



UMCS

UNIWERSYTET MARII CURIE-SKŁODOWSKIEJ
W LUBLINIE

Wydział Matematyki, Fizyki i Informatyki

Kierunek: Matematyka
Specjalność: Informatyczna

Rafał Szczerski
nr albumu: 2524659

Równoległe algorytmy mnożenia macierzy Parallel matrix multiplication algorithms

Praca licencjacka
napisana w Zakładzie Informatyki
pod kierunkiem dr Beaty Byliny

Lublin rok 2015

Spis treści

1	Wstęp	4
2	Wiadomości wstępne	5
2.1	Ustalenia terminologiczne	5
2.2	Klasyfikacja algorytmów	8
2.2.1	Reprezentacja algorytmów	12
2.2.2	Przyspieszenie	13
2.2.3	Prawo Amdahla	15
2.2.4	Prawo Gustafsona i Barsisa	16
2.3	Teoretyczne modele obliczeń	18
2.3.1	Model RAM	19
2.3.2	Model PRAM	22
2.3.3	Model sieciowy	25
3	Klasyczne algorytmy mnożenia macierzy	29
3.1	Algorytm naiwny	29
3.2	Algorytm „dziel i rządź”	30
3.3	Algorytm Cannona	30
3.4	Algorytm Strassena	31
3.5	Algorytm Strassena-Winograda	31
4	Równoległe algorytmy mnożenia macierzy	33
4.1	Algorytm w modelu PRAM	33
4.2	Algorytmy w modelu sieciowym	36
4.2.1	Algorytm w dwuwymiarowym torusie	36
4.2.2	Algorytm systoliczny w sieci dwuwymiarowej	37
4.2.3	Algorytm w topologii hipersześcianu	38
4.2.4	Algorytm Cannona w dwuwymiarowej sieci	40
4.2.5	Algorytm 2.5D	40

4.2.6	Algorytm 3D	40
4.2.7	Równoległy algorytm Strassena CAPS	40

Rozdział 1

Wstęp

Rozdział 2

Wiadomości wstępne

2.1 Ustalenia terminologiczne

Definicja 2.1 (Zbiór przechodni). Zbiór A nazywamy **przechodnim**, wtedy i tylko wtedy, gdy $\forall x (x \in A \wedge y \in x \implies y \in A)$.

Definicja 2.2 (Domknięcie przechodnie zbioru). Domknięciem przechodnim zbioru X nazywamy najmniejszy w sensie inkluzji zbiór przechodni, który zawiera X .

Definicja 2.3 (Graf skierowany (DG)). Powiedzmy, że:

1. $V \neq \emptyset$ jest zbiorem
2. $E \subseteq V \times V$

Grafem skierowanym G nazwiemy dwójkę (V, E) .

Definicja 2.4 (Acykliczny graf skierowany (DAG)). Acyklicznym grafem skierowanym nazywamy graf skierowany nie zawierający cykli.

Definicja 2.5 (Domknięcie przechodnie grafu). Niech $G = (V, A)$ będzie grafem skierowanym. Graf skierowany $G^+ = (V, A^+)$ nazywamy **domknięciem przechodnim** grafu G , gdy A^+ jest zbiorem wszystkich takich par (a, b) wierzchołków zbioru V , że w grafie G istnieje droga z a do b .

Definicja 2.6 (Graf zależności). Niech dane będą zbiór $S \neq \emptyset$, relacja przechodnia $R \subseteq S \times S$. **Grafem zależnościami** nazywamy graf $G = (S, T)$ i $T \subseteq R$, gdzie R jest przechodnim domknięciem T .

Definicja 2.7 (Ścieżka). **Ścieżką** łączącą v_0 z v_n o długości n nazywamy ciąg wierzchołków (v_0, v_1, \dots, v_n) taki, że dla każdego $k \in \{0, 1, \dots, n-1\}$ istnieje krawędź z v_k do v_{k+1} .

Definicja 2.8 (Droga). **Drogą** w grafie G nazywamy ścieżkę, której wierzchołki są różne.

Definicja 2.9 (Długość drogi). **Długością** drogi w grafie G nazywamy liczbę krawędzi, które zawiera droga.

Definicja 2.10 (Cykl). Drogę zamkniętą długości co najmniej 1 z ciągiem wierzchołków $x_1 x_2 \dots x_n x_1$ nazywamy **cyklem**, jeśli wszystkie wierzchołki $x_1, x_2 \dots x_n$ są różne.

Definicja 2.11 (Stopień wierzchołka). **Stopień** $d_G(v)$ **wierzchołka** v definiujemy jako liczbę incydentnych z v krawędzi. Każdemu wierzchołkowi v grafu

skierowanego G możemy przypisać stopień wyjściowy (ang. *indegree*) $d_G^+(v)$ i stopień wejściowy (ang. *outdegree*) $d_G^-(v)$:

$$d_G^+(v) = \#\{w \mid (v, w) \in E\}$$

$$d_G^-(v) = \#\{w \mid (w, v) \in E\}$$

Definicja 2.12 (Macierz). Niech \mathbb{K} będzie ciałem. Macierzą o m wierszach i n kolumnach i wartościami w \mathbb{K} (krótko: macierzą $m \times n$) nazwamy każde odwzorowanie $A : \{1, \dots, m\} \times \{1, \dots, n\} \rightarrow \mathbb{K}, (i, j) \mapsto A_{ij}$

2.2 Klasyfikacja algorytmów

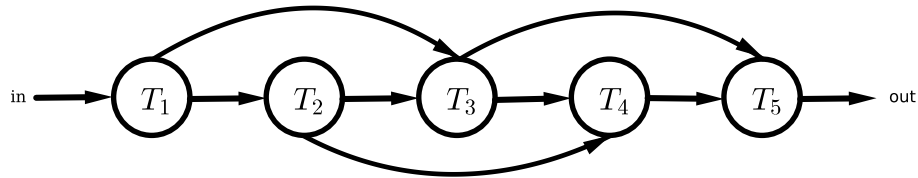
Definicja 2.13 (Algorytm). Zbiór jednoznacznie określonych reguł lub zadań obliczeniowych prowadzących w skończonej ilości kroków do rozwiązania pewnego problemu [1].

Określone w ten sposób zadania obliczeniowe są z reguły względem siebie niezależne. Pewne z nich mogą być wykonywane równolegle, inne – sekwencyjnie, jedno po drugim. Wobec tego algorytm może być określony częściowo równolegle, częściowo sekwencyjnie.

