# MIDDLE EAST TECHNICAL UNIVERSITY

## Project Report

### CENG478 - Introduction to Parallel Programming
### Parallellizing Matrix Multiplication

June 26, 2021

**Student Name/Number:** Ersel Hengirmen / 2468015

# Contents

# 1 Introduction

In this problem we will focus on N×N matrix multiplication with 2 N×N matrices. Let's say we have A and B matrices (A and B being N×N matrices) and want to find A.B=C matrix. For example:

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} * \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \\ b_{31} & b_{32} & b_{33} \end{bmatrix} = \begin{bmatrix} c_{11} & c_{12} & c_{13} \\ c_{21} & c_{22} & c_{23} \\ c_{31} & c_{32} & c_{33} \end{bmatrix}$$

To find any $c_{xy}$ in a matrix multiplication like the one above, we would need to calculate the below equation:

$$c_{xy} = \sum_{z=1}^{N} a_{xz} * b_{zk}$$

Which would require $O(N^3)$ amount of time complexity.

There are also algorithms like Strassen's algorithm which has a complexity near $O(N^{2.8074})$ and as of December 2020 there has been a new algorithm of complexity $O(N^{2.3728596})$ which is said to be best complexity ıntil now. This algorithm was created by Josh Alman and Virginia Vassilevska Williams. However this algorithm is a galactic algorithm.[1]

In this project we will not focus on the other algorithms but try to design a Paralel model for $O(N^3)$ model to reduce its complexity because of its convenience to adapt into parallel models. We will use Cannon's algorithm to parallelize it[2].

Note: A galactic algorithm is one that outperforms any other algorithm for problems that are sufficiently large, but where "sufficiently large" is so big that the algorithm is never used in practice.[3]

# 2 Proposed Model

## 2.1 Assumptions

- We are asssuming that we have an hpc network with cut through routing.

- We are asssuming that only one of the processes has the Initial A and B matrices stored in it.

- Because of the first assumption we will assume that message transfer cost is $t_s + mt_m$

- We are assuming that the network is in 2D-mesh topology.

- We are assuming that multiplication and addition operations take only 1 unit of time.

## 2.2 Restrictions

- To use cannon's algorithm easily we will restrict process count to be a square number(like 1 4 9 16)

- Also to make the matrix blocks homogenous we will only accept N when it is square root of process count's multiple.(for example when N is 12 the process count can be 4,9,16,36,144 because 12 is divisible to their square roots (2,3,4,6,12))

## 2.3 Parallelizing Matrix Multiplication

For parallel matrix multiplication of 2 $N^2$ matrices A and B to create a $N^2$ matrix C there are 3 basic methods that comes to mind.

1. Every process will calculate rows of the C matrix

2. Every process will calculate columns of the C matrix

3. Every process will calculate a block of the C matrix

Even though in the end there will be exactly the same number of calculations the first 2 methods has an inefficient side opposed to the 3rd method. The problem with the first 2 methods comes from the way the matrix multiplication works. If we are using matrix row method then every process will only need that corresponding row from the A matrix but since the elements in a rows are in every column then every process will need all of the B matrix to process their rows.

so every process(assume there are P number of processes) would need to get N/P rows from matrix A and all of the B matrix which means to send the matrices we ould need to make 1 Broadcast operation and 1 scatter operation from process 0. Broadcast would take $O(log(p)(t_s + N^2 t_m))$ time (since it is the whole matrix the message size is $N^2$). On the other hand scatter operation would be done on rows and would take $O(log(p)t_s + (p-1)\frac{N^2}{p}t_m)$ time ($\frac{N^2}{p}$ is the message size every process receives) so in the end just sending the matrices would take $O(N^2 log(p))$ operations(because we need to broadcast the whole matrix)

On the other hand for the block approach we will only need to send rows from A matrix and columns from B matrix normally and since we chose cannon's algorihm we will only send **blocks** of $\frac{N}{\sqrt{p}} \times \frac{N}{\sqrt{p}}$ instead. Which would need 2 scatter operations of complexity $O(log(p)t_s + (p-1)\frac{N^2}{p}t_m))$. so the complexity of initialization would be $O(N^2)$ for this case.

