criterii	97
4.2	Tehnici generale de paralelizare	99
4.2.1	Partiționarea în taskuri	99
4.2.2	Grafuri de precedență	100
4.2.3	Planificarea taskurilor	102
4.3	Câțiva algoritmi paraleli fundamentali	108
4.3.1	Suma a n numere	108
4.3.2	Suma globală	113
4.3.3	Suma prefixelor	114
4.3.4	Recurențe de ordinul I	119
4.3.5	Produsul matrice-vector	121
5	Sortare și probleme conexe	129
5.1	Găsirea minimului	131
5.2	Selecția și găsirea rangului	133
5.2.1	Rangul într-o secvență sortată	133
5.2.2	Selecție pe PRAM	134
5.2.3	Selecție pe MIMD cu memorie distribuită	138
5.3	Interclasarea	140
5.3.1	Complexitate paralelă teoretică	141
5.3.2	Un algoritm pe CREW PRAM, pentru număr mic de procesoare	142
5.3.3	Mediana secvenței $A \perp B$ (partiționare mediană)	143
5.3.4	Interclasare prin multipartiționare	147
5.3.5	Interclasare pe arhitecturi cu memorie distribuită	150
5.4	Sortare	153
5.4.1	Sortare par-impar	154
5.4.2	Quicksort paralel	159
5.4.3	Sortare echilibrată	165
5.4.4	Sortare bitonică	169
5.4.5	Sortare prin interclasare	175
5.4.6	Ce algoritm de sortare alegem ?	176
6	Algoritmi numerici	181
6.1	Înmulțirea de matrice	182
6.1.1	Complexitatea problemei	183
6.1.2	Algoritmi pe inel	184
6.1.3	Algoritmi pe tor	186
6.1.4	Un algoritm pentru memorie partajată	190

6.2	Sisteme liniare triunghiulare	192
6.2.1	Algoritmi secvențiali	192
6.2.2	Complexitatea paralelă	193
6.2.3	Algoritmi cu comunicație globală	195
6.2.4	Algoritmi în front de undă	200
6.2.5	Algoritmi ciclici	205
6.2.6	Comparații ale algoritmilor precedenți	208
6.2.7	Un algoritm pe tor	208
	Răspunsuri la probleme	213

Capitolul 1

Arhitecturi paralele

Deși această carte este dedicată algoritmilor paraleli, primul capitol se ocupă de arhitecturi. Motivul este foarte simplu. Așa cum între programatorii care rezolvă o problemă utilizând un același calculator secvențial și un același algoritm apar mari diferențe în eficacitatea programelor scrise—datorită modului în care fiecare știe să folosească particularități ale hardware-ului, ale sistemului de operare sau ale limbajului de programare folosit—departajarea este și mai netă în cazul utilizării calculatoarelor paralele; cauza principală este în primul rând neadecvarea algoritmilor la arhitectură. Dacă a face un program să "meargă" pe o arhitectură paralelă nu este, până la urmă, o dificultate prea mare, a-l face eficient este cu siguranță mult mai greu decât în cazul secvențial. De aceea, un prim pas în abordarea paralelismului, chiar pentru programatorul obișnuit (cel care rezolvă probleme generale) este trecerea în revistă a tipurilor de arhitecturi paralele. Nu vom intra în amănunte, ci ne vom mulțumi cu o descriere generală, suficientă însă pentru a releva implicațiile asupra algoritmilor, deci asupra programării; pentru o informare mai completă vezi, de exemplu, cartea scrisă de Hockney și Jesshope [15], apărută și în limba română; de asemenea, Internetul conține vaste resurse pentru aflarea celor mai noi tendințe în domeniu; un bun punct de start este de exemplu <http://www.top500.org>.

Tipuri de paralelism. La ora actuală, practic orice calculator, inclusiv micul PC pe care-l aveți aproape zilnic în față, conține elemente de paralelism. Ele sunt încă suficient de bine ascunse pentru a nu obliga utilizatorul să-și schimbe semnificativ felul de a gândi. Apariția lor a fost determinată de necesitatea creșterii vitezei de calcul, prin alte mijloace decât creșterea frecvenței ceasului (permisă de permanente inovații tehnologice în conceperea și producerea circuitelor electronice). Cel mai evident exemplu este cel al coprocesorului pentru PC-uri; este, în fond, un alt procesor specializat într-un singur gen de calcule, cele în virgulă mobilă; el lucrează în paralel cu procesorul, dar la comanda acestuia; dacă acum 10 ani coprocesorul era un circuit separat, acum el este parte a procesorului. Paralelismul acesta se numește *functional*; există unități de calcul pentru diverse tipuri de operații—unele pentru întregi, altele

pentru numere reale—care lucrează în paralel. Astfel, procesorul poate executa mai multe instrucțiuni, sau faze diferite (decodificare, citire/scriere operanzi, calcul efectiv) ale mai multor instrucțiuni, în același timp. Totuși, din punctul de vedere al programării, un calculator secvențial are prin excelență un singur procesor.

La extrema cealaltă se află calculatoarele care au zeci, sute sau chiar mii de procesoare identice, care pot executa fiecare un flux propriu de instrucțiuni, cooperând la rezolvarea unei probleme. Astfel de calculatoare se numesc *masiv paralele*. Deoarece procesoarele sunt ele însele extrem de puternice, un calculator masiv paralel este de sute/mii de ori mai rapid decât calculatoarele secvențiale din aceeași generație. Acestor calculatoare le sunt destinați cei mai mulți dintre algoritmi prezentați în această carte.

Paralel și concurent. Nu este poate prea devreme pentru a face deosebirea între noțiunile de *paralel* și *concurrent*. Presupunând că există mai multe procese (taskuri, părți de program, sau chiar programe) care trebuie să se execute, modul concurent implică numai o eventuală senzație de paralelism, pentru că există o singură unitate centrală, la ocuparea căreia concurează procesele; dacă mecanismul după care se permite alocarea unității centrale este bun, procesele pot părea că se execută în paralel datorită vitezei mari de execuție. Este situația întâlnită în cazul sistemelor de operare multitasking. O execuție cu adevărat paralelă a proceselor se face folosind mai multe procesoare, în cazul ideal câte unul pentru fiecare proces. În multitasking, problemele dificile apar în a face să coopereze mai multe procese ce utilizează aceeași resursă; în paralelism, problema importantă e de a face un program, format din mai multe procese care cooperează, să se execute pe mai multe procesoare. Aceste probleme au, totuși, destule puncte comune; de aceea, nu în puține cazuri, modurile de descriere paralel, respectiv concurent, sunt asemănătoare; diferența esențială, repetăm, apare la execuție.

1.1 O clasificare generală

Una dintre marile probleme ale algoritmicii paralele este că, spre deosebire de cea secvențială, unde modelul de calculator—von Neumann—era unic și foarte simplu conceptual, există o multitudine de clase de arhitecturi care trebuie avute în vedere. Secvențial se gândește într-un singur fel (sau aproape), paralel în foarte multe. Un algoritm excelent pe un anumit tip de calculator poate fi ineficient pe un altul. De aceea nu se poate vorbi despre varianta paralelă a unui algoritm secvențial, ci despre adaptări paralele pe diverse arhitecturi.

Principalele categorii de arhitecturi sunt departajate—într-o clasificare datorată lui Flynn, probabil cea mai populară datorită simplității ei—prin numai două criterii: al fluxului de instrucțiuni și al fluxului de date; ambele pot fi sau unice, sau multiple. Conform acestei clasificări calculatoarele pot fi:

- SISD (Single Instruction Single Data): la un moment dat se execută o singură instrucțiune asupra unei singure date; este calculatorul secvențial obișnuit.

- SIMD (Single Instruction Multiple Data): la un moment dat se execută o singură instrucțiune, dar asupra mai multor date; este un calculator paralel, în care procesoarele execută sincron același program, dar fiecare asupra altui set de date.
- MIMD (Multiple Instruction Multiple Data): fiecare procesor execută un program propriu, folosind date proprii.
- MISD (Multiple Instruction Single Data): o categorie deocamdată vidă, în care aceleași date ar fi prelucrate simultan în mod diferit, de către mai multe procesoare.

Se rețin deci principalele două categorii de calculatoare paralele—SIMD și MIMD—de care ne vom ocupa în continuare.

1.1.1 Arhitecturi SIMD

Arhitecturi pipeline. Din punct de vedere istoric, primul tip de paralelism a fost cel *pipeline*, în care prelucrarea datelor se face ca pe o bandă rulantă, pe care există mai multe posturi de lucru; la fiecare dintre acestea se execută o operație asupra produsului; la un moment dat există pe bandă atâtea produse câte posturi, în diverse faze de execuție. Așadar, modul de prelucrare pipeline este un caz particular de paralelism funcțional. Ideea de bază este de a descompune operațiile în cele mai mici părți componente; de exemplu, pentru a aduna două numere reprezentate în format virgulă mobilă, se parcurg secvențial următoarele etape:

E1: compararea exponenților;

E2: aducerea la același exponent cu numărul mai mare a celui mai mic, a cărui mantisă se modifică corespunzător;

E3: adunarea mantiselor;

E4: normalizarea rezultatului.

Ilustrăm prelucrarea datelor în figura 1.1. O pereche de date D_k trebuie să parcurgă etapele secvențial, dar pot exista patru perechi prelucrate la un moment dat; practic se execută aceeași operație asupra mai multor date, dar prin suboperații diferite. Aceasta permite execuția unei adunări—sau a unei operații, în general—la fiecare perioadă de ceas, cu condiția asigurării unui flux de date suficient; desigur, este necesar un timp de amorsare, timpul în care prima pereche de date parcurge cele patru "posturi de lucru".

Pe acest principiu s-au bazat primele calculatoare *vectoriale*, a căror arhitectură este reprezentată extrem de schematic în figura 1.2. Denumirea provine din eficiența cu care se executau, pe aceste calculatoare, operații vectoriale, ca de exemplu $y = \alpha x + y$, cu x, y vectori n -dimensionali și α un scalar.

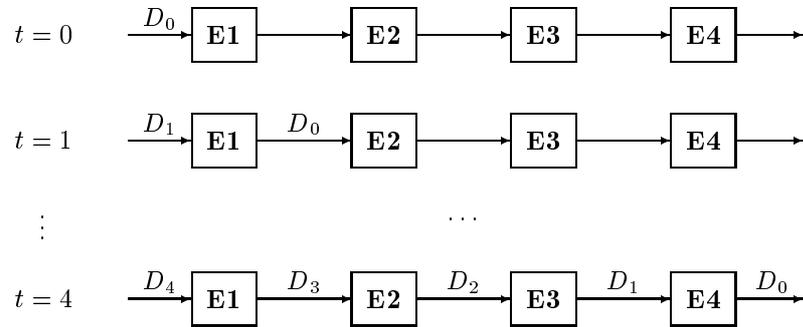


Figura 1.1: Evoluția datelor într-o arhitectură pipeline. Notăm generic D_k datele în toate fazele prelucrării.

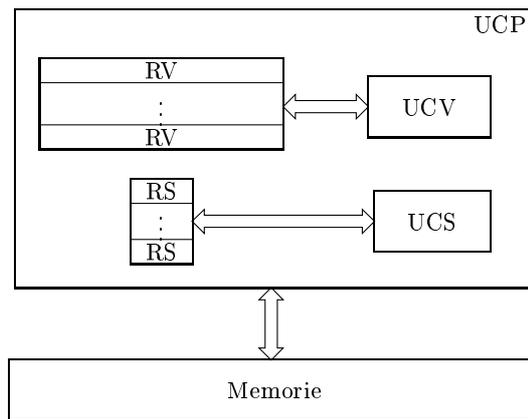


Figura 1.2: Structura unui calculator vectorial monoprocesor. UCP - unitate centrală de prelucrare; RV - registru vectorial; UCV - unitate de calcul vectorial; RS - registru scalar; UCS - unitate de calcul scalar.

Unitatea centrală conține unități de calcul și registre, atât scalare cât și vectoriale. Prin stocarea operanzilor de tip vector în registrele vectoriale, unitatea de calcul vectorial poate fi alimentată foarte rapid, asigurându-i-se funcționare în regim pipeline. Desigur că de foarte mare importanță este și modul de conectare a unității centrale cu memoria, astfel încât să se poată asigura traficul dintre aceasta și registrele vectoriale. Există și unități arhitecturale pentru operații scalare, care pot funcționa (uneori) în paralel cu cele vectoriale.

Consecința acestei organizări arhitecturale este rescrierea algoritmilor în așa fel încât ei să conțină cât mai multe operații vectoriale cu vectori de dimensiune cât mai apropiată de cea a registrelor; astfel, timpul de execuție este mai mic decât pentru aceiași algoritmi în variantă scalară.

La această oră toate procesoarele—inclusiv cele dintr-un simplu PC—conțin elemente de paralelism pipeline; însă este vorba de obicei despre paralelism între decodificarea unei instrucțiuni și execuția celei anterioare, sau/și între acestea și accesul la memorie. Se poate spune că ideea de pipeline a migrat de la nivelul calculatorului, la nivelul procesorului.

SIMD pur paralel. O altă variantă de arhitectură SIMD este cea prezentată în figura 1.3, aici fiind vorba cu adevărat de paralelism. Procesoarele P_i execută simultan aceeași instrucțiune, dictată de unitatea de comandă—cea care deține programul. Operațiile se execută asupra unor date aflate în memoriile locale M_i ; acestea sunt de capacitate mică, dar rapide, și sunt alimentate cu date de la o memorie externă. Uneori există o rețea de interconectare între procesoare și memorii (figurată printr-un dreptunghi cu linie întreruptă), astfel încât, la un moment dat, un procesor poate accesa alte memorii, dar nu pot exista două memorii accesate simultan de același procesor; de exemplu, fiecare procesor accesează memoria vecinului său din dreapta. Acest tip de arhitectură—SIMD veritabil, spre deosebire de cazul pipeline—este adecvat tot calculelor vectoriale; de aceea, astfel de calculatoare sunt numite și ele vectoriale.

De exemplu, adunarea a doi vectori de lungime egală cu numărul de procesoare și cu elementele stocate câte două în fiecare memorie locală, în locații aflate la aceeași adresă relativă la începutul memoriei respective, durează tot atât cât adunarea a două numere reale, adică, în general, un singur tact.

Arhitecturile SIMD de tipul celor din figura 1.3 au fost foarte populare spre sfârșitul anilor '80; exemple de astfel de arhitecturi sunt calculatoarele Connection Machine 2, DAP, MasPar MP-1, Zephir. După 1990, ele au devenit din ce în ce mai puțin competitive, arhitecturile MIMD dovedind posibilități superioare de progres.

1.1.2 Arhitecturi MIMD

Calculatoarele MIMD se împart în două mari categorii, după modul de repartizare a memoriei: cu memorie *partajată (comună)* și cu memorie *distribuită*. În ambele cazuri vom presupune că există p procesoare, notate P_i , $i = 0 \dots p - 1$; vom numi i adresa (sau identitatea) procesorului; vom utiliza în egală măsură reprezentările zecimală

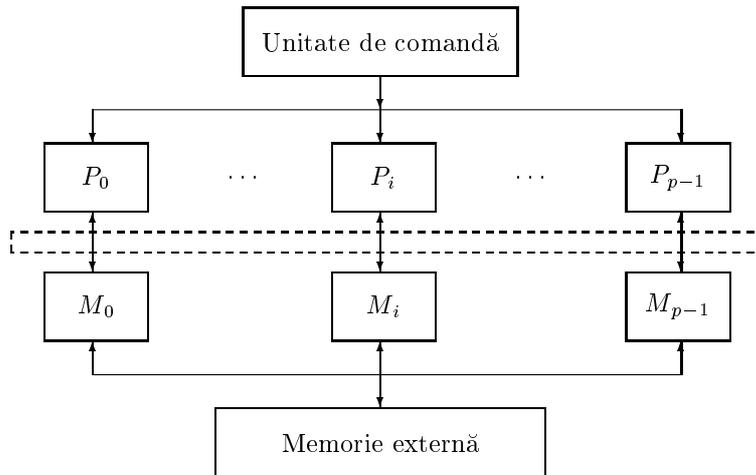


Figura 1.3: Structura unui calculator SIMD cu rețea de procesoare.

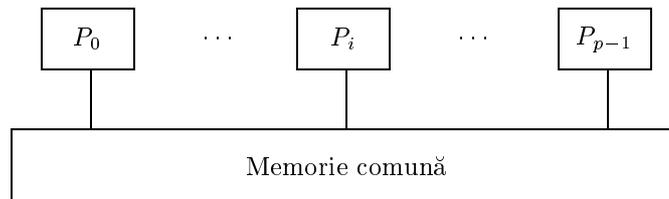


Figura 1.4: Arhitectura de principiu a unui calculator MIMD cu memorie comună.

și binară ale numărului i . Fiecare procesor poate executa programul său propriu, independent de cel al celorlalte procesoare; aceasta nu înseamnă că procesoarele au întotdeauna programe diferite; diferența esențială față de arhitecturile SIMD este că execuția programelor nu se face sincron, ci independent pe fiecare procesor.

MIMD cu memorie partajată. La arhitectura cu memorie partajată, cum este cea prezentată în figura 1.4, există o memorie comună de dimensiuni mari și, eventual, memorii locale (nerepresentate în figură) de dimensiuni modeste, dar mai rapide; în figura 1.5 este reprezentată o configurație cu un singur procesor și cu o astfel de memorie, numită de obicei *cache*; în cazul multiprocesor, accesul la memoria comună (globală) este coordonat de o unitate logică specială, care are rolul de a arbitra conflictele apărute în urma mai multor cereri simultane de acces la memorie.

Principala problemă a acestor arhitecturi este accesul la memoria comună; posibilele conflicte introduc întârzieri, ceea ce nu permite o bună eficiență atunci când

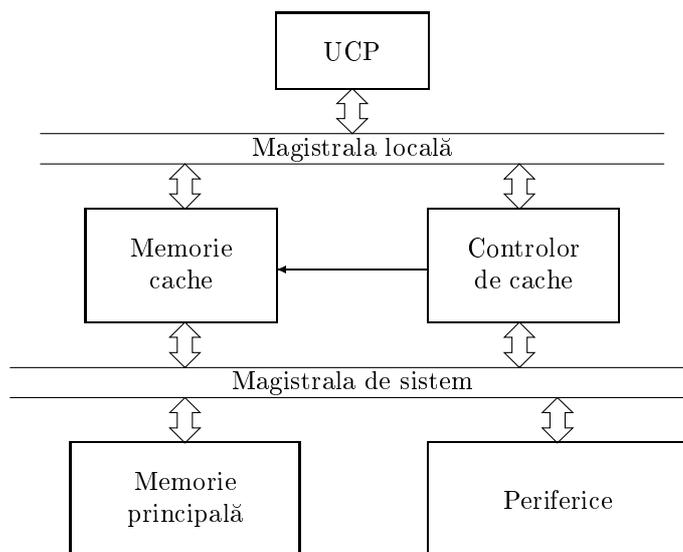


Figura 1.5: Structură monoprocesor cu memorie cache.

numărul de procesoare este mare (peste câteva zeci), adică o dată cu creșterea probabilității apariției conflictelor. Tehnologia actuală permite realizarea unor calculatoare cu memorie comună performante cu cel mult 32 de procesoare.

Se știe că viteza procesoarelor este mai mare, uneori semnificativ, decât cea a memoriilor. De aceea este naturală utilizarea de către un procesor a mai multor blocuri (pagini) de memorie, la care să facă pe rând cereri de citire/scriere, la fiecare tact câte una; de exemplu, dacă un bloc are nevoie de patru tacte pentru a satisface cererea, atunci patru blocuri vor fi, în general, suficiente pentru a nu se obține timpi morți în funcționarea procesorului. Ori în arhitecturile MIMD cu memorie partajată pare a fi exact invers. În fapt nu e chiar așa; imensa memorie comună este împărțită într-un număr relativ mare de blocuri accesibile independent; există acea unitate funcțională a calculatorului care gestionează cererile de acces la memorie ale procesoarelor și care permite simultaneitatea, dar numai în cazul accesului la blocuri diferite; practic, toate blocurile de memorie pot lucra în paralel, reușindu-se alimentarea procesoarelor; de exemplu, dacă un acces la memorie durează patru tacte și sunt 16 procesoare, e nevoie de minimum 64 de blocuri de memorie. În plus, memoriile cache permit micșorarea traficului cu memoria comună, prin reutilizarea unor date depuse acolo, deci deja locale.

Așadar, accesul la memoria comună nu este secvențial. Implicația algoritmică este organizarea accesului la memorie astfel încât să se evite conflictele, să se acceseze deci simultan locații de memorie aflate în blocuri diferite. Esențial în caracterizarea

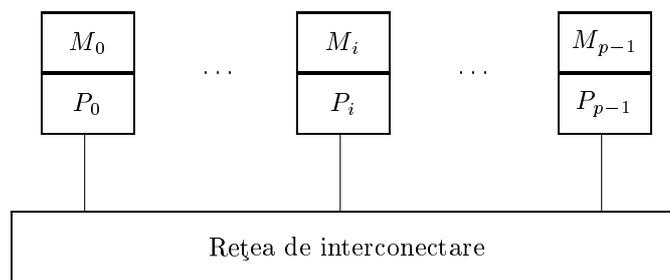


Figura 1.6: Arhitectura de principiu a unui calculator MIMD cu memorie distribuită.

arhitecturii rămâne faptul că fiecare procesor are acces la întreaga memorie. Printre calculatoarele actuale cu memorie comună se numără Sun 10000 Enterprise.

MIMD cu memorie distribuită. În arhitectura cu memorie distribuită, așa cum este prezentată în figura 1.6, fiecare procesor P_i are o memorie locală M_i ; nu există memorie comună. Cooperarea procesoarelor se realizează prin intermediul legăturile asigurate de o rețea de interconectare, cu topologie fixată sau programabilă; de altfel, configurația acestei rețele este una dintre cele mai importante trăsături ale unui astfel de calculator MIMD. Fiecare procesor P_i execută un program propriu aflat în M_i și prelucrează date aflate tot în M_i . Vom dezvolta prezentarea acestui tip arhitectural în secțiunea următoare, deoarece el pare a avea cel mai bun viitor, limitarea numărului de procesoare venind deocamdată mai mult din motive economice.

Arhitecturi mixte. Deoarece arhitecturile cu memorie comună sunt foarte eficiente pentru un număr mic de procesoare, o tendință actuală este de a realiza calculatoare în care clustere de 4-8 procesoare cu memorie comună sunt conectate printr-o rețea de comunicație. Putem spune astfel că, într-un calculator cu memorie distribuită, procesoarele sunt înlocuite cu astfel de clustere. Se îmbină astfel calitățile celor două tipuri de calculatoare MIMD. Este interesant că modul de programare a acestor arhitecturi mixte este tipic fie memoriei comune, fie celei distribuite, fie amândurora (dar nu simultan); așadar, nu s-a creat și un mod de programare specific arhitecturilor mixte. Exemple de astfel de arhitecturi sunt calculatoarele HP 9000 Superdome, SGI Origin 3000.

Accesul la un calculator paralel. O problemă destul de delicată este modul de realizare a legăturilor între calculatoarele MIMD și exterior. De obicei aceasta se face prin intermediul unui calculator gazdă, pe care utilizatorii lucrează în regim multitasking; pe acest calculator se realizează pregătirea programului, adică editarea și compilarea. Cum acest tip de legătură este lent, transmitându-se printr-un singur canal informațional programe și date pentru mai multe procesoare, de către un singur utilizator la un moment dat, tendința actuală este de a permite conectarea utilizatorilor direct la procesoare ale calculatorului paralel, pe care se execută nucleul

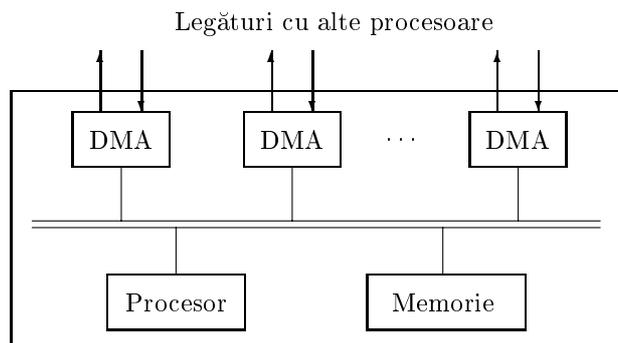


Figura 1.7: Structura simplificată a unui nod într-o arhitectură MIMD cu memorie distribuită.

unui sistem de operare, căruia îi revine sarcina de a ocupa mai multe procesoare la momentul execuției.

Probleme

P 1.1.1 Presupunem că o unitate pipeline realizează o operație în trei etape, pentru fiecare existând câte un "post de lucru". Timpul de calcul în fiecare dintre cele trei etape A , B și C este diferit: 1, 2, respectiv 1 (tacte). La câte tacte produce unitatea pipeline un rezultat? Cum se poate face astfel încât la fiecare tact să se producă un rezultat? În cât timp sunt prelucrate 100 de date?

1.2 Arhitecturi MIMD cu memorie distribuită

1.2.1 Caracteristici ale nodurilor

Revenind la figura 1.6, o pereche procesor-memorie (P_i, M_i) este denumită *nod*. Fiecare nod are un număr de canale de comunicație (porturi) prin care schimbă informații cu alte noduri. Canalele sunt fie prevăzute on-chip—ca la transputerul T800 sau la procesorul actual Compaq Alpha EV67/68, care au câte patru canale—fie asigurate de procesoare (circuite) specializate pentru comunicație. În ambele cazuri, comunicația poate decurge în paralel cu alte operații (calcul) efectuate de către procesor; structura unui nod poate fi simplificată la modelul prezentat în figura 1.7, în care canalele de comunicație sunt realizate prin interfețe DMA (Direct Memory Acces); pentru a comunica pe un canal, procesorul indică interfeței DMA adresa din memoria locală la care se găsește sau la care trebuie depus mesajul, precum și lungimea acestuia; apoi, DMA-ul transmite mesajul pe legătura fizică, accesând memoria concurrent cu procesorul și, eventual, alte interfețe DMA.

Procesoarele utilizate în calculatoarele paralele cu memorie distribuită actuale sunt de obicei de tip RISC (Reduced Instruction Set Computer); exemple actuale sunt Hewlett Packard PA/RISC 8600, MIPS R14000, Sun UltraSPARC III; setul de instrucțiuni al acestor procesoare este redus, obiectivele avute în vedere fiind:

- posibilitatea de a executa o instrucțiune pe ciclu mașină;
- reducerea suprafeței de circuit ocupate de partea de decodificare a instrucțiunilor, eliberând loc pentru memoria cache (o parte din M_i fiind deci pe aceeași pastilă cu P_i) sau pentru partea de comunicație.

S-a ajuns aici deoarece s-a constatat că procesoarele cu set mare de instrucțiuni (CISC – Complex Instruction Set Computer) au două dezavantaje: compilatoarele nu folosesc la întreaga valoare acest set de instrucțiuni și, în plus, partea de decodificare a instrucțiunilor ocupă foarte mult din suprafața circuitului. Simplificând tipurile instrucțiunilor și modurile de adresare, se câștigă în viteza de execuție. Suprafața eliberată astfel poate fi folosită și în scopul introducerii paralelismului în interiorul procesorului. Execuția instrucțiunilor se face, așa cum am menționat, în mod pipeline, cu următoarele faze: citirea instrucțiunii, decodarea ei și citirea operanzilor, execuția propriu-zisă de către unitatea aritmetică și logică, depunerea rezultatului. Aceasta funcționare decurge la viteză maximă atâta vreme cât nu există conflicte sau dependențe de date între instrucțiuni succesive, instrucțiuni de salt sau întreruperi.

O altă tendință este aceea de încerca execuția mai multor instrucțiuni pe tact (procesoare VLIW – Very Long Instruction Word); pentru aceasta e nevoie de mai multe unități funcționale pentru diverse tipuri de calcule (întregi, în virgulă mobilă, sau chiar mai specializat: adunare, înmulțire), fiecare putând funcționa în mod pipeline; procesoarele menționate mai sus au câte două unități de calcul în virgulă mobilă. De asemenea, magistrala de comunicație cu memoria are dimensiuni mari, de exemplu 64 de biți, un cuvânt putând conține mai multe instrucțiuni, care se vor executa în paralel.

1.2.2 Câteva noțiuni din teoria grafurilor

Putem asocia imediat o arhitectură cu memorie distribuită unui graf neorientat: procesoarele sunt nodurile iar liniile de comunicație între procesoare (considerate bidirecționale)—arcele. De aceea este necesară reamintirea unor noțiuni din teoria grafurilor și reliefaarea semnificațiilor lor specifice acestui context. (De asemenea, cum se va vedea mai târziu, și unui algoritm îi putem asocia un graf, important pentru detectarea paralelismului.)

Un *graf* neorientat $G = (V, E)$ este constituit dintr-o mulțime V , a nodurilor, și o alta, $E \subset V \times V$, conținând perechi de noduri, (în care nu distingem ordinea, deci $(x, y) = (y, x)$), adică arcele care unesc aceste noduri¹.

Numărul de noduri ale unui graf este numit *ordin*; îl vom nota cu p ; în cazul arhitecturilor paralele, p este numărul de procesoare.

¹Notăția (V, E) este oarecum tradițională; în engleză, *vertex* înseamnă nod, iar *edge* arc.

Un arc de forma (v, v) se numește *bucă*. Un graf fără bucle se numește *simplu*. Grafurile care ne interesează vor fi toate simple, o cale de comunicație de la un procesor la el însuși neavând sens.

Două noduri sunt *adiacente* (sau vecine), dacă există un arc între ele. Numărul de vecini ai unui nod se numește *gradul* nodului respectiv; este deci numărul de arce plecând (sau venind, sensul nu contează) dintr-un nod. Altfel zis, gradul este numărul de canale de comunicație ale unui procesor.

Gradul unui graf (Δ) este maximul gradelor nodurilor componente. Graful este *regulat* dacă gradele tuturor nodurilor sunt egale. Observăm că numărul de arce al unui graf regulat este $p\Delta/2$.

Se numește *drum* (sau *cale*) între două noduri $x, y \in V$ o secvență de noduri x_0, x_1, \dots, x_k , astfel încât două noduri consecutive sunt adiacente și $x_0 = x, x_k = y$. Drumul se numește *elementar* dacă nodurile componente x_1, \dots, x_{k-1} sunt distincte. Într-o arhitectură cu memorie distribuită, un drum este o cale de comunicație între două procesoare; vor fi utilizate doar drumuri elementare, după cum este logic.

Lungimea unui drum este numărul de arce ale sale, k în notația de mai sus. *Distanța* dintre două noduri, $dist(x, y)$, este lungimea celui mai scurt drum între acestea; în contextul memorie distribuită, vom nota $dist(P_i, P_j)$ sau chiar $dist(i, j)$ distanța între procesoarele cu adresele i și j .

Diametrul grafului (D) este maximul tuturor distanțelor între perechi de noduri ale grafului, adică distanța dintre cele mai depărtate noduri. *Excentricitatea* unui nod este cea mai mare distanță dintre acesta și un alt nod; pentru orice nod, excentricitatea este mai mică sau egală cu diametrul (o putem privi ca diametrul din punct de vedere al nodului respectiv).

Vom numi *ciclu* un drum elementar pentru care nodul inițial și cel final sunt identice. Un *ciclu hamiltonian* este un ciclu trecând (exact o dată) prin fiecare nod al grafului; lungimea sa este egală cu ordinul grafului. Un graf posedând un ciclu hamiltonian se numește *graf hamiltonian*.

Un graf se numește *conex* dacă există (cel puțin) un drum între oricare două noduri ale sale.

Un *arbore* este un graf conex fără cicluri.

1.2.3 Topologii curente

Calitățile topologiei rețelei de interconectare. Din punctul de vedere al algoritmilor pentru arhitecturi cu memorie distribuită, este de mare importanță topologia rețelei de interconectare, deoarece aceasta restrânge posibilitățile de comunicare. Să vedem întâi care ar trebui să fie calitățile acestei topologii.

Evident, graful trebuie să fie conex. Apoi, pe de o parte, e de dorit ca gradul grafului să fie cât mai mic și, dacă se poate, graful să fie regulat; deci, un procesor să nu aibă prea multe canale de comunicație, pentru a fi ușor de realizat (un număr mare de canale înseamnă complexitate mare a procesorului și număr ridicat de pini), și toate procesoarele—identice, în mod obișnuit—să-și folosească integral canalele;

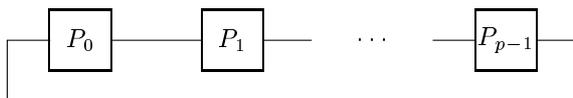


Figura 1.8: Calculator MIMD cu topologie de inel.

strâns legată de aceasta este cerința ca numărul total de arce să fie relativ mic, fiecare arc însemnând o conexiune fizică realizată printr-un traseu electric pe o placă și, eventual, între plăci. Pe de altă parte, e cu atât mai bine cu cât diametrul este mai mic, astfel încât distanțele între procesoare să fie limitate, iar comunicarea cât mai facilă.

Este clar că aceste două deziderate sunt antagonice; cu cât gradul e mai mic, cu atât sunt mai puține arce, deci scad șansele să existe drumuri scurte între oricare două procesoare. Micșorarea gradului și a numărului de arce este o cerință hardware, cea a diametrului provine din software.

La o extremă se află graful *total interconectat* (sau *complet*), în care există arce între oricare două noduri; diametrul său este 1, gradul $O(p)$ (mai exact² $p - 1$) și numărul de arce $O(p^2)$ (exact $p(p - 1)/2$); dacă diametrul este perfect, gradul și numărul de arce sunt imense; o astfel de topologie este realizabilă practic doar pentru un număr infim de procesoare. La cealaltă extremă se află inelul, despre care vom vorbi imediat.

În continuare vom enumera câteva tipuri curente de configurații, împreună cu unele proprietăți ale lor. Este vorba despre grafuri cu foarte bune proprietăți de simetrie, pentru care descrierea algoritmilor se face într-un mod mult mai simplu decât în cazul unor asimetrii pronunțate.

Inelul. Într-un inel (figura 1.8), procesorul P_i este legat de P_j , unde $j = (i - 1) \bmod p$ sau $j = (i + 1) \bmod p$; în primul caz vom spune că P_j este vecinul din stânga, în al doilea, cel din dreapta. Gradul inelului este $\Delta = 2$, numărul de arce p (ambele valori fiind bune), dar diametrul este $D = \lfloor p/2 \rfloor$, adică mare; de aceea, în astfel de arhitecturi se utilizează puține procesoare; de altfel, nici unul dintre calculatoarele performante actuale nu are această topologie (deși ea a fost utilizată la începutul anilor '80 în diverse modele experimentale). Din punct de vedere al notației procesoarelor, pentru simplificare, vom considera că notația P_i este corectă, indiferent de valoarea întregului i , deci $P_i \equiv P_{i \bmod p}$; de exemplu, P_{-1} și P_{p-1} indică același procesor.

Grila de dimensiune 2 (figura 1.9a) are $p = m \cdot n$ procesoare. Pentru simplitate, convenim să asociem procesorului o adresă cu două coordonate, e.g. (i, j) ; o adresă formată dintr-un număr în intervalul $0 : p - 1$ se poate obține banal, de pildă $(i - 1) * p + j$ corespunde lui (i, j) , dacă numerotăm procesoarele în ordinea liniilor. Procesorul

²Se spune despre o funcție $f(n)$ că are valori de ordinul $O(g(n))$ dacă, pentru toți n mai mari decât un n_0 , există două constante c_1 și c_2 astfel încât $c_1 g(n) \leq f(n) \leq c_2 g(n)$. Cu alte cuvinte, funcția $g(n)$, de obicei simplă (monom sau logaritm), poate fi folosită pentru aproximarea ordinului de mărime al funcției mai "complicate" $f(n)$.

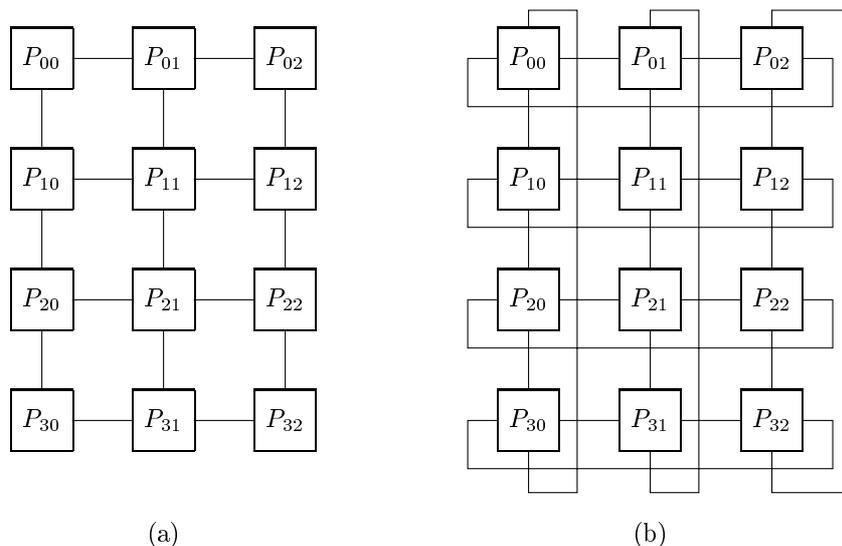


Figura 1.9: Exemple de calculator MIMD cu topologie de (a) grilă 4×3 (b) tor 4×3 .

P_{ij} este legat de procesorul P_{kl} (cu $0 \leq i, k \leq m-1$ și $0 \leq j, l \leq n-1$), una singură din condițiile următoare fiind îndeplinită: $|k-i|=1$ sau $|l-j|=1$. Un procesor are cel mult patru vecini ($\Delta = 4$), numiți stânga, dreapta, sus, jos sau vest, est, nord, sud; diametrul grafului este $D = m + n - 2$, iar numărul de arce $2mn - m - n$. Un exemplu de mare succes la mijlocul deceniului trecut este calculatorul Intel Paragon.

Torul de dimensiune 2 (figura 1.9b), este o grilă 2D, cu procesoarele de la margini legate între ele astfel încât fiecare linie sau coloană de procesoare să formeze un inel. Torul este regulat, cu $\Delta = 4$, $D = \lfloor m/2 \rfloor + \lfloor n/2 \rfloor$ (vezi diametrul inelului) și număr de arce $2mn$. Ca și la grilă, procesoarele vor fi notate P_{ij} ; în afară de aceasta, se va folosi artifiциul de notație de la inel, adică $P_{ij} \equiv P_{i \bmod m, j \bmod n}$. Un exemplu de calculator cu topologie tor 2D este Fujitsu AP3000. Grila și torul se pot generaliza imediat în 3 dimensiuni. Calculatorul Cray T3E, unul din cele mai puternice actualmente, are o topologie de tor 3D.

Hipercubul (figura 1.10) a fost cea mai populară arhitectură în perioada 1985-1990; marele său succes a venit atât din calitățile intrinseci ale grafului, cât și din numărul mare de algoritmi ce i se potrivesc (vom vedea mai târziu care e sensul acestei potriviri); calitățile algoritmice sunt cele care îl fac în continuare demn de un interes special. Notăm \mathcal{H}_d hipercubul de dimensiune d , care are $p = 2^d$ procesoare. Fiecare procesor este legat de d vecini; P_i este conectat de P_j dacă reprezentările binare ale

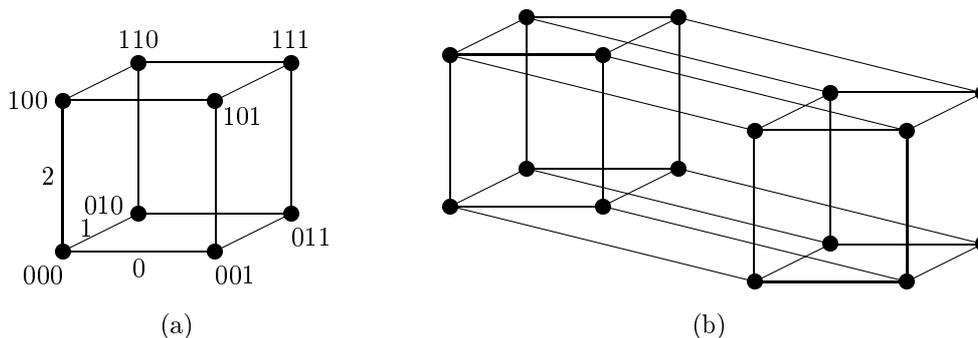


Figura 1.10: Două exemple de hipercub: (a) de dimensiune 3 (b) de dimensiune 4.

numerelor i și j diferă printr-un singur bit³; de exemplu, în hipercubul de dimensiune 3, P_0 este conectat cu P_1 , P_2 și P_4 (vom mai spune, cu același înțeles, că 000 este conectat cu 001, 010 și 100 sau că 0 este conectat cu 1, 2 și 4). Deci $\Delta = d = \log p$, numărul de arce $p(\log p)/2$ și diametrul $D = \log p$ (vezi problema 1.2.3). Un exemplu de calculator cu topologie de hipercub este familia nCUBE.

Hipercubul poate fi definit și recursiv; \mathcal{H}_d poate fi obținut luând două copii ale lui \mathcal{H}_{d-1} și adăugându-le arcele care unesc nodurile aflate în aceeași poziție în cele două copii; la adresele nodurilor se va adăuga un bit în poziția cea mai semnificativă: 0 pentru prima copie a lui \mathcal{H}_{d-1} , 1 pentru a doua.

Vom spune că arcul unind P_i și P_j se găsește în dimensiunea l dacă i și j diferă în bitul din poziția l ; ne vom referi întotdeauna la reprezentări binare: dacă $i = i_{d-1}i_{d-2}\dots i_l\dots i_0$, atunci $j = i_{d-1}i_{d-2}\dots \bar{i}_l\dots i_0$; se mai spune că P_i și P_j sunt vecini în dimensiunea l ; vezi, de exemplu, figura 1.10, în care sunt precizate dimensiunile în care se află arcele plecând din P_0 . Se vede deci că eliminând toate arcele dintr-o anumite dimensiune a lui \mathcal{H}_d , se vor obține două hipercuburi \mathcal{H}_{d-1} ; eliminând apoi arcele dintr-o altă dimensiune, vor rezulta patru hipercuburi \mathcal{H}_{d-2} etc. Notăm \mathcal{H}_d^r familia de 2^{d-r} hipercuburi \mathcal{H}_r , obținute dintr-un hipercub \mathcal{H}_d prin eliminarea arcelor din dimensiunile $d-1, d-2, \dots, d-r$; altfel zis, nodurile din fiecare din aceste hipercuburi au primii (i.e. cei mai semnificativi) $d-r$ biți identici. Notăm $\mathcal{H}_d^r(q)$ un astfel de hipercub, anume pe cel care conține procesorul P_q . În mod analog, notăm \mathcal{H}_d^r hipercuburile de dimensiune r , în fiecare nodurile având *ultimii* (i.e. cei mai puțin semnificativi) $d-r$ biți identici. De exemplu, în \mathcal{H}_3 , hipercubul $\mathcal{H}_3^2(0)$ are dimensiune 2 și conține nodurile 000, 001, 010, 011, iar hipercubul $\mathcal{H}_3^2(0)$ conține nodurile 000, 010, 100, 110.

Simetrie. Pentru inel, tor și hipercub, simetria poate fi subliniată foarte simplu.

³Se mai spune că distanța Hamming dintre cele două noduri este egală cu 1; distanța Hamming este numărul pozițiilor în care reprezentările binare diferă: $\sum_{i=0}^{d-1} i_i \oplus j_i$; notăm cu i_i bitul din poziția i , în reprezentarea binară a lui i .

Topologie	Δ	D	N
Inel	2	$O(p)$	$O(p)$
Grilă	4	$O(\sqrt{p})$	$O(p)$
Tor	4	$O(\sqrt{p})$	$O(p)$
Tor 3D	6	$O(\sqrt[3]{p})$	$O(p)$
Hipercub	$\log p$	$O(\log p)$	$O(p \log p)$

Tabelul 1.1: Gradul Δ , diametrul D și numărul de arce N pentru cele mai importante topologii.

Notând un procesor oarecare cu P_0 , și pe celelalte în funcție de acesta, se obține aceeași construcție în toate cazurile. Cu alte cuvinte, toate procesoarele au același context local, toate "văd" ansamblul în același fel; aceasta este o proprietate esențială în dezvoltarea algoritmilor: un procesor nu cunoaște decât adresa sa și direcțiile pe care se află canalele sale de comunicație, deci vecinii săi; vom discuta aceste aspecte în capitolul următor. În plus, toate nodurile au aceeași excentricitate, egală cu diametrul. Din aceste motive, unii algoritmi vor fi prezentați în forma pentru P_0 , pentru inel sau hipercub, sau P_{00} , pentru tor, fără ca particularizarea să fie simplificatoare; vom numi *central* acest procesor.

Produs cartezian. Topologiile de mai sus pot fi descrise alternativ sub forma unui produs cartezian. Notăm $\mathcal{I}(p) = 0 : p-1$; asociem elementelor lui $\mathcal{I}(p)$, văzute ca noduri ale unui graf, legături care formează un inel. Cu această convenție, identificăm $\mathcal{I}(p)$ cu un inel. Se observă imediat că torul $m \times n$ este produsul cartezian $\mathcal{I}(m) \times \mathcal{I}(n)$. De asemenea, hipercubul este $\mathcal{I}^d(2) = \{0, 1\}^d$.

Comparații. Recapitulăm principalele caracteristici ale topologiilor prezentate în tabelul 1.1, unde presupunem că grila și torul au același număr de procesoare în orice dimensiune; lăsăm cititorului demonstrația proprietăților torului 3D. În tabel, topologiile sunt ordonate crescător după grad și descrescător după diametru.

Rețele configurabile. Rețelele de interconectare pot fi și **configurabile**, de exemplu prin intermediul unor comutatoare în cruce (crossbar). Un astfel de comutator permite conectarea a $2n$ linii, putându-se realiza conectarea între orice linie dintr-un grup de n cu orice alta din celălalt grup de n . Se poate concepe o arhitectură ce permite oricăror două procesoare să fie vecine. Complexitatea rețelei este însă foarte mare, de aceea nu se construiesc crossbar cu mai mult de 32 de intrări; se folosesc în schimb diverse combinații de crossbar-uri, care permit realizare unui număr mare de topologii. Ca exemple, amintim calculatoarele din familia supernodurilor (din anii '90), în care procesoarele erau transputere și care puteau realiza orice topologie cu grad mai mic sau egal cu patru (numărul de canale al transputerului); de asemenea, calculatoarele actuale Sun 10000 Starfire, HP 9000 Superdome și SGI Origin 3000 (ultimele două sunt formate din clustere conectate prin comutatoare în cruce).

Alte topologii. Pornind de la constrângerile constructive legate de grad și di-

ametriu, se poate formula *problema* (Δ, D) : care este numărul maxim de noduri al unui graf cu grad Δ și diametru D ? Altfel spus, să se găsească topologia cu grad și diametru impuse, dar cu cel mai mare număr de procesoare; problema este pusă oarecum pe dos: ar fi fost poate mai natural să se impună gradul și numărul de procesoare și să se minimizeze diametrul; în fond, este exact același lucru. Problema nu e decât parțial rezolvată. De exemplu, pentru $\Delta = 2$, inelul este cel care atinge optimul. Este curios că grafurile obținute în mod aleator au de multe ori foarte bune proprietăți din acest punct de vedere; ele sunt însă inutilizabile, datorită neregularității lor. Oricum, formularea problemei a condus la o multitudine de propuneri de posibile topologii, care, deși nu au fost puse în practică la scară mare, dovedesc bune proprietăți teoretice, nu numai din punctul de vedere (Δ, D) , ci și din cel al adecvării diversilor algoritmi la aceste topologii. Menționăm aici doar grafurile *de Bruijn* și *Cube-Connected-Cycle*.

Tendențe actuale. Până în 1990, hipercubul părea a fi topologia cea mai promițătoare. După aceea, grila și torul au beneficiat de avantajul ca li se pot adăuga foarte ușor procesoare, cu modificări minime în configurația celor existente; grila reprezintă un graf planar (se poate reprezenta în plan fără ca două arce să se intersecteze), iar torul se poate construi în două planuri, de aceea sunt ușor de realizat fizic. Torul 3D beneficiază de avantaje asemănătoare. În prezent, tendința este de a face topologia transparentă pentru utilizator; comunicația este lăsată sistemului de operare sau unor biblioteci specializate, care permit comunicarea între oricare două noduri. Totuși, cunoașterea topologiei este importantă din cel puțin două motive:

- mulți algoritmi folosesc în mod natural o topologie virtuală în care comunicațiile se desfășoară numai între vecinii în graf; conceperea algoritmilor se bazează deci pe o topologie de comunicație;
- performanța obținută cu un program paralel poate fi îmbunătățită sensibil atunci când topologia virtuală a algoritmului corespunde topologiei reale a calculatorului.

1.2.4 Scufundări

Așa cum am sugerat mai sus, în proiectarea de algoritmi paraleli se preferă o oarecare distanțare față de arhitectură, pentru a nu fi necesară rescrierea periodică a algoritmilor, la eventuale schimbări de arhitectură. O cale de a realiza portabilitatea este de a gândi algoritmi pe topologii simple (această simplitate poate fi văzută în două feluri: topologia în sine este simplă sau topologia este adecvată algoritmului, deci simplă din punctul de vedere al acestuia; avem în vedere ambele sensuri) și de a găsi o metodă de "suprapunere" a acestei topologii cu cea a arhitecturii reale; sunt, așadar, două nivele: unul al legăturilor logice, tipice algoritmului, altul al legăturilor fizice, specifice mașinii; se încearcă o potrivire a legăturilor logice cu cele fizice; vom reveni mai pe larg în capitolul 4.

Problema este interesantă și din alte motive, de exemplu posibilitatea emulării unei arhitecturi prin alta și determinarea efortului cerut de această emulare; să ne postăm deocamdată în acest punct de vedere. Această problemă este cunoscută sub numele de *scufundare a grafurilor*.

Vom prezenta întâi formularea matematică a problemei, explicând apoi semnificația noțiunilor, de altfel destul de intuitive. Un studiu mai general poate fi găsit [22].

Definiții. Fie G și H două grafuri simple neorientate. O scufundare a grafului G în graful H este definită printr-o funcție f definită pe mulțimea nodurilor din G , cu valori în mulțimea nodurilor din H și o funcție S_f definită pe mulțime arcelor din G , cu valori în mulțimea căilor din H ; prin S_f se asociază unui arc xy din G , o cale în H care unește $f(x)$ cu $f(y)$.

De obicei se vorbește despre o scufundare doar prin intermediul funcției f , lăsând S_f subînțeleasă (unui arc xy asociindu-i-se una din căile de lungime minimă între $f(x)$ și $f(y)$). Pentru a măsura eficiența unei scufundări se folosesc mai mulți parametri, dintre care amintim pe cei mai importanți:

- *dilatatarea* unei scufundări, notată $dil(f)$, este lungimea maximă a unei căi $S_f(xy)$, cu $x, y \in V(G)$, unde $V(G)$ este mulțimea nodurilor grafului G . În cazul în care se consideră căile cele mai scurte—deci lungimea căii $S_f(xy)$ este $dist_H(f(x), f(y))$, adică distanța pe H dintre $f(x)$ și $f(y)$ —atunci dilatarea se exprimă unic în funcție de f prin (se notează cu $E(G)$ mulțimea arcelor din G):

$$dil(f) = \max_{xy \in E(G)} dist_H(f(x), f(y)).$$

- *congestia* unei scufundări f a lui G în H este maximul, după toate arcele $e \in E(H)$, a numărului de căi $S_f(xy) \in E(H)$, imagini de arce din G , care conțin e .

Putem privi H ca fiind arhitectura reală, iar G topologia emulată (fie că ea provine dintr-un algoritm, fie că este realmente a unei arhitecturi). Dilatarea exprimă distanța reală maximă între două noduri emulate; este de dorit să fie cât mai mică, valoarea ideală fiind 1. Congestia dă informații despre numărul de canale emulate pe un canal de comunicație fizică; cu cât acest număr este mai mic, conflictele de utilizare a canalului vor fi mai reduse; valoarea ideală este 1; dacă dilatarea este egală cu 1, atunci și congestia este tot 1.

Exemple. Vom prezenta doar câteva rezultate fundamentale, cu aspect intuitiv foarte pregnant, de aceea fără detalii. Ne vom referi la legăturile între topologiile cele mai curențe, anume inel, grilă, tor și hipercub.

O primă problema este de a afla dacă un inel poate fi scufundat cu dilatare 1 în topologiile de mai sus sau, cu alte cuvinte, dacă se poate găsi pentru acestea un *ciclu hamiltonian*.

O grilă 2D de dimensiuni $m \times n$ are ciclu hamiltonian dacă și numai dacă m sau n sunt pare ($m, n > 2$). Un exemplu semnificativ în ce privește construirea ciclului este prezentat în figura 1.11a, pe o grilă 4×5 . Într-o grilă cu ambele dimensiuni impare se poate scufunda un inel cu dilatare 2; se observă din figura 1.11b că linia oblică nu

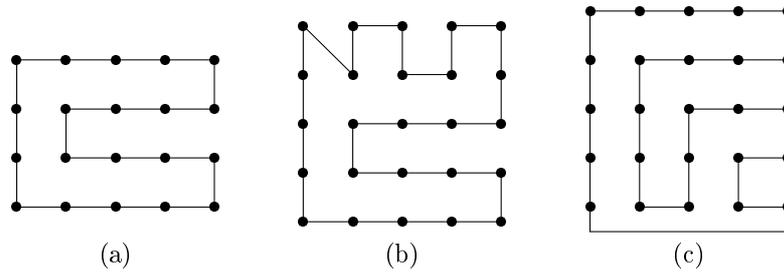


Figura 1.11: (a) Ciclu hamiltonian într-o grilă 2D 4×5 ; (b) scufundarea unui inel într-o grilă 5×5 , cu dilatare 2; (c) ciclu hamiltonian într-un tor 5×5 .

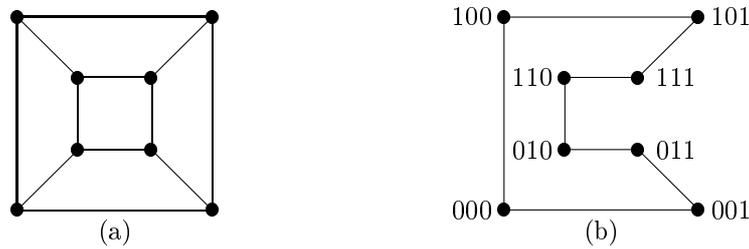


Figura 1.12: (a) un alt mod de desenare a unui hiper cub de dimensiune 3 (b) ciclu hamiltonian în hiper cubul de dimensiune 3.

corespunde unui arc al grilei; această unică conexiune trebuie realizată prin două arce ale grilei; în ambele variante posibile, dreapta-sus sau stânga-jos, congestia este 2.

Un tor 2D are întotdeauna ciclu hamiltonian; dacă m sau n sunt pare, se aplică rezultatul precedent; dacă ambele dimensiuni sunt impare, construcția este (de exemplu) asemănătoare cu cea din figura 1.11c, pentru un tor 5×5 .

De asemenea, un hiper cub are întotdeauna un ciclu hamiltonian (de fapt se pot construi d cicluri disjuncte pentru un hiper cub $2d$ -dimensional, așa cum se arată în [19]). Ordinea în care se parcurg nodurile este cea din codul Gray reflectat. Un cod Gray de dimensiune d este o permutare a celor 2^d numere binare pe d biți, astfel încât două numere consecutive diferă printr-un singur bit (nodurile cu numerele vecine sunt vecine în hiper cub). De exemplu, codul Gray reflectat se construiește recursiv, după cum urmează: $G_1 = (0, 1)$, $G_d = (0G_{d-1}, 1G_{d-1}^R)$, i.e codul de dimensiune d se obține din codul de dimensiune $d - 1$, căruia i se adaugă cifra 0 pe poziția cea mai semnificativă, urmat de același cod de dimensiune $d - 1$, dar în ordine inversă, cu cifra 1 adăugată pe poziția cea mai semnificativă. Deci, $G_2 = (00, 01, 11, 10)$, iar G_3 poate fi văzut, în ciclul hamiltonian al unui hiper cub de dimensiune 3, în figura 1.12.

O grilă 2D de dimensiuni $2^{n/2} \times 2^{n/2}$ poate fi scufundată cu dilatare 1 într-un

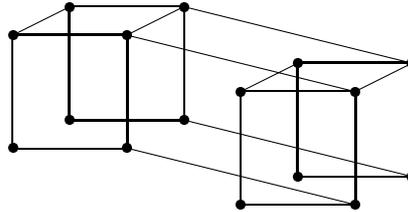


Figura 1.13: Grilă 4×4 mapată într-un hipercube de dimensiune 4; un tor 4×4 este chiar hipercube !

hipercub n -dimensional (cu n par), ca și, de altfel, orice grilă de dimensiuni $2^{n_1} \times 2^{n_2}$, cu $n_1 + n_2 = n$; un exemplu este prezentat în figura 1.13. Și pentru torul 2D este valabil acest rezultat.

Probleme

P 1.2.1 Să se arate că numărul de arce dintr-un arbore de ordin p este $p - 1$.

P 1.2.2 Care este numărul de arce într-un graf regulat de ordin p și grad Δ ?

P 1.2.3 Să se arate că distanța între două noduri oarecare cu adresele x, y din hipercubele \mathcal{H}_d este egală cu numărul de poziții în care cifrele din x și y diferă (distanța Hamming între x și y).

P 1.2.4 Câte noduri se află la distanță x de P_0 , într-un hipercube \mathcal{H}_d ?

P 1.2.5 Care sunt vecinii procesorului P_{ij} într-un ciclu hamiltonian pe o grilă $n \times n$, cu n par ?

P 1.2.6 Cum se face conversia între reprezentarea binară a unui număr și codul său Gray (reflectat) ? Care sunt vecinii procesorului P_g într-un ciclu hamiltonian în hipercubele \mathcal{H}_d ? Caz particular: vecinii lui P_4 într-un ciclu hamiltonian în \mathcal{H}_4 .

1.3 Criterii de performanță

Vom prezenta pe scurt câteva criterii prin care se apreciază posibilitățile hardware ale unui calculator paralel:

- *Viteza de calcul* maximă, măsurată în Mflop/s (Million Floating point Operation per Second) sau, mai rar, în MIPS (Million Instructions Per Second). În primul caz se măsoară numărul de operații în virgulă mobilă efectuate într-o secundă (alte prefixe: G pentru giga— 10^9 operații, sau T pentru tera— 10^{12} operații), în al doilea numărul de instrucțiuni. În descrierea vitezelor de calcul, se folosesc două valori: una de vârf, care presupune că procesoarele lucrează la încărcare maximă (ignorând

conflictele de acces la memorie, comunicația, sincronizările, timpii morți), și alta practică, evaluată cu programe speciale de test (benchmark), care rezolvă probleme des întâlnite în calculele științifice. Unul dintre cele mai utilizate teste este LINPACK, care conține rezolvarea unor probleme de algebră matriceală, ca de exemplu rezolvarea de sisteme de ecuații liniare. Recordul mondial actual este de peste 7 Tflop/s, realizat de calculatorul IBM ASCI White de la Lawrence Livermore National Laboratory, în anul 2000, și măsurat la rezolvarea unui sistem de 518000 ecuații liniare.

- *Numărul de procesoare.* Un calculator paralel, în special în clasa MIMD cu memorie distribuită, nu are un număr fixat definitiv de procesoare. Datorită puterii de cumpărare și necesităților diferite ale potențialilor clienți, se proiectează variante arhitecturale cu număr de procesoare de la câteva zeci la câteva mii, diferențele conceptuale între acestea fiind minime.

- *Dimensiunea memoriei locale* a unui procesor, sau a memoriei comune, dacă există. Acesta e un parametru de mare importanță, deoarece dictează dimensiunea problemelor care pot fi rezolvate. Pentru calculatoarele MIMD cu memorie distribuită actuale, memoria locală a unui procesor este de ordinul Go.

- *Frecvența ceasului unui procesor*, care dă informații despre puterea unui nod. În prezent, procesoarele utilizate în calculatoarele paralele lucrează la frecvențe de sub 1 GHz; totuși, datorită paralelismului lor funcțional, ele pot furniza viteze de calcul de peste 1 Gflop/s.

- Caracteristici legate de *viteza de comunicație* și în special banda de trecere a unei legături, măsurată în număr de octeți transferați pe secundă. Valorile curente sunt de ordinul 500 Mo/s. Totuși, ca și la viteza de calcul, trebuie să diferențiem între viteza de vârf, valabilă doar pentru mesaje foarte lungi (când nu se ia în considerare timpul de amorsare a comunicației etc.) și cea atinsă efectiv în aplicații reale. Vom discuta despre caracteristicile comunicației în capitolul 3.

- *Viteza de acces la memorie* a fiecărui procesor, esențială pentru calculatoarele cu memorie comună. Valorile actuale ale acestei viteze sunt de ordinul sutelor de Mo/s.

- Viteza de comunicație cu exteriorul, care poate produce un efect de strangulare a vitezei efective de calcul, prin introducerea unor timpii morți datorati nealimentării cu programe și date a calculatorului paralel.

- Suportul software oferit o dată cu calculatorul. În ultimii ani, se constată un efort intens de standardizare, materializat în apariția de biblioteci implementate eficient pe toate arhitecturile (sau pe clase de arhitecturi). Aceasta permite portabilitatea programelor paralele, ceea ce acum 10 ani părea un țel îndepărtat. Vom prezenta câteva biblioteci importante în această carte.

- Raportul între costul calculatorului și viteza sa de calcul sau, mai popular, *prețul* unui Mflop/s. Acesta, ca peste tot, e de multe ori hotărâtor. Tendința actuală este de a menține cât mai scăzut acest raport, chiar dacă nu se obține o viteză de calcul deosebit de ridicată.

Este clar că trebuie să existe o armonie între valorile acestor parametri. De exemplu, nu are sens ca viteza de comunicație să fie mult mai mare decât viteza de calcul:

un element în virgulă mobilă să poată fi mult mai repede transmis decât calculat. Evident că nici situația inversă nu e de dorit. Există, de asemenea, corelații între viteza de calcul și capacitatea memoriei; aceasta din urmă trebuie gândită după dimensiunea problemelor pe care calculatorul le poate rezolva în timp rezonabil. Criteriul principal de apreciere a unui calculator paralel rămâne totuși viteza de calcul.

Lista celor mai rapide 500 de calculatoare din lume, actualizată semestrial, poate fi găsită la <http://www.top500.org>. Pe primele locuri se găsesc calculatoare de construcție specială, în majoritate unicate, cu mai mult de 1000 de procesoare. Pentru ilustrarea criteriilor de mai sus, prezentăm fișa tehnică a două calculatoare comerciale.

Cray T3E este un calculator MIMD cu memorie distribuită, cu o topologie de tor 3D; nodurile sunt construite cu procesoare DEC Alpha 21164, cu o frecvență a ceasului de 675 MHz și o performanță de vârf de 1.35 Gflop/s (procesorul are două unități de calcul în virgulă mobilă), și cu o memorie de 2Go. Numărul de procesoare poate varia între 40 și 2176. Viteza maximă de comunicație pe fiecare legătură între noduri este de 325 MB/s, bidirecțional. Viteza maximă de calcul atinsă este de 1.127 Tflop, pentru rezolvarea unui sistem liniar de dimensiune 148800, folosind 1488 procesoare. Mai multe informații se pot obține la <http://www.cray.com>.

Sun Enterprise E10000 este un calculator MIMD cu memorie comună, având 4 până la 64 procesoare UltraSPARC, funcționând la o frecvență de 400 MHz și având o performanță de vârf de 800 Mflop/s. Capacitatea memoriei comune este de până la 64 Go, iar viteza de acces la memorie de 200-400 Mo/s pe procesor (valoarea minimă este pentru numărul maxim de procesoare). Viteza maximă raportată este 43.8 Gflop/s (test LINPACK). Mai multe informații se pot obține la <http://www.sun.com>, iar pentru exemplarul instalat în Universitatea Politehnica București, la <http://hpc.ss.pub.ro>.

Capitolul 2

Modele de programare

În acest capitol vom prezenta modele de programare pentru calculatoarele MIMD, cu memorie distribuită sau comună. Aceste modele ne vor permite descrierea precisă a algoritmilor într-un pseudocod care poate fi translatat cu ușurință în formă implementabilă pe un calculator paralel.

Paradigma unică de programare folosită în această carte și în calculatoarele MIMD actuale este SPMD (Single Program Multiple Data): toate procesoarele execută același program, dar fiecare utilizează un set propriu de date. Acest mod de lucru combină versatilitatea arhitecturilor MIMD cu facilitatea programării.

Organizarea memoriei în arhitectura paralelă determină caracteristicile secundare ale modelului de programare, referitoare la modul în care cooperează procesoarele. În cazul memoriei distribuite se lucrează cu modelul de comunicație prin mesaje; mai precis, procesoarele comunică numai prin intermediul unor mesaje transmise de la un procesor sursă la unul destinație; așadar, comunicația este explicită. În cazul memoriei partajate, comunicația se face implicit prin memoria comună; modelul de programare se bazează pe arhitectura ideală PRAM (Parallel Random Access Machine).

Toate aceste modele vor fi detaliate în continuare. De asemenea, dedicăm o secțiune consistentă standardului de comunicație prin mesaje MPI (Message Passing Interface).

2.1 Descrierea algoritmilor SPMD

Deși în stilul SPMD toate procesoarele execută același program, trebuie făcute două precizări importante:

- Execuția instrucțiunilor nu este sincronă, adică fiecare procesor execută instrucțiunile în ritmul său propriu.
- Nu toate procesoarele execută aceleași instrucțiuni, adică deși programele încărcate în memoriile procesoarelor sunt identice, totuși pot exista diferențe la execuție.

Prima precizare este justificată de însăși arhitectura MIMD, în care nu există sincronism. A doua necesită detaliere, deoarece pare în contradicție cu definiția modelului SPMD.

Adresa unui procesor. Un program se execută pe p procesoare ale unui calculator paralel; în general, valoarea lui p este mai mică decât numărul total de procesoare al calculatorului, în funcție de dorința utilizatorului, care ia în seamă dimensiunile problemei pe care o are de rezolvat și gradul de încărcare a calculatorului. (În general, noi nu vom face distincție între p și numărul total de procesoare, deoarece nu este important din punct de vedere algoritmic.) Procesoarele sunt numerotate de la 0 la $p - 1$, dar această numerotare nu este statică, ci se face în momentul lansării în execuție. Totuși, procesoarele pot să-și afle adresa proprie, printr-o primitivă a sistemului de operare sau de bibliotecă. În descrierea algoritmilor, vom preciza că, de exemplu, adresa este conținută de variabila k , spunând că algoritmul e valabil pentru procesorul P_k ; implicit, variabila conținând adresa procesorului curent va fi numită id . Procesoarele pot executa instrucțiuni diferite în funcție de adresa lor. De exemplu, dacă programul unic al procesoarelor arată astfel

1. **dacă** $id = 3$ **atunci**
 1. $a_{id} \leftarrow 1$
2. **altfel**
 1. $a_{id} \leftarrow 0$

atunci doar procesorul P_3 execută instrucțiunea $a_{id} \leftarrow 1$, în timp ce toate celelalte execută $a_{id} \leftarrow 0$. În consecință, în final, vectorul $a \in \mathbb{R}^p$ are elementele $a_k = 0$, $k \neq 3$, și $a_3 = 1$.

Variabile. Indici locali și indici globali. În descrierile algoritmilor, variabilele nu vor fi declarate; utilizarea unei variabile înseamnă implicit declararea ei; tipul variabilei (întreg, real etc.) va rezulta din context. În privința locului în care se află o variabilă în memorie și a modului în care ne referim la ea, trebuie să facem o distincție importantă.

- În cazul memoriei *comune*, fiecare variabilă se află în aceasta, iar referirea la variabile se face la fel ca în programele secvențiale. Micul exemplu de mai sus este elocvent, elementele vectorului a fiind referite cu indicii lor reali.
- În cazul memoriei *distribuite*, un procesor execută programul unic utilizând date din memoria sa *locală*. Așadar, o variabilă a unui program SPMD este multiplicată: fiecare procesor deține un exemplar, asupra căruia are control complet. Un procesor nu poate modifica variabile din memoria altui procesor. Prin urmare, în exemplul de mai sus, *fiecare* procesor deține doar un element al vectorului $a \in \mathbb{R}^p$, anume elementul al cărui indice este egal cu adresa procesorului. În mod natural, vom nota acest element tot cu a , însă indicii vor dispărea, iar programul de mai sus va fi transformat astfel:

1. **dacă** $id = 3$ **atunci**

1. $a \leftarrow 1$
2. **altfel**
 1. $a \leftarrow 0$

De unde în primul caz a era un vector aflat într-o singură zonă de memorie, acum vectorul este distribuit în memoriile locale ale procesoarelor.

Dacă pentru memoria comună lucrurile sunt foarte clare, în cazul memoriei distribuite avem de ales între cele două variante de mai sus, în *prezentarea* algoritmilor. Varianta didactică este de a folosi indici *globali*, ca și cum vectorii s-ar afla în întregime în memoriile procesoarelor, dar fiecare procesor efectuează numai adunarea sa; în acest caz scriem că P_k execută $a_k \leftarrow 0$, la fel ca pentru memorie comună. Dacă ținem seama că fiecare procesor are de fapt în memorie un singur element al vectorului, atunci scriem că P_k execută $a \leftarrow 0$; variabila este scalară; indicele a dispărut. Dacă vectorul a ar avea dimensiune mai mare decât p , atunci fiecare procesor ar avea mai multe elemente în memoria locală, stocate contiguu; accesul la ele se va face cu indici *locali*.

În general, vom folosi indicii globali, pentru că permit o descriere mai clară a algoritmilor. Când indicii sunt locali, adică se referă la partea dintr-un vector (sau matrice) aflată în memoria locală a unui procesor, vom menționa explicit acest lucru; acesta este modul de lucru în programele reale, de aceea vom da suficiente exemple de programe în ambele forme.

Pseudocod. Limbajul în care vor fi descriși algoritmii va fi mai mult sau mai puțin formal. În general vom folosi un pseudocod format din câteva instrucțiuni obișnuite în limbajele de programare de nivel înalt, așa cum am și făcut-o mai sus. Rezumăm aici aceste instrucțiuni, presupunând cititorul familiarizat cu semnificația lor.

- Atribuirea se va simboliza cu \leftarrow ; deci $a \leftarrow 1$ înseamnă că variabila a capătă valoarea 1 (a nu se confunda cu $=$, operatorul semnificând egalitatea). Semnul \leftrightarrow se va folosi pentru interschimbarea valorilor a două variabile.

- Instrucțiunea de decizie va avea forma **dacă** *expresie logică* **atunci** *instrucțiuni*; instrucțiunile se execută dacă condiția este adevărată. Poate exista și o ramură care se execută atunci când condiția este falsă, ramură care va urma cuvântului cheie **altfel**.

- Instrucțiunea **pentru** $i = vi : pas : vf$, *instrucțiuni* reprezintă ciclul cu număr cunoscut de pași; contorul i ia valori de la valoarea inițială vi la valoarea finală vf , la fiecare iterație adăugându-i-se valoarea pas . Pasul poate lipsi, fiind implicit 1, ca de exemplu în $i = 1 : n$.

- Instrucțiunea **cât timp** *expresie logică*, *instrucțiuni* este ciclul cu număr necunoscut de pași și test inițial. Instrucțiunile se execută atât timp cât expresia logică este adevărată.

- Instrucțiunea **stop** termină execuția programului; ea nu e neapărat necesară, dar ne permite să scurtăm descrierile algoritmilor.

Pentru descrierea expresiilor logice sau a celor aritmetice vom folosi semnele matematice uzuale.

Algoritmii vor fi numerotați, eventual denumiți după autorul lor sau cu un nume consacrat, și vor avea câteva cuvinte explicative la început, rezumând problema pe care o rezolvă și condițiile de funcționare, incluzând, dacă e cazul, adresa procesorului pentru care este valabil algoritmul.

O instrucțiune va ocupa, în general, o singură linie; pentru cele care pot ocupa mai multe linii, adică **dacă**, **pentru** și **cât timp**, liniile următoare primeia vor fi indentate (aliniată mai spre dreapta), prin această metodă putând fi determinat sfârșitul instrucțiunii; algoritmii vor avea descrieri suficient de scurte pentru a nu mai introduce un cuvânt cheie special pentru terminarea unei instrucțiuni de acest tip. Liniile algoritmului sunt numerotate după o regulă evidentă în exemplul de mai sus: liniile aliniată la același nivel sunt numerotate în ordine, începând de la 1, până când se întâlnește o instrucțiune de nivel superior (adică aliniată mai la stânga). Pentru referire se folosește numărul liniilor tuturor instrucțiunilor în interiorul cărora se află instrucțiunea curentă, la care se concatenează numărul liniei instrucțiunii. Așadar, în ultimul exempl de mai sus, instrucțiunea $a \leftarrow 0$ are numărul 2.1.

Descrierea paralelismului. O problemă oarecum delicată este modul în care se descriu acțiuni ce au loc în paralel. Modelul SPMD descrie implicit paralelismul; pentru facilitarea studiului unui algoritm, putem presupune că toate cele p procesoare încep simultan execuția programelor lor, ceea ce în realitate nu se întâmplă de obicei; a vedea însă în ce măsură operațiile se desfășoară într-adevăr în paralel ține de buna înțelegere a algoritmului; vom încerca să explicăm cât mai des ce anume se întâmplă în paralel și ce acțiuni sunt condiționate de altele (capitolul 4 va oferi unele instrumente suplimentare). În afara paralelismului implicit, vom folosi totuși și instrucțiunea **în paralel**, în două contexte destul de diferite:

1. Primul, coerent cu cele spuse până acum, este cel al unor acțiuni care se desfășoară în paralel pe un același procesor—de exemplu mai multe comunicații și, eventual, calcul.
2. Al doilea va fi cel al unor formulări la nivel global ale algoritmilor (mai puțin riguroase), descriindu-se acțiuni paralele desfășurate pe procesoare diferite; va fi utilizat mult mai rar, în special atunci când lucrurile vor fi suficient de clare.

În primul caz, instrucțiunea **în paralel** nu înseamnă obligatoriu execuție paralelă, ci numai potențialul pentru aceasta. Instrucțiunea va fi utilizată mai ales pentru noduri ale unui calculator cu memorie distribuită (vezi figura 1.7), acolo unde comunicația pe diverse canale poate decurge în bună măsură în paralel cu activitatea (de calcul) a procesorului (și unde este destul de greu de decis în ce măsură este vorba despre paralelism sau concurență). Execuția unei instrucțiuni **în paralel** se termină o dată cu terminarea execuției ultimei instrucțiuni componente.

Posibile ambiguități. Combinația instrucțiunilor **în paralel** și **pentru** poate da naștere la ambiguități: trebuie precizat dacă instrucțiunile din bucla **pentru** se

execută secvențial sau paralel. În secvența de instrucțiuni

1. **în paralel**
 1. **pentru** $i = 0 : n - 1$
 1. *instrucțiuni 1*
 2. *instrucțiuni 2*

instrucțiunea **pentru** 1.1 și grupul *instrucțiuni 2* din linia 1.2 se execută în paralel; în interiorul instrucțiunii **pentru**, iterațiile după contorul i se execută secvențial: se execută grupul *instrucțiuni 1* pentru $i = 0$, apoi aceleași instrucțiuni pentru $i = 1$ etc. În schimb, în secvența

1. **pentru** $i = 0 : n - 1$, **în paralel**
 1. *instrucțiuni(i)*

se execută în paralel grupul *instrucțiuni(i)* pentru toate valorile lui i între 0 și $n - 1$. Confuzii cu cazul precedent pot fi în următoarea situație:

1. **pentru** $i = 0 : n - 1$
 1. **în paralel**
 1. *instrucțiuni 1*
 2. *instrucțiuni 2*

când, pentru fiecare iterație în i a buclei **pentru**, instrucțiunile 1.1.1 și 1.1.2 se execută în paralel, dar iterațiile se execută secvențial.

2.2 Modelul de comunicație prin mesaje

Discutăm în această secțiune numai despre programarea calculatoarelor MIMD cu memorie distribuită. Deoarece fiecare procesor are memoria sa proprie, singura modalitate de comunicare între procesoare este transmiterea de mesaje. Operația de bază este cea în care un procesor sursă P_s transmite un mesaj M conținând date din memoria sa M_s unui procesor destinație P_a , care stochează datele în memoria sa M_a . În general, vom presupune că procesoarele sursă și destinație sunt vecine, adică există o legătură fizică între ele, pe care circulă mesajul. (Cazurile în care această regulă nu este respectată vor fi menționate explicit.)

Primitive de bază. În modelul de programare descris mai sus, numit *model de comunicație prin mesaje*, orice operație de comunicație poate fi descrisă utilizând doar două primitive de comunicație: una de transmitere executată de procesorul sursă, pe care o vom numi **send**, alta de recepție executată de destinatar și numită **recv**. Sintaxa lor este următoarea

```
send(date, dest)
recv(date, sursă)
```

Pentru ambele moduri de apel, *date* este un nume (o variabilă) din care să rezulte clar locul în memoria locală în care se află datele de transmis (sau în care se stochează datele recepționate), precum și lungimea mesajului. Al doilea parametru este un nume ce identifică vecinul cu care se comunică și deci portul care va efectua comunicația; pentru simplitate, atunci când va fi posibil, vom folosi drept nume al vecinului direcția pe care se găsește acesta: stânga, dreapta, sus, jos, vest, est, nord, sud; aceste denumiri sunt potrivite pentru inel, grilă sau tor; pentru hiper cub va fi folosită dimensiunea pe care se comunică, un număr între 0 și $d - 1$. Așadar,

send($A(i, :)$, dreapta)

înseamnă operația de transmitere a liniei i a matricei A (folosim notația în stil Matlab), de dimensiune presupusă cunoscută, către vecinul din dreapta al procesorului ce execută operația.

Vom presupune întotdeauna că datele formează o zonă contiguă în memorie, ceea ce nu este neapărat adevărat. Aceasta reprezintă o problemă practică de cea mai mare importanță în eficiența comunicației; necontiguitatea datelor în memorie implică fie transmiterea mai multor mesaje, ceea ce înseamnă de regulă un consum mai mare de timp, fie copierea datelor într-o zonă contiguă, deci folosind operații suplimentare. De aceea, modul în care sunt organizate matricele în memorie—pe linii sau pe coloane—nu poate fi neglijat.

Corectitudinea comunicației. Este evident că primitivele **send** și **recv** trebuie să apară în perechi, pe ansamblul procesoarelor; orice operație de transmisie a unui procesor este însoțită de una de recepție a unui vecin al său. Acesta este un prim criteriu de corectitudine a comunicației într-un algoritm. Iată un singur exemplu, foarte simplu; să presupunem că procesorul P_k trebuie să transmită un același mesaj M vecinilor săi pe un inel, P_{k-1} și P_{k+1} . Algoritmul are forma următoare

1. **dacă** $\text{id} = k$ **atunci**
 1. **send**(M , dreapta)
 2. **send**(M , stânga)
2. **altfel dacă** $\text{id} = (k - 1) \bmod p$ **atunci** **recv**(M , dreapta)
3. **altfel dacă** $\text{id} = (k + 1) \bmod p$ **atunci** **recv**(M , stânga)

Sunt în total două apeluri **send** și două **recv**, deci totalul este corect. Primul **send** este efectuat în instrucțiunea 1.1 de către P_k , spre dreapta, adică spre procesorul $P_{(k+1) \bmod p}$; la rândul său, acest procesor, în linia 3, execută **recv** din stânga, adică recepționează de la P_k . Se verifică în mod analog că a două pereche **send-recv** este corectă.

Terminarea locală a comunicației. Să presupunem că procesorul P_s transmite un mesaj procesorului P_a . Deoarece fiecare procesor are fluxul său de instrucțiuni, momentele în care P_s apelează primitiva **send**, iar P_a **recv**, sunt în general diferite. Transmiterea efectivă a mesajului se face doar după ce ambele procesoare au apelat primitivele respective. Se pune întrebarea: ce face procesorul care a apelat primul

primitiva sa de comunicație, din momentul apelului și până când celălalt procesor apelează primitiva pereche? Sunt două răspunsuri posibile:

1. așteaptă (fără să facă nimic); este cazul comunicației *blocante*;
2. poate executa alte operații; comunicația este *non-blocantă*.

Comunicație blocantă. Ilustrăm acest tip de comunicație în figurile 2.1 și 2.2, presupunând că procesorul sursă este primul care apelează primitiva de comunicație (cazul în care destinatarul face primul apel este obținut prin simetrie). În figura 2.1, la apelul primitivului **send** de către sursă zona de memorie numită generic "mesaj" conține datele care trebuie transmise (ceea ce este marcat printr-un dreptunghi hașurat); până în momentul apelării **recv** de către destinatar, procesorul sursă nu execută nici o operație, deci este blocat; după apelul **recv** are loc comunicația efectivă (marcată în figură cu o săgeată), după care zona "mesaj" devine liberă, adică poate fi utilizată în alte scopuri. Deoarece unul din procesoare așteaptă fără a face altceva, iar rutinele **send** și **recv** se termină doar după transferul mesajului, comunicația se numește blocantă și *sincronă*.

În figura 2.2, comunicația blocantă are loc printr-un *buffer*. Primitiva **send** mută mesajul din zona sa originală într-un buffer; după aceasta, execuția **send** se termină, zona de memorie a mesajului putând fi refolosită. Mesajul trece din buffer la destinație atunci când se execută **recv**; de comunicația efectivă se ocupă o rutină de bibliotecă sau a sistemului de operare (deci la nivel ierarhic inferior), care lucrează în paralel (sau concurent) cu programul principal, utilizând facilitățile arhitecturale pentru comunicație, de exemplu DMA în figura 1.7. Procesorul sursă este blocat doar în timpul în care mesajul este copiat în buffer. Observăm că terminarea rutinei **send** la sursă nu este condiționată de apelul **recv** la destinație.

Atât în varianta sincronă cât și în cea prin buffer, la terminarea execuției primitivului **send**, zona de memorie în care era memorat mesajul poate fi refolosită, iar la terminarea execuției **recv**, mesajul este recepționat în zona de memorie alocată acestui scop.

Comunicație non-blocantă. În modul non-blocant, primitivile **send** și **recv** se termină *imediat* după apel, fără ca zona de memorie a mesajului să fie liberă la **send** sau să conțină mesajul la **recv**. Primitivile au doar rolul de a iniția comunicația, transferul efectiv fiind realizat mai târziu la un nivel inferior. După terminarea execuției **send** sau **recv**, procesorul poate executa alte operații, în paralel sau concurent cu comunicația. Și în modul non-blocant comunicația se poate desfășura sincron sau prin buffer; prima variantă este ilustrată în figura 2.3; cealaltă este lăsată ca exercițiu cititorului.

Este evident că primitivile de comunicație non-blocante trebuie să fie însoțite de posibilitatea de a afla când se termină comunicația efectivă, mai exact de a afla când zona de memorie a mesajului poate fi refolosită de procesorul sursă sau poate fi utilizată de procesorul destinație. În acest scop se introduce primitiva **wait** (pe care o

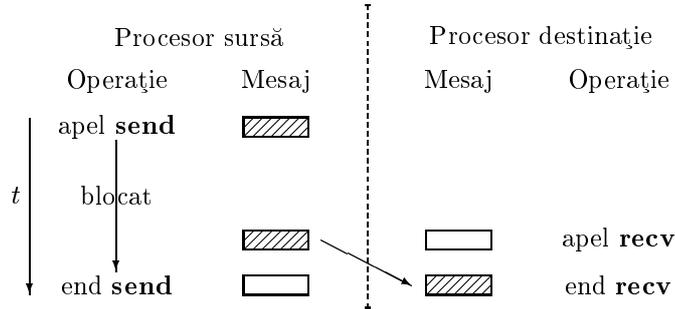


Figura 2.1: Comunicație blocantă între două procesoare, modul sincron.

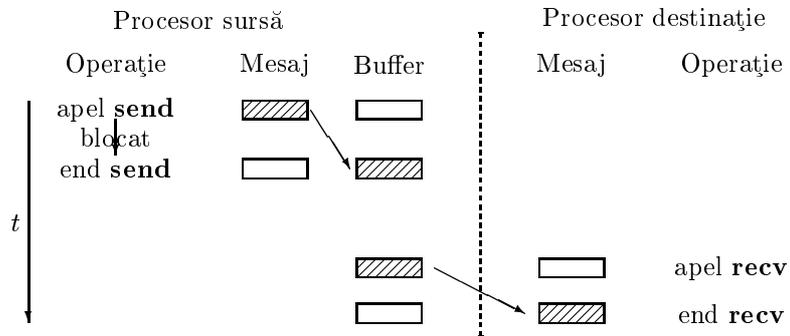


Figura 2.2: Comunicație blocantă prin buffer.

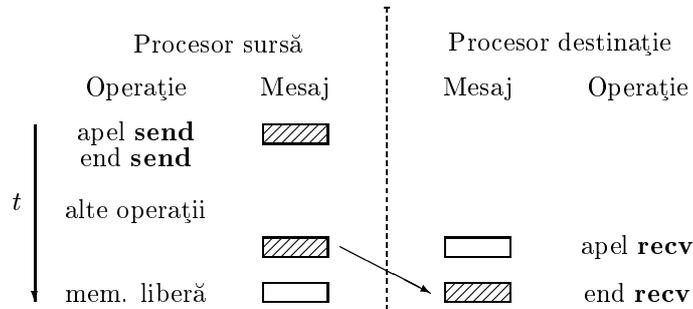


Figura 2.3: Comunicație non-blocantă sincronă.

vom folosi mai ales în traducerea **așteaptă**), de așteptare a terminării comunicației. Așadar, utilizarea tipică a rutinelor non-blocante este

1. **send**(*date*, *dest*)
2. execută operații care nu modifică *date*
3. **așteaptă** terminarea **send**
4. modifică *date*

În mod evident, dacă **așteaptă** este executată imediat după **send** sau **recv**, efectul este același ca în cazul primitivelor blocante.

Comparații. Comunicația non-blocantă este mai flexibilă și este mai folosită, în special în cazul tipic când există coprocesoare specializate pentru comunicație; ea are avantajul de a permite execuția altor operații, de către procesor, în timpul comunicației, și deci economisirea timpului unității centrale. Desigur, programarea concurentă a comunicației cu alte operații este uneori mai dificilă; pot apărea erori de gestionare a zonelor de date, cu manifestări diferite în funcție de desfășurarea în timp a comunicației.

Convenții de apel. Păstrăm un singur nume pentru rutinele de comunicație, indiferent de modul blocant sau non-blocant. Contextul va decide ce mod este folosit. De asemenea, vom folosi relativ rar primitiva **așteaptă**; în general, vom presupune, pentru claritatea algoritmilor, că operațiunile de comunicație nu se desfășoară în paralel cu calculele. Atunci când vom dori să evidențiem acest paralelism o vom face explicit. Să presupunem, de exemplu, că un procesor dintr-un hipercub cu $p = 2^d$ noduri transmite un mesaj tuturor vecinilor săi. În cazul comunicației blocante vom scrie (numai pentru procesorul sursă)

1. **pentru** $i = 0 : d - 1$
 1. **send**(*mesaj*, *i*)

(Reamintim că i reprezintă dimensiunea pe care se transmite.) Pentru comunicația non-blocantă, vom evidenția desfășurarea în paralel a mai multor transmisii, scriind

1. **pentru** $i = 0 : d - 1$, **în paralel**
 1. **send**(*mesaj*, *i*)
2. **așteaptă** terminarea tuturor operațiilor **send**

În general, pentru economie de spațiu, vom *ignora* apelul primitivei **așteaptă**, presupunându-l implicit.

Erori frecvente în comunicația blocantă. Să presupunem că, pe un inel de procesoare, se efectuează următoarea operație de comunicație colectivă: fiecare procesor transmite vecinului său din dreapta un mesaj. Echivalent, fiecare procesor trebuie să recepționeze un mesaj de la vecinul său din stânga. Notăm generic cu $d1$ și $d2$ variabilele locale alocate mesajului propriu, respectiv celui recepționat. În modul non-blocant, cu convențiile de mai sus, programul unui procesor poate fi scris în oricare din cele două forme de mai jos.

1. **în paralel**

1. **send**($d1$, dreapta)
2. **recv**($d2$, stânga)

1. **send**($d1$, dreapta)2. **recv**($d2$, stânga)3. **așteaptă** terminarea **send** și **recv**

Vom prefera scrierea din stânga.

