

Final report

Distributed matrix multiplication

March 2019

The work has been done by MSc. 1st year students Artem Vasilev and Mahmud Allahverdiyev

1 Introduction

Matrix multiplication is one the central operations in a number of numerical algorithms. Up to now, several different approaches have been proposed to improve the efficiency of matrix multiplication algorithms. If we apply the original mathematical definition of matrix multiplication, it gives an algorithm that takes time on the order of n^3 to multiply two $n \times n$ matrices ($\Theta(n^3)$ in big \mathcal{O} notation). Better asymptotic bounds on the time required to multiply matrices have been known since the work of Strassen in the 1960s, but it is still unknown what the optimal time is. The definition of matrix (1) multiplication is that if $C = AB$ for an $n \times m$ matrix A and an $m \times p$ matrix B , then C is an $n \times p$ matrix with entries

$$c_{ij} = \sum_{k=1}^m a_{ik}b_{kj} \quad (1)$$

2 General and distributed matrix multiplication algorithms

Strassen's algorithm which runs in $\mathcal{O}(N^{2.81})$ time was an improvement to the naive algorithm. In the following decades a number of theoretical and practical improvements have been made, including CopperSmith-Winograd algorithm which has a complexity of $\mathcal{O}(n^{2.3754})$. This algorithm is based on approaching the problem via Arithmetic progressions. Despite of the smaller exponent in the complexity, this algorithm has a large constant; therefore, is not used in practice.

In general, algorithms with better asymptotic running time than the Strassen algorithm are rarely used in practice, due to the same factor in their running times which make them impractical. Theoretically, it is possible to improve the exponent of the algorithm further; however, the exponent must be at least 2.

The reason for the latter is obvious; the matrices contains $\mathcal{O}(n^2)$ entries and all of them have to be read at least once to calculate the exact result.

Strassen's algorithm improves the naive multiplication by grouping the operations in a clever manner. For convenience, let's assume that matrices A, B have sizes $2^n \times 2^n, n \geq 0$; if it's not the case, we can pad the matrices with zeros to get them have the sizes equal to 2-s powers. Let's partition matrices A, B and C into block form

$$A = \begin{bmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{bmatrix}$$

,

$$B = \begin{bmatrix} B_{1,1} & B_{1,2} \\ B_{2,1} & B_{2,2} \end{bmatrix}$$

and

$$C = \begin{bmatrix} C_{1,1} & C_{1,2} \\ C_{2,1} & C_{2,2} \end{bmatrix}$$

where the smaller matrices have dimensionality of $2^{n-1} \times 2^{n-1}$. The naive algorithm does not exploit the grouping possibilities for submatrices to reduce the number of possible matrix multiplications. Strassen introduces the result via new matrices:

$$\begin{aligned} c_{11} &= f_1 + f_4 - f_5 + f_7, \\ c_{12} &= f_3 + f_5, \\ c_{21} &= f_2 + f_4, \\ c_{22} &= f_1 - f_2 + f_3 + f_6, \end{aligned} \tag{2}$$

where

$$\begin{aligned} f_1 &= (a_{11} + a_{22})(b_{11} + b_{22}), \\ f_2 &= (a_{21} + a_{22})b_{11}, \\ f_3 &= a_{11}(b_{12} - b_{22}), \\ f_4 &= a_{22}(b_{21} - b_{11}), \\ f_5 &= (a_{11} + a_{12})b_{22}, \\ f_6 &= (a_{21} - a_{11})(b_{11} + b_{12}), \\ f_7 &= (a_{12} - a_{22})(b_{21} + b_{22}). \end{aligned} \tag{3}$$

Using only 7 multiplications and 18 matrix additions (2), (3) it has the lower complexity than the naive algorithm. In order to get a reasonable speedup via Strassen's algorithm, large sized matrices are needed since the constant factor in the naive algorithm is comparably smaller.

Block algorithms have become building blocks for distributed algorithm routines, including the world-ruling BLAS-1, BLAS-2 and BLAS-3 modules are crucial factors in fast developing parallel programming world.

The three loops in iterative matrix multiplication can be arbitrarily swapped with each other without an effect on correctness or asymptotic running time. However, the order can have a considerable impact on practical performance due to the memory access patterns and cache use of the algorithm; which order is best also depends on whether the matrices are stored in row-major order, column-major order, or a mix of both.

In a perfect setting, we can distribute the data of matrix in blockwise fashion between p servers (1).

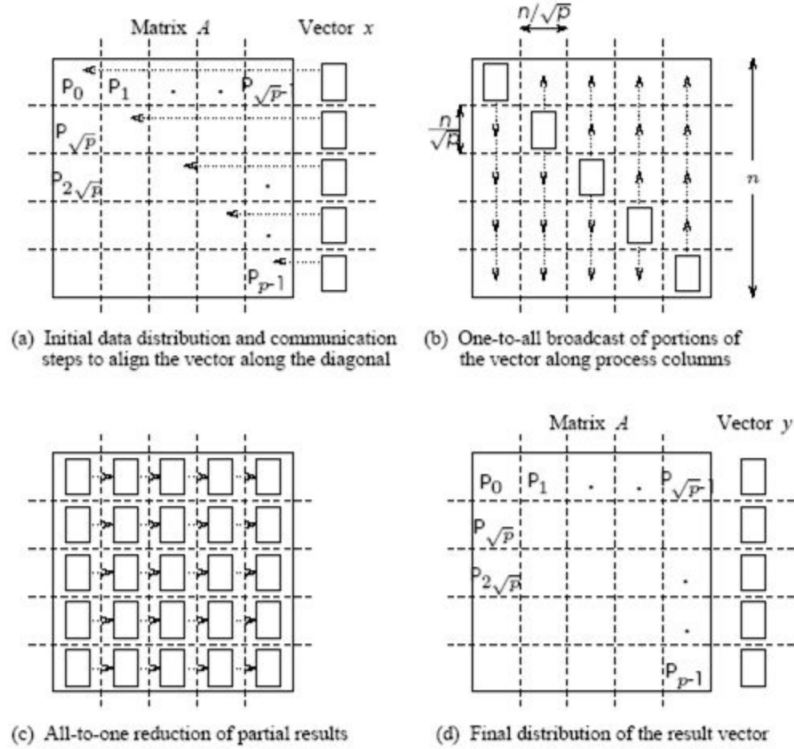


Figure 1: Distributing operations between servers

3 GASP codes for Distributed Matrix Multiplication

Now, consider a setting where two users have matrices $A \in \mathcal{F}_q^{r \times s}$ and $B \in \mathcal{F}_q^{r \times s}$, respectively. The common goal is to calculate AB in such a distributed way that the operations are done by N independent servers while assuming that no information is about A or B is revealed to any of the servers. Among all

servers, it is guaranteed that at most T of them may collude. Moreover, the communication costs, thus total amount of work should be minimized.

4 Conclusion

Aha

References