Na podstawie niezależności zadań obliczeniowych algorytmy możemy podzielić na pięć klas [2]:

1. Algorytmy szeregowo
2. Algorytmy równoległe
3. Algorytmy szeregowo-równoległe (SPA, Serial-Parallel Algorithms)
4. Algorytmy nieszeregowo-równoległe (NSPA, Nonserial-Parallel Algorithms)
5. Algorytmy regularno-iteracyjne (RIA, Regular-Iterative Algorithms)

Definicja 2.14 (Algorytm sekwencyjny). **Algorytm sekwencyjny** (rys. 2.1) jest ciągiem dokładnie sprecyzowanych zadań obliczeniowych T_i , $i \in \mathbb{N}$ rozwiązujących dany problem, tj. wyznaczających dane wyjściowe na podstawie danych wejściowych. Zakłada się, że w algorytmie sekwencyjnym zadania wykonywane są przez jeden procesor.

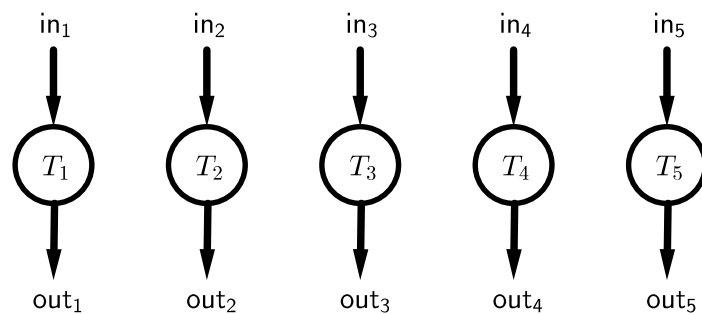


Rysunek 2.1: Algorytm sekwencyjny

W celu rozwiązania problemu za pomocą większej liczby procesorów należy go zdekomponować na podproblemy, które mogą być rozwiązane równoległe. Każdy z podproblemów rozwiązywany jest przez odrębny algorytm będący składową algorytmu równoległego.

Definicja 2.15 (Równoległość). **Równoległość** w odniesieniu do oprogramowania jest to symultaniczny transfer, występowanie albo przetwarzanie poszczególnych części pewnej całości, takich jak bity składające się na znak albo znaki pewnego słowa, używając osobnych urządzeń dla ich różnych części [1].

Definicja 2.16 (Algorytm równoległy). **Algorytmem równoległym** (rys. 2.2) nazywamy każdy algorytm w którym spośród określonych w nim zadań T_1, T_2, \dots, T_n co najmniej dwa zadania $T_i, T_j, i \neq j$ dzięki ich wzajemnej niezależności, mogą być wykonane równocześnie [2].



Rysunek 2.2: Algorytm równoległy

Definicja 2.17 (Architektura równoległa). **Architektura równoległa** jest to architektura wieloprocessorowa, na której można wykonywać przetwarzanie równoległe [1].

Algorytmy równoległe i architektury równoległe są ze sobą blisko spokrewnione. Równoległość może być zaimplementowana na wielu poziomach używając technik sprzętowych i programowych[2]

1. Równoległość na poziomie danych (*Data-level parallelism*), gdzie pracujemy na wielu bitach danych lub na wielu danych jednocześnie.
2. Równoległość na poziomie instrukcji (*Instruction-level parallelism*, ILP), gdzie jednocześnie procesor może wykonać więcej niż jedną instrukcję.
3. Równoległość na poziomie wątków (*Thread-level parallelism*, TLP). Wątkiem jest część programu, która współdzieli zasoby procesora z innymi wątkami. W TLP wiele programowych wątków jest uruchamianych jednocześnie na jednym bądź wielu procesorach.

4. Równoległość na poziomie procesów (*Process-level parallelism*). Proces to program, który jest uruchomiany na komputerze. Rezerwuje on własne zasoby komputera, takie jak przestrzeń pamięciową i rejestry.

Przykład 2.1. Prostym przykładem algorytmu równoległego jest serwer sieciowy, który każde zapytanie przychodzące przetwarza niezależnie od innych zapytań. Innym przykładem są wielozadaniowe systemy operacyjne, radzące sobie z jednoczesną obsługą kilku uruchomionych programów.

2.2.1 Reprezentacja algorytmów

2.2.2 Przyspieszenie

Potencjalną korzyść z równoległego wykonania zadania obliczeniowego możemy zmierzyć licząc czas jaki zajmuje wykonanie go na jednym procesorze i porównanie wyniku z wykonaniem tego samego zadania równolegle na N procesorach. Współczynnik przyspieszenia $S(p, n)$ możemy zdefiniować jako

$$(2.1) \quad S(p, n) = \frac{T_p^*(1)}{T_p(N)}$$

gdzie $T_p^*(1)$ jest pesymistyczną złożonością czasową najszybszego znanego algorytmu sekwencyjnego rozwiązującego problem Z na jednym procesorze, $T_p(N)$ jest pesymistyczną złożonością algorytmu równoległego. Wyrażenie 2.1 nazywamy **przyspieszeniem bezwzględnym**.

Maksymalną wartością przyspieszenia $S(p, n)$ jest p , ponieważ używając p procesorów można przyspieszyć obliczenia najlepszego algorytmu sekwencyjnego co najwyżej p razy. Zwykle uzyskiwane przyspieszenie jest mniejsze niż p . Przyczyną tego może być niewystarczający stopień współbieżności w problemie Z .

W złożoności $T_1^r(n)$ można wyróżnić operacje obliczeniowe, które muszą być wykonane sekwencyjnie, $T_1^s(n)$, oraz obliczenia, które mogą być wykonane równolegle, $T_1^r(n)$. Wobec tego mamy, że:

$$(2.2) \quad T_1(n) = T_1^s(n) + T_1^r(n)$$

Zakładając, że obliczenia $T^r(n)$ da się równomiernie rozdzielić między p procesorami, przyspieszenie $S(p, n)$ wyraża się wzorem:

$$(2.3) \quad S(p, n) = \frac{T_1(n)}{T_p(n)} \leq \frac{T_1^s(n) + T_1^r(n)}{T_1^s(n) + T_1^r(n)/p + T_p^o(n)}$$

gdzie $T_p^o(n)$ jest złożonością dodatkową wynikającą z organizacji obliczeń równoległych. W jej skład wchodzi m.in. operacje komunikacji między procesorami.