În modul blocant, banala secvență

1. **send**($d1$, dreapta)
2. **recv**($d2$, stânga)

reprezintă un program incorect, dacă transmisia se face sincron. În acest caz, fiecare procesor începe prin a transmite și apoi așteaptă procesorul din dreapta să execute recepția, ceea ce nu se întâmplă: toate procesoarele se blochează. În schimb, dacă transmisia se face prin buffer, atunci operația **send** se termină local la transferul mesajului în buffer (independent de acțiunile procesorului destinație), deci recepția poate să înceapă.

O soluție simplă de a evita blocarea în cazul blocant sincron este ca parte din procesoare să înceapă prin a transmite, iar celelalte prin a recepționa; este interesant că oricum se face împărțirea procesoarelor în două mulțimi nevide, comunicația va funcționa. Cel mai rapid algoritm este următorul:

1. **dacă** id este par
 1. **send**($d1$, dreapta)
 2. **recv**($d2$, stânga)
2. **altfel**
 1. **recv**($d2$, stânga)
 2. **send**($d1$, dreapta)

Dacă numărul de procesoare p este par, atunci au loc efectiv două etape de comunicație; în prima transmit procesoarele cu adresă pară și recepționează cele cu adresă impară, în a doua lucrurile se petrec invers. În schimb, dacă p este impar, au loc trei etape; de exemplu, când $p = 5$, în prima etapă, transmisiile sunt $P_0 \rightarrow P_1$ și $P_2 \rightarrow P_3$; observăm că P_4 nu poate transmite, deoarece vecinul său din dreapta este P_0 , care este ocupat. În a doua etapă, au loc transmisiile $P_4 \rightarrow P_0$ și $P_1 \rightarrow P_2$; procesorul P_3 nu poate transmite, deoarece P_4 este ocupat. În fine, în ultima etapă, rămâne o singură transmisie de efectuat, indiferent de numărul de procesoare; în exemplul nostru este vorba despre $P_3 \rightarrow P_4$.

Pentru a evita complicații de tipul celor de mai sus, vom presupune că în modul blocant transmisia se face prin buffer.

Probleme

P 2.2.1 Descrieți în pseudocod următoarea operație de comunicație într-un hipercub: un (unic) procesor P_k transmite un mesaj M tuturor vecinilor săi pe hipercub. Aceeași problemă pentru tor.

P 2.2.2 Descrieți în pseudocod operația de comunicație în care, pe un hipercub, fiecare procesor transmite un același mesaj M tuturor vecinilor săi. Scrieți un algoritm general, precum și unul care să funcționeze în cazul blocant sincron. Idem pentru tor.

2.3 Standardul MPI

Message Passing Interface (MPI) este un standard descriind primitive de comunicație—în contextul unui model de programare SPMD cu comunicație prin mesaje—apărut la mijlocul anilor '90, ca urmare al unui efort comun al cercetătorilor din mediul academic și din industrie. La ora actuală, MPI este disponibil pe toate calculatoarele MIMD cu memorie distribuită, sub formă de bibliotecă de rutine optimizate pe calculatoarele respective. Scrierea unui program folosind rutinele de comunicație MPI asigură portabilitatea; eficiența comunicației este garantată de implementările performante disponibile ale rutinelor MPI. Vom prezenta aici doar câteva rutine, împreună cu concepte specifice MPI; cititorul interesat poate găsi o documentație completă la <http://www.mpi-forum.org>; în afara implementărilor pe calculatoare paralele, MPI este disponibil și pentru alte medii, de exemplu Linux sau Windows NT; astfel, scrierea și testarea unui program utilizând MPI se pot face extrem de ușor pe o rețea de câteva PC-uri—desigur, pe probleme de dimensiuni reduse. Standardul MPI este descris pentru limbajele FORTRAN și C. În cele ce urmează, vom folosi numai convențiile de apelare a funcțiilor C. Numim program MPI un program care utilizează MPI pentru comunicație.

Generalități. O rutină MPI se execută în cadrul unui grup de procesoare, numit *comunicator*. Inițial, există un singur comunicator, numit `MPI_COMM_WORLD`, conținând toate procesoarele pe care se execută programul. Pentru multe aplicații, acest comunicator este suficient; există rutine MPI care permit crearea dinamică a grupurilor de procesoare. În cadrul unui comunicator, procesoarele sunt numerotate de la 0 la $p - 1$, unde p este numărul de procesoare din comunicator. Un procesor poate afla numărul de procesoare din comunicator și adresa proprie (numită *rang* în MPI) apelând următoarele rutine:

```
int MPI_Comm_size(MPI_Comm com, int *p)
int MPI_Comm_rank(MPI_Comm com, int *my_id)
```

În ambele rutine, `MPI_Comm` este tipul predefinit pentru comunicatoare; variabila `com` este comunicatorul în care se află procesorul, în cazul cel mai general `MPI_COMM_WORLD`. În variabilele `p` și `my_id`, rutinele returnează valorile numărului de procesoare, respectiv a adresei procesorului apelant. Toate rutinele MPI (care reprezintă funcții C) întorc un întreg caracterizând succesul execuției: `MPI_SUCCES` înseamnă execuție corectă, alte valori reprezintă coduri de eroare predefinite.

Structura unui program MPI. Încadrându-se în modelul SPMD, un program MPI are forma obișnuită a unui program secvențial. Rutinele `MPI_Init` și `MPI_Finalize` trebuie apelate înainte, și respectiv după, orice alte rutine MPI (menționăm ca `MPI_Init` primește argumentele `argc` și `argv`, acestea având sau nu o

Comunicație	blocantă	non-blocantă
standard	MPI_Send MPI_Recv	MPI_Isend MPI_Irecv
prin buffer	MPI_Bsend	MPI_Ibsend
sincronă	MPI_Ssend	MPI_Issend

Tabelul 2.1: Numele rutinelor MPI, pentru diverse tipuri de comunicație.

semnificație în funcție de implementarea MPI). Aceasta este singura constrângere, în rest programatorul este liber să folosească rutinele MPI potrivite algoritmului pe care îl implementează.

Rutine de comunicație. În modelul MPI, oricare două procesoare își pot transmite mesaje. Așadar, topologia virtuală este graful complet conectat. Calea pe care o parcurge efectiv mesajul este transparentă pentru utilizator. Principalele rutine de comunicație sunt `MPI_Send` și `MPI_Recv`, a căror sintaxă completă este următoarea

```
MPI_Send(void *mes, int lung, MPI_Datatype tipdate, int dest,
         int eticheta, MPI_Comm com)
MPI_Recv(void *mes, int lung, MPI_Datatype tipdate, int sursa,
         int eticheta, MPI_Comm com, MPI_Status *stare)
```

La execuția `MPI_Send` de către un procesor, mesajul aflat în memoria locală la adresa `mes`, de lungime `lung` și având tipul `tipdate` (MPI are identificatori predefiniți pentru practic toate tipurile uzuale de date; alte tipuri se pot construi cu rutine speciale) este transmis procesorului `dest`. Mesajului îi este asociată o etichetă (tag), care îl personalizează. Transmisia mesajului se face în cadrul comunicatorului `com`, care trebuie să conțină și procesorul sursă, și cel destinație. Argumentele rutinei `MPI_Recv` au semnificații similare, cu următoarele observații: la execuție, se recepționează un mesaj de lungime cel mult `lung`, trunchiind eventual mesajul transmis; tipul de date poate diferi la sursă și destinație; în plus, după recepție, variabila `stare` dă informații suplimentare despre mesajul primit, ca de exemplu lungimea efectivă a mesajului.

Moduri de comunicație. MPI prevede rutine pentru toate tipurile de comunicație descrise în secțiunea 2.2; ne referim la comunicație blocantă, respectiv non-blocantă, precum și la modurile de transmisie prin buffer, respectiv sincron. Numele rutinelor sunt prezentate în tabelul 2.1, regulile de formare a numelor fiind evidente. Trebuie subliniat că modul de implementare a rutinelor de bază (standard) `MPI_Send`, `MPI_Recv` nu este precizat; de exemplu, este posibil ca mesajele scurte să fie transmise prin buffer, iar cele lungi sincron; de aceea programatorul trebuie să-și ia precauțiile necesare, discutate anterior. Pentru comunicația non-blocantă, rutina `MPI_Wait` așteaptă terminarea transmisiei sau recepției, iar `MPI_Test` verifică dacă transmisia sau recepția s-au terminat sau nu.

Exemple. Vom prezenta aici două exemple de programe MPI, implementând

două operații simple de comunicație pe inel discutate în secțiunea 2.2. Prima este transmiterea unui mesaj de către procesorul P_0 vecinilor săi, P_1 și P_{p-1} ; programul este prezentat în figura 2.4. A doua operație este transmiterea unui mesaj de către fiecare procesor vecinilor săi, realizată de programul din figura 2.5. În ambele cazuri, am prezentat un program MPI complet pentru a ilustra structura acestuia.

Programele sunt suficient de simple pentru a fi clare, așadar prezentăm numai scurte comentarii. În programul din figura 2.4, lungimea mesajului transmis este aleasă astfel încât să fie transmis și zeroul de terminare a unui șir de caractere; astfel, procesoarele care recepționează pot afișa mesajul fără a ști lungimea lui exactă (ci doar că aceasta este mai mică decât 50); o altă variantă ar putea fi transmiterea șirului de caractere fără terminator și aflarea lungimii la receptori prin intermediul variabilei `stare`. Eticheta asociată mesajului este `99`, aceeași la sursă și destinație. În cazul e.g. $p = 6$, programul afișează următorul text:

```
0: am trimis:  mesaj de la 0
1: am primit:  mesaj de la 0
5: am primit:  mesaj de la 0
```

Atragem atenția că ordinea în care aceste linii sunt afișate este mai mult sau mai puțin aleatoare. Chiar dacă P_0 apelează `printf` înaintea celorlalte două procesoare, nu există garanția că și afișarea se face înainte. Cu atât mai mult între P_1 și P_5 nu poate exista o ordine fixată de afișare. Mai multe execuții ale programului pot produce ordini diferite la afișare.

În programul din figura 2.5, observăm folosirea rutinei de transmisie prin buffer, care garantează corectitudinea programului; în general, pentru mesaje atât de scurte, `MPI_Send` ar fi putut înlocui cu succes `MPI_Bsend` în practică; recomandăm totuși utilizarea variantei de program a cărei corectitudine nu depinde de implementarea rutinelor MPI. Rutina `MPI_Buffer_attach` este destinată alocării unei zone de memorie folosită pentru buffere; utilizatorul alege dimensiunea acestei zone în funcție de o estimare a numărului de mesaje care se pot găsi simultan în buffer (și a lungimii lor); în exemplul nostru, lungimea 100 ar fi fost suficientă. Calculul adresei procesorului din stanga se face ca în program, adică adunând p , deoarece operația modulo din C diferă de cea matematică pentru numere negative.

Alte categorii de rutine. Standardul MPI conține numeroase alte rutine. Enumerăm aici câteva categorii interesante.

- **Comunicație globală.** MPI conține rutine pentru diverse operații de comunicație implicând toate procesoarele dintr-un grup (comunicator). Este vorba despre sincronizare, difuzare, distribuție etc. Aceste operații vor fi discutate în capitolul următor, de aceea ne mărginim să semnalăm existența lor.
- **Alte operații globale,** ca de exemplu calculul maximului unor valori distribuite tuturor procesoarelor unui grup, sau calculul sumei prefixelor. Vom prezenta în capitolele 4 și 5 algoritmi pentru aceste operații.

```
#include <stdlib.h>
#include <mpi.h>

void main(int argc, char ** argv)
{
    int my_id, p;
    MPI_Status stare;
    char mesaj[100];

    MPI_Init(&argc, &argv);           // initializare MPI
    MPI_Comm_rank(MPI_COMM_WORLD, &my_id); // afla adresa proprie
    MPI_Comm_size(MPI_COMM_WORLD, &p);    // afla numar de procesoare

    if ( p < 3 )
    { if ( my_id == 0 )                // doar proc. 0 afiseaza mesaj eroare
      printf("Sunt necesare cel putin 3 procesoare !\n");
      exit(0);
    }

    if ( my_id == 0 )                  // proc. 0 transmite mesajul
    { sprintf(mesaj, "mesaj de la %d", my_id);
      MPI_Send(mesaj, strlen(mesaj)+1, MPI_CHAR, 1, 99, MPI_COMM_WORLD);
      MPI_Send(mesaj, strlen(mesaj)+1, MPI_CHAR, p-1, 99, MPI_COMM_WORLD);
      printf("%d: am trimis:  %s\n", my_id, mesaj);
    }
    else if ( my_id == 1 || my_id == p-1 ) // proc. 1 si p-1 receptioneaza
    {
      MPI_Recv(mesaj, 50, MPI_CHAR, 0, 99, MPI_COMM_WORLD, &stare);
      printf("%d: am primit:  %s\n", my_id, mesaj);
    }
    MPI_Finalize();                    // terminare MPI
}
```

Figura 2.4: Program MPI pentru operația: procesorul P_0 transmite un mesaj procesoarelor P_1 și P_{p-1} .

```
#include <stdlib.h>
#include <mpi.h>

void main(int argc, char ** argv)
{
    int my_id, p, dreapta, stanga;
    MPI_Status stare;
    char msend[100], mrecv[100], buffer[1000];

    MPI_Init(&argc, &argv);           // initializare MPI
    MPI_Comm_rank(MPI_COMM_WORLD, &my_id); // afla adresa proprie
    MPI_Comm_size(MPI_COMM_WORLD, &p);    // afla numar de procesoare

    dreapta = (my_id + 1) % p;          // adresa procesor dreapta
    stanga = (my_id - 1 + p) % p;       // adresa procesor stanga
    sprintf(msend, "mesaj de la %d", my_id);
    MPI_Buffer_attach(buffer, 1000);
    MPI_Bsend(msend, strlen(msend)+1, MPI_CHAR, dreapta, 99, MPI_COMM_WORLD);
    MPI_Recv(mrecv, 100, MPI_CHAR, stanga, 99, MPI_COMM_WORLD, &stare);
    printf("%d: am primit:  %s\n", my_id, mrecv);

    MPI_Finalize();                    // terminare MPI
}
```

Figura 2.5: Program MPI pentru operația: fiecare procesor transmite un mesaj procesorului din dreapta sa, pe un inel.

- **Topologii virtuale.** Pentru creșterea performanțelor unui program MPI pe un anumit calculator, programatorul poate descrie o topologie virtuală care să corespundă—exact, sau în sensul unei dilatări minime—topologiei reale a calculatorului. Rutinele MPI dedicate descrierii topologiei asociază procesoarelor adresele corespunzătoare poziției în topologia reală. MPI permite descrierea unei topologii arbitrare, modelată printr-un graf, așa cum am văzut în capitolul 1. Rutine speciale sunt dedicate topologiilor ce pot fi descrise printr-un produs cartezian, aici incluzându-se toate topologiile curente: inel, tor (2D sau 3D), hipercub. Programele astfel scrise sunt în continuare portabile, chiar dacă eventual au performanțe mai slabe pe arhitecturi diferite de cea pentru care au fost concepute.
- **Procese.** În versiunea 1 a MPI, fiecare procesor executa un singur proces, reprezentând programul MPI; așadar, procesele erau create static, de către sistemul de operare, la încărcarea programului în procesoare. Versiunea 2 a adus ca principală noutate posibilitatea de creare dinamică a proceselor; se permite astfel o mai mare flexibilitate în descrierea paralelismului, precum și execuția concurentă a mai multor procese pe același procesor.

2.4 Modelul PRAM

Modelul teoretic al arhitecturii MIMD cu memorie partajată este PRAM (Parallel Random Access Machine), o extensie a modelului RAM pentru arhitectura secvențială; pentru amănunte, vezi [16]. În acest model, într-o unitate de timp, fiecare procesor poate să citească o dată din memoria comună, să execute o operație și să scrie într-o locație de memorie. Se trece deci peste implicațiile pe care le au structura și viteza de acces ale memoriei, presupunându-se că timpul de acces la memorie este constant. De asemenea, se presupune că procesoarele încep sincron execuția fiecărei instrucțiuni a lor; aceasta este o abatere de la modelul SPMD, dar permite descrierea simplă a algoritmilor; vom reveni asupra sincronizării cu exemple.

Modelul PRAM poate avea mai multe variante, în funcție de posibilitatea accesului simultan la o locație de memorie:

- EREW PRAM (Exclusive Read Exclusive Write)—în care o locație poate fi citită sau scrisă de către un singur procesor la un moment dat.
- CREW PRAM (Concurrent Read Exclusive Write)—în care, în plus, este permisă citirea unei locații de către mai multe procesoare simultan.
- CRCW PRAM (Concurrent Read Concurrent Write)—în care, în plus, mai multe procesoare pot scrie simultan în aceeași locație; aceasta se poate realiza în mai multe feluri:
 - toate procesoarele scriu aceeași valoare—modelul COMMON.

- procesoarele au asociate o prioritate; scrie procesorul cu prioritate maximă; acesta este modelul PRIORITY.
- procesorul care scrie este ales aleator—modelul ARBITRARY.

Din punct de vedere teoretic, modelul este foarte puternic, el permițând descrierea simplă a algoritmilor și evidențierea paralelismului în sine al problemelor, făcând abstracție de detaliile arhitecturii țintă (cea pe care se urmărește implementarea). Trecerea de la forma PRAM a algoritmilor la implementări practice este relativ simplă, de aceea arhitectura PRAM va fi folosită intens în această lucrare.

Exemplul 1: suma a doi vectori. Să presupunem că vrem să efectuăm adunarea $z \leftarrow x + y$, unde variabilele x, y, z sunt vectori de lungime p . Această operație este perfect paralelă, fiecare sumă a două elemente putând fi efectuată independent de câte un procesor. În general, vom scrie că procesorul P_k efectuează operația $z_k \leftarrow x_k + y_k$. Vectorii se află în memoria comună, deci fiecare procesor trebuie să acceseze locația adecvată, așa că indicii sunt obligatorii. Cu totul excepțional, putem scrie

1. **pentru** $k = 0 : p - 1$, **în paralel**
 1. $z_k \leftarrow x_k + y_k$

Aici este vorba despre o descriere globală a paralelismului, în care nu se precizează adunarea efectuată de fiecare procesor. Acest mod mai vag de descriere este util doar când evidențiem paralelismul unor operații; va fi utilizat foarte rar în această lucrare.

Exemplul 2. Pentru ilustrarea sincronizării presupuse de modelul PRAM, considerăm calculul de către două procesoare a sumei a patru numere, $u_k, k = 0 : 3$. Un mod rapid de a face calculele este ca procesorul P_0 să efectueze suma $u_0 + u_1$ și, în paralel, P_1 suma $u_2 + u_3$, după care P_0 să adune cele două sume parțiale. Notând $id = i$, programul corespunzător acestei idei este

1. $v_i \leftarrow u_{2i} + u_{2i+1}$
2. **dacă** $i = 0$
 1. $s \leftarrow v_0 + v_1$

Algoritmii scriși pentru modelul PRAM presupun (tacit) că procesoarele execută simultan instrucțiunile, adică lucrează sincron; în cazul nostru, avem garanția că ambele sume parțiale sunt calculate în instrucțiunea 1, înainte de a se trece mai departe la instrucțiunea 2. În lipsa sincroniei, este posibil ca P_0 să utilizeze v_1 înainte ca P_1 să pună în această variabilă suma $u_3 + u_4$. Astfel de sincronizări vor fi presupuse implicit întotdeauna, chiar dacă ele nu sunt strict necesare.

Observăm că, în ambele exemple de mai sus, modelul este de fapt EREW PRAM, deoarece nu există citiri sau scrieri simultane în aceeași locație de memorie. Vom da mai târziu exemple de algoritmi și pentru alte modele.

Probleme

P 2.4.1 Se poate simula cu un PRAM un calculator MIMD cu memorie distribuită ?

P 2.4.2 Scrieți un program de care calculează suma numerelor u_k , $k = 0 : 3$, fără a utiliza variabile suplimentare.

Capitolul 3

Algoritmi de comunicație

În calculatoarele paralele, comunicația este o problemă ce are o cu totul altă natură decât la calculatoarele secvențiale; la acestea din urmă, comunicația se face îndeosebi între un proces (unic, de obicei) ce ocupă unitatea centrală și exterior. La primele, comunicația are loc între procese care se execută pe procesoare diferite și care colaborează la rezolvarea aceleiași probleme; în același timp există posibilitatea dialogului între procesoare și exterior, în mod analog cu cazul secvențial, dar de acest tip de comunicație nu ne vom ocupa. Vom studia deci comunicația doar atunci când aceasta este parte intrinsecă a algoritmului paralel, transferul de informație (cooperarea între procesoare) fiind absolut necesar rezolvării problemei și oferind cheia eficienței, adică a execuției rapide pe un calculator paralel.

În acest capitol, ne vom ocupa de cazul arhitecturilor MIMD cu memorie distribuită, la care există linii de comunicație directe între anumite procesoare, după o topologie specifică fiecărui calculator. Pentru arhitecturile MIMD cu memorie partajată, comunicația se face prin intermediul memoriei comune; problemele corespunzătoare acestui caz sunt mai vechi și au apărut o dată cu calculatoarele secvențiale cu sistem de operare multitasking; există multă literatură pe această temă, de exemplu [?].

Vom detalia o serie de algoritmi, cu precădere în domeniul comunicațiilor globale—i.e. atunci când mai multe procesoare participă la realizarea comunicației. După cum veți vedea, tehnicile folosite depind fundamental nu numai de topologia rețelei de interconectare, ceea ce era clar de la început, dar și de modelul de comunicație adoptat; de aceea alegerea unui model adecvat calculatorului țintă și sistemului său de operare este o etapă extrem de importantă în proiectarea sau alegerea unui algoritm de comunicație. Vom încerca, în acest capitol, să oferim o serie de idei utile nu numai atunci când este vorba exclusiv de comunicație, ci și în alte cazuri, în probleme în care comunicația este totuși preponderentă iar calculele se pot greșa în mod natural pe structura algoritmului de comunicație.

3.1 Modele de comunicație

Așa cum am sugerat în figura 1.7, fiecare nod al unei arhitecturi MIMD cu memorie distribuită are unități specializate în comunicație. În cele ce urmează nu ne vom preocupa prea mult de ceea ce se întâmplă la nivelul acestor unități: modul în care sunt legate fizic unități aparținând unor procesoare vecine, protocolul de comunicație folosit, etc. Ne vom mulțumi cu faptul că există legături între anumite procesoare, după topologii de genul celor descrise în capitolul anterior și vom prezenta doar modele generale simple, dar care sunt suficiente în proiectarea algoritmilor.

Toate modelele prezentate în această secțiune se referă la arhitecturi cu memorie distribuită. Principalele aspecte abordate sunt: modul de transmitere a mesajului între procesoare, modul de folosire a canalelor de către un procesor, modelarea timpului necesar pentru comunicarea unui mesaj de lungime dată. Reamintim că o clasificare a comunicației după modul de execuție local—blocant sau non-blocant—a fost discutată în capitolul anterior.

Modul de transmitere a unui mesaj între procesoare. Problema de bază în comunicație este de a transmite un mesaj, de la un procesor oarecare la un altul. În general, pentru a ajunge de la procesorul sursă la cel destinație, mesajul trebuie să parcurgă mai multe legături între procesoare vecine (în afara cazului simplu când sursa și destinația sunt vecine). Nu ne punem deocamdată problema cum se găsește o cale între cele două procesoare—o presupunem cunoscută—ci doar a modului în care este transmis mesajul de către procesoarele de pe această cale.

1. O primă posibilitate are la bază doar comunicația între fiecare procesor și vecinii săi direcți, pe liniile de legătură cu aceștia. Fiecare mesaj este transmis integral de către posesorul său și recepționat de către vecin într-o zonă de memorie înainte de a fi prelucrat sau, eventual, trimis mai departe către destinatar; acest mod de transmitere este numit *store and forward* (memorează și apoi trimite mai departe) sau *comutare de pachete*. Ilustrăm în figura 3.1 această tehnică; succesiunea etapelor de comunicație este reprezentată de sus în jos; mesajul M —având un antet h , care conține adresa destinatarului și alte informații referitoare la mesaj, de exemplu lungimea acestuia—este transmis de către sursă primului procesor de pe calea spre destinație; după recepție, acesta îl retransmite celui de-al doilea, ș.a.m.d. până la destinație. O analogie simplă este cea cu un lanț de oameni care transmite din mână în mână o cărămidă.

2. Un alt model posibil, frecvent întâlnit în calculatoarele actuale este cel cu *comutare de circuit* (*circuit switched*), în care întâi se stabilește fizic o cale între sursă și destinatar, indiferent de distanța dintre aceștia, și apoi se transmite direct mesajul destinatarului, după cum se sugerează în figura 3.2. Deci, dacă în prima variantă mesajul se transmitea în mai multe etape identice, în fiecare între două procesoare vecine, fiecare procesor memorând pe rând mesajul, acum mesajul va fi transmis în două etape principial diferite; întâi este trimis antetul h (ca un deschizător de drum) care parcurge calea între sursă și destinație, rezervând legăturile de pe această cale; apoi, procesorul sursă este informat de rezervarea întregii căi, și deci că poate începe

transmisia; apoi, mesajul este transmis direct destinatarului, pe calea acum rezervată, ca și cum acesta ar fi vecin cu expeditorul; procesoarele intermediare nu memorează mesajul. Această tehnică este utilizată în telefonie, de exemplu.

3. În fine, o a treia variantă posibilă, asemănătoare ca idee cu comutarea de circuit, este cea *wormhole*. Mesajul M este împărțit în mai multe pachete (*flit*), notate M^i în figura 3.3, care au fiecare lungimea egală cu cea a cozii (zonei de memorie, bufferului) atașate unui canal; primul flit, pe care-l vom nota h conține și adresa destinatarului. Antetul mesajului avansează spre destinație de câte ori aceasta este posibil, iar restul mesajului îl urmează; ultimul flit, M^f , eliberează coada canalului pe care o ocupa. Mesajul se deplasează deci ca un vierme, care nu poate fi "tăiat" de alte mesaje, deci ocupă mai multe canale la un moment dat. Avantajul acestui mod de lucru este că un flit nu trebuie stocat în întregime de un procesor; acesta, dacă canalul de ieșire este liber, introduce doar o mică întârziere—de câțiva octeți—cea ce conduce la o deplasare rapidă a trupului viermelui, dacă primul flit, capul, se deplasează el însuși rapid. Dacă, în schimb, canalul de ieșire este ocupat, fiecare flit este memorat de către procesorul la care a ajuns și transmis mai departe atunci când este posibil. Față de comutarea de circuit, se elimină deci etapa de confirmare a rezervării căii de comunicație. Din nou comunicația este transparentă pentru procesoarele de pe traseu, ele nememorând întregul mesaj.

Modul de utilizare a canalelor. Un procesor are mai multe canale de comunicație; el poate folosi un singur canal la un moment dat, sau pe toate în paralel; primul model se numește *1-port*, al doilea *n-port* sau *multiport* (sau *link-bound*, adică numărul maxim de comunicații este mărginit numai de numărul de canale fizice ale procesorului). Cele două modele sunt extreme: în primul canalele sunt folosite pur secvențial, în al doilea în perfect paralelism. În modelul 1-port se presupune de obicei că un procesor poate emite și recepționa simultan pe câte un canal, nu neapărat același. Modelul multiport este idealizat; de exemplu, pentru un procesor cu arhitectură de genul celei din figura 1.7, la care canalele pot transmite sau recepționa în paralel, modelul nu este totuși perfect adecvat: programarea canalelor (cu adresa mesajului, lungimea sa) se face în mod secvențial, de-abia apoi ele lucrând în paralel, dar concurând în accesul la memorie, ceea ce poate eventual întârzia comunicația. Cel mai realist, dar și mai greu de manipulat algoritmic, ar fi un model în care k porturi pot fi folosite simultan, unde $k > 1$, dar mai mic decât numărul de canale al unui procesor.

Canalele vor fi presupuse bidirecționale, ceea ce este cazul pentru toate calculatoarele actuale; dacă transmisia se poate face în ambele sensuri simultan, modelul va fi numit *full duplex*; dacă se poate folosi un singur sens la un moment dat, modelul este *half duplex*. În modelul full duplex se presupune că transmisia simultană în cele două sensuri ale unei legături se face cu aceeași viteză ca în cazul folosirii unui singur sens, ceea ce nu este întotdeauna adevărat.

Ca unitate de măsură pentru lungimea mesajelor se alege de obicei octetul. Totuși, deoarece în această lucrare ne vom ocupa în special de probleme de calcul numeric, în care datele sunt numere reale reprezentate întotdeauna în virgulă mobilă la imple-

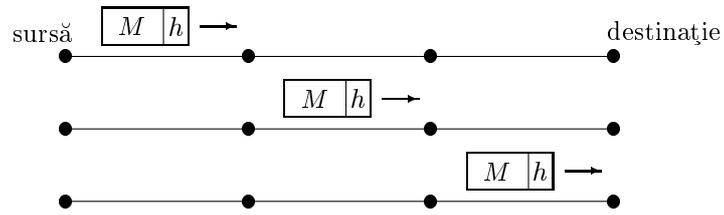


Figura 3.1: Comunicație store and forward (comutare de pachete).

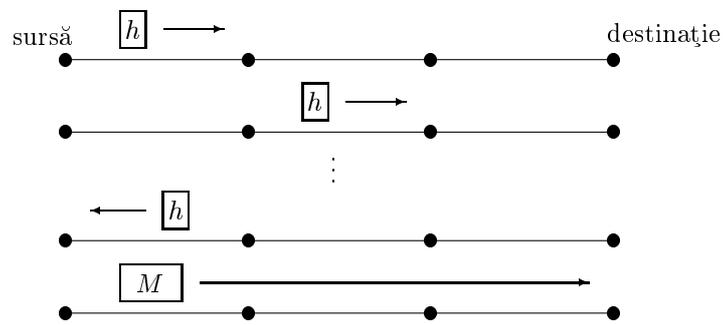


Figura 3.2: Comunicație cu comutare de circuit.

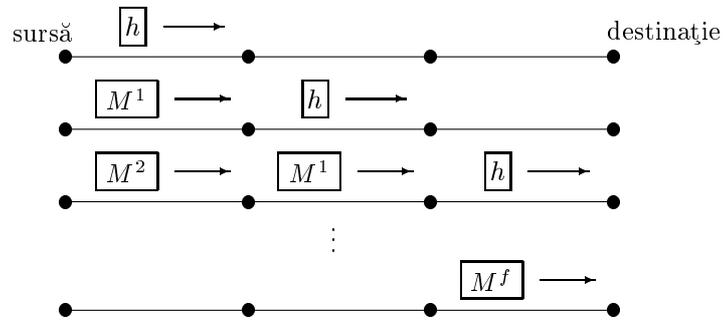


Figura 3.3: Comunicație wormhole.

mentare, vom presupune în general dimensiunea mesajelor măsurată nu în octeți, ci în elemente în virgulă mobilă, făcând abstracție de numărul de octeți necesar pentru reprezentarea unui element (care este, pentru multe calculatoare și limbaje de programare, de 4 octeți pentru reprezentarea în simplă precizie și de 8 octeți pentru reprezentarea în dublă precizie). De asemenea, în alte probleme, vom alege ca unitate de măsură lungimea obiectului fundamental pentru acel algoritm (șir de caractere, de exemplu, la sortare).

Modelarea timpului de comunicație. În ce privește timpul necesar pentru transmiterea unui mesaj între două procesoare vecine (implicit, deci, presupunem modul de transmisie store and forward) sunt mai multe modele posibile, fiecare adecvat unor anume tipuri de mesaje și de arhitecturi. Cel mai simplu model este acela în care durata de transmisie a unui mesaj este *constantă*, indiferent de lungimea mesajului; modelul este adecvat cazului în care mesajele se transmit în pachete de lungime fixă—dacă mesajul efectiv este mai scurt, el se completează până la lungimea cerută cu elemente ne semnificative. De asemenea, modelul este bun dacă mesajele sunt relativ scurte (zeci de octeți, să zicem).

Mai aproape de realitate este modelul *proporțional*, în care timpul de transmisie a unui mesaj de lungime m este

$$t_c(m) = m\beta,$$

unde β este timpul cât durează transmiterea unui element; valoarea $1/\beta$ este denumită *rată de transmisie*.

În fine, cel mai potrivit model, și deci cel mai folosit, este cel *liniar*:

$$t_c(m) = \sigma + m\beta, \tag{3.1}$$

în care σ este timpul de inițializare (start-up) a transmisiei; acesta apare, de exemplu, datorită necesității de a programa canalul de comunicație, a procedurilor de începere a comunicației, a confirmării recepției etc. În cazul în care mesajele au lungime mare, timpul de inițializare poate fi neglijat, folosindu-se modelul proporțional; de asemenea, dacă σ este mai mic sau de același ordin de mărime cu β , atunci σ poate fi neglijat chiar pentru mesaje relativ scurte. Pe de altă parte, dacă σ este relativ mare și m mic, modelul liniar este apropiat de cel cu timp constant. Se remarcă așadar că modelul liniar este o generalizare a modelelor precedente. Pentru majoritatea calculatoarelor actuale, relația între constantele din (3.1) este $\sigma \gg \beta$; ca ordin de mărime, valori ca $\sigma = 50\mu\text{s}$ și $\beta = 0.03\mu\text{s}$ sunt caracteristice (amintim că valoarea timpului de inițializare σ depinde și de suportul software pentru comunicație).

În cazul modurilor de transmisie prin comutare de circuit sau wormhole, timpul de transmisie are forma

$$t_c(m) = \sigma + d\tau + m\beta. \tag{3.2}$$

În acest caz procesoarele nu sunt neapărat vecine, ci se află la distanță d ; termenul τ reflectă timpul necesar antetului să străbată legătura între doi vecini. Se observă

că factorul ce înmulțește dimensiunea m a mesajului nu depinde de lungimea căii de comunicație; cum τ este relativ mic, aceasta reflectă faptul că modul wormhole este, în general, mai eficient decât cel store and forward. Expresia (3.2) este valabilă în ipoteza că nici una dintre legăturile de pe calea dintre sursă și destinație nu este ocupată.

Alte observații. Alegerea modelului pentru timpul de transmisie a unui mesaj este foarte importantă; există probleme de comunicație care se rezolvă optim, pe anumite arhitecturi, în mod destul de diferit în funcție de modelul adoptat. În această lucrare vom folosi îndeosebi modelul liniar și transmisia store and forward; de asemenea, modelul multiport va fi cel mai des adoptat. În ipoteza utilizării unui alt model, acesta va fi specificat în mod explicit acolo unde este cazul.

3.2 Operații de comunicație globală

În algoritmi de calcul numeric (în care datele se reprezintă sub formă matriceală) dar nu numai în aceștia, sunt deseori întâlnite operații de comunicare globală, în care sunt implicate mai mult de două procesoare. Unele dintre operații apar frecvent, de aceea au beneficiat de un dublu interes. În primul rând, ele au fost studiate intens pentru a se găsi algoritmi care să le execute rapid, pe fiecare dintre arhitecturile curente; amintim aici doar câteva din articolele de pionierat, care au propus algoritmi valabili și astăzi, ale căror autori sunt Saad & Schultz [23], Johnsson & Ho [19], Stout & Wagar [25], Bertsekas *et al* [3]. În al doilea rând, aceste operații de comunicație au fost standardizate, în prezent toate bibliotecile de comunicație oferind implementări ale lor.

Cele mai uzuale operații de schimbare a datelor între procesoare sunt enumerate mai jos. Operațiile sunt ilustrate în figura 3.4.

1. Transmiterea unui mesaj de la un procesor la altul (one-to-one). Includem această operație de comunicație în categoria celor globale pentru că, în general, procesoarele sursă și destinație nu sunt vecine, deci mesajul transmis circulă pe la mai multe procesoare în drumul său.
2. Transmiterea aceluiași mesaj de către un procesor tuturor celorlalte procesoare, denumită *difuzare* (broadcast, one-to-all).
3. Transmiterea câte unui mesaj de către fiecare procesor tuturor celorlalte procesoare, denumită *difuzare generală* (all-to-all broadcast, multinode broadcast, gossiping¹). Așadar este vorba de câte o operație de difuzare efectuată de fiecare nod.
4. Transmiterea de către un procesor a câte unui mesaj pentru fiecare din celelalte procesoare, denumită *difuzare personalizată* sau *distribuție* (distributing,

¹Dacă vă pasionează engleza și nu cunoașteți cuvântul, dați fuga la dicționar.

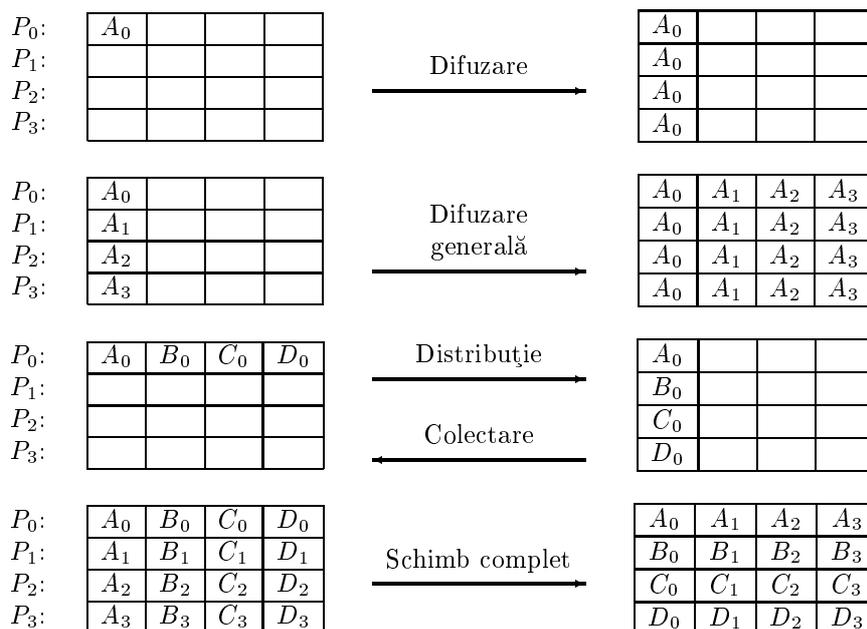


Figura 3.4: Operații de comunicație colectivă, exemplificate pentru $p = 4$ procesoare participante. Un procesor deține datele de pe orizontala sa: în stânga cele dinainte de comunicație, în dreapta cele de după.

scattering, personalized one-to-all). Operația inversă, de transmitere de către fiecare procesor a câte unui mesaj pentru un anumit procesor, denumită *colectare* (collecting, gathering), este similară difuzării personalizate și nu trebuie tratată separat; inversând sensul căilor pe care se efectuează transmisia într-un algoritm de difuzare personalizată, se obține un algoritm de colectare. Ambele operații pot apărea în cazul în care unul dintre procesoare are o poziție privilegiată: de exemplu, el împarte seturi de date celorlalte procesoare, printr-o distribuție, și apoi culege rezultatele printr-o colectare.

5. Transmiterea de către fiecare procesor a câte unui mesaj pentru fiecare dintre celelalte procesoare, denumită *schimb complet*, *multidistribuție* sau *transpunere* (multiscattering, personalized all-to-all, complete exchange, data transposition). Această operație consistă din difuzări personalizate efectuate de către fiecare procesor. Dacă fiecare procesor deține inițial o linie a unei matrice $p \times p$, el va avea în final coloana cu același indice, operația efectuată fiind așadar transpunerea matricei respective.

O altă operație extrem de utilă ce poate fi inclusă în categoria comunicațiilor

globale este **sincronizarea** (termenul englezesc echivalent este **barrier**). Ea este echivalentă cu o barieră în calea execuției programului unui procesor, care se oprește la întâlnirea ei și nu trece mai departe decât împreună cu toate celelalte procesoare. Este evident că, dacă este utilizată, primitiva **sincronizare** trebuie să apară în programul fiecărui procesor; după execuția ei, toate procesoarele vor începe, în paralel, execuția instrucțiunii următoare. Operația de sincronizare poate apărea atât în calculatoarele MIMD cu memorie comună, cât și în cele cu memorie distribuită. Multe calculatoare oferă suport hardware pentru execuția ei, reducând semnificativ timpul necesar sincronizării. Să remarcăm însă că este vorba de timpul efectiv, scurs între momentul în care ultimul procesor apelează local primitiva de sincronizare și momentul în care procesoarele încep execuția primei instrucțiuni ulterioare sincronizării. În sfârșit, să mai observăm că în arhitecturile MIMD cu memorie distribuită nu este neapărat necesară o sincronizare globală a procesoarelor; sincronizarea se face implicit prin transmisia-recepția de mesaje.

3.3 Comunicații globale

În această secțiune vor fi prezentați algoritmi pentru operațiile de comunicație globală. Deși în bibliotecile de comunicație actuale există rutine implementând toate operațiile de comunicație globală, algoritmi descriși în continuare sunt interesanți deoarece oferă idei ce pot fi aplicate în algoritmi de calcul; în plus, studiul algoritmilor de comunicație permite o mai bună cunoaștere a modului în care se poate profita de o topologie dată. Topologiile considerate vor fi inelul, hipercubul și torul de dimensiuni egale, deci $\sqrt{p} \times \sqrt{p}$; pentru grilă vom particulariza unii dintre algoritmi pentru tor; de asemenea, vom ține seama că unii algoritmi de la inel se pot generaliza ușor pentru tor (deși există algoritmi specifici mai performanți). Să precizăm de la bun început că, de exemplu, orice algoritm valabil pentru inel va funcționa și pe un hipercub, deoarece un inel poate fi scufundat în hipercub; totuși, pentru hipercub se pot găsi algoritmi specifici mult mai performanți, care folosesc conectivitatea mai bună a acestuia.

Se va presupune că lungimea tuturor mesajelor este egală, anume m , ceea ce, în general, este cazul în practică, atunci când există mai multe mesaje de transmis, adică în difuzarea generală, distribuție și schimb complet. În algoritmi care implică situația specială a unui procesor, adică difuzarea și distribuția, se va presupune P_0 acest procesor (sau P_{00} pentru tor), ceea ce nu este o restrângere, arhitecturile considerate fiind simetrice, după cum am văzut în primul capitol. Modelul de comunicație va fi cel *store and forward*, pentru comutarea de circuit rezervând o secțiune specială.

Timpii de execuție ai algoritmilor vor fi notați cu doi până la cinci indici. Primul este numărul de ordine al tipului de comunicație, după numerotarea din secțiunea precedentă. Al doilea reprezintă tipul arhitecturii: R pentru inel (ring), T pentru tor, G pentru grilă, H pentru hipercub. Al treilea reprezintă felul comunicației: H pentru half duplex, F pentru full duplex; acest indice va fi omis dacă timpii de execuție sunt aceiași pentru ambele modele. Al patrulea precizează numărul canalelor

de comunicație ce pot transmite simultan; valorile pot fi 1 pentru 1-port și * pentru multiport (acesta din urmă fiind cazul implicit). Al cincilea se referă la modelul timpului de comunicație: L pentru timp liniar, C pentru timp constant, P pentru timp proporțional; dacă lipsește, se presupune modelul liniar. Primii doi indici vor fi prezenți întotdeauna; ultimii trei, având valori diferite, vor putea apare în orice combinație. Deci $T_{2,H,F}$ reprezintă timpul de execuție pentru difuzare, pe un hiper-cub, model full duplex, multiport, timp liniar; la fel și $T_{2,H,F,C}$, dar pentru timp constant; nu apare nici o ambiguitate, chiar dacă indicele C se află în a patra poziție, și nu într-a cincea.

Un mesaj va fi notate M_a , unde a este procesorul destinație; indicele a poate lipsi, dacă destinația este evidentă din context. Dacă mesajul este împărțit în pachete, notăm M_a^i pachetul cu numărul de ordine i . Dacă este necesară precizarea procesorului sursă, atunci un pachet va fi notat cu $M_{s,a}^i$, unde s este procesorul sursă. În loc de s și a pot apărea liste de procesoare, semnificând concatenarea unor mesaje provenind de la, sau cu destinația procesoarele din lista respectivă. Prin aceste notații lășăm a se înțelege că un mesaj poate fi separat în mai multe submesaje, sau, la fel de bine, mai multe mesaje pot fi concatenate; aceste operații vor fi presupuse a se efectua în timp neglijabil.

3.3.1 Comunicație de la un procesor la altul

Deși această operație poate părea banală, este totuși interesantă tratarea ei pentru a pune în evidență mijloacele prin care se poate crește viteza de comunicație a unui mesaj, pentru fiecare din modelele de timp expuse. Motivația includerii ei între operațiile de comunicație globală, deși doar două procesoare sunt beneficiarele efective ale comunicației, se justifică prin aceea că multe alte procesoare, eventual toate celelalte, participă la transmisia între sursă și destinație.

Presupunem că procesorul sursă P_s transmite un mesaj procesorului destinație P_a și că aceasta este singura operație de comunicație efectuată pe ansamblul procesoarelor.

Comunicație punct-la-punct pe inel

Comunicație pe calea cea mai scurtă. În cazul inelului, distanța între sursă și destinație este $\delta = \min(|a - s|, p - |a - s|)$. Cea mai simplă metodă este de a transmite mesajul pe cea mai scurtă cale dintre P_s și P_a , din procesor în procesor, după algoritmul următor.

ALGORITM 3.1 (transmiterea mesajului M de la P_s la P_a , pe calea cea mai scurtă, pe inel)

1. $\delta = \min(|a - s|, p - |a - s|)$
2. **dacă** $dist(id, s) + dist(id, a) = \delta$ **atunci**
 1. **dacă** $id = s$

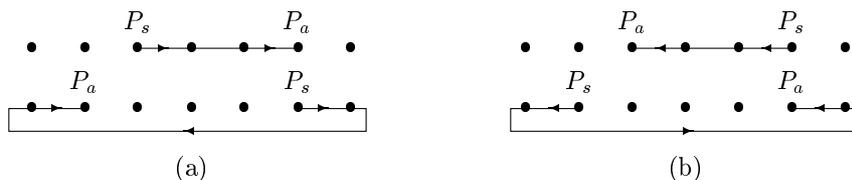


Figura 3.5: Pozițiile relative posibile pentru sursă și destinație, pe un inel: (a) P_a la dreapta lui P_s (b) P_a la stânga lui P_s .

1. **dacă** $(a - s) \bmod p = \delta$ **atunci** $\text{send}(M, \text{dreapta})$
2. **altfel** $\text{send}(M, \text{stânga})$
2. **altfel dacă** $(id - s) \bmod p \leq \delta$ **atunci**
 1. $\text{recv}(M, \text{stânga})$
 2. **dacă** $id \neq a$ **atunci** $\text{send}(M, \text{dreapta})$
3. **altfel dacă** $\delta < p/2$
 1. $\text{recv}(M, \text{dreapta})$
 2. **dacă** $id \neq a$ **atunci** $\text{send}(M, \text{stânga})$

Așadar, un procesor determină în primul rând, în instrucțiunile 1 și 2, dacă se află sau nu pe calea cea mai scurtă între sursă și destinație. Apoi, procesoarele de pe calea cea mai scurtă trebuie să determine în ce sens se efectuează transmisia; condiția din instrucțiunea 2.2 arată când un astfel de procesor se găsește la dreapta lui P_s ; dacă $id \geq s$, atunci $(id - s) \bmod p = id - s \leq \delta$; dacă $id < s$, atunci $(id - s) \bmod p = p - (s - id) < \delta$, pentru că $s - id > p - \delta$; aceste două cazuri sunt cele din figura 3.5a; în celelalte două cazuri, P_a se găsește în stânga lui P_s , deci și procesoarele de pe calea cea mai scurtă sunt în stânga. Să mai observăm că atunci când ambele căi între P_s și P_a au aceeași lungime, atunci este preferat convențional sensul de la stânga la dreapta. În cazul în care sensul de transmisie este de la stânga la dreapta, comunicația este inițiată de P_s în instrucțiunea 2.1.1. Celelalte procesoare participante realizează retransmiterea mesajului, în instrucțiunile 2.2.1 și 2.2.2; evident, procesorul destinație doar recepționează. Circulația în sensul celălalt este realizată de 2.1.2, 2.3.1 și 2.3.2 (testul din 2.3 interzice procesoarelor să transmită de la dreapta spre stânga în cazul în care există două căi de lungime minimă, ceea ce e posibil doar când p este par și $\delta = p/2$). Poate că am detaliat excesiv acest algoritm, însă am făcut-o pentru a arăta drumul de la simpla enunțare a ideii, până la o implementare; de aici înainte vor fi prezentați algoritmi doar în cazurile mai complicate.

Timpul necesar execuției algoritmului 3.1 este:

$$T_{1,R,C} = \delta; \quad T_{1,R,P} = \delta m \beta; \quad T_{1,R,L} = \delta(\sigma + m\beta).$$

Pentru modelul timp constant nu se poate face mai bine, mesajul trebuind să parcurgă oricum δ canale. Pentru celelalte două modele, două idei conduc la creșterea vitezei comunicației.

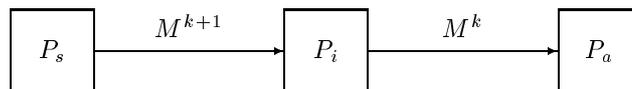


Figura 3.6: Transmisie pipeline între procesoarele P_s și P_a , prin intermediul procesorului P_i .

Utilizarea mai multor căi. În primul rând, mesajul se poate împărți în două și transmite pe cele două căi distincte dintre P_s și P_a , una de lungime δ , cealaltă de lungime $p - \delta$; se străbate deci inelul în cele două sensuri, pornind de la P_s . Dacă mesajul se împarte într-un pachet de dimensiune γm trimis pe calea cea mai scurtă și unul de lungime $(1 - \gamma)m$ trimis pe calea mai lungă, cu $\gamma \in (0, 1)$, timpul optim este

$$T'_{1,R,L} = \delta(\sigma + \gamma m\beta),$$

atunci când ambele pachete ajung simultan la P_a . Se obține

$$\gamma = \frac{p - \delta}{p} + \frac{p - 2\delta}{p} \frac{\sigma}{m\beta}.$$

Pentru modelul timp proporțional se ia $\sigma = 0$ în relațiile de mai sus.

Transmisie pipeline. În al doilea rând, pentru a fi folosite cât mai multe canale de comunicație în paralel (spre deosebire de simpla transmitere a mesajului, în care o singură legătură este ocupată la un moment dat), un mesaj transmis pe o cale poate fi împărțit în mai multe pachete (submesaje), care sunt transmise succesiv, unul câte unul. Un exemplu este prezentat în figura 3.6, în care între procesoarele P_s și P_a se află un singur alt procesor P_i ; după ce P_i primește primul pachet din mesaj, va trimite spre P_a acest pachet, în paralel cu recepția celui de-al doilea pachet de la P_s . De fapt se efectuează o transmisie *pipeline* a mesajului. În cazul unor mesaje foarte lungi, timpul pentru transmiterea mesajului între P_s și P_a poate fi aproximat cu cel din cazul în care procesoarele ar fi fost vecine; diagrama ocupării legăturilor dintre procesoare din figura 3.7 arată că, în transmisia pipeline, legăturile sunt aproape permanent ocupate, mai puțin unii mici timpi morți inițiali și/sau finali. Se poate observa că aceasta este o abordare software a modului de comunicare wormhole, cu diferența, totuși, că acolo pachetele nu erau memorate în nodurile intermediare decât dacă era absolut necesar.

Numărul optim de pachete. Aplicăm metoda pipeline pentru transmisia mesajului de la sursă la destinație, pe calea cea mai scurtă; considerăm mesajul împărțit în ν pachete notate $M^0, M^1, \dots, M^{\nu-1}$, fiecare de lungime m/ν ; presupunem, pentru simplitate, că ν divide m . Se constată că primul pachet ajunge la P_a după δ pași, iar ultimul după $\nu - 1 + \delta$ pași; deci timpul total este, pentru modelul timp proporțional:

$$T''_{1,R,P} = (\nu - 1 + \delta) \frac{m}{\nu} \beta = m\beta + (\nu - 1) \frac{m}{\nu} \beta.$$

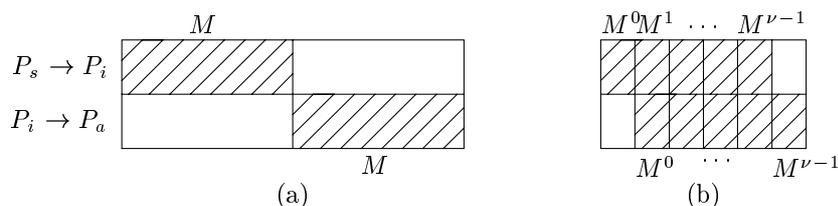


Figura 3.7: Ocuparea legăturilor pentru comunicarea din figura anterioară, în cazurile: (a) transmisie integrală a mesajului; (b) transmisie pipeline. Dreptunghiurile hașurate reprezintă timpii necesari transmiterii unui pachet; zonele albe reprezintă timpii morți.

Pentru a minimiza acest timp se observă că ν trebuie să fie cât mai mare; deci, $\nu = m$, adică fiecare pachet este format dintr-un singur octet (element de mesaj).

Pentru modelul timp liniar, timpul de comunicație este

$$T''_{1,R,L} = (\nu - 1 + \delta)\left(\sigma + \frac{m}{\nu}\beta\right).$$

Numărul de pachete care minimizează acest timp este

$$\nu_{\text{opt}} = \sqrt{\frac{m\beta(\delta - 1)}{\sigma}},$$

iar timpul minim corespunzător este

$$T''_{1,R,L\text{min}} = \left(\sqrt{m\beta} + \sqrt{(\delta - 1)\sigma}\right)^2. \quad (3.3)$$

Acest minim se poate încadra astfel:

$$(\delta - 1)\sigma + m\beta < T''_{1,R,L\text{min}} \leq 2(\delta - 1)\sigma + 2m\beta$$

și el se apropie de termenul din stânga pe măsură ce m crește.

Pe de altă parte, pentru a se atinge optimul, este necesar să fie îndeplinită condiția $1 \leq \nu_{\text{opt}} \leq m$, ceea ce se întâmplă dacă:

$$\frac{\sigma}{m} \leq (\delta - 1)\beta \leq m\sigma.$$

Pentru valori suficient de mari ale lungimii mesajului m , această condiție este satisfăcută. Parametrul σ variază, pentru diverse calculatoare, de la valori de ordinul lui β până la valori de sute de ori mai mari; în plus, distanța între sursă și destinație δ , fiind mai mică decât diametrul grafului reprezentat de topologia rețelei de interconectare, are o valoare de ordinul zecilor, să spunem; deci m trebuie să fie de cel puțin câteva zeci, ceea ce este cazul aproape întotdeauna; practic, nu se folosesc mesaje prea scurte

datorită ineficienței comunicației; oricum, se poate concluziona că tehnica pipeline nu se poate aplica eficient mesajelor scurte.

Pipeline pe mai multe căi. Metoda pipeline se poate aplica tuturor algoritmilor de comunicație de la un procesor la altul, combinată cu folosirea tuturor căilor de comunicație *arc-disjuncte* (care nu au nici un arc comun) între sursă și destinație; căile pot avea noduri comune, deoarece lucrăm în modelul multiport, deci se poate transmite permanent în mod simultan prin acele noduri. De exemplu, în cazul inelului, transmitând pipeline în ambele sensuri, de la sursă spre destinație, timpul de transmisie are o expresie similară cu (3.3), anume

$$T_{1,R,L}''' = \left(\sqrt{\frac{m}{2}}\beta + \sqrt{(p-\delta-1)\sigma} \right)^2. \quad (3.4)$$

Acest timp este de aproape două ori mai mic decât $T_{1,R,L}''_{\min}$, atunci când m este suficient de mare. El poate fi încă micșorat dacă se folosesc pachete de dimensiuni diferite pe cele două căi, dar nu cu mult. Expresia (3.4) se poate imediat generaliza pentru un graf oarecare.

PROPOZIȚIA 3.1 *Presupunem că, pe un graf oarecare, există c căi arc-disjuncte între sursă și destinație, cea mai lungă dintre ele având lungime h . Atunci timpul de comunicație a unui mesaj de lungime m , folosind toate aceste căi în regim pipeline, este*

$$T_1 = \left(\sqrt{\frac{m}{c}}\beta + \sqrt{(h-1)\sigma} \right)^2. \quad (3.5)$$

Se poate observa că sunt valabile relațiile $c \leq \Delta$, unde Δ este gradul grafului, i.e. numărul de canale al unui procesor, și $h \geq \delta$.

Concluzii. Pentru modelul timp constant cea mai eficientă transmitere a unui mesaj are loc pe calea cea mai scurtă între sursă și destinație. Pentru modelul timp liniar (și m suficient de mare), trebuie găsite cât mai multe căi arc-disjuncte între sursă și destinație, lungimea lor fiind mai puțin importantă; transmisia eficientă pe aceste căi are loc în modul pipeline. În continuare, pentru tor și hiper-cub, nu vom mai preciza timpii obținuți în urma aplicării tehnicii pipeline, ei rezultând direct din formula (3.5), ci ne vom rezuma la prezentarea căilor pe care se poate face comunicația.

Mai observăm că timpul (3.5) este același pentru modelele half duplex și full duplex, canalele fiind folosite întotdeauna într-un singur sens.

Comunicație punct-la-punct pe tor

Pe un tor există cel mult patru căi arc-disjuncte între sursă și destinație, dat fiind gradul acestui graf. În această secțiune vom presupune că \sqrt{p} este par, deci diametrul torului este chiar $D = \sqrt{p}$. Să considerăm întâi cazul în care procesoarele P_s și P_a se află pe linii și coloane diferite. O construcție imediată este cea din figura 3.8a, în care cele patru căi sunt construite folosind pentru fiecare o singură linie și o singură

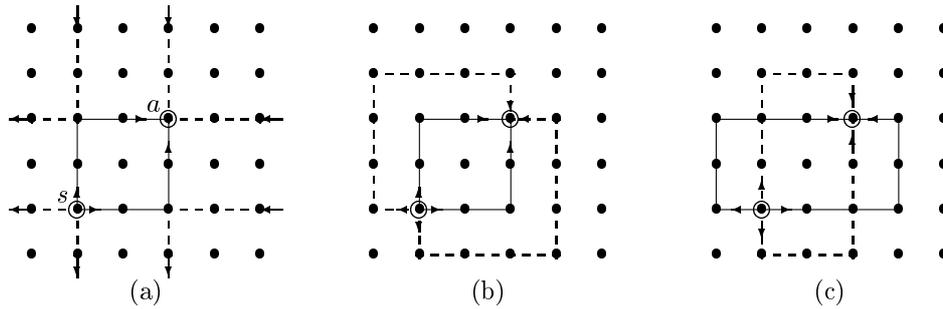


Figura 3.8: Trei variante de construcție a patru căi arc-disjuncte între două procesoare aflate pe linii și coloane diferite, într-un tor.

coloană; pentru prima se merge întâi la dreapta, pe linia lui P_s , până se întâlnește coloana lui P_a , apoi în sus; pe a doua, întâi în sus, apoi la dreapta; celelalte două sunt construite folosind direcțiile stânga și jos, în cele două ordini posibile. Se observă că două căi, cele reprezentate cu linie continuă, au lungimea egală cu distanța δ dintre noduri. Celelalte, reprezentate cu linie întreruptă, au în general lungimea mai mare, anume $2\sqrt{p} - \delta$. Aceste patru căi pot fi privite ca o generalizare a celor două căi posibile pe un inel.

Această variantă are în avantajul său simplitatea, dar e posibil ca lungimea căilor să fie prea mare; de asemenea, nu este posibilă o generalizare pentru grilă. Se pot imagina două alte soluții, care încearcă minimizarea lungimii maxime a celor patru căi. În figura 3.8b, celor două căi de lungime δ li se adaugă două de lungime $\delta + 4$, construite în jurul dreptunghiului format de primele două căi. Se poate demonstra optimalitatea celei de-a treia variante, cea din figura 3.8c, în care cele patru căi au toate lungime $\delta + 2$; ele sunt construite folosind una din cele două căi cele mai scurte între fiecare dintre vecinii din aceeași poziție ai lui P_s , respectiv P_a (vecinul din dreapta al lui P_s cu cel din dreapta al lui P_a , etc.); toate aceste căi au lungime δ , la care se adaugă cele două arce între P_s și vecin, respectiv între P_a și vecinul corespunzător.

Când cele două procesoare, P_s și P_a se află pe aceeași linie sau aceeași coloană, în oricare dintre variantele din figura 3.8, două căi se suprapun, deci trebuie găsită o altă soluție. O posibilitate este cea din figura 3.9, în care fiecare cale are lungime $\delta + 4$; căile au fost desenate câte două, pentru mai multă claritate; închipuiți-vă că torul este ceva mai mare, altfel există o variantă mai bună, dar particulară, P_a fiind la aceeași distanță de P_s pe orizontală, și prin stânga și prin dreapta.

Și această variantă cade atunci când procesoarele sunt vecine. Vă rămâne să găsiți d-voastră soluția (problema 3.3.4).

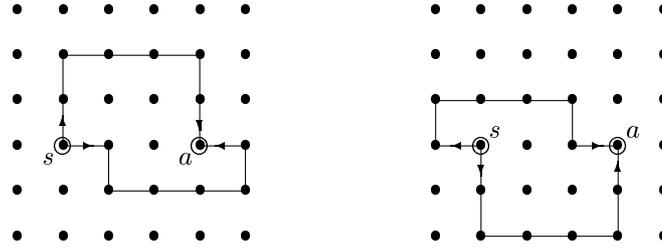


Figura 3.9: Patru căi arc-disjuncte între procesoare aflate pe aceeași linie.

Comunicație punct-la-punct pe hipercub

Pe un hipercub, distanța dintre P_s și P_a este egală cu numărul de biți de 1 din $s \oplus a$, adică distanța Hamming dintre s și a , pe care o notăm cu $H(s, a) = \delta$. O cale de lungime minimă între P_s și P_a se obține complementând unul câte unul biții diferiți din reprezentările binare ale adreselor s și a . Numim *principală* calea obținută prin complementarea biților diferiți în ordine *de la dreapta la stânga* (de la cel mai puțin la cel mai semnificativ).

Pentru simplitate, descriem comunicarea între procesoarele 0 și $2^d - 1 = 0 \dots 01 \dots 11$ (număr ce are 1 pe ultimele δ poziții). Generalizarea la adrese s și a oarecare este imediată, vezi problema P3.3.5.

Comunicație între colțuri opuse. Vom începe cu cazul particular al comunicării între două colțuri opuse ale hipercubului, deci $s = 0$ și $a = 2^d - 1$. În acest caz, e.g. pentru \mathcal{H}_3 , calea principală este formată din nodurile 000, 001, 011, 111, și are lungimea $d = 3$. Pentru obținerea altor $d - 1$ căi, complementăm biții adresei $s = 0$ tot de la dreapta la stânga, dar începând cu biții din pozițiile 1, \dots , $d - 1$; după bitul $d - 1$, se continuă ciclic, adică se complementează bitul 0. De exemplu, pentru \mathcal{H}_3 , cele două noi căi sunt 000, 010, 110, 111 și 000, 100, 101, 111, după cum se vede și în figura 3.10a. Aceste căi pot fi obținute din adresele nodurilor de pe calea principală, prin rotiri (deplasări ciclice) la stânga cu 1, \dots , $d - 1$ poziții. Numerotăm căile după poziția primului bit care este modificat (deci calea principală are numărul 0). Introducem notația $\text{Ro}(c, l)$, care semnifică numărul obținut prin rotația numărului binar c cu l biți la stânga. Așadar, dacă c este adresa unui nod de pe calea principală, atunci $\text{Ro}(c, l)$ este adresa nodului de pe calea l , aflat la aceeași distanță de nodul 0 ca și c .

PROPOZIȚIA 3.2 *Cele d căi astfel construite între nodurile 0 și $2^d - 1$ sunt arc-disjuncte.*

Demonstrație. Notăm cu C_k mulțimea nodurilor aflate la distanță k de nodul 0 (și deci având k biți de 1). Prin construcție, fiecare cale are exact un arc între C_k și

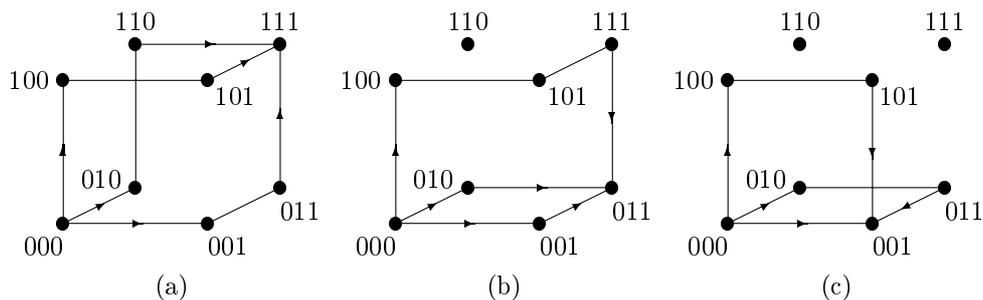


Figura 3.10: Cele trei căi arc-disjuncte într-un hipercube de dimensiune 3, între nodurile: (a) 000 și 111; (b) 000 și 011; (c) 000 și 001.

C_{k+1} ; aceste d arce sunt în dimensiuni diferite, mai precis arcul de pe calea l se află în dimensiunea $(k + l) \bmod d$. Deci căile sunt arc-disjuncte. ■

Cazul general. Trecem acum la cazul în care $\delta < d$. Calea principală între 0 și $2^\delta - 1$ are lungime δ . Rotind acum numai grupul primilor δ biți ale adreselor nodurilor de pe calea principală, obținem δ căi arc-disjuncte de lungime δ ; sunt, de fapt, căile între colțurile opuse ale unui hipercube de dimensiune δ , construite toate cu arce din acel hipercube. Mai multe căi de lungime δ nu se pot obține, deoarece modificând biți din poziții mai mari decât $\delta - 1$ ne îndepărtăm de nodul $2^\delta - 1$. Scopul nostru este însă să construim d căi, maximul ce se poate spera într-un hipercube, acesta fiind gradul grafului. Celelalte $d - \delta$ căi se construiesc astfel; în calea l , primul arc se obține prin negarea bitului l , deci se află în dimensiunea l ; este deci arcul între 0 și 2^l ; apoi se construiește calea principală de la nodul 2^l la nodul $2^\delta - 1$. În figura 3.10b,c sunt ilustrate astfel de construcții pentru noduri aflate la distanță 2, respectiv 1, în \mathcal{H}_3 .

PROPOZIȚIA 3.3 *Cele d căi între nodurile 0 și $2^\delta - 1$ sunt arc-disjuncte.*

Demonstrație. Primele δ căi sunt arc-disjuncte conform Propoziției 3.2. Celelalte—în afara primului și ultimului arc, aflate pe dimensiunea l pentru calea l —unesc noduri ale căror adrese au bitul din poziția l egal cu 1, deci nu se pot suprapune cu nici una dintre celelalte căi, care conțin noduri având acest bit egal cu 0. ■

Putem sintetiza toate aceste considerații într-un algoritm. Notăm cu y numărul format din ultimii δ biți din adresa procesorului pentru care scriem algoritmul (adică id), și cu x pe cel format din primii $d - \delta$ biți; cu u_k notăm numărul reprezentat pe δ biți, având ultimii k biți egali cu 1, iar pe ceilalți 0 (de exemplu, dacă $\delta = 4$, $u_3 = 0111$). În întreg algoritmul, pachetul M^i , de lungime m/d , va circula pe calea i , folosind numerotarea descrisă mai sus. Algoritmul este valabil și pentru cazul $\delta = d$.

ALGORITM 3.2 (comunicație de la P_0 la $P_{2^\delta - 1}$ a mesajului M , pe d căi arc-disjuncte, în hipercubele \mathcal{H}_d)

1. $b \leftarrow$ numărul de biți de 1 din y
2. $e \leftarrow$ numărul de biți de 1 din x
3. **dacă** $id = 0$ **atunci** {procesorul sursă}
 1. **pentru** $i = 0 : d - 1$, **în paralel** {emite pe toate căile}
 1. **send**(M^i, i)
4. **altfel dacă** $id = 2^\delta - 1$ **atunci** {procesorul destinație}
 1. **pentru** $i = 0 : \delta - 1$, **în paralel** {primele δ căi}
 1. **recv**($M^i, (i - 1) \bmod \delta$)
 2. **pentru** $i = \delta : d - 1$, **în paralel** {celelalte $d - \delta$ căi}
 1. **recv**(M^i, i)
5. **altfel dacă** $e = 0$ **atunci** {proc. poate fi într-una din primele δ căi}
 1. **dacă** $y = \text{Ro}(u_k, i)$ **atunci** {implicit $0 < k < \delta, 0 \leq i < \delta$ }
 1. **recv**($M^i, (k + i - 1) \bmod \delta$)
 2. **send**($M^i, (k + i) \bmod \delta$)
6. **altfel dacă** $e = 1$ **atunci** {proc. poate fi într-una din ultimele $d - \delta$ căi}
 1. $i \leftarrow$ poziția bitului de 1 din x , în id
 2. **dacă** $y = u_k$ **atunci**
 1. **dacă** $k = 0$ **atunci** **recv**(M^i, i) {de la P_0 }
 2. **altfel** **recv**($M^i, k - 1$)
 3. **dacă** $k = \delta$ **atunci** **send**(M^i, i) {spre $P_{2^\delta - 1}$ }
 4. **altfel** **send**(M^i, k)

În loc de suplimentarea explicațiilor oferite în comentarii, vă propunem ca exercițiu aplicarea algoritmului pentru fiecare din nodurile hipercubului de dimensiune 3 pentru comunicația între nodurile 000 și 011. Se obțin într-adevăr căile din figura 3.10b ?

3.3.2 Difuzare

Dacă în comunicarea între două procesoare nu am avut nevoie de a preciza sensul în care se transmit mesajele pe o legătură, acesta fiind evident, acum și mai departe sensul devine important. De aceea e necesar să considerăm grafurile orientate, între fiecare două noduri vecine existând două arce, câte unul în fiecare sens. Vom da în continuare câteva noi definiții referitoare la arbori.

Arbori. Reamintim că un arbore este un graf conex fără bucle. Într-un arbore, nodurile unite de un arc se numesc *tată* și *fiu*, sensul arcului fiind de la tată spre fiu; un tată poate avea mai mulți fii; dacă are cel mult doi, arborele se numește *binar*; în schimb, un fiu are întotdeauna un singur tată. Nodul care nu are tată se numește *rădăcină*; nodurile care nu au fii se numesc *frunze*. *Înălțimea* (numită și *adâncime*) unui arbore este distanța cea mai mare între rădăcină și o frunză. *Nivelul* unui nod este distanța dintre acesta și rădăcină; așadar, înălțimea grafului este nivelul maxim al unui nod; nivelul unui arc va fi considerat egal cu nivelul nodului tată.

Un *arbore de acoperire* (spanning tree) al unui graf X este un arbore care conține toate nodurile din X și arce din X (este arbore și graf parțial al lui X). O *familie* de

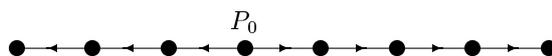


Figura 3.11: Un arbore de acoperire de înălțime $\lfloor p/2 \rfloor$, pe un inel.

arbori de acoperire este formată din mai mulți arbori de acoperire cu aceeași rădăcină (care pot fi arc-disjuncti sau nu).

Difuzare pe arbori de acoperire. Legătura dintre difuzare și arborii de acoperire este imediată. Procesorul care începe difuzarea se află în rădăcină. El își transmite mesajul către fii săi, care transmit la rândul lor fiilor lor, ș.a.m.d până când mesajul ajunge la toate frunzele; evident, fiecare nod păstrează pentru sine o copie a mesajului. Pe o familie de arbori de acoperire, mesajul este împărțit în mai multe pachete, fiecare difuzat pe unul dintre arbori. Pe un singur arbore de acoperire se poate utiliza întotdeauna tehnica pipeline. Cum, la un moment dat, este posibil ca toate arcele unui astfel de arbore să fie ocupate, se poate face pipeline eficient pe o familie de arbori de acoperire numai dacă arborii sunt arc-disjuncti. Pentru timpul de difuzare este deci valabil următorul rezultat analog Propoziției 3.1.

PROPOZIȚIA 3.4 *Utilizând transmisia pipeline cu dimensiune optimă a pachetelor pentru difuzare pe o familie de c arbori de acoperire arc-disjuncti, având înălțimea maximă h , se obține un timp de difuzare care respectă formula (3.5).*

Atragem atenția că este vorba despre modelul full duplex, deoarece ambele arce orientate corespunzând unui canal de comunicație pot aparține unor arbori (distincti) din familie. Fiindcă arborii de acoperire din familie au aceeași rădăcină (nodul care difuzează), avem $c \leq \Delta$. În plus, înălțimea arborilor este limitată inferior de diametrul grafului, i.e. $h \geq D$. Deci, în cel mai bun caz, timpul de difuzare este egal cu timpul de comunicație punct-la-punct între două noduri situate la distanță maximă (D) în graf. Așadar, în cele ce urmează vom evidenția numărul de arbori de acoperire și înălțimea lor maximă, timpul de difuzare rezultând din (3.5).

Difuzare pe inel

Un arbore de înălțime minimă. Presupunem fără pierderea generalității că P_0 este procesorul care difuzează. Un arbore de acoperire evident util pentru difuzare este cel din figura 3.11, având înălțimea egală cu diametrul inelului, adică $\lfloor p/2 \rfloor$. Algoritmul de difuzare este intuitiv: P_0 trimite pe ambele canale mesajul, vecinii săi îl retransmit, etc. Orice alt arbore de acoperire conduce la un timp mai mare de difuzare, deoarece înălțimea arborelui este mai mare decât diametrul inelului.

Doi arbori arc-disjuncti. Putem încerca construirea unei familii formată din doi arbori de acoperire arc-disjuncti; mai mult de atât nu se poate, deoarece din rădăcină pornesc doar două arce. Singura soluție este cea din figura 3.12, în care primul arbore conține $p-1$ arce orientate spre stânga, iar al doilea $p-1$ arce orientate spre dreapta;

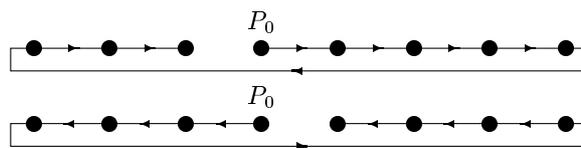


Figura 3.12: O familie formată din doi arbori de acoperire arc-disjuncti de înălțime $p - 1$, pe un inel.

înălțimea arborilor este $p - 1$; singurele arce nefolosite din inel sunt cele orientate spre P_0 . Se împarte așadar mesajul M în două pachete de dimensiuni egale, trimise de P_0 în cele două sensuri, și retransmise de-a lungul întregului inel, conform următorului algoritm:

ALGORITM 3.3 (P_0 difuzează mesajul M pe inel, pe doi arbori de acoperire arc-disjuncti, model multiport, full duplex)

1. **dacă** $\text{id} = 0$ **atunci**
 1. **în paralel** $\text{send}(M^0, \text{dreapta})$, $\text{send}(M^1, \text{stânga})$
2. **altfel, în paralel**
 1. $\text{rcv}(M^0, \text{stânga})$, **dacă** $\text{id} \neq p - 1$ **atunci** $\text{send}(M^0, \text{dreapta})$
 2. $\text{rcv}(M^1, \text{dreapta})$, **dacă** $\text{id} \neq 1$ **atunci** $\text{send}(M^1, \text{stânga})$

Așadar, pachetul M^0 este trimis pe primul arbore de acoperire din figura 3.12, iar M^1 pe al doilea. Dacă urmărim puțin în timp "traectoria" celor două pachete, observăm că instrucțiunile 2.1 și 2.2 nu se execută în general în paralel, ci numai în cazul unui singur procesor (cu $\text{id} = p/2$, când p este par); acesta nu e un motiv pentru a nu le programa în paralel; dimpotrivă, nefiind nici o corelație între operațiile din cele două instrucțiuni, algoritmul 3.3 este general și simplu. Mai mult, algoritmul se poate generaliza pentru un caz de transmisie pipeline. Fiecare procesor va executa de ν ori instrucțiunile 1.1 (sursa), respectiv 2.1 și 2.2 (celelalte procesoare), unde ν este numărul de pachete.

Pentru modelul half duplex, timpul de difuzare este de două ori mai mare decât pentru full duplex, deci nu se obține nici un câștig prin utilizarea a doi arbori de acoperire.

Translația unui arbore. Ce se întâmplă dacă nodul care face difuzarea nu este P_0 , ci un P_s oarecare? Este clar că arborii de acoperire se construiesc în același mod ca pentru P_0 . Vom numi *translație* operația prin care dintr-un arbore de acoperire $\mathcal{A}(0)$ cu rădăcina în P_0 se obține un altul, $\mathcal{A}(s)$, cu aceeași formă, dar cu rădăcina în P_s . Pentru inel, un nod oarecare $x \in \mathcal{A}(0)$ se transformă simplu în $x' \in \mathcal{A}(s)$ prin relația $t(x) = (x + s) \bmod p$; un arc unind nodurile x și y în $\mathcal{A}(0)$, se transformă în arcul ce unește nodurile $t(x)$ și $t(y)$; de acum înainte vom vorbi doar despre transformările aplicate nodurilor, rămânând implicit că un arc ce unește două noduri se transformă în arcul ce unește transformatele nodurilor.

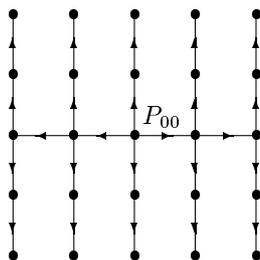


Figura 3.13: Un arbore de acoperire pentru tor.

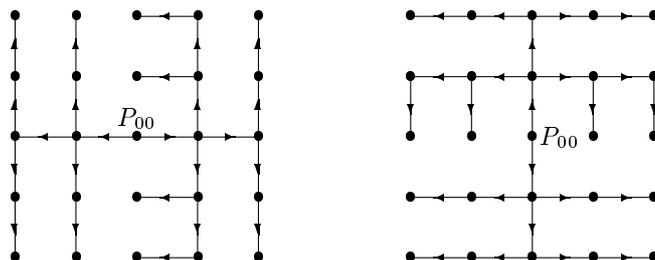


Figura 3.14: O familie formată din doi arbori de acoperire arc-disjuncti, pe un tor.

Difuzare pe tor

Vom trece repede peste tor; P_{00} este procesorul care difuzează. În figura 3.13 este prezentat cel mai evident arbore de acoperire, de înălțime egală cu diametrul torului. Distanța, pe acest arbore, între P_{00} și orice nod, este cea mai mică posibil (egală cu distanța pe tor între cele două noduri). Algoritmul de difuzare pe acest arbore este prezentat în rezolvarea problemei 3.3.12.

Cum el conține toate cele patru arce care pornesc din rădăcină, este de presupus că se poate mai bine. În figura 3.14 sunt prezentați doi arbori de acoperire arc-disjuncti; privind cu puțină atenție, se observă că cel de-al doilea este obținut din primul printr-o rotație cu 90 de grade în sens trigonometric, în jurul procesorului P_{00} , pe care-l identificăm cu punctul de coordonate $(0,0)$; o astfel de rotație transformă punctul (x,y) în $(-y,x)$ (pentru verificare, considerați un exemplu, să zicem $(2,1)$); evident, așa cum am mai spus-o, toate coordonatele trebuie gândite modulo \sqrt{p} . Înălțimea arborilor este egală cu diametrul torului. Construcția decurge aproape imediat din cea din figura 3.13; arcele de pe verticala nodului P_{00} sunt singurele stâmpenitoare: după rotație, s-ar suprapune cu cele de pe orizontala lui P_{00} ; ele sunt eliminate și înlocuite cu arce orizontale.

Cum ideea rotației pare benefică și este și deosebit de simplă, să o aplicăm pentru

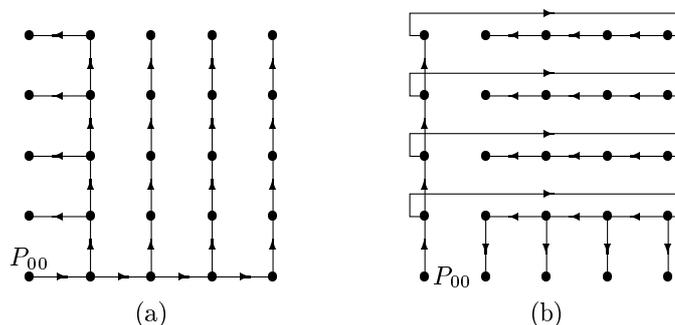


Figura 3.15: (a) Arbore de acoperire pe tor, din care se poate genera prin rotații o familie de patru arbori arc-disjuncti; (b) arbore de acoperire obținut prin rotația cu 90 de grade în sens trigonometric a arborelui de la (a).

a construi patru arbori de acoperire arc disjuncti. Este evident că fiecare va folosi un singur arc plecând din P_{00} . Primul arbore este cel construit în figura 3.15a, gândind aproximativ după cum urmează; nodurile de pe orizontala lui P_{00} sunt toate unite prin arce mergând spre dreapta; din fiecare se construiesc arce în sus, pe toată verticala; în acest fel se acoperă toate nodurile torului, mai puțin cele de pe verticală lui P_{00} ; se vede imediat că, între arborii obținuți prin rotație, nu există conflicte de utilizare a arcelor. În fine, nodurile de pe verticala lui P_{00} sunt legate de cele de pe verticala lui P_{01} , deci de la dreapta spre stânga, din nou fără conflicte la rotație. Pentru edificare, prezentăm și arborii obținuți prin rotație cu 90 de grade în sens trigonometric, în figura 3.15b. Înălțimea acestor arbori este $2(\sqrt{p} - 1)$, deci de aproximativ două ori mai mare decât a arborilor din figura 3.14.

Trebuie menționat că se pot construi patru arbori de acoperire arc disjuncti de înălțime $D + 1$, unde $D = 2\lfloor\sqrt{p}/2\rfloor$ este diametrul torului; acesta este optimul, dar construcția este relativ complicată.

Pentru tor, translația se face cu aceeași operație ca pentru inel, însă aplicată pentru ambele coordonate ale adresei unui procesor.

Difuzare pe hipercub

Arbore de acoperire binomial. În figura 3.16a este prezentat un arbore de acoperire binomial; construcția sa poate fi descrisă în mai multe feluri; de exemplu, în mod recursiv: arcul între nodurile 0 și 1 face parte din arbore; se continuă analog construcția în cele două hipercuburi de dimensiune $d - 1$, obținute după eliminarea ultimului bit al adreselor nodurilor, adică $\mathcal{H}_d^{d-1}(0)$ și $\mathcal{H}_d^{d-1}(1)$; o observație imediată este că fiecare subarbore al unui arbore binomial este la rândul său un arbore binomial. Sau, dacă i este poziția pe care se află cel mai semnificativ bit de 1 (-1, pentru nodul 0), atunci fiii acestui nod sunt vecinii în hipercub pe dimensiunile $i + 1, \dots$,

$d - 1$. Arborele se numește binomial deoarece numărul nodurilor pe nivelul k este $\binom{d}{k}$ (totuși, nu e singurul arbore de acoperire cu această proprietate).

O formă a algoritmului de difuzare pe acest arbore este sugerată într-un mod mai intuitiv în figura 3.17. La fiecare pas, se transmite pe o singură dimensiune a hipercubului, de la toate procesoarele aflate în posesia mesajului M , către vecinii din acea dimensiune; fiecare procesor folosește un singur canal la un moment dat. După k pași, difuzarea într-un hipercub k -dimensional conținând P_0 este încheiată. Algoritmul este formulat mai riguros în continuare:

ALGORITHM 3.4 (P_0 difuzează mesajul M pe hipercub, după un arbore binomial, model 1-port)

1. $l \leftarrow$ poziția celui mai semnificativ bit de 1 din id (-1 pentru P_0)
2. **dacă** id $\neq 0$ **atunci**
 1. **recv**(M, l) {recepționează de la tată}
3. **pentru** $i = l + 1 : d - 1$
 1. **send**(M, i) {trimite tuturor fiilor}

Cu acest algoritm, procesoarele nu primesc mesajul în cel mai scurt timp fiecare. Aceasta se poate realiza în modelul multiport, adăugând o instrucțiune **în paralel** instrucțiunii 3; astfel, fiecare procesor va expedia imediat mai departe mesajul tuturor fiilor. Timpul global de difuzare nu se va modifica însă, el fiind cel după care procesorul $2^d - 1$ primește mesajul.

Familie de arbori binomiali. Deoarece în algoritmul 3.4 la fiecare pas se comunică pe o singură dimensiune, o idee de îmbunătățire a timpului de comunicare este de a folosi toate dimensiunile simultan. Se împarte mesajul în d pachete M^0, M^1, \dots, M^{d-1} ; la pasul 0, procesorul P_0 trimite câte un pachet fiecărui vecin, i.e. pachetul M^i vecinului din dimensiunea i ; la pasul k , pachetul M^i va fi trimis pe direcția $(i + k) \bmod d$, de către toate procesoarele care îl dețin. Funcționarea algoritmului este ilustrată în figura 3.18.

Aceasta revine la a utiliza o familie de arbori de acoperire construiți din cel din figura 3.16a, pe care-l vom nota $\mathcal{A}_0(0)$. Ceilalți $d - 1$ arbori se construiesc prin rotații: arborele $\mathcal{A}_i(0)$, cu $1 \leq i \leq d - 1$, se obține transformând orice nod $x \in \mathcal{A}_0(0)$ în $\text{Ro}(x, i)$. Un exemplu de astfel de arbori este prezentat în figura 3.16b,c, în care se află $\mathcal{A}_1(0)$, respectiv $\mathcal{A}_2(0)$, pentru un hipercub 3-dimensional. Faptul că arborii sunt obținuți prin rotații a condus la denumirea algoritmului, numit *rotativ*. Acești arbori nu sunt arc-disjuncti, dar mecanismul de difuzare descris mai sus evită conflictele în ocuparea arcelor. Totuși, deși timpul de difuzare este de aproximativ de d ori mai mic decât al algoritmului 3.4, nu se poate utiliza transmisia pipeline. Acest defect este remediat în continuare.

O familie de arbori de acoperire arc-disjuncti. Pornim de la următoarea idee algoritmică intuitivă: P_0 trimite mesajul tuturor vecinilor săi; apoi, aceștia difuzează mesajul după câte un arbore de acoperire binomial, cu singurul amendament că nu îl mai retrimite lui P_0 .

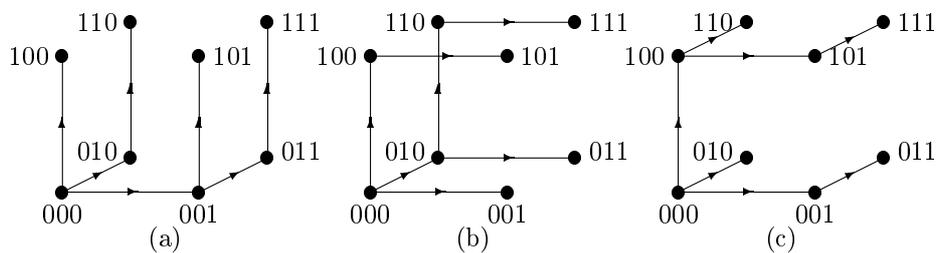


Figura 3.16: (a) Arbore de acoperire binomial într-un hipercub de dimensiune 3 (b), (c) arbori binomiali obținuți prin rotația primului cu una, respectiv două poziții la stânga.

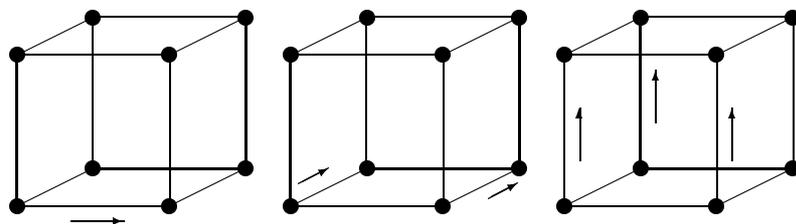


Figura 3.17: Difuzare după un arbore binomial, la fiecare pas pe o direcție.

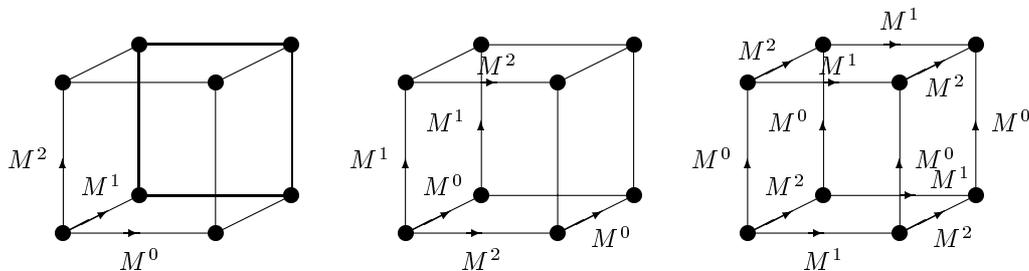


Figura 3.18: Difuzare într-un hipercub de dimensiune 3, algoritmul rotativ (pe o familie de arbori de acoperire binomiali).

