# Memory Hierarchies and Matrix-Matrix Multiplication

Pratyuksh Bansal

Università della Svizzera italiana

September 22, 2021

# What is Matrix multiplication?

If $A$ is $m \times n$ matrix and $B$ is $n \times p$ matrix, then $C = AB$ is $m \times p$ matrix.



$$c_{ij} = \text{row i} \times \text{col j} = \sum_{k=1}^{n} a_{ik} b_{kj}$$

$$\underset{m \times n}{A} \quad \times \quad \underset{n \times p}{B} \quad = \quad \underset{m \times p}{C}$$

Also known as **GEMM** - General matrix multiplication.

# Naive algorithm for GEMM

**Algorithm 1** Matrix multiplication

1: **for** $i = 1$ to $m$ **do**
2:   **for** $j = 1$ to $p$ **do**
3:     **for** $k = 1$ to $n$ **do**
4:       $C(i,j) = \sum_{k=1}^{n} A(i,k) * B(k,j)$
5:     **end for**
6:   **end for**
7: **end for**

Question: Is this the most optimal way to do it?

# Performance



NxN matrix-matrix-multiplication on Quad-Core AMD Opteron(tm) Processor 2344@1.7GH

Theoretical Peak Performance (TPP) :
$(1.7 \times 10^9 cycles/sec) \times (4 flops/cycle) = 6.8$ GFlops/sec.

More than $50\times$ difference in performance!!
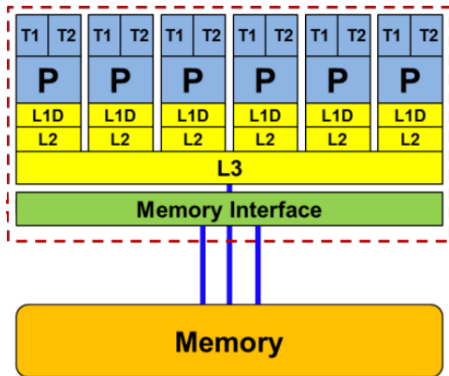
# Memory hierarchy



Figure: Memory hierarchy of a multi-core architecture

# Why is it important?

- Basic linear algebra operation, appears in several applications in physics, engineering, etc.
- Benchmark to compare the performance of processors.

# Why is it important?

- Basic linear algebra operation, appears in several applications in physics, engineering, etc.
- Benchmark to compare the performance of processors.

Learning objective :

Fundamental concepts and ideas used to optimize GEMM.

# Thank you