

Updates on Practical Strassen's Algorithms

Rodrigo Brandão
Devangi N. Parikh

UT Austin

BLIS Retreat 2023

Outline

Introduction

Fast Matrix Multiplication

Practical FMM

Performance Results

Conclusions

Introduction

Introduction

Consider, $C+ = AB$, where C , A and B are $m \times n$, $m \times k$ and $k \times n$ matrices respectively and m, n, k are even.

Introduction

Consider, $C+ = AB$, where C , A and B are $m \times n$, $m \times k$ and $k \times n$ matrices respectively and m, n, k are even.

Consider the following partition of C , A and B :

$$\begin{pmatrix} C_0 & C_1 \\ C_2 & C_3 \end{pmatrix} + = \begin{pmatrix} A_0 & A_1 \\ A_2 & A_3 \end{pmatrix} \begin{pmatrix} B_0 & B_1 \\ B_2 & B_3 \end{pmatrix}$$

Direct Computation

$$C_0 := A_0 B_0 + A_1 B_2 + C_0$$

$$C_1 := A_0 B_1 + A_1 B_3 + C_1$$

$$C_2 := A_2 B_0 + A_3 B_2 + C_2$$

$$C_3 := A_2 B_1 + A_3 B_3 + C_3$$

Direct Computation

$$C_0 := A_0B_0 + A_1B_2 + C_0$$

$$C_1 := A_0B_1 + A_1B_3 + C_1$$

$$C_2 := A_2B_0 + A_3B_2 + C_2$$

$$C_3 := A_2B_1 + A_3B_3 + C_3$$

8 multiplications, 8 additions

Fast Matrix Multiplication

Strassen's Algorithm¹

$$M_0 = (A_0 + A_3)(B_0 + B_3); \quad C_0 += M_0; C_3 += M_0;$$

$$M_1 = (A_2 + A_3)B_0; \quad C_2 += M_1; C_3 -= M_1;$$

$$M_2 = A_0(B_1 - B_3); \quad C_1 += M_2; C_3 += M_2;$$

$$M_3 = A_3(B_2 - B_0); \quad C_0 += M_3; C_2 += M_3;$$

$$M_4 = (A_0 + A_1)B_3; \quad C_1 += M_4; C_0 -= M_4;$$

$$M_5 = (A_2 - A_0)(B_0 + B_1); \quad C_3 += M_5;$$

$$M_6 = (A_1 - A_3)(B_2 + B_3); \quad C_0 += M_6;$$

¹V. Strassen, "Gaussian elimination is not optimal," in Numerische Mathematik, 1969

Strassen's Algorithm¹

$$M_0 = (A_0 + A_3)(B_0 + B_3); \quad C_0 += M_0; C_3 += M_0;$$

$$M_1 = (A_2 + A_3)B_0; \quad C_2 += M_1; C_3 -= M_1;$$

$$M_2 = A_0(B_1 - B_3); \quad C_1 += M_2; C_3 += M_2;$$

$$M_3 = A_3(B_2 - B_0); \quad C_0 += M_3; C_2 += M_3;$$

$$M_4 = (A_0 + A_1)B_3; \quad C_1 += M_4; C_0 -= M_4;$$

$$M_5 = (A_2 - A_0)(B_0 + B_1); \quad C_3 += M_5;$$

$$M_6 = (A_1 - A_3)(B_2 + B_3); \quad C_0 += M_6;$$

7 multiplications, 22 additions

¹V. Strassen, "Gaussian elimination is not optimal," in Numerische Mathematik, 1969

Fast Matrix Multiplication (FMM)

A $\langle \tilde{m}, \tilde{k}, \tilde{n} \rangle$ FMM algorithm² by partitioning

$$C = \left(\begin{array}{c|c|c} C_0 & \cdots & C_{\tilde{n}-1} \\ \hline \vdots & & \vdots \\ \hline C_{(\tilde{m}-1)\tilde{n}} & \cdots & C_{\tilde{m}\tilde{n}-1} \end{array} \right), A = \left(\begin{array}{c|c|c} A_0 & \cdots & A_{\tilde{k}-1} \\ \hline \vdots & & \vdots \\ \hline A_{(\tilde{m}-1)\tilde{k}} & \cdots & A_{\tilde{m}\tilde{k}-1} \end{array} \right),$$
$$\text{and } B = \left(\begin{array}{c|c|c} B_0 & \cdots & B_{\tilde{n}-1} \\ \hline \vdots & & \vdots \\ \hline B_{(\tilde{k}-1)\tilde{n}} & \cdots & B_{\tilde{k}\tilde{n}-1} \end{array} \right)$$

where A_i , B_j , and C_p are the submatrices of A , B and C , with a single index in the row major order.

²A. R. Benson and G. Ballard, “A framework for practical parallel fast matrix multiplication,” PPOPP 2015.

Fast Matrix Multiplication (FMM) (contd.)