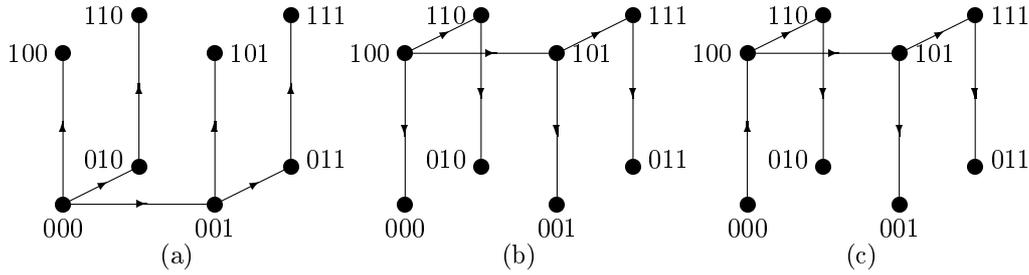


Figura 3.19: (a) Arbore de acoperire binomial; (b) Arbore obținut prin translatarea primului în nodul 100; (c) Arborele traslatat, cu sensul arcului din P_0 schimbat.

Definim întâi operația de *translație pe hipercub*. Un arbore oarecare $\mathcal{A}(0)$ este traslatat cu rădăcina în nodul cu adresa s prin transformarea fiecărui nod $x \in \mathcal{A}(0)$ în $x \oplus s \in \mathcal{A}(s)$ (și transformarea corespunzătoare a arcelor). Un exemplu este prezentat în figura 3.19, în care desenul (a) prezintă arborele de acoperire binomial $\mathcal{A}_0(000)$, iar desenul (b), arborele $\mathcal{A}_0(100)$. Să observăm că, după translație, orice arc rămâne în aceeași dimensiune (dacă nodurile x și y sunt vecine în arbore, sunt și în hipercub, deci adresele lor diferă într-un singur bit; atunci și $x \oplus s$ și $y \oplus s$ diferă în același bit).

Acum putem construi câte un arbore de acoperire binomial în fiecare dintre vecinii nodului 0. Pe care însă, dintre cei d arbori posibili, rădăcina fiind fixată și scopul obținerea unei familii arc-disjuncte? Vom încerca să traslatăm câte un arbore de acoperire binomial cu rădăcina în nodul 0, în fiecare din vecini. Ne bazăm pe faptul că în orice arbore de acoperire binomial există o singură frunză vecină cu rădăcina; de exemplu, pentru $\mathcal{A}_0(0)$, aceasta este 100 (în \mathcal{H}_3); vom traslata arborele în această frunză; rezultatul a fost deja prezentat în figura 3.19b. Cum, la difuzare, comunicația între nodurile 0 și 100 se desfășoară dinspre 0, vom schimba sensul arcului dintre aceste două noduri, obținând un arbore cu rădăcina în nodul 0; rezultatul este prezentat în figura 3.19c; acest arbore va fi notat $\mathcal{A}_0^D(0)$. Să generalizăm: în $\mathcal{A}_i(0)$, frunza vecină cu rădăcina este $2^{(i-1) \bmod d}$; atunci $\mathcal{A}_i^D(0)$ se obține din $\mathcal{A}_i(2^{(i-1) \bmod d})$ schimbând sensul unicului arc adiacent nodului 0. Cei d arbori $\mathcal{A}_i^D(0)$, cu $0 \leq i \leq d-1$, sunt arc-disjuncti (lăsăm fără demonstrație această afirmație) și au fiecare înălțime $d+1$. Pentru \mathcal{H}_3 , această familie este prezentată în figura 3.20.

3.3.3 Difuzare generală

Actul de naștere al difuzării generale l-a constituit următoarea problemă frivolă: p persoane știu fiecare câte o parte dintr-o istorie mondenă; de câte telefoane este nevoie pentru ca toate persoanele să afle întreaga istorie? Se presupune că, într-o comunicare telefonică, fiecare persoană spune tot ce știe în momentul respectiv. Recunoaștem imediat modelul: 1-port, timp constant, full duplex (deși nu se comunică simultan,

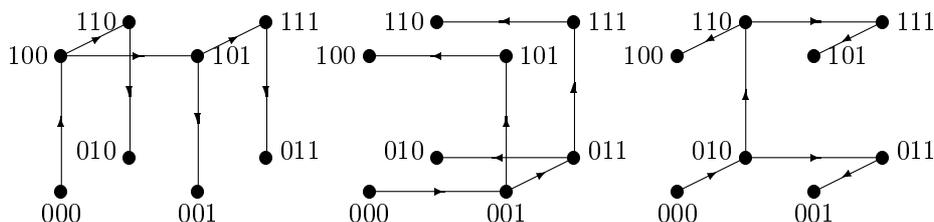


Figura 3.20: Familie de arbori de acoperire arc-disjuncti într-un hipercube de dimensiune 3.

totuși se folosesc ambele sensuri în *aceeași* convorbire); comunicarea se desfășoară pe un graf complet; pentru aprofundare, vezi problema 3.3.19. Și acum să trecem la formulări serioase.

O abordare naivă a problemei ar fi următoarea: un procesor colectează toate mesajele, apoi le difuzează celorlalte procesoare; după ce vom fi discutat despre colectare se va vedea că, de multe ori, aceasta durează cât o difuzare generală; în plus, difuzarea ar fi făcută cu un mesaj imens, de dimensiune pm . Deci, sunt slabe speranțe de performanță.

Idee eficientă generală. O modalitate naturală și eficientă de a obține algoritmi de difuzare generală este de a aplica algoritmi de difuzare din fiecare nod, aceștia executându-se în paralel. Dacă, la un moment dat, un procesor are de trimis mai multe mesaje pe un canal, atunci le va concatena într-un mesaj mai lung, procesorului care recepționează acest mesaj revenindu-i implicit sarcina de a extrage mesajele inițiale și de a le retransmite pe căile corespunzătoare. Dacă graful este simetric, nu vor exista dezechilibre de încărcare a canalelor. Se va vedea în continuare că algoritmi simpli de difuzare duc la bune performanțe în difuzarea generală; de asemenea, deoarece multe arce (eventual toate) vor fi ocupate simultan, nu mai este necesară utilizarea tehnicii pipeline.

Difuzare generală pe inel

Pe un inel, difuzarea generală durează tot atât ca și o difuzare; dacă fiecare procesor își transmite mesajul pe ambele canale, apoi retransmite mesajele primite de la vecini, în ambele sensuri, toate canalele sunt folosite în paralelism complet, iar situația de a transmite două mesaje pe același canal în aceeași etapă nu apare. Practic, fiecare procesor își difuzează mesajul M_i după câte un arbore de acoperire de tipul celui din figura 3.11, procesorul fiind situat în rădăcină. Sunt $\lfloor p/2 \rfloor$ etape de comunicare, cât este înălțimea arborelui, după algoritmul următor:

ALGORITM 3.5 (difuzare generală pe inel, fiecare procesor pe un arbore de acoperire de înălțime minimă)

1. pentru $k = 0 : \lfloor p/2 \rfloor - 1$

1. **în paralel**

1. **send**($M_{(\text{id}-k) \bmod p}$, dreapta)
2. **recv**($M_{(\text{id}-k-1) \bmod p}$, stânga)
3. **dacă** p este impar sau $k \neq \lfloor p/2 \rfloor - 1$ **atunci, în paralel**
 1. **send**($M_{(\text{id}+k) \bmod p}$, stânga)
 2. **recv**($M_{(\text{id}+k+1) \bmod p}$, dreapta)

În etapa k , un procesor retransmite vecinilor mesajele celor două procesoare aflate la distanță k de el (propriile mesaje când $k = 0$) și recepționează de la vecini mesajele procesoarelor aflate la distanță $k + 1$; evident, procesorul memorează separat aceste mesaje. Condiția din linia 1.3 se explică prin faptul că atunci când p este par, cele două ramuri ale arborelui din figura 3.11 nu au lungimi egale; în acest algoritm am presupus că ramura din dreapta are lungime $p/2$, iar cea din stânga $p/2 - 1$, de aceea există o etapă în plus de comunicație spre dreapta. În regim full duplex, timpul de execuție este

$$T_{3,R,F} = \lfloor p/2 \rfloor (\sigma + m\beta).$$

În modelul half duplex, circulația mesajelor se face într-un singur sens, de exemplu fiecare procesor transmite spre dreapta și recepționează simultan din stânga. Difuzarea generală se termină atunci când mesajul fiecărui procesor a ajuns la vecinul său din stânga. Sunt deci $p - 1$ etape de comunicație.

Difuzare generală pe tor

Pe tor, se poate generaliza algoritmul 3.5 de la inel. Într-o primă fază, se face difuzare generală pe fiecare inel orizontal în timp $(\sqrt{p}/2)(\sigma + m\beta)$; apoi, se face din nou o difuzare generală, dar acum pe fiecare inel vertical, fiecare procesor comunicând toate mesajele procesoarelor de pe linia sa; timpul necesar este $(\sqrt{p}/2)(\sigma + \sqrt{p}m\beta)$.

Se observă că, în orice moment, nu se folosesc decât jumătate dintre arce, cele orizontale în prima fază, cele verticale în a doua. Pentru a le utiliza pe toate recurgem la un truc simplu și deja cunoscut; fiecare procesor își împarte mesajul în două pachete; primul este difuzat cum am spus mai sus, iar al doilea, după aceeași idee, dar întâi pe verticală și apoi pe orizontală. În acest fel toate arcele sunt ocupate tot timpul, iar timpul de difuzare generală este

$$T_{3,T,F} = \sqrt{p}\sigma + \left(\frac{mp}{4} + \frac{m\sqrt{p}}{4}\right)\beta. \quad (3.6)$$

Difuzare generală pe hipercub

Pe un hipercub, simetrizarea algoritmului de difuzare pe un arbore de acoperire binomial conduce la o difuzare generală ce se desfășoară după cum se sugerează în figura 3.21. La fiecare pas se comunică pe o singură dimensiune, în ordinea $0, 1, \dots, d - 1$, pe toate arcele posibile; fiecare procesor transmite vecinului său din acea dimensiune

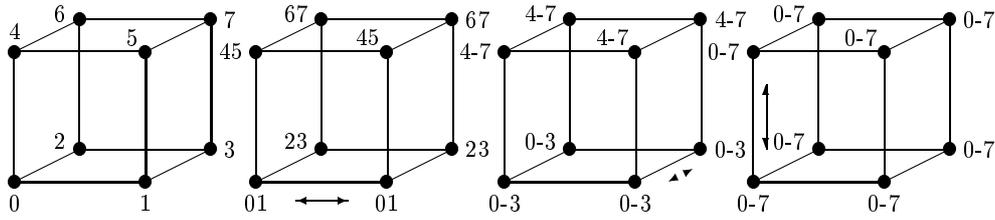


Figura 3.21: Difuzare generală într-un hipercub de dimensiune 3, la fiecare pas pe o direcție (în noduri sunt notate mesajele ajunse acolo după fiecare etapă).

toate mesajele pe care le posedă. După fiecare pas dimensiunea mesajelor transmise se dublează; în figură, fiecare nod are asociate numerele mesajelor deja primite; de exemplu, după doi pași, nodul 1 posedă mesajele 0, 1, 2, 3. La pasul k este încheiată difuzarea generală în cele 2^{d-k+1} hipercuburi formate cu arce din dimensiunile $0, 1, \dots, k-1$ (\mathcal{H}_d^{d-k+1} , în notația de până acum), ceea ce arată natura recursivă a algoritmului (care poate fi descris scurt astfel: se face difuzare generală în două subhipercuburi componente, apoi se schimbă toate datele locale între procesoare aflate pe aceeași poziție în cele două subhipercuburi).

Ca și la difuzare, defectul acestui algoritm este că nu se folosesc, la un moment dat, decât arcele dintr-o singură dimensiune. De aceea încercăm simetrizarea algoritmului de difuzare rotativ, care conduce la folosirea în paralel a tuturor canalelor de comunicație. Procesorul P_j împarte mesajul propriu în d pachete M_j^i , cu $0 \leq i \leq d-1$. La pasul k , procesorul P_j transmite pe direcția $l = (i+k) \bmod d$ pachetul (modificat prin concatenare) M_j^i ; de asemenea, concatenează pachetele primite pe direcția l cu M_j^i ; în final, toate pachetele M_j^i , pentru $0 \leq j \leq p-1$, au același conținut; pentru procesorul P_j , pachetele M_j^i , cu $0 \leq i \leq d-1$, conțin mesajele tuturor procesoarelor, adică întreaga cantitate de informație pe care trebuie să o primească P_j . Deși principalul algoritm e destul de simplu, concatenarea pachetelor conduce la amestecarea lor; de aceea, în final (sau, mai bine, după fiecare etapă de comunicație) fiecare procesor trebuie să pună în ordine pachetele primite; aceasta se poate face simplu dacă se atașează fiecărui pachet M_j^i numerele j (procesorul sursă) și i (numărul de ordine) ca antet. Pentru a exemplifica amestecarea, să ne mai uităm la figura 3.21; presupunem că se concatenează mesajul primit la finalul celui local; procesorul 5 are după prima etapă mesajele 5,4, după a doua 5,4,7,6, și după a treia 5,4,7,6,1,0,3,2; aici ordonarea se face numai după numărul procesorului sursă; evident, când fiecare mesaj e împărțit în pachete, lucrurile se complică și mai mult. Timpul de difuzare generală este

$$T_{3,H,F} = (\log p)\sigma + \frac{(p-1)m}{\log p}\beta.$$

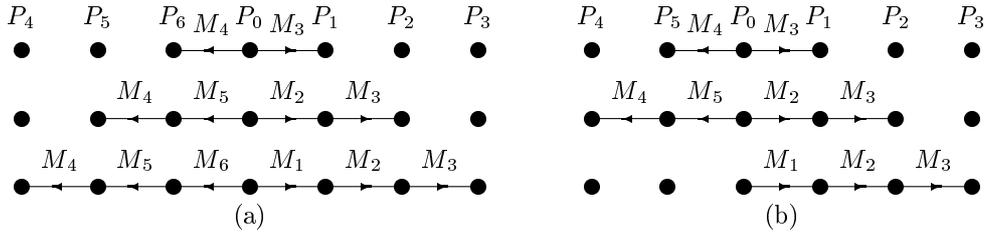


Figura 3.22: Distribuție în trei etape, pe un inel: (a) cu 7 procesoare; (b) cu 6 procesoare.

3.3.4 Difuzare personalizată (distribuție)

Deoarece procesorul P_0 trebuie să trimită mesaje tuturor celorlalte procesoare, utilizarea unui arbore de acoperire sau a unei familii de arbori de acoperire arc-disjuncti este naturală. Principiul de utilizare a unui arbore de acoperire pentru distribuție este ca P_0 să trimită mesajele *întâi celor mai îndepărtate* procesoare. Astfel, în timp ce primele mesaje se deplasează mai departe prin arbore spre destinațiile lor, procesorul P_0 poate trimite mesaje procesoarelor mai apropiate; se poate spera la utilizarea în paralel a multor arce ale arborelui. Vom vedea, în cazul hipercubului, ce calitate este bine să aibă arborii de acoperire.

Distribuție pe inel

Pe un inel, procesorul sursă P_0 transmite spre stânga și spre dreapta simultan, deci folosind arborele de acoperire din figura 3.11, întâi celor mai îndepărtate procesoare, i.e. cele cu adresele $\lfloor p/2 \rfloor$, prin dreapta, și respectiv $\lfloor p/2 \rfloor + 1$, prin stânga. La pasul următor, P_0 trimite mesajele pentru procesoarele aflate cu un arc mai aproape etc. Toate celelalte procesoare retransmit mesajele primite, până în momentul în care îl primesc pe cel destinat lor. Algoritmul este prezentat în continuare și exemplificat în figura 3.22.

ALGORITM 3.6 (P_0 distribuie pe inel)

1. $j \leftarrow \lfloor p/2 \rfloor - \text{dist}(\text{id}, 0) + 1$
2. **dacă** p este par și $\text{id} > \lfloor p/2 \rfloor$ **atunci** $j \leftarrow j - 1$
3. **pentru** $k = 0 : j - 1$
 1. **dacă** $\text{id} = 0$ **atunci, în paralel**
 1. **send**($M_{\lfloor p/2 \rfloor - k}$, dreapta)
 2. **dacă** p este impar sau $k < j - 1$ **atunci** **send**($M_{\lfloor p/2 \rfloor + 1 + k}$, stânga)
 2. **altfel dacă** $\text{id} \leq \lfloor p/2 \rfloor$ **atunci, în paralel**
 1. **recv**($M_{\lfloor p/2 \rfloor - k}$, stânga)
 2. **dacă** $k > 0$ **atunci** **send**($M_{\lfloor p/2 \rfloor - k - 1}$, dreapta)
 3. **altfel, în paralel** $\{\text{id} > \lfloor p/2 \rfloor\}$

1. **recv**($M_{\lfloor p/2 \rfloor + 1 + k}$, dreapta)
2. **dacă** $k > 0$ **atunci** **send**($M_{\lfloor p/2 \rfloor + k}$, stânga)

Numărul de etape la care participă un procesor este calculat în instrucțiunile 1 și 2; el este cu atât mai mic cu cât distanța de procesorul care distribuie este mai mare; când p este par, arborele de acoperire din figura 3.11 are ramurile de lungimi inegale, deci procesoarele din stânga lui P_0 participă la o etapă de comunicație mai puțin, după cum se vede și în figura 3.22b. Procesorul P_0 transmite în ambele sensuri simultan în fiecare etapă, mai puțin în ultima, dacă p este par, după cum se vede din 3.1.2. În fine, celelalte procesoare doar recepționează în prima lor etapă de comunicație, apoi recepționează și transmit simultan. Instrucțiunile 3.2.1 și 3.2.2 sunt identice cu 3.3.1 și respectiv 3.3.2, doar sensul de comunicație este schimbat; este descrisă comunicația pentru procesoarele din dreapta, respectiv stânga, lui P_0 . Remarcăm că un procesor poate folosi doar două zone de memorie, alternativ pentru recepție și transmisie; ultimul mesaj primit este cel destinat lui. Complexitatea acestui algoritm este

$$T_{A,R} = \lfloor p/2 \rfloor (\sigma + m\beta).$$

Aplicarea tehnicii pipeline în distribuție nu face decât să înrăutățească performanțele. La inel, aceasta se vede imediat, deoarece procesorul P_0 este în permanentă ocupat. Spargerea mesajelor în pachete n-ar modifica timpul de propagare, dar ar mări coeficientul lui σ , deci și timpul total; așadar, nu se utilizează pipeline.

Distribuție pe tor

Pe tor, putem să ne inspirăm din algoritmul de la inel, folosind implicit arborele de acoperire din figura 3.13 după cum urmează. În prima fază, nodul P_{00} distribuie fiecărui nod de pe orizontala sa toate cele \sqrt{p} mesaje ce trebuie trimise pe verticala acelui nod; timpul necesar este (aproximativ) $(\sqrt{p}/2)(\sigma + m\sqrt{p}\beta)$. În a doua fază, fiecare procesor de pe linia 0 distribuie pe inelul vertical pe care se află; de data aceasta talia mesajelor este m , deci timpul necesar este $(\sqrt{p}/2)(\sigma + m\beta)$.

Pentru utilizarea simultană a mai multor canale, se împarte fiecare mesaj în două pachete; primul pachet se distribuie ca mai sus; al doilea, după aceeași idee, dar întâi pe verticală și apoi pe orizontală, adică după un arbore de acoperire obținut prin rotația celui din figura 3.13 cu 90 de grade în sens trigonometric. Timpul de distribuție este identic cu cel de difuzare generală dat de (3.6).

Distribuție pe hipercub

Distribuție pe un arbore de acoperire binomial. Pe un hipercub, o primă idee este de a transmite folosind un arbore de acoperire binomial; să încercăm, pentru a vedea care sunt dezavantajele. În figura 3.23 este prezentată evoluția comunicațiilor, dacă în fiecare etapă se transmite pe câte o dimensiune; în prima etapă, P_0 transmite (concatenate) toate mesajele pentru procesoarele din subhipercubul $\bar{\mathcal{H}}_d^{d-1}(1)$, în a

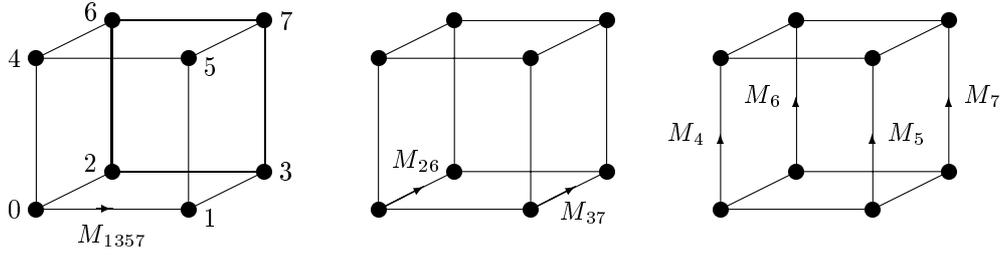


Figura 3.23: Difuzare personalizată, într-un hiper cub de dimensiune 3, pe un arbore de acoperire binomial, la fiecare pas pe o direcție.

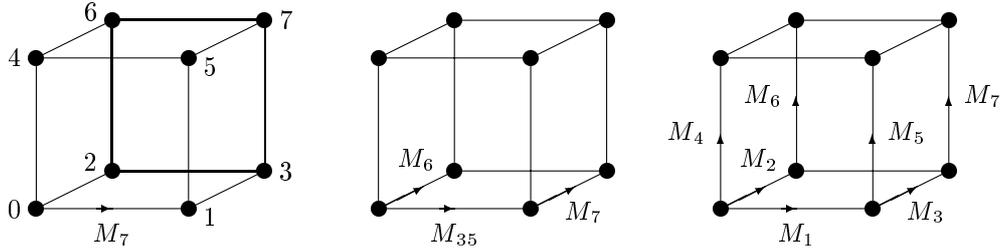


Figura 3.24: Difuzare personalizată pe un arbore de acoperire binomial, după principiul "întâi celui mai îndepărtat".

două, toate mesajele pentru $\bar{\mathcal{H}}_d^{d-2}(2)$; în general, în etapa k , cu $0 \leq k \leq d-1$, trimite pe dimensiunea k toate mesajele pentru subhipercubul $\bar{\mathcal{H}}_d^{d-1-k}(2^k)$; fiecare procesor care primește un mesaj păstrează porțiunea destinată lui și transmite restul mai departe în subarborii a cărui rădăcină este, într-un mod asemănător cu P_0 . Se observă că dimensiunea mesajelor se înjumătățește la fiecare pas, inițial fiind $2^{d-1}m$. Complexitatea acestui algoritm 1-port este:

$$T_{4,H,1} = (\log p)\sigma + (p-1)m\beta.$$

Tot pe un arbore de acoperire binomial, se poate proceda trimițând mesajele după principiul "întâi celui mai îndepărtat", așa cum este sugerat în figura 3.24. În etapa k ($0 \leq k \leq d-1$), P_0 trimite, simultan pe toate canalele sale, mesajele pentru nodurile aflate la distanță $d-k$, mesajele pentru nodurile din același subarbor fiind concatenate; timpul de transmisie este mai bun de maximum două ori decât precedentul:

$$T'_{4,H} = (\log p)\sigma + \frac{p}{2}m\beta.$$

El se justifică prin faptul că, într-un arbore binomial, unul dintre subarborii rădăcinii, cel mai "stufos", are jumătate din noduri; timpul de distribuție este cel necesar comunicărilor de la rădăcină spre acest subarbor; nu mai sunt introduse întârzieri

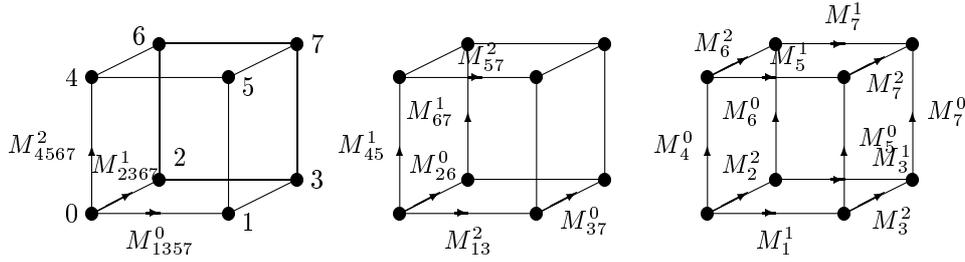


Figura 3.26: Difuzare personalizată pe o familie de arbori de acoperire binomiali rotativi.

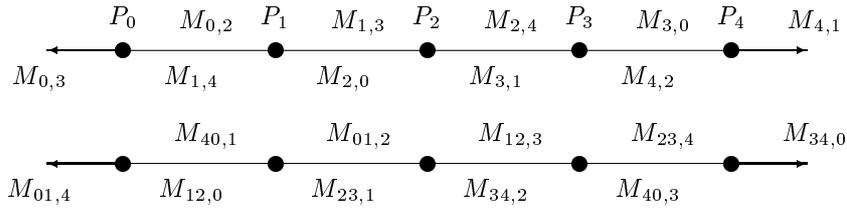


Figura 3.27: Schimb complet pe un inel cu 5 procesoare, în două etape; mesajele de deasupra liniilor orizontale circulă spre dreapta, iar cele de dedesubt, spre stânga.

3.26 și este o generalizare acum banală a celui din figura 3.23.

3.3.5 Schimb complet

Deoarece în schimbul complet fiecare procesor efectuează o distribuție, putem obține algoritmi simpli prin simetrizarea celor de distribuție. Ca și la difuzarea generală, această strategie poate asigura ocuparea în paralel a tuturor canalelor, fără a fi necesară transmisia pipeline. Vom ilustra ideea doar în cazul inelului, cititorul fiind invitat să scrie algoritmi pentru tor și hiper-cub.

Schimb complet pe inel. Dacă fiecare procesor trimite mesaje pe ambele canale, începând cu cele adresate celor mai depărtate procesoare, iar în etapele următoare transmite atât mesajele sale, cât și pe cele primite de la vecini, evoluția comunicației este de genul celei prezentate în figura 3.27; în etapa k , cu $0 \leq k < p/2$, fiecare procesor transmite k mesaje concatenate (cu un mesaj mai mult decât în etapa precedentă) pe fiecare canal. Descriem mai jos algoritmul pentru p impar; în cazul p par se întâmplă același fenomen ca la distribuție: în ultima etapă se comunică într-un singur sens; nu mai detaliem, pentru simplitate.

ALGORITM 3.7 (schimb complet pe inel, pentru p impar)

1. $T_s \leftarrow \emptyset, T_d \leftarrow \emptyset$
2. **pentru** $k = 0 : \lfloor p/2 \rfloor - 1$
 1. $jd \leftarrow (\text{id} + \lfloor p/2 \rfloor - k) \bmod p$

2. $js \leftarrow (id - \lfloor p/2 \rfloor + k) \bmod p$
3. $L_d \leftarrow T_s$ concatenat cu $M_{id, jd}$
4. $L_s \leftarrow M_{id, js}$ concatenat cu T_d
5. **în paralel**
 1. **send**(L_d , dreapta), **send**(L_s , stânga), **recv**(T_d , dreapta), **recv**(T_s , stânga)

Variabilele folosite au următoarea semnificație: jd și js sunt adresele celor două procesoare cărora le sunt destinate mesajele curente, aflate la dreapta procesorului curent, respectiv la stânga (pe calea cea mai scurtă); T_d și T_s sunt zonele de memorie în care se depun mesajele primite, în etapa curentă, din dreapta, respectiv stânga; în fine, în L_d și L_s se formează mesajele destinate transmisiei, în etapa curentă, spre dreapta, respectiv stânga. Deci, în fiecare etapă, un procesor concatenează câte un mesaj propriu mesajelor primite în etapa anterioară și apoi le trimite mai departe; mesajele destinate lui sunt recepționate în ultima etapă, deci se vor afla în variabilele T_d și T_s . Datorită modului de scriere a instrucțiunilor 2.3 și 2.4, T_d va conține mesajele de la procesoarele, în ordine, $(id + 1) \bmod p, \dots, (id + \lfloor p/2 \rfloor) \bmod p$, iar T_s de la, tot în ordine, $(id - \lfloor p/2 \rfloor) \bmod p, \dots, (id - 1) \bmod p$; deci, mesajele sunt în ordinea crescătoare a adreselor expeditorilor, mai puțin o rotație.

Ca și la distribuție, au loc $\lfloor p/2 \rfloor$ etape de comunicație, deci timpul de execuție al algoritmului, în regim full duplex, este

$$T_{5,R,F} = \sum_{k=1}^{p/2} (\sigma + km\beta) = \frac{p}{2}\sigma + \frac{p}{4}\left(\frac{p}{2} + 1\right)m\beta.$$

Este interesant că se poate aplica o strategie într-un fel inversă, fiecare procesor începe prin a trimite jumătate din mesaje spre dreapta (cele pentru procesoarele mai apropiate prin dreapta) și jumătate spre stânga. În a doua etapă, fiecare procesor reține mesajul destinat lui, și retransmite restul mai departe; deci, în etapa k ($0 \leq k < \lfloor p/2 \rfloor$), un procesor trimite $\lfloor p/2 \rfloor - k$ mesaje concatenate (cu unul mai puțin decât în etapa precedentă). Timpul de comunicație va fi exact același.

Schimb complet pe tor. Pe tor generalizăm algoritmi de la inel. Se face întâi câte un schimb complet pe fiecare inel orizontal, cu mesajele destinate tuturor procesoarelor de pe aceeași verticală; un procesor trimite altuia câte \sqrt{p} mesaje. Urmează câte un schimb complet pe fiecare verticală, fiecare procesor având de transmis câte \sqrt{p} mesaje fiecărui alt procesor de pe același inel vertical (câte unul de la fiecare procesor de pe inelul orizontal, inclusiv el însuși). Timpii pentru cele două faze sunt deci egali (iată totuși o diferență față de difuzarea generală și distribuție).

Din nou, dacă se separă fiecare mesaj în două pachete și se folosește algoritmul precedent pentru primul pachet, și același algoritm, dar întâi pe verticală, pentru al doilea, timpul de propagare se înjumătățește.

Transmisie permutată. Vom aminti acum, pe scurt, o utilizare a schimbului complet într-o problemă de comunicație aparent fără mare legătură cu acesta. Să presupunem că fiecare procesor dorește să transmită un singur mesaj, unui alt procesor,

astfel încât fiecare procesor are de transmis și de recepționat câte un mesaj. Problema se numește a transmisiei permutate (permuted send). Un exemplu de astfel de comunicație este *rotația*, în care fiecare procesor transmite vecinului său din dreapta, pe un inel; alt exemplu este *inversarea*, în care P_i schimbă câte un mesaj cu P_{p-1-i} . În general, pot apărea destule situații în care comunicația nu este deloc regulată.

O metodă de a regulariza transmisia permutată este următoarea: fiecare procesor își împarte mesajul în p pachete, pe care le transmite, câte unul fiecăruia dintre celelalte procesoare, unul păstrându-l; acesta este un schimb complet. Acum, fiecare procesor are tot p pachete, dar destinate câte unul fiecărui procesor; fiecare procesor păstrează pachetul său, trimițându-le pe celelalte $p - 1$ celorlalte procesoare; acesta este, din nou, un schimb complet. Dacă privim din punctul de vedere al transmiterii mesajului între procesoarele P_s și P_a , atunci P_s distribuie cele p pachete, după care P_a le colectează. În acest fel, transmisia permutată se transformă în două schimburi complete. Aceasta poate fi o soluție excelentă, mai puțin în cazurile în care există o variantă simplă, ca, de exemplu, pentru rotație.

3.3.6 Câteva concluzii

În tot restul lucrării vom utiliza primitive cu numele operațiilor de comunicație globală prezentate până acum, folosind o descriere nu foarte formală, dar suficient de explicită, de exemplu

difuzare: P_x trimite mesajul A celorlalte procesoare.

Aici poate apărea o confuzie destul de gravă. Această primitivă nu este apelată doar de procesorul P_x , ci de *toate* procesoarele, în aceeași formă. În funcție de adresa sa, fiecare procesor își îndeplinește rolul în difuzare. Evident, adresa x trebuie să aibă aceeași valoare în toate apelurile, deci, în programul unui anume procesor, să nu depindă în nici un fel de adresa procesorului.

În ce privește estimarea timpilor pentru executarea acestor primitive, vom considera, în general, timpii cei mai buni dintre cei prezentați. Există totuși un amendament, pentru cazurile în care se poate folosi pipeline. Dacă mesajele au dimensiuni mari (din păcate, noțiunea de "mare" nu poate fi explicitată prea mult, ea depinzând de modelul de comunicație și de valorile parametrilor comunicației), vom lucra cu timpii pentru pipeline, deci într-adevăr cei mai buni. Dacă însă mesajele sunt scurte, atunci vom presupune că nu se utilizează pipeline (chiar dacă s-ar utiliza, câștigul ar fi nesemnificativ). Oricum, în multe cazuri ne vom mulțumi cu o apreciere mai vagă a timpului, precizând, de exemplu, că o difuzare pe hipercub durează $O(m)$ (sau $O(m/\log p)$, cu pipeline).

Probleme

P 3.3.1 Deduceți o limită inferioară pentru timpul de comunicație între două noduri aflate la distanță j , într-un graf regulat oarecare de grad Δ . La fel, dacă graful este neregulat.

P 3.3.2 Considerăm un graf complet, în care oricare două noduri sunt conectate. Câte căi arc-disjuncte există între două noduri oarecare ?

P 3.3.3 Când sunt optime cele patru căi din figura 3.8a ?

P 3.3.4 Găsiți patru căi arc-disjuncte între două noduri vecine pe o grilă. Indicație: lungimile lor sunt 1, 3, 5, respectiv 7.

P 3.3.5 Cum se construiesc căile de comunicare între nodurile P_s și P_a dintr-un hipercub, cunoscând numai căile dintre nodurile 0 și $2^j - 1$, unde $j = H(s, a)$?

P 3.3.6 a) Să se arate că cele d căi construite prin rotații ale căii principale între nodurile 0 și $2^d - 1$, în hipercubul \mathcal{H}_d , nu au nici un nod comun (în afara capetelor). b) La fel, pentru cele d căi între nodurile 0 și $2^j - 1$, construite în secțiunea 3.3.1.

P 3.3.7 Câte arce sunt între C_k și C_{k+1} , în hipercubul \mathcal{H}_d ?

P 3.3.8 În tabelul 3.1 sunt prezentate margini inferioare ale timpului necesar unor operații de comunicație globală, în modelele full și half duplex, timp liniar, pentru un graf oarecare cu diametrul D și gradul Δ . Demonstrați valabilitatea acestor margini.

P 3.3.9 Care este numărul maxim de arbori de acoperire arc-disjuncți, cu aceeași rădăcină, care se pot construi într-un graf oarecare, regulat și de grad Δ ?

P 3.3.10 Scrieți algoritmul de difuzare pe inel, după arborele de acoperire din figura 3.11.

P 3.3.11 Modificați algoritmul din problema anterioară astfel încât să fie corect indiferent de procesorul P_s care face difuzarea. La fel, pentru algoritmul 3.3.

P 3.3.12 Scrieți algoritmul de difuzare pe tor, după arborele de acoperire din figura 3.13.

P 3.3.13 Dacă un arc este în dimensiunea k în $\mathcal{A}_0(0)$, în ce dimensiune se va afla după rotația prin care se obține $\mathcal{A}_i(0)$?

P 3.3.14 În ce etapă de comunicație va primi mesajul M procesorul P_j , în algoritmul 3.4 ? Dar mesajul M^t , în algoritmul de difuzare rotativ ?

P 3.3.15 Scrieți algoritmul de difuzare rotativ pentru hipercub.

P 3.3.16 Dintr-un algoritm de difuzare se poate obține întotdeauna un algoritm de comunicare de la procesorul care difuzează (sursa) la un oricare altul (destinația), eliminând toate transmisiile inutile ale mesajului; gândind în termeni de arbori de acoperire, se elimină toate arcele care nu sunt pe căile dintre sursă și destinație. Ce se obține din particularizarea algoritmului de difuzare după familia de arbori arc-disjuncți $\mathcal{A}_i^D(0)$?

P 3.3.17 Care este timpul de difuzare într-un graf oarecare de diametru D , considerându-se modelul multiport, cu timp constant ?

P 3.3.18 Care este timpul de difuzare într-un graf complet, în modelul 1-port, cu timp constant ? Dar într-un hipercub ? Care este limita inferioară pentru un graf oarecare ?

P 3.3.19 Cât durează difuzarea generală pe graful complet, în modelul timp constant, multiport, full duplex ? Dar în half duplex ? Dar în modelul 1-port ?

P 3.3.20 Scrieți algoritmul de difuzare în regim half duplex pe inel.

Operația	Full duplex	Half duplex
Difuzare	$D\sigma + (D - 1 + \frac{m}{\Delta})\beta$	$\frac{2(p-1)}{p\Delta}m\beta$
Difuzare generală	$(p - 1)\frac{m}{\Delta}\beta$	$2(p - 1)\frac{m}{\Delta}\beta$
Distribuție	$(p - 1)\frac{m}{\Delta}\beta$	$(p - 1)\frac{m}{\Delta}\beta$
Schimb complet, inel	$\frac{p}{4}(\frac{p}{2} + 1)m\beta$	$\frac{p}{2}(\frac{p}{2} + 1)m\beta$
Schimb complet, hipercub	$\frac{p^m}{2}\beta$	$pm\beta$

Tabelul 3.1: Margini inferioare pentru timpii necesari operațiilor de comunicație globală.

P 3.3.21 Care este timpul de distribuție pe un graf oarecare în modelul cu timp constant ? Sugerati algoritmi de distribuție în modelul 1-port (și timp constant), pentru inel și hipercub. Indicație: timpi optimi: $\lfloor (p + 1)/2 \rfloor$ pentru inel și $\log p$ pentru hipercub.

P 3.3.22 Sugerati un algoritm simplu de distribuție pe tor, care să aibă un timp de execuție asemănător cu $T_{4,T}$, deci să aibă coeficientul lui β aproape optim, dar în care mesajele să nu fie împărțite în pachete.

P 3.3.23 Demonstrați că, în algoritmul de schimb complet pe tor, fiecare mesaj parcurge calea cea mai scurtă între sursă și destinație.

P 3.3.24 Demonstrați că $\sum_{i=0}^d i \binom{d}{i} = d2^{d-1}$.

P 3.3.25 Se consideră o matrice pătrată A , partiționată în blocuri egale, și ele pătrate, notate A_{ij} . a) Presupunem că blocurile se distribuie celor p procesoare astfel încât P_k are blocurile P_{kj} , pentru $0 \leq j < p$ (are o linie de blocuri). Cum se efectuează transpunerea matricii, pe un inel ? b) Dar pe un hipercub ? c) Pe un tor, presupunem că procesorul P_{ij} deține blocul A_{ij} . Care este acum algoritmul de transpunere ?

P 3.3.26 Cât durează schimbul complet pe un graf complet, în modelul timp constant și multiport, respectiv 1-port ?

3.4 Comunicații în modelul comutare de circuit

Dedicăm o secțiune separată acestui model, dar fără amploarea celei pentru store and forward. Ne vom mulțumi cu câteva exemple care să sublinieze filosofia de abordare a problemelor. Notațiile sunt cele deja utilizate. Trebuie să facem o mică modificare, totuși; la primitiva **send**, în modelul store and forward era suficientă indicarea canalului pe care era transmis un mesaj; acum suntem obligați să precizăm și procesorul destinație a mesajului, astfel:

send(*date*, *canal*, P_a)

Deci, mesajul *date* va fi trimis pe canalul *canal*, către procesorul P_a . Vă întrebați poate de ce mai trebuie precizat canalul; răspunsul e simplu: pentru că, în general,

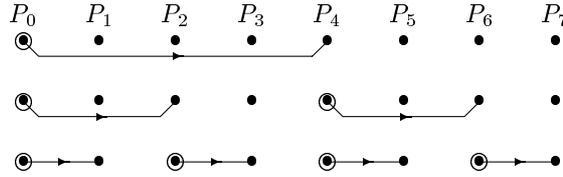


Figura 3.28: Difuzare într-un inel cu 8 procesoare, model comutare de circuit, în 3 etape.

nu există o cale unică până la P_a ; alte argumente vor fi oferite peste câteva pagini, când vom vorbi despre dirijare (acolo se va lămuri și cum se alege calea dintre sursă și P_a , care nu e unică nici după ce sursa a trimis mesajul pe un canal precizat). Vom introduce și pentru **recv** acest al treilea argument, cu semnificația de adresă a expeditorului unui mesaj; el este necesar pentru a putea selecta un mesaj din mai multe sosite pe același canal; în plus, clarifică scrierea algoritmilor.

Comunicație punct-la-punct. În comunicația între două procesoare, timpul este cu atât mai mic cu cât există mai multe căi arc-disjuncte între sursă și destinație. Dacă există c căi, de lungime cel mult h , atunci mesajul este împărțit în c pachete, fiecare trimis pe o cale. Timpul total este dictat de calea cea mai lungă:

$$T_1 = \sigma + h\tau + (m/c)\beta.$$

Pentru comparație, revedeți formula (3.5), de la store and forward în regim pipeline.

Difuzare pe inel, model 1-port. Dacă în comunicarea între două procesoare modelul comutare de circuit este evident mai performant, în difuzare lucrurile nu sunt chiar simple. Să studiem puțin difuzarea pe inel, în modelul 1-port. În primul pas de comunicație, procesorul sursă P_0 informează un alt procesor; în al doilea, cele două procesoare mai pot informa (simultan) alte două, ș.a.m.d., numărul procesoarelor informate dublându-se la fiecare pas. Deci, numărul minim de etape de comunicație este $\lceil \log p \rceil$. Desigur, pentru a se atinge acest minim, e nevoie ca, în fiecare etapă, căile de comunicație să fie arc-disjuncte. Din fericire, pentru inel, aceasta se obține cu ușurință, după cum se vede în figura 3.28. Să presupunem că sunt $p = 2^r$ procesoare. În prima etapă P_0 comunică cu $P_{2^{r-1}}$; în a doua, cele două comunică cu procesoarele aflate la distanță 2^{r-2} spre dreapta; în fiecare etapă distanța de comunicare se înjumătățește față de precedentă; în etapa k , cu $0 \leq k \leq r-1$, se comunică la distanță 2^{r-k-1} . Algoritmul poate fi descris astfel:

ALGORITHM 3.8 (P_0 difuzează mesajul M pe inel, model comutare de circuit, 1-port, $p = 2^r$, în r etape)

1. $j \leftarrow$ numărul de biți finali de 0 din id
2. **dacă** id $\neq 0$ **atunci** **recv**(M , stânga, id $- 2^j$)
3. **pentru** $l = j - 1 : -1 : 0$ $\{l = r - k - 1\}$
 1. **send**(M , dreapta, id $+ 2^l$)

Fiecare procesor primește o singură dată mesajul M și îl transmite de j ori; valoarea j se calculează în instrucțiunea 1; echivalent, se poate scrie $id = z2^j$, cu z un număr impar (convenim că $j = r$, pentru $id = 0$). Timpul total pentru acest algoritm este

$$T_{2,R} = \sum_{k=0}^{r-1} (\sigma + 2^{r-k-1}\tau + m\beta) = (\log p)\sigma + (p-1)\tau + (\log p)m\beta.$$

Comparând acum cu timpii pentru difuzare pe inel în modelul store and forward, observăm că, pentru mesaje scurte, timpii erau $O(pm/2)$, dar pentru mesaje lungi, $O(m/2)$ în regim pipeline. Timpul în model comutare de circuit se află între aceștia; avantajul modului pipeline provine din faptul că toate procesoarele de pe o cale memorau conținutul mesajelor, în timp ce acum ele nu o fac (am încercat să sugerăm aceasta în figura 3.28, liniile pe care "circulă" mesajele, netrecând prin procesoarele intermediare, deși în realitate o fac, desigur).

Difuzare în model multiport și generalizare. În modelul multiport, P_0 poate informa simultan două procesoare, de exemplu pe cele aflate la distanță $p/3$; apoi aceste trei procesoare informează alte 6, la distanță $p/3^2$; în general, la pasul k (primul pas este notat 0), se comunică la distanță $p/3^{k+1}$, în total existând 3^{k+1} procesoare informate. Numărul de pași este deci $\lceil \log_3 p \rceil$.

Aceasta strategie poate fi încercată pe orice graf regulat cu grad Δ . După prima etapă, sunt $\Delta + 1$ procesoare informate, după a doua $(\Delta + 1)^2$ etc. Se poate spera la un număr de $\lceil \log_{\Delta+1} p \rceil$ etape. Problema este de a găsi căi arc-disjuncte și apoi, eventual, de a găsi distanța optimă între procesoarele care comunică (de a minimiza coeficientul lui τ). De exemplu, pe tor, un algoritm destul de regulat poate decurge în modul schițat în figura 3.29. Se împarte torul în 9 "pătrate" egale, P_{00} fiind în centrul pătratului central; în prima etapă P_{00} informează simultan procesoarele aflate în centrele celor patru pătrate din colțuri; în a doua, procesoarele din centrele celorlalte patru pătrate. Apoi, procesoarele din centrele tuturor celor 9 pătrate continuă în același mod, fiecare pătrat fiind împărțit în 9 pătrate mai mici; evident că în pătratele din colțuri difuzarea poate începe mai devreme, după cum se vede în figura 3.29. La fiecare două etape, numărul procesoarelor informate crește de 9 ori; un mic calcul arată că numărul total de etape este $2\lceil \log_9 p \rceil$.

Alte operații de comunicație globală. În difuzarea generală nu se poate spera la îmbunătățiri datorate comutării de circuit, deoarece în cei mai eficienți algoritmi în modelul store and forward nu se utilizează pipeline, iar arcele sunt tot timpul ocupate. E interesant că nici la distribuție nu se poate face mai bine; într-un algoritm în care procesorul sursă emite permanent pe toate canalele, este inutil ca el să-și transmită mesajele mai departe de vecini, ceea ce ar costa mai mulți τ , atâta vreme cât vecinii nu au altceva de făcut decât să retransmită mesajele primite. Totuși, dacă distribuția are loc în paralel cu alte operații de comunicație, este de așteptat ca modelul comutare de circuit să fie mai eficient, utilizând algoritmul banal în care procesorul sursă trimite pe rând mesajele direct la destinație. Nici la schimbul complet lucrurile nu sunt diferite. Prezentăm totuși un algoritm pe hiper-cub care are avantajul simplității.

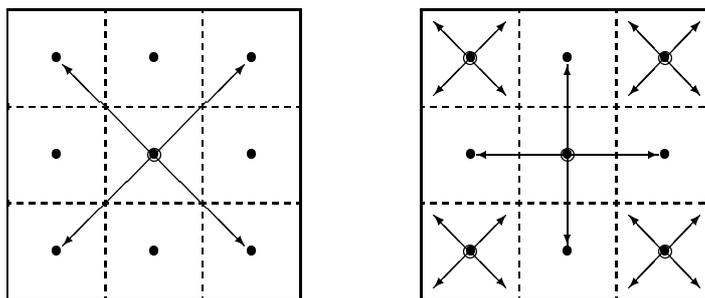


Figura 3.29: Ideea difuzării într-un tor.

Schimb complet pe hipercub. Folosind comutarea de circuit, cea mai banală idee de a realiza schimbul complet este următoarea: vor fi $p - 1$ etape în care un anumit procesor comunică, pe rând, cu celelalte $p - 1$ procesoare, trimițând fiecăruia mesajul său, și primind în același timp mesajul de la procesorul respectiv. Întrebarea este dacă se poate ca toate procesoarele să facă acest lucru în paralel. Altfel zis, dacă în fiecare etapă se pot găsi $p/2$ perechi de procesoare (în fiecare etapă altele) astfel încât cele p căi ce leagă fiecare procesor cu perechea lui să fie arc-disjuncte.

Ideea fericită de formare a perechilor e ca în etapa k , cu $1 \leq k \leq p - 1$, procesorul s să trimită mesajul pentru $s \oplus k$; se vede imediat că, în aceeași etapă, $s \oplus k$ trimite mesajul pentru s , pentru că $(s \oplus k) \oplus k = s$. De asemenea, în etape diferite, s comunică cu procesoare diferite (funcția $f(k) = s \oplus k$ este bijectivă!). Mai rămâne problema nesuprapunerii căilor. Să presupunem că fiecare mesaj circulă pe calea principală dintre două procesoare; în etapa k , calea principală între procesoarele s și $s \oplus k$ depinde numai de $s \oplus s \oplus k = k$; deci biții de 1 din k decid dimensiunile pe care circulă mesajele, indiferent de sursă. Să demonstrăm că aceste căi sunt arc-disjuncte pentru cazul în care $k = p - 1$, adică are toți biții de 1; primul arc de pe fiecare cale se află în dimensiunea 0 și, evident, nu există nici o suprapunere; mulțimea nodurilor aflate la distanță 1 pe fiecare cale este compusă din toate nodurile hipercubului (fiecare cale a ajuns într-un alt nod); al doilea arc este pe dimensiunea 1; cum se pleacă din noduri diferite, arcele sunt diferite și ajung în noduri diferite; de asemenea sunt diferite de oricare prim arc de pe o cale, pentru că acesta era în dimensiunea 0; ș.a.m.d., fiecare cale având câte un arc în fiecare dimensiune. Pentru alți k demonstrația e similară, cu atât mai mult cu cât căile au mai puține arce. În figura 3.30 sunt prezentate toate cele 7 etape de comunicație, cu căile respective, pentru un hipercub de dimensiune 3. Algoritmul este deci următorul:

ALGORITM 3.9 (Schimb complet în hipercub, model comutare de circuit)

1. **pentru** $k = 1 : p - 1$
 1. $ds \leftarrow$ poziția celui mai puțin semnificativ bit de 1 din k

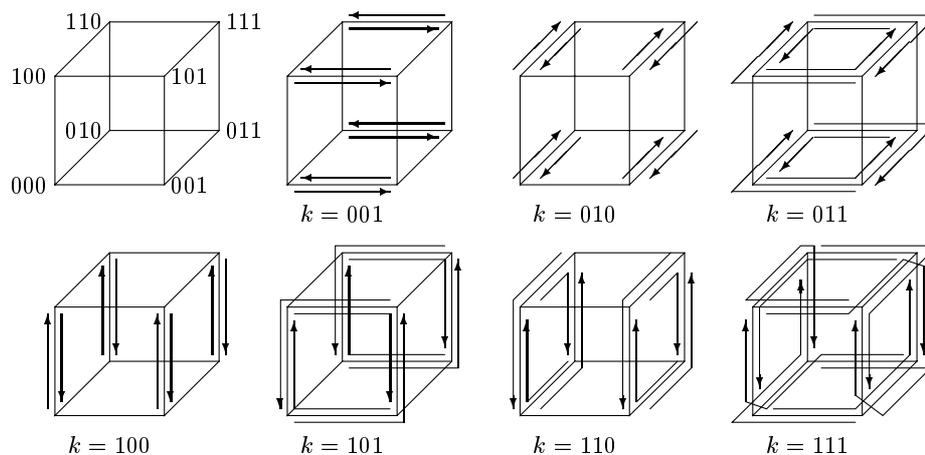


Figura 3.30: Schimb complet în hipercub, model comutare de circuit.

2. $dr \leftarrow$ poziția celui mai semnificativ bit de 1 din k
3. $\text{send}(M_{id, id \oplus k}, ds, id \oplus k)$
4. $\text{recv}(M_{id \oplus k, id}, dr, id \oplus k)$

Trebuie reamintit că algoritmul nu este performant decât dacă mesajele parcurg calea principală (sau alte căi, pentru care arcele de la același nivel sunt în aceeași dimensiune); nu se precizează nicăieri în algoritm acest lucru; calea este decisă de modul de dirijare; cum se face aceasta, vom vorbi ceva mai târziu. Variabilele ds și dr conțin dimensiunile pe care pleacă, respectiv sosește, un mesaj în etapa k . Algoritmului i se poate adăuga o operație de sincronizare efectuată în finalul fiecărei iterații; aceasta previne rămânerea în urmă a unor procesoare, care se poate agrava prin blocarea unor căi; totuși, pentru a se ajunge aici, e necesar ca unele procesoare să ajungă cu două iterații înaintea altora, ceea ce este destul de improbabil.

Timpul necesar execuției algoritmului 3.9 este

$$T_{5,H} = \sum_{k=1}^{p-1} (\sigma + \text{dist}(0, k)\tau + m\beta) = (p-1)\sigma + \frac{dp}{2}\tau + (p-1)m\beta.$$

Coefficientul lui τ este suma $\sum_{i=1}^d i \binom{d}{i}$ pe care am calculat-o deja în problema 3.3.24 și provine din faptul că un procesor comunică o dată cu fiecare alt procesor, și sunt $\binom{d}{i}$ procesoare la distanță i .

3.5 Dirijare (routing)

Dacă un procesor trimite un mesaj unui alt procesor, care nu este vecin cu el, mesajul trebuie să străbată o anume cale pentru a ajunge la destinație, trecând pe la procesoarele de pe această cale. Pentru procesoarele intermediare se pune problema alegerii unui canal de comunicație pe care să expedieze mesajul, cunoscând numai destinația și, desigur, canalul pe care au primit mesajul; aceasta este problema dirijării (routing).

S-ar părea că ne-am ocupat deja de chestiune, în secțiunea 3.3.1. Acolo însă am presupus că *toate* procesoarele dintr-o arhitectură paralelă au cunoștința de operația de comunicație. Aici însă adoptăm ipoteze diferite: *numai sursa și destinația* știu că va avea loc comunicația; celelalte procesoare, prin unitățile de comunicație de pe fiecare canal, așteaptă eventuale mesaje din orice direcție (și aici e ceva nou, până acum n-am avut nevoie de o astfel de primitivă de interogare); dacă apare un mesaj, unitatea de comunicație îl recepționează și, pe baza destinației lui, îl trimite mai departe sau îl reține. Aceste ipoteze sunt mai realiste; în primul rând, pe o arhitectură paralelă se pot executa simultan mai multe programe, pe diverse seturi de procesoare, disjuncte sau nu; în acest caz mesajele pot trece pe la procesoare care nu lucrează la același program ca sursa și destinația, deci nu au de unde ști de comunicația respectivă; în al doilea, în modurile comutare de circuit sau wormhole, unitatea de comunicație a unui procesor intermediar trebuie să rezolve rapid problema găsirii canalului pe care să transmită mai departe un mesaj și, dacă se poate, asincron cu unitatea centrală a nodului; de aceea, practic, unitățile de comunicație au asociate funcții de dirijare, care rezolvă problema.

Definiții. Asociem unei arhitecturi MIMD cu memorie distribuită un graf orientat, $X = (V, E)$, cu V mulțimea nodurilor (reprezentând procesoarele) și E mulțimea arcelor (reprezentând canalele de comunicație). Funcția de dirijare $f : E \times V \rightarrow E$ asociază, pentru un nod oarecare $v \in V$, unui canal de intrare $c = (u, v) \in E$, un canal de ieșire $e = (v, w) \in E$. Există mai multe atribute posibile ale funcțiilor de dirijare, dintre care enumerăm:

- dirijare *distribuită*: decizia de alegere a căii de ieșire este luată de procesorul intermediar prin care trece mesajul; în acest mod nu se poate cunoaște drumul deja parcurs de mesaj. Modul de lucru opus este dirijarea *centralizată*, în care un singur nod gestionează toate traseele.

- dirijare *statică*: există o funcție de dirijare care asociază fiecărui cuplu (sursă, destinație) o cale unică. Pot apărea probleme de congestie, dacă traficul nu e omogen.

- dirijare *adaptivă*: funcția de dirijare permite alegerea între mai multe căi (funcția se poate modifica în raport cu contextul), atunci când unele sunt deja ocupate, pentru a echilibra traficul. Funcția de dirijare se poate defini în acest caz astfel: $f : E \times V \times B \rightarrow E$, unde B reprezintă stările canalelor—ocupat sau nu.

- dirijare *pe calea cea mai scurtă*: traseul mesajului de la sursă la destinație are un număr minim de canale (atenție, aceasta nu implică unicitate). Funcția de dirijare poate fi văzută ca o scufundare de dilatare minimă a grafului complet în graful X .

Noțiunea de dirijare apare nu numai în calculatoarele paralele, ci și în orice rețea eterogenă de calculatoare. De regulă, în rețele, funcția de dirijare este realizată prin tabele, care asociază canale unor destinații; dirijarea este realizată prin program. În calculatoarele paralele sunt interesante funcțiile simple de dirijare, care pot fi implementate hardware, și pot deci găsi în timp foarte scurt canalul de ieșire.

Se pune acum problema dacă tot ceea ce am făcut în secțiunea 3.3.1 este util. Vom vedea imediat că aproape toate operațiile descrise acolo global, suportă foarte bine o descriere la nivel *local*. Ne vom ocupa de dirijarea distribuită și statică. Pe scurt, întrebarea la care vrem să răspundem este: putem construi o funcție de dirijare astfel încât, dacă procesorul sursă trimite un pachet pe fiecare canal, aceste pachete urmează căi arc-disjuncte până la destinație?

Dirijarea pe inel este foarte simplă: un procesor retransmite spre dreapta mesajele primite dinspre stânga și invers. Se vede imediat că orice algoritm de comunicare între două procesoare funcționează corect, indiferent dacă se utilizează o singură cale (ca în algoritmul 3.1), sau amândouă, indiferent dacă se utilizează pipeline sau nu.

Dirijare pe tor. Pentru a transmite pe cele patru căi din figura 3.8a, e suficient ca procesorul sursă să trimită cele patru pachete pe canalele sale de ieșire, iar celelalte procesoare să aplice următorul algoritm de dirijare (notăm cu id_x și id_y poziția procesorului curent pe orizontală, respectiv verticală):

ALGORITM 3.10 (dirijare pe tor; procesorul id determină canalul de ieșire co , pentru un mesaj venit pe canalul ci)

1. recepționează un mesaj pentru P_{xy} , pe canalul ci
2. **dacă** $id = (x, y)$ **atunci** păstrează mesajul
3. **altfel dacă** $id_x \neq x$ și $id_y \neq y$ **atunci**
 1. **dacă** $ci = \langle \text{stânga, dreapta, sus, jos} \rangle$ **atunci**
 1. $co \leftarrow \langle \text{dreapta, stânga, jos, sus} \rangle$
4. **altfel dacă** $id_x = x$ **atunci**
 1. **dacă** $ci = \langle \text{stânga, dreapta, sus, jos} \rangle$ **atunci**
 1. $co \leftarrow \langle \text{sus, jos, jos, sus} \rangle$
5. **altfel dacă** $id_y = y$ **atunci**
 1. **dacă** $ci = \langle \text{stânga, dreapta, sus, jos} \rangle$ **atunci**
 1. $co \leftarrow \langle \text{dreapta, stânga, stânga, dreapta} \rangle$

Unele instrucțiuni **dacă** au în condiție toate valorile posibile pentru ci ; variabilei co i se atribuie valoarea din poziția din lista din instrucțiunea sa identică cu poziția valorii actuale a lui ci ; de exemplu, dacă în 3.1 avem $ci = \text{sus}$, atunci în 3.1.1 se execută atribuirea $co \leftarrow \text{jos}$.

Desigur că algoritmul este executat în buclă infinită, mai sus fiind descrisă o singură iterație.

Din păcate, în cazul în care sursa și destinația sunt pe aceeași linie sau coloană, algoritmul 3.10 trebuie îmbunătățit; în forma aceasta el funcționează corect, însă

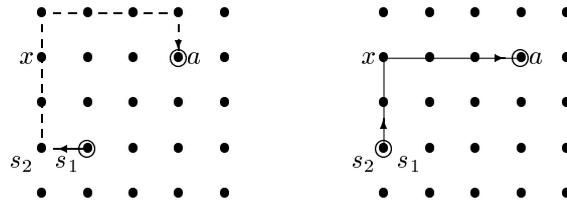


Figura 3.31: Două căi pentru care nu se poate construi o funcție de dirijare în nodul x , pe tor.

ineficient (două mesaje din patru ar face întâi un tur complet de inel înainte de a schimba direcția, pe căi ce s-ar suprapune cu primele două). Pentru că algoritmul 3.10 este practic foarte bun, e mai bine ca procesorul sursă să rezolve el problema, netrimțând decât două mesaje.

Și tot din păcate, pentru variantele de căi arc-disjuncte din figura 3.8bc, nu se pot construi funcții de dirijare. Un contraexemplu imediat pentru cazul 3.8b este prezentat în figura 3.31; în primul desen, nodul s_1 trimite un mesaj către a , începând spre stânga; în al doilea, s_2 trimite un mesaj, tot spre a , începând în sus; căile pe care circulă aceste mesaje trec prin nodul x , folosind același canal de intrare, dar canale de ieșire diferite. Numai că, printr-o funcție de dirijare, se asociază unui canal de intrare și unei destinații, un singur canal de ieșire. Evident, singura soluție aici ar fi ca funcția de dirijare să țină seama și de adresa sursei mesajului; în acest fel ne îndepărtăm de modul clasic de a privi dirijarea.

Dirijare pe hipercub. Pe hipercub, ca întotdeauna, soluțiile sunt mai elegante. O primă idee e de a dirija orice pachet pe calea principală. Un procesor oarecare trimite un mesaj pe dimensiunea cu cel mai mic număr de ordine în care adresa sa diferă de cea a destinatarului; deci, co este poziția bitului de 1 cel mai nesemnificativ din $id \oplus a$ (P_a fiind destinația). De exemplu, dacă nodul 10001 recepționează un mesaj pentru 01011, indiferent pe ce canal, îl trimite mai departe pe canalul 1. Faptul că nu se ține seama de canalul de intrare sugerează că nu întotdeauna se pot obține d căi disjuncte între sursă și destinație. De exemplu, dacă 000 trimite către 110, începând cu canalul 0, deci lui 001, acesta va retrimite înapoi mesajul lui 000, apoi acesta circulând pe ruta 000, 010, 110, deci pe calea începând din 000 pe dimensiunea 1; vor fi numai două căi arc-disjuncte între 000 și 110.

Defectul acesta este remediat în dirijarea *rotativă*. Dacă un procesor primește un mesaj pe ci , atunci îl va expedia pe primul canal aflat într-o dimensiune mai mare, în care adresa sa diferă de cea a destinatarului, sau pe cel în cea mai mică astfel de dimensiune, dacă nu există una mai mare; altfel zis, co este poziția bitului de 1 din $id \oplus a$, cea mai aproape prin stânga (ciclic) de ci . De exemplu, dacă 10001 primește un mesaj pentru 01011 pe canalul 2, îl va retrimite pe canalul 3; dacă $ci = 4$, atunci $co \leftarrow 0$. Căile sunt întotdeauna arc-disjuncte, dacă sursa trimite mesaje spre aceeași

destinație pe toate canalele. Dacă 000 trimite spre 110, căile vor fi: prima 000, 001, 011, 111, 110, a doua 000, 010, 110, a treia 000, 100, 110.

Pentru comunicarea dintre procesoarele 0 și $2^j - 1$ se regăsesc exact căile din algoritmul 3.2; primele j căi sunt calea principală și rotațiile ei; pentru celelalte, dacă se pornește pe dimensiunea $l \geq j$, următorul arc, pornind din nodul 2^l , se află în dimensiunea 0 (nu există biți diferiți în poziții mai mari ca l , în adresele 2^l și $2^j - 1$), al treilea în dimensiunea 1, etc., deci se urmează calea principală dintre 2^l și $2^j - 1$. În cazul general, căile dintre sursă și destinație nu sunt însă identice cu cele din observația din rezolvarea problemei 3.3.5.

Probleme

P 3.5.1 Scrieți un algoritm de dirijare pe grilă.

P 3.5.2 Considerând dirijarea rotativă la dreapta, în care co este poziția bitului de 1 cel mai apropiat spre dreapta (ciclic) de poziția ci , în $id \oplus a$, asigură aceasta d căi arc-disjuncte între sursă și destinație? Care sunt căile între 000 și 110?

P 3.5.3 Propuneți un algoritm de dirijare pe graful complet și analizați-l.

Capitolul 4

Algoritmi paraleli: concepție și apreciere

Acest capitol este dedicat modului în care se studiază un algoritm paralel, de la formularea problemei până la aprecierea performanțelor obținute. Nu trebuie însă să vă așteptați la un panaceu universal, ci numai la unele instrumente, în general folosite, dar cu aplicabilitate diferită de la caz la caz.

Formularea problemei de rezolvat se face la fel ca în cazul secvențial, adăugând acum specificațiile arhitecturii țintă. Proiectarea unui algoritm de rezolvare pornește de obicei de la algoritmi secvențiali existenți. Se studiază în ce măsură aceștia oferă un grad de paralelism satisfăcător și, dacă e cazul, se scrie algoritmul paralel corespunzător. Se estimează apoi performanțele algoritmului, după care acesta se implementează și se testează experimental; desigur, e de dorit ca estimarea să fie cât mai adecvată, pentru a economisi din timpul de calculator necesar testării (încă relativ scump pentru calculatoarele paralele). Dacă algoritmi secvențiali clasici nu dau satisfacție, trebuie încercată găsirea unor noi, eventual nu atât de valoroși secvențiali, dar cu potențial mare de paralelizare.

De obicei, algoritmi secvențiali au forme consacrate; apar însă și o serie de variante ale lor, obținute printr-o altă *ordonare* a aceluiași operații; ceea ce e important însă, și comun tuturor variantelor, este ideea unică ascunsă în spatele lor. Ceea ce numim algoritm paralel este tot o formă a aceleiași idei, dar în care sunt evidențiate operațiile ce pot decurge în paralel; de asemenea, există o multitudine de variante ale unui algoritm paralel, chiar pe o aceeași arhitectură; toate aceste forme permit descrieri secvențiale, de multe ori echivalente. Pentru a face distincția terminologică necesară, vom numi algoritm secvențial, așa cum se utilizează curent, combinația idee + descriere secvențială; la fel, algoritm paralel = idee + descriere paralelă pe o arhitectură precizată (eventual idealizată). Am vrut astfel să subliniem că ideea este esențialmente comună, deosebirea între secvențial și paralel apărând îndeosebi în descriere. În această lucrare vor fi prezentate atât idei noi, născute având ca scop

implementarea paralelă, cât și descrieri noi ale unor idei vechi. În ambele cazuri se vorbește însă despre algoritmi paraleli noi; de multe ori există justificarea: o descriere nouă poate încorpora multă originalitate și poate da un aspect nebănuț ideii originale. Vom conserva această mică ambiguitate de limbaj, intrată de-acum în uz.

4.1 Performanțele unui algoritm paralel

Vom începe cu etapa finală, aprecierea unui algoritm, pentru că astfel putem fixa mai bine obiectivele de avut în vedere la proiectare. În tot ce urmează, ne vom referi la o unică problemă, rezolvată prin unul sau mai mulți algoritmi, paraleli și secvențiali. Trebuie deci să fie clar că nu se apreciază un algoritm în sine, izolat de context, ci eficacitatea cu care el rezolvă problema în cauză.

4.1.1 Criterii de bază

Țimpul de execuție. Principalul criteriu de apreciere a performanțelor unui algoritm paralel este Țimpul de execuție. Vom nota acest Țimp cu $T(n, p)$, arătând prin aceasta dependența sa de cel puțin doi factori; n este dimensiunea problemei, reprezentând într-un anumit fel numărul datelor de intrare, de exemplu (cazul cel mai curent), dimensiunea unor matrice sau vectori; p este, ca și până acum, numărul de procesoare. Vom adăuga acestei notații diverși indici, arătând de obicei algoritmul pentru care este valabil Țimpul respectiv. Toate celelalte criterii sunt derivate ale Țimpului de execuție.

Aprecierea Țimpului de execuție se poate face teoretic sau experimental; în ambele cazuri se consideră ca Țimp de execuție intervalul dintre momentul în care *primul* procesor începe lucrul și cel în care *ultimul* procesor termină lucrul.

Pentru un procesor, Țimpul de lucru este consumat nu numai pentru calculele pe care le efectuează, ci și pentru comunicație (cu alte procesoare sau pentru transferuri între memoria locală și cea comună); de asemenea, pot exista și eventuale Țimpuri morți, necesari sincronizării cu alte procesoare (așteptarea unor mesaje, a accesului la secvențe critice). Pentru calcule vom considera, în general, ca unitate de măsură Țimpul necesar execuției unui flop, notat cu α . Pentru Țimpii de comunicație sunt prezentate mai multe modele în capitolul anterior.

În ce privește măsurarea experimentală a Țimpilor de execuție există probleme de implementare care țin de arhitectura calculatorului. În general, în arhitecturile MIMD cu memorie distribuită, se atribuie unui procesor sarcina de a trimite mesaje de pornire tuturor celorlalte și apoi, după terminarea execuției părții din algoritm atribuite lui, de a recepționa mesaje de terminare de la celelalte procesoare; el va măsura Țimpul dintre emiterea primului mesaj și recepționarea ultimului, o aproximare superioară Țimpului de execuție. Sau, soluție convenabilă și pentru arhitecturile cu memorie comună, algoritmul începe și se termină cu câte o sincronizare, Țimpii inițial și final

fiind măsurate imediat după acestea. Dacă există un ceas global, fiecare procesor memorează momentul de început și pe cel de sfârșit.

Menționăm că, de cele mai multe ori, timpul necesar traficului datelor între o memorie externă și memoriile locale ale procesoarelor (în cazul memoriei distribuite) sau memoria comună este neglijat. Toate aprecierile timpului de execuție se fac începând dintr-un moment în care datele se află într-un loc accesibil direct procesoarelor și până când rezultatele sunt depuse tot acolo. Aceasta este de obicei corect atunci când se compară doi algoritmi paraleli, dar poate fi incorect la comparația între un algoritm paralel și unul secvențial.

Accelerarea. Modul uzual în care este gândită valoarea unui algoritm paralel poate fi formulat astfel: pe un calculator cu p procesoare va fi acesta mai rapid de p ori, sau de aproape p ori, decât algoritmul secvențial (executat pe un singur procesor)? Cu alte cuvinte: se obține un câștig semnificativ de viteză în rezolvarea problemei dacă se utilizează algoritmul paralel? Este acest câștig pe măsura puterii de calcul folosite? Dacă răspunsul este negativ—ceea ce înseamnă fie că algoritmul este profund secvențial, fie că nu a fost găsită o variantă paralelă bună—atunci cercetarea trebuie continuată. Pe de altă parte, găsirea unui algoritm paralel foarte avantajos, comparat cu varianta sa secvențială (cum am spus, orice algoritm paralel se poate descrie secvențial), nu înseamnă și sfârșitul căutării; deoarece se urmărește rezolvarea unei probleme, compararea algoritmului paralel se face cu *cel mai rapid* algoritm secvențial din clasa celor care rezolvă problema.

Parametrul care cuantifică acest mod de a privi performanța este *accelerarea* (speed-up, creștere de viteză) algoritmului paralel X , definită prin relația

$$S_X(n, p) = T_S(n)/T_X(n, p),$$

în care $T_S(n)$ este timpul de execuție al celui mai rapid algoritm secvențial care rezolvă problema respectivă, iar $T_X(n, p)$ timpul de execuție al algoritmului paralel. Accelerarea depinde, deseori semnificativ, de n și p ; în general, pentru p constant, aceasta crește o dată cu n (vom vedea ceva mai târziu cum); de asemenea, pentru n constant, crește o dată cu p , dar fără a depăși o valoare maximă. În general, accelerarea este mărginită astfel:

$$1 \leq S(n, p) \leq p,$$

valoarea din dreapta fiind considerată ideală. Sunt posibile și valori mai mici ca 1, de exemplu în cazul în care algoritmul paralel se desfășoară de fapt secvențial pe ansamblul celor p procesoare, pierzându-se și timp suplimentar pentru comunicație; cazul $S(n, p) > p$, aparent imposibil, poate apărea în practică; vom exemplifica puțin mai târziu.

Eficiența. Parametrul cel mai folosit în aprecierea unui algoritm paralel X este *eficiența*, definită astfel:

$$\varepsilon_X(n, p) = \frac{S_X(n, p)}{p} = \frac{T_S(n)}{pT_X(n, p)}. \quad (4.1)$$

Eficiența are (aproape) întotdeauna o valoare subunitară; un algoritm este cu atât mai valoros cu cât eficiența sa este mai aproape de 1. În general, eficiența crește o dată cu n și scade o dată cu p . Vom spune că un algoritm paralel are eficiența *asimptotic egală cu 1*, dacă, pentru orice p constant, avem

$$\lim_{n \rightarrow \infty} \varepsilon(n, p) = 1.$$

Desigur că, pentru un anumit algoritm, limita de mai sus poate fi dovedită numai teoretic; rezultatele experimentale o pot doar sugera.

Vom numi *cost* (sau volum de lucru) al unui algoritm paralel mărimea

$$W(n, p) = pT(n, p),$$

reflectând timpul de execuție cumulat pe cele p procesoare. El poate fi comparat cu costul algoritmului secvențial; diferența $W(n, p) - W(n, 1)$ (unde $W(n, 1) \equiv T_S(n)$) arată timpul total suplimentar consumat în execuția paralelă.

Pentru varianta paralelă XP a unui algoritm secvențial XS , o formulă de același tip cu (4.1) definește *eficiența relativă* a paralelizării:

$$\varepsilon'_X(n, p) = \frac{T_{XS}(n, p)}{pT_{XP}(n, p)}.$$

Aceasta este utilă pentru a aprecia cât de bună este paralelizarea algoritmului, fără a ne interesa cât de bine este rezolvată problema.

În general, vom nota $S(p)$, $\varepsilon(p)$ etc., ignorând dimensiunea problemei, ori de câte ori va fi posibil.

4.1.2 Despre overhead

Definiție și clasificare. Am menționat deja că pentru orice algoritm paralel eficiența este subunitară; este interesantă studierea cauzelor care contribuie la ineficiență ($1 - \varepsilon$). Denumirea generică a factorilor care micșorează eficiența este de *overhead*. Prezentăm în continuare câteva tipuri generale de overhead, valabile pentru toți algoritmii (se pot defini și tipuri specifice, în funcție de natura problemei):

- overhead *algoritmice*, datorat părților intrinsec secvențiale ale algoritmului—părți care nu se pot paraleliza, sau datorat unor operații suplimentare efectuate în algoritmul paralel față de cel mai rapid algoritm secvențial.
- overhead datorat *încărcării inegale* a procesoarelor; cazul ideal este cel în care calculele sunt repartizate uniform procesoarelor, astfel încât fiecare să aibă același număr de operații de executat; dacă unele procesoare vor avea mai mult de lucru decât altele, este evident că timpul total de execuție va crește.
- overhead datorat *comunicației*; aceasta reprezintă o activitate cu totul nouă față de cazul secvențial, deci orice comunicație contribuie la ineficiența algoritmului.

- overhead *software*; acesta apare datorită repetării unor operații ce se execută o singură dată în cazul secvențial, pe fiecare procesor al arhitecturii paralele. Exemplul cel mai simplu este cel al deciziilor. Există, în plus, anumite operații suplimentare datorate implementării software a algoritmului—decizii în funcție de adresa procesorului, aritmetică întreagă repetată pe mai multe procesoare etc.

Estimarea overhead-ului. Dacă estimarea teoretică a fiecărui tip de overhead poate fi făcută relativ ușor, în schimb determinarea experimentală nu este deloc simplă; aceasta e necesară deseori, pentru a confirma estimările teoretice, care nu pot ține cont de toți factorii posibili. Prezentăm în continuare un mod de a face determinarea overhead-urilor. Pentru a simplifica notația vom omite în continuare dependența de n și de p , presupuse fixe, a diverselor mărimi definite mai sus.

Notăm cu T_{nc} timpul de execuție al unui algoritm în care nu se efectuează comunicațiile și folosim un model simplificat pentru expresia lui T_p —timpul de execuție pentru un algoritm paralel:

$$T_p = T_a + T_c + T_l,$$

unde T_a este timpul pentru operații aritmetice (în cazul unei încărcări echilibrate a procesoarelor, atunci $T_a \approx T_s/p$ unde T_s este timpul secvențial—de execuție pe un singur procesor), $T_c = T_p - T_{nc}$ este timpul consumat pentru comunicație, iar $T_l = T_{nc} - T_a$ timpul pierdut în execuția paralelă din cauza încărcării inegale a procesoarelor, specificului implementării software, etc. Astfel se pot calcula două tipuri de overhead (pentru detalii vezi [12]):

1. overhead datorat comunicației: $O_c = 1 - T_{nc}/T_p$.

De fapt, deoarece uneori protocolul de comunicație implică sincronizare, o parte nemăsurabilă (pe care o neglijăm) a acestui overhead se datorează timpului mort de așteptare a sincronizării. De asemenea, alți timpi morți pot apărea atunci când un procesor așteaptă de la celelalte date de prelucrat și nu mai are nici o operație de efectuat; în acest caz, cauza overhead-ului este în primul rând algoritmică.

2. overhead datorat încărcării inegale, software, algoritmic: $O_l = (T_{nc} - \frac{T_s}{p})/T_p$.

4.1.3 Legea lui Amdahl

Legea lui Amdahl. Să analizăm puțin comportarea generală a unui algoritm în funcție de numărul de procesoare p , făcând o ipoteză simplificatoare; presupunem că algoritmul poate fi împărțit în două părți: una pur secvențială, al cărei cost este W_1 , și o alta perfect paralelă, indiferent de numărul de procesoare, al cărei cost este W_P . Atunci timpul de execuție este

$$T(p) \equiv T_1 + T_P = W_1 + W_P/p,$$

unde $T_1 = W_1$ și $T_P = W_P/p$, doar operațiile din partea perfect paralelă fiind repartizate între procesoare. Timpul de execuție are aspectul din figura 4.1a. Accelerarea

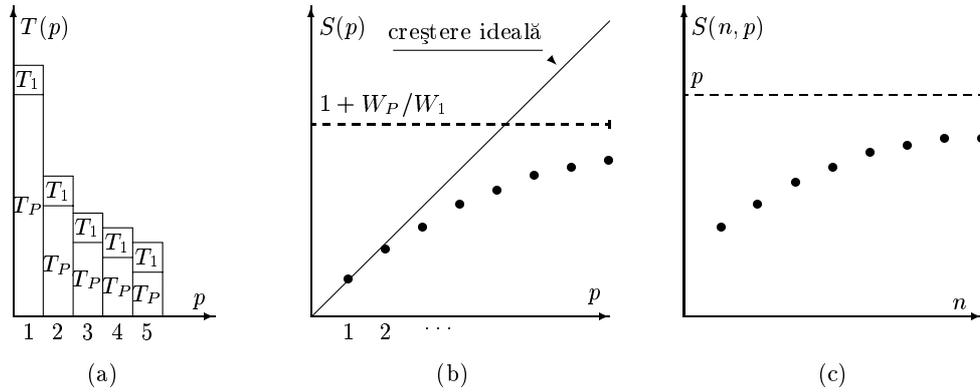


Figura 4.1: Evoluții specifice pentru: (a) Timpul de execuție al unui algoritm paralel funcție de numărul de procesoare p ; (b) Accelerarea funcție de p ; (c) Accelerarea funcție de dimensiunea problemei.

are expresia

$$S(p) = \frac{W_1 + W_P}{W_1 + W_P/p} < 1 + \frac{W_P}{W_1}, \tag{4.2}$$

unoscută sub numele de *legea lui Amdahl*, și o evoluție de genul celei din figura 4.1b. Așadar, accelerarea are o limită inerentă; din fericire, pentru mulți algoritmi W_1 este modest.

Discuție. În cele de mai sus am presupus implicit că raportul W_1/W_P nu depinde de dimensiunea problemei, ceea ce nu se întâmplă de obicei în practică. Să presupunem acum că W_1 depinde de n altfel decât W_P , anume raportul W_p/W_1 crescând o dată cu n ; de exemplu, $W_1 = O(n)$, dar $W_P = O(n^2)$; deci, cu cât crește n , cu atât valoarea maximă a creșterii de viteză din (4.2) este mai mare. Așadar, pentru p fixat, o alură tipică a creșterii de viteză în funcție de dimensiunea problemei este cea din figura 4.1c; o evoluție asemănătoare are și eficiența, valoarea maximă p fiind înlocuită cu 1. Dacă eficiența algoritmului este asimptotic egală cu 1, atunci accelerarea tinde către orizontala de ordonată p .

Să definim acum o altă mărime destul de utilizată, *performanța temporală*, $R(n, p) = 1/T(n, p)$; ea se folosește de multe ori pentru că este cu atât mai mare cu cât algoritmul se execută mai rapid, ceea ce e mai aproape de obiceiurile noastre de a aprecia performanța. De asemenea, pentru n fixat, îi este asociată o altă formă a legii lui Amdahl, anume

$$R(p) = \frac{R_\infty}{1 + p_{1/2}/p}, \tag{4.3}$$

în care $R_\infty = 1/W_1$ este performanța asimptotică, iar $p_{1/2} = W_P/W_1$ numărul de procesoare pentru care se atinge jumătate din performanța asimptotică (înlocuiți

aceste valori în (4.3), pentru verificare). Și pentru accelerare se poate scrie o relație de tipul (4.3) (calculați valorile parametrilor!).

Superliniaritate. Cum se poate ajunge la o accelerare mai mare decât p ? Algoritmii cu această proprietate, deci cu $T(n, p) < T(n, 1)/p$, sunt impropriu numiți *superliniari*; existența lor contrazice intuiția și, desigur, legea lui Amdahl (în (4.2), ar trebui ca $W_1 < 0$, ceea ce e imposibil, și $W_1 + W_P/p > 0$). Până acum ne-am referit numai la timp de execuție, atunci când a fost vorba de evaluarea unui algoritm. Să ne referim acum la numărul de operații; în principiu, n-ar trebui să existe nici o diferență; mai sus am definit timpul de execuție α al unui flop ca pe o mărime fixă (și așa o vom face și după acest paragraf).

Superliniaritatea poate proveni din două motive: fie numărul total de operații (costul), în paralel, este mai mic decât cel secvențial; fie o operație se execută mai rapid în contextul paralel, față de cel secvențial. Ambele par obscure și greu de imaginat. Pentru primul, nu vom da detalii, lăsând tema deschisă; de altfel, există destule controverse; ca exemplu (destul de dubios), imaginați-vă simularea unui mediu de programare paralel: secvențial, paralelismul trebuie simulat, în paralel el este parțial real (pe un procesor se pot totuși simula mai multe), deci e simulat mai rapid. Pentru al doilea însă, lucrurile sunt mai clare; timpul de execuție al unui algoritm nu depinde numai de numărul de flop, ci și, de exemplu, de numărul de referințe la memorie (pentru citirea datelor și a instrucțiunilor); presupunând că un nod are o arhitectură cu memorie ierarhică, precum cea din figura 1.5, timpul necesar unui acces la memorie este mai mic dacă data se află în memoria cache; ori, atunci când un algoritm se execută secvențial, există o singură memorie cache; la execuția paralelă, pentru aceeași cantitate totală de date, memoria cache este de p ori mai mare; deci, mult mai multe date vor avea loc în memoria cache, și astfel un acces la memorie va dura în medie mai puțin. Aceleași argumente sunt valabile pentru nivelul ierarhic următor de memorie, adică atunci când se pune problema transferurilor (repetate) între memoria principală și memoria externă (în contextul utilizării memoriei virtuale).

4.1.4 Alte criterii

Complexitate paralelă. Încercăm să ne apropiem acum de caracterizări mai teoretice ale unui algoritm. În general, vom înțelege prin *complexitate* a unui algoritm numărul de operații necesare execuției lui; pentru un algoritm paralel, este vorba de numărul de operații care pot fi executate de un singur procesor, între începutul și sfârșitul execuției algoritmului; deci, mai multe operații executate simultan de mai multe procesoare sunt contorizate (după cum e normal), ca una singură a algoritmului paralel; de asemenea, sunt echivalați în număr de operații eventualii timpi morți ai unui procesor. Vom face distincție între complexitatea aritmetică și cea a comunicației, cea de-a doua referindu-se la numărul de elemente comunicate (în aceleași condiții ca pentru numărul de operații, după cum am făcut-o deja în capitolul precedent). Deoarece între complexitate și timpul de execuție al algoritmului este, în principiu, o relație de proporționalitate, cele două noțiuni vor fi câteodată (dar cât

mai rar) folosite una în locul celeilalte.

Fiecare problemă are un număr minim de operații necesar rezolvării ei în paralel, cu un număr stabilit de procesoare, fie că există, fie că nu există un algoritm care să atingă această limită. Vom numi complexitate paralelă a *problemei*, numărul de operații necesar rezolvării ei, cu un număr oricât de mare de procesoare și neglijând comunicația (folosind deci toate mijloacele posibile); notația consacrată este $\Omega(f(n))$, unde $f(n)$ este o funcție simplă; $\Omega(\log n)$, de exemplu, poate însemna la fel de bine $\log n$ sau $3 \log n$; important este ordinul de mărime, și nu eventualii coeficienți. Complexitatea unui *algoritm* a fost definită mai sus și se notează cu $O(f(n))$; evident, ea este mai mare sau egală cu complexitatea problemei.

Vom folosi într-un mod mai neriguros, dar sugestiv notația $O(\cdot)$. De exemplu, deși $O(n)$ și $O(n/p)$ pot avea teoretic aceeași expresie, vom sublinia astfel influența lui p asupra complexității; sau, deși $O(n)$ și $O(n + \log n)$ sunt identice, a doua notație evidențiază mărimea termenului al doilea dintr-o sumă, pentru a compara sume în care primul termen e identic.

Algoritmi optimați, algoritmi eficienți. Fie $\text{polylog}(n) = \bigcup_{k>0} O(\log^k n)$. Fie S o problemă pentru care cel mai rapid algoritm secvențial se execută în timp $T(n)$.

DEFINIȚIA 4.1 Un algoritm paralel care rezolvă S și care se execută în timpul $t(n)$ pe $p(n)$ procesoare se numește *optimal* dacă:

- a) $t(n) = \text{polylog}(n)$, și
- b) volumul de lucru $W(n) = p(n) \cdot t(n)$ este $O(T(n))$.

Cu alte cuvinte, algoritmul este optimal dacă timpul său de execuție este logaritmic și dacă volumul său de lucru este de același ordin de mărime cu cel al algoritmului secvențial (nu se execută semnificativ mai multe operații în algoritmul paralel, pe ansamblul procesoarelor, față de cel secvențial). După cum se vede, notând $p(n)$ am lăsat să se înțeleagă că pentru atingerea optimalității, numărul de procesoare poate depinde de dimensiunea problemei (aceasta e specific unei abordări teoretice; practic, desigur că p este cel oferit de arhitectura paralelă).

DEFINIȚIA 4.2 Un algoritm paralel este *eficient* dacă condiția b) din definiția 4.1 se înlocuiește cu:

- b') volumul de lucru este $W(n) = T(n) \cdot \text{polylog}(n)$.

Un algoritm eficient are un înalt grad de paralelism, dar accelerarea nu este atinsă în mod optimal, ci cu prețul unor operații suplimentare, ceea ce se reflectă într-un cost de un factor logaritmic de ori mai mare decât cel secvențial.

Alte criterii în afara complexității. În toată această secțiune nu ne-am referit decât la un singur criteriu de apreciere (și comparație) a algoritmilor paraleli: timpul de execuție (sau, echivalent, numărul de operații, complexitatea). Nu trebuie uitate alte criterii, utilizate la fel și în algoritmica secvențială; poate fi vorba despre precizia cu care se obțin soluțiile, stabilitatea numerică a algoritmului, memoria necesară etc.

În continuare, vor fi menționate și acestea; decizia alegerii unui algoritm trebuie să depindă și de ele.

4.2 Tehnici generale de paralelizare

În drumul de la un algoritm descris în formă secvențială, până la scrierea algoritmului pentru un anumit calculator paralel, trebuie parcurse întotdeauna următoarele etape:

- gruparea în taskuri a operațiilor ce trebuie executate.
- identificarea ordinii în care se pot executa taskurile, precum și a taskurilor ce se pot executa în paralel; de un mare ajutor în această etapă este construirea grafului de precedență a taskurilor.
- planificarea (scheduling) execuției taskurilor, adică stabilirea unui mod de a descrie procesorul pe care se execută fiecare task.

Uneori, în prezentarea algoritmilor, etapele nu sunt separate în mod foarte explicit, în special a doua de cea de-a treia; de asemenea, de multe ori, nu se prezintă constrângerile de precedență a taskurilor. Vom expune, totuși, pe larg, noțiunile generale legate de etapele de mai sus, pentru a crea un cadru sistematic de abordare a problemei paralelizării unui algoritm. (Trebuie precizat că există în literatură numeroase articole ce conțin abordări brutale ale problemelor, justificate finalmente prin buna eficiență experimentală a algoritmului studiat. Vom folosi și noi această tactică, pentru a nu lungi, de multe ori inutil, expunerea.)