## 2.4   Cannon's Algorithm

The problem with Block based matrix multiplication is the fact that when we try to calculate an entry of the result matrix's value(for example $c_{ik}$) we tend to start from $a_{i0}$ and $b_{0k}$. Which means that every process would need to start with left-most A group and upper-most B group. (For example in the matrix below the Group 1 2 and 3 of C matrixes would need to start with Group1 of A.)

$$\begin{bmatrix} Group1 & Group2 & Group3 \\ Group4 & Group5 & Group6 \\ Group7 & Group8 & Group9 \end{bmatrix}$$

Cannon's algorithm uses a way to bypass this problem. By their method instead of sending every group of block of matrix to their correspongind process we would instead rotate every block of A by its block row and every block of B by its block column according to their block row number and block column number. In this way at any given time, each process would be using a different block in every iteration(while keeping the computation accurate).

For example assuming the above matrix is A it would become

$$\begin{bmatrix} Group1 & Group2 & Group3 \\ Group5 & Group6 & Group4 \\ Group9 & Group7 & Group8 \end{bmatrix}$$

and if it was the B matrix it would become

$$\begin{bmatrix} Group1 & Group5 & Group9 \\ Group4 & Group8 & Group3 \\ Group7 & Group2 & Group6 \end{bmatrix}$$

**Note:** We are not actually swapping the blocks then sending them. Instead we are sending corresponding A and B blocks to neccessary processes.

After Initial block sending from process 0 is done every program will need to calculate their given blocks. But to calculate C block of for example group 9 the program would need to make $\sqrt{p}$ block matrix multiplications (for Group9 it would calculate AGroup8*BGroup6+AGroup7*BGroup3+AGroup9*BGroup9). To do this every process will send its received blocks(A and B blocks) to the next process and receive 2 blocks from the previous one for $\sqrt{p}-1$ times($\sqrt{p}$ blocks in every row and column but since the first one was received from process 0 we only need to make $\sqrt{p}-1$ iterations.)

Sending these messages would take $t_s + \frac{N^2}{P}t_m$ for every iteration so in total sending and receiving messages to next and from previous partners would take $2(\sqrt{p}-1)(t_s + \frac{N^2}{P}t_m)$

And for every $\frac{N^2}{P}$ entry we would need to make N multiplications and N-1 additions so in total $\frac{2N^3-N^2}{P}$ operations would be done paralelly for every matrix and the time complexity of computations would become $O(\frac{N^3}{P})$

The last part is receiving these matrices. To do this the program would only need to make a gather operations to receive blocks from every process to one of the specified processes. And that would take $O(log(p)t_s + (p-1)\frac{N^2}{p}t_m))$ time

So in the end Cannon's Algorithm(Parallel) would make 3 scatter/gather operations(2 scatters for distributing A and B matrixes and 1 gather to

gather calculated C matrix) in $3log(p)t_s + 3(p-1)\frac{N^2}{p}t_m$ time. Parallely sends blocks to each partner processes which takes $2(\sqrt{p}-1)(t_s+\frac{N^2}{P}t_m)$ time and calculates the C block of the matrix in $\frac{2N^3-N^2}{P}$ time which comes to $\frac{2N^3-N^2}{P} + t_s(3log(p) + 2\sqrt{p} - 1) + t_m\frac{N^2}{P}(3p + 2\sqrt{p} - 4)$ or $O(max(\frac{N^3}{P}, N^2))$

Ant the sequential algorithm's time only depends on computations which comes up to $N^3$-$N^2$ or $O(N^3)$

**Note:** Because of the scatter and gather operations time complexity will never exceed $O(N^2)$. And increasing the number of processes is not always efficient.The problem is that just sending the matrix itself would take $O(N^2)$ operations so thinking that this is the best case.