Then, $C := C + AB$ is computed by, for $r = 0, \dots, R - 1$,

$$M_r := \left(\sum_{i=0}^{\tilde{m}\tilde{k}-1} u_{ir} A_i \right) \times \left(\sum_{j=0}^{\tilde{k}\tilde{n}-1} v_{jr} B_j \right);$$
$$C_p += w_{pr} M_r \quad (p = 0, \dots, \tilde{m}\tilde{n} - 1)$$

where (\times) is a matrix multiplication that can be done recursively, u_{ir} , v_{jr} , and w_{pr} are entries of a $(\tilde{m}\tilde{k}) \times R$ matrix U , a $(\tilde{k}\tilde{n}) \times R$ matrix V , and a $(\tilde{m}\tilde{n}) \times R$ matrix W , respectively.

One-level Strassen $\langle 2, 2, 2 \rangle$

The set of coefficients that determine the $\langle 2, 2, 2 \rangle$ algorithm is denoted as $\mathcal{T} = \llbracket U, V, W \rrbracket$, where

$$\begin{aligned} U &= \begin{matrix} & M_0 & M_1 & M_2 & M_3 & M_4 & M_5 & M_6 \\ \begin{matrix} A_0 \\ A_1 \\ A_2 \\ A_3 \end{matrix} & \left(\begin{array}{cccccc} 1 & 0 & 1 & 0 & 1 & -1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & -1 \end{array} \right) \end{matrix} \\ V &= \begin{matrix} \begin{matrix} B_0 \\ B_1 \\ B_2 \\ B_3 \end{matrix} & \left(\begin{array}{cccccc} 1 & 1 & 0 & -1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 1 \\ 1 & 0 & -1 & 0 & 1 & 0 & 1 \end{array} \right) \end{matrix} \\ W &= \begin{matrix} \begin{matrix} C_0 \\ C_1 \\ C_2 \\ C_3 \end{matrix} & \left(\begin{array}{cccccc} 1 & 0 & 0 & 1 & -1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 & 0 \\ 1 & -1 & 1 & 0 & 0 & 1 & 0 \end{array} \right) \end{matrix} \end{aligned}$$

Finding other FMM

- ▶ For $\langle \tilde{m}, \tilde{k}, \tilde{n} \rangle$ partitioning of the matrices, the associated \mathcal{T} has dimension $(\tilde{m}\tilde{k}, \tilde{k}\tilde{n}, \tilde{m}\tilde{n})$ represents the underlying matrix multiplication.
- ▶ Seek a rank- R decomposition of tensor \mathcal{T} , where $R < \tilde{m}\tilde{k}\tilde{n}$

Reinforcement Learning + FMM³

- ▶ Agent “AlphaTensor” is trained to play a single-player game where the objective is finding tensor decompositions within a finite factor space.
- ▶ For matrices in \mathbb{R} reduces number of multiples:
 1. $\langle 3, 4, 5 \rangle$ from 48 to 47.
 2. $\langle 4, 4, 5 \rangle$ from 64 to 63.
 3. $\langle 4, 5, 5 \rangle$ from 80 to 76.

³A. Fawzi et. al, “Discovering faster matrix multiplication algorithms with reinforcement learning,” Nature 2023.

Practical FMM

Practical Strassen's Algorithm

Conventional Implementations

Matrix Size	must be large
Matrix Shapes	must be square
No Additional Workspace	×
Parallelism	usually task parallelism

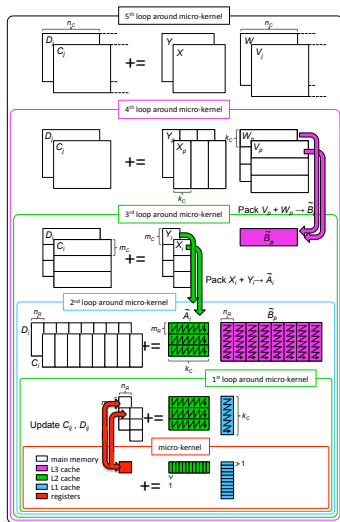
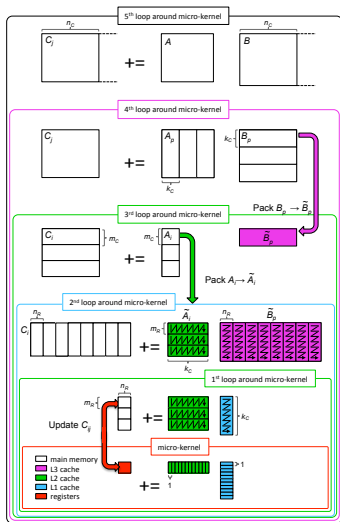
⁴J. Huang et. al, "Strassen's algorithm reloaded," in SC 16. IEEE, 2016.

Practical Strassen's Algorithm

	Conventional Implementations	Strassen Reloaded Impl ⁴
Matrix Size	must be large	speed up at smaller sizes
Matrix Shapes	must be square	speed up for rank-k
No Additional Workspace	×	✓
Parallelism	usually task parallelism	data parallelism

⁴J. Huang et. al, "Strassen's algorithm reloaded," in SC 16. IEEE, 2016.