În plus față de aceste etape, mai apar probleme legate de comunicație: trecerea de la comunicarea între taskuri, operație logică, la comunicarea între procesoare, operație fizică.

4.2.1 Partiționarea în taskuri

Task. Un task este definit ca o unitate indivizibilă de operații, specificată numai în termeni ai comportării sale exterioare (intrări, ieșiri, timp de execuție etc.). Noțiunea de operație este elementară și nu va fi definită riguros; operațiile luate în considerare ca durată—mai ales în algoritmii de calcul numeric—sunt cele aritmetice (adunare, înmulțire etc.).

Operațiile din interiorul unui task se execută secvențial.

Formarea taskurilor. Gruparea operațiilor în taskuri este o problemă delicată; nu există rețete universale pentru a face aceasta. Sunt o serie de criterii care pot ajuta în formarea taskurilor, dintre care unele euristice; de multe ori validarea se face experimental, alegându-se din mai multe variante.

Un prim criteriu se referă la ordinea în care se pot executa taskurile: părțile eminentemente secvențiale vor fi grupate în același task, în timp ce operații ce se pot

executa în paralel se vor introduce, în general, în taskuri diferite. Mai multe despre aceasta în secțiunea următoare.

Arhitectura unui procesor din componența calculatorului paralel poate influența și ea; dacă procesorul este vectorial, taskurile vor fi alese astfel încât să cuprindă cât mai multe operații cu vectori; dacă procesorul are memorie ierarhică, se vor crea taskuri pentru care datele să aibă dimensiuni apropiate de cele ale memoriei rapide și să fie folosite cât mai intens.

În fine, de mare însemnătate este ponderea comunicațiilor în raport cu calculele în algoritmul ce se implementează; în general, trebuie minimizate cantitatea de informație comunicată și, eventual, numărul de mesaje. Trebuie luată în considerare eficiența comunicațiilor în raport cu calculele pe arhitectura țintă. De asemenea, se vor utiliza la maximum primitivele de comunicație oferite de către calculator și sistemul de operare.

BLAS. Pentru algoritmi de calcul numeric, filosofia BLAS (Basic Linear Algebra Subprograms) este foarte utilă și flexibilă în alegerea taskurilor; ideea fundamentală este de a izola câteva rutine de bază—care implementează foarte eficient operațiile cele mai curențe, de exemplu înmulțirea matrice-vector sau matrice-matrice, sau rezolvarea de sisteme liniare triunghiulare—folosite apoi în descrierea algoritmilor; taskurilor vor fi compuse atunci din apeluri la astfel de rutine, și nu din operații elementare.

În 1990, Dongarra et al [10] au propus rutinele de bază pentru nivelul 3 BLAS. Acestea conțin operații matrice-matrice, de genul $C = \alpha AB + \beta C$, care implică $O(n^3)$ flop, și pentru care datele au dimensiune $O(n^2)$; deci, rutinele de nivel 3 sunt dense în calcule, în raport cu comunicația (sau cu transferurile între memoria comună și cele locale). Folosirea, în interiorul taskurilor, a operațiilor la nivel de bloc de matrice, pe lângă faptul că permite o mai mare claritate a algoritmului, oferă posibilitatea alegerii dimensiunilor blocurilor în concordanță cu posibilitățile oferite de arhitectura țintă. În capitolul 6, vor fi prezentați și algoritmi la nivel de bloc.

Granularitate. Noțiunea fundamentală în ce privește dimensiunea taskurilor este cea a *granularității* algoritmilor. Un același algoritm poate fi scris folosind operații la nivel de element, caz în care este denumit cu granularitate *fină* (fine grain) sau la nivel de blocuri, fiind denumit cu granularitate *mare* (coarse grain). În primul caz taskurile au dimensiuni mici, fiind formate din câteva operații, în al doilea dimensiunile sunt mari. Subliniem încă o dată importanța caracteristicilor calculatorului țintă; de exemplu, pentru calculatoarele MIMD cu memorie distribuită uzuală, algoritmi cei mai eficienți au, în general, granularitate mare; motivul principal este structura ierarhică a memoriei, care favorizează operațiile cu blocuri de matrice; de asemenea, un alt motiv este faptul că o aceeași cantitate de date se comunică mai repede în puține mesaje de lungime mare, decât în mai multe de lungime mică.

4.2.2 Grafuri de precedență

Noțiunea de graf de precedență a taskurilor—introdusă în [8]—modelează dependențele dintre taskuri, la nivel de date de intrare-ieșire.

Definiții. Datele de intrare ale unui task T_j pot fi chiar date de intrare ale întregului algoritm sau pot proveni din datele de ieșire ale altor taskuri. În primul caz, T_j poate fi executat oricând; în al doilea, T_j trebuie să aștepte execuția taskurilor care îi furnizează datele de intrare. Deci, ordinea în care se execută taskurile trebuie să respecte anumite constrângeri de precedență, impuse de circulația datelor. Relația de precedență este notată cu \prec ; dacă, pentru taskurile T_i, T_j , $T_i \prec T_j$, atunci operațiile din T_j folosesc date prelucrate de T_i (și apoi, eventual, și de alte taskuri); deci, execuția taskului T_j poate începe numai după terminarea execuției taskului T_i . În graful de precedență a taskurilor, fiecărui nod îi este asociat un task și, implicit, un număr reprezentând durata de execuție a acestuia; arcele (orientate) reprezintă constrângerile de precedență; dacă $T_i \prec T_j$, atunci există o cale de la T_i la T_j . Acest graf este orientat și aciclic (DAG - Directed Acyclic Graph); existența unui ciclu ar însemna o relație de genul $T_i \prec T_i$, evident absurdă.

Un DAG seamănă cu un arbore; există totuși diferențe importante. Un DAG are mai multe "rădăcini", adică taskuri fără predecesor; acestea folosesc exclusiv date de intrare ale algoritmului. De asemenea, un nod poate avea mai mulți tați. Ieșirile taskurilor frunze sunt ieșirile întregului algoritm.

Vom prezenta în continuare câteva definiții privind DAG și, în particular, grafurile de precedență; pentru simplitate, vom presupune implicit că duratele de execuție ale tuturor taskurilor sunt egale (pentru mulți algoritmi aceasta este sau adevărat sau o bună aproximație). Multe dintre definiții au un echivalent asemănător pentru arbori.

- *lungimea* unei căi este suma timpilor de execuție asociați nodurilor de pe calea respectivă.

- *nivelul* se definește astfel: nodurile terminale (care nu au nici un succesor) au nivel 0; nivelul $l(v)$ al unui nod v este $l(v) = \max_i l(v_i) + 1$, unde v_i sunt succesorii imediați (fii) ai lui v . Uneori se asociază nivelul 0 nodurilor inițiale (care nu au nici un predecesor) și, în relația de mai sus, v_i sunt predecesorii imediați ("tații") ai nodului v .

- *înălțimea* (adâncimea) unui DAG este lungimea celei mai lungi căi.

- *lățimea* unui DAG este $\max |L_k|$, unde L_k este mulțimea nodurilor având nivelul k , iar $|L_k|$ este cardinalul acestei mulțimi.

Reducerea grafului de precedență. Dacă un graf de precedență ar cuprinde arce pentru toate perechile de noduri care sunt în relație de precedență, atunci acest graf ar fi foarte dificil de folosit. De aceea, în măsura în care este posibil, din graf trebuie eliminate arcele care nu conțin informație necesară. Această operație se face de multe ori în mod natural, în analiza unui algoritm, când nu se studiază decât succesorii imediați ai unui task (cei care folosesc direct datele de ieșire ale aceluși task) și nu se rețin decât relațiile de precedență cu aceștia. O regulă de eliminare a arcelor inutile se bazează pe tranzitivitatea relației de precedență și pe o idee cât se poate de evidentă; dacă există o cale de la T_i la T_j trecând prin cel puțin un alt nod T_k , atunci arcul de la T_i la T_j se poate elimina; atenție, dacă există două căi între T_i și T_j , amândouă conținând încă cel puțin câte un alt nod, nu se poate elimina cea mai

scurtă dintre ele sau fie și numai un arc de pe aceasta.

Complexitatea și graful de precedență. Cu ajutorul grafului de precedență asociat unei partiționări în taskuri a operațiilor dintr-un algoritm se poate determina paralelismul intrinsec al partiționării. Presupunând disponibil un număr suficient de mare de procesoare și ignorând comunicațiile între taskuri, conflictele de acces la memoria comună și orice alte întârzieri, atunci timpul cel mai scurt în care se poate executa algoritmul este dat de înălțimea grafului. Într-adevăr, taskurile de pe orice cale trebuie executate secvențial, iar înălțimea este dată de cea mai lungă cale. Pentru a obține paralelismul intrinsec (teoretic) al algoritmului, se alege partiționarea în care fiecare operație elementară constituie un task. Înălțimea grafului este, în acest caz, *complexitatea paralelă* a algoritmului.

Lățimea grafului dă informații despre numărul necesar de procesoare pentru a atinge timpul cel mai scurt de execuție. Între taskurile de pe același nivel nu există constrângeri de precedență (dacă $T_i \prec T_j$, atunci $l(T_i) > l(T_j)$); deci toate taskurile de pe același nivel se pot executa în paralel. Deci, numărul de procesoare necesar pentru obținerea timpului minim de execuție este mai mic sau egal cu lățimea grafului. Dacă timpii de execuție ai taskurilor nu sunt egali, atunci problema numărului necesar de procesoare se complică, fiind nevoie de altă definiție pentru lățimea grafului.

În final subliniem încă o dată că noțiunea de graf de precedență este asociată unei partiționări în taskuri, și nu unui algoritm. Deci, fără o alegere adecvată a taskurilor, graful de precedență în sine nu poate aduce singur bunele performanțe.

4.2.3 Planificarea taskurilor

Prin definiție, fiind indivizibil, un task se execută de către un singur procesor. Se pune imediat întrebarea: pe ce procesor se va executa un anume task ?

DEFINIȚIA 4.3 Dată fiind o arhitectură cu p procesoare notate P_i , cu $i = 0 \dots p - 1$, *problema planificării taskurilor* este aceea de a găsi o funcție s definită pe mulțimea taskurilor, cu valori în mulțimea procesoarelor, semnificația relației $s(T_j) = P_i$ fiind că taskul T_j se va executa pe procesorul P_i .

În construirea funcției de planificare apar numeroase condiții legate de minimizarea timpului de execuție, a comunicațiilor între procesoare și, evident, de respectarea constrângerilor de precedență.

Tipuri de planificare. Planificarea poate fi de două feluri:

- *statică*, în care funcția s este cunoscută înaintea începerii execuției algoritmului.
- *dinamică*, în care funcția s se construiește pe măsură ce algoritmul se execută, atribuindu-se taskuri spre execuție unor procesoare libere. Nu se poate preciza, la începutul execuției, pe ce procesor se va executa un anume task.

Planificarea dinamică este mai flexibilă, ea fiind recomandată în condițiile unui sistem de operare concurrent; de asemenea, este unica viabilă în cazul defectării unui procesor (în acest din urmă caz execuția este blocată în cazul unei planificări statice). Pe de altă parte, planificarea dinamică este consumatoare de timp suplimentar, datorită programului care determină funcția s ; cu cât granularitatea este mai fină, cu atât costul planificării este mai important. În cazul arhitecturilor cu memorie partajată se folosește de obicei planificarea dinamică, iar pentru cele cu memorie distribuită planificarea statică. Motivul este legat de comunicație: în primul caz procesoarele sunt egal distanțate între ele, în timp ce în al doilea distanțele diferă; de exemplu se poate găsi o planificare statică în care comunicațiile să fie numai între vecini; cu o planificare dinamică, nu ar exista garanția eficienței comunicației. Totuși, o situație în care planificarea dinamică este evident mai utilă este cea în care durata taskurilor nu poate fi estimată dinainte; în acest caz, la o planificare statică pot apărea timpi morți mari. În această lucrare vom folosi exclusiv planificarea statică, ea reliefând mai bine paralelismul unui algoritm; de altfel, e mai ușor de obținut o planificare dinamică din una statică, decât invers; în plus, în majoritatea cazurilor, timpul de execuție al taskurilor va putea fi aproximat destul de bine.

Planificarea după listă. Într-un mod devenit clasic, problema planificării se reduce la problema construirii unei liste cuprinzând toate taskurile; apoi se alocă taskurile, în ordinea din listă, procesoarelor libere, cu condiția de a nu se planifica simultan decât taskuri independente (din punct de vedere al precedenței).

Prin reformulare, problema nu se simplifică deloc; din păcate ea este NP-completă pentru mai mult de 2 procesoare (nerezolvabilă în timp polinomial). Există numeroși algoritmi de planificare suboptimală—care determină rapid funcția de planificare s , dar care nu asigură timp minim de execuție al algoritmului planificat; se pot obține rezultate foarte bune pentru anumiți algoritmi, dar slabe pentru alții.

Nici pentru timpi egali de execuție a taskurilor problema planificării nu e simplă. Construirea listei, de la sfârșit spre început, în ordinea nivelelor, este o euristică foarte folosită, dar care nu duce obligatoriu la optim. Deci, ultimele în listă sunt nodurile de pe nivelul 0, imediat înaintea lor cele de pe nivelul 1 etc. Pentru a ordona nodurile de pe același nivel se folosește, de exemplu, criteriul numărului de succesori direcți: primul în listă este pus taskul cu cei mai mulți succesori direcți. Aceste criterii de ordonare sunt foarte intuitive; ordinea inversă a nivelelor e sugerată de relațiile de precedență, automat respectate în acest caz; e bine ca un task cu mai mulți succesori să fie executat cât mai repede posibil, pentru a da și acestora posibilitatea de execuție. Nu vom da alte detalii despre algoritmii de planificare, ci vom prezenta un exemplu simplu care sugerează dificultatea problemei.

Exemplu. Să considerăm graful de precedență a taskurilor din figura 4.2, în care cele șapte taskuri au durate egale de execuție (egale cu unitatea, pentru simplitate); taskul 7 are nivelul 0, taskurile 4, 5, 6 au nivelul 1, iar taskurile 1, 2, 3 au nivelul 2. Vom încerca planificarea taskurilor din acest graf pe o arhitectură cu două procesoare. Dacă utilizăm algoritmul sugerat mai sus, nu avem nici o prioritate între taskurile 1, 2, 3, sau între taskurile 4, 5, 6; le vom ordona deci la întâmplare. Lista 1234567

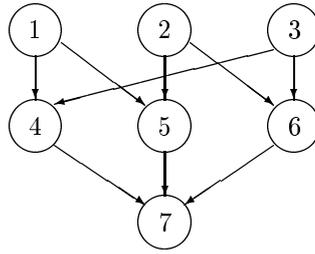


Figura 4.2: Un exemplu de graf de precedență.

a		b		c	
P_0	P_1	P_0	P_1	P_0	P_1
1	2	1	3	1	2
3		2	4	3	5
4	5	5	6	4	6
6		7		7	
7					

Tabelul 4.1: Planificări pe 2 procesoare pentru grafurile de precedență din figura 4.2.

conduce la planificarea din tabelul 4.1a, în care coloana pe care se află un task indică procesorul pe care se va executa acesta, iar linia momentul de timp la care va avea loc execuția; planificarea se construiește astfel: în prima etapă se execută taskurile 1 și 2; în a doua, 3, dar 4 nu, pentru că e succesori al lui 3 în grafurile de precedență; deci 4 se execută în etapa a treia, împreună cu 5; în etapa a patra taskul 6, dar nu și 7, care este succesor al lui 6; deci taskul 7 rămâne să se execute în etapa a cincea; timpul de execuție al întregului algoritmul este deci 5. Se vede imediat că ordinele 1324567 sau 1235467, care conduc la planificările din tabelele 4.1b și 4.1c, au timpul de execuție 4, care este, de altfel, cel optim. Dar dacă pe un astfel de graf sunt puține variante posibile, pe unul cu mii (milioane) de noduri, planificarea nu mai e deloc banală.

Presupunând că execuția algoritmului paralel are loc exact după planificările din tabelul 4.1, un mod curent de a prezenta situația ocupării procesoarelor este diagrama Gantt. În figura 4.3 sunt prezentate astfel de diagrame pentru primele două planificări din tabel; pe abscisă este reprezentat timpul (de multe ori omis, însă); zonele hașurate reprezintă perioadele în care procesoarele sunt ocupate, iar cele albe, timpurile morți.

Compromisul eficiență/timp de planificare. În ce privește efortul de calcul efectuat în vederea obținerii unei bune planificări trebuie făcută o precizare de multe ori neglijată. Pentru un algoritmul care va fi executat de puține ori, e suficientă și o planificare mediocră, dar obținută rapid. În schimb pentru un algoritmul pentru care se estimează un număr mare de execuții merită făcută investiția inițială a căutării unei

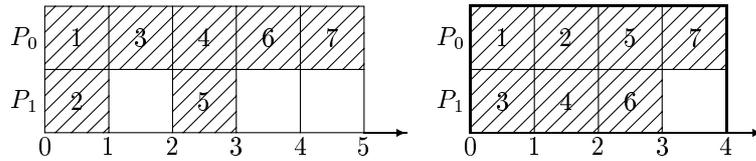


Figura 4.3: Diagrame Gantt pentru planificările din tabelul 4.1ab.

planificări apropiate de optim; aceasta cu atât mai mult cu cât execuția algoritmului se face pe un calculator paralel, unde timpul de calcul este scump, în timp ce planificarea se poate face pe orice calculator.

Principiul planificării (Brent). Am vorbit în secțiunea precedentă despre paralelismul intrinsec al unui algoritm (sau al unei partiționări în taskuri asociată algoritmului). Considerăm cazul în care toate taskurile au durată 1. Să presupunem că pentru a obține timpul minim de execuție t sunt necesare p_m procesoare. Care va fi timpul de execuție atunci când sunt disponibile doar $p < p_m$ procesoare? (De observat că acesta este cazul curent, numărul de procesoare al calculatorului țintă neavând legătură cu problema de rezolvat.)

Răspunsul imediat—cunoscut și sub numele de *principiul planificării al lui Brent*—este că timpul de execuție va fi de cel mult $\lceil p_m/p \rceil t$. Într-adevăr la un moment de timp oarecare se execută cel mult p_m taskuri, care sunt independente; folosind doar p procesoare, cele p_m taskuri se pot planifica oricum pe aceste procesoare, execuția lor durând maximum $\lceil p_m/p \rceil$; cum sunt t etape (corespunzând celor t niveluri din graficul de precedentă) în execuția algoritmului pe p_m procesoare, se obține rezultatul de mai sus.

Rafinând puțin prezentarea, să presupunem că pentru execuția etapei i sunt necesare x_i procesoare, deci $p_m = \max_i x_i$. Cu p procesoare timpul de execuție al pasului i va fi $\lceil x_i/p \rceil \leq x_i/p + 1$ și deci timpul total va fi mai mic decât $\lceil x/p \rceil + t$, unde $x = \sum_i x_i$.

Pe de altă parte, în ambele cazuri, o limită inferioară a timpului de execuție pe p procesoare este $t - 1 + \lceil p_m/p \rceil$; aceasta se obține când, în algoritmul cu paralelism maxim, într-o singură etapă se folosesc efectiv cele p_m procesoare (o astfel de etapă trebuie să existe, altfel nu se justifică numărul de procesoare), iar în restul mai puțin de p procesoare.

Vom ilustra principiul lui Brent în secțiunea următoare.

Efectul repartizării datelor. În arhitecturile cu memorie distribuită, la algoritmii de calcul numeric, de multe ori planificarea pornește de la o repartizare inițială a datelor în așa fel încât fiecare procesor să lucreze cât mai mult cu datele pe care le are (deci să se minimizeze comunicația) și încărcarea procesoarelor să fie echilibrată (overhead-ul de încărcare inegală să fie cât mai mic). Această abordare, care se bazează pe doar câteva variante—între care este ușor de ales—este facilitată de

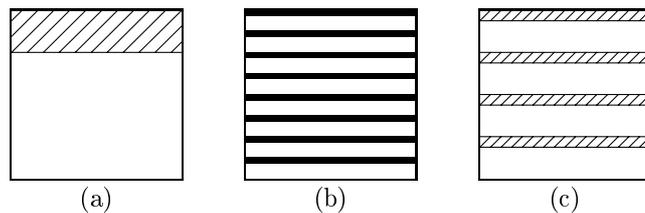


Figura 4.4: Elementele matricei A locale procesorului P_0 în repartizările: (a) bloc linii; (b) ciclic pe linii; (c) bloc ciclic pe linii.

aspectul algoritmilor, compuși în principal din bucle cu număr cunoscut de pași și omogeni din punct de vedere al organizării calculului.

De aceea considerăm necesară prezentarea aici a modurilor uzuale de repartizare a elementelor unei matrice între procesoarele unei arhitecturi paralele. Se dorește în general ca datele să fie repartizate procesoarelor în mod *echitabil* și *fără redundanțe* (aceeași dată în memoria mai multor procesoare), astfel încât acestea să ocupe cât mai puțin spațiu în memoriile locale și deci să se poată rezolva probleme de dimensiuni cât mai mari. Fie A o matrice cu n linii; fiecare dintre cele p procesoare va trebui deci să primească $m = n/p$ linii; pentru simplitate presupunem că p divide n . Repartizările următoare sunt naturale pe un inel (dar se pot avea în vedere și alte arhitecturi) și sunt ilustrate în figura 4.4:

- *bloc linii* în care primele m linii sunt locale procesorului P_0 , următoarele m lui P_1 , etc. În general, procesorul P_i posedă liniile cu numerele de la im la $(i+1)m - 1$ (liniile sunt numerotate începând de la 0).
- *ciclic pe linii* în care procesorului P_i îi aparțin liniile cu numerele $jp + i$, cu $j = 0, 1, \dots, m - 1$; sau, altfel spus, liniile l pentru care $l \bmod p = i$ (o imagine mai plastică: liniile sunt împărțite procesoarelor precum cărțile de joc dintr-un pachet unor jucători așezați în cerc).
- *bloc ciclic pe linii* în care fiecare procesor are câte d linii consecutive; matricea inițială se împarte în submatrice de câte pd linii, fiecare submatrice repartizându-se bloc linii pe procesoare.

Desigur că cele trei moduri descrise mai sus se pot folosi pe coloane. Particularizările în cazul unui vector sunt și ele evidente (matricea are o singură coloană).

Pe un tor sau o grilă, aceste moduri se pot generaliza, repartizarea făcându-se simultan pe linii și pe coloane. Vom discuta pe larg despre acest subiect în capitolul 6.

Pe hipercub, se folosesc uneori repartizări după codul Gray; notăm cu $g(i)$ codul Gray al numărului i . Se trece imediat de la cele trei repartizări pe linii de mai sus (bloc, ciclic și bloc ciclic), la alte trei, astfel: dacă o linie aparținea acolo procesorului P_i , acum ea va fi repartizată procesorului $P_{g(i)}$. Un avantaj al acestei tehnici este că, întotdeauna, linii vecine în matrice sunt repartizate aceluiași procesor sau unor

procesoare vecine (pe un ciclu hamiltonian); deci, atunci când se va adapta de la inel la hipercub un algoritm în care au loc doar comunicații între vecini, pentru a păstra această proprietate trebuie aleasă o repartizare conform codului Gray (aceeași, dintre cele trei, ca și pentru inel).

Indiferent de modul de repartizare, atunci când linia (coloana) i este repartizată în întregime unui singur procesor, atunci acest procesor va fi notat și cu $P(i)$.

Din nou despre indici globali și locali. Presupunem, de exemplu, o repartizare oarecare pe linii a unei matrice A ; aceasta înseamnă că un procesor deține un bloc $m \times n$. Într-un program, matricea A ca atare nu este referită printr-o singură variabilă (globală), ci prin p variabile locale fiecărui procesor, purtând același nume, să zicem B (evident, de obicei se folosește chiar numele A). Atunci când vom descrie un algoritm, ne vom referim de regulă la valorile indicilor ca și cum A ar fi memorată global; deci, a_{ij} este un element pe linia i a matricei A , repartizată procesorului $P(i)$; în variabila locală B a acestui procesor, acest element se numește $b_{\ell j}$, unde $\ell \in 0 : m - 1$ este indicele local al liniei respective; de pildă, pentru repartizare bloc linii, $\ell = i \bmod m$ (iar $P(i) = \lfloor i/m \rfloor$). Nu vom efectua decât rareori această transformare din indici globali în indici locali, pentru a nu îngreuna citirea algoritmilor. La implementare, însă, ea este obligatorie; altfel, memoria locală nu e utilizată eficient, ci aproximativ ca și cum întreaga matrice A ar fi memorată de fiecare procesor, dar un procesor nu ar utiliza efectiv decât cele m linii "ale sale".

Planificare și comunicație. Tot în arhitecturile cu memorie distribuită, ce efect are planificarea asupra comunicației? Dacă datele de ieșire ale unui task T_1 sunt date de intrare pentru un task T_2 , spunem că există o operație de comunicație logică între aceste taskuri; se poate construi un graf orientat al comunicațiilor logice, în noduri fiind taskurile, iar arcele legând taskuri care comunică. Dacă T_1 este planificat pe procesorul P_i , iar T_2 pe procesorul P_j , atunci, la execuția algoritmului, va avea loc o comunicație între aceste două procesoare; comunicația la acest nivel va fi numită fizică.

Are loc, așadar, o operație de scufundare a grafului legăturilor logice dintre taskuri, în grafurile legăturilor fizice dintre procesoare. O bună planificare trebuie să țină seama și de acest aspect. Noțiunile definite în primul capitol capătă acum semnificații noi. Dilatarea va reprezenta distanța maximă între procesoare care comunică; ea este foarte importantă, deoarece multe comunicații la distanțe mari vor mări mult timpul de execuție al algoritmului. Idealul este de a reuși o planificare în care să se comunice doar între vecini; de multe ori, în concepția algoritmilor, vom porni de la un astfel de criteriu. Congestia arată numărul legăturilor logice care se suprapun pe o aceeași legătură fizică; o congestie mare nu este neapărat nesatisfăcătoare, deoarece comunicațiile pot avea loc la momente diferite, deci, în realitate, nu vor avea loc conflicte de utilizare a legăturii fizice; astfel, minimizarea congestiei este un obiectiv secundar, în general.

Concluzii. Am enumerat până acum câteva criterii de planificare: planificarea cât mai multor taskuri în paralel, pentru a obține o încărcare egală, o distribuție

echitabilă și regulată a datelor, reducerea distanțelor între procesoarele care comunică. Este practic imposibil ca toate aceste criterii să fie utilizate într-o manieră explicită; adică, să se enumere toate (sau multe) variantele posibile de planificare, să se aprecieze performanțele lor și apoi să se aleagă varianta optimă. De cele mai multe ori sunt alese doar 2-3 planificări, pe baza unor particularități ale problemei (și a spiritului de observație al celui care o rezolvă), dintre care se alege varianta cea mai bună. După dobândirea unei oarecare experiențe, metoda este într-adevăr eficientă; sperăm că lectura acestei lucrări vă va ajuta în dezvoltarea unui bun spirit de observație specific.

4.3 Câțiva algoritmi paraleli fundamentali

Vom prezenta în această secțiune câteva exemple de studiu al paralelismului unor algoritmi rezolvând probleme simple, dar care apar foarte des; acestea ilustrează în același timp limitele ideilor de paralelizare descrise până acum, dar și procedee intuitive de obținere a paralelismului. Tratarea va fi aici mai completă decât în capitolele următoare, pentru că suntem la început; în plus, acești algoritmi vor fi folosiți efectiv mai târziu.

4.3.1 Suma a n numere

Prima problemă de care ne ocupăm este cea a calculului sumei a n numere x_0, x_1, \dots, x_{n-1} , elemente ale unui vector x . Ea ilustrează în primul rând importanța modului în care se formează taskurile.

Notăm $y_i^j = \sum_{l=i}^j x_l$. Pentru a putea determina paralelismul intrinsec al algoritmilor, vom presupune taskurile formate dintr-o singură operație. Vom utiliza relația (de recurență):

$$y_i^k = y_i^j + y_{j+1}^k, \text{ pentru } i \leq j < k$$

și cazul ei particular $y_i^{k+1} = y_i^k + x_{k+1}$ (care spune că suma a q numere se poate calcula adunând suma primelor $q-1$ cu cel de-al q -lea număr).

Algoritm secvențial. Într-o primă variantă notăm cu T_i , $0 < i \leq n-1$, operația de calcul al sumei parțiale $y_0^i = y_0^{i-1} + x_i$; calculul, în ordine crescătoare a indicilor, al sumelor parțiale, este algoritmul secvențial uzual de calcul al sumei. Se observă însă cu ușurință că, în acest fel, calculele se desfășoară pur secvențial; într-adevăr, cum $T_{i-1} < T_i$, graful de precedentă este o listă—deci fiecare nivel este format dintr-un singur task. Complexitatea paralelă este de $(n-1)$ operații, identică cu cea secvențială.

Paralelismul maxim. O a doua variantă se poate defini recursiv: se efectuează numărul maxim de adunări care pot decurge în paralel, e.g. se adună $x_i + x_{i+1}$, pentru toți i pari; se obțin astfel $n/2$ numere, pentru care se repetă algoritmul; la fiecare pas lungimea șirului de numere se înjumătățește, deci complexitatea este $\log n$

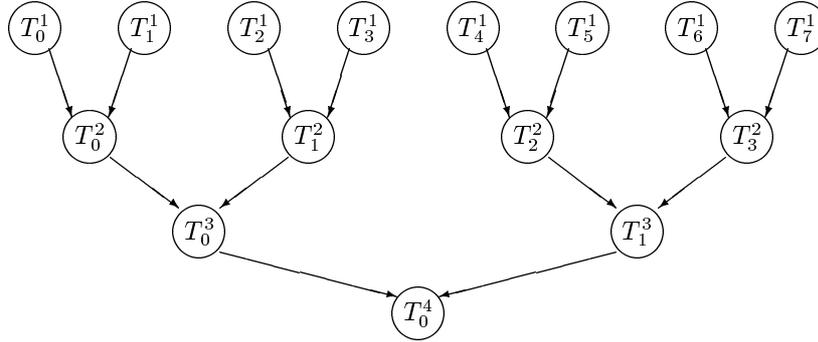


Figura 4.5: Graf de precedență pentru calculul sumei a 16 numere.

(presupunem, pentru simplitate, că n este o putere a lui 2). Altfel spus, pentru $k = 1, \dots, \log n$, se calculează în taskul T_i^k suma parțială

$$y_{i2^k}^{(i+1)2^k-1} = y_{i2^k}^{i2^k+2^{k-1}-1} + y_{i2^k+2^{k-1}}^{(i+1)2^k-1}, \text{ pentru } 0 \leq i \leq \frac{n}{2^k} - 1.$$

Relațiile de precedență sunt de tipul $T_{2i}^{k-1} \prec T_i^k$ și $T_{2i+1}^{k-1} \prec T_i^k$, ceea ce conduce la un graf de precedență de forma unui arbore binar (cu frunzele în sus!) echilibrat, deci având o înălțime de $\log n$; un astfel de arbore este prezentat în figura 4.5. Lățimea acestui graf este $n/2$. Așadar, putem spune că algoritmul are complexitate $O(\log n)$, folosind $O(n)$ procesoare; de altfel, aceasta este și complexitatea problemei calculului sumei. Costul algoritmului este $O(n \log n)$, dar efectiv nu se execută decât $n/2 + n/4 + \dots + 1 = n - 1$ adunări, exact cât și pentru algoritmul secvențial precedent. Din punct de vedere teoretic, algoritmul nu este optimal, dar este eficient.

Se observă diferența enormă între cele două variante (care, de altfel, sunt extreme), provenind exclusiv din împărțirea în taskuri a operațiilor, adică din ordonări diferite ale operațiilor; secvențial, algoritmi au același număr de operații.

Algoritmi pentru PRAM

Cazul $p = n/2$. Folosind $n/2$ procesoare, planificarea taskurilor din grafurile de precedență este simplă: în etapa k , procesorul P_i execută taskul T_i^k . Pe un EREW PRAM, algoritmul este următorul.

ALGORITHM 4.1 (suma numerelor x_i , pe EREW PRAM, $p = n/2$, $\text{id} \equiv i$)

1. **pentru** $k = 1 : \log n$
 1. **dacă** $i < n/2^k$ **atunci**
 1. execută T_i^k ($x_{i2^k} \leftarrow x_{i2^k} + x_{i2^k+2^{k-1}}$ sau $x_i \leftarrow x_{2i} + x_{2i+1}$)

a				b			
P_0	P_1	P_2	P_3	P_0	P_1	P_2	P_3
T_0^1	T_1^1	T_2^1	T_3^1	T_0^1	T_2^1	T_4^1	T_6^1
T_4^1	T_5^1	T_6^1	T_7^1	T_1^1	T_3^1	T_5^1	T_7^1
T_0^2	T_1^2	T_2^2	T_3^2	T_0^2	T_1^2	T_2^2	T_3^2
T_0^3	T_1^3			T_0^3	T_1^3		
T_0^4				T_0^4			

Tabelul 4.2: Două planificări, respectând principiul lui Brent, pentru graficul de precedentă din figura 4.5.

În instrucțiunea 1.1.1 sunt date două implementări posibile ale taskului T_i^k , în care elementele inițiale sunt modificate, astfel încât suma este conținută, în final, de x_0 . Cele două variante diferă prin locația în care se depune rezultatul unei sume. În prima, suma se depune în prima dintre cele două variabile sumate; astfel, după prima iterație, sumele parțiale sunt depuse în locațiile pare, după a doua în locațiile cu număr de ordine divizibil cu 4 etc. În a doua variantă, sumele parțiale se depun în primele locații ale vectorului x : prima jumătate, după prima iterație, primul sfert, după a doua iterație etc. Prima variantă este recomandată, deoarece, într-o iterație, fiecare locație de memorie este accesată doar de un singur procesor; în a doua, de exemplu, în același pas, x_1 este citit de P_0 și apoi scris de P_1 . Așa cum am anunțat în secțiunea 2.4, presupunem că procesoarele se sincronizează după fiecare iterație; mai mult, în a doua variantă presupunem că între citirea și scrierea din același pas, procesoarele se sincronizează. Accelerarea și eficiența au valorile

$$S(n) = \frac{n-1}{\log n}; \quad \varepsilon(n) = \frac{2(n-1)}{n \log n} \approx \frac{2}{\log n}.$$

Deci, deși timpul de execuție al algoritmului e foarte scurt, procesoarele nu sunt folosite eficient; în general, $\varepsilon(n) \ll 1$; aceasta se datorează timpilor morți ai majorității procesoarelor; doar P_0 lucrează permanent, în timp ce jumătate dintre procesoare nu fac decât o singură operație.

Cazul $p \ll n$. De obicei, lungimea vectorului x este semnificativ mai mare decât numărul de procesoare, adică $p \ll n$; cum planificăm acum taskurile? De exemplu, pentru $p = 4$, dacă aplicăm principiul lui Brent pentru graficul de precedentă din figura 4.5 (unde $n = 16$), atunci planificarea arată ca în tabelul 4.2a. Desigur, aceasta nu este unica variantă posibilă; același timp de execuție îl are și planificarea din tabelul 4.2b, care respectă și ea principiul lui Brent; deși aceasta din urmă nu pare a avea o calitate deosebită, totuși e mai naturală: fiecare procesor calculează suma a 4 numere consecutive, după care se folosește ideea de planificare din cazul $p = n/2$. Deci, în formă generală, algoritmul este următorul.

ALGORITM 4.2 (suma numerelor x_i , pe EREW PRAM, $p \ll n$, $\text{id} \equiv i$)

1. calculează (secvențial) $z_i \leftarrow \sum_{j=in/p}^{(i+1)n/p-1} x_j$
2. apelează algoritmul 4.1, pentru numerele $z_l, l = 0 : p-1$, cu $p/2$ procesoare

Am presupus că p divide n și că fiecare procesor calculează suma a n/p numere consecutive (și deci, o repartizare bloc linii a elementelor din vectorul x , pe procesoare); rezultatul final se va afla în z_0 . După calculul secvențial al sumelor, se obțin p sume parțiale; pentru aplicarea algoritmului 4.1 sunt suficiente $p/2$ procesoare, cele cu adrese $i = 0 : p/2 - 1$. Complexitatea algoritmului 4.2 este $T(n, p) = n/p - 1 + \log p$, iar eficiența

$$\varepsilon(n, p) = \frac{n-1}{p(n/p-1+\log p)} \approx \frac{n}{n+p \log p}$$

este asimptotic egală cu 1 și, în general, apropiată de 1. În acest caz, putem deci considera algoritmul valoros.

Algoritmi pentru MIMD cu memorie distribuită

Trecem acum la arhitecturile MIMD cu memorie distribuită, încercând să vedem care sunt implicațiile ideilor algoritmilor de mai sus asupra comunicațiilor. Presupunem că elementele din x sunt inițial distribuite egal celor p procesoare. Vom aborda întâi cazul în care $n = p = 2^d$, deci procesorul P_i are în memoria locală un singur element, să zicem x_i (acesta este cel mai mic număr de elemente de care merită să ne ocupăm; dacă rămâneam la o analogie perfectă cu cazul memorie comună, am fi ales $n = 2p$; cazul $n = p$ e însă cel mai defavorabil, dar care implică toate procesoarele). Păstrăm planificarea în care P_i execută taskul T_i^k ; dacă încercăm să utilizăm una dintre cele două forme ale instrucțiunii 1.1.1 din algoritmul 4.1, vom constata că P_i nu are în memoria proprie nici unul dintre elementele de sumat (prima formă e totuși adaptabilă, vezi problema 4.3.5). Vom încerca o a treia formă, în care P_i să calculeze de fiecare dată suma dintre elementul pe care îl posedă și un altul (pentru simplificarea scrierii am inversat și ordinea de ciclare):

1. **pentru** $k = \log n - 1 : -1 : 0$
 1. **dacă** $i < n/2^{d-k}$ **atunci** $x_i \leftarrow x_i + x_{i+2^k}$

În descriere recursivă, această formă seamănă cu cea din algoritmul 4.1, doar că se formează altfel perechile de elemente a căror sumă se calculează; se fac întâi sumele $x_0 + x_{n/2}, x_1 + x_{n/2+1}, \text{etc.}$; se obțin $n/2$ numere, pentru care se repetă procedeul. Atunci se observă imediat că procesorul P_i efectuează suma dintre numărul său (în care se acumulează o sumă parțială) și numărul (suma parțială) primit de la procesorul P_{i+2^k} ; condiția din instrucțiunea 1.1 spune că bitul k al numărului binar i este 0, adică P_i și P_{i+2^k} sunt vecini dintr-un hiper cub. În scrierea algoritmului în stil SPMD, vom ține seama de faptul că elementul x_i ocupă o singură locație în memoria locală a procesorului P_i , pe care o vom numi generic x (folosim așadar indici locali, în timp ce mai sus erau folosiți indici globali). Deci, algoritmul are forma:

ALGORITM 4.3 (suma a n numere, pe hipercub, $p = n$; numărul din memoria locală este x ; suma se calculează pe loc în x și se va afla în memoria procesorului P_0)

1. **pentru** $k = d - 1 : -1 : 0$
 1. **dacă** $\text{id} < 2^k$ **atunci**
 1. **recv**(z, k), $x \leftarrow x + z$
 2. **altfel dacă** $\text{id} < 2^{k+1}$ **atunci** **send**(x, k)

Se observă că după iterația k rămân active doar procesoarele dintr-un hipercub de dimensiune k , cel notat $\mathcal{H}_d^{d-k}(0)$; fiecare astfel de procesor P_i calculează în această etapă suma elementelor inițial locale procesoarelor din hipercubul de dimensiune $d - k$ notat $\bar{\mathcal{H}}_d^k(i)$. În final, suma este calculată de P_0 .

Timpul de execuție al algoritmului 4.3 este

$$T(p) = (\log p)\alpha + (\log p)(\sigma + \beta),$$

primul termen fiind pentru adunări, celălalt pentru comunicații. Desigur că eficiența algoritmului este modestă.

Dacă $n \gg p$, atunci fiecare procesor are în memoria locală n/p elemente, le calculează suma, după care se utilizează algoritmul 4.3. Timpul total este acum

$$T(n, p) = (n/p - 1)\alpha + (\log p)(\alpha + \sigma + \beta),$$

iar eficiența $\varepsilon(n, p) \approx \frac{n\alpha}{n\alpha + p(\log p)(\alpha + \sigma + \beta)}$, asimptotic egală cu 1.

Generalizare: colectare cu sumare. Să încercăm să facem acum legătura cu un algoritm de comunicație. Cu puțină atenție, se vede că, strict din punct de vedere al ordinii și dimensiunilor pe care se comunică, algoritmul 4.3 este inversul algoritmului sugerat în figura 3.23, de distribuție după un arbore de acoperire binomial, la fiecare pas după o direcție. Deci, seamănă cu un algoritm de colectare.

Într-adevăr, o primă idee naivă de a calcula suma pe o arhitectură cu memorie distribuită ar putea fi simpla colectare a elementelor x_i într-un singur procesor, urmată de calculul secvențial al sumei, în $p - 1$ operații. Aceasta este probabil o variantă mai puțin eficientă; ea poate fi luată în considerare dacă există o funcție hardware de colectare, mult mai rapidă decât una realizată prin program. Dar algoritmul 4.3 îmbină ideea colectării cu cea a calculului sumelor parțiale pe parcurs, astfel încât se transmite un singur element la un moment dat; în acest fel se reduce și numărul de operații aritmetice, și timpul de comunicație. Ideea se poate aplica pentru orice topologie și algoritm de colectare. Acesta este motivul pentru care algoritmi de sumare pe arhitecturi cu memorie distribuită sunt denumiți, de multe ori, algoritmi de *colectare cu sumare*.

O ultimă remarcă: toți algoritmi din această secțiune sunt valabili nu numai pentru sumă, ci și pentru orice altă operație, cu singura condiție ca aceasta să fie asociativă și, doar pentru algoritmul 4.3, comutativă (vezi problema 4.3.5 pentru a elimina această condiție). Calculul $x_0 \bullet x_1 \bullet \dots \bullet x_{n-1}$, unde \bullet este o operație asociativă mai este numit și *reducere*.

4.3.2 Suma globală

Problema sumei globale e o mică variație pe tema sumei, pe arhitecturi cu memorie distribuită. Fiind date n numere x_0, x_1, \dots, x_{n-1} , elemente ale unui vector x , să se calculeze suma lor astfel încât *toate* procesoarele să o cunoască (să o posedă în memoria locală). Desigur, calculul sumei ca mai sus, urmat de o difuzare efectuată de P_0 ar rezolva problema. Să căutăm însă o metodă mai rapidă. Ne vom mărgini discuția la hipercub. Presupunem întâi $p = n$, deci procesorul P_i deține în memoria locală numărul x_i .

Dacă algoritmul de calcul al sumei implica o comunicație identică celei de la colectare, formularea problemei sumei globale ar sugera să ne orientăm spre difuzarea generală. Să mai privim o dată figura 3.21, în care este prezentat un algoritm de difuzare generală, la fiecare pas pe o dimensiune. Modificăm algoritmul de difuzare generală astfel: după ce două procesoare vecine într-o anumită dimensiune își transmit valorile lor locale, ele efectuează și suma acestor valori, pe loc. Așadar, numerele asociate nodurilor în figura 3.21 reprezintă indicii inițiali și finali ai sumelor parțiale deținute de nodurile respective. Pentru exemplificare, să ne referim la procesorul P_4 ; în prima etapă, schimbă elementul său (x_4) cu cel al lui P_5 , deci pe dimensiunea 0, și calculează suma lor y_4^5 ; apoi, schimbă această sumă cu cea a lui P_6 , pe dimensiunea 1, și calculează noua sumă, adică y_4^7 ; în fine, schimbă cu P_0 , pe dimensiunea 2, și adună, obținând suma totală y_0^7 . Dacă $p < n$ generalizarea este imediată iar algoritmul este următorul.

ALGORITM 4.4 (suma globală a n numere, pe hipercub)

1. $x \leftarrow$ suma elementelor locale din vector
2. **pentru** $k = 0 : d - 1$
 1. **în paralel** $\text{recv}(z, k)$, $\text{send}(x, k)$
 2. $x \leftarrow x + z$

Operațiile efectuate sunt aceleași pentru toate procesoarele. Timpul său de execuție este identic cu cel pentru sumă, ceea ce justifică efortul de elaborare a algoritmului de sumă globală. Este interesant că am obținut un timp bun, chiar dacă unele operații sunt dublate; de exemplu, în primul pas, P_4 și P_5 calculează amândouă y_4^5 ; mai mult, în ultimul pas, toate procesoarele efectuează aceeași adunare. Avantajul este că prin aceste operații redundante se elimină necesitatea comunicației (a difuzării finale a sumei); în plus, ele sunt făcute de procesoare care altfel ar fi stat; timpii morți din algoritmul 4.3 sunt eliminați, fără ca P_0 , singurul care lucra permanent acolo, să aibă acum ceva în plus de făcut.

Calculul mai multor sume globale. Dacă trebuie calculate mai multe sume, aceasta se poate face comunicându-se în paralel. După calculul sumelor locale (dacă $n > p$), pentru fiecare sumă globală se începe comunicația pe altă dimensiune; de exemplu, pentru prima pe dimensiunea 0, pentru a doua pe dimensiunea 1, etc.; la pasul al doilea, pentru suma j se comunică pe dimensiunea $(j + 1) \bmod d$; ș.a.m.d., la pasul k , cu $0 \leq k \leq d - 1$, pentru suma j se comunică pe dimensiunea $(j +$

k) mod d , comunicația desfășurându-se în paralel, dacă sunt cel mult d sume globale de calculat. Deci, pentru d sume, timpul de comunicație este același ca pentru una singură; desigur, timpii necesari adunărilor se cumulează. Sursa de inspirație a acestei idei este algoritmul rotativ de difuzare generală !

4.3.3 Suma prefixelor

Aceasta este o generalizare a problemei sumei, cerându-se acum calculul tuturor sumelor parțiale $s_i \equiv y_0^i = \sum_{l=0}^i x_l$, cu $i = 0, \dots, n-1$. Numele de suma prefixelor (prefix sum) provine din faptul că aceste sume parțiale sunt formate întotdeauna cu primele elemente din vectorul x . Secvențial, calculul este banal, recurența $s_{i+1} = s_i + x_{i+1}$ implicând $n-1$ operații; însă, după cum am văzut la problema sumei, acest mod de calcul este absolut secvențial.

Algoritmi pentru PRAM

Sumare în cascadă. Să încercăm întâi generalizarea algoritmului 4.1; acolo, după prima iterație, $n/2$ procesoare calculau suma a câte două elemente consecutive din x ; după a doua, $n/4$ procesoare calculau suma a câte patru elemente; și tot așa, din ce în ce mai puține procesoare calculând suma tot mai multor elemente. Apare imediat ideea de a încerca să lăsăm *toate* procesoarele să lucreze, după același principiu: la pasul k , fiecare să calculeze suma a 2^k elemente consecutive. Pentru aceasta, trebuie modificată puțin prima formă a instrucțiunii 1.1.1 din algoritmul 4.1; se observă că se efectua suma între elemente aflate la distanță 2^{k-1} în vectorul x ; vom schimba doar locul în care se va depune rezultatul, anume în locația elementului cu indice mai mare. În plus, vom face aceasta pentru toate elementele x_i , de fiecare element ocupându-se un procesor; mai trebuie să avem grijă ca sumele să cuprindă doar elemente valide ale vectorului x , adică să nu efectuăm nici o operație dacă poziția $i - 2^{k-1}$ este în afara vectorului (mai mică decât 0). Mersul calculelor este exemplificat în figura 4.6, fiecare procesor efectuând o sumă parțială, și se desfășoară conform algoritmului următor, numit și de sumare în cascadă, în care sumele sunt calculate pe loc în elementele x_i .

ALGORITM 4.5 (*sumare în cascadă*) (suma prefixelor, pe CREW PRAM, $p = n$, $\text{id} \equiv i$)

1. **pentru** $k = 0 : \log n - 1$
 1. **dacă** $i \geq 2^k$ **atunci**
 1. $x_i \leftarrow x_i + x_{i-2^k}$

În etapa k , locația i e citită de două procesoare, P_i și P_{i+2^k} ; pentru a adapta algoritmul pe EREW PRAM, citirile ar trebui să aibă loc în două etape (folosind n locații de memorie în plus, pentru copii ale sumelor parțiale). De asemenea, chiar și pe CREW PRAM, funcționarea corectă a instrucțiunii 1.1.1 presupune că procesoarele se sincronizează între citiri și scrieri. Complexitatea paralelă a algoritmului este $\log n$, dar costul său de $n \log n$ îl face neoptimal (dar eficient, totuși).

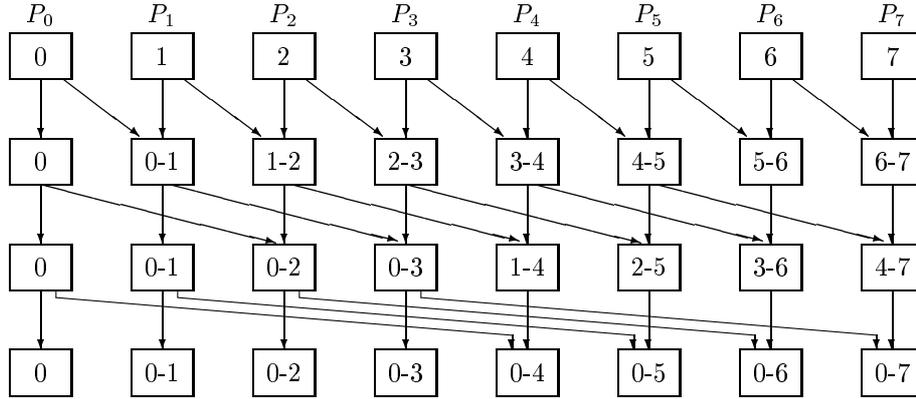


Figura 4.6: Calculul sumei prefixelor conform algoritmului 4.5 (sumare în cascadă); fiecare dreptunghi corespunde unei locații de memorie, conținutul reprezentând indicii limită ai sumelor parțiale (y_j^l este figurat ca $j-l$); etapele de calcul au loc de sus în jos; săgețile indică dependențele de date (se poate privi acest desen și ca o reprezentare a grafului de precedentă).

În cazul $n \gg p$, aplicarea principiului lui Brent algoritmului 4.5 nu duce la un rezultat prea fericit; costul algoritmului nu se modifică, deci timpul este $(n \log n)/p$. Ori, ne așteptăm să obținem ceva de genul n/p , pentru a avea eficiență aproape de 1. Pentru ameliorarea eficienței, putem proceda astfel: fiecare procesor P_i calculează suma elementelor "sale", cu indici de la in/p la $(i+1)n/p-1$, într-o variabilă z_i ; pentru aceste p variabile se aplică algoritmul 4.5; apoi, fiecare procesor calculează secvențial cele n/p sume ale sale, adunând pe rând elementele sale la suma $z_{i-1} = s_{in/p-1}$ (suma tuturor elementelor cu indice mai mic decât cel al elementelor sale).

ALGORITM 4.6 (suma prefixelor, pe CREW PRAM, $p \ll n$, $id \equiv i$)

1. $z_i \leftarrow \sum_{j=in/p}^{(i+1)n/p-1} x_j$ {calculul sumelor locale}
2. calculează suma prefixelor cu algoritmul 4.5, pentru elementele z_i
3. **dacă** $i \neq 0$ **atunci** $x_{in/p} \leftarrow x_{in/p} + z_{i-1}$
4. **pentru** $j = in/p + 1 : (i+1)n/p - 1$
 1. $x_j \leftarrow x_j + x_{j-1}$

Numărul de adunări este de n/p în instrucțiunea 1, $\log p$ pentru suma prefixelor, și n/p pentru instrucțiunea 4; deci, în total, sunt $2n/p + \log p$ adunări. Iată deci un prim caz în care eficiența asimptotică nu mai este 1, ci $1/2$; aceasta se datorează efectuării de două ori a adunărilor cu x_i , prima dată doar pentru a putea reduce dimensiunea problemei de la n la p ; este un procedeu destul de des întâlnit, acela de a face calcule suplimentare pentru ajungerea la o situație cu bun paralelism. În orice caz, rezultatul este mai bun decât cel obținut cu principiul lui Brent.

Algoritmul Ladner-Fischer. Vom prezenta acum un alt algoritm pentru care dependențele de date sunt mai simple; algoritmul se va implementa lesne pe un EREW PRAM. Ideea este de natură divide et impera; se reduce dimensiunea problemei sumei prefixelor la jumătate, se apelează recursiv algoritmul, apoi, din cele $n/2$ sume astfel calculate se deduc cele n sume cerute inițial. O descriere generală are forma următoare:

ALGORITM 4.7 (*Ladner-Fischer*) (suma prefixelor în paralel, n putere a lui 2)

1. **pentru** $i = 0 : n/2 - 1$, **în paralel** $\xi_i \leftarrow x_{2i} + x_{2i+1}$
2. apelează recursiv acest algoritm, pt. suma prefixelor valorilor ξ_i ($z_i \leftarrow \sum_{\ell=0}^i \xi_\ell$)
3. **pentru** $i = 0 : n - 1$, **în paralel**
 1. **dacă** i e impar **atunci** $s_i \leftarrow z_{(i-1)/2}$
 2. **dacă** i e par **atunci** $s_i \leftarrow z_{i/2-1} + x_i$

În instrucțiunea 1 se efectuează sumele a $n/2$ perechi de valori; după calculul recursiv al sumei prefixelor pentru aceste sume ξ_i , jumătate din sumele dorite sunt deja calculate (cele pentru i impar); pentru celelalte sume trebuie efectuată încă o adunare, cea din recurența $s_k = s_{k-1} + x_k$ (cu k par); convenim că $s_{-1} = z_{-1} = 0$. Se poate remarca o oarecare asemănare cu algoritmul 4.6, diferența esențială fiind că acela nu era recursiv. Complexitatea paralelă a algoritmului 4.7 este de $T(n) = 2 \log n$, câte $\log n$ înaintea, respectiv după apelurile recursive (se respectă relația $T(n) = T(n/2) + 2$). Pentru a prezenta o planificare a operațiilor, vom trece la o formă iterativă a algoritmului. Vom nota cu w_i^k variabila ce conține suma parțială i , la nivelul de recursie k ($w_i^0 \equiv x_i$).

ALGORITM 4.8 (suma prefixelor în paralel, pe EREW PRAM, n putere a lui 2, $p = n/2$, $\text{id} \equiv i$)

1. **pentru** $k = 1 : \log n$
 1. **dacă** $i < n/2^k$ **atunci** $w_i^k \leftarrow w_{2i}^{k-1} + w_{2i+1}^{k-1}$
2. **pentru** $k = \log n : -1 : 1$
 1. **dacă** $i < n/2^k$ **atunci**
 1. $w_{2i+1}^{k-1} \leftarrow w_i^k$, $w_{2i}^{k-1} \leftarrow w_{i-1}^k + w_{2i}^{k-1}$

După prima buclă **pentru**, $w_0^{\log n}$ conține suma totală; în a doua, în iterația k , procesorul P_i execută echivalentul instrucțiunilor 3.1 și 3.2 din algoritmul 4.7; deoarece acum privim din punctul de vedere al adresei procesorului, iar acesta se ocupă de elementele cu indicii $2i$ și $2i + 1$, se înlocuiește în 3.1 i cu $2i + 1$, și în 3.2 i cu $2i$. În plus, convenim că $w_{-1}^k = 0$. Într-o iterație, o locație de memorie nu este citită decât de un singur procesor; sunt necesare $n/2 + n/4 + \dots + 1 = n - 1$ locații suplimentare de memorie. Costul algoritmului este $O(n \log n)$.

În cazul $n \gg p$, se aplică aceeași idee ca la algoritmul 4.6, după sumarea secvențială apelându-se însă algoritmul 4.8; complexitatea paralelă este $2n/p + 2 \log p$. Eficiența asimptotică este deci $1/2$, pentru orice $p \leq O(n/\log n)$ (vezi problema 4.3.2, demonstrația fiind asemănătoare celei a algoritmului 4.2), iar costul algoritmului $O(n)$.

Algoritmi pentru MIMD cu memorie distribuită

Ne ocupăm doar de cazul $n = p$, pe hipercub; dacă $n \gg p$, se aplică ideea din algoritmul 4.6.

Sumare în cascadă. O primă încercare este de a vedea ce comunicații implică algoritmul 4.5; dacă elementul x_i este repartizat procesorului P_i , atunci, la pasul k , procesorul i comunică cu $i - 2^k$; din păcate, aceste două procesoare sunt vecine doar atunci când bitul k din i este 1; dacă însă $i_k = 0$, atunci toți biții din pozițiile mai mari decât k (inclusiv) sunt diferiți, deci distanța este $d - k$; concluzia este că algoritmul e inaplicabil în această formă.

Dacă se utilizează o repartizare după codul Gray, deci dacă x_i este în memoria locală a procesorului $g(i)$, atunci lucrurile se schimbă. Se poate demonstra (vezi problema 4.3.10) că distanța pe hipercub între procesoarele $g(i)$ și $g(i - 2^k)$ este cel mult 2; aceasta conduce la comunicații destul de performante. Nu detaliem algoritmul; se pot alege două implementări. Una este cea obișnuită, în care elementele rămân pe locurile lor, și se utilizează comunicația (cu eventualele funcții de dirijare oferite de sistemul de operare). A doua se bazează pe interschimbări ale elementelor, care se fac astfel încât să se comunice doar între vecini; un exemplu este prezentat în figura 4.7; în hipercubul \mathcal{H}_3 , au loc trei etape de comunicație, fiecare în doi pași (mai puțin ultima, după care nu mai e nevoie de interschimbări); fiecare desen reprezintă un pas de comunicație; limitele sumelor parțiale rezultate sunt cele din desenul următor; în primul pas se comunică la distanță 1, pe ciclul hamiltonian; urmează două interschimbări, reprezentând al doilea pas, care fac ca la al treilea pas să se comunice tot între vecini, etc. Nu este întâmplător faptul că, în final, repartizarea este cea în ordinea uzuală (și nu după codul Gray).

În ambele implementări, volumul comunicației poate fi apreciat la $(2 \log p)(\sigma + \beta)$, ceea ce este destul de convenabil; pentru adunări se consumă timpul $(\log p)\alpha$. Prima implementare este comodă, dacă există funcții accesibile de dirijare; a doua implică descoperirea și detalierea mecanismului general de interschimbare (la care vă urăm succes!).

Suma prefixelor și suma globală. Întrucât varianta pe hipercub a algoritmului 4.7 (Ladner-Fischer) nu este foarte simplă, vom prezenta o altă posibilitate, aceea de a adapta algoritmul de sumă globală astfel încât să fie calculate și sumele parțiale de care avem nevoie. Împărțim hipercubul în două jumătăți: inferioară (adresele mai mici decât $p/2$) și superioară; calculăm *atât* suma globală z , *cât și* suma prefixelor s , în fiecare dintre aceste hipercuburi, printr-un apel recursiv (z și s sunt variabile locale); apoi fiecare nod din hipercubul inferior schimbă suma sa globală cu cea a nodului echivalent din hipercubul superior; fiecare nod actualizează suma sa globală ca în algoritmul 4.4, adică adună suma z locală cu cea primită; pentru nodurile din hipercubul inferior, suma prefixelor este deja corectă; cele din hipercubul superior adună la suma prefixelor s locală, suma globală z primită. Justificarea este simplă: pentru $i \geq p/2$, avem $y_0^i = y_0^{p/2-1} + y_{p/2}^i$; dar $y_0^{p/2-1}$ este suma globală în hipercubul inferior, iar $y_{p/2}^i$ este suma prefixelor în hipercubul superior. În algoritmul următor, forma

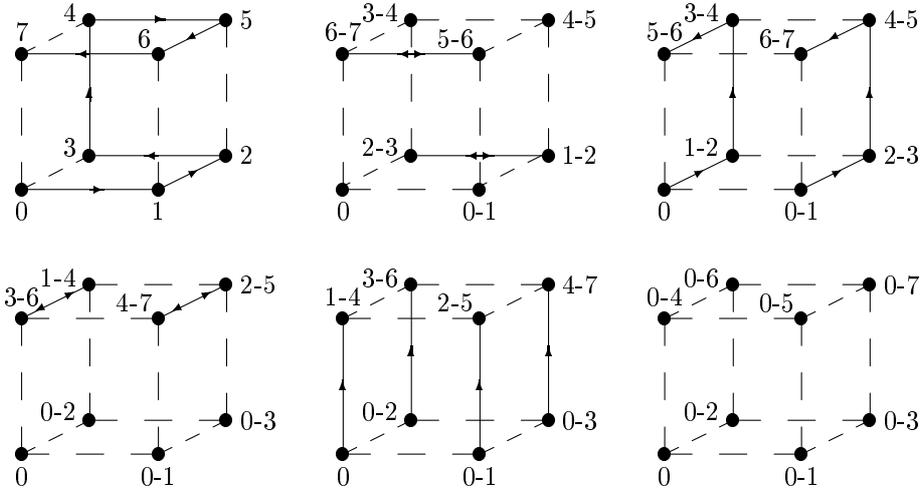


Figura 4.7: O implementare a algoritmului 4.5, pe hipercub, cu repartizare inițială după codul Gray și interschimbări la fiecare etapă.

recursivă capătă aspect iterativ datorită proprietăților topologice ale hipercubului.

ALGORITM 4.9 (suma prefixelor pe hipercub, adaptare a sumei globale, $p = n$, numărul din memoria locală este x)

1. $z \leftarrow x, s \leftarrow x$ { z – suma globală, s – suma prefixelor }
2. **pentru** $k = 0 : d - 1$
 1. **în paralel** $\text{recv}(v, k), \text{send}(z, k)$
 2. **dacă** $\text{id}_k = 1$ **atunci** $s \leftarrow s + v$
 3. $z \leftarrow z + v$

În 2.2 se actualizează suma prefixelor, în hipercubul superior, iar în 2.3, suma globală, întotdeauna. Timpul necesar calculelor este $(2 \log p)\alpha$ (P_{p-1} face la fiecare iterație câte două adunări), deci de două ori mai mare decât pentru adaptarea algoritmului 4.5 prezentată puțin mai sus. În schimb, timpul de comunicație este de două ori mai mic, adică $(\log p)(\sigma + \beta)$. Cum, în general, calculele sunt mai rapide decât comunicația, algoritmul 4.9 este de preferat. Totuși, dacă în locul adunării este o altă operație, trebuie avut în vedere că algoritmul 4.9 ca atare necesită și comutativitatea operației, spre deosebire de ceilalți algoritmi, pentru care asociativitatea este suficientă; o mică modificare face ca și acest algoritm să nu depindă de comutativitate (vezi problema 4.3.11).

4.3.4 Recurențe de ordinul I

Un șir $\{x_i\}$, recurent (liniar) de ordinul I, se definește prin relația $x_{i+1} = a_i x_i + b_i$, cu $0 \leq i \leq n-1$, unde a_i și b_i sunt coeficienți reali. Fiind dați acești coeficienți și valoarea x_0 , se pune problema determinării valorii termenului x_n . Calculul sumei a n numere este un caz particular al acestei probleme; dacă $s_i = \sum_{l=0}^i x_l$, atunci în recurența $s_{i+1} = a_i s_i + b_i$, coeficienții sunt $a_i = 1$ și $b_i = x_i$. Aplicarea formulei de recurență ca atare conduce la un algoritm perfect secvențial, cu complexitate de $2n$ operații.

Reducere ciclică. Ideea de paralelizare este clasică, și va mai fi întâlnită în această lucrare. În relația de calcul al termenului x_{i+2} , se înlocuiește x_{i+1} prin relația ce-l caracterizează, i.e.

$$x_{i+2} = a_{i+1}x_{i+1} + b_{i+1} = a_{i+1}a_i x_i + a_{i+1}b_i + b_{i+1} \equiv a'_i x_i + b'_i,$$

unde

$$a'_i = a_{i+1}a_i; \quad b'_i = a_{i+1}b_i + b_{i+1}. \quad (4.4)$$

Se obține astfel o nouă relație de recurență, de data aceasta între x_{i+2} și x_i , tot liniară și tot de ordinul I. Dacă se calculează noii coeficienți ai relației, ca în (4.4), pentru toți i pari, cu $0 \leq i < n$, se obține un nou șir recurent, conținând doar $n/2$ termeni. Se procedează analog în continuare, în fiecare etapă înjumătățindu-se lungimea șirului; în etapa k , lungimea șirului este $n/2^k$; dacă nu renumerotăm elementele șirului, după etapa k rămân în cursă elementele din pozițiile de forma $i = j2^k$, cu $0 \leq j < n/2^k$, ale șirului inițial; relația de recurență este acum de forma $x_{i+2^{k-1}} = a_i^k x_i + b_i^k$, iar noii coeficienți sunt

$$a_i^k = a_{i+2^{k-1}}^{k-1} a_i^{k-1}; \quad b_i^k = a_{i+2^{k-1}}^{k-1} b_i^{k-1} + b_{i+2^{k-1}}^{k-1}, \quad (4.5)$$

unde am notat în mod natural $a_i^0 = a_i$, $b_i^0 = b_i$. Presupunând că $n = 2^r$, după r etape se obține o unică relație de recurență: $x_n = a_0^r x_0 + b_0^r$, din care se calculează x_n .

Acest algoritm poartă numele de *reducere ciclică*, datorită reducerii permanente a numărului de elemente din șir, și, ca mecanism de desfășurare, este asemănător cu 4.1, cu modificarea din problema 4.3.5; totuși, trebuie făcută o observație importantă; deși graful de precedentă a taskurilor este, ca și pentru sumă, un arbore binar echilibrat cu sensurile tuturor arcelor schimbate, semnificația este cu totul alta; la sumă, o operație constă în efectuarea unei sume parțiale; la reducerea ciclică, însă, nu se calculează decât coeficienții unei noi relații, și nu ceva în legătură cu termenii x_i . Dacă sunt $n/2$ procesoare, calculele de tipul (4.5) se pot efectua în paralel, în fiecare etapă k ; pentru a planifica operațiile, ținem seama de faptul că un procesor P_j se ocupă de două perechi de coeficienți; deci, în etapa k , P_j va calcula a_i^k , cu $i = j2^k$ ($k \geq 1$). Forma algoritmului este următoarea:

ALGORITM 4.10 (reducere ciclică, pe EREW PRAM, $p = n/2$, $n = 2^r$)

1. **pentru** $k = 1 : \log n$
 1. **dacă** $\text{id} < n/2^k$ **atunci**
 1. $i \leftarrow \text{id} \cdot 2^k$
 2. $a_i^k \leftarrow a_{i+2^{k-1}}^{k-1} a_i^{k-1}$, $b_i^k \leftarrow a_{i+2^{k-1}}^{k-1} b_i^{k-1} + b_{i+2^{k-1}}^{k-1}$
2. **dacă** $\text{id} = 0$ **atunci**
 1. $x_n \leftarrow a_0^r x_0 + b_0^r$

Se observă că indicele superior al variabilelor a_i^k și b_i^k poate fi omis, calculele efectuându-se pe loc. Deoarece formulele (4.5) implică 3 operații, înseamnă că numărul total de operații va fi de $2 + 3 \log n$. Costul algoritmului este de $O(n \log n)$. Dar, ca și la calculul sumei, majoritatea procesoarelor nu lucrează decât o mică parte din timp; varianta secvențială a algoritmului 4.10 are numai $3(n/2 + n/4 + \dots + 1) + 2 = 3n - 1$ operații.

Cazul $n \gg p$. Aplicându-se principiul lui Brent, algoritmul obținut va avea o destul de bună eficiență. Dacă procesorul P_i se ocupă de calculul coeficienților a_l și b_l , cu $in/p \leq l \leq (i+1)n/p - 1$, atunci, prin $3n/p - 3$ operații, el reduce numărul de elemente din acest subsir la unul singur (aceasta se poate face și în altă ordine decât în algoritmul 4.10; de exemplu, P_0 poate calcula astfel: $a_0 \leftarrow a_1 a_0$, $b_0 \leftarrow a_1 b_0 + b_1$, apoi $a_0 \leftarrow a_2 a_0$, $b_0 \leftarrow a_2 b_0 + b_2$, ș.a.m.d., deci se elimină elementele din șir în ordine crescătoare); în acest fel, rămân în total p elemente din șir, pentru reducerea cărora se aplică algoritmul 4.10 cu $p/2$ procesoare, ceea ce va consuma $2 + 3 \log p$ operații; pentru o descriere mai riguroasă a algoritmului, vezi problema 4.3.13. Deci, în total, aproximativ $3(n/p + \log p)$ operații, ceea ce conduce la o eficiență

$$\varepsilon(n, p) \approx \frac{2n}{3(n + p \log p)} \xrightarrow{n \rightarrow \infty} \frac{2}{3}.$$

Prefixele unei recurențe. O altă problemă legată de șirurile recurente de ordinul I este calculul *tuturor* valorilor x_i , cu $0 \leq i \leq n$. Dacă analogul calculului elementului final x_n era calculul sumei, pentru această nouă problemă ne gândim, în mod natural, la o analogie cu suma prefixelor. Într-adevăr, toți algoritmi prezentați acolo pot fi ajustați cu ușurință; ne vom ocupa aici numai de adaptarea algoritmului 4.5, pentru ceilalți propunându-vă câteva probleme.

Ideea de bază este de a aplica algoritmul 4.10 din punctul de vedere al fiecărui procesor, după o schemă de tipul celei din figura 4.6. Se folosesc tot formulele (4.5), dar memorând acum noile valori ale coeficienților a și b în pozițiile cu indice superior și efectuând calculele pentru toți indicii i pentru care au sens; cu alte cuvinte, se folosesc formulele

$$a_i^k = a_i^{k-1} a_{i-2^{k-1}}^{k-1}; \quad b_i^k = a_i^{k-1} b_{i-2^{k-1}}^{k-1} + b_i^{k-1}, \quad (4.6)$$

pentru toți $i > 2^{k-1}$. Se obține așadar algoritmul:

ALGORITM 4.11 (reducere ciclică complet paralelă, pe CREW PRAM, $p = n$, $\text{id} \equiv i$)

1. pentru $k = 0 : \log n - 1$
 1. dacă $i \geq 2^k$ atunci
 1. $b_i \leftarrow a_i b_{i-2^k} + b_i$, $a_i \leftarrow a_i a_{i-2^k}$
 2. $x_{i+1} \leftarrow a_i x_0 + b_i$

După cum este normal, procesorul P_i calculează valoarea x_{i+1} (x_0 fiind cunoscut). Atribuirile din instrucțiunea 1.1.1 sunt scrise în această ordine pentru a se putea respecta relațiile (4.6). Trecerea la un algoritm EREW PRAM se face așa cum am discutat la suma prefixelor; totuși, aici apare o problemă în plus: accesul simultan al mai multor procesoare la x_0 ; chestiunea poate fi rezolvată foarte simplu, printr-un artificiu de notație; se renumerează x_i prin x_{i+1} , și se alege $x_0 = 0$; astfel, dificultatea dispare pur și simplu.

4.3.5 Produsul matrice-vector

Deși locul problemei calculului vectorului $z = Ax + y$, unde A este o matrice $n \times n$, y și z vectori de dimensiune n , ar fi mai degrabă în capitolul de algoritmi numerici, o vom prezenta totuși aici, în chip de avanpremieră; vom avea astfel o imagine a caracteristicilor algoritmilor implicând matrice și vectori. Operația $z = Ax + y$ este cunoscută și sub numele de gaxpy (General Ax plus y , matricea A fiind generală, i.e. neavând o structură specială); desigur, problema mai dificilă este de a calcula produsul matrice-vector, suma a doi vectori fiind o operație perfect paralelă.

Algoritmi pentru PRAM

Complexitatea paralelă. Notăm prin a_{ij} elementul din linia i , coloana j , al matricei A ; desigur, x_i este elementul din poziția (linia) i a vectorului x ; vom numera liniile și coloanele începând de la zero. Calculul vectorului z se face conform relației

$$z_i = y_i + \sum_{j=0}^{n-1} a_{ij} x_j, \text{ cu } 0 \leq i \leq n-1, \quad (4.7)$$

care implică $2n^2$ operații; pentru un i fixat, această relație este un produs scalar între linia i din A și vectorul x , ceea ce este sugerat în figura 4.8, în care termenul y_i nu a mai fost reprezentat. Se observă că, pe un CREW PRAM, fiecare z_i poate fi calculat independent; dacă alocăm câte $O(n)$ procesoare pentru calculul fiecărui z_i , acesta poate fi efectuat în $O(\log n)$ operații, fiind vorba despre o sumă; deci, cu $p = O(n^2)$ procesoare, complexitatea paralelă este $O(\log n)$. Desigur, deoarece costul este $O(n^2 \log n)$ și numărul de procesoare foarte mare, această abordare este utilă doar pentru a stabili care este complexitatea problemei.

Cazul $p = n$. Până la a ajunge la cazul practic $p \ll n$, să ne oprim la situația $p = n$, care permite o bună analiză a dificultăților ce apar. O primă idee e de a utiliza planificarea în care procesorul P_i calculează z_i ; din (4.7), se observă că P_i

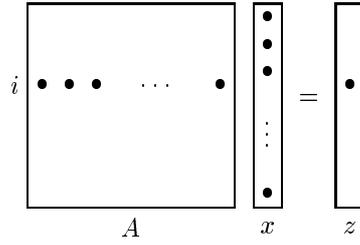


Figura 4.8: Produsul matrice-vector calculat prin produse scalare.

folosește linia i din A și vectorii x și y , în întregime. Pe un CREW PRAM, algoritmul decurge în felul următor.

ALGORITM 4.12 (gaxpy pe CREW PRAM, $p = n$, $\text{id} \equiv i$)

1. $z_i \leftarrow y_i$
2. **pentru** $j = 0 : n - 1$
 1. $z_i \leftarrow z_i + a_{ij}x_j$

Pentru a obține o versiune EREW PRAM, singurul impediment este accesul procesoarelor la vectorul x ; în ce privește matricea A , nu există nici un conflict la citire. Deoarece în instrucțiunea 2.1 ordinea de sumare este indiferentă, putem planifica operațiile astfel încât, la un moment dat, fiecare element din x să fie citit de un alt procesor; cea mai la îndemână modalitate este ca procesorul P_i să înceapă sumarea din 2.1 de la $j = i$ și să continue crescând j ; după $j = n - 1$, se continuă de la $j = 0$, deci ciclic.

ALGORITM 4.13 (gaxpy pe EREW PRAM, $p = n$, $\text{id} \equiv i$)

1. $z_i \leftarrow y_i$
2. **pentru** $j = i : n - 1$, $z_i \leftarrow z_i + a_{ij}x_j$
3. **pentru** $j = 0 : i - 1$, $z_i \leftarrow z_i + a_{ij}x_j$

Desigur că eficiența acestui algoritm depinde în mod critic de sincronismul executării instrucțiunilor.

Cazul $p \ll n$. Numărul de operații efectuate fără necesitatea sincronizării este acum mai mare. Dacă notăm cu \mathcal{L}_i mulțimea liniilor de care se ocupă procesorul P_i (de exemplu, de la in/p la $(i + 1)n/p - 1$), atunci calculele se desfășoară în felul următor:

ALGORITM 4.14 (gaxpy pe EREW PRAM, $p \ll n$, $\text{id} \equiv i$)

1. **pentru** $i \in \mathcal{L}_i$, $z_i \leftarrow y_i$
2. **pentru** $j = i : n - 1$
 1. **pentru** $i \in \mathcal{L}_i$, $z_i \leftarrow z_i + a_{ij}x_j$
3. **pentru** $j = 0 : i - 1$
 1. **pentru** $i \in \mathcal{L}_i$, $z_i \leftarrow z_i + a_{ij}x_j$

Este de remarcat ordinea buclor, care face ca un procesor să execute succesiv toate citirile unui element din x ; pentru a micșora și mai mult riscul conflictelor, se poate pune imediat înainte de instrucțiunile 2.1 și 3.1 un transfer de tipul $w_i \leftarrow x_j$, care mută elementul curent din x într-o locație de memorie "proprie" procesorului P_i .

Toți algoritmi de mai sus decurg în perfect paralelism; deci, dacă $p = n$, atunci timpul de execuție este $T(n) = 2n$; dacă $p \ll n$, atunci $T(n) = 2n^2/p$.

Algoritmi pentru MIMD cu memorie distribuită

Repartizare pe linii a matricei A . Ne ocupăm în special de cazul $p = n$, generalizarea fiind simplă. O primă idee e de a utiliza planificarea în care procesorul P_i calculează z_i ; la fel am procedat și în cazul memoriei comune. Implicația cea mai importantă este asupra repartizării datelor. Așa cum rezultată din relația (4.7) și din figura 4.8, procesorul P_i ar trebui să aibă în memoria locală linia i din matricea A , elementul y_i și vectorul x în întregime. Pe de altă parte, este naturală nu numai o repartitie echilibrată a datelor, cum e cea de mai sus, dar și una fără alocarea unei aceleiași date mai multor procesoare. Pentru a satisface și acest din urmă considerent, repartizarea inițială a datelor este una pe linii, pentru A , x și y ; de exemplu, P_i deține linia i din A , x_i și y_i . Rezultatul va fi și el distribuit în aceeași manieră, deci P_i va deține în final z_i .

Un algoritm imediat este următorul: fiecare procesor participă la o operație de comunicație globală în urma căreia va avea întregul x în memoria locală. Apoi, având toate datele necesare, aplică (4.7) pentru a calcula z_i . Deoarece fiecare P_i posedă x_i , comunicația globală constă în trimiterea acestui element tuturor celorlalte procesoare, deci este vorba de o difuzare generală. Pe scurt, algoritmul este următorul (peste tot în această secțiune folosim indici globali):

ALGORITM 4.15 (gaxpy pe MIMD cu memorie distribuită, $p = n$, $\text{id} \equiv i$, repartizare pe linii pentru A , x , y)

1. **difuzare generală:** P_i trimite x_i tuturor celorlalte procesoare și recepționează restul elementelor din x
2. $z_i \leftarrow y_i + \sum_{j=0}^{n-1} a_{ij} x_j$ (sau $z(i) \leftarrow y(i) + A(i, :) \cdot x$)

În instrucțiunea 2 am prezentat și o alternativă de scriere a relației (4.7), în stilul limbajului Matlab; $x(i)$ este același cu x_i , iar $A(i, :)$ reprezintă linia i din matricea A (indicele ':' are semnificația de "tot", în cazul nostru toate coloanele $j = 1 : n$).

Operațiile din algoritmul 4.15 decurg în perfect paralelism; apare însă un overhead datorat comunicației, a cărui mărime depinde de topologia arhitecturii distribuite (cu atât mai mic cu cât conectivitatea este mai bună). Timpul de execuție este de $2n\alpha$ pentru operații aritmetice, și, de exemplu pentru inel, de aproximativ $n/2(\sigma + \beta)$ pentru comunicație (vezi timpii de difuzare generală din secțiunea 3.3.3).

Acest algoritm are defectul că un procesor trebuie să memoreze întregul vector x ; în unele cazuri (pentru dimensiuni mari ale datelor), este de dorit evitarea acestei situații. Ne propunem deci să găsim un algoritm în care fiecare procesor să nu aibă la

un moment dat decât un singur element din vectorul x . Ideea este destul de simplă: pe parcursul operației de difuzare generală a vectorului x , fiecare procesor efectuează operațiile care implică elementele din x primite la un moment dat, după care trimite mai departe acele elemente către procesoare care nu le-au primit încă, însă *fără a le memora*; acest proces se repetă până la recepționarea tuturor elementelor din x .

Vom exemplifica aceasta pe inel, în cazul celui mai simplu algoritm cu putință, cel adecvat modelului half duplex. Deci, difuzarea generală decurge astfel: P_i posedă un element din x , pe care-l trimite vecinului din stânga; aceasta operație se repetă de $p - 1$ ori (până când fiecare element ajunge la vecinul din dreapta al posesorului său inițial, după ce a parcurs întreg inelul). Procesorul P_i , primind la un moment dat x_j , execută operația $z_i \leftarrow z_i + a_{ij}x_j$ (partea care implică x_j din 4.7) și trimite x_j mai departe, adică celui alt vecin. În formă compactă, aceasta se poate scrie astfel:

ALGORITHM 4.16 (gaxpy pe inel, $p = n$, $\text{id} \equiv i$, repartizare pe linii pt. A, x, y)

1. $z_i \leftarrow y_i, j \leftarrow i$
2. **pentru** $k = 0 : p - 1$
 1. $z_i \leftarrow z_i + a_{ij}x_j$
 2. **dacă** $k < p - 1$ **atunci, în paralel**
 1. **send**(x_j , stânga)
 2. **recv**($x_{(j+1) \bmod p}$, dreapta)
 3. $j \leftarrow (j + 1) \bmod p$

Un procesor nu trebuie să folosească decât două variabile locale pentru valorile x_j recepționate și retransmise (deoarece transmisia și recepția se fac simultan); memoria locală folosită pentru vectorul x este mult mai mică decât în algoritmul 4.16. Pentru o versiune a algoritmului 4.16 folosind numai indici locali, vezi problema 4.3.19.

Ultimii doi algoritmi se pot ușor generaliza pentru cazul în care $n \gg p$. Fiecare procesor deține n/p elemente din x și toate operațiile făcute mai sus cu un singur element x_i se vor face cu toate aceste n/p elemente. Vezi problema 4.3.20. De asemenea, generalizarea este banală și în cazul matricelor dreptunghiulare.

Repartizare pe coloane a matricei A . Ce se întâmplă dacă presupunem o repartizare pe coloane a matricei A ? (Evident, x și y , fiind vectori, își păstrează repartizarea.) Întrebarea nu este fără sens; produsul matrice-vector apare de obicei într-un context mai general, ca etapă în rezolvarea unei probleme și nu ca scop în sine; de aceea, repartizarea matricei A poate fi impusă de împrejurări.

Rămânând la cazul $n = p$, procesorul P_k va avea coloana k din A . Pentru o linie i oarecare, $0 \leq i \leq n - 1$, procesorul P_k poate contribui la calculul $z_i = y_i + \sum_{j=0}^{n-1} a_{ij}x_j$ din relația (4.7) doar cu efectuarea operației $s_{ik} = a_{ik}x_k$, restul elementelor din A nefiindu-i locale. Dacă fiecare procesor procedează la fel, calculul elementului z_i se reduce la efectuarea unei sume pe ansamblul procesoarelor. Modul de calcul al produsului matrice-vector este sugerat în figura 4.9 și este denumit prin *sumă vectorială*, deoarece, după efectuarea produselor locale, calculul lui z se reduce la o sumă de vectori; în figură este reprezentată doar partea ce revine procesorului P_k .

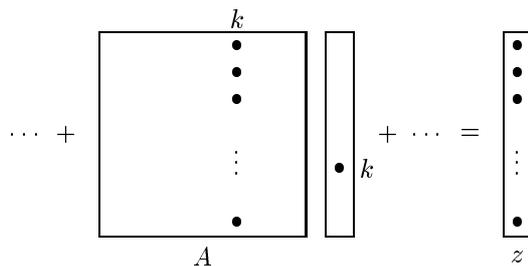


Figura 4.9: Produsul matrice-vector calculat prin sumă vectorială.

O metodă simplă de a realiza suma $z_i = y_i + \sum_{j=0}^{n-1} s_{ij}$ este colectarea tuturor valorilor s_{ij} în P_i (care trebuie să posede z_i , în final) și apoi efectuarea locală a sumei. Cum au loc p colectări simultane (pentru calculul tuturor elementelor din z), operația de comunicație globală este un schimb complet: procesorul P_i primește valoarea s_{ij} de la P_j , pentru orice $i, j \in 0 : p - 1$. Algoritmul este următorul.

ALGORITM 4.17 (gaxpy pe MIMD cu memorie distribuită, $p = n$, $\text{id} \equiv k$, repartizare pe coloane pentru A)

1. **pentru** $i = 0 : n - 1$
 1. $s_{ik} \leftarrow a_{ik} x_k$
2. **schimb complet**: P_k trimite s_{ik} lui P_i , cu $i = 0 : p - 1$
3. $z_k \leftarrow y_k + \sum_{j=0}^{n-1} s_{kj}$

Fiecare procesor efectuează n produse și n sume, ca în algoritmul 4.15. Volumul comunicației este însă mai mare, în general, deoarece aici este vorba de un schimb complet, fiecare procesor având n elemente, pe când acolo se făcea o difuzare generală, fiecare procesor având un singur element.

În plus, acest algoritm are același dezavantaj ca și 4.15: memoria locală ocupată este relativ mare; este necesară, pe lângă datele inițiale, și memorarea tuturor valorilor s , adică încă n pentru fiecare procesor.

Dar schimbul complet se poate face cu acumulare (procesorul P_i nu are nevoie de toate elementele s_{ij} , ci de suma lor); două valori, să zicem s_{iq} și s_{ir} , ambele cu destinația P_i (transmise de P_q și, respectiv, P_r), pot fi adunate pe traseu, într-un procesor în care se află amândouă la un moment dat, mai departe fiind transmisă doar suma lor. În fond, este același lucru cu a efectua p sume în paralel, pentru fiecare sumă destinația fiind un alt procesor. Aceasta reduce și memoria locală ocupată, și volumul comunicației, după cum se va vedea mai jos.

De exemplu, pe inel, o sumă se poate calcula foarte simplu transmitând datele într-un singur sens; fiecare procesor recepționează din dreapta o valoare (sumă parțială), o adună la cea proprie și trimite rezultatul mai departe spre stânga; altfel zis, suma circulă pe inel, îmbogățindu-se la fiecare pas cu o nouă valoare. În acest fel se pot

calcula p sume în paralel, fără nici un conflict de comunicație; la fiecare pas, pentru sume diferite sunt utilizate legături între perechi diferite de procesoare.

Algoritmul este echivalent cu a ”plimba” vectorul z pe inel, fiecare procesor deținând la un moment dat un element. La fiecare pas de comunicație, un element z_i se actualizează pe baza datelor oferite de procesorul la care se află; pentru a micșora memoria locală ocupată, procesorul P_k nu mai calculează de la început toate valorile s_{ik} (cu $0 \leq i < n$), ca în algoritmul 4.17, ci pe parcurs, pe măsură ce sunt necesare în actualizarea unui element z_i . Mai detaliat, forma algoritmului este:

ALGORITM 4.18 (gaxpy pe inel, $p = n$, $\text{id} \equiv k$, repartizare pe coloane pentru A)

1. $i \leftarrow (k + 1) \bmod p$, $z_i \leftarrow 0$
2. **pentru** $l = 0 : p - 1$
 1. $z_i \leftarrow z_i + a_{ik}x_k$
 2. **dacă** $l < p - 1$ **atunci, în paralel**
 1. **send**(z_i , stânga)
 2. **recv**($z_{(i+1) \bmod p}$, dreapta)
 3. $i \leftarrow (i + 1) \bmod p$
3. $z_k \leftarrow z_k + y_k$

Deoarece z_k parcurge inelul spre stânga, procesorul de la care pornește (pentru a ajunge în final la P_k) este vecinul din stânga al lui P_k , care are adresa $(k - 1) \bmod p$. Privind din punctul de vedere al procesoarelor, P_k deține inițial elementul din z care va ajunge în cele din urmă la vecinul său din dreapta, i.e. la procesorul $z_{(k+1) \bmod p}$, ceea ce explică instrucțiunea 1. Datorită acestui mod de comunicare, operația de adunare cu y din instrucțiunea 3 trebuie amânată la sfârșit, după ajungerea elementului z_k în P_k .

Algoritmii 4.18 și 4.16 au același timp de execuție atât pentru operații aritmetice, cât și pentru comunicație.

Pe alte topologii, ideile sunt asemănătoare, doar implementarea este mai complicată. Este de subliniat că, atunci când există rutine eficiente de comunicație globală și când memoria locală nu este o problemă, atunci algoritmii generali 4.15 (în special) și 4.17 sunt recomandați, datorită simplității lor.

Repartizare bidimensională a matricei A . Vom încheia acest paragraf, prezentând o variantă de înmulțire matrice-vector pe tor. Pentru simplitatea discuției, considerăm cazul în care fiecare procesor deține un element al matricei A , i.e. avem $p = n^2$ și procesorul P_{ij} deține elementul a_{ij} . Vectorii x , y și z sunt repartizați pe o singură dimensiune, de exemplu pe prima linie sau prima coloană de procesoare din tor. O repartizare echivalentă, la care se ajunge imediat din cele anterioare, dar care prezintă avantaje în cazul înmulțirii matrice-vector, este pe diagonala torului, i.e. procesorul P_{ii} deține elementele x_i , y_i și, în final, z_i .