2.2.3 Prawo Amdahla

2.2.4 Prawo Gustafsona i Barsisa

Niech p oznacza liczbę procesorów, σ – część czasu obliczeń algorytmu równoległego przypadającą na wykonanie obliczeń w sposób sekwencyjny, a ρ – część czasu obliczeń algorytmu równoległego przypadającą na wykonywanie obliczeń w sposób równoległy takie, że $\sigma + \rho = 1$. Czas wykonania tego samego algorytmu w hipotetycznym komputerze sekwencyjnym jest proporcjonalny do sumy $\sigma + p\rho$, gdzie wyrażenie $p\rho$ odpowiada czasowi wykonania części równoległej obliczeń przez jeden procesor. Przyspieszenie, które zostałoby uzyskane, gdyby obliczenia równoległe zostały przeprowadzone w komputerze sekwencyjnym wyraża się przez:

$$(2.4) \quad \Psi(p, n) \leq \frac{\sigma + p\rho}{\sigma + \rho} = \sigma + p\rho = \sigma + p(1 - \sigma) = p + (1 - p)\sigma$$

Wzór (2.4) jest znany jako **prawo Gustafsona i Barsisa**. Definiuje ono tzw. **skalowane przyspieszenie**, ponieważ wraz ze zmianą liczby procesorów skaluje się odpowiednio rozmiar problemu, tak aby utrzymać stały czas obliczeń równoległych [3].

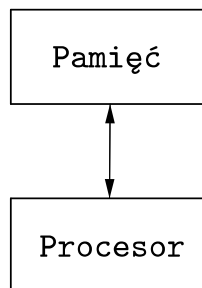
2.3 Teoretyczne modele obliczeń

2.3.1 Model RAM

Nim przejdziemy do omówienia modeli obliczeń równoległych zajmiemy się omówieniem modelu RAM zwanego również architekturą von Neumanna.

Model RAM (*Random Access Machine*) odpowiada rozważaniom zawartym w ???. Zakłada on:

1. Istnienie pewnego procesora wyposażonego w:
 - (a) skończoną listę instrukcji, które może on realizować
 - (b) pewną liczbę rejestrów arytmetycznych procesora R_1, R_2, \dots, R_n , $n > 1$ które mogą przechowywać dowolne skończone liczby w zapisie binarnym
 - (c) specjalny rejestr sterujący L zwany licznikiem programu.
2. Istnienie pamięci złożonej z potencjalnie nieskończonej liczby komórek M_i , $i = 1, 2, 3, \dots$ (Rys. 2.3) w których można przechowywać dowolną skończoną liczbę w zapisie binarnym.
3. Stały czas zapisu i odczytu wartości do/z komórki pamięci (inaczej *dostęp swobodny*).



Rysunek 2.3: Model obliczeń sekwencyjnych RAM

Tabela 2.1: Przykładowa lista instrukcji procesora[3]

Instrukcja	Argument	Znaczenie
LOAD	k, a	$R_k := w(a)$
STORE	k, b	$M_{w(b)} := R_k$
ADD	k, c	$R_k := R_k + w(c)$
SUB	k, c	$R_k := R_k - w(c)$
MULT	k, c	$R_k := R_k \times w(c)$
DIV	k, c	$R_k := \lfloor R_k / w(c) \rfloor$
JUMP	i	$L := i$
JPOS	k, i	if $R_k > 0$ then $L := i$ else $L := L + 1$
JZERO	k, i	if $R_k == 0$ then $L := i$ else $L := L + 1$
JNEG	k, i	if $R_k < 0$ then $L := i$ else $L := L + 1$
READ	k	Wczytaj daną z urządzenia zewnętrznego do rejestru R_k
WRITE	k	Wydrukuj daną z rejestru R_k
HALT		Zakończ obliczenie

2.3.2 Model PRAM

Model wspólnej pamięci składa się z pewnej liczby procesorów, z których każdy posiada własną pamięć i może lokalnie wykonywać programy. Wszystkie procesory mogą komunikować się za pomocą wspólnej globalnej pamięci. Każdemu procesorowi przyporządkowany jest niepowtarzająca się liczba naturalna. Jest to lokalnie dostępny indeks, numer procesora lub jego identyfikator.

W modelu wspólnej pamięci wyróżniamy dwa podstawowe tryby operacji. W pierwszym trybie, synchronicznym, wszystkie procesory działają synchronicznie według wspólnego zegara. Model ten nazywamy równoległą maszyną o dostępie swobodnym (PRAM, parallel random-access machine).

W drugi trybie, asynchronicznym, każdy procesor pracuje według osobnego zegara. W tym trybie programista jest odpowiedzialny za odpowiednią synchronizację procesorów, jeśli zachodzi taka potrzeba. Dokładniej mówiąc, jeśli procesor ma pobrać dane, to odpowiedzialnością programisty jest upewnienie się, że odpowiednie dane są już uzyskane, ponieważ wartości wspólnych zmiennych są określane dynamicznie w trakcie wykonania programu na różnych procesorach.

Ponieważ każdy procesor może uruchomić swój program lokalnie, ten model jest typu MIMD w klasyfikacji Flynna. Znaczy to tyle, że każdy procesor może wykonać pewną instrukcję lub operację na danych niezależnie od tych wykonanych na jakimkolwiek innym procesorze w trakcie danej jednostki czasu.

Dla danego algorytmu, rozmiar danych wymienionych pomiędzy pamięcią globalną i pamięcią lokalną różnych procesorów wyraża rozmiar **komunikacji** wymaganej przez algorytm.

Możemy wyróżnić kilka wariantów modelu PRAM w zależności od wymagań jakie postawimy odnośnie jednoczesnego dostępu kilku procesorów do tego samego adresu w pamięci globalnej.

Definicja 2.18. Klasyfikacja PRAM ze względu na dostęp do pamięci

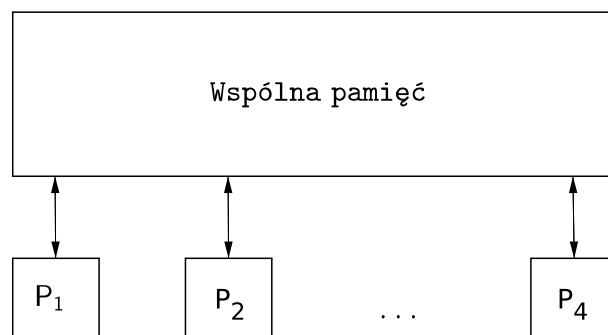
- **EREW** – algorytmy z wyłącznym odczytem i wyłącznym zapisem; nie pozwala na jednoczesny zapis do pamięci
- **CREW** – algorytmy z jednoczesnym odczytem i wyłącznym zapisem;

pozwała na jednoczesny dostęp do pamięci dla instrukcji odczytu

- **CRCW** – algorytmy z jednoczesnym odczytem i jednoczesnym zapisem;
- **ERCW** – algorytmy z wyłącznym odczytem i jednoczesnym zapisem.