Strassen Reloaded⁴



⁴J. Huang et. al, "Strassen's algorithm reloaded," in SC 16. IEEE, 2016.

Generating High-Performance Implementations of FMM⁵

- ▶ Generates code that takes as input $\langle \tilde{m}, \tilde{k}, \tilde{n} \rangle$ and $\llbracket U, V, W \rrbracket$ and as output generates implementations that build upon the primitives that combine taking linear combinations of matrices with the packing routines and/or micro-kernels that underlie BLIS.
- ▶ Provides a model of cost for each implementation that can then be used to choose the best FMM algorithm for a matrix of given size and shape.

⁵J. Huang et. al, “Generating Families of Practical Fast Matrix Multiplication Algorithms,” in IPDPS 2017.

Performance Results

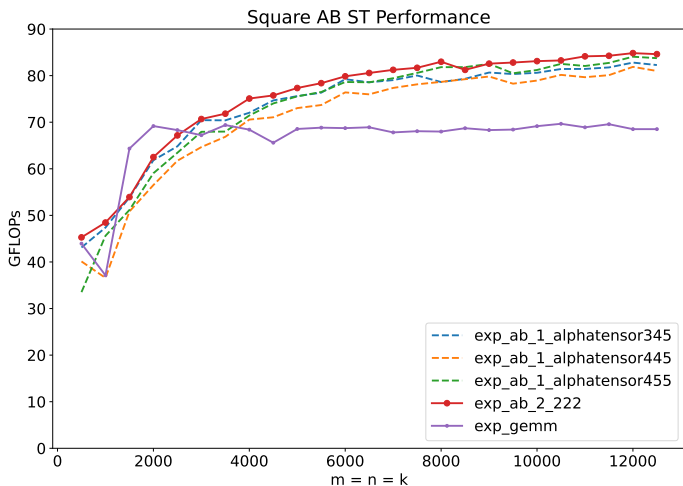
Setup

- ▶ Updated the code generator to support Intel's Haswell Architecture
- ▶ Intel(R) Xeon(R) CPU E3-1270 v6 @ 3.80/4.20GHz processor
- ▶ Three-level cache: L1 128 KB, L2 1 MB, L3 8 MB

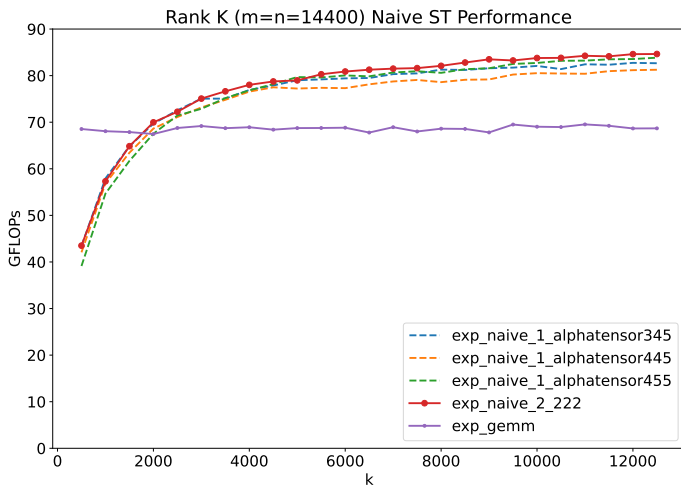
Results - Single-Threaded Naive ($m=k=n$)



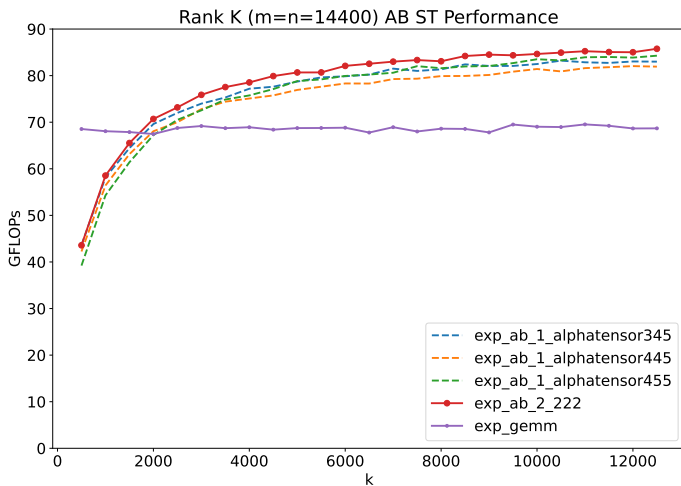
Results - Single-Threaded AB ($m=k=n$)



Results - Single-Threaded Naive ($m=n=14400$)



Results - Single-Threaded AB ($m=n=14400$)



Conclusions

Conclusions

- ▶ For a single thread, Naive and AB Strassen continues to perform generally better than the other FMM algorithms.

Future Work

- ▶ Create a BLIS plugin that implements Strassen and other FMM using the BLIS framework.
- ▶ Perform experiments with multi-threading as well as on other architectures.
- ▶ Investigate the necessary conditions for an FMM algorithm to outperform Strassen.

Thank you!