Pentru a vedea structura comunicației necesară efectuării produsului Ax , să observăm că, în mod natural, procesoarele de pe linia i a torului cooperează pentru calculul elementului z_i . Fiecare procesor efectuează produsul $s_{ij} = a_{ij}x_j$. Pentru

aceasta, procesorul P_{ij} are nevoie de elementul x_j , aflat la P_{jj} ; în concluzie, este necesară o difuzare a lui x_j pe coloana j a torului, efectuată de procesorul diagonal de pe acea coloană. După calculul produsului s_{ij} , procesoarele de pe linia i a torului cooperează în calculul sumei $x_i = \sum_{j=0}^{n-1} s_{ij}$, cu destinația în procesorul P_{ii} . Așadar, algoritmul are structura următoare.

ALGORITM 4.19 (gaxpy pe tor, $p = n^2$, pentru procesorul P_{ij} , care deține a_{ij} ; procesorul diagonal P_{ii} deține x_i , y_i și, în final, z_i)

1. **difuzare** pe coloane: P_{ij} recepționează x_j de la P_{jj}
2. $s_{ij} \leftarrow a_{ij}x_j$
3. **sumă** pe linii: procesorul P_{ii} colectează $x_i = \sum_{j=0}^{n-1} s_{ij}$
4. **dacă** $i = j$ **atunci** $x_i \leftarrow x_i + y_i$

Se observă că operațiile sunt echilibrate între procesoare. Comunicația se desfășoară pe grupuri de procesoare, aflate pe aceeași linie sau coloană a torului.

Cazul $p \ll n$. Generalizarea algoritmilor anteriori pentru cazul $p \ll n$ este extrem de simplă. Vom exemplifica doar pentru repartizarea bidimensională. În loc să repartizăm individual fiecare element al matricei A , considerăm blocuri A_{ij} de dimensiune $r \times r$; fiecare procesor primește unul sau mai multe blocuri. De asemenea, vectorii x , y și z sunt împărțiți în blocuri de lungime r , notate X_i , Y_i , respectiv Z_i . Se poate scrie o formulă perfect analoagă cu (4.7), dar conținând nu operații elementare, ci cu blocuri (matrice sau vectori) de dimensiune r , i.e.

$$Z_i = Y_i + \sum_{j=0}^m A_{ij}X_j, \text{ cu } 0 \leq i \leq m,$$

unde $m = n/r$ este numărul de blocuri pe orizontală sau verticală. Orice algoritm scris la nivel de element se poate generaliza imediat la nivel de bloc, atât în privința calculelor cât și a comunicației. Capitolul 6 va prezenta numeroase exemple de algoritmi la nivel de bloc, care sunt preferați în practică datorită eficienței lor.

Probleme

P 4.3.1 Cum se poate calcula produsul scalar a doi vectori, x și y ?

P 4.3.2 Care este numărul de procesoare pentru care algoritmul 4.2 este optimal ? (Adică, se execută în timp $O(\log n)$ și are costul de $O(n)$.)

P 4.3.3 După iterația k a algoritmului 4.3, suma a căror elemente este locală procesorului P_i ?

P 4.3.4 Scrieți un program MPI care să implementeze algoritmul 4.3.

P 4.3.5 Adaptați algoritmul 4.1 astfel încât, folosind prima variantă a instrucțiunii 1.1.1, comunicația implicată de implementarea pe un hipercub să aibă loc numai între vecini.

P 4.3.6 Aplicați ideea colectării după o familie de arbori rotativi, la fiecare pas după o direcție, pentru calculul sumei a $n \gg p$ numere.

P 4.3.7 Sugați un algoritm de calcul al sumei pe inel.

P 4.3.8 Cum se poate calcula suma globală pe un inel ? Dar mai multe sume ?

P 4.3.9 Se poate elimina o operație din algoritmul 4.6 ?

P 4.3.10 Dacă $g(i)$ este codul Gray al numărului i , demonstrați că, pe un hipercub, distanța dintre nodurile cu adresele $g(i)$ și $g(i + 2^l)$ este cel mult 2.

P 4.3.11 Modificați algoritmul 4.9 a.î. să fie corect și pentru operații necomutative.

P 4.3.12 Se mai poate reduce timpul de execuție al algoritmului 4.10 ?

P 4.3.13 Scrieți versiunea algoritmului reducerii ciclice 4.10 pentru cazul $n \gg p$.

P 4.3.14 Fie șirul recurent definit prin relația $x_{i+1} = (a_i x_i + b_i)/(c_i x_i + d_i)$, în care a_i, b_i, c_i și d_i sunt coeficienți cunoscuți; de asemenea, se cunoaște valoarea x_0 . Indicați un algoritm paralel pentru calculul termenului x_n .

P 4.3.15 Propuneți un algoritm de rezolvare a unui sistem de ecuații liniare $Ax = b$, atunci când A este inferior bidiagonală ($a_{ij} = 0$ dacă $i < j$ sau $i > j + 1$), de dimensiune $n \times n$.

P 4.3.16 Adaptați algoritmul 4.8, pentru calculul tuturor termenilor unui șir recurent de ordinul I.

P 4.3.17 Scrieți un algoritm pentru calculul termenilor unui șir recurent pe un hipercub, inspirându-vă din algoritmul 4.9.

P 4.3.18 Să se adapteze algoritmul de difuzare generală 3.5 pentru realizarea comunicației necesare în calcularea produsului matrice-vector pe un inel.

P 4.3.19 Să se scrie o versiune a algoritmului 4.16 folosind doar indici locali.

P 4.3.20 Să se generalizeze algoritmul 4.16, de calcul al produsului matrice-vector pe un inel, la cazul $n \gg p$; considerați și cazul în care n nu se divide cu p .

P 4.3.21 Care dintre algoritmi 4.16 și 4.18 este mai eficient pentru înmulțirea matrice-vector pe un inel, atunci când matricea A este dreptunghiulară, de dimensiune $m \times n$; se presupune $m \neq n$ și $m, n \gg p$.

Capitolul 5

Sortare și probleme conexe

Fie $A = \langle a_0, a_1, \dots, a_{n-1} \rangle$ o secvență de n elemente ale unei mulțimi care admite o relație de ordine; a_i pot fi numere de diverse tipuri, șiruri de caractere etc.; a_i vor fi uneori numite și *chei*, deoarece ele pot face parte din înregistrări cuprinzând mai multă informație—numai cheile fiind folosite pentru ordonare. Pentru simplitate, vom presupune că cele n elemente sunt distincte¹; mai mult, de multe ori vom presupune că A este o permutare a secvenței $\langle 0, 1, \dots, n - 1 \rangle$ (atenție însă, această ipoteză este făcută doar pentru a ușura prezentarea; dacă se cunosc informații suplimentare asupra elementelor secvenței A , sortarea se poate simplifica mult; noi nu vom profita de aceasta). Dacă nu este specificat altfel, vom considera că n este o putere a lui 2. Algoritmii prezentați sunt valabili, uneori cu modificări minore (de "bucătărie", lăsate cititorului, de obicei), și în absența tuturor acestor simplificări.

Problemele de care ne vom ocupa în acest capitol sunt următoarele:

- problema sortării—aceea de a ordona crescător (sau descrescător) secvența A , deci de a găsi $A' = \langle a_{i_0}, a_{i_1}, \dots, a_{i_{n-1}} \rangle$ astfel încât $a_{i_j} < a_{i_k}$ (sau $a_{i_j} > a_{i_k}$) dacă $j < k$;
- problema găsirii minimumului, deci a găsirii indicelui i_m pentru care $a_i > a_{i_m}$, oricare $i \neq i_m$;
- problema selecției—a găsirii elementului din poziția j în secvența sortată; altfel zis, a elementului pentru care există fix $j - 1$ alte elemente mai mici decât el (am precizat aceasta pentru a nu rămâne cumva impresia că secvența trebuie sortată pentru aflarea celui de-al j -lea element);

¹Matematic vorbind, e ușor să diferențiem elementele egale dintr-o secvență; în loc de elementul a_i vom considera perechea (a_i, i) ; în orice secvență vor fi atunci perechi distincte. Algoritm vorbind, procedeul este neplăcut, implicând calcule și memorie suplimentare; e de preferat să ne descurcăm fără acest subterfugiu.

- problema gășirii rangului unui element, inversa precedentei; dându-se elementul $v \in A$, să se găsească poziția sa în secvența sortată (sau, dacă vreți, câte elemente sunt mai mici decât v);
- problema interclasării (fuziunii, engl. *merge*): dându-se două secvențe crescătoare A și B , de lungimi n , respectiv m , să se găsească secvența sortată $A \perp B$ care cuprinde cele $n + m$ elemente din A și B .

Evident, vom insista asupra tratării paralele a problemelor. O imensă sursă de algoritmi secvențiali este Knuth [17], apărută în 1976 și în traducere românească. Surse de algoritmi paraleli vor fi citate pe parcurs; o bună lucrare pentru algoritmi pe PRAM este [16].

Pentru aprecierea complexității algoritmilor vom contoriza numai comparațiile între chei (deci nu vom număra operațiile de interschimbare a două elemente, aritmetica indicilor etc.); în acest capitol α va însemna timpul necesar pentru o comparație (el nu va fi folosit în formule decât atunci când este imperios necesar).

Mulți dintre algoritmii prezentați în acest capitol au complexități care depind, uneori semnificativ, de ordinea elementelor din A ; de aceea expresiile pe care le vom prezenta pentru acestea au o nuanță probabilistă; va fi vorba deci despre timpi de execuție în medie, în cel mai favorabil caz sau în cel mai defavorabil caz, cu precizarea cât de aproape de valorile respective se poate spera să fie timpul de execuție real. Trebuie făcută totuși o ipoteză (probabilistă) asupra ordinii din A : vom presupune că indicii elementelor din A sunt repartizați uniform, adică un element se poate afla, cu aceeași probabilitate, pe orice poziție din A ; aceasta este, de altfel, ipoteza uzuală.

Notății. În scrierea algoritmilor vom folosi notația $exch(a_i, a_j)$ sau—cu aceeași semnificație— $exch(i, j)$ (când va fi evident care este secvența implicată) pentru operația de comparație a elementelor a_i și a_j , urmată de interschimbare dacă $a_i > a_j$; deci, în prima poziție va ajunge elementul mai mic. Pe un calculator cu memorie comună, această operație va fi efectuată de un procesor (deci secvențial). În cazul unui calculator cu memorie distribuită sunt două posibilități; dacă ambele elemente aparțin aceluiași procesor, operația este cea banală; dacă elementele aparțin unor procesoare diferite, atunci fiecare procesor comunică elementul său celuilalt, apoi fiecare face comparația și reține o valoare: primul minimul, al doilea maximul; se efectuează un pas de comunicație și o comparație. Să observăm că în cazul memorie distribuită interschimbarea se face prin comunicație (și se comunică indiferent de valorile elementelor), în timp ce pentru memoria comună prin trei atribuiri succesive (făcute însă numai dacă e cazul).

Dacă A_i și A_j sunt două secvențe, vom folosi următoarele notații: $\langle A_i, A_j \rangle$ este secvența obținută prin concatenarea, în ordinea precizată, a secvențelor A_i și A_j ; $A_i < A_j$ semnifică faptul că orice element din A_i este mai mic decât orice element din A_j ; ca un caz particular, $v < A_i$, unde v este o cheie, înseamnă că v este mai mic decât orice element din A_i .

5.1 Găsirea minimului

Problema aflării elementului minim dintr-o secvență constituie un caz particular de reducere și poate fi rezolvată prin algoritmi din secțiunea 4.3.1. Vom relua aici unii algoritmi și vom prezenta și alții noi, care nu se pot generaliza pentru reducere.

Algoritmul secvențial uzual este următorul: se compară primele două elemente; minimul se compară cu al treilea; noul minim cu al patrulea etc. Sunt necesare $n - 1$ comparații, această valoare fiind optimă (fiecare comparație elimină un candidat la minim, și trebuie eliminați $n - 1$). În această formă, nu este nici urmă de paralelism.

Algoritm pe EREW PRAM. Găsirea minimului se poate face utilizând un arbore binar echilibrat de decizie, ca în figura 4.5, care are în cele n frunze elementele din A și în care fiecare nod tată are asociată valoarea minimă a celor doi fiu; altfel zis, algoritmul urmează principiul de organizare a unui turneu de tenis. Așadar, se formează $n/2$ perechi de elemente și se extrage din fiecare minimul; apoi, din aceste minime se formează $n/4$ perechi, ș.a.m.d. La pasul i , găsirea minimului din fiecare dintre cele $n/2^i$ perechi se poate face într-o singură operație de către $n/2^i$ procesoare. Algoritmul ce urmează este asemănător cu algoritmul 4.1 de calcul al sumei.

ALGORITM 5.1 (găsire minim pe EREW PRAM, $p = n/2$, $\text{id} \equiv k$)

1. **pentru** $i = 1 : \log n$
 1. **dacă** $k < n/2^i$ **atunci**
 1. $\text{exch}(k2^i, k2^i + 2^{i-1})$

În final minimul se găsește în a_0 . Deci, $T(n) = T(n/2) + 1$, adică $T(n) = \log n$, și se folosesc $O(n)$ procesoare. Acest timp este optim; de altfel și în modelul CREW PRAM, găsirea minimului necesită $\Omega(\log n)$ comparații.

Algoritm pe COMMON CRCW PRAM. În acest model, lucrurile nu mai sunt atât de evidente, chiar dacă simple. Se poate imagina un algoritm care să găsească minimul cu o singură comparație !

Să observăm întâi că putem compara simultan toate perechile posibile (a_i, a_j) , cu $i \neq j$; pentru aceasta este nevoie de $\binom{n}{2} = n(n-1)/2$ procesoare. Ce facem însă cu rezultatele comparațiilor ?

Fie un vector c de dimensiune n , inițializat cu valoarea 1 peste tot. Procesorul care compară a_i cu a_j scrie 0 în poziția din c corespunzătoare elementului mai mare dintre cele două. Nu există nici un conflict de scriere: dacă mai multe procesoare scriu în aceeași locație, ele scriu toate valoarea 0. În urma acestei operații, un singur element din c va păstra valoarea inițială 1, anume cel din poziția minimului; toate celelalte sunt superioare cel puțin unui alt element, deci cel puțin un procesor va scrie 0 în poziția corespunzătoare din c .

Presupunând deci că sunt $p = n(n-1)/2$ procesoare, să scriem algoritmul pentru procesorul notat P_{ij} ($i > j$):

ALGORITM 5.2 (găsire minim pe COMMON CRCW PRAM)

1. $c_i \leftarrow 1; c_j \leftarrow 1$
2. **dacă** $a_i > a_j$ **atunci** $c_i \leftarrow 0$
3. **altfel** $c_j \leftarrow 0$

Deci, o singură comparație e suficientă; e drept că numărul de procesoare e foarte mare.

Algoritmul de mai sus poate fi generalizat. La fiecare pas căutarea minimumului va fi limitată la o mulțime S , care inițial coincide cu A , iar în final conține un singur element—minimumul.

Fie x cel mai mic întreg astfel încât atunci când S este partiționată în x blocuri de dimensiune $\lfloor m/x \rfloor$ sau $\lceil m/x \rceil$ (unde $m = |S|$), sunt suficiente p comparații pentru a compara elementele fiecărei perechi dintr-un același bloc. Astfel, în x pași (folosind algoritmul 5.2 cu p procesoare, în fiecare bloc), este determinat minimumul din fiecare bloc și S se reduce la x elemente.

Pentru a compara într-un pas toate elementele dintr-un bloc sunt necesare $\binom{m/x}{2}$ procesoare, deci $p = x \frac{m/x(m/x-1)}{2}$ și, în fine, $x = \frac{m^2}{m+2p}$.

Obținem deci relația de recurență $T(m) = T(m^2/(m+2p)) + x$. Se poate arăta că, dacă $n/2 \leq p \leq \binom{n}{2}$, atunci numărul de operații pentru a reduce cardinalitatea lui S de la n la 1 este $T(n, p) = \log \log n - \log \log(2p/n) + O(1)$. Mai mult decât atât, acesta este și optimul pentru găsirea minimumului.

Algoritm pentru MIMD cu memorie distribuită. Găsirea minimumului este asemănătoare colectării. Întâi, fiecare nod calculează minimumul între cele $m = n/p$ elemente locale (în $m-1$ operații). Apoi, pe arborele de acoperire pe care se realizează colectarea, fiecare procesor așteaptă minimele care vin dinspre frunze, le compară cu minimumul local și trimite mai departe elementul minim dintre acestea.

De exemplu, într-un hipercub, algoritmul de găsire a minimumului de către P_0 folosind un arbore de acoperire binomial, este următorul (asemănător cu algoritmul 4.3):

ALGORITM 5.3 (găsire minim în hipercub, $p = 2^d$, $n \gg p$, $\text{id} \equiv k$)

1. găsește μ , minimumul elementelor locale
2. **pentru** $i = d-1 : -1 : 0$
3. **dacă** $k_i = 1$ **atunci**
 1. **send**(μ, i), **stop**
4. **altfel** **recv**(μ', i)
5. **dacă** $\mu' < \mu$ **atunci** $\mu \leftarrow \mu'$

Să observăm că procesorul P_k calculează elementul minim în subhipercubul $\tilde{\mathcal{H}}_d^{d-r}(k)$, unde r este numărul de biți consecutivi de 0 de la începutul reprezentării binare a lui k pe d biți.

Timpul său de execuție este $T(n) = \left(\frac{n-1}{p} + \log p\right)\alpha + \log p(\sigma + \beta)$. Se constată că $T(p) = O(\log p)$, deci algoritmul este optimal și în ce privește numărul de comparații și în ce privește comunicația.

Dacă se dorește ca fiecare procesor să cunoască minimumul, procesorul care l-a calculat va face o difuzare, sau se aplică ideea de la calculul sumei globale, i.e. algoritmul 4.4.

Probleme

P 5.1.1 Care este principalul motiv pentru care algoritmul secvențial curent de găsim a minimului este cel amintit în text și nu unul asemănător cu cel paralel pentru EREW PRAM (tip turneu de tenis) ?

P 5.1.2 Să se scrie algoritmul de găsim a minimului în modelul EREW PRAM atunci când sunt doar $p \ll n$ procesoare. Care este timpul de execuție ? Dar eficiența algoritmului ?

P 5.1.3 Scrieți algoritmul 5.1 astfel încât, după pasul i , minimul să fie între primele $n/2^i$ elemente ale secvenței A .

P 5.1.4 În algoritmul 5.2, în linia 1, sunt necesare două operații pentru inițializarea vectorului c (e drept că ne putem permite, câtă vreme nu numărăm decât comparațiile !); se poate cu una singură ?

P 5.1.5 a) Dacă în secvența A pot fi elemente egale, algoritmul 5.2 funcționează corect ?
b) Dar dacă în linia 2 se înlocuiește $a_i > a_j$ cu $a_i \geq a_j$?

P 5.1.6 Demonstrați că pentru $T(n, p) = \log \log n - \log \log(2p/n) + O(1)$ este adevărat $T(n, n(n-1)/2) = O(1)$.

P 5.1.7 Scrieți algoritmul de găsim a minimului pe un inel.

P 5.1.8 Cu care dintre cei doi algoritmi de găsim a minimului pe un EREW PRAM—5.1 sau cel din problema **P5.1.3**—seamănă algoritmul 5.3 ? Scrieți un algoritm care să semene cu celălalt.

5.2 Selecția și găsim a rangului

Aceste probleme apar de obicei în legătură cu sortarea și mult mai rar de sine stătătoare; de aceea vom considera mai ales cazurile în care despre ordinea din A există informații suplimentare.

5.2.1 Rangul într-o secvență sortată

Ne vom ocupa întâi de găsim a rangului unei valori v într-o secvență A sortată crescător (v nu este neapărat un element din A). Secvențial, problema este simplă și se rezolvă printr-o căutare binară în $\log n$ comparații (se compară v cu elementul aflat la mijlocul secvenței; se află astfel în care jumătate de secvență se găsește v ; se continuă în acea jumătate).

Rangul pe CREW PRAM. Pe un astfel de model arhitectural, rangul poate fi găsit în timp constant, cu n procesoare. Se folosește o variabilă logică suplimentară c , de lungime n . Procesorul P_k compară elementul a_k cu v și scrie în c_k valoarea 0 dacă a_k este mai mic și 1 dacă este mai mare. Cum secvența A este ordonată, în acest moment în c se vor afla r valori consecutive de 0, urmate de $n - r$ valori de 1. Procesorul P_k citește acum valorile c_k și c_{k-1} ; un singur procesor, anume P_r , va găsi două valori diferite (facem convenția că $c_{-1} = 0$). Deci algoritmul este

ALGORITM 5.4 (rangul lui v în secvența sortată A pe CREW PRAM, $p = n$, $\text{id} \equiv k$; în final r conține rangul)

1. **dacă** $a_k < v$ **atunci** $c_k \leftarrow 0$
2. **altfel** $c_k \leftarrow 1$
3. **dacă** $c_k \neq c_{k-1}$ **atunci** $r \leftarrow k$

Rangul pe MIMD cu memorie distribuită. În acest context, modificăm puțin problema, astfel: procesorul P_k deține o secvență ordonată crescător A_k de lungime $m = n/p$, dar despre ordinea în $A = \langle A_0, \dots, A_{p-1} \rangle$ nu se mai poate spune nimic altceva. Presupunem că valoarea v este locală procesorului P_0 , ceea ce nu este o particularizare.

Ideea algoritmului este simplă; procesorul P_0 difuzează valoarea v tuturor procesoarelor; fiecare procesor calculează rangul local r_k , deci rangul lui v în secvența locală A_k ; apoi se execută un algoritm de sumă globală, după care fiecare procesor va deține rangul lui v . Algoritmul se bazează pe faptul evident că rangul (global) este suma rangurilor locale $r = \sum_{k=0}^{p-1} r_k$.

ALGORITM 5.5 (calculează r , rangul lui v în secvența A sortată local, pe MIMD cu memorie distribuită, $mp = n$, $\text{id} \equiv k$)

1. **difuzare:** P_0 trimite tuturor valoarea v
2. $r_k \leftarrow$ rangul lui v în A_k (prin căutare binară)
3. $r \leftarrow$ suma globală a valorilor r_k (cu algoritmul 4.4)

Algoritmul determină corect rangul chiar dacă A conține chei egale (se calculează numărul de elemente strict mai mici decât v).

5.2.2 Selecție pe PRAM

Trecem acum la problema selecției, găsirea elementului din poziția j în secvența A sortată, pe care o vom rezolva în cazul cel mai general: nici o informație inițială despre ordinea din A .

Algoritmul secvențial utilizat de obicei este inspirat din cel de sortare rapidă (quicksort). Ideea este de separa A în două părți, de a vedea în care din ele se găsește elementul căutat, apoi de a continua căutarea numai în acea parte; o strategie divide et impera, așadar. Întâi se alege un element $v \in A$, numit pivot în quicksort, dar pe care l-am putea numi candidat. Apoi se împarte secvența A în două: în prima parte (stânga) vor fi elementele mai mici decât v , în a doua (dreapta), elementele mai mari sau egale cu v ; ca o mică observație, numărul de elemente din stânga este chiar rangul lui v în A ; dacă în stânga sunt exact j elemente, v este elementul căutat; dacă sunt mai puțin de j elemente, căutarea va fi continuată în dreapta (candidatul la selecție a fost prea mic), altfel se va căuta în stânga. Forma generală a algoritmului este (valabilă și la implementarea paralelă):

ALGORITM 5.6 (algoritm general de selecție a elementului din poziția j)

1. $s \leftarrow 0, d \leftarrow n - 1$
2. **cât timp** $s \leq d$ {secvența este nevidă}
 1. $alege_pivot(s, d, j)$
 2. $i_v = partiționare(s, d)$ {formează cele două subsecvențe}
 3. **dacă** $i_v = j$ **atunci stop**
 4. **altfel dacă** $i_v < j$ **atunci** $s \leftarrow i_v + 1$ {continuă în dreapta}
 5. **altfel** $d \leftarrow i_v - 1$ {continuă în stânga}

În acest algoritm, s și d reprezintă indicii primului element, respectiv cel al ultimului, în secvența în care se caută. Candidatul v este calculat de procedura $alege_pivot$ și adus în poziția s , printr-o eventuală interschimbare, deci $a_s = v$; în general se procedează astfel: se alege un eșantion de t elemente din secvența în care se caută, se ordonează acest eșantion și se alege de obicei mediana²; există multe euristici de alegere a pivotului; mediana e un candidat bun pentru că se speră împărțirea secvenței de căutare în două părți de dimensiuni apropiate, înjumătățind la fiecare pas dimensiunea problemei; cu cât t este mai mare, cu atât mediana eșantionului aproximează mai bine mediana secvenței; timpul de calcul al medianei crește, dar scade mai repede dimensiunea problemei; practic se alege $t = 3$ sau chiar $t = 1$. Procedura $partiționare$ calculează poziția i_v a pivotului, permutând elementele în A astfel încât $a_i < v$, când $i < i_v$ și $a_i \geq v$, când $i > i_v$; algoritmul secvențial de partiționare merită atenția: vezi problema 5.2.4. În cazul înjumătățirii dimensiunii, numărul de operații respectă relația $T(n) = T(n/2) + n$, deci $T(n) = 2n$. În cel mai rău caz se poate ajunge la $O(n^2)$ operații, dar aceasta este foarte improbabil; se poate conta, în medie, pe $O(n)$.

Selecție pe CREW PRAM, $p = n$. Putem să aflăm într-un singur pas ce elemente sunt mai mari sau mai mici decât pivotul; în schimb, este dificil să punem unele din aceste elemente într-o zonă contiguă de memorie, pentru a continua căutarea. Prezentăm întâi algoritmul și apoi explicațiile.

ALGORITM 5.7 (selecție a elementului din poziția j în secvența A , pe CREW PRAM, $p = n$, $id \equiv k$)

1. $n' \leftarrow n, r \leftarrow 0$
2. **cât timp** $n' \neq 1$
 1. se alege v și se pune în a_0 (deci $v \equiv a_0$)
 2. **dacă** $a_k < v$ **atunci** $c_k \leftarrow 1$ ($P_0 : c_0 \leftarrow 0$), $e_k \leftarrow 1, f_k \leftarrow 0$
 3. **altfel** $c_k \leftarrow 0, e_k \leftarrow 0, f_k \leftarrow 1$ ($P_0 : f_k \leftarrow 0$)
 4. calculează $c_k \leftarrow \sum_{i=0}^k c_i$, cu algoritmul 4.8 (suma prefixelor)
 5. calculează $f_k \leftarrow \sum_{i=0}^k f_i$, cu algoritmul 4.8 (suma prefixelor)

²Sau elementul din poziția $\lfloor t(j-s)/(d-s+1) \rfloor$, cel care se află într-o poziție proporțională, în eșantion, cu elementul din poziția j în A , în secvența de căutare, presupunând că ar fi sortată; adică $i/(j-s) \approx t/(d-s+1)$; primul este raportul pozițiilor, al doilea cel a lungimilor.

6. **dacă** $r + c_{n'-1} = j$ **atunci stop**
7. **dacă** $r + c_{n'-1} > j$ **atunci**
 1. **dacă** $e_k = 1$ **atunci** $b_{c_k-1} \leftarrow a_k$
($P_0 : n' \leftarrow c_{n'-1}$)
8. **altfel** $\{r + c_{n'-1} < j\}$
 1. **dacă** $e_k = 0$ **atunci** $b_{f_k-1} \leftarrow a_k$
($P_0 : n' \leftarrow f_{n'-1}, r \leftarrow r + c_{n'-1}$)
9. **dacă** $k \geq n'$ **atunci stop**
10. $a_k \leftarrow b_k$

Pentru a nu lungi prea mult descrierea, am apelat la următoarea convenție: instrucțiunile care încep cu " $P_0 :$ " sunt executate de P_0 în locul celor imediat precedente. n' este numărul de elemente rămase în zona de căutare; r este numărul de elemente mai mici decât candidatul curent v , care nu mai sunt în zona de căutare (au fost eliminate la iterații anterioare); acestea sunt variabile globale, citite de toată lumea și modificate doar de P_0 . Instrucțiunea 2.1, alegerea pivotului, nu este explicitată; important e să dureze $O(1)$; pentru simplitate ne putem gândi că v e chiar a_0 . În 2.2, c_k , copia sa e_k și negatul său f_k exprimă sensul comparației între a_k și v . În 2.4 se aplică suma prefixelor pentru c_k ; elementele mai mici decât v au acum valori c_k diferite, și anume chiar pozițiile lor într-o subsecvență conținând doar valorile mai mici decât v (prima poziție fiind 1). Să exemplificăm pentru $n = 8$:

k	0	1	2	3	4	5	6	7
a_k	5	1	8	9	3	2	6	4
c_k, e_k	0	1	0	0	1	1	0	1
f_k	0	0	1	1	0	0	1	0
$\sum_{i=0}^k c_k$	0	1	1	1	2	3	3	4
$\sum_{i=0}^k f_k$	0	0	1	2	2	2	3	3

În mod analog, suma prefixelor pentru f_k aduce în aceste variabile (când $a_k > v$) poziția într-o subsecvență conținând valorile mai mari decât v . În exemplu, secvența elementelor mai mici decât v este $\langle 1, 3, 2, 4 \rangle$ (rangurile sunt în penultima linie), iar $\langle 8, 9, 6 \rangle$ cea a celor mai mari (rangurile în ultima linie). Să observăm că, după instrucțiunea (3), $c_{n'-1}$ este chiar rangul lui v în secvența de căutare³; rangul în A este obținut adunând la $c_{n'-1}$ numărul de elemente mai mici decât v deja eliminate, adică r . Dacă rangul este j algoritmul se termină. Dacă este mai mare, rămân în cursă elementele mai mici decât v , cele pentru care $e_k = 1$, în număr de $c_{n'-1}$; ele sunt copiate într-un zonă de memorie suplimentară B . În linia 2.8.1 se procedează analog în cazul în care rangul este mai mic decât j ; în plus, se actualizează r , pentru că sunt $c_{n'-1}$ elemente mai mici decât oricare dintre cele între care se va continua căutarea,

³În treacăt am rezolvat deci și problema rangului atunci când secvența A nu este sortată; timpul necesar este $O(\log n)$, cât pentru suma prefixelor; de fapt e suficientă o simplă sumă.

acestea din urmă fiind mutate începând de la poziția 0 în A . La fiecare iterație sunt necesare n' procesoare; în 2.9 își încheie activitatea procesoarele cu adresă $\geq n'$. În 2.10, se recopiază elementele în A . Sunt necesare câteva sincronizări: înainte de 2.6, 2.9 și 2.10.

Timpul de execuție este dictat de suma prefixelor; în cazul favorabil al înjumătățirii lui n' la fiecare iterație, avem $T(n) < T(n/2) + O(\log n)$, adică $T(n) = O(\log^2 n)^4$. Costul algoritmului este de $O(n \log^2 n)$. Numărul minim de operații pentru selecție (de fapt pentru aflarea medianei) s-a demonstrat a fi $\Omega(\log n / \log \log n)$, pe un CREW PRAM, cu un număr mărginit polinomial de procesoare.

Selecție pe CREW PRAM, $p \ll n$. Schițăm mai jos modificările aduse algoritmului 5.7, păstrând structura lui.

ALGORITM 5.8 (selecție pe CREW PRAM, $p \ll n$, $\text{id} \equiv k$)

1. $n' \leftarrow n$, $r \leftarrow 0$, $m \leftarrow n/p$, \tilde{A} se află între pozițiile km și $(k+1)m - 1$
2. **cât timp** $n' \neq 1$
 1. se alege v și se pune în a_0 (deci $v \equiv a_0$)
 2. aplică algoritmul din **P5.2.4** pt. partiționarea $\tilde{A} = \langle \tilde{A}_0, \tilde{A}_1 \rangle$, $\tilde{A}_0 < v \leq \tilde{A}_1$
 3. $c_k \leftarrow |\tilde{A}_0|$, numărul de elemente din \tilde{A} mai mici decât v
 4. $f_k \leftarrow |\tilde{A}_1|$, numărul de elemente din \tilde{A} mai mari decât v
 5. $c_{k+1} \leftarrow \sum_{i=0}^k c_i$ și $f_{k+1} \leftarrow \sum_{i=0}^k f_i$, cu algoritmul 4.8 (puțin modificat)
 6. **dacă** $r + c_p = j$ **atunci stop**
 7. **dacă** $r + c_p > j$ **atunci**
 1. copiază \tilde{A}_0 în B , în pozițiile de la c_k la dreapta ($c_0 = 0$)
 2. **dacă** $k = 0$ **atunci** $n' \leftarrow c_p$
 8. **altfel** $\{r + c_p < j\}$
 1. copiază \tilde{A}_1 în B , în pozițiile de la f_k la dreapta ($f_0 = 0$)
 2. **dacă** $k = 0$ **atunci** $n' \leftarrow f_p$, $r \leftarrow r + c_p$
9. $m \leftarrow \lceil n'/p \rceil$, \tilde{A} se află între indicii (din A) km și $\min((k+1)m - 1, n' - 1)$
10. copiază în \tilde{A} elementele corespunzătoare din B

Să exemplificăm valorile c_k și f_k dintr-o iterație, în tabelul 5.1, pentru $p = 4$. Se observă că c_4 este numărul de elemente mai mici decât pivotul 13, iar f_4 numărul de elemente mai mari; cum c_4 este rangul pivotului în A , merită să ne amintim că rangul global este egal cu suma rangurilor locale c_i , cu $0 \leq i \leq 3$. În ultimele două linii se găsesc pozițiile de scriere în B .

Într-o iterație, fiecare procesor se ocupă de aproximativ $m = n'/p$ elemente succesive din A (notăm secvența acestor elemente cu \tilde{A}). De menționat că, în 2.9, unele

⁴Timpul poate fi redus dacă mutăm 2.5 imediat înainte de 2.8.1; în acest caz, a doua sumă a prefixelor se va face doar dacă este nevoie. De altfel, a doua sumă a prefixelor nu este necesară; notând $d_k = \sum_{i=0}^k c_i$ și $g_k = \sum_{i=0}^k f_i$, observăm că avem întotdeauna doar una din situațiile (i) $d_k = d_{k-1} + 1$ și $g_k = g_{k-1}$, sau (ii) $d_k = d_{k-1}$ și $g_k = g_{k-1} + 1$. În consecință, g_k rezultă din d_k , d_{k-1} și g_{k-1} .

k	0	1	2	3	4
\tilde{A}	$\langle 13, 1, 23, 25 \rangle$	$\langle 18, 9, 7, 16 \rangle$	$\langle 11, 6, 10, 0 \rangle$	$\langle 15, 20, 17, 5 \rangle$	
\tilde{A}_0, \tilde{A}_1	$\langle 1 \rangle, \langle 23, 25 \rangle$	$\langle 7, 9 \rangle, \langle 18, 16 \rangle$	$\langle 11, 6, 10, 0 \rangle, \langle \rangle$	$\langle 5 \rangle, \langle 20, 17, 15 \rangle$	
c_k	1	2	4	1	
f_k	2	2	0	3	
$\sum_{i=0}^k c_k$	0	1	3	7	8
$\sum_{i=0}^k f_k$	0	2	4	4	7

Tabelul 5.1: Exemplu de execuție pentru algoritmul 5.8.

procesoare pot rămâne cu o secvență \tilde{A} vidă, spre finalul execuției algoritmului; altfel, numai ultimul procesor are ceva mai puține elemente decât celelalte; în orice caz, procesoarele sunt aproximativ la fel de încărcate. Numărul de elemente mai mari sau mai mici decât v deținute de un procesor se calculează secvențial în algoritmul din **P5.2.4**. Suma prefixelor (algoritmul 4.8, modificat pentru că se scriu rezultatele în elementul din poziția următoare) durează acum $O(\log p)$. Timpul mare de calcul este consumat în partiționarea secvențială; dacă pivotul este bine ales, avem $T(n) \leq T(n/2) + n/p + \log p$, deci $T(n) = 2n/p + \log p \log n$, ceea ce este foarte bine, al doilea termen (cel suplimentar unei perfecte distribuiri a operațiilor celor p procesoare) fiind nesemnificativ.

5.2.3 Selecție pe MIMD cu memorie distribuită

Modificăm puțin problema selecției, în sensul că presupunem secvențele locale ca fiind ordonate; dorim în continuare aflarea celei de-a j -a valori din secvența A ordonată, s-o numim mai departe x ; presupunem că elementele din A sunt distincte. Schițăm întâi ideea algoritmului, după care vom prezenta detalii.

Un procesor propune un candidat v pentru x , după care, cu algoritmul 5.5, se calculează rangul acestui candidat în secvența A . Fiecare procesor posedă două variabile care delimitează spațiul în care se caută x , ele conținând indicii s și d în secvența locală \tilde{A} , pentru care este sigur că $\tilde{a}_s \leq x \leq \tilde{a}_d$; valorile inițiale sunt $s = 0$, $d = m - 1$; o dată calculat rangul candidatului, fiecare procesor poate micșora spațiul de căutare. Apoi este propus un nou candidat, ș.a.m.d., până când se găsește x ; condiția de terminare este aceeași ca la PRAM: rangul candidatului să fie chiar j .

Avem de dat răspunsuri mai multor întrebări.

Cine propune candidații? Este clar că dacă un procesor are $s > d$, el nu deține x , deci nu are sens să propună un candidat. Așadar, toate procesoarele trebuie să participe la propuneri; unul singur deține x , dar nu se știe care. Pot fi luate în considerare mai multe variante. De exemplu, la iterația i , propune procesorul P_i ; există o valoare specială, să-i spunem NIL , care este propusă atunci când $s > d$, și care înseamnă "nu am candidat". Mecanismul pare destul de greoi; dacă rămân

doar puține procesoare pentru care $s < d$, ele trebuie să le aștepte pe celelalte, care difuzează valori *NIL*. Vom alege o altă metodă: fiecare procesor propune câte un candidat la fiecare iterație; procesorul P_0 colectează acești candidați, elimină valorile *NIL*, apoi selectează unul singur, v , de exemplu mediana (elementul din poziția din mijloc în secvența sortată a candidaților), pe care-l difuzează. Această metodă nu e evident superioară anterioarei; are totuși avantajul că la fiecare iterație se propune un candidat; o adoptăm din motive ce se vor vedea abia în secțiunea dedicată sortării rapide.

Ce valori propune un procesor drept candidați? Inițial, cea mai plauzibilă valoare este cea din poziția $\lfloor j/n \rfloor$ (reamintim că $n = mp$ este lungimea totală a secvenței. Apoi, în iterațiile următoare, fiecare procesor va propune valoarea mediană din spațiul său de căutare, cea cu indicele $\lfloor (s+d)/2 \rfloor$ în secvența locală.

Cum se reduce spațiul de căutare? În procesul de detectare a rangului candidatului, un procesor determină un indice l astfel încât $\tilde{a}_{l-1} < v \leq \tilde{a}_l$; l este chiar rangul local. După calculul rangului global r , sunt posibile trei cazuri. Fie $r = j$ și atunci soluția a fost găsită. Fie $r < j$ și atunci candidatul a fost prea mic; se elimină toate elementele mai mici decât candidatul, deci $s \leftarrow l$ (sau $l+1$ dacă $v = \tilde{a}_l$; mai concis, elementul cu cel mai mic indice, mai mare decât v). Fie $r > j$ și deci candidatul a fost prea mare; atunci $d \leftarrow l-1$.

Care este elementul din poziția j ? Ultimul candidat v (deci valoarea v la sfârșitul algoritmului).

Vom numi P_0 procesorul *conducător* al selecției, datorită rolului său special în acest algoritm. Este evident că oricare procesor poate avea rolul de conducător într-o arhitectură simetrică (inel, tor, hipercub; dar nu și grilă), fără a se modifica numărul de operații.

Să sintetizăm toate acestea într-un algoritm; r este rangul global, fiecare procesor deținând câte o copie a acestei valori.

ALGORITM 5.9 (determină elementul din poziția j în secvența A sortată local, pe MIMD cu memorie distribuită, $mp = n$, $\text{id} \equiv k$, procesor conducător P_0)

1. $s \leftarrow 0$, $d \leftarrow m-1$, $c \leftarrow \tilde{a}_{\lfloor j/p \rfloor}$, $r \leftarrow -1$ {Inițializare}
2. **cât timp** $r \neq j$
 1. **colectare**: P_0 primește candidații c de la celelalte procesoare
 2. P_0 elimină valorile *NIL*, ordonează restul candidaților și alege v mediana lor
 3. $r \leftarrow$ rangul global al lui v (se aplică algoritmul 5.5)
 4. **dacă** $r \neq j$ **atunci** {reduce spațiul local de căutare}
 1. **dacă** $r < j$ **atunci** {candidat prea mic}
 1. $s \leftarrow$ cel mai mic indice pentru care $\tilde{a}_s > v$
 2. **altfel** $d \leftarrow$ cel mai mare indice pentru care $\tilde{a}_d < v$
 2. {propune un nou candidat}
 3. **dacă** $s \leq d$ **atunci**

1. $c \leftarrow \tilde{a}_{\lfloor (s+d)/2 \rfloor}$ {mediana în spațiul de căutare}
4. **altfel** $c \leftarrow NIL$

În instrucțiunea 2.2, deși nu am explicat, este clar că procesorul P_k nu face nimic, când $k \neq 0$; în schimb, la 2.1 au loc comunicații globale, deci P_k participă.

Să apreciem calitativ timpul de execuție pentru o singură iterație, pe un hipercub. Instrucțiunile 1, 2.4.1.1 și 2.4.2 se execută în timp constant, instrucțiunea 2.1 în $O(p)$, 2.2 în $O(p \log p)$ (dar vom vedea că găsirea secvențială a mediane necesită doar $O(p)$) și 2.3 în $O(\log(n/p))$. Adunând, putem aproxima timpul necesar unei iterații la $O(\log n)$ pentru comparații și $O(p)$ pentru comunicație.

A determina exact numărul de iterații este imposibil. La fiecare iterație, cel puțin un procesor își înjumătățește spațiul de căutare; deci, în cel mai defavorabil caz, au loc $p \log(n/p)$ iterații; aceasta este însă o apreciere foarte pesimistă. În general, spațiul de căutare al fiecărui procesor este redus la fiecare iterație, deci se poate aprecia că, în medie, numărul de iterații este de $O(\log n)$. Să remarcăm că în cazul algoritmului 5.8, problema dezechilibrului de încărcare se rezolvă la fiecare iterație; aici ar fi dificil, comunicația implicată fiind foarte complexă.

Complexitatea algoritmului este așadar de $O(\log^2 n)$ pentru comparații și de $O(p \log n)$ pentru comunicație.

Probleme

P 5.2.1 Cu ce amendamente este valabil algoritmul 5.4 pe un EREW PRAM ?

P 5.2.2 Scrieți algoritmul de găsim a rangului unei valori v într-o secvență ordonată A de lungime n , în contextul MIMD cu memorie distribuită. Se presupune că procesorul P_i are în memoria locală secvența A_i de lungime $m = n/p$ și că $A_i < A_j$, dacă $i < j$.

P 5.2.3 Scrieți un algoritm pentru calculul rangului pe o arhitectură cu memorie distribuită, atunci când A este o secvență oarecare (nu este ordonată local).

P 5.2.4 Detaliați procedura *partiționare* din algoritmul 5.6.

P 5.2.5 Care este condiția de terminare a algoritmului 5.9 dacă în secvența A pot exista elemente cu aceeași valoare ? Ce modificări mai trebuie aduse algoritmului ?

P 5.2.6 Cum trebuie modificat algoritmul 5.9 dacă secvențele locale nu sunt ordonate ?

5.3 Interclasarea

Fie secvențele $\langle a_0, a_1, \dots, a_{n-1} \rangle$ și $\langle b_0, b_1, \dots, b_{m-1} \rangle$, cu $n \leq m$. Algoritmul secvențial de interclasare a lor decurge în felul următor; se folosesc doi indici, unul în A , unul în B , inițializați cu 0; se compară elementele având acești indici; cel mai mic dintre ele este adăugat secvenței $A \perp B$ (inițial vidă); indicele corespunzător secvenței în care s-a găsit elementul minim este incrementat; se continuă până când primul dintre indici depășește lungimea secvenței. Se fac cel puțin $\min(n, m)$ comparații și cel mult $n + m - 1$. (Deși răspunsul va fi dat pe parcursul acestei secțiuni, gândiți-vă de acum la un algoritm paralel de interclasare pentru două procesoare.)

5.3.1 Complexitate paralelă teoretică

Se pot interclasa două secvențe în timp constant, pe un CREW PRAM ? Da, dacă folosim algoritmul 5.4 de găsim a rangului într-o secvență ordonată. Dacă se calculează rangul r_i al fiecărui element a_i din A între elementele din B , atunci poziția lui a_i în $A \perp B$ este $i + r_i$ (sunt i elemente în A și r_i elemente în B mai mici decât a_i); numărul de procesoare necesar este nm , câte m pentru calculul rangului fiecărui element din A . Se procedează analog cu elementele din B .

ALGORITM 5.10 (interclasare pe CREW PRAM, $p = nm$, în final $D = A \perp B$)

1. **pentru** $i = 0 : n - 1$, **în paralel**
 1. câte m procesoare calculează r_i , rangul lui a_i în B , cu algoritmul 5.4
 2. un procesor din fiecare grup de m execută: $d_i \leftarrow a_{i+r_i}$
2. **pentru** $j = 0 : m - 1$, **în paralel**
 1. câte n procesoare calculează r_j , rangul lui b_j în A , cu algoritmul 5.4
 2. un procesor din fiecare grup de n execută: $d_j \leftarrow b_{j+r_j}$

Dacă notăm cu P_{ij} procesoarele, o alocare a taskurilor pentru acest algoritm este simplă. Costul algoritmului este foarte mare, $O(n^2)$ comparații când $m = n$, ceea ce lasă acestui algoritm doar importanță teoretică.

Algoritmul Valiant. Vom prezenta acum un algoritm de interclasare a două secvențe, folosind cel mult $(n + m)/2$ procesoare, pe un CREW PRAM. Sursa este Kruskal [18].

ALGORITM 5.11 (*Valiant*) (interclasare $A \perp B$ pe CREW PRAM, $p = (n + m)/2$)

1. Împarte A în \sqrt{n} secvențe A_i de lungime \sqrt{n} fiecare și B în \sqrt{m} secvențe B_j de lungime \sqrt{m} fiecare (ignorăm micile modificări necesare atunci când n și m nu sunt pătrate perfecte).
2. Fie a_i^* ultimul element al secvenței A_i și b_j^* ultimul element al secvenței B_j . **În paralel**, găsește rangul fiecărui a_i^* în secvența $B^* = \langle b_0^*, \dots, b_{\sqrt{m}-1}^* \rangle$, folosind algoritmul 5.4; altfel spus, pentru fiecare a_i^* , găsește unicul indice $j(i)$ pentru care $b_{j(i)-1}^* < a_i^* < b_{j(i)}^*$ (convenim că $b_{-1}^* = -\infty$ și $b_{\sqrt{m}}^* = \infty$).
3. **În paralel**, găsește rangul fiecărui a_i^* în secvența $B_{j(i)}$.
4. repetă recursiv pașii 1-3, pentru interclasarea (în paralel) a fiecărei secvențe A_i cu o secvență din B .

Să detaliem puțin. Scopul pașilor 1-3 este de găsi poziția exactă a elementelor a_i^* în secvența B ; în figura 5.1 secvențele A și B sunt simbolizate prin cele două linii orizontale; liniile oblice începând în a_i^* au capătul celălalt în b_i^a , care este cel mai mic element din B mai mare decât a_i^* ; în figură, am considerat exemplul $a_0^* < b_0^* < a_1^* < b_1^* < b_2^* < a_2^* < a_3^* < b_3^*$.

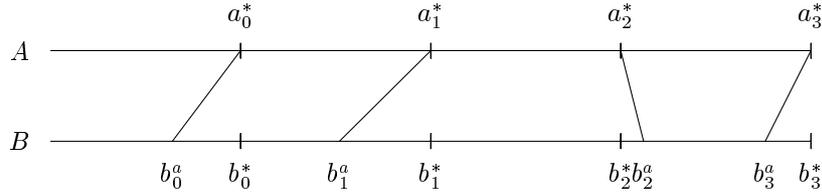


Figura 5.1: Separarea în subprobleme în algoritmul Valiant.

Găsirea fiecărui b_i^a se face în două etape; în prima se determină secvența B_j în care se găsește b_i^a (o căutare la nivel de bloc, mai grosieră); este pasul 2 al algoritmului; el se execută în timp constant folosind $\sqrt{n}\sqrt{m}$ procesoare (se aplică de \sqrt{n} ori, în paralel, algoritmul 5.4, care necesită \sqrt{m} procesoare). În a doua (pasul 3 al algoritmului) se rafinează căutarea la nivel de element; efortul este identic cu cel de la pasul 2. În acest moment, problema inițială a fost spartă în \sqrt{n} probleme, în care fiecare A_i trebuie interclasat cu o secvență din B care începe cu elementul b_{i-1}^a (convenim că $b_{-1}^a = b_0$) și se termină cu elementul precedent lui b_i^a ; vom nota această secvență cu B'_i .

Numărul de comparații efectuate respectă regula recursivă $T(n) \leq T(\sqrt{n}) + O(1)$, ceea ce conduce la $T(n) = O(\log \log n)$. La primul nivel de recursie se folosesc $\sqrt{n}\sqrt{m} \leq (n+m)/2$ procesoare. În continuare, pentru fiecare subproblemă este adevărată o relație de acest gen, ceea ce conduce la un necesar de maximum $(n+m)/2$ procesoare.

Algoritmul acesta are dezavantajul de a nu împărți în mod egal decât secvența A atunci când sunt formate subproblemele. Dacă aceasta nu afectează foarte mult timpul teoretic total (în general, dar nu și în cel mai defavorabil caz), în schimb pune mari probleme în alocarea procesoarelor.

Se poate demonstra că problema interclasării a două secvențe de lungime n se poate rezolva cu $O(n)$ procesoare în cel puțin $\Omega(\log \log n)$ operații, deci algoritmul Valiant este optim. În plus, costul lui este $O(n \log \log n)$, ceea ce e acceptabil.

5.3.2 Un algoritm pe CREW PRAM, pentru număr mic de procesoare

Să vedem cum se pot aplica ideile din algoritmul Valiant pentru cazul $p \ll n$. Să presupunem că împărțim A și B în secvențe de lungime egală, n/p , respectiv m/p ; păstrăm notațiile de mai sus; evident că acum $A^* = \langle a_0^*, \dots, a_{p-1}^* \rangle$ și $B^* = \langle b_0^*, \dots, b_{p-1}^* \rangle$. Algoritmul are două etape.

ALGORITM 5.12 (*Valiant-Kruskal*) (interclasare $A \perp B$ pe CREW PRAM, $p \ll n$)

1. Găsește poziția exactă a fiecărui a_i^* în secvența B (formează secvențele B'_i).
2. Fiecare procesor P_i interclasează secvențial secvențele A_i și B'_i .

Efortul mare de calcul se află în etapa a doua; putem spune că aceasta necesită $(n + m)/p$ comparații; numărul total de comparații care trebuie efectuat pentru interclasarea celor p perechi de secvențe $A_i \perp B_i'$ este $n + m$ în cel mai defavorabil caz; presupunând secvențele B_i' de lungimi aproximativ egale, rezultă complexitatea de mai sus.

Pentru prima etapă sunt mai multe variante posibile. Dacă fiecare P_i ar căuta separat poziția lui a_i^* în B , aceasta s-ar face printr-o căutare binară, în $\log m$ comparații.

Se poate însă și mai bine. Se interclasează A^* și B^* folosind algoritmul Valiant, deci în timp $O(\log \log p)$ (cele două secvențe au împreună $2p$ elemente; cu p procesoare se poate aplica algoritmul Valiant fără modificare). În acest moment s-a identificat, pentru fiecare a_i^* , segmentul $B_{j(i)}$ în care trebuie continuată căutarea. Fiecare procesor va căuta acum poziția exactă pentru un a_i^* într-o secvență de lungime m/p , deci făcând $\log(m/p)$ comparații. Cum până acum am presupus că $n \leq m$, putem inversa rolurile între A și B , și deci timpul necesar primei etape va fi în acest caz $\log(n/p) + O(\log \log p)$, iar în total:

$$T(n, p) = (n + m)/p + \log(n/p) + O(\log \log p). \quad (5.1)$$

Vom face acum o mică paranteză, pentru a fi a mai clar cum se găsește, după interclasarea secvențelor A^* și B^* , indicele $j(i)$, pentru fiecare a_i^* . Privind figura 5.1, reamintindu-ne că acolo $A^* \perp B^* = \langle a_0^*, b_0^*, a_1^*, b_1^*, b_2^*, a_2^*, a_3^*, b_3^* \rangle$ și revăzând încă o dată pasul 2 al algoritmului Valiant, observăm că $j(0) = 0$, $j(1) = 1$, $j(2) = 3$, $j(3) = 3$; în general, $j(i)$ este indicele celui mai mic element din B^* mai mare decât a_i^* sau $j(i)$ este indicele primului element din B^* care se află în dreapta lui a_i^* în $A^* \perp B^*$; să notăm cu $rp(a_i^*)$ acest indice (rp vine de la *right pair*, perechea din dreapta). În mod analog, vom nota cu $lp(a_i^*)$ (*left pair*) indicele primului element din B^* care se află în stânga lui a_i^* în $A^* \perp B^*$. Se vede imediat că $rp(a_i^*) = lp(a_i^*) + 1$. Convenim ca $lp(a_i^*) = -1$ dacă nu există element din B^* mai mic decât a_i^* . Să mai observăm o relație interesantă (a cărei demonstrație e lăsată cititorului):

$$rp(a_i^*) = rang(a_i^*) - i \quad (5.2)$$

unde rangul se calculează în secvența $A^* \perp B^*$.

Deci, în algoritmul Valiant, calculăm, pentru fiecare a_i^* rangul acestuia în $A^* \perp B^*$, ceea ce e echivalent cu sortarea acestei secvențe (deci cu interclasarea), apoi $j(i) \equiv rp(a_i^*)$ se calculează cu (5.2). Același lucru se face și în algoritmul Valiant-Kruskal.

5.3.3 Mediana secvenței $A \perp B$ (partiționare mediană)

Vom încerca să găsim un algoritm pentru calculul paralel al medianei secvenței $A \perp B$, fără a face însă interclasarea; definim mediana secvenței $A \perp B$ ca fiind valoarea cu rangul $\lceil (m + n)/2 \rceil - 1$. Soluția prezentată este inspirată de [27].

Problema se poate formula altfel și se numește partiționare mediană. Să se partiționeze A și B în câte două secvențe, $A = \langle A_0, A_1 \rangle$, respectiv $B = \langle B_0, B_1 \rangle$,

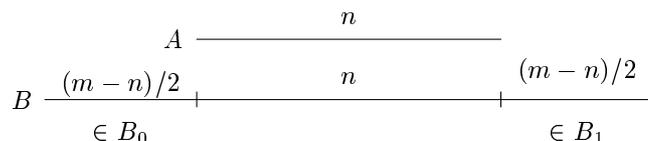


Figura 5.2: Reducerea căutării la $n = |A|$ elemente din B , când $m = |B| > n$ (aici $m - n$ este par).

astfel încât $|A_0| + |B_0| = \lceil (m+n)/2 \rceil$, $|A_1| + |B_1| = \lfloor (m+n)/2 \rfloor$ și orice element din A_0 sau B_0 este mai mic decât orice element din A_1 sau B_1 . În consecință, mediana secvenței $A \perp B$ va fi ultima valoare (cea mai mare) din $A_0 \perp B_0$.

În acest fel se face legătura cu interclasarea. După aflarea partiționării mediane, se pot găsi, în paralel, $A_0 \perp B_0$ și $A_1 \perp B_1$, problema de interclasare fiind spartă în două probleme de dimensiuni mai mici și *egale*.

Vom prezenta două idei care contribuie la găsirea rapidă a partiționării mediane. Întâi, că se poate lucra cu numai n elemente din B , deci cu două secvențe de dimensiuni egale. Apoi, că se pot utiliza comparații paralele.

Reducerea la secvențe egale. Dacă ne uităm la figura 5.2, în care secvențele A și B sunt reprezentate prin segmente de lungimi proporționale cu numărul de elemente, observăm că primele $m - \lfloor (m+n)/2 \rfloor$ din B , adică cele din $B_l = \langle b_0, b_1, \dots, b_{m - \lfloor (m+n)/2 \rfloor - 1} \rangle$ aparțin în mod sigur lui B_0 ; aceasta rezultă imediat din faptul că în secvența B există $\lfloor (m+n)/2 \rfloor$ elemente mai mari decât cele din B_l , adică jumătate dintre elementele din $A \perp B$. În mod analog, este sigur că elementele din $\langle b_{\lfloor (m+n)/2 \rfloor}, \dots, b_{m-1} \rangle$ aparțin secvenței B_1 . În acest fel este redusă dimensiunea problemei partiționării.

Să încorporăm această idee într-o procedură. Vom nota cu la și ua doi indici în secvența A , astfel încât toate elementele la stânga lui la sunt în A_0 , iar toate elementele la dreapta lui ua sunt în A_1 . Cele între la și ua sunt încă în A ; inițial $la = 0$ și $ua = n - 1$; la terminarea partiționării (dar mai avem până acolo), $la = ua + 1$. Se definesc în mod analog indicii lb și ub pentru secvența B (inițial $ub = m - 1$). Algoritmul de reducere a dimensiunii este (pentru rigurozitate, uităm aici că am presupus $n < m$ și considerăm cazul general):

ALGORITHM 5.13 (reducerea dimensiunii problemei partiționării mediane)

procedură *reduce*(la, ua, lb, ub)

1. **dacă** $n > m$ **atunci**
 1. $ua \leftarrow la + \lceil (m+n)/2 \rceil - 1$
 2. $la \leftarrow la + n - \lfloor (m+n)/2 \rfloor$
2. **altfel dacă** $n < m$ **atunci**
 1. $ub \leftarrow lb + \lceil (m+n)/2 \rceil - 1$
 2. $lb \leftarrow lb + m - \lfloor (m+n)/2 \rfloor$

Apelul procedurii se face prin *reduce*($0, n - 1, 0, m - 1$).

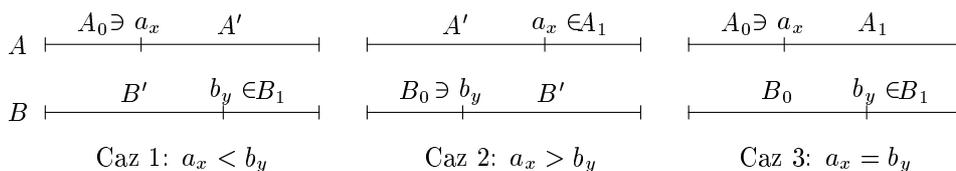


Figura 5.3: Comparație între elemente cu indici complementari; concluzii posibile.

Pentru a introduce ideea care permite efectuarea mai multor comparații în paralel în căutarea partiționării mediane, vom putea considera $n = m$.

Reducerea căutării prin utilizarea indicilor complementari. Vom încerca să extragem maximum de informație din comparații între a_x și b_y , unde $0 \leq x \leq n-1$ și $y = n-1-x$; numim y indicele *complementar* al lui x ; se observă că a_x și b_y sunt egal distanțate de începutul secvenței A , respectiv sfârșitul secvenței B .

Dacă $a_x < b_y$, atunci putem pune $\langle a_0, \dots, a_x \rangle$ în A_0 și $\langle b_y, \dots, b_{n-1} \rangle$ în B_1 ; într-adevăr, sunt $n-1-x$ elemente în A mai mari decât a_x și $n-y = x+1$ elemente în B mai mari decât a_x , în total n elemente din $A \perp B$ mai mari decât a_x , deci mediana este mai mare decât a_x . În mod analog se demonstrează că mediana este mai mică decât b_y . În figura 5.3, cazul 1, este ilustrată această situație; cu A' și B' sunt notate subsecvențele din A , respectiv B , în care trebuie continuată căutarea medianei.

Dacă $a_x > b_y$, situația se inversează: $\langle a_x, \dots, a_{n-1} \rangle$ aparține lui A_1 și $\langle b_0, \dots, b_y \rangle$ aparține lui B_0 (vezi figura 5.3, cazul 2).

În fine, dacă $a_x = b_y$, am descoperit partiționarea mediană și mediana însăși, care este a_x .

Reducere paralelă a căutării. Pentru a face mai multe procesoare să participe la reducerea secvențelor A' și B' , vom face comparații în mai multe puncte. Fie x_0, \dots, x_{t-1} un număr de t indici diferiți, în ordine crescătoare, și indicii complementari y_0, \dots, y_{t-1} . Comparațiile între a_{x_i} și b_{y_i} pot fi făcute în paralel. Dacă notăm $c_i = 0$ pentru $a_{x_i} < b_{y_i}$ și $c_i = 1$ altfel, atunci în secvența $\langle c_0, \dots, c_{t-1} \rangle$ există cel mult două elemente vecine diferite. Deci există un singur j pentru care $a_{x_j} < b_{y_j}$ și $a_{x_{j+1}} \geq b_{y_{j+1}}$. Această afirmație este evident adevărată deoarece secvența $\langle a_{x_0}, \dots, a_{x_{t-1}} \rangle$ este crescătoare, în timp ce $\langle b_{y_0}, \dots, b_{y_{t-1}} \rangle$ este descrescătoare; ele se pot "intersecta" cel mult o dată.

În concluzie, dacă notăm A' și B' părțile din A și B care mai pot conține mediana după efectuarea celor t comparații, se poate ajunge într-unul din următoarele patru cazuri.

1. Dacă există cel puțin un j pentru care $a_{x_j} = b_{y_j}$, atunci:

$$\begin{array}{lll} A_0 = \langle a_0, \dots, a_{x_j} \rangle & B_0 = \langle b_0, \dots, b_{y_{j-1}} \rangle & A' = \emptyset \\ A_1 = \langle a_{x_{j+1}}, \dots, a_{n-1} \rangle & B_1 = \langle b_{y_j}, \dots, b_{n-1} \rangle & B' = \emptyset \end{array}$$

2. Dacă $a_{x_0} < b_{y_0}$ și $a_{x_{t-1}} < b_{y_{t-1}}$, atunci ($a_{x_{t-1}}$ și $b_{y_{t-1}}$ sunt ca în figura 5.3, cazul 1):

$$\begin{array}{lll} A_0 = \langle a_0, \dots, a_{x_{t-1}} \rangle & B_0 = \emptyset & A' = \langle a_{x_{t-1}+1}, \dots, a_{n-1} \rangle \\ A_1 = \emptyset & B_1 = \langle b_{y_{t-1}}, \dots, b_{n-1} \rangle & B' = \langle b_0, \dots, b_{y_{t-1}-1} \rangle \end{array}$$

3. Dacă $a_{x_0} < b_{y_0}$, și există un indice j pentru care $a_{x_{j-1}} < b_{y_{j-1}}$, dar $a_{x_j} > b_{y_j}$, atunci:

$$\begin{array}{lll} A_0 = \langle a_0, \dots, a_{x_{j-1}} \rangle & B_0 = \langle b_0, \dots, b_{y_j} \rangle & A' = \langle a_{x_{j-1}+1}, \dots, a_{x_{j-1}} \rangle \\ A_1 = \langle a_{x_j}, \dots, a_{n-1} \rangle & B_1 = \langle b_{y_{j-1}}, \dots, b_{n-1} \rangle & B' = \langle b_{y_j+1}, \dots, b_{y_{j-1}} \rangle \end{array}$$

4. Dacă $a_{x_0} > b_{y_0}$ (a_{x_0} și b_{y_0} sunt ca în figura 5.3, cazul 2), atunci:

$$\begin{array}{lll} A_0 = \emptyset & B_0 = \langle b_0, \dots, b_{y_0} \rangle & A' = \langle a_0, \dots, a_{x_0-1} \rangle \\ A_1 = \langle a_{x_0}, \dots, a_{n-1} \rangle & B_1 = \emptyset & B' = \langle b_{y_0+1}, \dots, b_{n-1} \rangle \end{array}$$

Să mai observăm că partiționarea mediană a fost găsită în cazul 1 întotdeauna, în cazul 2 dacă $x_{t-1} = n - 1$, în cazul 3 dacă $x_{j-1} = x_j - 1$ și în cazul 4 dacă $x_0 = 0$.

Procedura următoare înglobează toate aceste rezultate. Toate notațiile se păstrează; vectorul c este folosit pentru memorarea sensului inegalităților.

ALGORITHM 5.14 (partiționare mediană, pe CRCW PRAM, $\text{id} \equiv k$)

procedură *partiționare_mediană*(la, ua, lb, ub)

1. **dacă** $k = 0$ **atunci** *reduce*(la, ua, lb, ub)

2. **cât timp** $la \leq ua$

1. $c_k \leftarrow 1$, $exit \leftarrow 0$, $d \leftarrow -1$

2. $x_k \leftarrow la + (k + 1)[(ua - la + 1)/(p + 1)]$ {se aleg indicii x_i echidistanți}

3. $y_k \leftarrow la + ub - x_k$ {complementarul lui x_k }

4. **dacă** $a_{x_k} > b_{y_k}$ **atunci**

1. **dacă** $k = 0$ {cazul 4}

1. $ua \leftarrow x_0 - 1$, $lb \leftarrow y_0 + 1$, $exit \leftarrow 1$

5. **altfel dacă** $a_{x_k} < b_{y_k}$ **atunci**

1. **dacă** $k = p - 1$ {cazul 2}

1. $la \leftarrow x_{p-1} + 1$, $ub \leftarrow y_{p-1} - 1$, $exit \leftarrow 1$

2. **altfel** $c_k \leftarrow 0$ {cazul 3}

6. **altfel** {cazul 1, $a_{x_k} = b_{y_k}$ }

1. $d \leftarrow k$

7. **dacă** $exit = 0$ **atunci**

1. **dacă** $d = k$ **atunci** {din nou cazul 1}

1. $la \leftarrow x_k + 1$, $ua \leftarrow x_k$, $lb \leftarrow y_k$, $ub \leftarrow y_k - 1$

2. **dacă** $c_k = 0$ și $c_{k+1} = 1$ **atunci** {din nou cazul 3}

1. $la \leftarrow x_k + 1$, $ua \leftarrow x_{k+1} - 1$, $lb \leftarrow y_{k+1} + 1$, $ub \leftarrow y_k - 1$

În final, se obțin $A_0 = \langle a_0, \dots, a_{ua} \rangle$, $B_0 = \langle b_0, \dots, b_{ub} \rangle$, $A_1 = \langle a_{1a}, \dots, a_{n-1} \rangle$, $B_1 = \langle b_{1b}, \dots, b_{n-1} \rangle$.

Pot apărea scrieri simultane în instrucțiunile 2.1 și 2.6.1; în prima, aceasta se evită foarte ușor lăsând un singur procesor să scrie variabilele *exit* și *d*; în 2.6.1, indiferent ce procesor scrie (în modelele ARBITRARY sau PRIORITY CRCW, desigur), algoritmul funcționează corect (vezi și problema 5.3.6). În variabila *exit* scrie cel mult un procesor (P_0 sau P_{p-1}). De asemenea, în *la*, *ua*, *lb*, *ub*, scrie un singur procesor, în general (vezi, totuși, problema 5.3.8).

Procesoarele lucrează în paralel. Totuși, unele pot sta ceva mai mult timp într-o iterație (în instrucțiuni **dacă**), de unde necesitatea unor sincronizări. De exemplu, înainte de 2.7, sincronizarea este necesară pentru a asigura actualizarea variabilelor *exit*, *C* și *d* în iterația curentă. De asemenea este necesară sincronizarea la începerea fiecărei iterații a buclei 2, pentru a permite actualizarea variabilelor *la*, *ua*, *lb*, *ub*.

Algoritmul poate fi ușor adaptat unui EREW PRAM. Se lasă căutarea să continue chiar dacă $a_{x_k} = b_{y_k}$ (de altfel nici nu se mai verifică egalitatea). În acest scop se elimină variabila *exit* și tot ce urmează după instrucțiunea 2.7 inclusiv. De asemenea, se elimină instrucțiunile 2.6, 2.6.1 și se înglobează cazul $a_{x_k} = b_{y_k}$ în instrucțiunea 2.5, transformând semnul $<$ în \leq .

Să vedem care este timpul de execuție al acestui algoritm; se observă că în fiecare iterație se execută în paralel o singură comparație. Cum *A* și *B* sunt împărțite în $p + 1$ secvențe de lungimi egale, iar în iterația următoare căutarea este continuată numai într-una din aceste perechi de secvențe, înseamnă că dimensiunea problemei se reduce de $p + 1$ ori. Dacă z este numărul total de iterații (după care $ua < la$), atunci z este cel mai mic întreg pentru care $n/(p + 1)^z < 1$. Adică se obține $z = T(n, p) = \lceil \log n / \log(p + 1) \rceil + 1$.

Cum complexitatea secvențială a acestui algoritm este $T(n, 1) = \log n / \log 2$, eficiența sa este $\varepsilon(n, p) \approx \log(p + 1) / p$, deci nu grozavă; algoritmul este bun pentru un număr mic de procesoare; pentru $p = 2$, eficiența este 0.79, însă pentru $p = 32$ ajunge la 0.16. (Acesta este însă cel mai bun algoritm paralel pentru partiționare mediană cunoscut deocamdată.)

5.3.4 Interclasare prin multipartiționare

Să mai prezentăm un algoritm de interclasare, tot pe CREW PRAM și tot pentru cazul $p \ll n$. Acesta are marele avantaj de a crea partiții egale, în număr de $t \leq p$, și pentru *A* și pentru *B*, ceea ce implică și un timp de execuție ceva mai scurt (aproximativ același, indiferent de valorile din *B*), și o planificare mai simplă a procesoarelor în comparație cu algoritmul 5.12 (Valiant-Kruskal). În acest scop este folosit algoritmul de partiționare mediană descris în secțiunea anterioară.

Problema multipartiționării este de a găsi secvențele A_i , B_i , astfel încât $|A_i| + |B_i| = \lceil (m + n) / t \rceil$ și $A \perp B = \langle A_0 \perp B_0, \dots, A_{t-1} \perp B_{t-1} \rangle$; deci, de a partiționa $A \perp B$ în secvențe de dimensiuni egale (eventual mai puțin ultima), fără însă a face interclasarea.

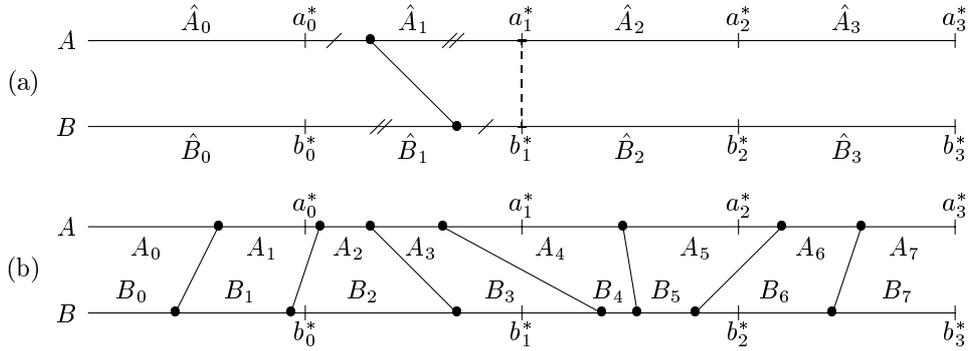


Figura 5.4: (a) linii de separare; exemplu de partiționare mediană pentru o pereche \hat{A}_i, \hat{B}_j ; (b) exemplu de multipartiționare în cazul $A^* \perp B^* = \langle a_0^*, b_0^*, a_1^*, b_1^*, b_2^*, a_2^*, a_3^*, b_3^* \rangle$.

Considerăm, pentru început, A și B compuse din subsecvențe de dimensiuni egale $s = \lceil (m+n)/t \rceil$ (ultimele fiind mai scurte, eventual):

$$\begin{aligned} \hat{A}_0 &= \langle a_0, \dots, a_{s-1} \rangle, \hat{A}_1 = \langle a_s, \dots, a_{2s-1} \rangle, \dots \\ \hat{B}_0 &= \langle b_0, \dots, b_{s-1} \rangle, \hat{B}_1 = \langle b_s, \dots, b_{2s-1} \rangle, \dots \end{aligned}$$

Definim, ca în algoritmul 5.11, $A^* = \langle a_i^* \rangle = \langle a_{is-1} \mid i = 1, \dots, \lfloor n/s \rfloor \rangle$, secvența formată din ultimele elemente ale subsecvențelor A_i ; B^* se definește analog. Secvența $A^* \perp B^*$ poate fi determinată în timp $\log \log t$, utilizând algoritmul Valiant (sunt cel puțin t procesoare, pentru a interclasa secvențe de lungime $\lfloor n/s \rfloor + \lfloor m/s \rfloor \approx (m+n)/s \approx t$, deci se poate atinge timpul optim).

Vom face multipartiționarea prin executarea în paralel a $t-1$ procese de partiționare mediană pe perechi de subsecvențe \hat{A}_i, \hat{B}_j , cu $j = rp(a_i^*)$ sau $i = rp(b_j^*)$ (pentru noțiunea de pereche din dreapta, $rp()$, vezi secțiunea 5.3.2; tot acolo se află modul de calcul al acestor entități, care se obțin imediat din algoritmul Valiant).

Să pregătim terenul cu câteva definiții și observații. În reprezentarea grafică din figura 5.4a, în care secvențele sunt figurate prin segmente de dreaptă, numim punct de separare un indice care separă o secvență în două subsecvențe (elementul cu indicele respectiv fiind în secvența din stânga sau în cea din dreapta), și linie de separare, o dreaptă care unește două puncte de separare a două secvențe; în figura 5.4a, sunt reprezentate două astfel de linii, una continuă, alta punctată.

OBSERVAȚIA 5.1 *Linia de separare prin a_i^* și b_j^* , cu $j = rp(a_i^*)$, lasă în stânga ei cel puțin $(i+j+1)s$ elemente din $A \perp B$ (exemplu: linia punctată din figura 5.4a).*

Demonstrație. Deoarece $a_i^* \leq b_j^*$, vom număra elementele din $A \perp B$ mai mici decât a_i^* . Sunt $(i+1)s$ elemente din A mai mici decât a_i^* . Cum $b_{j-1}^* < a_i^*$, atunci sunt js elemente din B mai mici decât a_i^* ; nu putem spune nimic despre elementele din B între b_{j-1}^* și b_j^* . ■

OBSERVAȚIA 5.2 *Linia de separare mediană a secvenței $\hat{A}_i \perp \hat{B}_j$, cu $j = rp(a_i^*)$, lasă în stânga sa exact $(i + j + 1)s$ elemente din $A \perp B$ (exemplu: linia continuă oblică din figura 5.4).*

Demonstrație. Imediată, din observația 5.1 și din figură, unde segmentele notate cu același număr de linii mici oblice au același număr de elemente (din cele $(i + j + 2)s$ elemente din stânga liniei de separare $a_i^* b_j^*$, se extrag s ; observația 5.1 este de ajutor pentru că arată că linia continuă este în stânga celei punctate). ■

Prin simetrie, cele două observații sunt valabile și pentru perechile b_i^* , a_j^* , cu $j = rp(b_i^*)$. În sfârșit, o ultimă precizare:

OBSERVAȚIA 5.3 *Este adevărată următoarea relație:*

$$\{i + rp(a_i^*) \mid i = 0, \dots, \lfloor n/s \rfloor - 1\} \cup \{i + rp(b_i^*) \mid i = 0, \dots, \lfloor m/s \rfloor - 1\} = \{0, 1, \dots, \lfloor n/s \rfloor + \lfloor m/s \rfloor - 1\} \quad (5.3)$$

Demonstrație. Prin inducție, vezi problema 5.3.9. ■

Această observație arată că sumele dintre indici și perechile lor din dreapta formează exact șirul primelor $\lfloor n/s \rfloor + \lfloor m/s \rfloor - 1$ numere naturale. Împreună cu celelalte observații, aceasta permite reducerea problemei multipartiționării în t segmente la $t - 1$ partiționări mediane. Figura 5.4b ilustrează acest lucru; pe scurt, observațiile de mai sus arată că între oricare două linii de separare consecutive se află exact s elemente din $A \perp B$. Algoritmul general de multipartiționare este următorul:

ALGORITM 5.15 (multipartiționare prin partiționare mediană)

1. găsește $A^* \perp B^*$ cu algoritmul 5.11
2. **în paralel**
 1. **pentru** $i = 0 : \lfloor n/s \rfloor - 1$, **în paralel**
 1. partiționează median perechea $\hat{A}_i, \hat{B}_{rp(a_i^*)}$ utilizând $\lfloor p/(t - 1) \rfloor$ procesoare
 2. **pentru** $j = 0 : \lfloor m/s \rfloor - 1$, **în paralel**
 1. partiționează median perechea $\hat{B}_j, \hat{A}_{rp(b_j^*)}$ utilizând $\lfloor p/(t - 1) \rfloor$ procesoare

Timpul necesar pentru multipartiționare este

$$\lceil \log((n + m)/t) / \log(\lfloor p/(t - 1) \rfloor + 1) \rceil + O(\log \log p).$$

Interclasarea prin multipartiționare este acum ușor de făcut. Se utilizează $t = p/2$ în algoritmul 5.15, pentru a obține $p/2$ perechi disjuncte de subsecvențe, care se vor interclasa, fiecare de către două procesoare (după un algoritm banal: fiecare procesor interclasează începând de la un capăt; eficiența este 1).

ALGORITM 5.16 (*Xiong*) (interclasare prin multipartiționare pe CREW PRAM)

1. apelează algoritmul 5.15 de multipartiționare cu $t = p/2$
2. **pentru** $i = 0 : p/2 - 1$, **în paralel**
 1. interclasează A_i și B_i utilizând două procesoare

Au loc citiri simultane deoarece procesoare diferite pot lucra cu aceeași subsecvență \hat{A}_i sau \hat{B}_j (în medie două procesoare operează pe o subsecvență). Complexitatea algoritmului este simplu de calculat:

$$T(n, p) = \frac{m+n}{p} + \frac{1}{\log 3} \log \frac{m+n}{p} + O(\log \log p). \quad (5.4)$$

Primul termen corespunde interclasării secvențiale (și e optim), iar următorii multi-partiționării. Comparând cu complexitatea algoritmului Valiant-Kruskal (5.1) se constată diferența doar la al doilea termen, ceva mai mic aici.

Marele avantaj al algoritmului rămâne însă faptul că se crează subprobleme de dimensiuni egale, ceea ce implică o încărcare echilibrată a procesoarelor.

5.3.5 Interclasare pe arhitecturi cu memorie distribuită

Titlul acestei secțiuni este foarte promițător; din păcate, cum veți vedea, conținutul nu va fi prea generos.

Începem cu partea mai consistentă și care ne va fi de mare folos mai departe, anume interclasarea cu două procesoare. S-ar părea că nu e foarte mult de discutat pe tema asta. Algoritmul cel mai simplu a fost deja sugerat și el este următorul: P_0 trimite secvența sa A_0 de lungime n lui P_1 și recepționează de la acesta secvența A_1 (tot de lungime n); apoi P_0 interclasează A_0 și A_1 , oprindu-se atunci când a obținut primele n elemente; în același timp, și P_1 interclasează A_0 și A_1 , dar începând de la sfârșitul secvențelor, oprindu-se și el atunci când a obținut n elemente; în acest moment secvența $A_0 \perp A_1$ este distribuită celor două procesoare. Cum P_0 are acum toate cele n elemente mai mici decât cele ale lui P_1 , putem asimila operația executată cu $exch(A_0, A_1)$, o generalizare a aceleiași operații la nivel de element. Presupunând A_i și A_j ordonate, după executarea $exch(A_i, A_j)$ se obține $A_i < A_j$, ca în procedura care urmează; ca și în cazul operației la nivel de element, vom nota $exch(i, j)$ atunci când nu există posibilitatea de confuzie; presupunem că procesoarele i și j sunt vecine.

ALGORITM 5.17 (interclasare cu două procesoare, pe arhitecturi cu memorie distribuită, $id \equiv k$, secvența locală este notată \tilde{A})

procedură $exch(i, j)$

1. **dacă** $i = k$ **atunci**

1. **în paralel** $send(\tilde{A}, j)$, $recv(B, j)$

2. $\tilde{A} \leftarrow$ primele n elemente din $\tilde{A} \perp B$, prin interclasare secvențială

2. **dacă** $j = k$ **atunci**

1. **în paralel** $send(\tilde{A}, i)$, $recv(B, i)$

2. $\tilde{A} \leftarrow$ ultimele n elemente din $\tilde{A} \perp B$, prin interclasare secvențială
(în ordine inversă)

Timpul de execuție este de $n\alpha$ pentru interclasarea propriu-zisă și de $(\sigma + n\beta)$ pentru comunicație. Numărul de comparații este jumătate din cel necesar, în cel mai

rău caz, interclasării secvențiale, și mai bine de atât nu se poate. Comunicația însă, poate fi îmbunătățită.

Să presupunem că procesoarele nu își comunică de la început secvențele în întregime, ci o fac în pachete de lungime ν . Ideea care se ascunde aici e de a transmite astfel de pachete până când procesoarele sunt sigure (și vor fi amândouă în același timp) că au fiecare suficiente elemente din secvența celuilalt pentru a reuși să obțină partea din rezultat pe care trebuie s-o calculeze. Speranța este de a ajunge aici înainte de a se comunica întregile secvențe.

Se începe cu P_0 trimițând lui P_1 ultimul pachet din A_0 și primind de la acesta primul pachet din A_1 ; este normal să se procedeze așa, din moment ce P_0 începe interclasarea de la început, adică cu elementele mai mici, în timp ce P_1 începe cu elementele mai mari. Să raționăm din punctul de vedere al lui P_0 (partenerul său va proceda în mod analog); ca în algoritmul 5.17, notăm secvența locală cu \tilde{A} ($= \langle \tilde{A}_0, \dots, \tilde{A}_{n/\nu-1} \rangle$) și cu B copia locală a lui A_1 , conținând deocamdată numai ν elemente. Să comparăm $\tilde{a}_{n-\nu}$ cu $b_{\nu-1}$, uitându-ne la figura 5.3 și amintindu-ne că indicii acestor două elemente sunt complementari; dacă $\tilde{a}_{n-\nu} \leq b_{\nu-1}$, atunci este sigur că primele n elemente din $\tilde{A} \perp B$ sunt deținute de P_0 . Dacă nu, atunci se continuă schimbând încă un pachet; se compară acum $\tilde{a}_{n-2\nu}$ și $b_{2\nu-1}$, ș.a.m.d.

ALGORITM 5.18 (interclasare cu două procesoare, $id \equiv k$, comunicație prin pachete)

procedură *exch1*(i, j)

1. $l \leftarrow 1$, $gata \leftarrow 0$, $B \leftarrow \emptyset$
2. **cât timp** $l \leq n/\nu$ și $gata = 0$ {schimbă un pachet}
 1. **dacă** $i = k$ **atunci**
 1. **în paralel** **send**($\tilde{A}_{n/\nu-l}, j$), **recv**(\tilde{B}, j)
 2. $B \leftarrow \langle B, \tilde{B} \rangle$ {concatenează în B , la dreapta}
 3. **dacă** $\tilde{a}_{n-l\nu} \leq b_{l\nu-1}$ **atunci** $gata \leftarrow 1$
 2. **dacă** $j = k$ **atunci**
 1. **în paralel** **send**(\tilde{A}_{l-1}, i), **recv**(\tilde{B}, i)
 2. $B \leftarrow \langle \tilde{B}, B \rangle$ {concatenează în B , la stânga}
 3. **dacă** $\tilde{a}_{l\nu-1} \geq b_0$ **atunci** $gata \leftarrow 1$
 3. $l \leftarrow l + 1$
- 3 **dacă** $i = k$ **atunci**
 1. $\tilde{A} \leftarrow$ primele n elemente din $\tilde{A} \perp B$, prin interclasare secvențială
4. **altfel**
 1. $\tilde{A} \leftarrow$ ultimele n elemente din $\tilde{A} \perp B$, prin interclasare secvențială (în ordine inversă)

Avantajele și dezavantajele unei astfel de abordări se pot evidenția pe cazul simplu $\nu = n/2$. Probabilitatea ca un singur pachet comunicat să fie suficient este $1/2$; într-adevăr: mediana secvenței $A_0 \perp A_1$ se află în jumătatea din stânga a lui A_1 sau în jumătatea din dreapta a lui A_0 , cu aceeași probabilitate de a se afla în una din

celelalte jumătăți. În primul caz e suficient un pachet, în al doilea trebuie comunicate ambele. Atunci, în medie, timpul de comunicație al variantei $exch1()$ este: $\frac{1}{2}(\sigma + \frac{1}{2}n\beta) + \frac{1}{2}(2\sigma + n\beta) = \frac{3}{2}\sigma + \frac{3}{4}n\beta$, deci, grosso modo, puțin peste trei sferturi din cel în care se comunicau secvențele întregi, și, în concluzie, metoda este benefică. Nu trebuie mizat însă foarte mult pe ea; dacă există multe perechi de procesoare care fac astfel de interclasări, iar apoi urmează, într-un fel sau altul, o sincronizare (se schimbă partenerii de comunicație, de pildă), timpul total de execuție va fi dictat mai degrabă de timpul cel mai mare de interclasare și nu de timpul mediu.

Pe tema algoritmului 5.18 se poate broda mult. Se poate trimite întâi un pachet mai mare, de lungime $n/2$, de exemplu, apoi altele mai mici. Probabil că problema nu merită un studiu foarte amănunțit; ideea merită pusă în practică, acordând experimental valoarea (valorile) lui ν .

Din nefericire, de îndată ce sunt mai mult de două procesoare, algoritmul de interclasare se complică, în sensul că nici una dintre topologiile curente nu se potrivește structurii comunicației. Adaptarea multipartiționării la arhitecturile cu memorie distribuită implică o structură neregulată a comunicației. Totuși, pentru o posibilă soluție paralelă, vezi problema 5.4.17.

Probleme

P 5.3.1 Scrieți (totuși...) algoritmul secvențial de interclasare a două secvențe.

P 5.3.2 Să se demonstreze că în algoritmul 5.11 sunt suficiente \sqrt{nm} procesoare.

P 5.3.3 Demonstrați că $T(n) = \log \log n$ verifică $T(n) = T(\sqrt{n}) + 1$.

P 5.3.4 Găsiți încă un argument pentru aproximarea cu $(n+m)/p$ a complexității etapei a doua din algoritmul Valiant-Kruskal.

P 5.3.5 Demonstrați relația (5.2).

P 5.3.6 Dacă x și y sunt indici complementari și $a_x = b_y$, demonstrați că valoarea medianei secvenței interclasate $A \perp B$ este a_x , chiar dacă în fiecare dintre secvențele A și B pot exista elemente egale.

P 5.3.7 Scrieți un algoritm de partiționare mediană a secvențelor ordonate A și B , de lungime egală n , pe un CREW PRAM, folosind $p = n$ procesoare.

P 5.3.8 În algoritmul 5.14, este posibil ca două procesoare să scrie simultan în variabilele la, ua, lb, ub ?

P 5.3.9 Dacă $X = \langle x_0, \dots, x_{n-1} \rangle$ și $Y = \langle y_0, \dots, y_{m-1} \rangle$ sunt două secvențe ordonate crescător, să se arate că $\{i + rp(x_i) \mid i = 0, \dots, n-1\} \cup \{j + rp(y_j) \mid j = 0, \dots, m-1\} = \{0, \dots, n+m-1\}$. Altfel zis, să se demonstreze (5.3).

P 5.3.10 O generalizare a noțiunii de partiționare mediană este k -separarea. Problema este definită astfel: dându-se secvențele ordonate A și B , să se găsească A_0, A_1, B_0, B_1 astfel

încât $|A_0| + |B_0| = k$, $|A_1| + |B_1| = n + m - k$ și $A \perp B = \langle A_0 \perp B_0, A_1 \perp B_1 \rangle$. Să se demonstreze că, pentru $n \leq m$, k -separarea este echivalentă cu partiționarea mediană pentru:

$$\begin{cases} A_{[k]} \perp B_{[k]}, & \text{dacă } 1 \leq k \leq n \\ A \perp (B_{[k]} - B_{[k-n]}), & \text{dacă } n < k \leq m \\ (A - A_{[k-n]}) \perp (B - B_{[k-m]}), & \text{dacă } m < k \leq n + m \end{cases}$$

unde am notat prin $A_{[k]}$ secvența primelor k elemente din A .

P 5.3.11 Presupunând comunicația realizată half duplex, care este cel mai bun mod de a implementa $exch(a_i, a_j)$, cele două elemente aparținând unor procesoare vecine ?