Jeśli nie poczyni się żadnych dodatkowych założeń, to nie jest jasno określone, co zostanie zapisane w komórce pamięci w wyniku jednoczesnego zapisywania do niej przez wiele procesorów w algorytmie typu CRCW. W literaturze można spotkać wiele typów maszyny PRAM, które różnią się sposobami rozwiązywania konfliktów zapisu. Można wśród nich wyróżnić[4]:

1. jednolity (ang. common) – procesory muszą zapisać do tej samej komórki pamięci jednolitą wartość
2. dowolny (ang. arbitrary) – zapamiętywana jest dowolna wartość z wartości zapisywanych do tej samej komórki pamięci
3. priorytetowy (ang. priority) – zapamiętywana jest wartość zapisywana przez procesor o najmniejszym numerze
4. (ang. combining) – zapamiętywana jest wartość jest pewną, jednak ściśle określoną kombinacją zapisywanych wartości



Rysunek 2.4: Model wspólnej pamięci

2.3.3 Model sieciowy

Sieć możemy przedstawić modelowo jako graf $G = (N, E)$, gdzie każdy węzeł $i \in N$ oznacza procesor, a każda krawędź $(i, j) \in E$ – dwukierunkową komunikację między procesorami i i j . Przyjmujemy, że każdy procesor ma swoją lokalną pamięć i nie ma żadnej pamięci współdzielonej przez procesory. Tak jak w przypadku modelu z pamięcią wspólną, operacje w sieci mogą być synchroniczne lub asynchroniczne.

W opisie algorytmów dla modelu sieciowego potrzebujemy zdefiniować dwie instrukcje do opisanania komunikacji między procesorami.

1. $\text{send}(X, i)$
2. $\text{receive}(X, j)$

Procesor P wykonujący instrukcję **send** wysyła kopię X do procesora P_i , następnie natychmiast przechodzi do wykonywania kolejnej instrukcji. Procesor P wykonujący instrukcję **receive** zatrzymuje wykonanie programu aż do chwili, gdy otrzyma dane z procesora P_j , a następnie przechowuje dane w Y i kontynuuje wykonanie programu.

Procesory pracujące w sieci asynchronicznej zarządzają swoimi zadaniami przez wymianę komunikatów. Schemat taki nazywamy modelem wymiany komunikatów. Procesory te niekoniecznie muszą być ze sobą sąsiadujące.

Charakteryzuje ją kilka parametrów:

1. średnica – maksymalna odległość (krawędziowa) między dowolną parą węzłów; im mniejsza, tym lepiej.
2. maksymalny stopień wierzchołka – maksymalna liczba łączy do dane procesora
3. szerokość połowienia sieci – minimalna liczba krawędzi, które muszą zostać usunięty, aby podzielić ją na dwie równe podsieci
4. spójność krawędziowa – minimalna liczba krawędzi, które muszą ulec awarii, aby sieć stała się niespójna
5. koszt sieci – koszt wykonania, zarządzania i utrzymania połączeń między procesorami; w najprostrzym przypadku mierzony liczbą krawędzi

Sieć liniowa

Definicja 2.19 (Sieć liniowa). Model składa się z p procesorów P_1, P_2, \dots, P_p połączonych ze sobą w ciąg, tzn. procesor P_i połączony jest z procesorem P_{i-1} i P_{i+1} , o ile takie istnieją. Średnica takiej sieci wynosi $p - 1$, jej maksymalny stopień wynosi 2.

Definicja 2.20 (Torus). Sieć liniowa z połączonymi końcami.

Sieć dwuwymiarowa

Dwuwymiarowa sieć jest dwuwymiarową wersją sieci liniowej. Składa się ona z $p = m^2$ procesorów ułożonych w siatkę $m \times m$ taką, że procesor $P_{i,j}$ jest połączony z procesorem $P_{i\pm 1,j}$ i $P_{i,j\pm 1}$.

Średnica takiej sieci złożonej z $p = m^2$ procesorów wynosi \sqrt{p} a jej maksymalny stopień 4

Sieć hipersześcienna

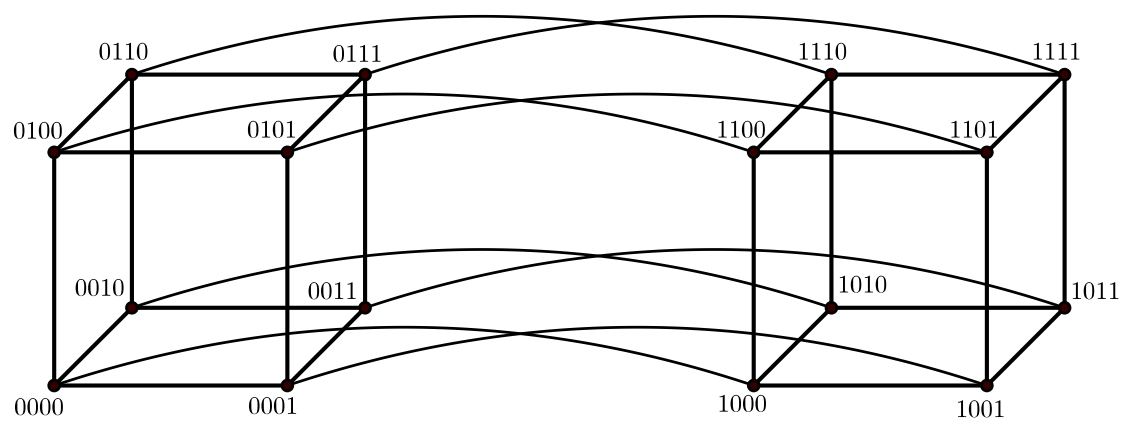
Definicja 2.21. Kostka Boola

Niech $i_{d-1}i_{d-2} \dots i_0$, gdzie $0 \leq i \leq p - 1$ będzie binarną reprezentacją i . Wówczas procesor i jest połączony z procesorem $P_{i^{(j)}}$, gdzie $i^{(j)} = i_{d-1} \dots \bar{i}_j \dots i_0$ i $\bar{i}_j = 1 - i_j$. Innymi słowy, dwa procesory są ze sobą połączone wtedy i tylko wtedy, gdy ich wskaźniki różnią się tylko jednym bitem.

Sieć w topologii hipersześcianu składa się z $p = 2^d$ procesorów połączonych w d -wymiarową kostkę Boola.

Hipersześcian ma strukturę rekursywną. Kostkę d -wymiarową możemy rozszerzyć do $d + 1$ wymiarów przez połączenie poszczególnych procesorów do d -wymiarowych kostek.

Średnica d -wymiarowego hipersześcianu wynosi $d = \log p$. Jest tak ponieważ odległość w grafie między dwoma procesorami P_i i P_j jest równa liczbie pozycji bitów, którymi wskaźniki i i j różnią się między sobą. Stąd jest ona mniejsza lub równa d , a ponadto odległość między P_0 a P_{2^d-1} wynosi d . Każdy węzeł jest stopnia $d = \log p$.



