

# GPU Multiplication Algorithms

J. C. Thomas

The University of Iowa, Dept. of Biostatistics

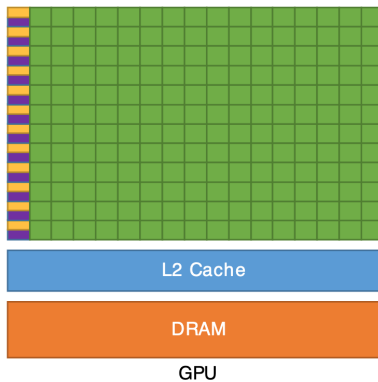
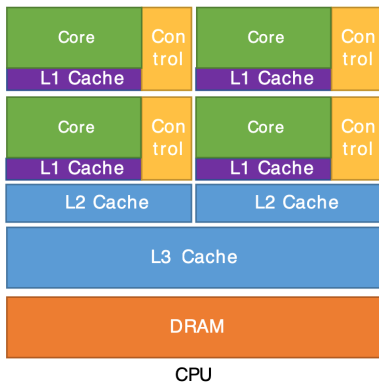
Dec 14 2025

GPUs

# Why GPUs?

- Why are GPUs so expensive?
  - It turns out everyone loves to do matrix multiplication
- Why is matrix multiplication so popular?
  - Deep learning and AI are all pretty much just a bunch of matrix multiplications. Analogous: think about the closed form solution of linear regression
- How do GPUs help with matrix multiplication?
  - Matrix multiplication is a bunch of independent arithmetic
  - This can be parallelized via a GPU

# How GPUs?



# Attacking a Village

- If we are a fantasy warlord and we need to destroy all the houses in a village, should we send in our 4 trolls or our army of 100 goblins?
  - Our trolls are really strong, but there are only a few of them. They will need to destroy a house, move to the next house, destroy it, ect.
  - Our goblins are not as strong but there are a lot of them and they are quick. Each goblin can concentrate on destroying an individual house.

**IOWA**

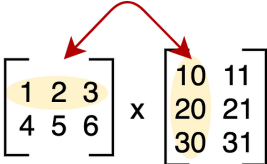
# GPU Computation

- Units are threads, blocks, and grids
- Little green boxes/goblins are threads
  - This is the primary computational unit
- Blocks are groups of threads, can share memory
- Grids are the entire workspace
- For matrix multiplication, this maps nicely. Each thread is an element of the output matrix, the entire output matrix is the grid. Blocks are less intuitive.

# Matrix Multiplication

# By Hand Approach

- Why is this a good candidate for parallel computing?
- How does this work within the thread, block, grid structure?


$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} \times \begin{bmatrix} 10 & 11 \\ 20 & 21 \\ 30 & 31 \end{bmatrix}$$
$$= \begin{bmatrix} 1 \times 10 + 2 \times 20 + 3 \times 30 & 1 \times 11 + 2 \times 21 + 3 \times 31 \\ 4 \times 10 + 5 \times 20 + 6 \times 30 & 4 \times 11 + 5 \times 21 + 6 \times 31 \end{bmatrix}$$
$$= \begin{bmatrix} 10 + 40 + 90 & 11 + 42 + 93 \\ 40 + 100 + 180 & 44 + 105 + 186 \end{bmatrix} = \begin{bmatrix} 140 & 146 \\ 320 & 335 \end{bmatrix}$$



# By Hand Approach Considerations

- Perfect use of parallel computing!
- Not very memory efficient. For multiplying two  $4 \times 4$  matrices, each value is loaded in 4 times
- Tile method improves on this, each value loaded just 2 times

# Tiled Approach

				B			
				0	2	3	1
				4	1	1	0
				2	1	3	0
				1	2	4	4
A	1	0	2	3			
	0	1	2	1			
	0	2	4	3			
	1	3	1	0			

For Block 1,  $i=0$

```
sA[tx, ty] = A[x, ty]
sB[tx, ty] = B[tx, y]
```

1	0
0	1

sA

0	2
4	1

sB

C = Grid

# Tiled Approach

For Block 1,  $i=1$

$$sA[tx, ty] = A[x, ty + TPB]$$
$$sB[tx, ty] = B[tx + TPB, y]$$

**A**

1	0	2	3				
0	1	2	1				
0	2	4	3				
1	3	1	0				

**B**

0	2	3	1
4	1	1	0
2	1	3	0
1	2	4	4

**sA**

2	3
2	1

**sB**

2	1
1	2

**C = Grid**


# Tiled Approach

**A**

1	0	2	3
0	1	2	1
0	2	4	3
1	3	1	0

**B**

0	2	3	1
4	1	1	0
2	1	3	0
1	2	4	4

**C = Grid**


For Block 1

```
for j in range(TPB):  
    tmp += sA[tx, j] * sB[j, ty]
```

$i=0$

1	0
0	1

sA

0	2
4	1

sB

$\times$

0	2
4	1

$+$

$i=1$

2	3
2	1

sA

2	1
1	2

sB

$\times$

7	8
5	4

$+$

7	10
9	5

$C[x, y] = tmp$

# Other Approaches

- Matrix multiplication is essential and this problem has been studied a ton
- There are many, MANY complex algorithms that are much faster than these. However, these give a flavor of how complex ones work with the computational system.

Computation

# Interfacing

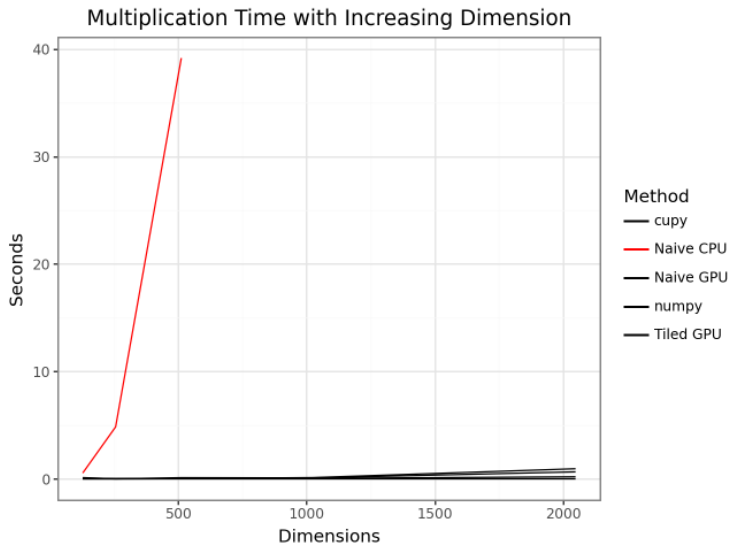
- GPUs (specifically Nvidia GPUs) can be communicated with via a C-like language called CUDA
- Instructions are passed to threads in things called kernels, function-like syntax
- Most softwares have tools and extensions for interfacing with GPUs (R, python, Julia, ect)
- Good python packages: cupy and numba

# Simulation 1

- Goal: compare speed with increasing matrix dimension for the following algorithms
  - naive method on CPU
  - pre-built in method on CPU (numpy)
  - naive method on GPU
  - tiled method on GPU
  - pre-built in method on GPU (cupy)
- Square matrices of dimension 128, 256, 512, 1024, and 2048
- Threads: 16



# Simulation 1 Results



# Simulation 1 Results