P 5.3.12 Pe un CREW PRAM, $exch(a_i, a_j)$ se poate implementa prin instrucțiunea: **dacă** $a_i > a_j$ **atunci** $b \leftarrow a_i$, $a_i \leftarrow a_j$, $a_j \leftarrow b$. Se poate reduce numărul de atribuiri ?

5.4 Sortare

Vom considera numai sortarea *internă*, în care datele se presupun stocate în memoria internă a calculatorului; spre deosebire, în sortarea externă datele se află într-o memorie mult mai lentă, externă calculatorului (disc, etc.).

Complexitatea secvențială a sortării este $\Omega(n \log n)$ și există mai multe metode care ating această limită: sortarea rapidă, pe care o vom numi *quicksort* așa cum e deja consacrat, sortarea prin interclasare (mergesort) sau sortarea de ansamble (heapsort). Spre deosebire de ultimele două, care au aceeași complexitate (ca ordin de mărime) în toate cazurile, quicksort se comportă excelent în medie, dar poate ajunge la $O(n^2)$ în cel mai defavorabil caz; totuși, cu unele mici precauții explicate în secțiunea dedicată ei, sortarea rapidă este într-adevăr cea mai rapidă și deci cea mai răspândită, fiind pe deasupra și relativ ușor de programat. În cele ce urmează vom considera pur și simplu $n \log n$ complexitatea sortării.

Vom prezenta algoritmi secvențiali la începutul fiecărei secțiuni următoare, ei fiind în majoritatea cazurilor sursa de inspirație, mai mult sau mai puțin directă, pentru algoritmi paraleli.

Atunci când spunem (pentru calculatoarele cu memorie distribuită, de obicei) că un procesor ordonează o secvență, vom presupune tacit că o face utilizând cel mai bun algoritm secvențial cu putință; acesta depinde de mulți factori, dar în primul rând de lungimea secvenței de sortat.

În cazul arhitecturilor secvențiale este ușor de spus unde se va găsi secvența A după sortare: în exact aceeași zonă de memorie ca la început, ordinea crescătoare fiind de la stânga la dreapta. Pentru arhitecturile paralele cu memorie comună situația este aceeași. În cazul memoriei distribuite, însă, ordinea finală nu mai este atât de evidentă; fiecare procesor deține o parte din secvență și nu e obligatoriu ca aceste părți să aibă dimensiuni egale, deși ar fi de preferat. Vom considera, în mod natural, că A este sortată dacă secvența A_i locală procesorului P_i este ordonată crescător, $0 \leq i \leq p - 1$, și dacă $A_i < A_j$, când $i < j$, $0 \leq i, j \leq p - 1$.

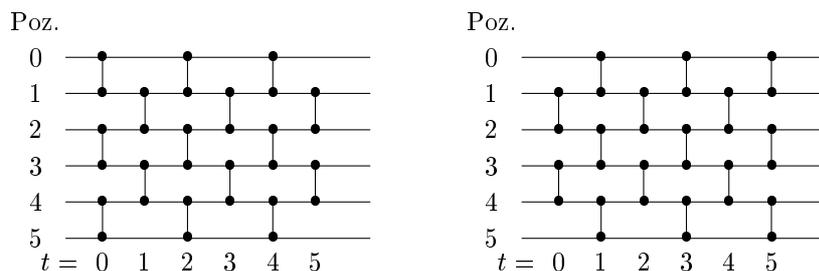


Figura 5.5: Rețele de sortare par-impair.

5.4.1 Sortare par-impair

Există o mulțime de algoritmi secvențiali de sortare în timp $O(n^2)$ (vezi, de exemplu, problema 5.4.1); ei sunt utili din două motive: sunt extrem de simplu de programat și sunt cei mai eficienți pentru număr mic de elemente. În continuare vom prezenta un astfel de algoritm, numit sortare *par-impair* din motive care vor fi imediat evidente.

Rețele de comparatoare. Introducem întâi noțiunea de rețea de comparatoare, cu ajutorul figurii 5.5; liniile orizontale reprezintă elementele din secvența de intrare A , iar segmentele verticale unind două orizontale reprezintă comparatoarele, care efectuează operații $exch(i, j)$, cu $i < j$ (i și j sunt numerele de ordine ale orizontalelor); așadar, circulând de la stânga la dreapta, elementele mai mici au tendința de a "urca", cele mai mari de a "cobori". Rețeaua se numește *de sortare* dacă la ieșirea ei secvența A este ordonată, oricare ar fi fost ordinea inițială. Puteți vedea o rețea ca o structură hardware, în care pe orizontale circulă datele, iar comparatoarele au două intrări și două ieșiri, sau pur și simplu ca pe o descriere grafică a unui algoritm, cu o bună scoatere în evidență a paralelismului.

Pentru a păstra tradiția, introducem câteva notații specifice. Notăm $[i : j]$ comparatorul care execută $exch(i, j)$; dacă α este o rețea, atunci $\alpha[i : j]$ este o rețea formată din concatenarea rețelei α cu comparatorul $[i : j]$ (acesta se adaugă la dreapta lui α); dacă X este o secvență de intrare, atunci $X\alpha$ este secvența la ieșirea din α , iar $(x\alpha)_i$ elementul din poziția i al acesteia.

Rețele de sortare par-impair. Figura 5.5 prezintă două astfel de rețele; se observă imediat ca la fiecare pas (sau moment de timp) se execută maximum de comparații posibile, alternând paritatea poziției elementului mai mic (în urma comparației). Rezultatul se obține în timp n , cu gradul de paralelism $n/2$. Se utilizează $\binom{n}{2} = n(n-1)/2$ comparatoare; altfel spus, costul algoritmului este $O(n^2)$. Se observă că rețeaua din dreapta, în figura 5.5 este cea din stânga văzută în sens invers; pentru n par, rețelele sunt diferite, pentru n impar identice. Deși e intuitiv că aceste rețele sortează, demonstrația nu e tocmai simplă; cititorul grăbit poate sări până la sfârșitul ei.

Proprietăți ale rețelelor de sortare. O rețea în care toate comparatoarele

sunt de forma $[i : i + 1]$ (compară elemente adiacente) se numește *simplă*; astfel sunt și rețelele de sortare par-impair. Multe din proprietățile demonstrate în continuare sunt valabile *numai* pentru rețele simple.

Într-o secvență A , numim *inversare* perechea (a_i, a_j) dacă $a_i > a_j$ și $i < j$. O secvență poate avea cel mult $\binom{n}{2}$ inversări (dacă e ordonată descrescător).

PROPOZIȚIA 5.1 *O rețea simplă de sortare are cel puțin $\binom{n}{2}$ comparatoare.*

Demonstrație. Un comparator $[i : i + 1]$ reduce numărul de inversări din secvență cu 0 sau 1. Secvența sortată nu are inversări; maximul de inversări este $\binom{n}{2}$. ■

Deci, din acest punct de vedere, rețeaua par-impair stă bine: conține numărul minim de comparatoare.

Notăm cu Z o secvență care conține toate inversările posibile: pentru orice $i < j$, avem $z_i > z_j$; sau, dacă vrei, putem considera $Z = \langle n, n - 1, \dots, 1 \rangle$.

PROPOZIȚIA 5.2 (Floyd) *Dacă o rețea α nu sortează o secvență oarecare X pentru două poziții $i < j$, adică $(x\alpha)_i > (x\alpha)_j$, atunci nu sortează nici Z pentru acele poziții, deci $(z\alpha)_i > (z\alpha)_j$.*

Demonstrație. Prin inducție după numărul de comparatoare din α . Dacă α are un singur comparator $[r : r + 1]$ și presupunem că $(z\alpha)_i < (z\alpha)_j$, atunci $i = r$, $j = r + 1$ și deci $(x\alpha)_i < (x\alpha)_j$, ceea ce contrazice ipoteza.

Dacă β este o rețea pentru care lema este adevărată, $\alpha = \beta[r : r + 1]$ și $(x\alpha)_i > (x\alpha)_j$, să demonstrăm că $(z\alpha)_i > (z\alpha)_j$.

Dacă $r = i$, atunci $j > r + 1$ ($j = r + 1$ implică $(x\alpha)_i < (z\alpha)_j$). Cum $(x\alpha)_r < (x\alpha)_{r+1}$ și $(x\alpha)_r > (x\alpha)_j$, atunci $(x\alpha)_{r+1} > (x\alpha)_j$, deci $(x\beta)_{r+1} > (x\beta)_j$ și $(x\beta)_r > (x\beta)_j$; din ipoteza de inducție rezultă că $(z\beta)_{r+1} > (z\beta)_j$ și $(z\beta)_r > (z\beta)_j$ și în final $(z\alpha)_i = \min((z\beta)_r, (z\beta)_{r+1}) > (z\alpha)_j$.

Dacă $r = i - 1$, atunci $(x\beta)_r$ sau $(x\beta)_{r+1}$ mai mare decât $(x\beta)_j$ (sau amândouă), deci $(z\beta)_r$ sau $(z\beta)_{r+1}$ mai mare decât $(z\beta)_j$ (din ipoteza de inducție) și atunci $(z\alpha)_i = \max((z\beta)_r, (z\beta)_{r+1}) > (z\alpha)_j$.

Lăsăm cititorului completarea demonstrației pentru $r = j$ și $r = j - 1$. Pentru alți i, j e banal, deoarece $[r : r + 1]$ nu afectează aceste poziții. ■

PROPOZIȚIA 5.3 *Dacă o rețea simplă sortează Z , atunci sortează orice secvență.*

Demonstrație. Consecință directă a Prop. 5.2. ■

Pentru a demonstra că rețeaua de sortare par-impair funcționează corect este suficient să demonstrăm următoarea afirmație.

PROPOZIȚIA 5.4 *Rețeaua par-impair sortează Z .*

Demonstrație. Cum rețeaua conține numărul minim de comparatoare, înseamnă că la fiecare comparație urmează interschimbarea, atunci când Z parcurge rețeaua. Să presupunem că se întâmplă așa. Ne vom referi la rețeaua din figura 5.5, dreapta; pentru cea din stânga, vezi P5.4.4.

Considerăm că efectul unei comparații la momentul t se manifestă la momentul $t + 1$; deci, la $t = 0$, z_0 e pe poziția 0, z_1 pe poziția 1; la momentul 1, z_1 e pe poziția 0, z_0 pe poziția 1. Prima comparație are loc la momentul $t = 0$, ultima la $t = n - 1$; poziția finală e cea de la momentul $t = n$.

Să încercăm să vedem care este poziția elementului z_k în momentul t dacă toate comparatoarele fac interschimbări.

$$\text{Pentru } k \text{ impar : } \begin{cases} k - t, & \text{pentru } t \leq k \\ t - 1 - k, & \text{pentru } t \geq k + 1 \end{cases}$$

z_k urcă de la $t = 1$ la $t = k$, ajungând în poziția 0; acolo stă un pas; apoi coboară de la $t = k + 2$ la $t = n$, ajungând pe poziția $n - 1 - k$, adică cea bună.

$$\text{Pentru } k \text{ par : } \begin{cases} k + t, & \text{pentru } t \leq n - 1 - k \\ 2n - 1 - k - t, & \text{pentru } t \geq n - k \end{cases}$$

z_k coboară de la $t = 1$ la $t = n - 1 - k$, ajungând în ultima poziție; stă acolo un pas; apoi urcă de la $t = n - k + 1$ la $t = n$, ajungând pe poziția $n - 1 - k$.

Să demonstrăm prin inducție după t că această descriere a poziției funcție de timp este corectă. Vom face doar un sfert din demonstrație, celelalte cazuri fiind asemănătoare ca tehnică.

Pentru $t=0$, funcția este evident corectă.

Presupunem că este corectă la momentul t și vom arăta că, în acest moment, au loc toate interschimbările, ceea ce face funcția corectă și la momentul $t + 1$. Presupunem k impar și t par. Atunci:

a) Dacă z_k se află în poziția $k - t$ (impară, $k < t$), comparația se face cu elementul din poziția $k - t - 1$, fie el z_j . Pentru ca interschimbarea să aibă loc este necesar ca $j < k$. Sunt posibile următoarele situații:

Pentru j impar: $k - t - 1 = j - t \Rightarrow j = k - 1$, imposibil (j ar fi par)

$$k - t - 1 = t - 1 - j \Rightarrow j = 2t - k < k$$

Pentru j par: $k - t - 1 = j + t \Rightarrow j = k - 2t - 1 < k$

$$k - t - 1 = 2n - 1 - j - t \Rightarrow j = 2n - k, \text{ imposibil}$$

b) Dacă z_k se află pe poziția $t - 1 - k$ (pară, $t \geq k + 1$), comparația se face cu elementul din poziția $t - k$, fie el z_j . Interschimbarea are loc dacă $j > k$.

Pentru j impar: $t - k = j - t \Rightarrow j = 2t - k > k$

$$t - k = t - 1 - j \Rightarrow j = k - 1, \text{ imposibil}$$

Pentru j par: $t - k = j + t \Rightarrow j = -k$, imposibil

$$t - k = 2n - 1 - j - t \Rightarrow j = 2n - 2t + k - 1 > k$$

Mai trebuie verificat că, din cele patru variante bune, exact una are loc la fiecare moment de timp (ceea ce e simplu) și că, la momente de timp diferite, nu se obțin

soluții j identice (ceea ce iarăși e simplu, dar trebuie coroborat cu cazul t impar). Din nou lăsăm cititorului aceste detalii⁵. ■

Aici se termină demonstrația corectitudinii rețelelor de sortare par-impar. Extensia acestui algoritim la arhitecturi MIMD cu memorie partajată este lăsată ca exercițiu (vezi P5.4.2).

Sortare par-impar pe MIMD cu memorie distribuită

Ne ocupăm direct de cazul $p \ll n$. Fiecare dintre cele p procesoare are în memoria locală $m = n/p$ elemente din A .

Aplicarea principiului lui Brent. Am putea aplica ideea din figura 5.5 fără nici o modificare; închipuiți-vă că $p = 2$ și că tragem o linie orizontală care separă trei poziții sus pentru procesorul P_0 și trei poziții jos pentru procesorul P_1 ; fiecare procesor execută comparațiile și interschimbările care îi sunt locale, pentru cele în care un element e la P_0 , iar celălalt la P_1 , efectuându-se comunicație (în cadrul unui $exch()$); se comunică de $n/2$ ori, de fiecare dată numai câte un element (deci vor apărea foarte mulți σ în timpul de comunicație); în plus, timpul de execuție rămâne $O(n^2)$. Nu am făcut altceva decât să aplicăm principiul lui Brent.

Sortare par-impar cu subsecvențe. O idee mult mai bună e aceea de a aplica sortarea par-impar la nivel de subsecvențe; ca imagine grafică, în figura 5.5, pe liniile orizontale nu ar mai circula elemente ci subsecvențe din A . Fiecare procesor se ocupă de o astfel de "linie". Un "comparator" va trebui acum, dintre cele $2m$ elemente de pe cele două linii, să le aleagă pe cele mai mici m și să le dirijeze pe prima linie, și pe cele mai mari m , pe cealaltă linie. Dacă subsecvențele sunt ordonate, atunci aceasta este, într-o arhitectură cu memorie distribuită, chiar o operație $exch(A_i, A_{i+1})$, așa cum a fost definită în secțiunea 5.3.5. Nu este greu de văzut că, dacă un algoritim sortează utilizând $exch()$ la nivel de element, el va sorta și cu $exch()$ la nivel de subsecvență, cu singura condiție ca subsecvențele să fi fost inițial sortate; implementarea $exch()$ din algoritmul 5.17 ne asigură că subsecvențele rămân permanent sortate, iar faptul că algoritmul sortează ne spune că, la sfârșitul execuției sale, subsecvențele sunt în ordine crescătoare; deci, întreaga secvență este sortată.

Deoarece în sortarea par-impar comparațiile sunt adiacente, topologia suficientă în cazul implementării pe o arhitectură cu memorie distribuită este inelul. Scrierea algoritmului necesită doar un pic de atenție la parități și la cazurile când procesorul este la una din extreme.

ALGORITM 5.19 (sortare par-impar pe inel, $p \ll n$, $id \equiv k$)

1. sortează crescător secvența locală (cu quicksort)
2. **pentru** $i = 0 : p - 1$
 1. **dacă** i este par **atunci**
 1. **dacă** k este par și $k \neq p - 1$ **atunci** $exch(k, k + 1)$

⁵Important e să fie bine înțeles drumul unui element prin rețea.

2. **dacă** k este impar **atunci** $exch(k-1, k)$
2. **altfel** $\{i$ este impar $\}$
 1. **dacă** k este par și $k \neq 0$ **atunci** $exch(k, k-1)$
 2. **dacă** k este impar și $k \neq p-1$ **atunci** $exch(k, k+1)$

În acest caz, numărul de comparații este (primul termen pentru sortarea locală, al doilea pentru cele p interclasări):

$$T_a(n, p) = \frac{n}{p} \log \frac{n}{p} + p \frac{n}{p} = \frac{n}{p} \log \frac{n}{p} + n.$$

Timpul de comunicație este $T_c(n, p) = p(\sigma + (n/p)\beta) = O(n)$. Acești timpi nu mai sunt atât de răi ca aceia rezultați din aplicarea principiului lui Brent; cauza este că metoda de sortare locală a fost cea mai bună posibil, și nu tot cea par-impar. Astfel, eficiența asimptotică a algoritmului este totuși 1, deși algoritmul secvențial de origine este în $O(n^2)$ operații.

Generalizare. Nu am fi insistat atât de mult dacă situația nu ar fi tipică; e prima dată când ne întâlnim cu așa ceva, deci merită să generalizăm. Să zicem că rezolvarea unei anume probleme, folosind un algoritm secvențial rapid \mathcal{A}_1 , necesită mult mai puțin timp decât dacă se folosește un alt algoritm secvențial \mathcal{A}_2 . Totuși, \mathcal{A}_2 se paralelizează foarte ușor; în plus are proprietatea că, într-o arhitectură multiprocesor, înainte sau/și după o fază de cooperare între procesoare, reduce dimensiunea problemei de la n , la n/p , în așa fel încât fiecare procesor se vede pus în situația de a rezolva secvențial problema sa. Este clar că el o va rezolva folosind cel mai bun algoritm secvențial cu putință, adică \mathcal{A}_1 . Dacă n/p este suficient de mare și \mathcal{A}_2 nu consumă prea mult timp pentru a reduce dimensiunea problemei, atunci algoritmul paralel inspirat de \mathcal{A}_2 are o complexitate de același ordin de mărime ca o ipotetică paralelizare ideală a algoritmului \mathcal{A}_1 , pentru că efortul cel mai mare este rezolvarea *locală* a problemei! Aparent am făcut dintr-un slab algoritm secvențial, un bun algoritm paralel. În fapt am folosit fiecare algoritm acolo unde era mai bun; din \mathcal{A}_2 am luat paralelismul, din \mathcal{A}_1 rapiditatea secvențială. Totuși, prin tradiție și, de ce să nu recunoaștem, pe merit, numele algoritmului paralel este dat de \mathcal{A}_2 .

În cazul nostru, \mathcal{A}_1 este quicksort și are o complexitate de $O(n \log n)$; \mathcal{A}_2 este sortarea par-impar, are complexitate $O(n^2)$, dar combină p sortări de dimensiune n/p în timp $O(n)$ pentru a face o sortare de dimensiune n . Hibridul lor are o complexitate paralelă de $O(n \log n)$.

Nu trebuie însă să rămâneți cu ideea că quicksort nu se pretează paralelizării. Va fi demonstrat imediat contrariul. Totuși, algoritmul 5.19 se distinge prin simplitate, chiar dacă eficiența lui este ceva mai slabă, cum se va dovedi după studierea altor metode de sortare paralelă.

Rețele sistolice. În final, o ultimă observație privitoare la rețelele de sortare în general, ca structuri hardware. Putem să ne închipuim că există un ceas global, la fiecare tact, comparatoarele citind datele de intrare, efectuând interschimbarea,

dacă e cazul, și scoțând datele pe liniile de ieșire, pentru a fi citite de următoarele comparatoare, la tactul următor; există o mică dificultate în realizarea sincroniei, anume că unele date nu sunt utilizate la tactul următor, ci mai târziu (în rețelele par-impair, așa ceva se întâmplă doar pentru p impar); se pot însă introduce module de întârziere cu un tact astfel încât, pe fiecare "linie" orizontală (a unui element din secvență) dintr-o rețea de sortare, la fiecare moment de timp (deci pe fiecare verticală), să se afle fie un comparator, fie un modul de întârziere. În acest datele vor circula "aliniat" pe verticală (frontul de undă va fi vertical).

Avantajul acestui mod de a privi lucrurile este că se pot sorta mai multe secvențe în regim pipeline. Dacă se introduce la fiecare tact câte o secvență la intrarea în rețea, nu va exista nici un conflict pe parcursul străbaterii rețelei, decalajul de un tact păstrându-se până la ieșire. Deci, dacă sunt suficiente date, o rețea poate realiza o sortare la fiecare tact, eficiența fiind foarte mare. Astfel se poate realiza o funcționare numită *sistolică*, deoarece datele înaintează prin rețea ca fluxul de sânge prin vase, câte puțin la fiecare bătaie a inimii—ceasul global al rețelei. (Termenul de sistolic este mult ulterior rețelelor de sortare; noutatea nu constă însă decât în ideea funcționării sincrone, în ritmul impus de un ceas global.)

5.4.2 Quicksort paralel

Algoritmul secvențial. Ideea sortării rapide e simplă și ingenioasă: se împarte secvența A în două părți; într-una elementele sunt mai mici decât o valoare pivot v , în celalaltă, mai mari; sortându-se aceste două secvențe, A va rezulta și ea sortată. Forma generală a algoritmului este cea care urmează.

ALGORITM 5.20 (quicksort, forma generală secvențială)

procedură *quicksort*(s, d)

1. **cât timp** $s < d$ {secvența este nevidă}
 1. *alege_pivot*(s, d)
 2. *partitionare*(s, d, i_v) {formează cele două subsecvențe}
 3. *quicksort*($s, i_v - 1$) {și le sortează}
 4. *quicksort*($i_v + 1, d$)

Variabilele s și d sunt indicii elementului cel mai din stânga, în secvență, respectiv al celui mai din dreapta. Despre alegerea pivotului nu spunem nimic nou față de ceea ce am discutat în secțiunea 5.2; reamintim că pivotul se va găsi în poziția s , adică $v \equiv a_s$. De asemenea, știm acum cum se realizează secvențial partiționarea, din **P5.2.4**, și chiar paralel, din algoritmul 5.8 (în care ocupă cea mai mare parte); i_v este poziția pivotului în secvența A . După realizarea partiționării, sortarea subsecvențelor rezultate poate decurge independent. Apelul se face cu *quicksort*($0, n - 1$).

Algoritm paralel simplu. O paralelizare naivă a acestui algoritm ar fi următoarea: un procesor partiționează secvența inițială A ; apoi i se mai asociază un procesor, pentru a partiționa cele două secvențe rezultate; se mai adaugă două procesoare, fiind

acum patru secvențe, etc., până când se obțin p secvențe și atunci fiecare procesor sortează secvențial (cu quicksort) o secvență. Cum pentru partiționarea de către un procesor a unei secvențe de lungime n sunt necesare n comparații, și presupunând că se obțin secvențe egale la fiecare partiționare, numărul de comparații necesar va fi

$$n + n/2 + n/4 + \dots + 2n/p + n/p \log(n/p) \approx n/p \log n + 2n$$

ceea ce nu e convenabil; gândiți-vă numai la cazul $p = \log n$, care e ușor de întâlnit în practică, și pentru care eficiența este $1/3$. De fapt va fi și mai rău, pentru că secvențele pot fi inegale, și atunci unele procesoare vor lucra mai mult.

Quicksort pe PRAM

Quicksort bazat pe partiționare paralelă. Defectul ideii de mai sus este că nu se efectuează partiționarea în paralel, de către mai multe procesoare. Pe o arhitectură cu memorie comună, o structura generală a unui algoritm quicksort paralel care să înlăture acest neajuns este foarte asemănătoare cu cea a algoritmului secvențial.

ALGORITHM 5.21 (quicksort, forma generală paralelă pe un PRAM)

procedură $pquicksort(s, d, q)$

1. **dacă** $q = 1$ **atunci**

1. sortează secvențial folosind algoritmul quicksort 5.20

2. **altfel**

1. **cât timp** $s < d$ {secvența este nevidă}

1. $alege_pivot(s, d, q)$

2. $partiționare(s, d, i_v, q)$ {formează cele două subsecvențe}

3. $q' \leftarrow \lceil q(i_v - s)/(d - s) \rceil$ {număr de proc. pentru secvența din stânga}

4. $pquicksort(s, i_v - 1, q')$ {sortează în stânga cu q' procesoare}

5. $pquicksort(i_v + 1, d, q - q')$ {sortează în dreapta cu $q - q'$ procesoare}

Față de forma secvențială a mai fost adăugat un parametru, numărul de procesoare q ; inițial, avem $q = p$; apoi, după crearea celor două subsecvențe, procesoarele se împart și ele în două grupe, de mărime proporțională cu lungimea secvențelor; de aici mai rezultă un avantaj al acestei abordări, anume că se echilibrează mai bine eforturile procesoarelor, cel puțin câtă vreme q nu devine foarte mic.

Partiționare paralelă. Întrucât alegerea pivotului e o operație care durează puțin și poate fi făcută în multe feluri, trecem direct la punctul cel mai important: partiționarea; de fapt implementarea sortării rapide în paralel se reduce la un bun algoritm paralel de partiționare, o dată adoptată forma 5.21. Cum am menționat, algoritmul 5.8 conține în mare măsură ideile de care avem nevoie; acolo se extrăgea doar una dintre cele două secvențe rezultate; aici avem nevoie de amândouă. Pe scheletul aceluiași algoritm vom obține (procesoarele participante sunt P_0, \dots, P_{q-1} , după o eventuală renumerotare; secvența de care se ocupă un procesor este notată A ; pivotul este $v \equiv a_0$):

ALGORITM 5.22 (algoritm de partiționare pe CREW PRAM, $p \ll n$, pentru procesorul P_k aparținând unui grup de q procesoare)

procedură *partiționare*(s, d, i_v, q)

1. $m \leftarrow \lceil (d - s)/q \rceil$
2. \tilde{A} se află între pozițiile $l = s + km$ și $u = \min[s + (k + 1)m - 1, d]$ din A
3. aplică algoritmul din **P5.2.4** pt. partiționarea $\tilde{A} = \langle \tilde{A}_0, \tilde{A}_1 \rangle$, $\tilde{A}_0 < v \leq \tilde{A}_1$
- 3'. (în plus, P_0 pune v pe prima poziție în \tilde{A}_1)
4. $c_k \leftarrow |\tilde{A}_0|$, numărul de elemente din \tilde{A} mai mici decât v
5. $f_k \leftarrow |\tilde{A}_1|$, numărul de elemente din \tilde{A} mai mari decât v
6. $c_{k+1} \leftarrow \sum_{i=0}^k c_i$ și $f_{k+1} \leftarrow \sum_{i=0}^k f_i$, cu algoritmul 4.8 (suma prefixelor)
7. copiază \tilde{A}_0 în B , în pozițiile de la $s + c_k$ la dreapta ($c_0 = 0$)
8. copiază \tilde{A}_1 în B , în pozițiile de la $s + c_q + f_k$ la dreapta ($f_0 = 0$)
9. copiază în \tilde{A} elementele corespunzătoare din B (de la l la u)
10. **dacă** $k = 0$ **atunci** $i_v \leftarrow s + c_q$

În instrucțiunea 2, procesorul calculează indicii care delimitează zona de care se ocupă. Partiționarea ei și calculul pozițiilor în B ale subsecvențelor locale \tilde{A}_0 și \tilde{A}_1 se face ca în algoritmul 5.8. La copiere, intervine o diferență; în 8, secvențele de tip \tilde{A}_1 se copiază într-o zonă din B care începe cu poziția $s + c_q$; în această poziție va fi chiar pivotul (în 3', P_0 l-a pus pe prima poziție în \tilde{A}_1 , iar acum se copiază \tilde{A}_1 începând de la $s + c_q$); cum sunt $c_q = \sum_{i=0}^{q-1} c_k$ elemente mai mici decât v , aceasta este poziția corectă. În 9, fiecare procesor recopiază aproximativ m elemente din B în A , dar ele nu sunt aceleași cu cele deținute inițial. Se observă că se poate folosi porțiunea din B cuprinsă între indicii s și d . Timpul de calcul este bine repartizat între procesoare, deci se poate estima la $(d - s)/q + O(\log q)$ primul termen fiind pentru partiționarea secvențială, al doilea pentru suma prefixelor.

Să considerăm acum algoritmul 5.21 în care se utilizează partiționarea din procedura 5.22. Repartizarea procesoarelor în funcție de lungimea subsecvențelor face ca termenul $O(\log q)$ să fie întotdeauna neglijabil față de $(d - s)/q$ (pentru $p \ll n$); astfel, timpul de execuție al algoritmului va fi $O((n/p) \log n)$, de unde o eficiență asimptotic egală cu 1; al doilea termen ca importanță în timpul de execuție este cel mult $O(\log p \log n)$, dar mai mic decât la algoritmul 5.8 (acolo toate cele p procesoare lucrau pe secvențe din ce în ce mai mici, aici numărul de procesoare scade o dată cu dimensiunea secvenței). Memoria suplimentară folosită este de n locații, adică o întregă copie a lui A , ceea ce poate însemna destul de mult.

Reducerea memoriei suplimentare. O metodă de a elimina parțial acest inconvenient va fi sugerată în continuare, pentru că este instructivă; ideea apare în [6]. Să presupunem că fiecare procesor nu se ocupă de o secvență contiguă, ca în algoritmul 5.22 pasul 2, ci elementele sunt împărțite ciclic procesoarelor, ca în exemplul din tabelul 5.6. În primul grup de linii marcate cu numele procesoarelor este arătată repartizarea elementelor pe procesoare; în al doilea, ordinea locală a elementelor după ce fiecare procesor aplică o partiționare secvențială ca în **P5.2.4** (pivotul este 13);

k	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
a_k	13	1	23	25	18	9	7	16	11	6	10	0	15	20	17	5
P_0	13				18				11				15			
P_1		1				9				6				20		
P_2			23				7				10				17	
P_3				25				16				0				5
P_0	11				13				18				15			
P_1		1				9				6				20		
P_2			10				7				23				17	
P_3				5				0				16				25
a_k	11	1	10	5	13	9	7	0	18	6	23	16	15	20	17	25

Figura 5.6: Reducerea memoriei suplimentare prin repartizarea ciclică a secvenței A .

ultima linie a tabelului arată ordinea indusă pe A de această partiționare locală. Se observă că există trei zone în A : una în care elementele sunt mai mici decât pivotul, cea de la valoarea **5** în stânga; o a doua în care elementele sunt mai mari decât pivotul, la dreapta de **23**; între ele, o zonă cu elemente amestecate. Doar pentru această zonă trebuie acum făcută o reordonare, într-un mod oarecum asemănător cu partiționarea din algoritmul 5.22 (partea de după (2')); nu intrăm în amănunte tehnice și ne mulțumim de data aceasta numai cu ideea că modul de repartizare a elementelor pe procesoare poate aduce neașteptate avantaje.

Quicksort pe hipercub

Să vedem acum modul de implementare a sortării rapide pe un hipercub, în care fiecare procesor deține $m = n/p$ elemente din A ($p \ll n$). Din capul locului se vede că algoritmul 5.21 nu poate fi folosit în litera lui datorită alocării unui număr diferit de procesoare pentru sortarea celor două secvențe rezultate după partiționare; pentru o arhitectură cu memorie distribuită aceasta implică o regrupare a elementelor care necesită multă comunicație; mai grav, comunicația într-un grup de procesoare care se ocupă de aceeași secvență trebuie să fie regulată; e greu de imaginat, de exemplu, cum vor coopera 7 procesoare dintr-un hipercub cu $p = 16$.

Va trebui deci să ne mulțumim cu o înjumătățire a procesoarelor la fiecare partiționare; cele p procesoare vor coopera la partiționarea secvenței A în două, apoi câte $p/2$ procesoare vor lucra pe fiecare din secvențele rezultate, ș.a.m.d. până când fiecare procesor va trebui să sorteze secvența pe care o deține. O mai slabă echilibrare a efortului va fi compensată de facilitatea comunicației; am ales totuși răul cel mai mic. Pentru ca nici acest rău să nu fie semnificativ, e necesară ceva mai multă grijă în alegerea pivotului. Este evident care ar fi cea mai bună alegere: chiar mediana secvenței A ; astfel, partiționarea s-ar face chiar la jumătate; din păcate, calcul exact al medianei e destul de costisitor, deci ne vom mulțumi cu o aproximare a ei.

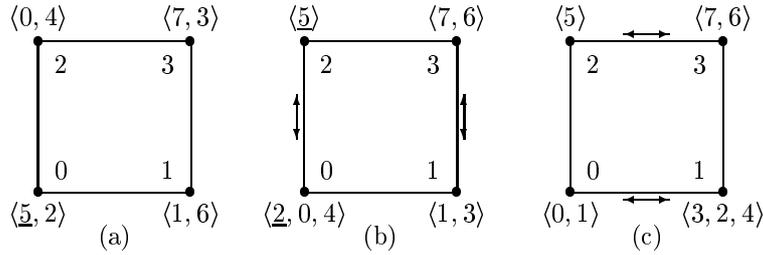


Figura 5.7: Exemplu de desfășurare a sortării rapide pe un hiper cub cu patru procesoare: (a) situația inițială; (b) după comunicația pe direcția 1; (c) după comunicația pe direcția 0.

Pe un hiper cub cu $p = 2^d$, acest mod de operare se poate pune ușor în practică. Să spunem că P_0 alege pivotul v dintre elementele locale și îl difuzează tuturor. Fiecare procesor particionează secvența proprie: $\tilde{A} = \langle \tilde{A}_0, \tilde{A}_1 \rangle$ astfel încât $\tilde{A}_0 < v \leq \tilde{A}_1$. Până aici totul e banal. Acum procesoarele vecine pe direcția $d - 1$ (care au diferit doar primul bit al reprezentării binare a adresei; pentru P_k acest bit este k_{d-1}) schimbă între ele câte o secvență, astfel încât secvențele mai "mici" \tilde{A}_0 să ajungă în jumătatea inferioară a hiper cubului (unde $k_{d-1} = 0$), iar cele mai "mari" \tilde{A}_1 în jumătatea superioară (unde $k_{d-1} = 1$); astfel procesoarele P_i , cu $0 \leq i \leq p/2 - 1$, au acum elemente mai mici decât pivotul și deci mai mici decât cele ale oricărui procesor P_j , cu $p/2 \leq j \leq p - 1$. În continuare, totul se repetă în cele două subhipercuburi ($\mathcal{H}_d^{d-1}(0)$, jumătatea inferioară și $\mathcal{H}_d^{d-1}(2^{d-1})$, jumătatea superioară), în paralel; se comunică pe direcția $d - 2$, de această dată. Se procedează similar mai departe, la fiecare iterație comunicând pe direcția imediat inferioară, până la direcția 0; de fiecare dată, numărul grupurilor de procesoare care cooperează se dublează, iar numărul procesoarelor dintr-un grup se înjumătățește. În final, sunt p grupuri formate dintr-un singur procesor; atunci secvența A este sortată la nivel mare (P_i are elementele mai mici decât cele ale lui P_j , dacă $i < j$) și fiecare procesor va avea doar de sortat secvența locală.

O ilustrare a algoritmului este prezentată în figura 5.7, pentru $p = 4$ și $n = 8$; elementele pivot sunt cele subliniate (primele din secvențele locale); ultima fază, de sortare locală, nu a mai fost reprezentată.

Această primă implementare a algoritmului quicksort are forma următoare.

ALGORITM 5.23 (quicksort pe hiper cub, $p = 2^d \ll n$, $\text{id} \equiv k$)

1. **pentru** $l = d - 1 : -1 : 0$
 1. **dacă** ultimii $l + 1$ biți din k sunt egali cu 0 **atunci**
 1. alege pivotul v dintre elementele locale
 2. **difuzează** v în $\mathcal{H}_d^{l+1}(k)$
 2. aplică algoritmul din **P5.2.4** pentru partiționarea $\tilde{A} = \langle \tilde{A}_0, \tilde{A}_1 \rangle$, $\tilde{A}_0 < v \leq \tilde{A}_1$
 3. **dacă** $k_l = 0$ **atunci** $\{P_k$ e în subhipercubul "inferior"}

1. **în paralel** $\text{send}(\tilde{A}_1, l)$, $\text{recv}(X, l)$
2. $\tilde{A} \leftarrow \langle \tilde{A}_0, X \rangle$ {concatenează secvența primită cu cea rămasă}
4. **altfel** { $k_l = 1$, subhipercubul "superior"}
 1. **în paralel** $\text{send}(\tilde{A}_0, l)$, $\text{recv}(X, l)$
 2. $\tilde{A} \leftarrow \langle \tilde{A}_1, X \rangle$
2. sortează secvența proprie \tilde{A}

Timpul de execuție poate fi estimat astfel, în ipoteza unei încărcări echilibrate (fiecare procesor are în permanență aproximativ n/p elemente):

$$T(n, p) = \frac{n}{p} \log \frac{n}{p} \alpha + \sum_{i=0}^{d-1} \left[\frac{n}{p} \alpha + \left(\sigma + \frac{n}{2p} \beta \right) \right] = \frac{n}{p} (\log n) \alpha + \left(\sigma + \frac{n}{2p} \beta \right) \log p. \quad (5.5)$$

Primul termen este pentru sortarea locală, iar în sumă, n/p este numărul de comparații necesare partiționării, iar celălalt termen corespunde comunicației; am neglijat timpul pentru difuzări, acestea făcându-se în paralel pe hipercuburi din ce în ce mai mici. Rezultatul este promițător, comparațiile fiind perfect distribuite procesoarelor, deci existând numai overhead de comunicație. În practică lucrurile nu stau chiar așa; dacă alegerea pivotului nu este făcută cu atenție, nu numai că încărcarea poate deveni foarte inegală, dar algoritmul poate nici să nu funcționeze. De exemplu, ca un caz extrem, dacă P_0 alege ca pivot cel mai mic element al său, care, întâmplător, este mai mic și decât orice element din P_{2^d-1} , după prima iterație P_0 rămâne fără nici un element, deci nu mai poate propune un pivot (ar putea propune o valoare aleatoare, dar asta nu are decât o logică: să nu se blocheze execuția; nici vorbă de performanță); vezi și problema 5.4.6.

Hyperquicksort. Un algoritm asemănător, dar cu o mult mai stabilă comportare, a fost propus de Wagar și numit hiperquicksort [26]. Să presupunem că fiecare procesor începe prin a sorta secvența proprie și să vedem ce avantaje decurg de aici. Alegerea pivotului este acum imediată: P_0 va alege elementul său median; este vorba deci de mediana unui eșantion de n/p elemente, care este o bună aproximare a medianei secvenței A . P_0 face difuzarea pivotului. Partiționarea este și ea simplă; trebuie găsită poziția pivotului în secvența locală, care este ordonată; aceasta se face prin căutare binară, deci rapid. Se comunică, exact ca în algoritmul precedent, secvențele "mici" în jumătatea inferioară a hipercubului și secvențele "mari" în jumătatea superioară. Fiecare procesor va avea acum două secvențe ordonate, pe care le interclasează, pentru a avea din nou o secvență ordonată. După care se continuă în cele două subhipercuburi, apoi în patru, etc., structura comunicației fiind aceeași ca la algoritmul dinainte. Toate aceste considerații ne conduc la:

ALGORITM 5.24 (*Wagar*) (hiperquicksort, pe hipercub, $p = 2^d \ll n$, $\text{id} \equiv k$)

1. sortează secvența locală \tilde{A}
2. **pentru** $l = d - 1 : -1 : 0$

1. **dacă** ultimii $l + 1$ biți din k sunt egali cu 0 **atunci**
 1. alege mediana secvenței locale \tilde{A} ca pivot v
 2. **difuzează** v în $\mathcal{H}_d^{l+1}(k)$
2. partiționează $\tilde{A} = \langle \tilde{A}_0, \tilde{A}_1 \rangle$, $\tilde{A}_0 < v \leq \tilde{A}_1$, prin căutare binară
3. **dacă** $k_l = 0$ **atunci** $\{P_k$ e în subhipercubul "inferior" $\}$
 1. **în paralel** $\text{send}(\tilde{A}_1, l)$, $\text{recv}(X, l)$
 2. $\tilde{A} \leftarrow \tilde{A}_0 \perp X$ $\{\text{interclasează secvența primită cu cea rămasă}\}$
4. **altfel** $\{k_l = 1$, subhipercubul "superior" $\}$
 1. **în paralel** $\text{send}(\tilde{A}_0, l)$, $\text{recv}(X, l)$
 2. $\tilde{A} \leftarrow \tilde{A}_1 \perp X$

Făcând din nou ipoteza unei încărcări echilibrate (mult mai aproape de realitate, acum), timpul de execuție este

$$T(n, p) = \frac{n}{p} \log \frac{n}{p} \alpha + \sum_{i=0}^{d-1} \left[\left(\log \frac{n}{p} + \frac{n}{p} \right) \alpha + \left(\sigma + \frac{n}{2p} \beta \right) \right] \approx \frac{n}{p} (\log n) \alpha + \left(\sigma + \frac{n}{2p} \beta \right) \log p. \quad (5.6)$$

O comparație rapidă cu relația (5.5) ne dezvăluie un singur termen în plus, în sumă, cel corespunzător căutării binare, de $\log(n/p)$, pe care-l putem și neglija, de altfel. În rest, sortarea locală este făcută tot o singură dată, comunicația este aceeași, și, coincidentă plăcută, partiționarea unei secvențe neordonate (de lungime n/p) din algoritmul 5.23 are aceeași complexitate ca interclasarea a două secvențe (de lungime aproximativ $n/(2p)$) în hiperquicksort. Astfel algoritmul hiperquicksort are aceeași complexitate, în cel mai bun caz, ca și algoritmul 5.23, dar o va atinge cu probabilitate mult mai mare decât acesta din urmă. Este și motivul pentru care a fost preferat; de asemenea, el este și relativ simplu.

5.4.3 Sortare echilibrată

Algoritmii de sortare pentru arhitecturi cu memorie distribuită inspirați de quicksort au, după cum am văzut, un dezavantaj: procesoarele nu au același număr de operații de efectuat, de unde existența unui overhead de încărcare inegală.

Vom prezenta în această secțiune un algoritm, preluat din [1] care înlătură acest dezavantaj, desigur, cu prețul unor operații suplimentare. Notăm cu A secvența de sortat, de lungime n și cu A_k secvența locală procesorului P_k , de lungime $m = n/p$. Primul lucru pe care-l face fiecare procesor este să-și ordoneze secvența locală.

Ideea: multiselecție și schimb complet. Să încercăm să facem astfel încât procesorul P_0 să aibă cele mai mici m elemente, P_1 pe următoarele m , etc. După aceea va fi suficientă o nouă sortare locală. Pentru aceasta este necesară parcurgerea următoarelor etape:

1. Se efectuează selecția elementelor din pozițiile im , cu $i = 1, \dots, p-1$, în secvența A ordonată; sunt necesare $p-1$ astfel de elemente, notate v_i , pentru a partiționa A în p secvențe.

2. Fiecare procesor împarte secvența locală în p părți, în concordanță cu cheile v_i ; i.e. P_k găsește subsecvențele A_k^i ale lui A_k astfel încât $A_k^{i-1} < v_i \leq A_k^i$.
3. Procesorul P_k colectează secvențele A_l^k de la celelalte procesoare; pe ansamblu, aceasta este o operație de schimb complet. În acest moment, P_k deține toate subsecvențele pentru care $v_{k-1} \leq A_l^k < v_k$, adică exact m elemente, cele dintre pozițiile km și $(k+1)m-1$, în secvența A ordonată (pentru rigurozitate, $v_0 = -\infty$).

Algoritm de multiselecție. Am văzut în secțiunea 5.2 un algoritm de selecție a unei valori pentru arhitecturi cu memorie distribuită. Acum însă este vorba de selecția a $p-1$ valori; desigur că aceste valori ar putea fi selectate una după alta, aplicând de $p-1$ ori succesiv algoritmul 5.9. Să vedem însă dacă nu se poate face totul în paralel.

În algoritmul 5.9, unul dintre procesoare avea rolul conducător: era cel care, în instrucțiunea 2.2, alegea unul singur dintre cei p candidați la selecție; în acest timp, restul procesoarelor erau inactice. Vom executa de $p-1$ ori în paralel algoritmul 5.9, procesorul P_k având rolul conducător în căutarea valorii v_k . Pentru o mai bună înțelegere, rescriem detaliile algoritmului 5.9 și, în cadrul său, ale algoritmului 5.5 de calcul al rangului. Notăm \tilde{A} variabila locală în care P_k stochează A_k .

ALGORITM 5.25 (*Abali*) (determină cheile de partiționare a secvenței A sortată local în p subsecvențe egale, pe MIMD cu memorie distribuită, $mp = n$, $\text{id} \equiv k$)

1. **pentru** $i = 1 : p-1$
 1. $s_i \leftarrow 0$, $d_i \leftarrow m-1$, $c_i \leftarrow \tilde{a}_{\lfloor im/p \rfloor}$ ($c_0 \leftarrow NIL$), $r_i \leftarrow -1$
2. **cât timp** $r_i \neq im$, pentru cel puțin un i , $1 \leq i \leq p-1$
 1. **schimb complet:** P_k primește candidații c_k de la celelalte procesoare
 2. P_k elimină valorile NIL , ordonează restul candidaților și alege v_k mediana lor
 3. **difuzare generală:** P_k trimite tuturor valoarea v_k
4. **pentru** $i = 1 : p-1$
 1. $l \leftarrow$ rangul lui v_i în \tilde{A}
 2. $r_i \leftarrow$ suma globală a valorilor l (cu algoritmul 4.4)
 3. **dacă** $r_i \neq im$ **atunci**
 1. **dacă** $r_i < im$ **atunci** $s_i \leftarrow$ cel mai mic indice pentru care $\tilde{a}_{s_i} > v_i$
 2. **altfel** $d_i \leftarrow$ cel mai mare indice pentru care $\tilde{a}_{d_i} < v_i$
 3. **dacă** $s_i \leq d_i$ **atunci** $c_i \leftarrow \tilde{a}_{\lfloor (s_i+d_i)/2 \rfloor}$
 4. **altfel** $c_i \leftarrow NIL$

Se observă transformări tipice în comunicația globală atunci când se simetrizează un algoritm: în instrucțiunea 2.1, colectarea s-a transformat în schimb complet, în 2.3 difuzarea în difuzare generală. La pasul 2.2 toate procesoarele lucrează acum în paralel. La ceilalți pași totul se repetă de $p-1$ ori. Să remarcăm că P_0 are o poziție specială; din cauză că nu se caută o cheie v_0 , el primește numai valori NIL , deci este practic inactiv.

Să vedem cum se modifică timpul de execuție. Ne referim tot la hipercub. Față de algoritmul 5.9, toate complexitățile se înmulțesc cu p , mai puțin la pasul 2.1, unde $O(p)$ se transformă în $O(p \log p)$, și la pasul 2.2, unde rămâne valabil $O(p)$. În total, cu aceleași considerații asupra numărului de iterații, putem aprecia timpul de lucru al algoritmului la $O(p \log^2 n)$ pentru operații aritmetice (comparații) și la $O(p \log p \log n)$ pentru comunicație. Cu toate că algoritmul 5.9 necesită $O(\log^2 n)$ operații, se poate spune că algoritmul 5.25, deși de complexitate $O(p \log^2 n)$, nu implică de p ori mai multe operații, ci ceva mai puține.

Algoritmul de sortare echilibrată este următorul.

ALGORITM 5.26 (sortare echilibrată, pe MIMD cu memorie distribuită, $mp = n$, $\text{id} \equiv k$)

1. sortează secvența locală A_k
2. aplică algoritmul 5.25 pentru găsirea cheilor de partiționare v_1, \dots, v_{p-1}
3. separă A_k în secvențe A_k^i , cu $0 \leq i \leq p-1$, a.î. $A_k^{i-1} < v_i \leq A_k^i$
4. **schimb complet:** P_k colectează A_l^k de la celelalte procesoare P_l ($0 \leq l \leq p-1$)
5. ordonează cele p secvențe A_l^k ($0 \leq l \leq p-1$),
prin interclasare după un arbore binar

La pasul 3 separarea secvenței locale în p părți se poate face fie prin $p-1$ căutări binare ale pozițiilor cheilor v_1, \dots, v_{p-1} , în timp $O(p \log(n/p))$, fie printr-o simplă parcurgere secvențială a secvenței locale (ca la interclasare, secvența $\langle v_i \rangle$ fiind și ea ordonată) în timp $O(n/p)$. La pasul 4, considerând că lungimea medie a mesajelor este n/p^2 , timpul necesar schimbului complet va fi $O((n/p) \log p)$.

La pasul 5, fiecare procesor P_k deține elementele din pozițiile de la km la $(k+1)m-1$ în secvența A sortată, deci exact cele care îi trebuiau, dar sub forma a p secvențe, fiecare ordonată crescător. Pentru a profita de această structură, se poate face o interclasare după un arbore binar: cele p secvențe se interclasează câte două, obținându-se $p/2$ secvențe care, la rândul lor, se interclasează câte două, etc. (este un algoritm asemănător sortării prin interclasare, doar că acolo secvențele sunt de lungimi egale). La fiecare înjumătățire a numărului de secvențe, deoarece numărul total de elemente este n/p , se fac cel mult n/p comparații; deci, în total, $O((n/p) \log p)$ comparații, semnificativ mai puțin decât pentru o sortare obișnuită.

Sumând, obținem pentru algoritmul 5.26 de sortare echilibrată un număr de comparații de $O((n/p) \log n + p \log^2 n)$ și un timp de comunicație de $O((n/p) \log p)$ (dacă $(n/p) \log p > p \log p \log n$, adică $p^2 < n/\log n$, ceea ce poate fi rezonabil); cea mai mare pondere o au sortarea locală inițială și găsirea cheilor de partiționare. Timpul de execuție al acestui algoritm depinde în destul de slabă măsură de ordinea din secvența A . Comparații experimentale între el și hiperquicksort au arătat că, pentru o repartiție uniformă a elementelor din A , hiperquicksort se comportă mult mai bine; pentru diverse ordini particulare, cum ar fi de pildă cea în care $A_i < A_j$ pentru $i < j$, dar A_i nesortate, sortarea echilibrată are o execuție mai rapidă.

Sortare cvasi-echilibrată. Trecem acum la un algoritm care, fără a realiza o încărcare perfect echilibrată a procesoarelor, furnizează mijloace pentru ca dezechilibrul să nu fie prea mare; el a fost publicat în [24].

Algoritmul 5.26 căuta cheile v_1, \dots, v_{p-1} aflate pe pozițiile $m, \dots, (p-1)m$ în secvența sortată A . În locul determinării precise a acestora se pot căuta chei care să se afle aproape de pozițiile dorite. Ideea prezentată în continuare este foarte intuitivă; vom arăta că ea conduce la rezultate promițătoare.

Descriere generală. Ca de obicei, fiecare procesor ordonează secvența locală. Apoi, din fiecare din aceste p secvențe sunt extrase $p-1$ eşantioane; în mod natural, sunt alese cele din pozițiile $[im/p]$, cu $i = 1, \dots, p-1$; dacă vă amintiți, acestea sunt chiar valorile propuse de procesoare ca prim set de candidați în pasul (1) al algoritmului 5.25. Se obțin în total $p(p-1)$ valori, într-o secvență pe care o vom boteza Y . Se sortează Y și se aleg cele $p-1$ chei de partiționare, iarăși destul de natural, ca fiind cele din pozițiile $i(p-1)$, cu $1 \leq i \leq p-1$ (se mai pot alege, de exemplu, cele din pozițiile $p/2, p/2 + p, \dots, p/2 + (p-2)p$, adică mediana celor mai mici p chei, mediana următoarelor p , etc.). După care totul decurge ca în algoritmul 5.26: fiecare procesor partiționează secvența proprie în acord cu v_1, \dots, v_{p-1} , urmează un schimb complet în care un procesor colectează subsecvențe cu aceeași poziție în secvențele locale și le sortează apoi prin interclasare după un arbore binar.

Proprietăți. Să vedem acum cât de echilibrată este încărcarea procesoarelor; aceasta este dictată de numărul de elemente care ajung în final la un anumit procesor.

PROPOZIȚIA 5.5 *In secvența A există cel puțin $[i(p-1) + 1]m/p$ elemente mai mici decât v_i și cel puțin $(p-i)(p-1)m/p$ elemente mai mari decât v_i .*

Demonstrație. v_i este elementul din poziția $i(p-1)$ în secvența eşantioanelor Y ; așadar sunt $i(p-1)$ elemente în Y mai mici decât v_i și $p(p-1) - i(p-1) - 1 = (p-i)(p-1) - 1$ elemente mai mari. Dar fiecare element din Y poate fi privit ca reprezentantul a m/p elemente mai mici decât el; în plus, există m/p elemente mai mici decât v_i în secvența locală din care provine acesta, al căror reprezentant este v_i ; cu totul, $[i(p-1) + 1]m/p$ elemente mai mici decât v_i . La fel de bine, un element din Y este reprezentantul a m/p elemente mai mari decât el, de unde $(p-i)(p-1)m/p$ elemente mai mari decât v_i . ■

PROPOZIȚIA 5.6 *Numărul de elemente x din A pentru care $v_{i-1} < x \leq v_i$ este de cel mult $2m$.*

Demonstrație. Pentru $1 \leq i \leq p-1$, numărul de elemente mai mari decât v_i este cel puțin $(p-i)(p-1)m/p$, iar numărul de elemente mai mici decât v_{i-1} este cel puțin $[(i-1)(p-1) + 1]m/p$; numărul de elemente a căror valoare este între v_{i-1} și v_i va fi cel mult $n - [(p-i)(p-1) + (i-1)(p-1) + 1]m/p = 2m - 2m/p < 2m$. Se poate verifica imediat că propoziția este adevărată și în cazurile $i = 0$ și $i = p$. ■

PROPOZIȚIA 5.7 *Propoziția 5.6 rămâne valabilă și dacă $v_i = y_{p/2+(i-1)p}$, cu $1 \leq i \leq p-1$.*

Cum elementele $x \in A$ pentru care $v_{i-1} < x \leq v_i$ au destinația finală P_{i-1} , din propoziția 5.6 rezultă că un procesor va deține mai puțin de $2n/p$ elemente; deci, în cel mai rău caz, de două ori mai multe decât a avut inițial.

Multiselecție aproximativă. Algoritmul de sortare va fi deci identic cu 5.26, cu singura mică (dar importantă) diferență că, la pasul 2, găsirea cheilor de partiționare se face cu un alt algoritm. Pe o arhitectură cu memorie distribuită, detalierea acestui pas ar putea fi următoarea.

ALGORITHM 5.27 (*Shi*) (determină cheile de partiționare a secvenței A sortată local în p subsecvențe de lungime $\leq 2m$, pe MIMD cu memorie distribuită, $mp = n$, $\text{id} \equiv k$, \tilde{A} secvența locală)

1. formează secvența $X = \langle \tilde{a}_{\lfloor im/p \rfloor} \mid i = 1, \dots, p-1 \rangle$
2. **colectare:** P_0 colectează secvențele X și le concatenează în secvența Y
3. P_0 sortează Y și extrage $v_i = y_{i(p-1)}$ (sau $v_i = y_{p/2+(i-1)p}$), cu $1 \leq i \leq p-1$
4. **difuzare:** P_0 trimite tuturor valorile v_i

Colectarea și difuzarea durează $O(p^2)$, respectiv $O(p \log p)$ (mesajele au lungime $p-1$), iar numărul de comparații este cel necesar sortării secvenței Y , deci $p^2 \log p^2$ dacă aplicăm un algoritm uzual de sortare, sau, dacă ținem seama de faptul că secvențele Y sunt ordonate, $p^2 \log p$ aplicând sortarea prin interclasare după un arbore binar. Se vede deci că sortarea este cea care consumă cel mai mult timp. Problema **P5.4.11** arată cum se poate îmbunătăți acest timp.

Considerând acum complexitatea algoritmului (să-l numim de sortare cvasi-echilibrată) 5.26 combinat cu algoritmul 5.27 (*Shi*), vom constata că ea nu se modifică decât la pasul 2, așa cum am arătat. La pasul 5, dacă partiționarea este aproape echilibrată, putem aprecia că se fac tot $(n/p) \log p$ comparații; în cel mai rău caz, când există secvențe de lungime aproape $2n/p$, vor fi de două ori mai multe operații (dar numai dacă există $p/2$ partiții de lungime $2n/p$ și $p/2$ de lungime 0, iar la interclasare se combină cele lungi între ele, la fiecare pas). În medie sunt $O((n/p) \log n + p \log n)$ comparații iar timpul de comunicație este $O((n/p) \log p)$.

Algoritmul de sortare cvasi-echilibrată (*Shi*) este ușor adaptabil pentru o arhitectură cu memorie comună.

5.4.4 Sortare bitonică

Algoritmul de sortare ce va fi prezentat în continuare este un exemplu tipic de algoritm paralel, care nu se inspiră în nici un fel din cele secvențiale (având, de altfel, o complexitate secvențială mai mare de $O(n \log n)$). El este datorat, ca idee, lui Batchier [2].

Definiții și proprietăți. Fie $A = \langle a_0, a_1, \dots, a_{n-1} \rangle$ secvența de ordonat. Pentru $i = 1, 2, \dots, n-2$ vom spune că a_i este un *minim local* dacă a_{i-1} și a_{i+1} sunt amândouă mai mari decât a_i și că a_i este un *maxim local* dacă a_{i-1} și a_{i+1} sunt amândouă mai mici decât a_i . Secvența A este *unimodală* dacă are cel mult un element care este minim

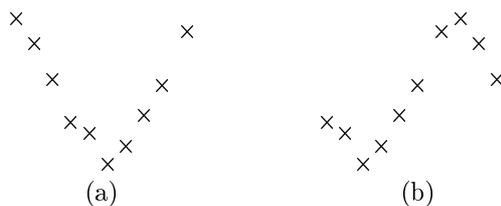


Figura 5.8: Secvențe (a) unimodală (cu minim) (b) bitonică.

local sau maxim local și *bitonică* dacă este o rotație a unei secvențe unimodale; vezi figura 5.8. Următoarea leamnă deschide calea către algoritm.

LEMA 5.1 (*Batcher*) Fie $A = \langle a_0, a_1, \dots, a_{2N-1} \rangle$ o secvență bitonică de lungime pară. Se definesc secvențele $m(A)$ și $M(A)$ astfel:

$$m(A) = \langle \min(a_0, a_N), \min(a_1, a_{N+1}), \dots, \min(a_{N-1}, a_{2N-1}) \rangle$$

$$M(A) = \langle \max(a_0, a_N), \max(a_1, a_{N+1}), \dots, \max(a_{N-1}, a_{2N-1}) \rangle.$$

Atunci secvențele $m(A)$ și $M(A)$ sunt bitonice și orice element din $m(A)$ este mai mic decât orice element din $M(A)$.

Demonstrație. Putem presupune că secvența A este unimodală și conține un minim local; într-adevăr, dacă A este rotită cu p pași la stânga (dreapta), atunci și $m(A)$, $M(A)$ vor fi rotite cu p pași la stânga (dreapta), deci își vor păstra proprietatea de bitonicitate; pe de altă parte, din orice secvență bitonică se poate obține prin rotație o secvență unimodală cu minim (vezi problema 5.4.12).

Fie a_j minimul local din A . Să presupunem că $j < N$ (dacă $j \geq N$, considerăm secvența A în ordine inversă, ceea ce inversează ordinea în $m(A)$, $M(A)$, dar nu modifică bitonicitatea sau unimodalitatea).

Fie k indicele pentru care: $a_0 > a_N, \dots, a_{k-1} > a_{N+k-1}, a_k < a_{N+k}, \dots, a_{N-1} < a_{2N-1}$. Este evident că $0 \leq k \leq j$. Atunci

$$m(A) = \langle a_N, \dots, a_{N+k-1}, a_k, a_{k+1}, \dots, a_j, a_{j+1}, \dots, a_{N-1} \rangle$$

este o rotație a secvenței $\langle a_k, a_{k+1}, \dots, a_j, a_{j+1}, \dots, a_{N-1}, a_N, \dots, a_{N+k-1} \rangle$ care este unimodală (cu a_j minim), deci $m(A)$ este bitonică.

$M(A)$ este o rotație a secvenței unimodale $\langle a_0, \dots, a_{k-1}, a_{N+k}, \dots, a_{2N-1} \rangle$ (cu a_{k-1} sau a_{N+k} minim).

Pe de alta parte, pentru orice $x \in M(A)$, $y \in m(A)$ sunt adevărate inegalitățile $x > a_{k-1}$ sau $x > a_{N+k}$, $y < a_k$ sau $y < a_{N+k-1}$; în plus, avem: $a_{k-1} > a_k$, $a_{N+k-1} < a_{N+k}$ (din unimodalitate) și $a_{k-1} > a_{N+k-1}$, $a_k < a_{N+k}$ (din ipotezele de mai sus); așadar obținem $x > y$. ■

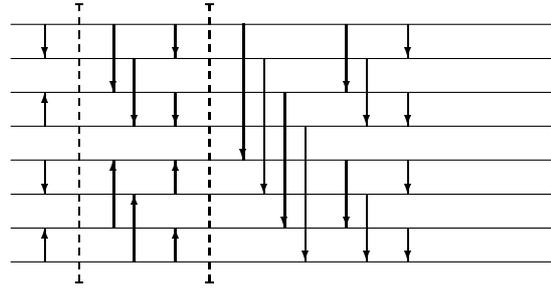


Figura 5.9: Sortare bitonică pentru 8 procesoare.

Ideea sortării bitonice. Algoritmul decurge în doi pași. Întâi se transformă A într-o secvență bitonică, apoi aceasta este sortată folosind rezultatul lemei: se formează secvențele bitonice $m(A)$ și $M(A)$, după care se sortează $m(A)$ și $M(A)$, ceea ce poate fi făcut în paralel.

Revenind la primul pas, să observăm că din concatenarea unei secvențe descrescătoare cu una crescătoare se obține o secvență bitonică. Pornind de la o secvență de n elemente, aceasta poate fi privită ca listă de $n/2$ secvențe bitonice de lungime 2, din care se poate construi o listă de $n/4$ secvențe de lungime 4, apoi o listă de $n/8$ secvențe de lungime 8, și așa mai departe până când se obține o secvență bitonică de lungime n . Pentru a transforma o listă de $n/2^i$ secvențe bitonice de lungime 2^i într-o listă de $n/2^{i+1}$ secvențe bitonice de lungime 2^{i+1} este suficient să se ordoneze secvențele de lungime 2^i alternativ descrescător și crescător, utilizând algoritmul de sortare a unei secvențe bitonice (sau un alt algoritm, pe care-l vom prezenta mai jos, mai eficient în execuția secvențială).

Rețea de sortare bitonică. În figura 5.9 sunt prezentate comparațiile care se efectuează într-o arhitectură cu 8 procesoare, sub forma unei rețele de sortare; o săgeată indică operația efectuată de un procesor: comparație urmată de o eventuală interschimbare; elementul mai mare va fi în poziția indicată de vârful săgeții. Se observă că permanent gradul de paralelism este $n/2$ (au loc $n/2$ comparații în paralel).

Să deducem complexitatea algoritmului. Dacă $B(n)$ este numărul de comparații necesar pentru sortarea unei secvențe bitonice folosind $n/2$ procesoare, atunci $B(2^i) = i$ (din lema lui Batcher). Atunci timpul de sortare a unei secvențe oarecare este

$$T(n) = \sum_{i=1}^{\log p} B(2^i) = 1 + 2 + \dots + \log p = O(\log^2 p)$$

pe un EREW PRAM (nu există conflicte nici la scriere, nici la citire). Costul algoritmului este deci $O(n \log^2 n)$; cum costul secvențial este $O(n \log n)$, se pune problema

dacă algoritmul lui Batcher poate fi îmbunătățit; răspunsul este pozitiv: se poate construi un algoritm paralel de sortare în $O(\log n)$, folosind $n/2$ procesoare, dar constanta ascunsă de $O(\cdot)$ este imensă. Practic, algoritmul lui Batcher rămâne deocamdată optim.

Sortare bitonică pe hiper cub

Cazul $p = n$. Să vedem cum se poate implementa algoritmul pe o arhitectură cu memorie distribuită, mai precis pe un hiper cub. Considerăm întâi cazul în care există un singur element pentru fiecare procesor: $p = n$. Dacă ne uităm puțin pe figura 5.9, imaginându-ne acum că fiecare linie reprezintă un procesor, observăm că un procesor k comunică doar cu procesoare $k \pm 2^i$, dar întotdeauna astfel încât într-un hiper cub au loc doar comunicații între vecini. Algoritmul este următorul:

ALGORITM 5.28 (sortare bitonică pe hiper cub, $p = n = 2^d$, $\text{id} \equiv k$)

1. **pentru** $l = 0 : d - 1$
 1. **pentru** $i = l : -1 : 0$
 1. **dacă** $k_{l+1} = 0$ **atunci**
 1. **dacă** $k_i = 0$ **atunci** $\text{exch}(k, k + 2^i)$
 2. **altfel** $\text{exch}(k - 2^i, k)$
 2. **altfel** $\{k_{l+1} = 1\}$
 1. **dacă** $k_i = 0$ **atunci** $\text{exch}(k + 2^i, k)$
 2. **altfel** $\text{exch}(k, k - 2^i)$

Singura problemă este stabilirea sensului în care se fac interschimbările; mecanismul din algoritm este ușor de înțeles dacă observăm că, la pasul l , elementul dintr-un procesor k face parte dintr-o secvență care trebuie ordonată crescător dacă bitul $l + 1$ al reprezentării binare a lui k este 0 și descrescător dacă acesta este 1 (prin convenție $k_{d+1} = 0$). Pentru un l fixat, interschimbările se fac astfel încât, pentru o pereche de procesoare P_k, P_j (cu $k = j \pm 2^i$) care face interschimbare, elementul mai mic să ajungă în procesorul cu număr mai mic dacă se ordonează crescător și în procesorul cu număr mai mare dacă se ordonează descrescător.

Complexitatea acestui algoritm este clar $\sum_{l=0}^{d-1} \sum_{i=0}^l O(1) = O(\log^2 p)$, la fiecare pas având loc o comunicație între vecini și o comparație.

Cazul $p \ll n$. Ar fi naiv să invocăm principiul lui Brent și deci (pe un PRAM) să păstrăm algoritmul neschimbat, modificând doar alocarea procesoarelor; timpul de execuție ar fi de $O((n/p) \log^2 p)$, ori trebuie să tindem către $O((n/p) \log p)$ pentru a spera o eficiență apropiată de 1.

Rămânem în cadrul arhitecturilor cu memorie distribuită. Să presupunem că fiecare procesor se ocupă de $m = n/p$ elemente. În primul rând, fiecare procesor P_k sortează elementele sale într-o secvență A_k crescătoare dacă k este par și descrescătoare dacă e impar; aceasta durează $O((n/p) \log(n/p))$. În continuare, se poate aplica algoritmul de sortare bitonică, înlocuind elementele a_i cu secvențe ordonate A_i , toate de aceeași lungime; formarea secvențelor de tipul $m(A)$ și $M(A)$ se face simplu,

$\min(A_i, A_j)$ fiind acum o secvență obținută prin alegerea fiecărui minim dintre elementele aflate pe aceeași poziție în secvențele A_i, A_j . Și totuși mai trebuie adăugat ceva.

Să exemplificăm pentru un grup de 4 procesoare; presupunem că secvența $\langle A_0, A_1 \rangle$ este ordonată crescător (adică A_0 și A_1 sunt fiecare ordonate crescător și, în plus, $A_0 < A_1$) și că secvența $\langle A_2, A_3 \rangle$ este ordonată descrescător. Deci secvența $A = \langle A_0, A_1, A_2, A_3 \rangle$ este bitonică (chiar unimodală); în algoritmul de sortare bitonică aceasta este situația după primul pas ($l = 0$ în algoritmul 5.28, vezi și figura 5.9). Urmează acum operațiile (pasul al doilea pentru cele patru procesoare):

$$\begin{array}{llll} A_0 \leftarrow \min(A_0, A_2) & A_1 \leftarrow \min(A_1, A_3) & A_2 \leftarrow \max(A_0, A_2) & A_3 \leftarrow \max(A_1, A_3) \\ A_0 \leftarrow \min(A_0, A_1) & A_1 \leftarrow \max(A_0, A_1) & A_2 \leftarrow \min(A_2, A_3) & A_3 \leftarrow \max(A_2, A_3). \end{array}$$

Lema lui Batcher ne asigură că $A_0 < A_1 < A_2 < A_3$; dacă ar fi vorba doar de elemente (adică $|A_i| = 1$) am fi satisfăcuți cu atât; problema e că elementele din nici un A_i nu sunt ordonate crescător. Deci, pentru ca secvența A să fie global ordonată, mai trebuie ordonate (încă o dată) elementele din fiecare A_i . Din fericire nu trebuie să mai pierdem $O((n/p) \log(n/p))$ pentru aceasta; aceeași leamnă ne spune că secvențele A_i sunt bitonice; ori ele pot fi ordonate mai repede (dar nu aplicând sortarea bitonică—aceasta conduce la un timp și mai mare: $O((n/p) \log^2(n/p))$); aici ne-ar aduce respectarea ad litteram a principiului lui Brent).

Sortarea secvențială a unei secvențe unimodale. (Sortarea unei secvențe bitonice este lăsată ca exercițiu.) O secvență unimodală cu minim este compusă din două subsecvențe, una descrescătoare, alta crescătoare. Dacă o privim pe cea descrescătoare în sens invers, obținem una crescătoare. Secvența sortată se obține prin interclasarea celor două subsecvențe crescătoare, deci în timp $O(n/p)$; mai detaliat: cel mult $\log(n/p)$ pentru găsirea minimumului și cel mult n/p pentru interclasare.

ALGORITHM 5.29 (sortarea secvenței unimodale A cu minim, secvențial)

1. găsește j , indicele elementului minim
2. interclasează $\langle a_j, a_{j-1}, \dots, a_0 \rangle$ cu $\langle a_{j+1}, a_{j+2}, \dots, a_{m-1} \rangle$

Sortare bitonică pe hipercub. Putem scrie acum algoritmul de sortare bitonică pe un hipercub, în cazul general; notăm A' secvența pe care o primește la un moment dat P_k de la un vecin al său (sunt necesare deci n/p locații de memorie suplimentare).

ALGORITHM 5.30 (sortare bitonică pe hipercub, $p = 2^d \ll n$, varianta 1, $\text{id} \equiv k$)

1. sortează elementele din A_k , crescător dacă $k_0 = 0$, descrescător dacă $k_0 = 1$
2. **pentru** $l = 0 : d - 1$
 1. **pentru** $i = l : -1 : 0$
 1. **în paralel** $\text{send}(A_k, i)$, $\text{recv}(A', i)$
 2. **dacă** $k_{l+1} = k_i$ **atunci** $A_k \leftarrow \min(A_k, A')$
 3. **altfel** $A_k \leftarrow \max(A_k, A')$
 4. sortează secvența bitonică A_k , crescător dacă $k_{l+1} = 0$, descrescător dacă $k_{l+1} = 1$

Numărul de comparații efectuat de acest algoritm este (termenul $2n/p$ corespunde complexității sortării unei secvențe bitonice, vezi problemele)

$$\frac{n}{p} \log \frac{n}{p} + \sum_{l=0}^{d-1} \left(\frac{2n}{p} + \sum_{i=0}^l \frac{n}{p} \right) \approx \frac{n}{p} \log n + \frac{n}{p} \log p + \frac{n}{2p} \log^2 p = O\left(\frac{n}{p} \log n\right) \quad (5.7)$$

și deci obiectivul propus a fost atins. În ce privește comunicația, aceasta ocupă timpul

$$\sum_{l=0}^{d-1} \sum_{i=0}^l \left(\sigma + \frac{n}{p} \beta \right) = \left(\sigma + \frac{n}{p} \beta \right) \log^2 p = O\left(\frac{n}{p} \log^2 p\right).$$

Se poate concluziona că algoritmul propus are o eficiență asimptotic egală cu 1. Vom vedea că se pot găsi algoritmi cu o complexitate mai bună a comunicației, deci cu un overhead datorat comunicației mai mic. În ce privește overhead-ul de încărcare inegală, algoritmul 5.30 stă foarte bine: procesoarele au aceeași cantitate de operații de efectuat; singura etapă la care pot interveni diferențe—însă ne semnificative—este cea de ordonare locală a secvențelor bitonice.

Alt algoritm de sortare bitonică. Să deducem acum o variantă mai simplă și mai eficientă a sortării bitonice pe hipercub. Pentru claritate, să precizăm câțiva termeni. Spunem că, în secvența $A = \langle A_0, \dots, A_{p-1} \rangle$, A_i este un minim local dacă $A_i < A_{i-1}$ și $A_i < A_{i+1}$; definiția este aceeași ca pentru cazul în care A_i ar fi elemente, doar semnificația semnului $<$ s-a schimbat. În continuare, noțiunile de unimodalitate și bitonicitate se definesc la nivel de subsecvență la fel ca la nivel de element.

În definirea secvențelor $m(A)$ și $M(A)$, subsecvențele componente au forma $\min(A_i, A_j)$, respectiv $\max(A_i, A_j)$; observația esențială este că, în acest fel, secvențele $m(A)$ și $M(A)$ sunt definite în mod *unic* din punct de vedere al elementelor componente (nu și al ordinii lor, evident). Indiferent de ordinea elementelor în interiorul unei subsecvențe, dacă A este bitonică la nivel de subsecvențe, atunci $m(A)$ și $M(A)$ din lema lui Batcher vor fi bitonice la nivel de subsecvențe. Atunci nu mai trebuie să avem grijă ca secvența A să fie în permanență bitonică la nivel de element, așa cum am făcut în algoritmul 5.30, ci doar bitonică la nivel de subsecvență.

Mai rămâne de stabilit metoda de obținere a subsecvențelor de genul $\min(A_i, A_j)$. Ideea cea mai la îndemână e de a ține în permanență secvențele A_i sortate crescător, și de a obține minimul prin interclasare. În acest fel, dacă secvențele A_i, A_j au fiecare câte m elemente, subsecvența minim se calculează cu doar m comparații și este ordonată crescător.

Toate aceste considerații ne arată că putem aplica foarte bine algoritmul 5.28 și în cazul $p \ll n$, cu două amendamente: elementele locale unui procesor să fie inițial sortate crescător și operațiile de tip $exch(i, j)$ să fie acum la nivel de subsecvență și să fie efectuate așa cum a fost descris în secțiunea dedicată interclasării.

ALGORITM 5.31 (sortare bitonică pe hipercub, $p = 2^d \ll n$, varianta 2, id $\equiv k$)

1. sortează crescător elementele din A_k (locale)
2. apelează algoritmul 5.28, cu $exch()$ la nivel de subsecvență

Numărul de comparații efectuat de acest algoritm este

$$T(n, p) = \frac{n}{p} \log \frac{n}{p} + \sum_{l=0}^{d-1} \sum_{i=0}^l \frac{n}{p} \approx \frac{n}{p} \log n + \frac{n}{2p} \log^2 p \quad (5.8)$$

și se dovedește mai mic decât cel din (5.7). Deci, nu numai că algoritmul 5.31 este mai simplu decât 5.30, dar are și un număr de comparații mai mic; în ce privește comunicația, algoritmi implică același efort.

Vom vedea, în problema 5.4.16, că efortul deducerii algoritmului 5.30 nu a fost totuși inutil. În cazul memoriei partajate, acesta este cel mai eficient.

5.4.5 Sortare prin interclasare

Algoritm secvențial. Sortarea prin interclasare (mergesort) are o descriere recursivă foarte scurtă: se împarte secvența A în două jumătăți, se sortează fiecare dintre ele și apoi se interclasează secvențele ordonate; înjumătățirea secvențelor are loc până când lungimea secvențelor este 1, deci ele sunt fiecare ordonate.

ALGORITM 5.32 (mergesort, forma generală secvențială)

procedură *mergesort*(s, d)

1. **cât timp** $s \neq d$ {secvența are mai mult de un element}
 1. $x \leftarrow \lfloor (s + d)/2 \rfloor$ {poziția mediană}
 2. *mergesort*(s, x) {sortează prima jumătate}
 3. *mergesort*($x + 1, d$) {sortează a doua jumătate}
 4. interclasează $A(s : x)$ cu $A(x + 1 : d)$ {cele două jumătăți}

Dacă nu vă place descrierea recursivă, puteți gândi altfel: se interclasează $n/2$ perechi de secvențe de un element, apoi $n/4$ perechi de secvențe de două elemente, etc. Complexitatea secvențială a algoritmului respectă $T(n) \leq T(n/2) + n$ (interclasarea a două secvențe de lungime $n/2$ se face cu cel mult n comparații), deci $T(n) \leq 2n \log n$. Sortarea prin interclasare implică de două ori mai multe operații decât quicksort, în cazul cel mai favorabil pentru acesta din urmă; marea deosebire în ce privește timpii de execuție este că pentru mergesort variațiile sunt destul de mici, deci $2n \log n$ este o bună aproximare, în timp ce pentru quicksort sunt mai mari. Totuși, în general, quicksort este mai rapid.

Algoritmi pe PRAM. Din forma algoritmului 5.32 se vede că o variantă paralelă a sa necesită numai un bun algoritm paralel de interclasare. Altfel, toate interclasările de secvențe de aceeași lungime pot decurge în paralel, începând de la cele $n/2$ interclasări de secvențe de lungime 1, și până la unica interclasare de secvențe de lungime $n/2$. Este același gen de paralelism ca și la quicksort, numai că aici gradul mare de paralelism este la început, dimensiunile secvențelor de interclasat crescând, în timp ce la quicksort era invers: prin partiționare se obțineau secvențe din ce în ce mai mici; mai există o diferență; aici dimensiunile se dublează după fiecare interclasare, în timp ce acolo subsecvențele obținute prin partiționare nu erau egale.

Din punct de vedere teoretic, pe un CREW PRAM, folosind $O(n)$ procesoare, putem aplica pentru interclasare algoritmul Valiant 5.11; cu acesta, pentru interclasarea a două secvențe de lungime t cu $O(t)$ procesoare este nevoie de $O(\log \log t)$ comparații. Deoarece pe tot timpul sortării prin interclasare suma elementelor din secvențele care se interclasează în paralel este n , rezultă că sunt necesare $O(n)$ procesoare. Numărul de comparații va fi $O(\log \log 1 + \log \log 2 + \dots + \log \log(n/4) + \log \log(n/2))$, unde fiecare termen corespunde unui nivel de recursie (primii sunt evident greșiți, primul chiar matematic absurd; ei pot fi oricum neglijați, ca fiind mici; am vrut să fie subliniată dimensiunea secvențelor de interclasat); efectuând suma (vezi **P5.4.20**) se obține o complexitate de $O(\log n \log \log n)$ pentru problema sortării. Costul acestui algoritm nu este totuși optim, ci $O(n \log n \log \log n)$; există un algoritm optim, care asigură sortarea în timp $O(\log n)$, cu $O(n)$ procesoare, dar el este destul de complicat; totuși, se bazează tot pe interclasare.

În cazul $p \ll n$, pe un CREW PRAM se poate aplica algoritmul de interclasare 5.16 (Xiong). Fiecare procesor sortează secvențial cele $m = n/p$ elemente ale sale. Apoi, procesoarele se grupează câte două pentru a interclasa două secvențe de lungime m , câte patru pentru a interclasa secvențe de lungime $2m$, câte opt pentru a interclasa secvențe de lungime $4m$, etc. Timpul necesar acestor operații este

$$\begin{aligned} T(n, p) &= \frac{n}{p} \log \frac{n}{p} + \sum_{i=1}^{\log p} \left[\frac{2 \cdot 2^{i-1} m}{2^i} + \frac{\log(2 \cdot 2^{i-1} m / 2^i)}{\log(\lfloor 2^i / 2^{i-1} \rfloor + 1)} + O(\log \log p) \right] \\ &= \frac{n}{p} \log n + \frac{1}{\log 3} \log \frac{n}{p} \log p + O(\log p \log \log p). \end{aligned} \quad (5.9)$$

Primul termen este pentru sortarea secvențială, iar cei din sumă sunt luați din (5.4). Pentru p mic, eficiența acestui algoritm va fi foarte bună.

Cazul memoriei distribuite. După cum am văzut în secțiunea dedicată interclasării, pe o arhitectură cu memorie distribuită nu există un algoritm eficient pentru a rezolva această problemă. Pe un hipercub, după sortarea secvențelor locale, doar prima etapă de interclasare decurge ușor, cea în care perechi procesoare își interclasează secvențele locale printr-un *exch()*; apoi ar urma ca grupuri de patru să interclaseze două secvențe ordonate distribuite pe câte două procesoare, grupuri de opt să interclaseze secvențe ordonate distribuite pe patru procesoare, etc. Nu se poate obține un mod eficient de comunicație (sau, cel puțin, el nu a fost găsit până în prezent). Dacă nu ați făcut-o deja, vedeți acum problema 5.4.17 pentru o frumoasă rețea de sortare prin interclasare, și pentru comunicația pe care ar implica-o într-un hipercub.

5.4.6 Ce algoritm de sortare alegem ?

Am prezentat o mulțime de algoritmi care rezolvă aceeași problemă, fiecare în alt mod; și să nu vă închipuiți că aceștia sunt toți. Pe care dintre ei îl alegem ? De ce nu ne-am mulțumit să descriem unul singur, cel mai bun ?

Să restrângem discuția la calculatoarele cu memorie distribuită. În tabelul 5.2 sunt recapitulate metodele paralele de sortare din acest capitol, împreună cu complexitățile

Algoritm	Comparații	Comunicație
Sortare par-impair (alg. 5.19)	$O(\frac{n}{p} \log n + n)$	$O(n)$
Hiperquicksort (Wagar, alg. 5.24)	$O(\frac{n}{p} \log n)$	$O(\frac{n}{2p} \log p)$
Sortare echilibrată (Abali, alg. 5.26+5.25)	$O(\frac{n}{p} \log n + p \log^2 n)$	$O(\frac{n}{p} \log p)$
Sortare cvasi-echilibrată (Shi, alg. 5.26+5.27)	$O(\frac{n}{p} \log n + p \log n)$	$O(\frac{n}{p} \log p)$
Sortare bitonică (Batcher, alg. 5.31)	$O(\frac{n}{p} \log n + \frac{n}{2p} \log^2 p)$	$O(\frac{n}{p} \log^2 p)$

Tabelul 5.2: Complexitățile metodelor de sortare pentru calculatoare MIMD cu memorie distribuită, prezentate în acest capitol.

lor—număr de comparații și număr de elemente comunicate, în medie. Dintre cele 5 metode, una singură este potrivită topologiei de inel, sortarea par-impair. Deci, dacă procesoarele calculatorului pe care îl avem sunt conectate în inel, aceasta e metoda aleasă; ceilalți algoritmi se pot și ei implementa pe inel, dar va crește complexitatea comunicației, ceea ce elimină interesul pentru ei. Pe de altă parte, dacă avem un hipercub, ar fi păcat să alegem sortarea par-impair, care nu profită decât de puține dintre conexiunile existente.

Alte criterii sunt dificil de generalizat; este vorba despre valoarea raportului viteză de comunicație/viteză de calcul pentru calculatorul respectiv; de comportarea algoritmilor în medie și în cazul cel mai defavorabil; de cât de aproape de medie se află în general timpul de execuție; de constantele ascunse de notația $O(\cdot)$; de caracteristici hardware sau software ale calculatorului, care pot fi speculate mai ușor sau mai greu de către un algoritm.

De obicei decizia vine pe baza posibilităților calculatorului și, nu în ultimul rând, a timpului disponibil pentru scrierea programului. Scopul acestei cărți e de a oferi o gama mare de variante, nu numai pentru a avea de unde alege, dar și pentru a dezvolta instrumente aplicabile altor probleme, netratate aici.

Probleme

P 5.4.1 Doi algoritmi secvențiali foarte populari de sortare cu complexitate $O(n^2)$ sunt metoda bulelor (bubblesort) și sortarea prin inserție. Îi prezentăm mai jos în forma lor cea mai simplă (se pot introduce unele teste care diminuează numărul de operații):

BUBBLESORT
pentru $i = n - 1 : -1 : 1$
 pentru $j = 0 : i - 1$
 $exch(j, j + 1)$

SORTARE PRIN INSERȚIE
pentru $i = 1 : n - 1$
 pentru $j = i - 1 : -1 : 0$
 $exch(j, j + 1)$

Desenați rețelele de sortare corespunzătoare acestor algoritmi și apreciați timpul de execuție, nivelul de paralelism, etc.

P 5.4.2 Scrieți algoritmul de sortare par-impair pentru un EREW PRAM, în cazurile $p = n/2$, $p \ll n$.

P 5.4.3 Pentru n impar, desenați rețelele analoage celor din figura 5.5. Sunt ele identice așa cum am afirmat în text ?

P 5.4.4 Fie α o rețea de sortare simplă cu $\binom{n}{2}$ comparatoare. Atunci α^R , rețeaua conținând aceleași comparatoare, dar în ordine inversă, este o rețea de sortare.

P 5.4.5 Ce modificări trebuie aduse algoritmului 5.21 și procedurii 5.22 pentru adaptarea la un EREW PRAM ?

P 5.4.6 Să presupunem că, în algoritmul 5.23, este aleasă ca pivot mediana unui eșantion de t elemente. Care este valoarea minimă pentru t astfel încât algoritmul să nu se oprească înainte de sortarea locală (aceasta se poate întâmpla dacă un procesor care trebuie să aleagă un pivot rămâne fără nici un element din secvență) ?

P 5.4.7 Este algoritmul de sortare echilibrată 5.26 într-adevăr echilibrat dacă în A pot exista elemente egale ?

P 5.4.8 Fie $X = \langle X_0, \dots, X_{p-1} \rangle$ o secvență de lungime n , fiecare X_i fiind ordonată crescător și având lungime n_i . Scrieți algoritmul (secvențial) de sortare a lui X , prin interclasare după un arbore binar. De câtă memorie suplimentară este nevoie ?

P 5.4.9 Scrieți algoritmul de sortare cvasi-echilibrată (Shi) pentru un PRAM.

P 5.4.10 Presupunem că trei procesoare dețin secvențele $\langle 0, 1, 2, 9, 16, 17, 24, 25, 27, 28, 30, 33 \rangle$, $\langle 7, 8, 11, 12, 13, 18, 19, 21, 23, 29, 34, 35 \rangle$, respectiv $\langle 3, 4, 5, 6, 10, 14, 15, 20, 22, 26, 31, 32 \rangle$, fiecare de lungime 12, deja sortate. Care sunt cheile de partiționare calculate de algoritmul 5.27 ? Câte elemente va deține în final fiecare procesor ?

P 5.4.11 Pe un hipercub, algoritmul 5.27 poate fi îmbunătățit după următoarea idee: interclasarea mesajelor X poate fi făcută în timpul colectării; pe un arbore de acoperire binomial, fiecare procesor care primește un X de la un procesor mai depărtat decât el de P_0 interclasează secvența primită cu cea proprie și trimite rezultatul mai departe. Scrieți algoritmul și apreciați timpul de execuție.

P 5.4.12 Să se arate că dintr-o secvență bitonică se pot obține, prin rotație, secvențe unimodale cu minim, respectiv cu maxim.

P 5.4.13 Scrieți algoritmul secvențial de găsim a minimului unei secvențe unimodale A , de lungime m , cu minim. Se poate adapta acest algoritm în cazul în care A este bitonică ?

P 5.4.14 Detaliați algoritmul 5.29 pentru cazul unei secvențe bitonice.

P 5.4.15 Găsiți o planificare a procesoarelor pentru sortarea bitonică pe un EREW PRAM, în cazul $p = n/2$.

P 5.4.16 Găsiți o variantă eficientă pentru sortarea bitonică pe un EREW PRAM, când $p \ll n$.

P 5.4.17 Un alt algoritm paralel de sortare propus de Batchier poate fi descris recursiv după cum urmează; în esență, este un algoritm de sortare prin interclasare, numit sortare cu interclasare par-impară (nici o legătură cu sortarea par-impar). Se împarte secvența de sortat C în două subsecvențe, $C = \langle A, B \rangle$; se sortează A și B ; apoi acestea se interclasează tot în mod recursiv: se interclasează secvențele "pare" $\langle a_0, a_2, \dots \rangle$ și $\langle b_0, b_2, \dots \rangle$, rezultatul fiind $\langle x_0, x_1, \dots \rangle$; se interclasează secvențele "impare" $\langle a_1, a_3, \dots \rangle$ și $\langle b_1, b_3, \dots \rangle$, obținându-se

$\langle y_0, y_1, \dots \rangle$; se consideră secvența $Z = \langle x_0, y_0, x_1, y_1, x_2, y_2, \dots \rangle$, căreia i se aplică operațiile de comparare-interschimbare $exch(y_0, x_1), exch(y_1, x_2), \dots$; în urma acestor operații $Z = A \perp B$.