Rysunek 2.5: Sieć w topologii hipersześcianu

Rozdział 3

Klasyczne algorytmy mnożenia macierzy

W swojej pracy „*Gaussian Elimination is not Optimal*” z 1969 roku VOLKER STRASSEN pokazał rekursywny algorytm mnożenia macierzy kwadratowych wymiaru $m2^k$ o złożoności $\mathcal{O}(n^{2.81})$ [5]. W przypadku macierzy 2×2 oznaczało to, że mnożenie można wykonać już za pomocą 7 mnożeń i 18 dodawań. Algorytm za sprawą SHMUELA WINOGRADA został zoptymalizowany[6][7] do najczęściej implementowanej dzisiaj postaci algorytmu STRASSENA-WINOGRADA (implementacja zawiera się na przykład w bibliotece GEMMW[8]). W przypadku macierzy 2×2 wykonuje on 7 operacji mnożenia i 15 dodawań[9].

3.1 Algorytm naiwny

Niech $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$. Rozważmy algorytm sekwencyjny wyznaczania macierzy $\mathbf{C} = \mathbf{AB}$ o złożoności $\mathcal{O}(n^3)$.

```
for  $i \in \{1, 2, \dots, n\}$  do
  for  $j \in \{1, 2, \dots, n\}$  do
     $\mathbf{C}(i, j) = 0$ ;
    for  $k \in \{1, 2, \dots, n\}$  do
       $\mathbf{C}(i, j) = \mathbf{C}(i, j) + \mathbf{A}(i, k)\mathbf{B}(k, j)$ ;
    end
  end
end
```

Algorithm 1: Sekwencyjny algorytm mnożenia macierzy

3.2 Algorytm „dziel i rządź”

Dla danych macierzy wejściowych $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$ oraz macierzy wyjściowej $\mathbf{C} \in \mathbb{R}^{n \times n}$ mamy:

$$\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} \quad \mathbf{C} = \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix}$$

gdzie

$$(3.1) \quad \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix}$$

(3.1) możemy wyrazić inaczej:

$$(3.2) \quad \begin{aligned} C_{11} &= A_{11}B_{11} + A_{12}B_{21} \\ C_{12} &= A_{11}B_{12} + A_{12}B_{22} \\ C_{21} &= A_{21}B_{11} + A_{22}B_{21} \\ C_{22} &= A_{21}B_{12} + A_{22}B_{22} \end{aligned}$$

Z powyższych konstatacji nasuwa się łatwy algorytm rekurencyjny.

multiply(\mathbf{A}, \mathbf{B});

Dane wejściowe : Macierze wejściowe \mathbf{A}, \mathbf{B} wymiaru n

Dane wyjściowe : Macierz \mathbf{C}

begin

if $n=1$ **then**

$C_{11} = A_{11}B_{11}$

else

$C_{11} = \text{multiply}(A_{11}B_{11}) + \text{multiply}(A_{12}B_{21});$

$C_{12} = \text{multiply}(A_{11}B_{12}) + \text{multiply}(A_{12}B_{22});$

$C_{21} = \text{multiply}(A_{21}B_{11}) + \text{multiply}(A_{22}B_{21});$

$C_{22} = \text{multiply}(A_{21}B_{12}) + \text{multiply}(A_{22}B_{22});$

end

return C

end

3.3 Algorytm Cannona

[10]

3.4 Algorytm Strassena

Niech A i B będą macierzami $m2^k \times m2^k$. Definiując następujące macierze pomocnicze

$$\begin{aligned} H_1 &= (A_{11} + A_{22})(B_{11} + B_{22}) & H_2 &= (A_{21} + A_{22})B_{11} \\ H_3 &= A_{11}(B_{12} + A_{22}) & H_4 &= A_{22}(B_{21} + A_{11}) \\ H_5 &= (A_{11} + A_{12})B_{22} & H_6 &= (A_{21} + A_{11})(B_{11} + B_{12}) \\ H_7 &= (A_{12} - A_{22})(B_{21} + B_{22}) \end{aligned}$$

otrzymujemy

$$(3.3) \quad C = \begin{bmatrix} H_1 + H_4 - H_5 + H_7 & H_3 + H_5 \\ H_2 + H_4 & H_1 + H_3 - H_2 + H_6 \end{bmatrix}$$

3.5 Algorytm Strassena-Winograda

Dla danych macierzy wejściowych \mathbf{A}, \mathbf{B} oraz macierzy wyjściowej \mathbf{C} mamy

$$\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} \quad \mathbf{C} = \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix}$$

Następnie ustalmy odpowiednio po siedem kombinacji liniowych T_i, S_i , $i \in \{1, 2, \dots, 7\}$ dla każdej z podmacierzy \mathbf{A} i \mathbf{B} .

$$\begin{aligned} T_0 &= A_{11} & S_0 &= B_{11} \\ T_1 &= A_{12} & S_1 &= B_{21} \\ T_2 &= A_{21} + A_{22} & S_2 &= B_{12} + B_{11} \\ T_3 &= T_2 - A_{12} & S_3 &= B_{22} - S_2 \\ T_4 &= A_{11} - A_{21} & S_4 &= B_{22} - B_{12} \\ T_5 &= A_{12} + T_3 & S_5 &= B_{22} \\ T_6 &= A_{22} & S_6 &= S_3 - B_{21} \end{aligned}$$

oraz

$$\begin{array}{ll} Q_0 = T_0 S_0 & U_1 = Q_0 + Q_3 \\ Q_1 = T_1 S_1 & U_2 = U_1 + Q_4 \\ Q_2 = T_2 S_2 & U_3 = U_1 + Q_2 \\ Q_3 = T_3 S_3 & C_{11} = Q_0 + Q_1 \\ Q_4 = T_4 S_4 & C_{12} = U_3 + Q_5 \\ Q_5 = T_5 S_5 & C_{21} = U_2 - Q_6 \\ Q_6 = T_6 S_6 & C_{22} = U_2 + Q_2 \end{array}$$

Jest to jeden krok metody Strassena-Winograda. Algorytm jest rekursywny ponieważ może być użyty ponownie dla wyznaczenia Q_i , $i \in \{0, 1, \dots, 6\}$

W praktyce stosuje się tylko kilka kroków algorytmu Strassena-Winegrada[9]. Złożoność obliczeniową $O(n^{w_0})$ algorytmu oznacza, że jego wykonanie zatrzymuje się po osiągnięciu macierzy wymiaru 1×1 .