**Important Note:** Instead of scatter and gather operations I sended and received blocks of A,B and C because of to send Blocks I needed to send them with derived data types[4] because Collective Communication Routines can only be used with primitive data types. Because of that the actual complexity of sending A and B and gathering C became $O(p * t_s + N^2t_m))$ instead of $O(log(p)t_s + (p-1)\frac{N^2}{p}t_m))$. So in the end the total complexity of my implementation is $\frac{2N^3-N^2}{P} + (3p + 2\sqrt{p} - 2)(t_s + \frac{N^2}{p}t_m)$ which stil comes up to $O(max(\frac{N^3}{P}, N^2))$

# 3    Discussion

I have used mpi libraries inc to check efficiency in my local system with maximum N being 1250 with p being 4

While the program is faster then the sequential by 2 times(with 4 processes) it has not reached the expected efficiency of 4 times efficiency.

my last 10 calculations with matrices of size 1250 and process count 4 are:

| calculation | Parallel time | Sequential time |
|:---:|:---:|:---:|
| 1 | 5.188443 | 10.520247 |
| 2 | 5.205665 | 10.476868 |
| 3 | 5.205504 | 11.147428 |
| 4 | 5.082834 | 11.127606 |
| 5 | 5.141081 | 10.719176 |
| 6 | 5.029180 | 10.508834 |
| 7 | 5.132676 | 10.439199 |
| 8 | 5.113407 | 10.424923 |
| 9 | 5.040144 | 11.151711 |
| 10 | 5.079504 | 10.848710 |

But according to my calculations I was thinking that there was a problem with my system. $\frac{parallel}{sequential}$ (assuming $t_s$ and $t_m$ are little numbers (1)) for N=1250 and p=4. I was thinking of the theoritical model of using scatter/-gather in my calculations which led me to believe that I should have been getting 4 times faster for the parallel algorithm(since both of their gamma time complexities were same).

$$\frac{\frac{2N^3-N^2}{p} + 3log_2(p) + 2\sqrt{p} - 1 + \frac{N^2}{p}(3p + 2\sqrt{p} - 4)}{2N^3 - N^2} =$$

$$\frac{81738282}{325390625} = 0.251200482496998$$

Even when I changed my calculations to the model I have implemented there wasn't a big difference in the calculations for N 1250 and p being 4 it should have been 4 times more efficient yet it remained being only 2 times efficient.

$$\frac{\frac{2N^3-N^2}{P} + (3p + 2\sqrt{p} - 2)(t_s + \frac{N^2}{p}t_m)}{2N^3 - N^2} =$$

$$\frac{981640639}{1951562500} = 0.25140056380952$$

# 4   Conclusion

Since the best known sequential algorithm for matrix multiplication has a time complexity of O($N^{2.3728596}$) with enough processors Cannon's algorithm can be a great method for square matrices.

# 5 How to use the code

There are 2 c files in the src directory **serial.c** and **solution.c**. **solution.c** has the parallel algorithm while **serial.c** has the sequential one.

To use the codes use **script.sh** inside the src directory.It works in this order:

- It will compile **solution.c**(with 4 processors u can change it by going inside script.sh) and **serial.c**

- Then it will start parallel algorithm

  Parallel algorithm will first create A and B matrices

  Then copy them inside **input** file(for sequential algorithm's use)

  At this point it starts Parallel timer

  After calculating the solutions it will first write the time it took to calculate C matrix in seconds.

  After that it will record the C matrix inside the **output** file

- Then it will run sequential algorithm with the inputs from parallel algorithm

  Sequential algorithm will start the timer after getting inputs and stop the timer and print it after calculation is done

  after printing the timer sequential algorithm will save the C matrix inside **output2** file.

- then the script will call diff command with output and output2

**Important Note:** Before using them make sure serial. **solution.c** and **serial.c** both has the same defined value for N or they will produce different results because of their different matrix sizes.

# References

[1]https://en.wikipedia.org/wiki/Matrix_multiplication_algorithm

[2]https://en.wikipedia.org/wiki/Cannon%27s_algorithm

[3]https://en.wikipedia.org/wiki/Galactic_algorithm

[4]https://hpc-tutorials.llnl.gov/mpi/collective_communication_routines/