Demonstrați corectitudinea acestui algoritm folosind următoarea teoremă, deosebit de utilă, formulată inițial pentru rețelele de sortare:

TEOREMĂ (principiul zero-unu). Dacă o rețea cu n linii de intrare (un algoritm cu $n = |C|$) sortează nedescrescător toate cele 2^n secvențe de zero și unu, atunci va sorta orice secvență de n numere în ordine nedescrescătoare.

P 5.4.18 Desenați o rețea de sortare conformă metodei din problema anterioară. Apreciați timpul de execuție. Studiați posibilitatea de adaptare a acestui algoritm la o arhitectură MIMD cu memorie distribuită, comparând apoi cu sortarea bitonică.

P 5.4.19 Algoritmul de sortare bitonică poate fi descris și altfel, tot sub formă recursivă. Dacă $A = \langle a_0, a_1, \dots \rangle$ este unimodală, atunci se sortează secvențele unimodale $\langle a_0, a_2, \dots \rangle$ ("pară") și $\langle a_1, a_3, \dots \rangle$ ("impară") și se efectuează operațiile $exch(a_0, a_1), exch(a_2, a_3), \dots$. Demonstrați corectitudinea algoritmului folosind principiul zero-unu. Verificați că rețeaua de sortare care rezultă este identică cu cea pentru sortarea bitonică.

P 5.4.20 Demonstrați că $\log \log n + \log \log(n/2) + \dots + \log \log 2 = O(\log n \log \log n)$.

Capitolul 6

Algoritmi numerici

Primii algoritmi pentru care s-a pus acut problema paralelizării au fost cei numerici, deoarece ei sunt folosiți în rezolvarea problemelor științifice de mari dimensiuni, pentru care nevoia de viteză de calcul e imperioasă. De aceea, pe primele calculatoare paralele, algoritmi numerici au fost implementați cu prioritate.

Să precizăm – deși e probabil inutil – că algoritmi numerici sunt cei în care se efectuează calcule aritmetice cu numere reale (sau complexe, caz pe care nu-l abordăm, el neimplicând mari diferențe), reprezentate în format virgulă mobilă. Exemplele de probleme numerice sunt numeroase: începând de la înmulțirea matrice-vector, tratată într-un capitol anterior, și de la rezolvarea de sisteme de ecuații liniare de diverse tipuri, până la, să zicem, rezolvarea de ecuații diferențiale, sau cu derivate parțiale. Spre deosebire, sortarea unor numere reale nu este o problemă numerică, deoarece, deși se efectuează operații (comparații, adică scăderi), rezultatul lor nu este folosit mai departe; valorile reale nu se modifică pe întreg parcursul algoritmilor de sortare, ci numai locurile lor în memorie, adică ordinea lor.

Scopul acestei lucrări fiind în primul rând de inițiere, ne vom opri la prezentarea de algoritmi pentru probleme simple din punct de vedere al tratărilor matematice și algoritmice. Este vorba de înmulțirea de matrice și de rezolvarea de sisteme liniare triunghiulare, dense (cele obișnuite; în acest context se va vorbi în mod natural și despre factorizarea LU) și tridiagonale; deci, capitole de bază ale algebrei liniare. Există, în plus, și un alt interes pentru algoritmi paraleli ce rezolvă aceste probleme; ei sunt folosiți în numeroase alte probleme, care se reduc la, sau conțin ca etape intermediare, acești algoritmi numerici fundamentali. Vom insista îndeosebi asupra algoritmilor dedicați arhitecturilor cu memorie distribuită, aici intervenind o serie de tehnici de paralelizare de mare interes.

Pentru a înțelege acest capitol cititorului îi sunt necesare cunoștințe de matematică la nivelul anului I de facultate (cel mult) și parcurgerea capitolelor referitoare la sisteme liniare ale unui manual de metode numerice, pentru familiarizarea cu algoritmi secvențiali. Nu vom aborda mai deloc problema stabilității numerice a algoritmilor

(cât de "aproape" de soluția exactă a problemei este cea calculată numeric) – care este, de fapt, piatra de încercare a algoritmilor numerici; aceasta nu înseamnă că am uitat-o, ci că punem accentul pe aspectele legate de paralelism.

Deoarece în toate problemele pe care le vom prezenta se lucrează cu matrice, să prezentăm întâi câteva notații, păstrate de-a lungul întregului capitol. Fie $A \in \mathbf{R}^{n \times n}$, deci o matrice pătrată, de dimensiune $n \times n$; nu vom lucra cu matrice dreptunghiulare; în singurul loc unde ar fi într-adevăr cazul, înmulțirea de matrice, generalizările sunt ușoare. Notăm prin a_{ij} , sau prin $A(i, j)$, elementul din linia i , coloana j , al matricei A ; vom numerota liniile și coloanele începând de la zero.

Pentru a descrie algoritmi la nivel de blocuri de matrice (și nu la nivel de element), se consideră următoarele partiționări ale matricei A :

$$A = \begin{bmatrix} A_{00} & A_{01} & \cdots & A_{0,n_b-1} \\ A_{01} & A_{11} & \cdots & A_{1,n_b-1} \\ \vdots & \vdots & & \vdots \\ A_{n_b-1,0} & A_{n_b-1,1} & \cdots & A_{n_b-1,n_b-1} \end{bmatrix} = [A_0 A_1 \dots A_{n_b-1}] = \begin{bmatrix} A_0^T \\ A_1^T \\ \vdots \\ A_{n_b-1}^T \end{bmatrix} \quad (6.1)$$

în care: $A_{ij} \in \mathbf{R}^{r \times r}$ este un bloc pătrat; $A_j \in \mathbf{R}^{n \times r}$ este o bloc-coloană; $A_i^T \in \mathbf{R}^{r \times n}$ este o bloc-linie. În notație MATLAB:

$$\begin{aligned} A_{ij} &\equiv A(ir : (i+1)r - 1, jr : (j+1)r - 1), \\ A_j &\equiv A(:, jr : (j+1)r - 1), \\ A_i^T &\equiv A(ir : (i+1)r - 1, :). \end{aligned}$$

Presupunem că n se divide cu r , și deci numărul de blocuri pe o dimensiune este $n_b = n/r$. Desigur că presupunerea că blocurile A_{ij} sunt pătrate este simplificatoare, dar ea ușurează înțelegerea, fără a diminua prea mult generalitatea. Atenție, notațiile A_i și A_i^T pot genera confuzii; aceste două matrice nu sunt una transpusa celeilalte; ele nu vor fi folosite niciodată împreună.

6.1 Înmulțirea de matrice

Înmulțirea de matrice este operația cea mai folosită în algoritmi din algebra liniară, în special în variantele bazate pe apeluri la rutine BLAS (Basic Linear Algebra Subroutines) de nivel 3 (în aceste variante, algoritmi sunt descriși la nivel de bloc, și reduși la câteva operații simple, implementate foarte eficient; vom prezenta o astfel de versiune pentru factorizarea LU). Ne vom ocupa în această secțiune de prezentarea câtorva algoritmi bazați pe idei diferite și având în vedere mai multe arhitecturi țintă.

Vom expune întâi problema și vom deduce complexitatea ei paralelă, apoi vom prezenta algoritmi pentru arhitecturi MIMD cu memorie distribuită, pentru topologiile de inel și tor, iar în final algoritmi pentru arhitecturi MIMD cu memorie partajată.

Referințele de bază în alcătuirea acestei secțiuni au fost [12, 13, 5].

6.1.1 Complexitatea problemei

Fie $A, B \in \mathbf{R}^{n \times n}$ două matrice. Notăm prin $MM(n)$ problema calculului matricei produs $C = A \cdot B \in \mathbf{R}^{n \times n}$. Prin definiție:

$$c_{ij} = \sum_{k=0}^{n-1} a_{ik} \cdot b_{kj}, \quad \text{pentru } i, j \in 0 : n - 1 \quad (6.2)$$

Denumim folosirea explicită a acestei formule *metoda standard* pentru rezolvarea $MM(n)$ (spre a o diferenția de metodele rapide, care nu vor fi însă studiate în această lucrare; vezi [11]).

La nivel de blocuri, calculul produsului $C = A \cdot B$ se poate face folosind partiționările din relația (6.1). Vom prezenta mai multe variante, fiecare sugerând o altă paralelizare. O primă formulă este similară cu cea la nivel de element (6.2):

$$C_{ij} = \sum_{k=0}^{n_b-1} A_{ik} B_{kj}, \quad \text{pentru } i, j \in 0 : n_b - 1 \quad (6.3)$$

în care pentru calculul fiecărui produs de blocuri se aplică formula la nivel de element.

O a doua formă reprezintă calculul unei bloc coloane din C :

$$C_j = A \cdot B_j = \sum_{i=0}^{n_b-1} A_i \cdot B_{ij}, \quad \text{pentru } j \in 0 : n_b - 1. \quad (6.4)$$

În fine, în a treia, se calculează o bloc linie din C :

$$C_i^T = A_i^T \cdot B = \sum_{j=0}^{n_b-1} A_{ij} \cdot B_j^T, \quad \text{pentru } i \in 0 : n_b - 1. \quad (6.5)$$

Și în ultimele două forme, produsele de blocuri se calculează cu formula la nivel de element, cu observația că acum se înmulțesc matrice (blocuri) dreptunghiulare.

Întorcându-ne la formula (6.2), aceasta arată numărul de operații necesare în rezolvarea prin metoda standard a problemei $MM(n)$, anume $T_1(n) = 2n^3$; calculul fiecărui element din C implică $2n$ operații (este un produs scalar), și sunt n^2 elemente. Dacă se calculează în paralel, se observă că toate cele n^3 (câte n pentru cele n^2 elemente din C) produse de tipul $a_{ik} b_{kj}$ care apar se pot calcula simultan, ele depinzând numai de datele de intrare; apoi, pentru obținerea fiecărui element al rezultatului trebuie calculată o sumă de n valori; această problemă a fost studiată anterior și se poate rezolva în timp $O(\log n)$, utilizând $O(n^3)$ procesoare (câte $O(n)$ pentru fiecare dintre cele n^2 elemente ale rezultatului); aceasta este și complexitatea algoritmului paralel standard de înmulțire de matrice.

Este interesant că și complexitatea problemei $MM(n)$ este $O(\log n)$, dar numărul de procesoare cu care se poate obține aceasta este mai mic. Nu s-a descoperit încă

limita inferioară a acestui număr de procesoare; se știe că este mai mare decât $O(n^2)$ – ceea ce se poate demonstra cu ușurință – și mai mică sau egală cu $O(n^{2.37})$, care este cel mai mic majorant descoperit până acum, în 1987 [9]. De remarcat legătura între complexitatea (în timp) a algoritmilor secvențiali rapizi și complexitatea ”în spațiu” (număr de procesoare) a algoritmilor paraleli ($O(n^{2.37})$ este cel mai bun timp secvențial pentru $MM(n)$); este așa numita *teză a calculului paralel*, vezi [16].

O soluție facilă de a construi algoritmi pentru produsul de matrice este aceea de a utiliza ca operație de bază produsul matrice-vector, operație deja discutată. De exemplu, folosind ca bază relația (6.4), cu $r = 1$ (deci $n_b = n$), și ținând seama că produsele de tipul AB_j se pot calcula în paralel, o descriere simplă este următoarea:

ALGORITM 6.1 ($MM(n)$ paralel, o soluție generală)

pentru $j = 0 : n - 1$, **în paralel**
 calculează $C(:, j) \leftarrow A \cdot B(:, j)$

Desigur că pentru acest algoritm mai trebuie detaliată planificarea operațiilor, descrierea comunicației, etc. O versiune asemănătoare se poate obține pornind de la relația (6.5).

În cele ce urmează, vom descrie algoritmi relativ apropiați de această idee, deși forma va putea părea mult diferită. În fond, aspectele specifice înmulțirii de matrice se referă la gruparea mesajelor comunicate, pentru a se micșora timpii de start-up, și la alternarea calcule-comunicație care este cea mai avantajoasă. Totuși, vor apare și idei noi.

6.1.2 Algoritmi pe inel

Vom porni analiza de la ipoteza că datele sunt repartizate echilibrat între procesoare, fără repetiții, și, de asemenea, că toate matricele – și cele de intrare și cea de ieșire – au aceeași repartizare.

O primă soluție, și singura pe care o detaliem, este aceea de a lucra pe linii, fiecare procesor deținând deci $r = n/p$ linii; așadar, în formulele anterioare, $n_b = p$. Presupunem o repartizare bloc linii, cea mai simplă de manipulat; procesorul P_i va avea în memoria sa blocurile A_i^T, B_i^T (sau, dacă ne referim la blocuri $r \times r$, blocurile A_{ij}, B_{ij} , cu $j = 0 : p - 1$); în final, P_i va avea bloc linia corespunzătoare a rezultatului, C_i^T .

Un algoritm simplu este sugerat de prima egalitate din relația (6.5), adică $C_i^T = A_i^T \cdot B$. De aici se vede că procesorul P_i are toate elementele din A necesare calculului părții sale din rezultat, dar are nevoie de întreaga matrice B . O difuzare generală în care fiecare procesor participă cu bloc linia sa din B ar aduce tuturor procesoarelor datele necesare în memoria proprie; acesta este echivalentul algoritmului 4.15. Din păcate, defectul pe care l-am subliniat în cazul produsului matrice-vector este încă și mai important acum. În loc de $3rn$ elemente (câte o bloc linie din A, B și C), un procesor ar trebui să aloce spațiu pentru mai mult de n^2 elemente în memoria locală, ceea ce este foarte mult. De aceea soluția este, în general, inacceptabilă.

Trebuie deci să imaginăm un mod de comunicație prin care bloc liniile matricii B să circule pe la toate procesoarele; desigur, A rămâne nemișcată, ca de altfel și C . De exemplu, se translează spre stânga bloc liniile din B : la fiecare etapă un procesor trimite către stânga bloc linia sa și recepționează o bloc linie din dreapta, aceasta până când bloc liniile revin la procesoarele care le-au deținut inițial; evident, în fiecare etapă procesoarele efectuează calcule cu datele locale din etapa respectivă. Se poate observa că paradigma de comunicație este cea studiată într-un exemplu din secțiunea ??; de asemenea, este modul de lucru din algoritmul 4.16.

Vom nota cu l indicele liniilor aparținând unui procesor, și cu c indicele liniilor din B primite la un moment dat. În fiecare dintre cele p etape, procesorul P_i calculează un produs de tipul $A_{ij}B_j^T$, ca în termenul din dreapta al relației (6.5), deci, între un bloc local din A (pătrat) și bloc linia din B recepționată anterior. În algoritmul ce urmează, instrucțiunile sunt scrise în stil MATLAB; în chip de comentariu este prezentată și o versiune folosind notația din (6.5).

ALGORITM 6.2 ($MM(n)$ pe inel, $n \gg p$, repartizare bloc linii pentru A, B, C ; pentru procesorul P_i)

```

 $l \leftarrow ir : (i+1)r - 1, c \leftarrow l, C(l,:) \leftarrow 0, j \leftarrow i \quad \{C_i^T \leftarrow 0\}$ 
pentru  $k = 0 : p - 1$ 
   $C(l,:) \leftarrow C(l,:) + A(l,c)B(c,:)$ 
  în paralel
     $\text{send}(B(c,:), \text{stânga})$ 
     $\text{recv}(D, \text{dreapta})$ 
   $j \leftarrow (j+1) \bmod p$ 
   $c \leftarrow jr : (j+1)r - 1, B(c,:) \leftarrow D$ 

```

Trebuie remarcat că, atunci când este vorba de o implementare efectivă, $A(l,:)$, $C(l,:)$ și $B(c,:)$ (ca și variabila auxiliară D) sunt matrice de dimensiune $r \times n$ în memoria locală, și nu blocuri ale unor matrice $n \times n$ memorate integral. Algoritmul a fost scris astfel pentru a fi mai sugestiv. Doar aici, pentru exemplificare, vom prezenta și forma efectivă, apropiată de sursa într-un limbaj de programare; A_{local} , B_{local} și C_{local} sunt variabilele locale conținând bloc linii ale matricelor A, B , respectiv C .

ALGORITM 6.3 ($MM(n)$ pe inel, $n \gg p$, pentru procesorul P_i)

```

 $j \leftarrow i, C_{local} \leftarrow 0$ 
pentru  $k = 0 : p - 1$ 
   $c \leftarrow jr : (j+1)r - 1$ 
   $C_{local} \leftarrow C_{local} + A_{local}(:,c)B_{local}$ 
  în paralel
     $\text{send}(B_{local}, \text{stânga}), \text{recv}(D, \text{dreapta})$ 
   $j \leftarrow (j+1) \bmod p, B_{local} \leftarrow D$ 

```

Calcululele sunt perfect echilibrate între procesoare; fiecare procesor transmite p mesaje de lungime nr fiecare. Timpul total de execuție este deci:

$$T_{6.2}(n, p) = \frac{2n^3}{p}\alpha + p\sigma + n^2\beta,$$

iar overhead-ul datorat comunicației: $O_c(n, p) = \frac{p\sigma + n^2\beta}{T_{6.2}(n, p)}$. Acest overhead tinde la zero – iar eficiența la 1 – pe măsură ce n crește; deci algoritmul este cu atât mai bun cu cât matricele sunt mai mari.

Un algoritm asemănător se poate obține adoptând o repartizare bloc coloane, și făcând să circule blocurile matricei A . Întrucât este simplu, îl lășăm ca problemă. Timpul de execuție va fi identic cu cel de mai sus.

6.1.3 Algoritmi pe tor

Pornim din nou de la o repartizare uniformă a matricelor; procesorul P_{ij} are în memoria locală blocurile matricelor de intrare A_{ij} , B_{ij} , și va calcula blocul C_{ij} al rezultatului; formula de calcul este cea din (6.3) și, presupunând torul de dimensiuni egale, dimensiunea unui bloc este $r = n/\sqrt{p}$ (deci $n_b = \sqrt{p}$, în formula amintită).

Rezultatul C poate rămâne local procesoarelor, neexistând nici o dependență între blocuri diferite ale rezultatului; se observă însă că, spre deosebire de cazul topologiei de inel, trebuie să circule acum și A și B , blocurile locale ale nici uneia dintre aceste matrice nefiind suficiente pentru calculul din (6.3). Din această formulă se mai observă că blocurile din A de care are nevoie P_{ij} se află pe linia i a torului, iar blocurile necesare din B pe coloana j , adică pe linia și coloana pe care se află procesorul. De aici rezultă că matricea A va circula pe orizontală (pe liniile torului), iar B pe verticală (pe coloanele torului).

Să presupunem că blocurile matricei B sunt transmise, la fiecare etapă, în sus de către toate procesoarele, până când revin la procesoarele care le-au deținut inițial; sunt deci \sqrt{p} etape; în acest fel, fiecare procesor va avea acces la o bloc coloană a matricei B , ceea ce va fi suficient pentru a efectua calculele, din punctul de vedere al lui B .

Cum va circula A ? La un moment dat, cum o linie de procesoare deține o bloc linie a matricei B , pentru formarea unor produse care să contribuie la formarea unei bloc linii a rezultatului este necesar un singur bloc al matricei A – deținut de un singur procesor; de exemplu, în poziția inițială, bloc linia B_{00} , B_{01} , ..., $B_{0, \sqrt{p}-1}$ trebuie înmulțită cu blocul A_{00} ; este vorba despre, în relația (6.5), primul termen al sumei, cel pentru $i = 0$. Apoi, după ce blocurile matricei B sunt comunicate vecinilor de sus (practic matricea se translează în sus), pe linia 0 de procesoare se va găsi acum bloc linia B_{10} , B_{11} , ..., $B_{1, \sqrt{p}-1}$; ea va trebui înmulțită cu blocul A_{01} .

Se observă deci că matricea A poate păstra repartizarea inițială (nu este necesar să circule toată, ca B) dar, la fiecare etapă, câte un bloc al lui A va trebui difuzat pe fiecare linie de procesoare de către deținătorul său. Modul de comunicație este exemplificat în figura 6.1, pentru un tor cu 3×3 procesoare. În fiecare căsuță se află

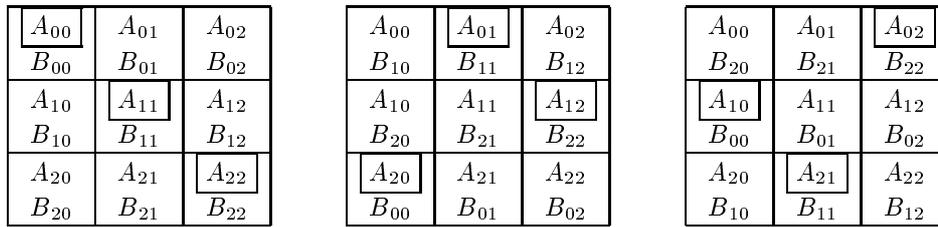


Figura 6.1: Comunicația în algoritmul pentru $MM(n)$ pe tor, cu difuzare pe linii a unor blocuri ale matricei A și cu translație pe verticală a matricei B .

blocurile locale procesorului din poziția corespondentă, la fiecare etapă a algoritmului, începând cu repartizarea inițială. Blocul din A aflat într-un mic chenar este cel difuzat pe linia procesorului respectiv în etapa curentă. După \sqrt{p} etape (3, în cazul nostru), se obține din nou repartizarea inițială.

Putem descrie adresa procesorului care difuzează la un moment dat, pe linia i ; în etapa k , cu $k = 0 : \sqrt{p} - 1$, difuzează procesorul din coloana j pentru care $(i - j + k) \bmod p = 0$; altfel spus, în prima etapă difuzează procesorul de pe diagonală, în a doua cel din dreapta sa, apoi următorul din dreapta, etc. Difuzarea pe o linie a torului este o operație banală, o linie de procesoare în tor fiind un inel. Vom nota cu c indicele de coloană al procesorului care face difuzarea: $c = (i + k) \bmod p$. Prezentăm algoritmul doar în notație la nivel de bloc; toate variabilele sunt blocuri $r \times r$.

ALGORITHM 6.4 ($MM(n)$ pe tor, algoritmul cu difuzare pe linii a matricei A , $n \gg p$, pentru procesorul P_{ij})

```

 $c \leftarrow i, C_{ij} \leftarrow 0$ 
pentru  $k = 0 : \sqrt{p} - 1$ 
(1)   difuzare (parțială):  $P_{ic}$  difuzează  $A_{ic}$  pe linia  $i$  a torului
       $C_{ij} \leftarrow C_{ij} + A_{ic}B_{cj}$ 
      în paralel
        send( $B_{cj}$ , sus)
        recv( $B_{(c+1) \bmod p, j}$ , jos)
       $c \leftarrow (c + 1) \bmod p$ 

```

Pentru claritate, să explicăm mai mult instrucțiunea (1), din punctul de vedere al procesorului P_{ij} ; în fiecare etapă k , pe linia i se face o difuzare; dacă $j = c$, atunci P_{ij} este procesorul care inițiază difuzarea; dacă $j \neq c$, atunci P_{ij} este unul dintre destinatarii difuzării, el recepționând blocul A_{ic} de la un vecin și, eventual, transmitându-l celuilalt, conform algoritmului de difuzare pe inel folosit.

Fiecare procesor are nevoie de cinci variabile de dimensiune $r \times r$: trei pentru C_{ij} , A_{ij} și blocul curent din B , plus două variabile auxiliare (care pot fi identice, de fapt), una pentru blocul difuzat din A , alta pentru recepția următorului bloc din B .

A_{00}	A_{01}	A_{02}
B_{00}	B_{11}	B_{22}
A_{11}	A_{12}	A_{10}
B_{10}	B_{21}	B_{02}
A_{22}	A_{20}	A_{21}
B_{20}	B_{01}	B_{12}

A_{01}	A_{02}	A_{00}
B_{10}	B_{21}	B_{02}
A_{12}	A_{10}	A_{11}
B_{20}	B_{01}	B_{12}
A_{20}	A_{21}	A_{22}
B_{00}	B_{11}	B_{22}

A_{02}	A_{00}	A_{01}
B_{20}	B_{01}	B_{12}
A_{10}	A_{11}	A_{12}
B_{00}	B_{11}	B_{22}
A_{21}	A_{22}	A_{20}
B_{10}	B_{21}	B_{02}

Figura 6.2: Comunicația în algoritmul de înmulțire de matrice pe tor, cu preskewing.

Ca și la inel, calculele sunt echilibrat repartizate între procesoare. În ce privește timpul necesar comunicației, evaluarea este ceva mai dificilă deoarece difuzarea este o operație a cărei eficiență depinde de algoritmul folosit și de caracteristicile fizice al calculatorului țintă. Considerând cel mai bun timp de difuzare pe inel (fără însă a face pipeline), adică $(\sqrt{p}/2)(\sigma + (n^2/p)\beta)$, timpul total va fi

$$T_{6.4} = \frac{2n^3}{p}\alpha + \sqrt{p}\left(\sigma + \frac{n^2}{p}\beta\right) + \sqrt{p}\frac{\sqrt{p}}{2}\left(\sigma + \frac{n^2}{p}\beta\right),$$

ultimul termen reprezentând timpul de difuzare, iar cel din mijloc cel pentru comunicarea matricei B (reamintim că un bloc B_{ij} are n^2/p elemente). După cum se vede, timpul de comunicație este $O(n^2)$, comparabil cu cel pentru inel. El poate fi redus prin optimizarea difuzării, dar nu va scădea sub, în cel mai bun caz, $O(2n^2/\sqrt{p})$, presupunând că difuzarea unui bloc durează tot atât cât comunicarea blocului între doi vecini. O ultimă îmbunătățire este executarea în paralel a difuzării și a translării matricei B , în fiecare etapă, deoarece ele decurg pe direcții diferite; vezi problema 6.1.2. Aceasta reduce timpul de comunicație la cel necesar difuzării.

Un alt algoritm pe tor provenind din metoda standard pentru rezolvarea $MM(n)$ se obține impunând o circulație asemănătoare pentru matricele A și B , dar pe direcții diferite – A spre stânga, B în sus (aceasta din urmă la fel ca în algoritmul 6.4); în acest mod se va reduce cu siguranță comunicația. Rămâne de văzut dacă această idee este valabilă și, dacă da, de rezolvat problema organizării întâlnirilor corespunzătoare între blocurile din A și din B , conform relației (6.3), astfel încât blocurile rezultatului să se formeze local; mai exact, să se găsească o dispunere inițială a matricelor A și B astfel încât, în cele \sqrt{p} etape de comunicație, în P_{ij} să se întâlnească, pe rând, toate blocurile A_{ik} , B_{kj} , cu $k \in 0 : \sqrt{p} - 1$; evident, ordinea întâlnirilor nu are importanță.

Se vede imediat că repartizarea ”obișnuită”, cu P_{ij} posedând A_{ij} și B_{ij} , nu este de folos; termenul $A_{ij}B_{ij}$ nu apare, în general, în sumele din (6.3). Dispunerea inițială a blocurilor matricelor A și B este cea prezentată în figura 6.2, în stânga, și a fost descoperită de Cannon (e cunoscută cu numele de ”preskewing” – repartizare ”oblică”).

Această repartizare poate fi dedusă relativ simplu, o dată stabilit sensul de mișcare

pentru A și B – stânga, respectiv sus (ideea algoritmului însă, nu e deloc banală!). Se începe prin a completa prima linie de procesoare cu prima bloc linie din A , în modul obișnuit: A_{0j} este local lui P_{0j} . Pe aceeași linie, blocurile din B sunt hotărâte de faptul că P_{0j} calculează C_{0j} , deci are nevoie de B_{jj} pentru a-l înmulți cu A_{0j} . Argumentăm în continuare doar pentru P_{10} , pentru celelalte procesoare procedându-se analog, din aproape în aproape; în a doua etapă, deoarece A se deplasează spre stânga, în P_{00} va ajunge A_{01} ; pentru a calcula C_{00} , acest bloc trebuie înmulțit cu B_{10} ; deci B_{10} este deținut inițial de P_{10} , care l-a trimis în sus în a doua etapă; tot inițial, în P_{10} se află A_{11} , singurul bloc care se înmulțește cu B_{10} pentru a contribui la C_{10} .

O dată completată dispunerea inițială, algoritmul este ușor de verificat, ceea ce este ilustrat și în figura 6.2, în care sunt prezentate cele trei etape ale înmulțirii pe un tor 3×3 . Într-o concluzie succintă, în repartizarea ”oblică” procesorul P_{ij} are inițial blocurile $A_{i,(j+i) \bmod \sqrt{p}}$ și $B_{(i+j) \bmod \sqrt{p},j}$.

Dată această repartizare inițială, algoritmul de rezolvare a problemei $MM(n)$ are o structură foarte simplă. Notând cu indicele *local* cele trei blocuri din A , B , C aflate la un moment dat în posesia unui procesor, și cu D și E două variabile auxiliare de aceeași dimensiune ($r \times r$), algoritmul va fi:

ALGORITHM 6.5 ($MM(n)$ pe tor, algoritmul cu preskewing, $n \gg p$, pentru proc. P_{ij})

```

 $C_{ij} \leftarrow 0$ 
pentru  $k = 0 : \sqrt{p} - 1$ 
   $C_{local} \leftarrow C_{local} + A_{local} \cdot B_{local}$ 
  în paralel
    send( $A_{local}$ , stânga)
    rcv( $D$ , dreapta)
    send( $B_{local}$ , sus)
    rcv( $E$ , jos)
   $A_{local} \leftarrow D, B_{local} \leftarrow E$ 

```

În această formă, timpul de execuție este următorul:

$$T_{6.5} = \frac{2n^3}{p}\alpha + \sqrt{p}\left(\sigma + \frac{n^2}{p}\beta\right).$$

Se observă dispariția întregului termen ce reflecta complexitatea difuzării, ceea ce conduce la un timp de comunicație de $O(n^2/\sqrt{p})$, mai mic decât cel corespunzător din algoritmul 6.4.

Nu trebuie să uităm însă că, de cele mai multe ori, matricele sunt disponibile în repartizarea prezentată la începutul secțiunii, adică procesorul P_{ij} are blocurile A_{ij} , B_{ij} . La dispunerea din figura 6.2 se poate ajunge ușor din această poziție prin translări, pe linii pentru A și pe coloane pentru B . Deoarece la procesorul P_{ij} trebuie să ajungă blocul $A_{i,(j+i) \bmod \sqrt{p}}$, se observă că sunt necesare i translări la stânga ale blocurilor din A , de pe linia i , deoarece coloana $(j+i) \bmod \sqrt{p}$ este cu i poziții la

dreapta (pe inelul format de linia i a torului) față de coloana j . În mod analog, pentru ca $B_{(i+j) \bmod \sqrt{p}, j}$ să ajungă în posesia procesorului P_{ij} , sunt necesare j translări în sus ale blocurilor din B , pe coloana j a torului. Un algoritm pentru preskewing este deci următorul (am omis variabilele auxiliare necesare ca tampon la recepție):

ALGORITM 6.6 (preskewing pe tor, $n \gg p$, pentru procesorul P_{ij})

```

în paralel
  pentru  $k = 0 : i - 1$ 
    în paralel
      send( $A_{local}$ , stânga)
      rcv( $A_{local}$ , dreapta)
  pentru  $k = 0 : j - 1$ 
    în paralel
      send( $B_{local}$ , sus)
      rcv( $B_{local}$ , jos)

```

Timpul de execuție al algoritmului 6.6 este apropiat de cel necesar unei singure difuzări pe un inel; ori, în algoritmul 6.4 se efectuau \sqrt{p} difuzări pe fiecare linie de procesoare; de unde o complexitate a comunicațiilor mai bună pentru algoritmi 6.5 și 6.6, luați împreună. Algoritmul 6.6 se poate îmbunătăți, dacă translarea se face în sensul drumului cel mai scurt; de exemplu, pe ultima linie, în loc de $\sqrt{p} - 1$ translări la stânga, e suficientă una singură la dreapta; astfel se reduce la jumătate timpul de comunicație; vezi problema 6.1.3.

Pentru ambii algoritmi prezentați, deci indiferent de repartizarea inițială a matricelor, eficiența tinde asimptotic la 1 deoarece, pentru n mare, termenul reprezentând complexitatea comunicațiilor este de $O(n^2)$, iar cel reprezentând calculele de $O(n^3)$.

6.1.4 Un algoritm pentru memorie partajată

Presupunem acum o arhitectură MIMD cu memorie partajată, cu p procesoare, fiecare având o mică memorie locală; cu minime modificări, algoritmul este valabil și în lipsa memoriei locale. Vom folosi ideea algoritmului 6.2 în ce privește modul de desfășurare a calculelor și repartizarea matricelor pe procesoare, cu deosebirea că acum matricele se află în memoria comună, și nu în cele locale, ceea ce simplifică sarcina.

Așadar, procesorul P_i se ocupă de o bloc linie din produs, pe care o va calcula pe baza relației (6.5), în care $r = n/p$ și $n_b = p$: $C_i^T = A_i^T \cdot B$. În primul rând, P_i va transfera în memoria sa locală bloc linia A_i^T , care, așa cum se vede, nu mai trebuie cunoscută de alte procesoare. În schimb matricea B trebuie citită în întregime de toate procesoarele, ceea ce pune două probleme: întâi, memoria locală este mică, și deci, în general, nu ajunge pentru stocarea matricei B ; apoi, dacă toate procesoarele citesc simultan aceleași elemente din matricea B , pot exista conflicte la citire care vor întârzia calculele efective. Soluția este sugerată de algoritmul 6.2: fiecare procesor

să acceseze câte o bloc linie din B la un moment dat, astfel încât, în p etape, un procesor să acceseze în întregime matricea B . Algoritmul care urmează este deci o simplă adaptare a algoritmului 6.2, al cărei aspect este apropiat de varianta 6.3.

ALGORITM 6.7 ($MM(n)$ pe arhitectură cu memorie comună, $n \gg p$, pt. proc. P_i)

```

 $j \leftarrow i, c \leftarrow ir : (i + 1)r - 1, C_{local} \leftarrow 0$ 
get( $A(c, :), A_{local}$ )           {citește  $A_i^T$ }
pentru  $k = 0 : p - 1$ 
    get( $B(c, :), B_{local}$ )       {citește  $B_j^T$ }
     $C_{local} \leftarrow C_{local} + A_{local}(:, c)B_{local}$     $\{C_i^T \leftarrow C_i^T + A_{ij}B_j^T\}$ 
     $j \leftarrow (j + 1) \bmod p, c \leftarrow jr : (j + 1)r - 1$ 
put( $C_{local}, C(ir : (i + 1)r - 1, :)$ )   {scrie  $C_i^T$ }

```

Variabilele cu indicele *local* sunt matrice $r \times n$ în memoria locală a fiecărui procesor. A_{local} conține permanent bloc linia A_i^T ; în B_{local} se transferă, din memoria comună, bloc linii din B ; în C_{local} se calculează bloc linia C_i^T a rezultatului, care este transferată, în final, în memoria comună, în poziția corespunzătoare; nu se pune problema nici unui conflict de scriere. Deși am presupus că după fiecare iterație dintr-o buclă **pentru** se face sincronizare, aici această ipoteză nu este neapărat necesară; e drept, nesincronizările pot introduce întârzieri, dar nu afectează corectitudinea algoritmului.

Algoritmul este echilibrat din punctul de vedere al încărcării; se apelează de $p + 2$ ori funcția **get** și o dată funcția **put**. Timpul de execuție al algoritmului va fi deci:

$$T_{6.7} = \frac{2n^3}{p}\alpha + (p + 2)\left(\sigma + \frac{n^2}{p}\beta\right).$$

Și pentru acest algoritm eficiența tinde asimptotic la 1.

Algoritmul 6.7 poate avea un singur inconvenient – insuficientă memorie locală pentru stocarea simultană a câte unei bloc coloane din A , B și C , ceea ce presupune spațiu pentru $3n^2/p$ elemente în virgulă mobilă. În acest caz se poate adopta, de exemplu, o partiționare bloc ciclică pe linii a matricelor; să spunem că fiecare procesor se ocupă de q bloc linii, fiecare bloc linie având r linii, și deci acum $n = rqp$. Fiecare procesor va calcula, pe rând, printr-un algoritm identic cu 6.7, fiecare dintre cele q bloc linii ce-i revin din C . Pentru fiecare bloc linie, trebuie citită întreaga matrice B , deci timpul de execuție al algoritmului va fi

$$T'_{6.7} = \frac{2n^3}{p}\alpha + q(qp + 2)\left(\sigma + \frac{n^2}{qp}\beta\right),$$

adică o comunicație de aproximativ q ori mai costisitoare decât în cazul precedent.

Probleme

P 6.1.1 Scrieți un algoritm pentru $MM(n)$ pe inel, pornind de la ipoteza unei repartizări bloc coloane pentru matricele A , B , C .

P 6.1.2 Modificați algoritmul 6.4, de înmulțire de matrice pe tor, astfel încât difuzarea unui bloc din A și comunicarea blocurilor din B vecinilor de sus să decurgă în paralel.

P 6.1.3 Modificați algoritmul 6.6 a. î. translările matricelor A și B să se facă întotdeauna pe drumul cel mai scurt.

P 6.1.4 Scrieți un algoritm de rezolvare a $MM(n)$ pe arhitecturi cu memorie comună, pornind de la egalitatea $C_{ij} = A_i^T \cdot B_j$. De ce n-a putut fi utilizată această egalitate în cazul inelului?

P 6.1.5 Scrieți un algoritm pentru $MM(n)$ pe inel, presupunând A repartizată bloc linii, iar B bloc coloane.

6.2 Sisteme liniare triunghiulare

Ne vom ocupa în această secțiune de problema rezolvării sistemelor liniare de forma $Ax = b$, cu $A \in \mathbf{R}^{n \times n}$, inferior triunghiulară și $b \in \mathbf{R}^n$; deci, $a_{ij} = 0$, pentru $i < j$; presupunem că problema are soluție unică, adică toate elementele diagonale a_{ii} sunt nenule. Deși această problemă, pe care o vom nota $TR(n)$, este banală în tratarea secvențială și are o complexitate de numai $O(n^2)$ – față de $O(n^3)$, pentru înmulțirea de matrice sau rezolvarea de sisteme dense – se impune studierea ei cu atenție în cazul implementării pe o arhitectură paralelă, deoarece, așa cum se va vedea, sunt necesare multe comunicații, care, dacă nu sunt bine organizate, pot scădea mult eficiența; de altfel, dacă pentru alți algoritmi de calcul numeric eficiențe de 80-90% sunt obișnuite, pentru rezolvarea sistemelor triunghiulare 50% este deja o bună performanță. De asemenea, rezolvarea de sisteme triunghiulare constituie etapa finală în rezolvarea sistemelor liniare dense, de unde și un interes practic deosebit.

După ce vom prezenta pe scurt considerente legate de rezolvarea secvențială a problemei, de complexitatea secvențială și de cea paralelă, vor fi discutați și evaluați mai mulți algoritmi paraleli pentru $TR(n)$. Ei vor fi grupați în perechi; algoritmi din aceeași pereche sunt duali, fiecare fiind inspirat dintr-una din cele două forme secvențiale clasice. Toți algoritmi sunt adecvați arhitecturilor cu memorie distribuită și va fi considerat doar cazul practic $n \gg p$. Referințele bibliografice cele mai folosite în alcătuirea acestei secțiuni au fost Heath & Romine [14], Chamberlain [7] și Li & Coleman [20, 21].

6.2.1 Algoritmi secvențiali

Există doi algoritmi consacrați pentru rezolvarea de sisteme triunghiulare. Ei se deosebesc numai prin ordinea operațiilor, având aceeași idee de bază – substituția înainte – și aceeași complexitate. Din linia i , cu $i \in 0 : n - 1$, a ecuației $Ax = b$ se deduce imediat că:

$$x_i = \left(b_i - \sum_{j=0}^{i-1} a_{ij} x_j \right) / a_{ii}, \quad (6.6)$$

cea ce arată că elementele din x se pot calcula în ordine crescătoare a indicilor, fiecare x_i depinzând numai de x_j cu $j < i$; astfel, în momentul calculării lui x_i , în relația de mai sus, toți termenii din dreapta sunt cunoscuți.

Să presupunem că x este inițializat cu elementele din b . În primul algoritm, după calcularea unui element x_j al soluției, se actualizează cu ajutorul acestuia toate elementele x_i , cu $i > j$, adică elementele încă necalulate din x . Practic, se adună vectorul x cu un multiplu al coloanei j a matricei A – de aceea această formă este numită algoritmul cu *sumă vectorială*:

ALGORITM 6.8 ($TR(n)$ secvențial, prin sumă vectorială)

```

 $x \leftarrow b$ 
pentru  $j = 0 : n - 1$ 
     $x_j \leftarrow x_j / a_{jj}$ 
    pentru  $i = j + 1 : n - 1$ 
         $x_i \leftarrow x_i - a_{ij}x_j$ 

```

De notat că bucla **pentru** i poate fi înlocuită prin suma vectorială $x(j + 1 : n - 1) \leftarrow x(j + 1 : n - 1) - A(j + 1 : n - 1, j)x(j)$.

În a doua formă, o necunoscută se calculează folosind toate elementele deja calculate ale soluției, adică aplicând formula (6.6) ca atare. Pentru a realiza acest calcul se efectuează produsul scalar între linia i a matricei A și vectorul x , mai exact $A(i, 0 : i - 1) \cdot x(0 : i - 1)$ (care poate înlocui bucla **pentru** j de mai jos) – de aceea această formă este numită algoritmul cu *produs scalar*:

ALGORITM 6.9 ($TR(n)$ secvențial, prin produs scalar)

```

 $x \leftarrow b$ 
pentru  $i = 0 : n - 1$ 
    pentru  $j = 0 : i - 1$ 
         $x_i \leftarrow x_i - a_{ij}x_j$ 
     $x_i \leftarrow x_i / a_{ii}$ 
     $\{x(i) \leftarrow A(i, 0 : i - 1) \cdot x(0 : i - 1)\}$ 

```

Numărul de operații, în ambii algoritmi, este n^2 .

6.2.2 Complexitatea paralelă

Pentru ambii algoritmi, cu sumă vectorială sau cu produs scalar, graful de precedență este același, și anume cel din figura 6.3, în care s-au folosit notațiile: T_{ij} pentru operația $x_i \leftarrow x_i - a_{ij}x_j$, și S_i pentru operația $x_i \leftarrow x_i / a_{ii}$. Relațiile de precedență sunt $S_i \prec T_{ki}$ și $T_{ij} \prec S_i$.

În primul algoritm, ordinea de execuție este pe coloane, adică $S_0, T_{10}, T_{20}, T_{30}, S_1, T_{21}$, etc., în timp ce în al doilea ordinea este pe linii, adică $S_0, T_{10}, S_1, T_{20}, T_{21}$, etc.

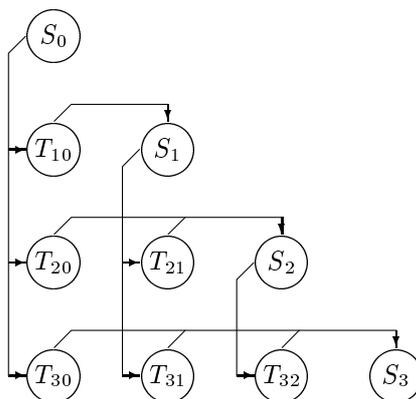


Figura 6.3: Graf de dependență pentru rezolvarea unui sistem inferior triunghiular cu $n = 4$.

Se observă că înălțimea grafului este dată de calea $S_0, T_{10}, S_1, T_{21}, \dots, S_{n-1}$, a cărei lungime este de $n + 2(n - 1) = 3n - 2$ operații. Lățimea grafului este de $n - 1$. Deci complexitatea paralelă a acestor doi algoritmi este $O(n)$, folosind $O(n)$ procesoare.

Dar complexitatea paralelă a problemei $TR(n)$ este mai bună, după cum se arată, de exemplu, în [4]. Se folosește următorul rezultat:

LEMĂ. Dacă A are diagonala unitate ($a_{ii} = 1$, ceea ce nu e o restricție; dacă elementele diagonale sunt diferite de 1, atunci se împarte cu ele întreaga linie, inclusiv elementul corespunzător din b), și descompunem $A = I - L$, cu L strict inferior triunghiulară (adică $l_{ij} = 0$, pentru $i \leq j$), atunci:

$$A^{-1} = I + L + L^2 + \dots + L^{n-1}.$$

Demonstrație (schită). Notând cu B termenul din dreapta al egalității de mai sus, se observă că $BA = I - L^n$. Se poate însă observa că L^2 are prima subdiagonală egală cu zero, L^3 are primele două subdiagonale zero, etc.; în concluzie $L^n = 0$. \square

Dar, după cum am văzut, o sumă se poate calcula cu $O(\log n)$ operații; calculul matricilor L, L^2, \dots, L^{n-1} este o problemă de tip prefix sum, și are complexitatea $O(\log^2 n)$, deoarece o înmulțire de matrici se poate calcula în $O(\log n)$, ceea ce se știe deja, iar calculul produselor parțiale comportă $O(\log n)$ etape la rândul său (vezi secțiunea dedicată sumei prefixelor); în fine, produsul $A^{-1}b$ se poate calcula în $O(\log n)$.

Soluția sistemului triunghiular se poate calcula așadar în timp $O(\log^2 n)$, cu $O(n^4)$ procesoare, adică de $O(n)$ ori (cât e nevoie la suma prefixelor) câte $O(n^3)$ (cât e nevoie la înmulțirea de matrici). Se pot obține și alte soluții cu aceeași complexitate, dar folosind $O(n^3)$ procesoare; ele pot fi găsite în [4].

În orice caz, acest algoritm optim din punct de vedere al timpului de execuție folosește un număr excesiv de mare de procesoare (deci este mai puțin eficient) și, în plus, este mai puțin stabil numeric. El este bun doar pentru determinarea teoretică a complexității problemei. Din punct de vedere practic, implementările curente se bazează pe algoritmi de tip substituție înainte.

6.2.3 Algoritmi cu comunicație globală

Pentru proiectarea unui algoritm paralel este importantă precizarea datelor care se folosesc la un moment dat (în bucla interioară a celor doi algoritmi secvențiali, în cazul nostru), precum și a datelor de intrare de care depinde un element al soluției. În ambii algoritmi, o necunoscută x_i depinde de elementul b_i corespunzător, dar această dependență nu implică decât inițializarea $x \leftarrow b$, deci nu are prea mare importanță. În bucla interioară a algoritmului cu sumă vectorială, pentru actualizarea elementelor x_i , cu $i > j$, e necesară coloana j a matricei A și soluția x_j (proaspăt calculată). În algoritmul cu produs scalar, pentru calculul soluției x_i se folosesc linia i a matricei A și soluțiile x_j , cu $j < i$.

Presupunem o repartizare pe linii sau pe coloane (cei doi algoritmi paraleli de un anume tip, rezultați din cei doi algoritmi secvențiali de bază, vor folosi repartizări diferite), fără a detalia însă, deocamdată, modul de repartizare; considerăm doar că fiecare procesor are $m = n/p$ linii (sau coloane) ale matricei A . Vom nota cu $P(k)$ adresa procesorului care deține linia (sau coloana, după caz) k . De asemenea, pentru un procesor oarecare, vom denumi mulțimea liniilor (coloanelor) care sunt locale acestuia *linii (coloane) proprii*. Elementele termenului liber b se repartizează la fel, indiferent dacă A se repartizează pe linii sau pe coloane, deoarece este vorba despre un vector; în orice caz, fiecare procesor va avea m elemente din b și e natural să stabilim că va poseda, în final, elementele corespunzătoare din x ; deci, $P(i)$ deține b_i și calculează x_i .

Cea mai simplă modalitate de a paraleliza algoritmi secvențiali prezentați este de a distribui procesoarelor calculele din bucla cea mai interioară, care se pot desfășura în paralel. În algoritmul cu sumă vectorială, pentru un anume j , acestea sunt operațiile T_{ij} , cu $i = j + 1 : n - 1$, adică cele de pe o "coloană" a grafului de precedentă din figura 6.3. Ele necesită coloana j a matricei A ; pentru a se executa calculele folosind doar date locale procesoarelor, coloana respectivă trebuie împărțită între procesoare; de aici concluzia că repartizarea naturală a matricei A este pe linii. Pe de altă parte, în toate calculele de tip T_{ij} apare x_j , care nu poate fi calculat decât de un singur procesor; pentru a comunica un x_j tuturor procesoarelor e deci necesară o difuzare. Prezentăm mai jos algoritmul, revenind după aceea cu alte explicații. Vom presupune, acum și în tot restul secțiunii dedicate rezolvării paralele a $TR(n)$, că vectorul x este inițializat cu b .

ALGORITM 6.10 ($TR(n)$ cu difuzare, pe arhitectură cu memorie distribuită, A repartizată pe linii)

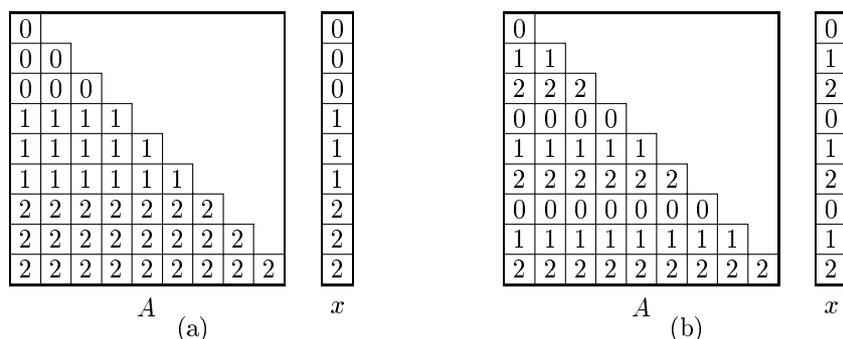


Figura 6.4: Repartizarea matricei inferioar triunghiulare A și a vectorului x , pentru $n = 9$ și $p = 3$, în cazurile: (a) bloc linii; (b) ciclic pe linii. Fiecare pătrat mic simbolizează un element, iar numărul din interiorul acestuia reprezintă adresa procesorului care deține elementul.

- pentru** $j = 0 : n - 1$
dacă $id = P(j)$ **atunci**
(1) $x_j \leftarrow x_j / a_{jj}$
(2) **difuzare:** $P(j)$ trimite x_j celorlalte procesoare
pentru $i = j + 1 : n - 1$
dacă $id = P(i)$ **atunci**
(3) $x_i \leftarrow x_i - a_{ij} x_j$

Prima parte a algoritmului, calculul unei soluții x_j , este efectuată de procesorul care deține inițial b_j , adică $P(j)$, în timp ce celelalte procesoare sunt inactice. Apoi, soluția x_j este comunicată de către $P(j)$ tuturor celorlalte procesoare, într-o operație de difuzare. În fine, în (3), fiecare procesor actualizează elementele proprii x_i , cu $i > j$, deci cele încă necalculate ale soluției, folosind valoarea x_j recepționată la difuzare și elementele a_{ij} de pe liniile care îi sunt locale.

Revenim acum la problema modului de repartizare pe linii. Dacă presupunem o repartizare bloc linii, cum este cea din figura 6.4a, atunci procesorul P_0 , care va deține liniile de la 0 la $m - 1$, își va termina complet treaba după primele m iterații, când va calcula elementele $x(0 : m - 1)$, adică exact cele de care se ocupă. Pentru $j \geq m$, calculele vor fi efectuate doar de $p - 1$ procesoare. În mod analog, după încă m iterații iese din cursă și P_1 , etc. Este clar că nu se poate spera o eficiență prea bună; chiar dacă instrucțiunile (3) sunt executate în paralel, doar o parte dintre procesoare sunt active; în figură, pe o coloană a matricei, numerele reprezintă adresele procesoarelor care execută în paralel operații de tip (3), pentru un j fixat. În plus, deși fiecare procesor deține m linii, numărul de elemente nenule este diferit pentru fiecare, matricea A fiind triunghiulară; de exemplu, P_0 are $m(m + 1)/2$ elemente,

în timp ce P_{p-1} are aproximativ mn , adică mult mai multe; deci, repartizarea, deși aparent echitabilă, nu este deloc așa.

Este de dorit, așadar, ca fiecare procesor să aibă aproximativ același număr de elemente (nenule) din *fiecare* coloană a matricei A , pentru a contribui în aceeași măsură la calculele de tip (3). Aceasta se poate realiza simplu printr-o repartizare ciclică pe linii, cum este cea din figura 6.4b; deci, $P(j) = j \bmod p$. Mai mult, spre deosebire de cazul anterior, acum fiecare procesor are și aproximativ același număr de elemente din A .

Pentru a evalua timpul de calcul al algoritmului 6.10, considerăm timpii pentru cele trei faze descrise:

- $n\alpha$ – timpul pentru calculul (secvențial) al soluțiilor x_j de către fiecare procesor.
- $\approx c_d n(\sigma + \beta)$ – timpul pentru cele $n - 1$ operații de difuzare (pentru $j = n - 1$ nu mai este necesară difuzarea); coeficientul c_d depinde de topologia arhitecturii și de algoritmul de difuzare ales; de exemplu, pentru hipercub $c_d = \log p$, iar pentru inel $c_d = p/2$.
- pentru actualizări (paralelism complet): $\frac{1}{p} \sum_{j=0}^{n-1} 2(n-1-j)\alpha \approx \frac{1}{p}(n^2 - n)\alpha$.

Timpul total este suma celor de mai sus:

$$T_{6.10} = \frac{1}{p}(n^2 + np - n)\alpha + c_d(n-1)(\sigma + \beta).$$

Precizăm că aceasta este o limită superioară a timpului de calcul; practic, pot apare valori mai mici, deoarece în timpul difuzării un procesor nu este permanent ocupat; deci, de îndată ce termină partea sa din operația de difuzare, el poate trece imediat la actualizări.

Pentru algoritmul cu produs scalar, executarea în paralel a operațiilor din bucla **pentru** cea mai interioară, înseamnă cooperarea procesoarelor pentru calculul sumei $\sum_{j=0}^{i-1} a_{ij}x_j$; deci, pentru un i fixat, suma implică elemente de pe linia i a matricei A , și elementele din x deja calculate. Pentru a repartiza echilibrat calculele, e normal să presupunem linia i împărțită aproximativ egal între procesoare, ceea ce semnifică o repartizare pe coloane a matricei A , și anume, din motive asemănătoare celor de la algoritmul precedent, ciclic pe coloane; repartizarea lui x nu pune probleme suplimentare, deoarece procesorul care deține a_{ij} , adică $P(j)$, deține și x_j . Deci, fiecare procesor calculează $\sum_{j < i} a_{ij}x_j$, cu j parcurgând coloanele proprii. Pentru a obține x_i ca în formula 6.6, este necesară adunarea acestor sume parțiale, după cum se procedează în algoritmul următor:

ALGORITM 6.11 ($TR(n)$ cu colectare, pe arhitectură cu memorie distribuită, A repartizată ciclic pe coloane)

```

pentru  $i = 0 : n - 1$ 
     $t \leftarrow 0$ 
    pentru  $j = 0 : i - 1$ 
        dacă  $id = P(j)$  atunci
(1)          $t \leftarrow t + a_{ij}x_j$ 
(2)      $s \leftarrow$  colectare cu sumare a valorilor  $t$  în  $P(i)$ 
        dacă  $id = P(i)$  atunci
(3)          $x_i \leftarrow (b_i - s)/a_{ii}$ 

```

Pentru un i fixat, calculele de tip (1) decurg în perfect paralelism; este vorba de operațiile T_{ij} de pe "linia" i a grafului de dependență din figura 6.3, care sunt independente; datorită repartizării ciclice, pentru orice i , toate procesoarele au aproximativ același număr de operații de efectuat. Toate sumele parțiale trebuie adunate, iar rezultatul comunicat procesorului $P(i)$; aceasta se poate realiza eficient printr-o operație de colectare cu sumare, așa cum am văzut în algoritmul 4.3, pentru hiper-cub, sau în problema 4.3.7, pentru inel. În sfârșit, în instrucțiunea (3), procesorul $P(i)$ calculează soluția x_i , timp în care celelalte procesoare sunt inactice (în realitate, celelalte procesoare continuă cu calculele (1) din iterația următoare, dar vor fi nevoite să-l aștepte pe $P(i)$ la următoarea colectare; până la urmă, timpul de inactivitate este același, chiar dacă în alt moment).

Timpul de execuție al acestui algoritm este aproximativ același cu cel al precedentului, singura diferență putând apare la complexitatea comunicației; dar difuzarea și colectarea sunt operații duale, deci vor consuma un timp identic; faptul că se face sumare în timpul colectării nu este semnificativ.

În concluzie, pentru ambii algoritmi prezentați, 6.10 și 6.11, se folosește comunicație globală (de unde și numele acestei secțiuni). În expresia $T_{6.10}$, partea reflectând calculele aritmetice este $O(n^2/p)$, iar cea descriind complexitatea comunicației este $O(pn)$, în cazul cel mai rău, al inelului; deci, pentru $n \rightarrow \infty$, ponderea comunicațiilor se devine nesemnificativă, și eficiența algoritmilor tinde către 1. Totuși, pentru valori uzuale ale lui n , termenul reflectând comunicația este important (să nu uităm că $\sigma \gg \alpha$), ceea ce va conduce la o creștere lentă a eficienței, pe măsură ce dimensiunea problemei crește.

Cauzele relativei ineficiențe sunt două. În primul rând, procesoarele participă la operații de comunicație globală, care implică sincronizare, și în care pierd timp în așteptarea mesajelor; în plus, există operații efectuate de un singur procesor, în timp ce celelalte sunt inactice (tocmai din cauza necesității sincronizării). În al doilea rând, numărul acestor comunicații globale este foarte mare, deoarece toate mesajele au dimensiune 1; algoritmi au granularitate fină; aceasta implică o pondere egală a coeficienților σ și β în timpul de comunicație.

Vom vedea în continuare cum se poate mări dimensiunea mesajelor, și deci micșora numărul lor, păstrând structura algoritmilor de mai sus. Altfel zis, vom deduce

algoritmi cu granularitate mare (la nivel de bloc de matrice); pentru a nu lungi prezentarea, ne vom ocupa doar de algoritmul cu difuzare.

Să presupunem că A este partiționată în blocuri pătrate ca în prima parte din relația (6.1); cum A este inferior triunghiulară, atunci $A_{ij} = 0_r$, pentru $i < j$, iar A_{jj} este inferior triunghiulară. De asemenea, ca în secțiunea referitoare la produsul matrice-vector, vom presupune vectorii x și b împărțiți fiecare în n_b vectori de dimensiune r , pe care-i vom nota cu X_i , respectiv B_i ($i \in 0 : n_b - 1$). Cu aceste notații, echivalentul relației (6.6), dedus acum dintr-o bloc linie a ecuației $Ax = b$, este:

$$A_{ii}X_i = B_i - \sum_{j=0}^{i-1} A_{ij}X_j,$$

și arată că X_i poate fi calculat prin rezolvarea unui sistem inferior triunghiular "mic", de dimensiune r , după calculul tuturor necunoscutelor din vectorii X_j , cu $j < i$.

Algoritmul secvențial asemănător cu 6.8 se deduce și el imediat; forma este aproape identică, numai că operațiile sunt acum la nivel de bloc, și nu de element.

ALGORITHM 6.12 ($TR(n)$ secvențial la nivel de bloc, prin sumă vectorială)

```

 $x \leftarrow b$ 
pentru  $j = 0 : n_b - 1$ 
    rezolvă sistemul inferior triunghiular cu necunoscuta  $Y$ :  $A_{jj}Y = X_j$ 
     $X_j \leftarrow Y$ 
    pentru  $i = j + 1 : n_b - 1$ 
         $X_i \leftarrow X_i - A_{ij}X_j$ 

```

Desigur că numărul de operații este același, adică n^2 .

Pentru a deduce echivalentul la nivel de bloc al algoritmului cu difuzare 6.10, trebuie să stabilim repartizarea matricei A . Deoarece se paralelizează bucla **pentru** cea mai interioară, adică operațiile $X_i \leftarrow X_i - A_{ij}X_j$, cu $i > j$ (j este fixat), este natural ca bloc coloana j a matricei A să fie distribuită echilibrat procesoarelor; aceasta înseamnă o repartizare pe linii, *bloc ciclică*; fiecare procesor va avea în memoria locală $m = n_b/p = n/(rp)$ bloc linii, adică n/p linii; notând (doar acum) cu $P(i)$ procesorul care deține blocul A_{ij} , și care calculează X_i , se obține:

ALGORITHM 6.13 ($TR(n)$ cu difuzare, la nivel de bloc, pe arhitectură cu memorie distribuită, A repartizată bloc ciclic pe linii)

```

pentru  $j = 0 : n_b - 1$ 
    dacă  $id = P(j)$  atunci
(1)        rezolvă sistemul inferior triunghiular cu necunoscuta  $Y$ :  $A_{jj}Y = X_j$ 
(1')        $X_j \leftarrow Y$ 
(2)       difuzare:  $P(j)$  trimite  $X_j$  celorlalte procesoare
    pentru  $i = j + 1 : n_b - 1$ 
        dacă  $id = P(i)$  atunci
(3)        $X_i \leftarrow X_i - A_{ij}X_j$ 

```

Din nou forma este foarte asemănătoare cu cea a algoritmului la nivel de element 6.10; numărul de operații este, însă, diferit de cel din expresia $T_{6.10}$. Defalcăt pe cele trei faze, timpul de execuție este:

- $n_b r^2 \alpha = nr\alpha$ pentru calculul (secvențial, deoarece un singur procesor lucrează la (1), restul așteptând difuzarea) al soluțiilor X_j de către fiecare procesor, pentru că se rezolvă n_b sisteme triunghiulare de dimensiune r .
- $\approx c_d(n_b\sigma + n\beta) = c_d((n/r)\sigma + n\beta)$ – pentru cele n_b difuzări, fiecare a unui mesaj de dimensiune r .
- timpul pentru actualizările (3) rămâne același ca în $T_{6.10}$.

Rezumând, timpul primei faze este de r ori mai mare ca în $T_{6.10}$, iar coeficientul lui σ din timpul de difuzare se diminuează de r ori. Cum $\sigma \gg \alpha$, aceasta conduce la reducerea timpului total. Totuși, timpul nu se îmbunătățește considerabil, deoarece am înlăturat doar una din cauzele ineficienței, cea cu pondere mai mică.

6.2.4 Algoritmi în front de undă

În algoritmi din secțiunea anterioară, la un moment dat, toate procesoarele execută operații de același tip, cooperând fie la actualizarea cu elementul x_j în algoritmul cu difuzare, fie la calculul sumelor parțiale necesare calculului unui element x_i în algoritmul cu colectare. O altă sursă de paralelism, exploatată în cele ce urmează, este lucrul simultan la mai multe elemente ale soluției, sau, echivalent, prelucrarea pipeline a unui singur element. Cu alte cuvinte, dacă în algoritmul cu difuzare se executau în paralel operațiile dintr-o "coloană" a grafului de precedență din figura 6.3, iar în algoritmul cu colectare cele dintr-o "linie", putem acum încerca combinarea acestor abordări. Evident, există mai multe posibilități, dintre care vom expune câteva în continuare.

Ne vom ocupa întâi de varianta corespunzătoare algoritmului secvențial cu sumă vectorială. Vom presupune acum o repartizare ciclic pe coloane a matricei A . După ce procesorul $P(j)$ calculează x_j , actualizarea termenilor x_i , cu $i > j$, nu poate fi împărțită cu alte procesoare deoarece $P(j)$ deține în întregime coloana j din A ; dacă toată actualizarea ar fi făcută de $P(j)$, atunci execuția algoritmului ar fi pur secvențială (deși distribuită), pentru că celelalte procesoare nu ar avea nimic de făcut. Pentru a obține paralelismul, se împarte actualizarea în mai multe etape.

Introducem un vector z , de dimensiune n , inițializat cu zero, în care se vor acumula actualizările; deci, în final, $z_i = \sum_{j=0}^{i-1} a_{ij}x_j$. De data aceasta nu elementele din x vor circula (ca în algoritmul 6.10, unde se difuzau), ci elementele din z . Acest vector se împarte în mai multe segmente, de dimensiune s ; evident, $1 \leq s \leq n$; în fiecare etapă $P(j)$ actualizează un segment, după care trimite acest segment către $P(j+1)$. În acest fel, $P(j+1)$ poate calcula x_{j+1} imediat după recepționarea primului segment, prelucrând apoi segmentele pe măsură ce le primește de la $P(j)$.

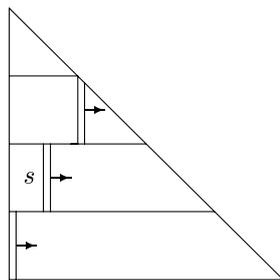


Figura 6.5: Poziția segmentelor la un moment dat, în algoritmul front de undă pe coloane.

Să detaliem. Procesorul $P(0)$ calculează x_0 , apoi începe calculul componentelor $z_i = a_{i0}x_0$, cu $i > 0$, ale vectorului de actualizare. După ce a calculat primele s astfel de elemente, adică primul segment, $P(0)$ le trimite lui $P(1)$, astfel încât acesta să-l poată calcula pe $x_1 \leftarrow (b_1 - z_1)/a_{11}$ și să-și înceapă actualizările, adică operații de genul $z_i \leftarrow z_i + a_{i1}x_1$, cu indicele i corespunzător liniilor primului segment. În acest timp, $P(0)$ continuă lucrul la următoarele s elemente ale lui z , al doilea segment. Pe scurt, după ce un procesor termină actualizarea unui segment, îl trimite mai departe, în așa fel încât, la un moment dat, poziția (spațială) a segmentelor este cea sugerată în figura 6.5, iar mișcarea acestora este asemănătoare înaintării unui front de undă, de unde și denumirea algoritmului.

Un exemplu este prezentat în figura 6.6, sub forma unei diagrame Gantt. Am notat cu S_j operația $x_j \leftarrow (b_j - z_j)/a_{jj}$ și cu T_{ij} operația $z_i \leftarrow z_i + a_{ij}x_j$, în concordanță cu graful de precedență din figura 6.3; cele două tipuri de operații au același timp de execuție. Timpii morți ai procesoarelor sunt figurați prin spații albe. Pentru a nu complica diagrama, am presupus comunicațiile instantanee; liniile verticale marchează momentele în care are loc comunicația. Alegând $s = 2$, un procesor execută asupra unui segment două, sau una, operații de tip T_{ij} .

Deși exemplul nu o sugerează pregnant, pentru $n \gg p$, în funcție de dimensiunea segmentului și de constantele specifice calculatorului gazdă (raportul timp comunicație/timp calcule, în special), este posibil ca toate procesoarele să fie active la un moment dat.

Vom prezenta în continuare o formă simplificată a algoritmului și apoi unele detalii. Variabila locală *segment* este o mulțime conținând cel mult s elemente ale vectorului de actualizare z ; un procesor nu lucrează la un moment decât asupra unui singur segment; *nr_seg* este numărul de segmente care trebuie prelucrate pentru un anume j . Deoarece comunicațiile au loc întotdeauna de la $P(j)$ la $P(j + 1)$ și repartizarea matricei A este ciclică, topologia minimală cerută de algoritm este inelul, pe care se comunică într-un singur sens.

P_0	S_0	T_{10}	T_{20}	T_{30}	T_{40}	T_{50}	T_{60}	T_{70}		S_4	T_{54}	T_{64}	T_{74}			
P_1				S_1	T_{21}	T_{31}	T_{41}	T_{51}	T_{61}	T_{71}			S_5	T_{65}	T_{75}	
P_2						S_2		T_{32}	T_{42}	T_{52}	T_{62}	T_{72}			S_6	T_{76}
P_3									S_3	T_{43}	T_{53}	T_{63}	T_{73}			S_7

Figura 6.6: Diagramă Gantt pentru execuția algoritmului în front de undă pe coloane, pentru $n = 8$, $p = 4$, $s = 2$.

ALGORITM 6.14 ($TR(n)$ – în front de undă pe coloane, pe inel, A repartizată ciclic pe coloane)

```

pentru  $j = 0 : n - 1$ 
  dacă  $id = P(j)$ 
    pentru  $k = 1 : nr\_segm$                                  $\{nr\_segm = \lceil (n - j)/s \rceil\}$ 
      dacă  $id = 0$  și  $j = 0$  atunci
        (1) crează  $segment$                                  $\{z((k - 1)s + 1 : ks) \leftarrow 0\}$ 
      altfel
        (2) recv( $segment$ , stânga)                         $\{de\ la\ P(j - 1)\}$ 
      dacă  $k = 1$ 
        (3)  $x_j \leftarrow (b_j - z_j)/a_{jj}$ 
        (4)  $segment \leftarrow segment - \{z_j\}$ 
      pentru  $z_i \in segment$ 
        (5)  $z_i \leftarrow z_i + a_{ij}x_j$ 
      dacă  $|segment| > 0$  atunci
        (6) send( $segment$ , dreapta)                         $\{c\grave{a}tre\ P(j + 1)\}$ 

```

Toate segmentele sunt create, în (1), de către P_0 , la prima iterație, și transmise permanent către dreapta. Numărul de segmente pe care le recepționează un procesor la un moment dat se calculează cu $nr_segm = \lceil (n - j)/s \rceil$, deoarece pe coloana j sunt $n - j$ elemente nenule în A . După recepția primului segment corespunzător unei coloane j , procesorul $P(j)$ poate calcula x_j , în (3) (care e echivalentul formulei (6.6)), după care z_j nu mai este necesar, deci poate fi eliminat din segment. Se observă deci că primul segment are cu câte un element mai puțin la fiecare necunoscută calculată și dispare după primele s etape, moment în care procesul se repetă cu al doilea segment (devenit acum primul). Celelalte valori din primul segment și celelalte segmente sunt actualizate în (5), după care segmentele sunt trimise către $P(j + 1)$, dacă sunt nevide. Se poate modifica algoritmul astfel încât instrucțiunile (6), dintr-o iterație, și (2), din următoarea, să se execute în paralel.

În acest fel segmentele parcurg inelul, fiecare până la dispariție; ultimul segment, de exemplu, este transmis de $n - 1$ ori; se poate aprecia, în general, că volumul comunicației este mare; totuși, nu trebuie uitat că multe din transmisiile pot avea loc în paralel. Dacă s , dimensiunea unui segment, este mic, crește gradul de paralelism, dar și timpul de comunicație, deoarece numărul segmentelor este mare; dacă s este mare, numărul de mesaje se reduce, dar scade și paralelismul. Faptul că volumul comunicației scade pe măsură ce calculele avansează compensează într-un fel trecerea de mai multe ori a unui segment pe la același procesor. De asemenea, se poate aprecia că în cazul comunicațiilor asincrone timpul de execuție poate fi redus, chiar semnificativ, față de cazul sincron.

Timpul de execuție al acestui algoritm este următorul (pentru detalii vezi [14]):

$$T_{6.14} = \frac{1}{p} \left(n^2 + np + \frac{s(s-1)}{2} p^2 \right) \left(\alpha + \frac{\sigma}{s} + \beta \right),$$

și se minimizează în funcție de s ; optimul se obține dintr-o ecuație de gradul 3 în s , ai cărei coeficienți depind de caracteristicile calculatorului gazdă. Deci, pentru fiecare calculator și pentru fiecare n , trebuie inițial calculată valoarea optimă a lui s . Calculul nu trebuie făcut la fiecare execuție, ci se poate crea, o dată, un tabel cu valorile optime, din care se extrage la execuție valoarea adecvată.

Varianta în front de undă a algoritmului cu produs scalar presupune o repartizare pe linii a matricei A . Procesorul $P(i)$ trebuie să calculeze singur produsul scalar ce implică linia i a matricei; dacă ar executa acest calcul în întregime, atunci calculele s-ar desfășura pur secvențial. Se folosește din nou ideea împărțirii în segmente, de data aceasta a calculului produsului scalar.

Având în vedere asemănarea pronunțată cu algoritmul în formă de undă orientat pe coloane, ne vom mulțumi cu o scurtă expunere comparativă. Dacă în acel algoritm circulau valorile lui z , care se actualizau pe baza valorilor locale ale lui x , acum lucrurile se petrec pe dos; pentru calculul $z_i = \sum_{j=0}^{i-1} a_{ij} x_j$, procesorul $P(i)$ are nevoie doar de valorile x_j , linia i din A fiindu-i proprie; deci ideea este de a face să circule segmente de dimensiune s , conținând elemente din x , în timp ce elementele din z rămân locale. Pentru a calcula x_i , procesorul $P(i)$ primește astfel de segmente, actualizând z_i , până când este obținută suma de mai sus; apoi poate efectiv determina x_i , pe care-l atașează ultimului segment; toate segmentele sunt trimise spre $P(i+1)$, pe rând, astfel încât un procesor lucrează cu un singur segment la un moment dat. Poziția segmentelor la un moment dat este cea sugerată în figura 6.7.

În algoritmul de mai jos, *segment* este o variabilă locală care conține cel mult s elemente din x ; suma z_i poate fi găzduită de o variabilă scalară; am păstrat notația numai pentru a nu strica asemănarea cu algoritmul pe coloane.

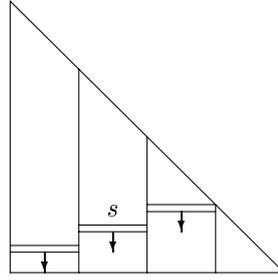


Figura 6.7: Poziția segmentelor la un moment dat, în algoritmul front de undă pe linii.

ALGORITHM 6.15 ($TR(n)$ – în front de undă pe linii, pe inel, A repartizată ciclic pe linii)

```

pentru  $i = 0 : n - 1$ 
     $z_i \leftarrow 0$ 
    dacă  $id = P(i)$  atunci
        pentru  $k = 1 : nr\_segm - 1$             $\{nr\_segm = \lceil i/s \rceil\}$ 
            (1) rcv( $segment$ , stânga)
                în paralel
                (2) send( $segment$ , dreapta)        $\{c\grave{a}tre P(i + 1)\}$ 
                    pentru  $x_j \in segment$ 
                        (3)  $z_i \leftarrow z_i + a_{ij}x_j$ 
                        (4) rcv( $segment$ , stânga)        $\{dac\grave{a} i \bmod s = 0\}$ 
                            pentru  $x_j \in segment$ 
                                (5)  $z_i \leftarrow z_i + a_{ij}x_j$ 
                                (6)  $x_i \leftarrow (b_i - z_i)/a_{ii}$ 
                                (7)  $segment \leftarrow segment \cup \{x_i\}$             $\{sau\ creaz\grave{a}\ un\ nou\ segment\}$ 
                                (8) send( $segment$ , dreapta)

```

Trebuie remarcate câteva detalii de implementare. Pentru un i fixat, tratarea segmentelor e despărțită în două etape, ultimul segment fiind prelucrat separat. Pentru primele, transmisia (2) are loc, în mod ideal, imediat după recepția (1), deoarece conținutul segmentului nu se modifică, deci e bine ca $P(i + 1)$ să-l primească cât mai devreme; pe de altă parte, pentru a nu se aștepta efectuarea acestei comunicații, actualizările (3) au loc în paralel (sau concurrent, dacă arhitectura procesorului nu permite altfel). Pentru ultimul segment se execută recepția, în (4), dar numai dacă este cazul, adică dacă nu trebuie creat un nou segment; acum actualizările (5) preced transmisia (8), deoarece trebuie calculat x_i , în (6), și adăugat segmentului, în (8). Se observă că numărul de segmente este acum în continuă creștere, până când se ajunge la numărul maxim de segmente.

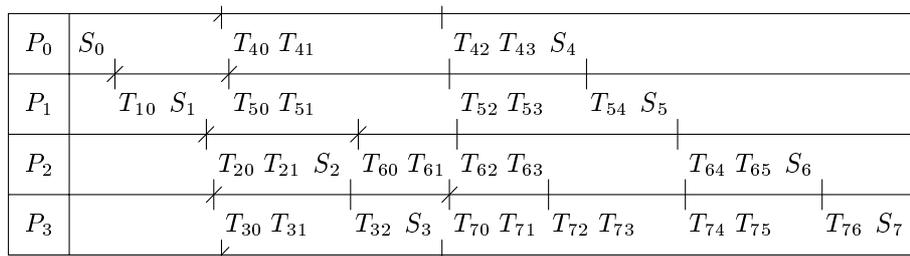


Figura 6.8: Diagramă Gantt pentru execuția algoritmului în front de undă pe linii, pentru $n = 8$, $p = 4$, $s = 2$.

Diagrama Gantt a unui exemplu este prezentată în figura 6.8, pentru aceiași parametri ca în cazul algoritmului în front de undă pe coloane. De data aceasta comunicația este simbolizată printr-o linie între zonele alocate unor procesoare vecine; am presupus că durata comunicațiilor este nenulă, dar foarte mică, pentru a sugera mai clar ordinea transmisiilor. O mică linie oblică marchează transmisiile primului segment. Comunicația a fost presupusă sincronă; pentru a nu complica diagrama, am presupus că transmisia (2) și actualizările (3) se pot efectua în orice ordine (cea mai avantajoasă), dar secvențial.

Timpul de execuție pentru acest algoritm are o expresie complicată, însă cu o alură asemănătoare celui pentru algoritmul în front de undă orientat pe coloane, adică cu un minim pentru $1 \leq s \leq n$.

6.2.5 Algoritmi ciclici

Algoritmii în front de undă au dezavantajul unui volum de comunicație relativ mare. Cei care vor fi prezentați în continuare vor folosi minimum de comunicație, atât ca număr de mesaje, cât și ca număr de elemente de matrice transmise.

Vom începe cu varianta paralelă a algoritmului cu sumă vectorială. Se presupune o repartizare ciclic pe coloane a matricei A . Există o oarecare asemănare cu algoritmul în front de undă corespondent, în sensul că actualizările se acumulează într-un vector z , de dimensiune n , inițializat cu zero și care, în final, conține valorile $z_i = \sum_{j=0}^{i-1} a_{ij}x_j$, pe baza cărora se calculează elementele din x . De data aceasta, însă, vectorul z nu se mai împarte în mai multe segmente; există un singur pachet (segment), de dimensiune $p - 1$, care circulă de-a lungul unui inel, conținând o parte din elementele lui z , sub forma unor sume parțiale.

Vom prezenta întâi algoritmul, și apoi unele explicații; în ce privește notațiile, SUM este pachetul de $p - 1$ elemente care circulă, și pentru care fiecare procesor trebuie să rezerve spațiu; de asemenea, fiecare procesor are o variabilă locală z , de dimensiune n (sunt necesare doar $n - p + 1$ valori, dar nu ne complicăm), în care

calculează partea sa din suma de mai sus; pentru toate procesoarele, vectorii SUM și z se inițializează cu 0. Algoritmul este următorul:

ALGORITM 6.16 ($TR(n)$ – ciclic pe coloane, pe inel, A repartizată ciclic pe coloane)

```

pentru  $j = 0 : n - 1$ 
  dacă  $id = P(j)$ 
    (1) recv( $SUM$ , stânga) {de la  $P(j - 1)$ }
    (2)  $z_j \leftarrow z_j + SUM(0)$ 
    (3)  $x_j \leftarrow (b_j - z_j)/a_{jj}$ 
    (4)  $SUM(0 : p - 3) \leftarrow SUM(1 : p - 2)$  {elimină  $SUM(0)$  din pachet}
    pentru  $i = 1 : p - 2$ 
      (5)  $SUM(i - 1) \leftarrow SUM(i - 1) + z_{j+i} + a_{j+i,j}x_j$ 
      (6)  $SUM(p - 1) \leftarrow z_{j+p-1} + a_{j+p-1,j}x_j$ 
      (7) send( $SUM$ , dreapta) {către  $P(j + 1)$ }
    pentru  $i = j + p : n - 1$ 
      (8)  $z_i \leftarrow z_i + a_{ij}x_j$ 

```

Deci, pachetul SUM circulă de-a lungul inelului de procesoare, fiecare procesor primindu-l din stânga în (1) și retransmițându-l spre dreapta în (7), după ce a făcut toate modificările necesare. Pachetul conține actualizări pentru următoarele $p - 1$ elemente din x de calculat; deoarece o linie a matricei A este repartizată între toate procesoarele, fiecare dintre acestea contribuie la calculul sumei $z_i = \sum_{j=0}^{i-1} a_{ij}x_j$ cu produsele $a_{ij}x_j$ pe care le poate efectua cu date din memoria locală; în loc de a transmite tot vectorul z , se poate folosi doar SUM , care la plecarea dintr-un procesor conține informația necesară celorlalte $p - 1$ procesoare, care e suficientă până la revenirea pachetului la procesorul curent. Instrucțiunea (2) reprezintă obținerea valorii finale a lui z_j , după ce în $SUM(0)$ s-au acumulat actualizările efectuate de celelalte procesoare. Deci, în (3), $P(j)$ calculează necunoscuta x_j folosind chiar formula (6.6). Apoi, în (4), elimină prima valoare din SUM (în algoritm, printr-o deplasare la stânga a vectorului SUM), care nu mai este necesară mai departe; actualizează, în (5), celelalte valori din SUM cu necunoscuta proaspăt calculată și adaugă, în (6), la pachetul SUM suma sa parțială (actualizată și cu necunoscuta curentă) destinată calculului variabilei x_{j+p-1} – care va fi calculată de $P(j + p - 1)$, vecinul din stânga al lui $P(j)$, dar cel mai depărtat în sensul de circulație al pachetului SUM . În acest moment pachetul SUM este complet prelucrat și poate fi expediat vecinului din dreapta. În fine, în (8), $P(j)$ actualizează cu necunoscuta proaspăt calculată sumele parțiale z pentru celelalte necunoscute (de la $j + p$ încolo), operație care ocupă, în general, mult mai mult timp decât toate cele precedente; doar pentru j suficient de mare, adică apropiat de n , aceste actualizări sunt de scurtă durată.

Se poate demonstra că, pentru acest algoritm, comunicația este minimală, atât ca număr de pachete cât și ca număr de elemente transmise. Aceasta nu înseamnă că algoritmul e optim ca timp total de execuție, pentru că apare altă problemă; procesorul $P(j)$ poate termina actualizările (8) ale vectorului z înainte de a reprimi

un nou pachet *SUM*; în acest caz, el nu are nimic de făcut și va trebui să aștepte; timpul mort astfel acumulat poate avea o pondere însemnată – cu atât mai mare cu cât raportul n/p este mai mic; acest lucru este absolut normal: dacă p este mare, un procesor va lucra mai mult la actualizarea pachetului *SUM*, care va circula astfel mai lent, în timp ce actualizarea vectorului z va dura relativ puțin.

Dacă $n > p(t_p + p)$ (unde t_p este costul, măsurat în flop, al transmisiei a $p - 1$ elemente în virgulă mobilă către un vecin¹) – ceea ce înseamnă o valoare destul de mare, atunci timpul de execuție al algoritmului va avea aproximativ expresia:

$$T_{6.16} = \left[\frac{1}{p}(n^2 + np) + p(t_p + p)^2 \right] \alpha.$$

Eficiența algoritmului tinde asimptotic spre 1, dar se apropie de această valoare numai pentru valori mari ale lui n . Se arată în [20] că o eficiență teoretică de aproximativ 50% se obține pentru $n = p(t_p + p)$.

Un algoritm asemănător se poate stabili și pentru o repartizare pe linii a matricei A . Ca și la algoritmul în front de undă pe linii, va circula x ; ca și la algoritmul ciclic pe coloane, va exista un singur pachet, numit acum *XSUB* (notația e consacrată) și conținând ultimele $p - 1$ necunoscute calculate. Pentru a păstra stilul de prezentare, am folosit și variabila z , cu semnificația de până acum; trebuie observat că z și x pot fi găzduite de aceeași variabilă locală și că z este repartizat ciclic între procesoare ($P(i)$ are elementul z_i). Algoritmul este următorul:

ALGORITM 6.17 (*TR*(n) – ciclic pe linii, pe inel, A repartizată ciclic pe linii)

```

pentru  $i = 0 : n - 1$ 
  dacă  $id = P(i)$ 
    (1)   recv(XSUB( $1 : p - 1$ ), stânga)
          pentru  $j = 1 : p - 1$ 
    (2)      $z_i \leftarrow z_i + a_{i,i-p+j} \cdot \textit{XSUB}(j)$ 
    (3)      $x_i \leftarrow (b_i - z_i)/a_{ii}$ 
    (4)     XSUB( $0 : p - 2$ )  $\leftarrow$  XSUB( $1 : p - 1$ )
    (5)     XSUB( $p - 1$ )  $\leftarrow$   $x_i$    {XSUB  $\leftarrow$  XSUB - { $x_{i-p+1}$ }  $\cup$  { $x_i$ } }
    (6)     send(XSUB( $1 : p - 1$ ), dreapta)
          pentru  $k = i + 1 : n - 1$ 
            dacă  $id = P(k)$  atunci
              pentru  $j = 1 : p$ 
    (7)        $z_k \leftarrow z_k + a_{k,i-p+j} \cdot \textit{XSUB}(j - 1)$ 

```

Am folosit un mic artificiu tehnic în utilizarea pachetului *XSUB*; cele $p - 1$ elemente transmise sunt cele din pozițiile $1 : p - 1$; vom explica puțin mai târziu la ce folosește poziția 0. Un procesor $P(i)$, după ce recepționează în (1) pachetul *XSUB*, completează suma parțială z_i în (2), pe baza soluțiilor aflate în *XSUB* (calculate în

¹Cu alte cuvinte $t_p = (\sigma + (p - 1)\beta)/\alpha$.

ultima tură de celelalte $p - 1$ procesoare), după care calculează x_i folosind clasică formulă (6.6). Apoi, în (4) și (5), șterge din pachetul $XSUB$ cea mai veche componentă a soluției și adaugă la acesta necunoscuta proaspăt calculată x_i ; după care, în (6), trimite $XSUB$ mai departe. Apoi actualizează elementele locale (de pe liniile care îi aparțin) ale vectorului z , folosind valorile din $XSUB$; în algoritm am ținut seama de deplasarea spre stânga din (4) a elementelor din $XSUB$; de aceea, în actualizările (8), apare și $XSUB(0)$, care nu trebuia trimis spre $P(i + 1)$, dar e încă necesar pentru actualizările locale.

Acest algoritm are proprietăți similare cu cele ale algoritmului ciclic pe coloane.

6.2.6 Comparații ale algoritmilor precedenți

Am prezentat până acum trei perechi de algoritmi paraleli de rezolvare a sistemelor liniare triunghiulare. Dintre aceștia, algoritmi cu comunicație globală sunt, teoretic, cei mai puțin performanți. Cel mai bun algoritm depinde de raportul n/p ; dacă acesta este mic, algoritmi în front de undă sunt superiori; dacă este mare (număr mic de procesoare față de dimensiunea problemei), atunci algoritmi ciclici sunt mai buni.

De exemplu, în [14] se raportează că, pe un hiper-cub de tip Ncube, pentru $n = 500$, pentru $p = 4, 8, 16$ algoritmi ciclici au avut cel mai bun timp de execuție; în schimb, pentru $p = 32, 64$ algoritmi în front de undă au fost superiori.

Se poate deci concluziona că, din punct de vedere practic, pentru rezolvarea paralelă a unui sistem liniar triunghiular se va alege unul dintre algoritmi în front de undă sau ciclic, cu repartizare a datelor pe linii sau pe coloane după cum e mai convenabil în context. Alegerea se face în funcție de valoarea n/p și de caracteristicile calculatorului pe care se face implementarea (în special raportul timp de comunicație/timp de calcul). Repartizarea pe linii sau pe coloane se alege în funcție de contextul în care se rezolvă $TR(n)$.

Nu trebuie neglijați, totuși, nici algoritmi cu comunicație globală, în special variantele la nivel de bloc (de genul algoritmului 6.13), datorită simplității lor și a folosirii unor operații de comunicație care sunt deseori gata implementate cu mare eficiență.

6.2.7 Un algoritm pe tor

Vom prezenta în continuare o modalitate de rezolvare a problemei $TR(n)$ pe un tor de dimensiuni egale. Primul lucru care trebuie stabilit este repartizarea matricei A ; de data aceasta, datorită "asemănării" între un tor și o matrice, există o soluție simplă și unică; profitând și de experiența anterioară, vom considera o repartizare ciclică atât pe linii, cât și pe coloane; altfel spus, elementul a_{ij} va fi deținut de procesorul de pe linia $i \bmod \sqrt{p}$ și coloana $j \bmod \sqrt{p}$ ale torului. Un exemplu este prezentat în figura 6.9; se observă că procesoarele au aproximativ același număr de elemente din A . În ce privește vectorii b și x , din motive ce vor deveni clare mai târziu, îi presupunem repartizați ciclic doar procesoarelor de pe diagonala torului; deci b_i aparține procesorului P_{tt} , cu $t = i \bmod \sqrt{p}$; același procesor va deține în final x_i .