Rozdział 4

Równoległe algorytmy mnożenia macierzy

4.1 Algorytm w modelu PRAM

Algorytm sumowania

Niech dany będzie tablica A o długości $n = 2^k$ liczb i maszyna PRAM z n -procesorami $\{P_1, P_2, \dots, P_n\}$. Każdy z procesorów wykonuje synchronicznie poniższy algorytm

Dane wejściowe: Tablica A o długości $n = 2^k$ przechowywana w pamięci wspólnej. Każdy procesor ma zainicjalizowane zmienne lokalne n oraz identyfikator i

Dane wyjściowe: Suma S wartości tablicy A . Tablica A nie ulega zmianie

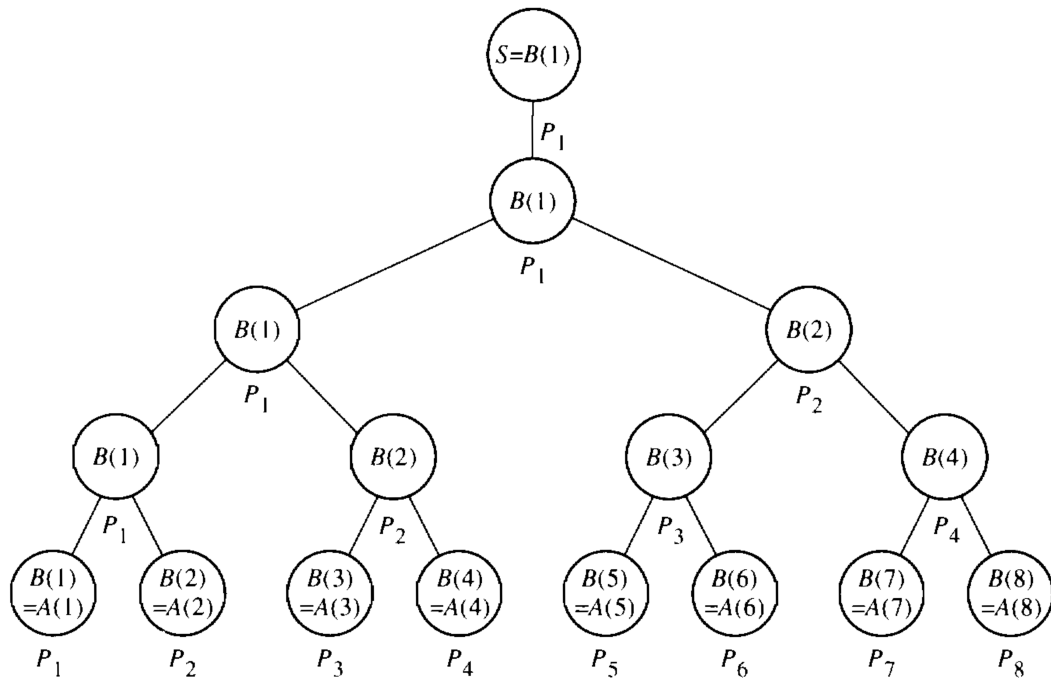
1. `global read($A(i)$, a)`
2. `global write (a , $B(i)$)`
3. `for $h = 1$ to $\log n$ do`
 `if $i \leq n/2^h$ then`
 `global read($B(2^{h-1}i)$, x)`
 `global read($B(2^hi)$, y)`
 `Set $z := x + y$`
 `global write (z , $B(i)$)`
4. `if ($i = 1$) then global write(z , S)`

Algorithm 2: Algorytm sumowania w PRAM

Przypadek dla $n = 8$ ilustruje rysunek 4.1. W pierwszym i drugim kroku kopia B tablicy A jest tworzona w pamięci wspólnej. Zadania obliczeniowe w kroku 3 są na podstawie wyważonego drzewa binarnego, którego liście odpowiadają elementom tablicy A. Procesor odpowiedzialny za wykonanie tej operacji jest określony przez indeks poniżej węzła reprezentującego tę operację. Zauważmy, że procesor P_1 , odpowiedzialny za ustawianie wartości $B(1)$ i zapisywanie sumy S , jest zawsze aktywny w trakcie wykonywania algorytmu, podczas gdy procesory P_5, P_6, P_7, P_8 są aktywne tylko podczas kroków 1 i 2.

Uwaga 4.1. Pomijamy szczegóły operacji dotyczących dostępu do pamięci. Operacje postaci **Ustaw** $A := B + C$, gdzie A, B i C są zmiennymi wspólnymi będziemy interpretować jako ciąg instrukcji

```
global read(B, x);
global read(C, y);
Ustaw z := x + y;
global write(z, A);
```



Rysunek 4.1: Algorytm sumowania ośmiu elementów w modelu PRAM z osmioma procesorami. Każdy wewnętrzny wierzchołek grafu reprezentuje operację sumowania.

Algorytm mnożenia

Rozważmy problem obliczenia iloczynu \mathbf{C} dwóch macierzy $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$, gdzie $n = 2^k$, dla pewnego $k \in \mathbf{N}$. Załóżmy, że dysponujemy n^3 procesorami $P_{i,j,l}$, $1 \leq i, j, l \leq n$ maszyny PRAM. Wówczas dla każdej pary (i, j) , n procesorów $P_{i,j,l}$, gdzie $1 \leq l \leq n$, oblicza sumę $\sum_{l=1}^n A(i, l)B(l, j)$ w myśl algorytmu 2.

wejście: Macierze $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$, gdzie $n = 2^k$, dla pewnego $k \in \mathbf{N}$
przechowywanych we wspólnej pamięci. Lokalnie
zainicjalizowane zmienne to n i trójka wskaźników (i, j, l)

wyjście: Iloczyn $\mathbf{C} = \mathbf{AB}$ w pamięci współdzielonej

1. Oblicz $C'(i, j, l) = A(i, l)B(l, j)$
2. for $h = 1$ to $\log n$ do
if $l \leq n/2^h$ then **Ustaw** $C'(i, j, l) := C'(i, j, 2l - 1) + C'(i, j, 2l)$
3. if $(l = 1)$ then **Ustaw** $C(i, j) := C'(i, j, 1)$

Algorithm 3: Algorytm mnożenia macierzy w modelu PRAM

Uwaga 4.2. Algorytm 4.1 wymaga równoległego odczytu ponieważ w trakcie wykonania kroku (1) procesory $P_{i,l,k}$ mogą równocześnie odczytywać te same dane. Przykładowo procesory $P_{i,1,l}, P_{i,2,l}, \dots, P_{i,n,l}$ w trakcie wykonywania kroku (1) wszystkie wymagają dostępu do $A(i, l)$.

4.2 Algorytmy w modelu sieciowym

4.2.1 Algorytm w dwuwymiarowym torusie

Dane wejściowe :

1. macierze $a[1..n, 1..n]$ i $b[1..n, 1..n]$, których elementy $a[i, j]$ umieszczone są w pamięciach lokalnych procesorów $P_{i,j}$
2. zmienna lokalna w każdym procesorze służąca do przechowywania rozmiaru n
3. zmienne lokalne i oraz j przechowujące numer procesora

Dane pomocnicze: zmienna lokalna k w każdym procesorze $P_{i,j}$

Dane wyjściowe :

1. iloczyn macierzy $c = ab$
2. elementy macierzy $c[i, j]$ umieszczone w lokalnych zmiennych c procesorów $P_{i,j}$

begin

```
  for  $k := 2$  to  $n$  do
    parfor  $P_{i,j}, 1 \leq i, j \leq n$  do
      if  $k \leq i$  then
         $a \leftarrow \text{prawy}(a)$  ;           // Przemieszczenie wierszy  $a$ 
        ;                               // macierzy cyklicznie w lewo
      end
      if  $k \leq j$  then
         $b := \text{dół}(b)$  // Przemieszczenie kolumn macierzy  $b$  ;
        // cyklicznie w górę
      end
    end
  end
```

end

// Obliczanie iloczynu macierzy a i b ;

```
  for  $k := 1$  to  $n$  do
    parfor  $P_{i,j}, 1 \leq i, j \leq n$  do
      if  $k = 1$  then
         $c := a * b$ ;
      else
         $a \leftarrow \text{prawy}(a)$ ;  $b \leftarrow \text{dół}(b)$ ;  $c := c + a * b$ ;
      end
    end
  end
```

end

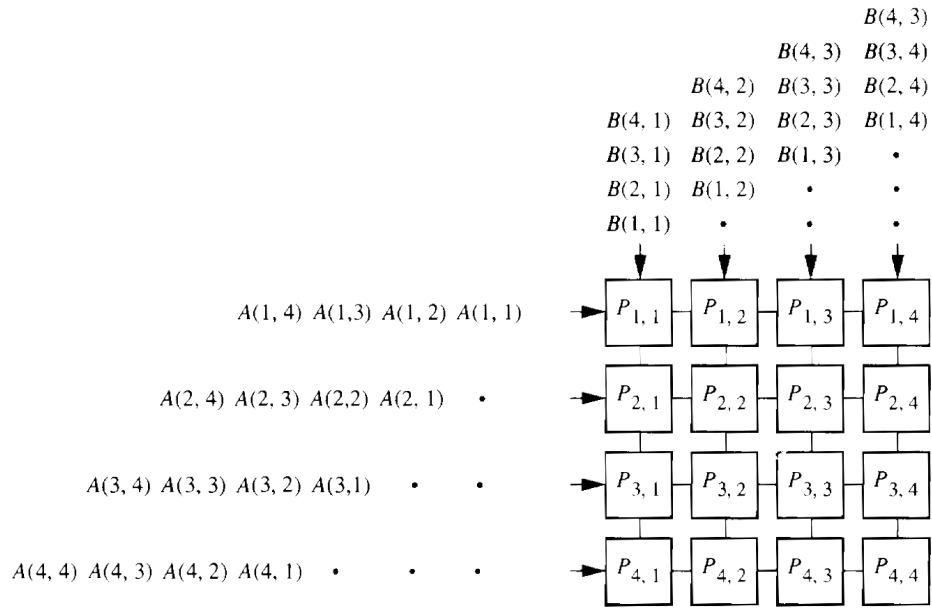
4.2.2 Algorytm systoliczny w sieci dwuwymiarowej

Rysunek 4.2 przedstawia możliwy schemat obliczenia iloczynu $AB = C$ w myśl paradygmatu systolicznego. Wiersze macierzy A są wprowadzane synchronicznie skośnie od lewej strony sieci; kolumny macierzy B są wprowadzane synchronicznie skośnie od góry sieci.

Gdy procesor $P_{i,j}$ odbierze dwie dane wejściowe $A(i, l)$ i $B(l, j)$, to przeprowadza operację $C(i, j) := C(i, j) + A(i, l)B(l, j)$; po tym przesyła $A(i, l)$ do swojego prawego sąsiada, a $B(l, j)$ do sąsiada poniżej.

Po $O(n)$ krokach, każdy procesor $P_{i,j}$ będzie miał szukaną wartość $C(i, j)$.

Algorytmy systoliczne pracują całkowicie synchronicznie; w każdej jednostce czasu, procesor otrzymuje dane od pewnego sąsiada, przeprowadza na nich lokalne obliczenia i następnie wysyła dane do któregoś swojego sąsiada.



Rysunek 4.2: Mnożenie macierzy w modelu sieciowym za pomocą algorytmu systolicznego. Wiersze macierzy A są synchronicznie umieszczane w sieci od lewej strony, podczas, gdy równocześnie „od góry” umieszczane są synchronicznie kolumny macierzy B . Gdy elementy $A(i, l)$ i $B(l, j)$ są dostępne na procesorze $P_{i,j}$, wykonywane jest działanie $C(i, j) = C(i, j) + A(i, l)B(l, j)$, $A(i, l)$ zostaje wysłany do procesora $P_{i,j+1}$ (o ile taki istnieje) oraz $B(l, j)$ zostaje wysłany do procesora $P_{i+1,j}$ (o ile taki istnieje).

4.2.3 Algorytm w topologii hipersześcianu

Rozgłaszanie „jednego do wszystkich” w hipersześcianie

Rozważmy problem rozgłaszania elementu X przechowywanego w rejestrze $D(0)$ procesora P_0 do wszystkich procesorów P_i p -procesorowego hipersześcianu, gdzie $p = 2^d$.

Algorytm polega na przechodzeniu z podkrostek najniższego wymiaru, kolejno do najwyższego w d iteracjach. Pierwsza iteracja polega na wysłaniu kopii X przez procesor P_0 do procesora P_1 . W drugiej iteracji procesor P_0 i P_1 wysyłają kopię X do P_2 i P_3 odpowiednio. Analogiczną operację przeprowadza się d -razy.

Dane wejściowe : Procesor P_0 z $p = 2^d$ -procesorowego synchronicznego hipersześcianu przechowuje element danych X w rejestrze $D(0)$.

Dane wyjściowe : X jest rogłaszany do wszystkich procesorów tak, że $D(i) = X$, gdzie $1 \leq i \leq p - 1$

; // Algorytm dla procesora P_i

begin

 for $l = 0$ to $d - 1$ do

 if $0 \leq i \leq 2^l - 1$ then

 Ustaw $D(i^l) := D(i)$;

 end

 end

end

Algorithm 4: Algorytm rozgłaszania w sieci hipersześciennej

Algorytm ?? ma złożoność równoległą $\mathcal{O}(\log p)$

Algorytm mnożenia

Rozważmy problem mnożenia macierzy $AB = C$ w synchronicznym hipersześcianie z $p = n^3$ procesorów, gdzie wszystkie macierze są wymiaru $n \times n$.