00										00
10	11									11
20	21	22								22
00	01	02	00							00
10	11	12	10	11						11
20	21	22	20	21	22					22
00	01	02	00	01	02	00				00
10	11	12	10	11	12	10	11			11
20	21	22	20	21	22	20	21	22		22

A
 x

Figura 6.9: Repartizare a matricei inferior triunghiulare A și a vectorului x pe un tor 3×3 , pentru $n = 9$.

Să încercăm să deducem un algoritm care să combine ideile algoritmilor cu difuzare 6.10 și cu colectare 6.11. E natural să dorim aceasta, din moment ce repartizarea pe tor îmbină repartizările din cei doi algoritmi amintiți, pe linii, respectiv pe coloane. Grupând procesoarele torului pe linii, algoritmul cu difuzare s-ar putea generaliza astfel: după calculul unui element x_j , grupul de procesoare $P(j)$ îl difuzează celorlalte grupuri, care efectuează actualizări. Pe de altă parte, grupând procesoarele pe coloane, se poate generaliza algoritmul cu colectare: la iterația i , fiecare grup calculează partea din suma din (6.6) care îi revine, după care se obține întreaga sumă printr-o colectare cu sumare cu destinația $P(i)$, grup care poate calcula acum x_i .

Deci, doar din argumente de bun simț, se observă că se vor face difuzări pe coloanele torului și colectări pe linii. De asemenea, deoarece calculul final al unui element din x este efectuat de un singur procesor, și trebuie să alegem câte un procesor reprezentativ pentru fiecare linie și fiecare coloană, e normal ca procesoarele de pe diagonala torului să fie cele "avantajate"; de aici, cerința repartizării vectorului b numai între aceste procesoare.

Să trecem acum la o abordare mai concretă. Introducem matricea $W \in \mathbf{R}^{n \times \sqrt{p}}$; fiecare coloană de procesoare din tor posedă o coloană din W , repartizarea pe linii fiind ciclică, la fel ca pentru A . Procesorul P_{st} , care deține elementul w_{it} (deci $s = i \bmod \sqrt{p}$), va avea de calculat suma (prin care definim W):

$$w_{it} = \sum_{j < i, P(j)=t} a_{ij}x_j, \quad (\text{unde } P(j) = j \bmod \sqrt{p}).$$

Se observă că această sumă necesită doar elemente a_{ij} locale procesorului. De asemenea, folosind consacratul vector z , deducem imediat că

$$z_i (= \sum_{j=0}^{i-1} a_{ij}x_j) = \sum_{t=0}^{\sqrt{p}-1} w_{it}.$$

Această relație sugerează că, o dată calculate toate valorile w_{it} de pe linia $s = P(i)$ a torului, $x_i = (b_i - z_i)/a_{ii}$ se poate calcula după o colectare cu sumare pe linia respectivă, cu destinația P_{ss} (care deține atât b_i , cât și a_{ii}); este, de fapt, ideea algoritmului cu colectare 6.11.

Pe de altă parte, pentru a calcula w_{it} , un procesor are nevoie de toate valorile x_j , pentru $P(j) = t$; ori acestea sunt calculate de procesorul P_{tt} . Deci, de fiecare dată când un procesor de pe diagonala torului calculează o valoare din x , el trebuie să o difuzeze tuturor procesoarelor de pe coloana sa, acestea urmând a actualiza valorile w_{it} care le sunt locale. Deoarece operația de actualizare va apare în mai multe ocazii, să descriem procedura pe care o urmează procesorul P_{st} :

procedură actualizare(k, x_j)
pentru $i = k + 1 : n - 1$
dacă $i \bmod \sqrt{p} = s$ **atunci** {pentru liniile proprii}
 $w_{it} \leftarrow w_{it} + a_{ij}x_j$

Introducerea parametrului suplimentar k se va dovedi necesară mai târziu; în mod normal $k = j$, deci un procesor actualizează toate elementele care îi aparțin, pentru $i > j$.

S-a conturat până acum următorul algoritm de principiu, schițat pentru ansamblul procesoarelor:

ALGORITM 6.18 ($TR(n)$ pe tor, schiță generală)

pentru $j = 0 : n - 1$
 $t \leftarrow j \bmod \sqrt{p}$
 colectare cu sumare (pe linia t a torului) în P_{tt} : $z_j \leftarrow \sum_{l=0}^{\sqrt{p}-1} w_{jl}$
 P_{tt} calculează $x_j \leftarrow (b_j - z_j)/a_{jj}$
 P_{tt} difuzează x_j pe coloana t a torului
 P_{st} , cu $s \in 0 : \sqrt{p} - 1$ (coloana t a torului), efectuează *actualizare*(j, x_j)

Acest algoritm arată plăcut, dar ascunde defecte majore; deci, înainte de a detalia, să le evidențiem și să încercăm să le eliminăm. Un exemplu simplu ne va fi de folos; după ce P_{00} calculează x_0 , îl difuzează pe coloana 0 a torului; printre alții, îl recepționează și P_{10} , care trece imediat la actualizări; în acest timp, P_{11} , care așteaptă de la P_{10} valoarea w_{10} , necesară pentru calculul lui x_1 , nu are nimic de făcut; de-abia după ce P_{10} își termină actualizările, este disponibil pentru comunicarea cu P_{11} (participarea la colectare). Se arată astfel că algoritmul nu are un grad înalt de paralelism; dacă toate procesoarele de pe o coloană a torului sunt simultan ocupate (mai puțin pe durata calculului lui x_j , care e nesemnificativă), în schimb coloanele lucrează aproximativ secvențial; timpul de execuție nu va fi cu mult mai bun decât cel al unui inel de numai \sqrt{p} procesoare, ceea ce e absolut neconvenabil.

Soluția e de a amâna o parte a actualizărilor până când se rezolvă toate problemele de comunicație ale etapei curente; de exemplu, P_{10} va calcula întâi w_{10} , îl va trimite către P_{11} , iar apoi va actualiza $w_{\sqrt{p}+1,0}$, $w_{2\sqrt{p}+1,0}$, etc. În acest fel P_{11} nu va aștepta

decât un timp minim, iar P_{10} va avea doar de schimbat ordinea unor operații. Ideea este deci de a face comunicațiile cât mai devreme cu putință, imediat ce datele de transmis sunt gata.

Această tactică ne obligă la spargerea în etape a operațiilor de difuzare și colectare, care vor avea loc într-un timp mai îndelungat (nu vor însemna practic o sincronizare a procesoarelor), dar care se vor suprapune cu alte operații. Pentru a obține o bună eșalonare în timp a comunicațiilor între vecini, vom efectua difuzarea prin transmitere în jos, pe inelul format de fiecare coloană a torului, iar colectarea prin transmitere la dreapta, pe fiecare linie a torului.

Vom prezenta acum algoritmul pentru procesorul P_{st} , într-o formă mai amănunțită. Pentru a clarifica oarecum lucrurile, să precizăm comunicațiile la care participă un procesor nediagonal; într-o etapă, P_{st} participă la o difuzare pe coloana t a torului și la o colectare pe linia s ; deci, are de recepționat câte un mesaj de sus și din stânga, și de transmis jos și spre dreapta.

ALGORITM 6.19 ($TR(n)$ pe tor, detaliat, A repartizată ciclic pe linii și coloane, pentru procesorul P_{st})

```

inițializează toți  $w_{it}$  locali cu 0
pentru  $j = 0 : n - 1$ 
  dacă  $t = j \bmod \sqrt{p}$  atunci
    dacă  $s = t$  atunci                                     {procesor diagonal}
      (1)   recv( $z_j$ , stânga)                               {dacă  $j \neq 0$ }
      (2)    $z_j \leftarrow z_j + w_{jt}$ 
      (3)    $x_j \leftarrow (b_j - z_j) / a_{jj}$ 
      (4)   send( $x_j$ , jos)                                  {dacă  $j \neq n - 1$ }
      (5)   actualizare( $j, x_j$ )
    altfel                                                 { $s \neq t$ , procesor nediagonal}
      (6)   recv( $x_j$ , sus)
      (7)   send( $x_j$ , jos)                                  {dacă  $s \neq (t - 1) \bmod \sqrt{p}$ }
      (8)    $i \leftarrow j + (s - t) \bmod \sqrt{p}$ 
      (9)    $w_{it} \leftarrow w_{it} + a_{ij}x_j$ 
      (10)  recv( $z_i$ , stânga)                             {dacă  $t \neq (s + 1) \bmod \sqrt{p}$ }
      (11)   $z_i \leftarrow z_i + w_{it}$ 
      (12)  send( $z_i$ , dreapta)
      (13)  actualizare( $j + \sqrt{p}, x_j$ )

```

Deci, după ce un procesor diagonal primește în (1) rezultatul colectării de pe linia sa, la care adaugă în (2) și contribuția sa, e pregătit să calculeze un element din x în (3), aplicând formula (6.6); apoi, în (4), inițiază difuzarea elementului proaspăt calculat și actualizează în funcție de acesta elementele din W proprii, în (5).

Un procesor nediagonal așteaptă întâi ultima necunoscută x_j calculată de procesorul diagonal de pe coloana sa, pe care o trimite imediat mai departe; (6) și (7) sunt contribuția procesorului la difuzarea pe coloana sa. Apoi actualizează, în (9), doar primul element local w_{it} ; indicele i , calculat anterior în (8), este cel mai mic indice al

unei linii proprii, mai mare decât j ; într-adevăr, dacă $s > t$ atunci $i = j + (s - t)$, iar dacă $s < t$ atunci $i = j + (s - t + \sqrt{p})$. Cu acest w_{it} , procesorul participă la colectarea cu sumare de pe linia sa, în instrucțiunile (10), (11) și (12); în sfârșit, poate continua actualizările, în (13). Ca detaliu, instrucțiunile (7), (9) și (10) pot decurge în paralel; toate valorile z_i locale unui procesor pot fi găzduite de o singură variabilă scalară.

O formulă exactă pentru timpul de execuție al acestui algoritm este greu de dedus. Să ne mulțumim cu aprecierea că operațiile aritmetice sunt echilibrat împărțite între procesoare, iar volumul comunicațiilor este aproximativ $4(n/\sqrt{p})(\sigma + \beta)$ (un procesor efectuează 4 comunicații pe etapă, și sunt n/\sqrt{p} etape). Desigur, pot exista timpi morți în activitatea procesoarelor, dar ei nu pot fi prea mari, deoarece actualizările ocupă mult mai mult timp decât celelalte operații.

Probleme

P 6.2.1 Scrieți analogul algoritmului cu difuzare 6.10 pentru rezolvarea sistemului *superior* triunghiular $Ax = b$ ($a_{ij} = 0$ pentru $i > j$).

P 6.2.2 Scrieți un algoritm pentru rezolvarea $TR(n)$ pentru o arhitectură cu memorie comună, pornind de la varianta secvențială cu sumă vectorială.

P 6.2.3 Desenați diagrame Gantt pentru evoluția algoritmilor cu difuzare 6.10 și cu colectare 6.11, în cazul $n = 6$, $p = 3$.

P 6.2.4 La fel, pentru evoluția algoritmilor ciclici de rezolvare a problemei $TR(n)$, pe coloane 6.16 și pe linii 6.17.

P 6.2.5 Ce instrucțiuni pot fi executate în paralel (sau concurent) de către un procesor, în algoritmul ciclic pe coloane 6.16 de rezolvare a problemei $TR(n)$? Dar în algoritmul 6.17?

P 6.2.6 Cum se poate ajunge de la o repartizare a unui vector pe tor ca aceea a lui x din figura ?? (pentru înmulțirea matrice-vector), la cea din algoritmul pentru $TR(n)$ pe tor?

Răspunsuri la probleme

P1.1.1 Numărul de tacte la care se produce un rezultat este dictat de timpul cel mai mare al unui "post"; răspunsul este deci 2. Se poate obține un rezultat la fiecare tact dacă se dublează postul B , ca în figura 7.1; comutatoarele dinaintea și de după cele două blocuri B își schimbă poziția la fiecare tact și sunt în fază (blocul C citește rezultatul furnizat de un bloc B și, în același moment blocul B respectiv trebuie să citească un nou operand). În primul caz 100 de rezultate finale erau produse în 203 tacte, în al doilea, în 103.

P1.2.1 Se procedează prin inducție după p . Pentru $p = 1$ este evident. Presupunem că orice arbore de ordin p are $p - 1$ arce. Fie un arbore de ordin $p + 1$; presupunem că nu există nici un nod cu gradul 1; deci, fiecare nod are cel puțin două arce. Vom construi un ciclu astfel: plecăm dintr-un nod oarecare pe unul din arce; în nodul următor, continuăm pe al doilea arc al său, primul fiind deja utilizat; și tot așa, până când se ajunge într-un nod deja parcurs; arcele sunt toate distincte, deci am obținut un ciclu, ceea ce e imposibil într-un arbore. Deci există un nod de grad 1; eliminând acest nod și arcul său se obține un arbore de ordin p , cu $p - 1$ arce; deci arborele inițial are p arce.

P1.2.2 $\Delta p/2$. Din fiecare nod pleacă Δ arce; fiecare arc unește două noduri.

P1.2.3 Un drum între cele două noduri se construiește din arce aflate pe dimensiunile egale cu pozițiile biților diferiți. De exemplu, de la 100 la 001 se poate ajunge în doi pași, prin 000 sau 101 (pe dimensiunile 2,0, respectiv 0,2).

P1.2.4 Nodurile la distanță x de P_0 au reprezentări binare conținând x biți de 1. Deci, în total, sunt $\binom{d}{x}$ noduri. Reamintim în treacăt o identitate binecunoscută: $\binom{d}{0} + \binom{d}{1} + \dots + \binom{d}{d} = 2^d$, care arată aici numărul total de noduri în hipercub, acestea fiind grupate după distanțele față de P_0 .

P1.2.5 Sunt șase cazuri posibile; ne referim la figura 1.11a; deci, vecinii lui P_{ij} sunt:

- $P_{i+1,j}$, $P_{i,j+1}$, dacă $i = 0$, $j = 0$ sau $j = 1$, i impar și $\neq n - 1$.

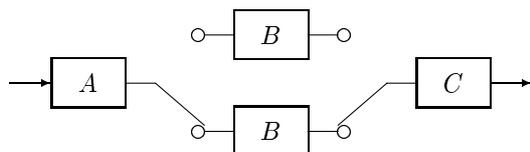


Figura 7.1: Schema pipeline pentru problema 1.1.1.

- $P_{i-1,j}, P_{i,j+1}$, dacă $i = n - 1, j = 0$ sau $j = 1, i$ par și $\neq 0$.
- $P_{i,j-1}, P_{i+1,j}$, dacă i par, $j = n - 1$.
- $P_{i,j-1}, P_{i-1,j}$, dacă i impar, $j = n - 1$.
- $P_{i-1,j}, P_{i+1,j}$, dacă $j = 0, 0 < i < n - 1$.
- $P_{i,j-1}, P_{i,j+1}$, altfel.

Iată deci că o figură simplă ascunde complicații destul de mari.

P1.2.6 Dacă i este numărul binar și g codul Gray, conversia binar-Gray se face cu $g_l = i_l \oplus i_{l+1}$, cu $0 \leq l \leq d - 1$ (cu convenția că $i_d = 0$). Invers, conversia Gray-binar se face cu $i_l = g_l \oplus g_{l+1} \oplus \dots \oplus g_{d-1}$, sau $i_l = g_l \oplus i_{l+1}$, cu $l = d - 1 : -1 : 0$.

În ciclul hamiltonian construit după codul Gray reflectat, care este un inel, procesorul g are vecinii g' și g'' , codurile Gray ale numerelor binare $i' = i - 1$ și $i'' = i + 1$, unde i are codul Gray g . Deci, pentru P_4 , parcurgem etapele: $g = 4 = 0100$ este codul Gray pentru $i = 0111$ (P_4 este pe poziția 7 în inel); $i' = 0110, i'' = 1000$, deci $g' = 0101 = 5$ și $g'' = 1100 = 12$.

P2.2.1 Pe hipercub, algoritmul este următorul:

1. dacă $\text{id} = k$
 1. pentru $i = 0 : d - 1$, în paralel
 1. send(M, i)
2. altfel dacă $\text{dist}_H(\text{id}, k) = 1$ {vecinii sunt la distanță Hamming 1}
 1. $i \leftarrow$ numărul de ordine al bitului de 1 din $\text{id} \oplus k$
 2. recv(M, i)

P2.2.2 În mod non-blocant, fiecare procesor transmite și recepționează pe toate canalele, în paralel. În mod blocant sincron, un procesor transmite întâi pe o dimensiune i dacă e.g. bitul i din adresa sa este 1, și recepționează întâi dacă acest bit este 0.

P2.4.1 Da. Fiecare procesor va avea alocată o parte din memoria comună, care va simula memoria locală. Pentru fiecare canal de comunicație între două procesoare va fi alocată câte o zonă, împărțită în două, fiecare parte pentru un sens de comunicație.

P2.4.2 Suma se calculează în u_0 . Considerăm $\text{id} \equiv i$.

1. $u_{2i} \leftarrow u_{2i} + u_{2i+1}$
2. dacă $i = 0$
 1. $u_0 \leftarrow u_0 + u_2$

P3.3.1 $j(\sigma + \beta) + \frac{m-\Delta}{\Delta}\beta$. Dacă graful este neregulat, Δ este gradul nodului sursă.

P3.3.2 Sunt $p - 1$ căi; una are lungime 1, iar celelalte lungime 2 și trec prin cele $p - 2$ noduri diferite de sursă și destinație.

P3.3.3 Când toate căile au aceeași lungime, deci atunci când $j = \sqrt{p}$; poate că merită să faceți două mici desene, cu \sqrt{p} par, respectiv impar.

P3.3.4 Un exemplu este cel din figura 7.2, în care sunt reprezentate câte două căi într-un desen; primele două au lungime 1, respectiv 7, iar celelalte lungime 3, respectiv 5.

P3.3.5 $b = a \oplus s$ are j biți de 1. Fie o permutare $(0, 1, \dots, d - 1) \xrightarrow{i} (i_0, i_1, \dots, i_{d-1})$, astfel încât $b_{i_0} = \dots = b_{i_{j-1}} = 1$ și $b_{i_j} = \dots = b_{i_{d-1}} = 0$; deci numărul binar $b_{i_{d-1}} \dots b_{i_1} b_{i_0}$ este egal cu $2^j - 1$. Fie acum o cale între nodurile 0 și $2^j - 1$ și fie c un nod de pe acea cale; transformăm această cale într-una între P_s și P_a , transformând numărul c după cum urmează. Se aplică permutarea inversă lui i biților numărului c , deci $c_0 \leftarrow c_{i_0}$, etc. (în



Figura 7.2: Patru căi arc-disjuncte între doi vecini, pe grilă.

acest fel se construiește o cale între 0 și b); apoi se face operația sau exclusiv $c \leftarrow c \oplus s$, acesta fiind rezultatul. Pentru sursă și destinație, transformările decurg astfel: $0 \xrightarrow{i^{-1}} 0 \xrightarrow{\oplus s} s$ și $2^j - 1 \xrightarrow{i^{-1}} b \xrightarrow{\oplus s} a$. Transformarea nu e unică, ci depinde de modul de alegere a permutării i . De exemplu, să vedem care este calea de lungime 4 dintre $s = 110$ și $a = 011$, într-un hipercub de dimensiune 3; $b = 101$ și putem alege i astfel: $(0, 1, 2) \xrightarrow{i} (0, 2, 1)$; în figura 3.10 se observă că între 0 și $2^2 - 1 = 011$, calea de lungime 4 este 000, 100, 101, 111, 011; aplicând inversa permutării i , aceasta se transformă în 000, 010, 011, 111, 101; în fine, după un sau exclusiv cu $s = 110$ se obține 110, 100, 101, 001, 011. Dacă se alege $(0, 1, 2) \xrightarrow{i} (1, 2, 0)$, atunci aplicând inversa lui i se obține 000, 010, 110, 111, 101; în final calea este 110, 100, 000, 001, 011.

Observație. Este poate mai simplu să se construiască direct căile între P_s și P_a . Se consideră mulțimea \mathcal{D} a pozițiilor în care s și a diferă și se construiesc j căi de lungime j prin modificări ale biților din acest grup (rotații, în \mathcal{D} , ale căii principale prin \mathcal{D}). Cele $d - j$ căi de lungime $j + 2$ se construiesc negând câte un bit care nu e în \mathcal{D} , mergând apoi pe calea principală prin \mathcal{D} , și negând la loc bitul inițial.

P3.3.6 a) Fiecare cale are un singur nod în C_k , obținut prin rotația lui $u_k = 0 \dots 0 \underbrace{1 \dots 1}_k$;

aceste d noduri sunt evident diferite. b) Cele j căi obținute prin rotirea căii principale nu au noduri comune. Nodurile de pe calea l , cu $l \geq j$, au bitul l egal cu 1, și diferă prin aceasta de toate celelalte noduri.

P3.3.7 În C_k sunt $\binom{d}{k}$ noduri; pentru fiecare nod, k arce merg către C_{k-1} și $d - k$ arce către C_{k+1} . Deci sunt $\binom{d}{k}(d - k) = \binom{d-1}{k}d$ arce între C_k și C_{k+1} .

P3.3.8 Difuzare, full duplex: primele Δ emise de P_s ajung la cel mai depărtat procesor cel puțin după timpul $D(\sigma + \beta)$, iar restul de $m - \Delta$, după cel puțin încă $\frac{m - \Delta}{\Delta}\beta$; se obține aceeași margine ca și în cazul comunicației între două procesoare aflate la distanță egală cu diametrul grafului.

Difuzare, half duplex: graful are cel mult $(p\Delta)/2$ arce folosite unidirecțional la un moment dat; numărul minim de valori care trebuie transmise global este $(p - 1)m$, mesajul trebuind să ajungă la toate celelalte procesoare; presupunând toate legăturile folosite permanent, se obține marginea inferioară cerută.

Difuzare generală: marginile sunt obținute din condiții de ocupare a arcelor grafului. Numărul minim total de valori care trebuie transmise este $p(p - 1)m$ (fiecare procesor trimite m valori pentru celelalte $p - 1$ procesoare; pentru a ajunge la câte un procesor, un mesaj trebuie să străbată cel puțin câte un arc), iar numărul maxim de mesaje care pot fi transmise simultan este de $p\Delta$ în cazul full duplex și $(p\Delta)/2$ în modelul half duplex.

Distribuție: procesorul care distribuie trebuie să transmită $(p - 1)m$ valori pe cele Δ

canale ale sale.

Schimb complet, inel: la distanță i de un procesor se află două procesoare. Numărul minim de arce parcurse de mesajele emise de procesor este $\sum_{i=1}^{p/2} 2i = \frac{p}{2}(\frac{p}{2} + 1)$. Pentru toate procesoarele, numărul total minim de treceri pe canale este deci $\frac{p^2}{2}(\frac{p}{2} + 1)$. Dacă toate cele $2p$ canale, în regim full duplex, respectiv p în half duplex, sunt folosite simultan, rezultă marginile din textul problemei.

Schimb complet, hiper cub: la distanță k de un procesor se află $\binom{d}{k}$ procesoare; cele $(p - 1)$ mesaje emise de către un procesor, spre fiecare dintre celelalte, au de parcurs cel puțin $\sum_{k=1}^d k \binom{d}{k} = \frac{d}{2}p - 1$ canale (pentru calculul sumei, vezi problema 3.3.24), deci în total se parcurg cel puțin $(dp^2)/2$ canale. Presupunând canalele folosite în paralelism complet, se obțin marginile din textul problemei.

P3.3.9 Deoarece din rădăcină pornesc Δ arce, se pot construi cel mult Δ arbori de acoperire. Graful are $p\Delta$ arce (orientate); cei Δ arbori ar avea în total $\Delta(p - 1)$ arce, deci nici o contradicție.

P3.3.10 P_0 difuzează mesajul M pe inel, după arborele de acoperire din figura 3.11:

1. **dacă** $id = 0$ **atunci**
 1. **în paralel** $\text{send}(M, \text{dreapta}), \text{send}(M, \text{stânga})$
2. **altfel** **dacă** $id \leq \lfloor p/2 \rfloor$ **atunci**
 1. $\text{recv}(M, \text{stânga})$
 2. **dacă** $id \neq \lfloor p/2 \rfloor$ **atunci** $\text{send}(M, \text{dreapta})$
3. **altfel**
 1. $\text{recv}(M, \text{dreapta})$
 2. **dacă** $id \neq \lfloor p/2 \rfloor + 1$ **atunci** $\text{send}(M, \text{stânga})$

P3.3.11 Toate expresiile e cu care este comparat id se modifică în $(s + e) \bmod p$. Evident că se pot face simplificări; de exemplu, în instrucțiunea 2.1 din algoritmul 3.3, $(s + p - 1) \bmod p \equiv (s - 1) \bmod p$.

P3.3.12 Algoritmul pentru procesorul P_{xy} este:

1. **dacă** $x = 0$ și $y = 0$ **atunci**
 1. **în paralel** $\text{send}(M, \text{dreapta}), \text{send}(M, \text{stânga}), \text{send}(M, \text{sus}), \text{send}(M, \text{jos})$
2. **altfel** **dacă** $y = 0$ **atunci**
 1. **dacă** $x \leq \lfloor \sqrt{p}/2 \rfloor$ **atunci**
 1. $\text{recv}(M, \text{stânga})$
 2. **în paralel** $\text{send}(M, \text{jos}), \text{send}(M, \text{sus}),$
 1. **dacă** $x \neq \lfloor \sqrt{p}/2 \rfloor$ **atunci** $\text{send}(M, \text{dreapta})$
 2. **altfel**
 1. $\text{recv}(M, \text{dreapta})$
 2. **în paralel** $\text{send}(M, \text{jos}), \text{send}(M, \text{sus}),$
 1. **dacă** $x \neq \lfloor \sqrt{p}/2 \rfloor + 1$ **atunci** $\text{send}(M, \text{stânga})$
3. **altfel** $\{y \neq 0\}$
 1. **dacă** $y \leq \lfloor \sqrt{p}/2 \rfloor$ **atunci**
 1. $\text{recv}(M, \text{jos}), \text{dacă } y \neq \lfloor \sqrt{p}/2 \rfloor$ **atunci** $\text{send}(M, \text{sus})$
 2. **altfel**
 1. $\text{recv}(M, \text{sus}), \text{dacă } y \neq \lfloor \sqrt{p}/2 \rfloor + 1$ **atunci** $\text{send}(M, \text{jos})$

P3.3.13 Adresele nodurilor adiacente arcului diferă în bitul k ; după rotația cu i biți la stânga (prin care se obține $\mathcal{A}_i(0)$), ele vor diferi în bitul $(k + i) \bmod d$.

P3.3.14 În etapa l , unde l este poziția celui mai semnificativ bit de 1 din j . În algoritmul de difuzare rotativ, mesajul M^i va fi primit în etapa l_i , unde l_i este poziția celui mai semnificativ bit de 1 din $\text{Ro}(j, -i)$; explicație: prin rotația cu i biți la dreapta se transformă nodurile din $\mathcal{A}_i(0)$ în noduri din $\mathcal{A}_0(0)$; ori, pentru acest arbore, știm deja răspunsul: l_i .

P3.3.15 Algoritmul de difuzare rotativ are forma următoare:

1. pentru $i = 0 : d - 1$
 1. $l_i \leftarrow$ poziția celui mai semnificativ bit de 1 din $\text{Ro}(\text{id}, -i)$ (-1 când $\text{id} = 0$)
2. pentru $k = 0 : d - 1$
 1. în paralel
 1. pentru $i = 0 : d - 1$, în paralel
 1. dacă $l_i = k$ atunci $\text{recv}(M^i, (l_i + i) \bmod d)$
 2. pentru $i = 0 : d - 1$, în paralel
 1. dacă $l_i < k$ atunci $\text{send}(M^i, (k + i) \bmod d)$

Se impun unele explicații, dată fiind forma condensată a algoritmului; k este etapa de comunicație; în fiecare etapă, fiecare procesor poate efectua în paralel, atât recepții, în 2.1.1, cât și transmisii, în 2.1.2. Condiția din 2.1.1.1 este justificată în rezolvarea problemei precedente: mesajul M^i sosește în etapa l_i . Canalul pe care sosește M^0 este l_0 (folosind arborele $\mathcal{A}_0(0)$); canalul pe care vine M^i se obține prin rotație, deci este în dimensiunea $(l_i + i) \bmod d$ (se folosește $\mathcal{A}_i(0)$; vezi și problema 3.3.13). După ce a primit un mesaj, un procesor îl emite în toate etapele, până la sfârșit, ceea ce se vede în condiția din 2.1.2.1, pe canalele din dimensiunile succesive celei pe care a fost primit (adică $(l_i + i + k - l_i) \bmod d = (k + i) \bmod d$).

P3.3.16 Se obține algoritmul 3.2 de comunicare după d căi disjuncte (vezi fig. 3.10).

P3.3.17 Timpul de difuzare este chiar D . Inițiatorul difuzării își informează vecinii, apoi aceștia vecinii lor, etc. La momentul k , toate nodurile la distanță $\leq k$ au primit mesajul.

P3.3.18 În prima etapă, inițiatorul difuzării comunică mesajul unui vecin oarecare. Apoi, ambele noduri informează alte două noduri; astfel numărul nodurilor care au primit mesajul se dublează la fiecare etapă. Timpul necesar este deci $\lceil \log p \rceil$.

Aceasta este și limita inferioară pentru un graf oarecare, care are mai puține arce decât graful complet. De asemenea, este și timpul necesar pe hipercub, în algoritmul 3.4. Deci, în modelul timp constant, 1-port, acest algoritm este optim (o bună dovadă a calităților hipercubului, care are mult mai puține arce decât graful complet).

P3.3.19 În modelul multiport, full duplex, timpul este 1 (fiecare comunică cu fiecare în ambele sensuri). În half duplex, timpul este 2; de exemplu, un procesor colectează simultan toate mesajele, apoi le difuzează; dacă nu vă place ideea de a difuza un mesaj atât de lung (deși în timp constant nu contează lungimea), se poate și altfel: în prima etapă toate perechile P_i, P_j comunică în acest sens dacă $i < j$; în a doua etapă se comunică în sens invers.

În 1-port, full duplex, timpul este $\lceil \log p \rceil$, ca pentru hipercub (algoritmul sugerat în figura 3.21). În half duplex, un timp de două ori mai mare este obținut imediat. Se poate însă și mai bine, dar cu complicații; până acum se știe că timpul se încadrează între limitele $1.44 \lceil \log p \rceil$ și $1.88 \lceil \log p \rceil$; prima valoare se poate demonstra ca limită teoretică (dar pe vreo două pagini); a doua este timpul de difuzare generală pe hipercub, după un algoritm complicat.

P3.3.20 Se transmite pe o singură direcție, spre dreapta, să zicem.

1. pentru $k = 0 : p - 2$

1. în paralel $\text{send}(M_{(\text{id}-k) \bmod p}, \text{dreapta}), \text{recv}(M_{(\text{id}-k-1) \bmod p}, \text{stânga})$

P3.3.21 Pe un graf oarecare, P_s trimite toate mesajele tuturor vecinilor; aceștia rețin mesajul lor și transmit restul mai departe; după k pași, toate procesoarele aflate la distanță k de P_s și-au primit mesajul. Timpul este deci D , diametrul grafului (mai exact, excentricitatea nodului P_s).

În modelul 1-port, pe inel, se folosește arborele de acoperire din figura 3.11. Se transmit întâi toate mesajele spre dreapta, apoi toate spre stânga, după care fiecare procesor retransmite mai departe. Când p este par, ramurile arborelui sunt inegale, deci timpul va fi $p/2$. Când p este impar, timpul va fi $(p + 1)/2$.

Pe hipercub, algoritmul este cel deja prezentat în figura 3.23.

P3.3.22 Mesajul pentru procesorul P_{xy} este trimis întâi pe orizontală, dacă $x + y$ este par, și întâi pe verticală dacă $x + y$ este impar. De exemplu, mesajul pentru procesorul $(2, 2)$ este trimis lui $(2, 0)$, împreună cu toate mesajele pentru care $x = 2$ și $x + y$ este par; apoi procesorul $(2, 0)$ distribuie pe verticala sa. În schimb, mesajul pentru $(2, 1)$ este trimis lui $(0, 1)$, împreună cu toate mesajele pentru care $y = 1$ și $x + y$ e impar; procesorul $(0, 1)$ distribuie apoi pe orizontala sa.

P3.3.23 E suficient să demonstrăm aceasta pentru distribuție. Dar distribuția decurge pe arborele din figura 3.13, construit pe principiul căii celei mai scurte (e intuitiv: se merge întâi pe orizontală pe drumul cel mai scurt, apoi la fel pe verticală).

P3.3.24 Considerăm suma termenilor egali distanțați de limitele de sumare: $i \binom{d}{i} + (d - i) \binom{d}{d-i} = i \binom{d}{i} + (d - i) \binom{d}{i} = d \binom{d}{i}$. Deci $\sum_{i=0}^d i \binom{d}{i} = (d/2) \sum_{i=0}^d \binom{d}{i} = d2^{d-1}$.

P3.3.25 a), b) Transpunerea se face printr-un schimb complet, indiferent de topologie, deoarece procesorul P_i trebuie să trimită blocul A_{ij} procesorului P_j , pentru orice i și j .

c) Se poate proceda astfel: fiecare procesor "diagonal" P_{kk} colectează mesajele de la procesoarele de pe linia sa, P_{kj} , cu $0 \leq j < \sqrt{p}$. Apoi, P_{kk} distribuie pe coloana sa, procesoarelor P_{ik} , cu $0 \leq i < \sqrt{p}$. Cum colectarea și distribuția au același timp de execuție și se desfășoară pe un inel, timpul total va fi $\sqrt{p}(\sigma + m\beta)$ (m fiind numărul de elemente dintr-un bloc).

Acum se poate aplica îmbunătățirea clasică; fiecare mesaj se împarte în două pachete, trimise unul pe orizontală, altul pe verticală (P_{ij} trimite primul pachet lui P_{ii} și pe al doilea lui P_{jj}). Procesoarele diagonale distribuie pe verticală primul pachet și pe orizontală pe cel de-al doilea. Se înjumătățește timpul de propagare de mai sus.

S-ar părea că mai bine nu se poate. Dar, să observăm că procesorul P_{ij} schimbă un mesaj cu P_{ji} , deci cu exact un alt procesor; transpunerea pe tor este deci o operație de transmisie permutată. Ea va dura cât două schimburi complete cu mesaje de lungime m/p , adică, dacă luăm cel mai eficient algoritm de schimb complet, $2\sqrt{p}\sigma + (\sqrt{p}/4)\beta$.

P3.3.26 În modelul multiport, timpul este 1; fiecare procesor trimite pe fiecare canal mesajul potrivit, recepționând în același timp. În modelul 1-port, transmite (și recepționează) pe rând pe câte un canal, deci timpul este $p - 1$.

P3.5.1 Un algoritm simplu poate fi schițat după cum urmează, procesorul destinație fiind P_{xy} : dacă $\text{id}_x \neq x$ și $\text{id}_y \neq y$ atunci co va fi canalul opus lui ci ; dacă $\text{id}_x = x$ sau $\text{id}_y = y$ atunci se alege co , canalul (unic) care apropie de destinație. Prima regulă poate să nu fie aplicabilă în cazul unui procesor de la periferie; pentru acestea, se alege orice canal

diferit de ci , care apropie mai mult de destinație (pot fi două astfel de canale). Conform acestor reguli se obțin, în general, două căi distincte, sau trei în cazul procesoarelor aflate pe aceeași linie sau coloană.

P3.5.2 Desigur. Căile între 000 și 110 sunt: prima 000, 001, 101, 111, 110, a doua 000, 010, 110, a treia 000, 100, 110.

P3.5.3 Indiferent de canalul de intrare, mesajul pentru P_a este trimis spre P_a ; astfel, se obțin o cale de lungime 1 și $p - 2$ căi de lungime 2.

P4.3.1 Produsul scalar $x^T y = \sum_{i=0}^{n-1} x_i y_i$, este în fond tot o sumă, după efectuarea produselor $x_i y_i$.

P4.3.2 $p = O(n/\log n)$. Fiecare procesor efectuează întâi suma a câte $O(n/p) = O(\log n)$ elemente, în timp $O(\log n)$; rămân astfel de sumat câte un element pentru fiecare procesor; aceasta se poate face în timp $O(\log(n/\log n)) = O(\log n - \log \log n) = O(\log n)$. Deci, timpul total de execuție este $O(\log n)$, iar costul, evident $O(n)$.

P4.3.3 $\sum_{i=0}^{n/2^k-1} x_{i+2^k}$.

P4.3.4 Funcția de mai jos calculează suma unor valori întregi, dar generalizarea este banală. Suma este calculată doar de P_0 , celelalte procesoare returnând o valoare nesemnificativă. Variabila $k2$ este egală cu 2^{k+1} (atenție: nu se verifică egalitatea lui p cu o putere a lui 2).

```
int suma_hip(int val)
{
    int my_id, p, d, vecin, x;
    int k, k2;
    MPI_Status status;

    MPI_Comm_rank (MPI_COMM_WORLD, &my_id);
    MPI_Comm_size (MPI_COMM_WORLD, &p);
    d = log(p)/log(2) + 0.5;

    for ( k = d-1, k2 = p ; k >= 0 ; k--, k2 >>= 1 )
    {
        vecin = my_id ^ ( 1 << k );
        if ( my_id < (k2>>1) )
        {
            MPI_Recv(&x, 1, MPI_INT, vecin, 0, MPI_COMM_WORLD, &status);
            val += x;
        }
        else if ( my_id < k2 )
            MPI_Send(&val, 1, MPI_INT, vecin, 0, MPI_COMM_WORLD);
    }

    return val;
}
```

P4.3.5 Trebuie modificate condiția din instrucțiunea **dacă**, din $i < n/2^k$ în $i \bmod 2^k = 0$, și adaptați corespunzător indicii din atribuirea 1.1.1. Algoritmul pentru P_i este:

1. pentru $k = 1 : \log n$
 1. dacă $i \bmod 2^k = 0$ atunci
 1. $x_i \leftarrow x_i + x_{i+2^{k-1}}$

Astfel, în etapa k , rămân în cursă doar procesoarele cu adrese de forma $j2^k$, comunicând cu procesoare $j2^k + 2^{k-1}$, deci vecine în hipercub în dimensiunea $k - 1$. Atenție, algoritmul este scris pentru $p = n$ procesoare, situația "normală" în cazul memoriei distribuite (spre deosebire de $p = n/2$ pe PRAM); procesoarele cu adresă impară nu efectuează nici o operație. Se va comunica în ordine inversă a dimensiunilor, față de algoritmul 4.3, dar tot după un arbore de acoperire binomial, la fiecare pas după o direcție.

Nu mai este necesară comutativitatea operației $+$; la pasul k , variabila locală procesorului P_i (cu adresa astfel încât $i \bmod 2^k = 0$) are valoarea $\sum_{l=i}^{i+2^k-1} x_l$, adică suma unor valori x_l cu indici consecutivi.

P4.3.6 Fiecare procesor își împarte elementele proprii în d grupuri și calculează suma pentru fiecare grup. Procesorul P_i are acum d valori: x_0^i, \dots, x_{d-1}^i . Apoi se calculează simultan sumele $s_k = \sum_{l=0}^{p-1} x_k^l$, fiecare comunicând pe un alt arbore din familia de arbori binomiali rotativi. Comunicațiile decurg simultan, deoarece pentru fiecare s_k se folosește la un moment dat o singură dimensiune. Din păcate, P_0 comunică la fiecare etapă pe toate direcțiile și face d adunări, câte una pentru fiecare s_k ; deci timpul consumat de P_0 dictează timpul total. Acesta este: $d[n/(pd) - 1]\alpha$ pentru sumele locale, $d^2\alpha$ pentru sumele efectuate de P_0 pe parcurs (d etape, câte d adunări) și $d(\sigma + \beta)$ pentru comunicație. Totalul este mai mare decât cel pentru algoritmul 4.3, deci ideea arborilor rotativi nu e eficient aplicată astfel. Dacă însă trebuie calculate mai multe sume, cel mult d , și pentru fiecare se folosește un arbore din familie, atunci timpul de comunicație va fi identic cu cel din cazul calculului unei singure sume.

P4.3.7 Fiecare procesor calculează suma elementelor locale, apoi o trimite spre P_0 , pe drumul cel mai scurt (arborele de acoperire din figura 3.11). P_0 , de îndată ce recepționează două sume, le adună la suma sa locală, continuând să recepționeze, în paralel; dacă $2\alpha \leq \sigma + \beta$, calculele nu întârzie comunicația! (Desigur, aceasta este valabil în cazul modelului de comunicație asincron; totuși, chiar în modelul sincron, execuția concurrentă a calculelor și comunicației pe P_0 va reduce posibilele întârzieri intervenite în comunicație.) Timpul total va fi $(n/p - 1)\alpha$ pentru sumele locale, $(p/2)(\sigma + \beta)$ pentru comunicație, și 2α pentru adunarea ultimelor două sume primite.

Dacă procedam ca la hipercub, efectuând sumele pe parcurs (fiecare procesor primește de la un vecin o valoare, o adună la valoarea proprie și transmite această sumă celui alt vecin), s-ar fi făcut $(p/2)\alpha$ adunări, nesuprapuse cu comunicația.

Ideea colectării în P_0 și a efectuării calculelor pe măsură ce vin sumele locale, în paralel cu comunicația este bună la inel pentru că la orice distanță de P_0 se află doar 2 procesoare. Pe un hipercub, colectând după un arbore de acoperire echilibrat, ar trebui efectuate mult mai multe calcule în timpul $\sigma + \beta$; dacă comunicațiile sunt foarte lente, merită încercat.

P4.3.8 Fiecare procesor calculează suma elementelor locale. Apoi, se comunică în $p/2$ etape, ca la difuzarea generală; în prima etapă, un procesor transmite în ambele direcții suma sa; apoi, pe măsură ce primește sumele altor procesoare, le adună la suma sa și le transmite mai departe, cu excepția ultimei etape. Practic, aceasta este simetrizarea algoritmului din problema 4.3.7; timpul de calcul este același; și acum retransmisia a două sume se poate face în paralel cu adunarea lor la suma locală.

Dacă sunt de calculat m sume, singurul truc aplicabil este de a grupa mesajele aparținând unor calcule diferite. Astfel timpul de comunicație va fi $(p/2)\sigma + (pm/2)\beta$. Coeficientul lui σ este același ca în cazul unei singure sume.

P4.3.9 Da, bucla **pentru** poate merge doar până la $(i+1)n/p-2$, pentru că $s_{(i+1)n/p-1}$ este chiar z_i .

P4.3.10 Fie $j = i + 2^l$; notăm gi și gj codurile Gray ale lui i , respectiv j . Reamintim regula de conversie binar-Gray: $gi_k = i_k \oplus i_{k+1}$ (vezi **P1.2.6**). Sunt două cazuri posibile.

a) $j_l = 1$, deci $i_l = 0$ și atunci $j_k = i_k$ pentru $k \neq l$. Atunci $gj_l = \overline{g_i}$ și $gj_{l-1} = \overline{g_{i-1}}$, iar în rest adresele coincid. Deci distanța este 2, mai puțin în cazul $l = 0$ (când nodurile sunt vecine în ciclul hamiltonian).

b) $j_l = 0$, deci $i_k = \overline{j_k}$, pentru $k \geq l$; pentru $k < l$, adresele coincid. Atunci $gj_{d-1} = \overline{g_{i-1}}$, dar $gj_{d-2} = g_{i_{d-2}}, \dots, gj_l = g_{i_l}$ (deoarece $a \oplus b = \overline{a} \oplus \overline{b}$); apoi, $gj_{l-1} = \overline{g_{i-1}}$, $gj_{l-2} = g_{i_{l-2}}, \dots, gj_0 = g_{i_0}$; deci, dimensiunile în care diferă gi și gj sunt $d-1$ și $l-1$, iar distanța între noduri este 2, făcând din nou excepție cazul $l = 0$.

P4.3.11 Instrucțiunile 2.2 și 2.3 se modifică în:

2.2. **dacă** $id_k = 1$ **atunci** $s \leftarrow v + s, z \leftarrow v + z$

2.3. **altfel** $z \leftarrow z + v$

Explicația e următoarea: pentru un procesor dintr-un hipercub "superior" (la un moment dat), suma globală v recepționată provine dintr-un hipercub "inferior", deci ea trebuie să fie termen stâng în operația notată cu $+$, iar variabila locală termen drept. Pentru procesoarele dintr-un hipercub "inferior", lucrurile se petrec pe dos.

P4.3.12 Da, dacă două procesoare calculează în paralel coeficienții a_i , respectiv b_i . Numărul de operații va fi $2 + 2 \log n$, cu n procesoare, dar pe un CREW PRAM.

P4.3.13 Păstrăm ipoteza simplificatoare $n = 2^r$; presupunem în plus că $p = 2^s$, deci $n/p = 2^{r-s}$. Pe EREW PRAM, algoritmul pentru procesorul P_i este următorul:

1. **pentru** $k = 1 : r - s - 1$
 1. **pentru** $l = in/p : 2^k : (i+1)n/p - 1$
 1. $a_l \leftarrow a_{i+2^{k-1}l}, b_l \leftarrow a_{i+2^{k-1}l} + b_{i+2^{k-1}l}$
 2. **pentru** $m = 1 : s + 1$
 1. **dacă** $i < 2^{s+1-m}$ **atunci**
 1. $l \leftarrow in/p, k \leftarrow r - s - 1 + m$
 2. $a_l \leftarrow a_{i+2^{k-1}l}, b_l \leftarrow a_{i+2^{k-1}l} + b_{i+2^{k-1}l}$
 3. **dacă** $i = 0$ **atunci** $x_n \leftarrow a_0^r x_0 + b_0^r$

În prima buclă **pentru**, se reduce numărul elementelor din șir la $2p$; a doua buclă **pentru** este transcrierea algoritmului 4.10.

P4.3.14 Substituția expresiei termenului x_{i+1} în relația de calcul a termenului x_{i+2} conduce la:

$$x_{i+2} = \frac{(a_{i+1}a_i + b_{i+1}c_i)x_i + (a_{i+1}b_i + b_{i+1}d_i)}{(c_{i+1}a_i + d_{i+1}c_i)x_i + (c_{i+1}b_i + d_{i+1}d_i)}$$

Deci dependența dintre x_{i+2} și x_i are aceeași formă ca relația de recurență inițială. Atunci se poate aplica algoritmul de reducere ciclică 4.10; bineînțeles, se înlocuiește calculul noilor coeficienți ai recurenței cu expresiile ce rezultă din relația de mai sus.

P4.3.15 Dată fiind forma sistemului, ecuația i are forma $a_{i,i-1}x_{i-1} + a_{ii}x_i = b_i$, pentru $i \geq 2$; deci, $x_1 = b_1/a_{11}$, și $x_i = (b_i - a_{i,i-1}x_{i-1})/a_{ii}$, pentru $i \geq 2$. Ori, aceste relații definesc un șir recurent de ordinul I; practic, este al mod de scriere a relațiilor de recurență, sub formă matriceală.

P4.3.16 Algoritm pentru procesorul P_i , pe EREW PRAM, n putere a lui 2, $p = n/2$:

1. **pentru** $k = 1 : \log n$
 1. **dacă** $i < n/2^k$ **atunci**
 1. $a_i^k \leftarrow a_{2i+1}^{k-1} a_{2i}^{k-1}$, $b_i^k \leftarrow a_{2i+1}^{k-1} b_{2i}^{k-1} + b_{2i+1}^{k-1}$
 2. **dacă** $i = 0$ **atunci** $w_0^{\log n} \leftarrow a_0^{\log n} x_0 + b_0^{\log n}$
3. **pentru** $k = \log n : -1 : 1$
 1. **dacă** $i < n/2^k$ **atunci**
 1. $w_{2i+1}^{k-1} \leftarrow w_i^k$, $w_{2i}^{k-1} \leftarrow a_{2i}^{k-1} w_{i-1}^k + b_{2i}^{k-1}$

Am notat cu w elementele șirului, indicele inferior fiind poziția în șir, iar cel superior nivelul de recurență; convenim că $w_{-1}^k = 0$; rezultatul este $w_i^0 \equiv x_{i-1}$. În prima buclă sunt calculați doar coeficienți, iar în a doua, doar elemente ale șirului; această a doua buclă trebuie amorsată prin calculul ultimului element $x_n \equiv w_0^{\log n}$.

P4.3.17 Presupunem că un procesor P_i are o singură pereche de coeficienți, notați a, b ; în aceste variabile, procesorul va calcula coeficienții pentru care $x_{i+1} = ax_0 + b$, aplicând de mai multe ori formule de tipul (4.5). În plus, fiecare procesor mai are o pereche de coeficienți, notați c și d , pentru care, în final, $x_n = cx_0 + d$. Deci, ca și în algoritmul 4.9, un procesor calculează coeficienții necesari doar lui (echivalentul sumelor prefixelor), dar și coeficienții necesari ultimului termen (echivalentul sumei totale). Algoritm este următorul:

1. **pentru** $k = 0 : d - 1$
 1. **în paralel** $\text{recv}((a', b'), k)$, $\text{send}((c, d), k)$
 2. **dacă** $\text{id}_k = 1$ **atunci**
 1. $b \leftarrow ab' + b$, $a \leftarrow aa'$
 3. $d \leftarrow cb' + d$, $c \leftarrow ca'$

Ca și în cazul algoritmului 4.9, e de așteptat o eficiență relativă sub 1/2, deoarece toate operațiile sunt dublate.

P4.3.18 La fiecare pas, procesorul P_i recepționează două valori din x și actualizează z_i ca în instrucțiunea 2.4 din algoritmul următor (corect doar pentru $p = n$ impar, dar ușor adaptabil și pentru p par):

1. $z_i \leftarrow y_i + a_{ii}x_i$, $j \leftarrow i$, $l \leftarrow i$
2. **pentru** $k = 0 : \lfloor p/2 \rfloor - 1$
 1. **în paralel**
 1. $\text{send}(z_j, \text{stânga})$
 2. $\text{recv}(z_{(j+1) \bmod n}, \text{dreapta})$
 3. $\text{send}(z_l, \text{dreapta})$
 4. $\text{recv}(z_{(l-1) \bmod n}, \text{stânga})$
 2. $j \leftarrow (j + 1) \bmod n$
 3. $l \leftarrow (l - 1) \bmod n$
 4. $z_i \leftarrow z_i + a_{ij}x_j + a_{il}x_l$

P4.3.19 Procesorul P_i are în memorie linia i din A , memorată sub forma unui vector a , și elementele x_i, y_i , memorate în variabilele scalare x, y . Rezultatul z_i este memorat în variabila scalară z .

1. $z \leftarrow y, j \leftarrow \text{id}$
2. **pentru** $k = 0 : p - 1$
 1. $z \leftarrow z + a_j x$
 2. **dacă** $k < p - 1$ **atunci, în paralel**
 1. **send**(x , stânga)
 2. **recv**(u , dreapta)
 3. $x \leftarrow u$
 4. $j \leftarrow (j + 1) \bmod p$

P4.3.20 În notație MATLAB, liniile pe care le posedă procesorul P_i , presupunând o repartizare bloc-linii, sunt $l = ip : (i + 1)p - 1$; datele locale lui P_i sunt $A(l, :), x(l), y(l)$. La fiecare pas se comunică n/p elemente din x și se actualizează n/p elemente din z ; vom nota cu w variabila locală care găzduiește elementele din x (un vector de dimensiune n/p) și cu c indicii corespunzători lui w în x . Cu aceste notații, algoritmul pentru P_i este următorul:

1. $z(l) \leftarrow y(l), j \leftarrow i, c \leftarrow l, w \leftarrow x(c)$
2. **pentru** $k = 0 : p - 1$
 1. $z(l) \leftarrow z(l) + A(l, c)w$
 2. **dacă** $i < p - 1$ **atunci, în paralel**
 1. **send**(w , stânga)
 2. **recv**(w' , dreapta)
 3. $j \leftarrow (j + 1) \bmod n$
 4. $c \leftarrow jp : (j + 1)p - 1$ {indicii elementelor din x }
 5. $w \leftarrow w'$ {acum $w = x(c)$ }

În cazul în care n nu se divide cu p , ultimul procesor are repartizate mai puține linii. Pentru a descrie aceasta, e suficient a se înlocui, în expresia lui l , $(i + 1)p - 1$ cu $\min((i + 1)p - 1, n)$. De asemenea, pentru ca înmulțirile să fie efectuate corect, în expresia lui c se înlocuiește $(j + 1)p - 1$ cu $\min((j + 1)p - 1, n)$. (O altă variantă, aproape echivalentă, e de a completa cu zerouri matricea A și vectorii x, y , până la cel mai apropiat multiplu de p .)

P4.3.21 Indiferent de relația între m și n , numărul de operații efectuate de un procesor este același, în algoritmi 4.16 și 4.18, repartizarea calculelor între procesoare fiind echilibrată. În ce privește comunicația, lucrurile nu mai stau așa. În algoritmul 4.16 se comunică elementele lui x , în fiecare dintre cele $p - 1$ etape de comunicație câte n/p (x are dimensiune n), ceea ce conduce la o complexitate a comunicației de aproximativ $p\sigma + n\beta$. În algoritmul 4.18 se comunică elementele lui z , în fiecare dintre cele $p - 1$ etape câte m/p (z are dimensiune m), deci, în total, $p\sigma + m\beta$. Acum devine clar care dintre algoritmi e mai bun: dacă $m > n$, atunci algoritmul 4.18 e mai rapid, dacă $m < n$, atunci 4.16.

P5.1.1 Algoritmul secvențial nu face nici o interschimbare de elemente, folosind o singură locație suplimentară de memorie (cea pentru minim); un algoritm paralel fără interschimbări ar avea nevoie de p locații suplimentare (de ce?).

P5.1.2 Bucla în i a algoritmului 5.1 va fi "spartă" în mai multe bucle de dimensiune p :

1. **pentru** $i = 1 : \log n$
 1. $j \leftarrow k$
 2. **cât timp** $j < n/2^i$
 1. $exch(2^i j, 2^{i-1}(2j+1))$, $j \leftarrow j+p$

Presupunând $p = 2^r$ se obține:

$$T(n) = \frac{n}{2p} + \frac{n}{4p} + \dots + \frac{n}{2^{\log n - (r-1)}} + \underbrace{1 + \dots + 1}_{\log p} = \frac{n}{p} + \log p - 1$$

și deci $\varepsilon(n, p) = \frac{n-1}{n-1+p \log p}$. Dacă $p \log p \ll n$ atunci $\varepsilon(n, p) \approx 1$.

O variantă mai simplă este aceea în care fiecare procesor calculează minimul între n/p elemente, apoi, pentru cele p elemente rămase, se aplică algoritmul 5.1.

P5.1.3 În algoritmul 5.1 minimul rămâne între elementele cu indice divizibil cu 2^i , după pasul i . Pentru a rămâne între primele 2^i elemente:

1. **pentru** $i = \log n - 1 : -1 : 0$ {se inversează ordinea aici}
 1. **dacă** $k < n/2^{\log n - i}$ **atunci** $exch(k, k+2^i)$

P5.1.4 De exemplu: $c_{(i+j) \bmod n} \leftarrow 1$ sau $c_{i \oplus j} \leftarrow 1$; ultimul mod e corect doar pentru p putere a lui 2.

P5.1.5 a) Da b) Da; dacă sunt mai multe elemente cu valoare minimă, în cazul testului $a_i > a_j$ va fi selectat cel cu indicele cel mai mare, iar în cazul testului $a_i \geq a_j$ cel cu indicele cel mai mic.

P5.1.8 Cu cel din problema **P5.1.3** (se elimină procesoarele cu numere mari). Pentru a semăna cu celălalt algoritmul se inversează ordinea în buclă: **pentru** $i = 0 : d-1$.

P5.2.1 Citirea valorilor c_k și c_{k-1} se face în doi pași în loc de unul. În plus, trebuie presupus că fiecare procesor cunoaște valoarea v , altfel e necesară difuzarea acesteia, în $\log n$ etape. Numărul de comparații rămâne însă tot 1.

P5.2.2 Se face difuzarea valorii v , dacă nu e deja cunoscută de toate procesoarele. Fiecare procesor face o căutare binară locală cu $\log m$ comparații. Un procesor P_k poate întâlni trei situații: $v > A_k$, $v < A_k$ sau v între două elemente din A_k ; să codificăm aceste situații cu 0, 1, 2. Fiecare procesor trimite codul său vecinului din dreapta (P_k lui P_{k+1}). Un singur procesor P_j va difuza valoarea rangului; fie unicul care posedă codul 2 și atunci rangul e $j+m+r$, r fiind rangul local; fie cel care deține codul 1 și primește din stânga codul 0 și atunci rangul e jm (completați algoritmul cu acțiunile suplimentare ale procesoarelor P_0 și P_{p-1}).

P5.2.3 Algoritmul este identic cu 5.5, cu diferența că evaluarea rangului local se face parcurgând secvența A_k de la un capăt la celălalt, deci cu n/p comparații în loc de $\log(n/p)$.

P5.2.4 Desigur că s-ar putea copia elementele mai mici decât v într-o zonă contiguă de memorie, la fel cele mai mari, apoi se va recopia totul în A . Se poate face însă totul pe loc în A , după cum urmează (presupunem că pivotul este a_s):

procedură *partiționare*(s, d, v, i_v)

1. $l \leftarrow s+1$, $u \leftarrow d$, $v \leftarrow a_s$
2. **cât timp** $l < u$
 1. **cât timp** $a_l < v$, $l \leftarrow l+1$
 2. **cât timp** $a_u \geq v$, $u \leftarrow u-1$

3. dacă $l < u$ atunci $a_l \leftrightarrow a_u, l \leftarrow l + 1, u \leftarrow u - 1$

3. $i_v \leftarrow u, a_s \leftrightarrow a_u$

Cei doi indici demarcează zona în care elementele nu au fost încă grupate; la stânga lui l sunt elemente mai mici decât v , la dreapta lui u , elemente mai mari; l este crescut până la găsirea primului element mai mare decât v ; u este micșorat până la găsirea unui element mai mic decât v ; apoi aceste două elemente sunt interschimbate. În final, u indică poziția celui mai din dreapta element mai mic decât pivotul; în această poziție este adus pivotul, prin interschimbare cu a_s . Numărul de comparații este $d - s$.

Acest algoritm este foarte secvențial. Algoritmii paraleli de partiționare (pentru arhitecturi cu memorie comună) folosesc copierea în altă zonă.

P5.2.5 Să considerăm un mic exemplu, în care două procesoare dețin secvențele $\langle 0, 1, 1, 2 \rangle$, respectiv $\langle 0, 0, 1, 3 \rangle$. Rangul lui 0 este 0, rangul lui 1 este 3. Ce se întâmplă dacă dorim aflarea elementului cu rangul 2 și aplicăm algoritmul 5.9? Iterația 1: se propun candidații 1 și 0, este ales 0, care are rang $0 < 2$; spațiile de căutare se reduc la $\langle 1, 1, 2 \rangle$ și $\langle 1, 3 \rangle$. Iterația 2: ambii candidați sunt 1, rangul este $3 > 2$; spațiile de căutare se reduc la mulțimea vidă și algoritmul nu a găsit soluția. Explicația este simplă: valoarea cu rangul 2 nu există, dacă definim rangul unei valori ca fiind egal cu numărul de elemente strict mai mici decât acea valoare. Ar trebui să rămânem la definiția rangului ca poziție în secvența sortată, dar aceasta nu ne este de mare ajutor din punct de vedere al calculului; P_1 din exemplul de mai sus n-are de unde ști ce rang local să atribuie valorii 0: 0 sau 1?

Poate că am divagat un pic. Iată o soluție. Să definim rang *minim* r al unei valori v , numărul de elemente din A strict mai mici decât v (ceea ce am numit simplu rang până acum) și rang *maxim* R numărul de elemente din A mai mici sau egale cu v (rangul primului element strict mai mare decât v în secvența A sortată). Să presupunem că, în algoritmul 5.9, se calculează și r și R . Atunci condiția de terminare este $r \leq j \leq R - 1$, ceea ce se explică foarte simplu: în secvența sortată A , v apare prima dată în poziția r și ultima dată în poziția $R - 1$, între aceste poziții fiind numai elemente cu valoarea v .

Calculul celor două ranguri în algoritmul 5.5 nu dublează decât partea de comunicație și calculul sumei globale. Găsirea rangurilor locale este simplă; după stabilirea rangului minim prin căutare binară, cel maxim se determină prin căutarea primului element mai mare decât v la dreapta rangului minim (inclusiv acesta).

În fine, o concluzie mai filosofică: găsirea rangului și selecția nu sunt una inversa celeilalte decât dacă elementele din A sunt distincte.

P5.2.6 Alegerea candidatului local, la pașii (1) și (6), se face luând mediana unui mic eșantion sau un element oarecare. Alegerea globală a candidatului decurge la fel (pașii (2) și (3)). P_0 difuzează candidatul global. Fiecare procesor aplică algoritmul de partiționare din 5.2.4, calculând și rangul local al lui v ; abia acum se continuă cu calculul rangului global (pasul (4), mai puțin difuzarea lui v , vezi algoritmul 5.5). La pasul (5), reducerea spațiului de căutare se face prin eliminarea uneia din cele două părți ale secvenței locale, mai mică sau mai mare decât candidatul după cum e cazul.

Complexitatea comunicației rămâne aceeași, adică $O(p \log n)$, iar numărul de comparații va fi $O(n/p)$ (se raționează ca la algoritmul 5.8).

În nici un caz nu se face sortarea locală, urmată de aplicarea algoritmului 5.9 ca atare. Aceasta înseamnă $O(n/p \log(n/p))$ operații (sortarea este cea mai costisitoare).

P5.3.1 Notăm $D = A \perp B$; notăm i indicele în A , j indicele în B și l indicele în D .

1. $i \leftarrow 0, j \leftarrow 0, l \leftarrow 0$

2. cât timp $i \leq n$ și $j \leq m$
 1. dacă $a_i < b_j$ atunci $d_i \leftarrow a_i, i \leftarrow i + 1$
 2. altfel $d_i \leftarrow b_j, j \leftarrow j + 1$
 3. $l \leftarrow l + 1$
3. cât timp $l \leq n + m$ {adaugă un sfârșit de secvență}
 1. dacă $i \leq n$ atunci $d_l \leftarrow a_i, i \leftarrow i + 1$
 2. altfel $d_l \leftarrow b_j, j \leftarrow j + 1$
 3. $l \leftarrow l + 1$

P5.3.2 La primul pas de recursie am văzut că sunt suficiente \sqrt{nm} procesoare. La al doilea pas, problema de interclasare este spartă în $q = \sqrt{n}$ probleme, fiecare de dimensiune n_i, m_i , unde $n_i = \lfloor \sqrt{n} \rfloor$ sau $n_i = \lceil \sqrt{n} \rceil$; despre fiecare m_i în parte nu știm nimic; în schimb $\sum_{i=1}^q n_i = n, \sum_{i=1}^q m_i = m$. Rezolvarea fiecărei subprobleme necesită $\sqrt{n_i} \sqrt{m_i}$ procesoare. Trebuie să demonstrăm că

$$\sum_{i=1}^q \sqrt{n_i} \sqrt{m_i} \leq \sqrt{n} \sqrt{m}$$

Dar aceasta este chiar inegalitatea Cauchy-Buniakowski-Schwarz pentru numerele $\sqrt{n_i}, \sqrt{m_i}$.

P5.3.4 Un procesor care termină de interclasat secvențele sale îl poate ajuta pe un altul, de exemplu pe cel care are cea mai lungă secvență B'_i și nu este încă ajutat; el începe să interclaseze începând de la coadă (și cu asta am sugerat răspunsul întrebării puse la începutul secțiunii). Evident că așa ceva nu-i deloc ușor de programat; deci nu e eficient, probabil.

P5.3.5 Sunt $\text{rang}(a_i^*)$ elemente mai mici decât a_i^* în secvența $A^* \perp B^*$ și sunt i elemente mai mici decât a_i^* în A^* . Deci, din cele $\text{rang}(a_i^*)$ mai mici în $A^* \perp B^*$, $\text{rang}(a_i^*) - i$ aparțin lui B^* . Cum $lp(a_i^*)$ este indicele celui mai mare element din B^* mai mic decât a_i^* , atunci $lp(a_i^*) = \text{rang}(a_i^*) - i - 1$ (al j -lea element din B^* este b_{j-1}^*), de unde (5.2).

P5.3.6 Sunt n valori mai mari sau egale cu a_x ($n - 1 - x$ în A și $n - y = x + 1$ în B) (a_x este cea mai mare valoare în $A_0 \perp B_0$, pentru că $b_{y-1} \leq a_x$). Este de subliniat că valoarea a_x este mediana, și nu elementul a_x este cel median (dacă $a_{x_1} = a_{x_2}$, doar unul va fi în poziție mediană în $A \perp B$).

P5.3.7 Procesorul P_k compară a_k cu b_{n-1-k} și scrie în c_k 0 sau 1 după cum a_k este mai mic sau mai mare; C este un vector logic, de n elemente; în el se vor afla j valori de 0 consecutive, urmate de valori de 1, până la sfârșit. Mediana este a_{j-1} sau b_{n-j} și se va găsi în m .

1. dacă $a_k < b_{n-1-k}$ atunci $c_k \leftarrow 0$
2. altfel $c_k \leftarrow 1$
3. dacă $c_k \neq c_{k-1}$ atunci $m \leftarrow \min(a_{k-1}, b_{n-k})$

Al doilea **dacă** este executat de toate procesoarele, dar unul singur va face atribuirea ce urmează. Algoritmul se adaptează imediat pe un EREW PRAM.

P5.3.8 Da. Dacă, de exemplu, $a_{x_0} < b_{y_0}$ și $a_{x_1} = b_{x_1}$, atunci, după instrucțiunea (4), $d = 1, c_0 = 0, c_1 = 1$. În acest caz și P_0 , în instrucțiunea (6), și P_1 , în instrucțiunea (5), modifică la, ua, lb, ub . Și totuși algoritmul funcționează corect în continuare. Dacă vrem neapărat (sau pentru eleganță...) să evităm acest conflict procedăm astfel: mutăm instrucțiunea (5) imediat după (3), la același nivel; tot acolo adăugăm o instrucțiune **stop**, pentru că, în fond, căutarea s-a terminat.

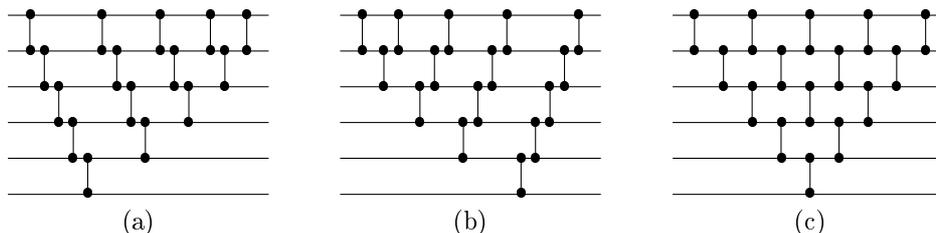


Figura 7.3: Rețele de sortare inspirate de (a) bubblesort, (b) sortarea prin inserție. (c) O redesenare arată că rețelele sunt identice.

P5.3.9 Fie $Z = X \perp Y = \langle z_0, \dots, z_{m+n-1} \rangle$. Demonstrăm prin inducție după indicele în Z . Argumentăm numai pentru X ; prin simetrie totul este valabil și pentru Y .

Dacă $z_0 = x_0$, atunci $rp(x_0) = 0$, deci $0 + rp(x_0) = 0$.

Dacă $z_k = x_l$ atunci fie $z_{k+1} = x_{l+1}$, și atunci $rp(x_{l+1}) = rp(x_l)$, fie $z_{k+1} = y_r$, și atunci $rp(x_l) = r$, $rp(y_r) = l + 1$.

P5.3.10 Demonstrăm numai primul caz. O k -separare constă de fapt în găsirea celor mai mici k elemente din $A \perp B$; acestea se află toate în $A_{[k]} \perp B_{[k]}$, pentru că a_i , cu $i \geq k$ are deja k elemente din A mai mici decât el, iar b_i , cu $i \geq k$ are k elemente din B mai mici. În plus, $A_{[k]} \perp B_{[k]}$ are $2k$ elemente.

P5.3.11 Se poate proceda în două feluri, teoretic echivalente. Primul procesor primește elementul celuilalt, îl compară cu al său, păstrează minimumul și trimite maximum celuilalt. Sau, procesoarele își transmit, pe rând, elementele și apoi amândouă fac comparația. Pentru operația cu secvențe, $exch(A_i, A_j)$, prima variantă este mai slabă.

P5.3.12 Da, dar cu două procesoare. Fiecare face o comparație și o atribuire, în paralel.

P5.4.1 Rețelele sunt prezentate în figura 7.3. Surpriza este că este vorba de aceeași rețea, așa cum se vede după redesenare. Se observă că, într-o rețea, un comparator poate glisa la stânga și la dreapta atâta vreme cât nu se "ciocnește" de un altul. Timpul de execuție este $2n - 3$, cu grad de paralelism cel mult $n/2$ și folosind $\binom{n}{2}$ comparatoare, adică minimum necesar.

P5.4.2 Cazul $p = \lfloor n/2 \rfloor$, pentru procesorul P_k :

1. pentru $i = 0 : n - 1$
 1. dacă k este par atunci $exch(2k, 2k + 1)$
 2. altfel dacă $k \neq \lfloor n/2 \rfloor$ atunci $exch(2k + 1, 2k + 2)$

Pentru cazul $p \ll n$ se pot scrie mai multe variante; utilizați principiul lui Brent pentru algoritmul de mai sus sau prin sortări locale urmate de $exch()$; comparați complexitățile.

P5.4.3 Desenele nu sunt identice; dacă ne gândim însă că o permutare de linii nu modifică proprietatea unei rețele de a sorta, rețelele sunt identice (a doua provine din prima prin inversarea liniilor—prin oglindire față de orizontală, dacă preferați).

P5.4.4 Acum este vorba despre oglindire față de verticală și aici nu e evident. Cum rețeaua are minimum de comparatoare, rezultă că pentru $Z = \langle n, \dots, 1 \rangle$ se fac interschimbări la fiecare comparație, rezultatul fiind $Z' = \langle 1, \dots, n \rangle$. Dacă se parcurge rețeaua în sens

invers, pornind de la Z' și schimbând sensul comparațiilor se obține Z , deci rețeaua inversată sortează descrescător, iar dacă se reface sensul comparațiilor, sortează crescător.

P5.4.5 Modificările sunt minore. Este vorba numai despre difuzări: ale parametrilor s , d , q , q' , a pivotului, imediat după alegerea lui, a poziției pivotului i_v , calculată de procedura *partiționare*. Deci $O(\log q)$ operații (neimplicând comparații), la fiecare iterație, ceea ce este ne semnificativ. În *partiționare*, suma prefixelor este deja pentru EREW PRAM.

P5.4.6 La iterația l ($l = d-1 : -1 : 0$), se alege $t = p/2^{d-l}$, ceea ce asigură că procesorul respectiv rămâne cu cel puțin $p/2^{d-l-1}$ elemente, din care va alege pivotul la următoarea iterație; în final poate rămâne fără nici un element, dar scopul – găsirea unui pivot – a fost atins. De menționat că celelalte cel puțin $p/2^{d-l-1}$ elemente sunt transmise unui procesor care va alege și el pivotul, începând din iterația următoare; nici acesta nu va rămâne fără elemente, până după ultima iterație. De obicei t se ține constant; aici am vrut doar să arătăm o limită inferioară; deci $t \geq p/2$.

P5.4.7 Nu. Cheile egale vor fi în final trimise aceluiași procesor. Evident că am făcut ipoteza că algoritmul 5.9 lucrează corect (vezi problema 5.2.5).

P5.4.8 Presupunem secvențele X_0, \dots, X_{p-1} concatenate într-o zonă contiguă de memorie X ; notăm i_0, \dots, i_{p-1} indicii de la care încep secvențele ($i_0 = 0$, $i_1 = n_0$, etc.).

1. pentru $j = 1 : \log p$
 1. pentru $l = 0 : p/2^j - 1$
 1. $s \leftarrow 2^j l$, $d \leftarrow 2^{j-1}(2l + 1)$

interclasează secvențele începând la i_s și i_d , de lungime n_s , respectiv n_d
 $i_s \leftarrow i_s + i_d$, $n_s \leftarrow n_s + n_d$

La fiecare iterație în j se pot face interclasările într-o zonă Y , de lungime n , care se copiază apoi în X (în fiecare interclasare în parte scriindu-se în Y începând de la valoarea i_s curentă). Oricum, e nevoie de n locații de memorie suplimentare.

P5.4.10 Fiecare procesor alege elementele din pozițiile $[im/p]$, $i = 1, 2$; deci P_0 : 16 și 27, P_1 : 13 și 23, P_2 : 10 și 22. Se ordonează aceste elemente: $Y = \langle 10, 13, 16, 22, 23, 27 \rangle$; din Y se aleg elementele din pozițiile $i(p-1)$, $i = 1, 2$, adică 16 și 23; acestea sunt cheile de partiționare. În final P_0 va avea 16 elemente, P_1 7 și P_2 13.

P5.4.11 Algoritmul pentru procesorul P_k este ($p = 2^d$):

1. formează secvența locală X
2. pentru $i = 0 : d-1$
 1. dacă $k_i = 1$ atunci
 1. **send**(X, i), **stop**
 2. **altfel** **recv**(X', i), $X \leftarrow X \perp X'$
3. **difuzare**: P_0 difuzează cheile de partiționare

După pasul i , rămân active numai procesoarele cu număr de ordine având ultimii $i+1$ biți egali cu 0. La fiecare pas mesajele transmise își dublează dimensiunea (presupunând că erau inițial egale). Numărul de operații aritmetice este $\sum_{i=0}^{d-1} 2^{d+1}(p-1) \approx 2p^2$; comunicația se face în timp $O(p^2)$, plus $O(p \log p)$ pentru difuzare. Față de algoritmul 5.27, timpul de sortare este redus de $d/2$ ori, pentru că acum toate procesoarele participă la această operațiune.

P5.4.13 Se folosește un algoritm de căutare binară; secvență unimodală este separată în două; se compară ultimul element din prima parte cu primul din a doua parte; se continuă căutarea numai în subsecvența care conține cel mai mic dintre aceste două elemente (o subsecvență a unei secvențe unimodale este și ea unimodală). De exemplu:

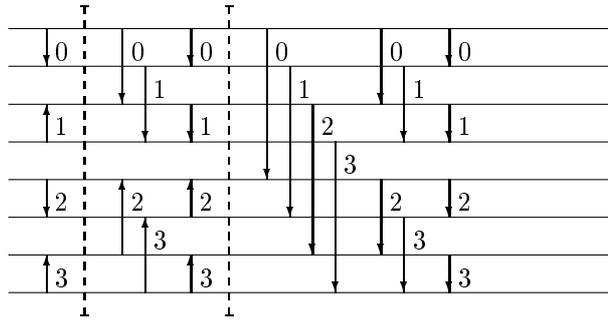


Figura 7.4: O planificare a procesoarelor pentru sortarea bitonică pe un PRAM.

1. $i \leftarrow 0, f \leftarrow m - 1$
2. **cât timp** $i \neq f$
 1. $j = \lfloor (i + f)/2 \rfloor$
 2. **dacă** $a_j < a_{j+1}$ **atunci** $f \leftarrow j$
 3. **altfel** $i \leftarrow j + 1$

Minimul se va găsi în a_j . În cazul unei secvențe bitonice, acest algoritm nu se poate aplica. Se folosește algoritmul secvențial clasic de căutare a minimului, eventual cu o condiție de terminare mai bună (care să profite de faptul că există un singur minim local). Oricum, în cel mai rău caz e nevoie de $n - 1$ comparații.

P5.4.14 Detaliem doar partea referitoare la interclasare; rezultatul va fi secvența sortată \tilde{A} ; număr de comparații: $m - 1$.

1. găsește j , indicele elementului minim
2. $\tilde{a}_0 \leftarrow a_j, l \leftarrow 1, i \leftarrow (j - 1) \bmod m, f \leftarrow (j + 1) \bmod m$
3. **cât timp** $i \neq f$
 1. $\tilde{a}_l \leftarrow \min(a_i, a_f), l \leftarrow l + 1, \tilde{f} \leftarrow \text{indmin}(a_i, a_f)$
 2. **dacă** $\tilde{f} = i$ **atunci** $i \leftarrow (i - 1) \bmod m$
 3. **altfel** $f \leftarrow (f + 1) \bmod m$

P5.4.15 Să mai desenăm o dată figura 5.9, amintindu-ne că fiecare săgeată reprezintă o comparație și deci este efectuată de un procesor (modelul este EREW PRAM); dacă asociem fiecărei săgeți numărul procesorului care efectuează comparația respectivă, rezultă figura 7.4 (evident că aceasta nu este singura planificare posibilă).

Ne vom inspira din algoritmul 5.28 în ce privește organizarea buclelor. Vă rămâne acum doar sarcina de a verifica buna funcționare a algoritmului ce urmează; k este numărul de ordine al procesorului, între 0 și $n/2$.

1. **pentru** $l = 0 : d - 1$
 1. **pentru** $i = l : -1 : 0$
 1. $j_1 \leftarrow k \bmod 2^i + 2^{i+1} \lfloor k/2^i \rfloor$
 2. $j_2 \leftarrow j_1 + 2^i$

3. dacă $k_i = 0$ atunci $exch(j_1, j_2)$ {ordonează crescător}
 4. altfel $exch(j_2, j_1)$

P5.4.16 Se împarte secvența A în $2p$ secvențe de $m = n/(2p)$ elemente; fiecare procesor sortează două astfel de secvențe. Apoi cele p procesoare aplică sortarea bitonică celor $2p$ secvențe, cu amendamentul descris pe larg pentru cazul memorie distribuită (ideea din algoritmul 5.30). Diferența majoră față de acest caz este că un singur procesor calculează $\min(A_i, A_j)$ și $\max(A_i, A_j)$, cu m comparații (în cazul distribuit două procesoare făceau fiecare câte m comparații). Numărul total de comparații este

$$2 \frac{n}{2p} \log \frac{n}{2p} + \sum_{i=0}^{d-1} \left(\frac{n}{p} + \sum_{i=0}^l \frac{n}{2p} \right) \approx \frac{n}{p} \log n + \frac{n}{p} (\log p - 1) + \frac{n}{4p} \log^2 p = O((n/p) \log n).$$

Comparând cu relația (5.7) se constată că aici coeficientul termenului în $(n/p) \log p$ (al doilea ca importanță) este de două ori mai mic.

Dacă aplicăm ideea algoritmului 5.31, de a obține subsecvențele minim și maxim prin interclasare, atunci numărul de comparații va fi

$$2 \frac{n}{2p} \log \frac{n}{2p} + \sum_{i=0}^{d-1} \sum_{i=0}^l 2 \frac{n}{2p} \approx \frac{n}{p} \log n - \frac{n}{p} + \frac{n}{2p} \log^2 p,$$

deci mai mare decât în cazul anterior; explicația e simplă: interclasarea necesită de două ori mai multe operații decât simpla extragere a minimului și maximului prin comparații între câte două elemente. De altfel, comparând cu (5.8), se constată aproape aceeași expresie, deci pe PRAM nu se reduce numărul de operații.

P5.4.17 Algoritmul este corect dacă demonstrăm că metoda de interclasare este corectă; o vom face în cazul general $|A| = n$, $|B| = m$. După sortare, secvența A va conține k zerouri, urmate de $n - k$ valori de unu; B va conține l zerouri urmate de $m - l$ de unu. Atunci $X = \langle x_0, x_1, \dots \rangle$ va începe cu $\lceil k/2 \rceil + \lceil l/2 \rceil$ zerouri, iar $Y = \langle y_0, y_1, \dots \rangle$ va începe cu $\lfloor k/2 \rfloor + \lfloor l/2 \rfloor$ zerouri. Diferența $(\lceil k/2 \rceil + \lceil l/2 \rceil) - (\lfloor k/2 \rfloor + \lfloor l/2 \rfloor)$ poate avea valorile 0, 1 sau 2 ($\lceil k/2 \rceil - \lfloor k/2 \rfloor$ poate fi 0 sau 1). Dacă este 0 sau 1, atunci secvența Z este deja sortată. Dacă este 2, atunci una din operațiile $exch(y_i, x_{i+1})$ va ordona Z . Cele trei cazuri sunt ilustrate de următorul exemplu:

$$\begin{array}{lll} A = \langle 0, 1, 1, 1, \dots \rangle & A = \langle 0, 0, 1, 1, \dots \rangle & A = \langle 0, 0, 0, 1, \dots \rangle \\ B = \langle 0, 1, 1, 1, \dots \rangle & B = \langle 0, 1, 1, 1, \dots \rangle & B = \langle 0, 1, 1, 1, \dots \rangle \end{array}$$

P5.4.18 Rețeaua de sortare este cea din figura 7.5 (secretul e de a desena dinspre dreapta spre stânga). Timpul de calcul este același cu cel al sortării bitonice, adică $\frac{1}{2} \log^2 n$; se folosesc mai puține procesoare decât la sortarea bitonică, dar gradul de paralelism este tot $n/2$; s-a demonstrat că numărul total de comparații folosite este asimptotic tot $O(n \log^2 n)$; este una dintre cele mai economicoase rețele de sortare. La implementarea pe o arhitectură MIMD cu memorie distribuită sunt importante gradul de paralelism (care e bun) și topologia de comunicație; dacă, la sortarea bitonică, pe un hipercub comunicau doar vecinii, acum, chiar dacă unele procesoare au timpi morți, există comunicație între procesoare aflate la orice distanță (și tocmai acestea sunt procesoarele care nu au timpi morți).

P5.4.19 Se presupune că A este o secvență conținând, în ordine, k de unu, l de zero și $n - k - l$ de unu (forma generală a unei secvențe unimodale cu minim formată din unu și

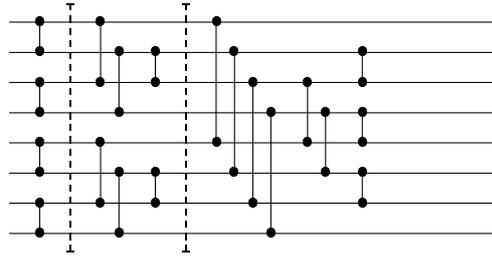


Figura 7.5: Alg. Batcher de sortare cu interclasare par-impară: rețea de sortare.

zero). Se pot întâlni patru cazuri, după paritățile lui k și l ; secvențele "pară" și "impară" rezultate sunt exemplificate mai jos:

$$\begin{array}{cccc}
 \langle 1, 1, 0, 1, 1, \dots \rangle & \langle 1, 1, 0, 0, 1, \dots \rangle & \langle 1, 1, 0, 0, 1, \dots \rangle & \langle 1, 1, 0, 0, 1, \dots \rangle \\
 \langle 1, 0, 0, 1, 1, \dots \rangle & \langle 1, 0, 0, 1, 1, \dots \rangle & \langle 1, 1, 0, 1, 1, \dots \rangle & \langle 1, 1, 0, 0, 1, \dots \rangle
 \end{array}$$

După sortare, numărul de zerouri inițiale din secvențele "pară" și "impară" diferă prin cel mult 1, deci una dintre operațiile $exch()$ va ordona secvența A .

Desenul rețelei de sortare se face, ca și în cazul sortării bitonice, înapoi (spre stânga) pe etape și înainte în cadrul unei etape (etapele sunt separate de linii punctate, în figura 7.4).

P5.4.20 Presupunând $n = 2^m$, suma va fi egală cu $\log m!$. Dar $m^{m/2} < m! < m^m$, și deci $\log m! = O(m \log m)$. Cum $m = \log n$, se obține rezultatul cerut.

Bibliografie

- [1] B. Abali, F. Özgüner, and A. Bataineh. Balanced Parallel Sort on Hypercube Multiprocessors. *IEEE Trans. Par. Distrib. Syst.*, 4(5):572–581, 1993.
- [2] K.E. Batcher. Sorting networks and their application. In *Proc. AFIPS 1968 SICC*, pages 307–314, 1968.
- [3] D.P. Bertsekas, C. Özveren, G.D. Stamoulis, P. Tseng, and J.N. Tsitsiklis. Optimal Communication Algorithms for Hypercubes. *J.Par.Distrib.Comput.*, 11:263–275, 1991.
- [4] D.P. Bertsekas and J.N. Tsitsiklis. *Parallel and distributed computation*. Prentice-Hall, 1989.
- [5] P. Bjørstad, F. Manne, T. Sørevik, and M. Vajteršic. Efficient Matrix Multiplication on SIMD Computers. *SIAM J. Matrix Anal. Appl.*, 13(1):386–401, 1992.
- [6] T. Brown and R. Xiong. A Parallel quicksort Algorithm. *J. Par. Distrib. Comput.*, 19:83–89, 1993.
- [7] R.M. Chamberlain. An Algorithm for LU Factorization with partial pivoting on a hypercube. Technical Report CCS 86/11, Dept. of Science and Technology, Chr. Michelson Institute, Bergen, Norway, June 1986.
- [8] E.G. Coffman and P.J. Denning. *Operating Systems Theory*. Prentice-Hall, 1973.
- [9] D. Coppersmith and S. Winograd. Matrix multiplication via arithmetic progressions. In *19 th. Annual ACM Symp. Theory Comp.*, pages 1–6, 1987.
- [10] J.J. Dongarra, J. Du Croz, S. Hammarling, and I. Duff. A Set of Level-3 Basic Linear Algebra Subprograms. *ACM Trans.Math.Software*, 16:1–17,18–28, 1990.
- [11] B. Dumitrescu, J-L. Roch, and D. Trystram. Fast matrix multiplication algorithms on mimd architectures. *Parallel Algorithms and Applications*, 4(2):53–70, 1994.
- [12] G. Fox et al. *Solving problems on concurrent processors*. Prentice-Hall, 1988.
- [13] G.H. Golub and C.F. Van Loan. *Matrix Computations. Second edition*. The John Hopkins University Press, 1989.
- [14] M.T. Heath and C.H. Romine. A Parallel Triangular Solver for a Distributed-memory Multiprocessor. *SIAM J. Sci. Stat. Comput.*, 9(3):558–587, May 1988.
- [15] R.W. Hockney and C.R. Jesshope. *Parallel Computers 2*. Adam Hilger Ltd., 1988.
- [16] R.M. Karp and V. Ramachandran. Parallel Algorithms for Shared-Memory Machines. In J. van Leeuwen, editor, *Handbook of Theoretical Computer Science*. Elsevier Science Publishers, 1990.

- [17] D.E. Knuth. *The Art of Computer Programming. Vol.3, Sorting and Searching.* Addison-Wesley, 1973.
- [18] C.P. Kruskal. Searching, merging, and sorting in parallel computation. *IEEE Trans.Comput.*, 32(10):942–946, 1983.
- [19] S. Lennart Johnsson and C.T. Ho. Optimum Broadcasting and Personalized Communication in Hypercubes. *IEEE Trans.Comput.*, 38(9):1249–1268, 1989.
- [20] G. Li and T.F. Coleman. A Parallel Triangular Solver for a Distributed-memory Multiprocessor. *SIAM J. Sci. Stat. Comput.*, 9(3):485–502, May 1988.
- [21] G. Li and T.F. Coleman. A New Method for Solving Triangular Systems on Distributed-memory Message-passing Multiprocessor. *SIAM J. Sci. Stat. Comput.*, 10(2):382–396, March 1989.
- [22] B. Monien and H. Sudborough. Embedding one Interconnection Network in Another. *Computing Suppl.*, 7:257–282, 1990.
- [23] Y. Saad and M.H. Schultz. Data Communication in Parallel Architectures. *Parallel Computing*, 11:131–150, 1989.
- [24] H. Shi and J. Schaeffer. Parallel Sorting by Regular Sampling. *J. Parallel Distrib. Comput.*, 14:361–372, 1992.
- [25] Q.F. Stout and B. Wagar. Intensive Hypercube Communications - Prearranged Communication in Link-Bound Machines. *J.Par.Distrib.Comput.*, 10:167–181, 1990.
- [26] B. Wagar. Hyperquicksort. In *Hypercube Multiprocessors 1987, Greece*, pages 292–299. SIAM, 1987.
- [27] R. Xiong and T. Brown. Parallel Median Splitting and k -splitting with Application to Merging and Sorting. *IEEE Trans.Par.Distrib.Syst.*, 4(5):559–565, 1993.