Niech $n = 2^q$ i stąd $p = 2^{3q}$. Przypiszmy procesorom indeksy (l, i, j) takie, że $P_{l,i,j}$ oznacza procesor P_r , gdzie $r = ln^2 + in + j$. Innymi słowy, rozkładając indeks r binarnie otrzymuje, że q najbardziej znaczących bitów odpowiada indeksowi l , następne q najbardziej znaczących bitów odpowiada indeksowi i i ostatecznie q najmniej znaczących bitów odpowiada wskaźnikowi j . W szczególności, jeśli ustalimy dowolną parę wskaźników spośród l, i oraz j oraz bę-

dziemy przechodzili z pozostałym wskaźnikiem po wszystkich jego możliwych wartościach, otrzymamy podkostkę wymiaru q .

Wejściowy ciąg A jest zapamiętany w podkostce wyznaczonej przez procesory $P_{l,i,0}$, gdzie $0 \leq l, i \leq n-1$, tak, że $A(i, l)$ jest zapamiętane w procesorze $P_{l,i,0}$.

Podobnie ciąg B jest zapamiętany w podkostce procesorów $P_{l,0,j}$, gdzie procesor $P_{l,0,j}$ zapamiętuje $B(l, j)$.

Celem jest obliczeniem $C(i, j) = \sum_{l=0}^{n-1} A(i, l)B(l, j)$ dla $0 \leq i, j \leq n-1$. Algorytm składa się z trzech etapów:

1. Dane wejściowe są rozdyskrebowane tak, że procesor $P_{l,i,j}$ pamięta $A(i, l)$ i $B(l, j)$ dla $0 \leq l, i, j \leq n-1$.
2. Procesor $P_{l,i,j}$ oblicza iloczyn $C'(l, i, j) = A(i, l)B(l, j)$ dla wszystkich $0 \leq i, j, l \leq n-1$.
3. Dla wszystkich $0 \leq i, j \leq n-1$ procesorów $P_{l,i,j}$, gdzie $0 \leq l \leq n-1$, obliczają sumę $C(i, j) = \sum_{l=0}^{n-1} C'(l, i, j)$

Implementacja pierwszego etapu składa się z dwóch części. W pierwszej z nich *rozgłaszamy*, dla każdego i, l , $A(i, l)$, z procesora $P_{i,l,0}$ do $P_{l,i,j}$ dla $0 \leq j \leq n-1$. Ponieważ zbiór procesorów $\{P_{l,i,j} | 0 \leq j \leq n-1\}$ wyznacza q -wymiarową kostkę dla każdej z par i oraz l , możemy użyć **algorytmu rozgłaszania**, żeby rozgłosić $A(i, l)$ od procesora $P_{i,l,0}$ do wszystkich procesorów $P_{l,i,j}$. W drugiej części każdy element $B(l, j)$ przechowywany w procesorze $P_{l,0,j}$ jest rozgłaszany do procesorów $P_{l,i,j}$ dla wszystkich $0 \leq i \leq n-1$. Pod koniec procesor $P_{l,i,j}$ będzie przechowywał dwie wartości: $A(i, l)$ i $B(l, j)$. Używając algorytmu (tutaj referencja) etap pierwszy ma równoległą złożoność czasową $O(\log n)$.

Drugi etap polega na wykonywaniu pojedynczych mnożeń na każdym z procesorów $P_{l,i,j}$. Stąd etap na etap ten składa się tylko jeden równoległy krok obliczeniowy. Pod koniec procesor $P_{l,i,j}$ przechowuje $C'(l, i, j)$.

Trzeci etap polega na obliczeniu n^2 sum $C(i, j)$. Wartości $C'(l, i, j)$ każdej z sum znajdują się w q -wymiarowym hipersześcianie $\{P_{l,i,j} | 0 \leq l \leq n-1\}$. Obliczanie takich sum ma równoległą złożoność czasową $O(\log n)$. Procesor $P_{0,i,j}$ będzie przechowywał wartość $C(i, j)$ iloczynu.

Wobec powyższych iloczyn macierzy wymiaru $n \times n$ w sieci hipersześciennej może być obliczony w czasie $O(\log n)$ dysponując n^3 procesorami.

4.2.4 Algorytm Cannona w dwuwymiarowej sieci

[11]

4.2.5 Algorytm 2.5D

[12]

4.2.6 Algorytm 3D

[11]

4.2.7 Równoległy algorytm Strassena CAPS

[9]

Bibliografia

- [1] Ed J. Radatz. *Standards Coordinating Committee 10, Terms and Definitions. The IEEE Standard Dictionary of Electrical and Electronics Terms*. IEEE, 1996.
- [2] Fayez Gebali. *Algorithms and Parallel Computing*. Wiley Publishing, 1st edition, 2011.
- [3] Zbigniew Czech. *Wprowadzenie do obliczeń równoległych*. Wydawnictwo Naukowe PWN, Warszawa, 2013.
- [4] Cormen Thomas H., Leiserson Charles E., Rivest Roland L., and Stein Clifford. *Wprowadzenie do algorytmów*. WNT, 2001.
- [5] Volker Strassen. Gaussian elimination is not optimal. *Numerische Mathematik*, 13:354–356, 1969.
- [6] Shmuel Winograd. On multiplication of 2×2 matrices. *Linear Algebra and its Applications*, 4:381–388, 1971.
- [7] Jacques Loeckx, editor. *Automata, Languages and Programming, 2nd Colloquium, University of Saarbrücken, July 29 - August 2, 1974, Proceedings*, volume 14 of *Lecture Notes in Computer Science*. Springer, 1974.
- [8] Craig C. Douglas, Michael Heroux, Gordon Slishman, Roger M. Smith, and Roger M. Gemm. A portable level 3 blas winograd variant of strassen’s matrix-matrix multiply algorithm, 1994.
- [9] Grey Ballard, James Demmel, Olga Holtz, Benjamin Lipshitz, and Oded Schwartz. Communication-optimal parallel algorithm for strassen’s matrix multiplication. *CoRR*, abs/1202.3173, 2012.
- [10] Lynn Elliot Cannon. *A Cellular Computer to Implement the Kalman Filter Algorithm*. PhD thesis, Bozeman, MT, USA, 1969. AAI7010025.

- [11] H. Gupta and P. Sadayappan. Communication efficient matrix-multiplication on hypercubes. Technical Report 1994-25, Stanford Infolab, 1994.
- [12] Edgar Solomonik and James Demmel. Communication-optimal parallel 2.5d matrix multiplication and lu factorization algorithms. Technical Report UCB/EECS-2011-72, EECS Department, University of California, Berkeley, Jun 2011.