

## Flops : Memory Cache Hierarchy

To assess the performance of an algorithm we have been counting arithmetic operations.

Around the 1980s, CPUs became fast enough that arithmetic costs were not the primary bottleneck. Instead, the costs of storing and manipulating numbers in memory became increasingly important.

### Matrix Multiplication

$$C = A B \quad \Rightarrow \quad \sim 2n^3 \text{ flops}$$

for  $i = 1:m$

for  $j = 1:m$

$$C_{ij} = \sum_{k=1}^m A_{ik} B_{kj}$$

"Naive Algorithm"

$$\begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}_{i\text{th entry}} = \begin{bmatrix} \text{\textit{i}th row} \\ \hline \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix} \begin{bmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \\ \cdot \end{bmatrix}_{j\text{th column}}$$

If flops are the dominant cost on a 2.6 GHz processor, we might expect

$$P_{nn} = \left[ \frac{2n^3 \text{ flops}}{t_{ns}} \right] \approx \left[ \frac{2 \text{ flops}}{1 \text{ cycle}} \right] \left[ \frac{2.6 \text{ cycles}}{1 \text{ ns}} \right]$$

$\uparrow$  time for  
mat-mat product  
in nanoseconds

= 5.2 flops/ns  
(Gigaflops)

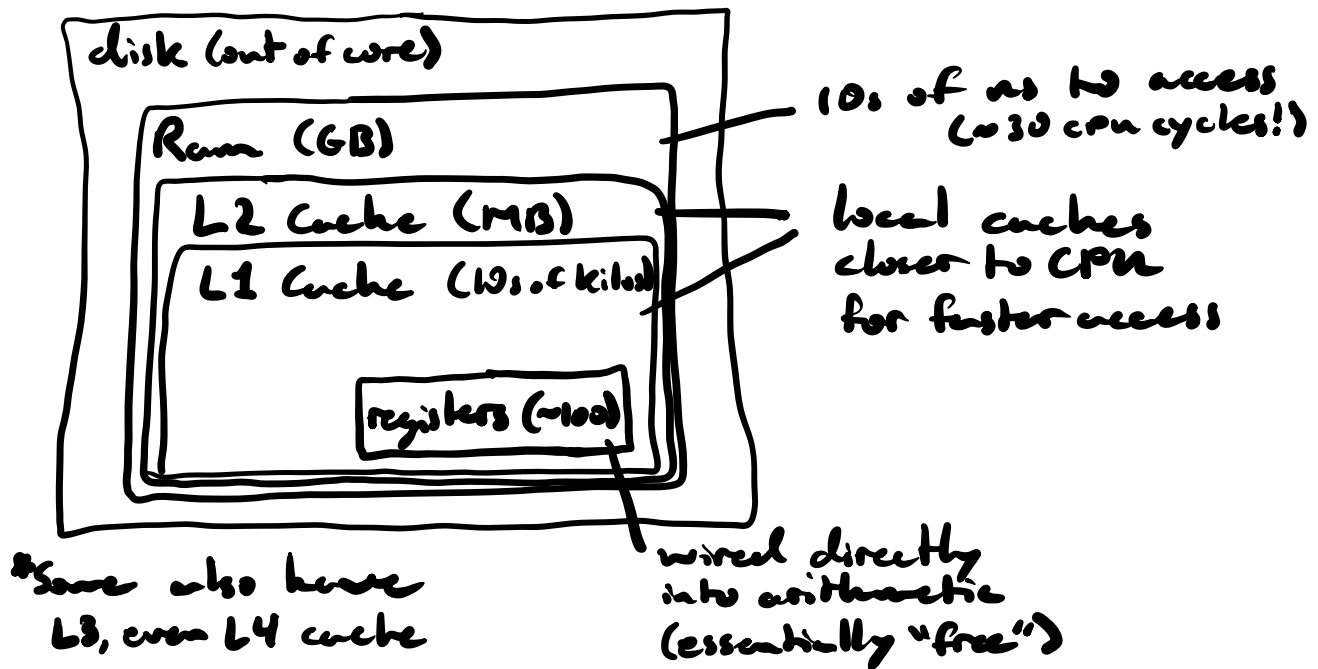
$\Rightarrow$  See Matmult.pdf for actual experiment.

$\Rightarrow$  Observe that  $P_{nn}$  is  $< 1$  gigaflop and decreases as  $n$  increases!

$\Rightarrow$  Highly optimized BLAS mat-mat product does achieve near peak theoretical flop rate using same # operations and mathematically equivalent algorithm.

How do we understand the difference in performance if flops are the same?

# Memory Hierarchy



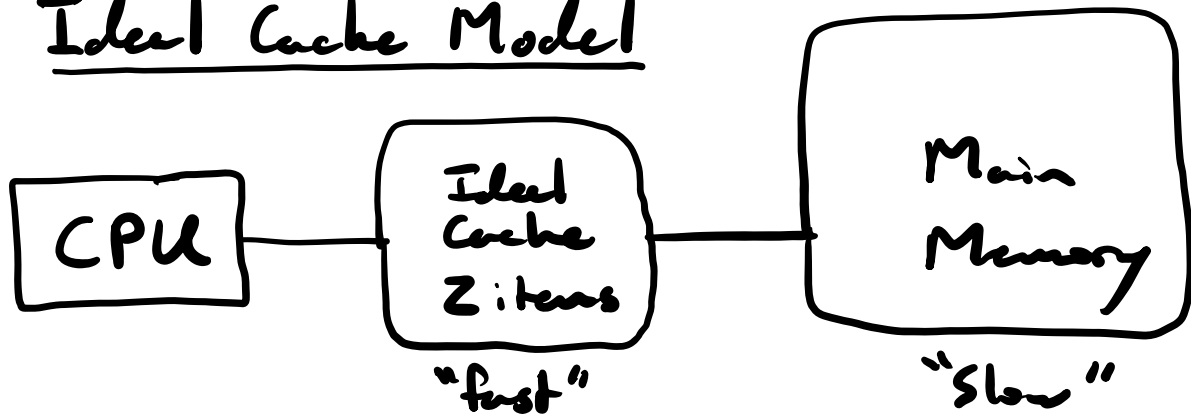
Not all flops are created equal!

⇒ Depends where : when data is accessed

Idea: do as much work as possible with data in "fast" memory before loading new data into local cache.

In practice, local caches are managed at hardware level, not by typical programmers. They replace data which has not been used recently with data needed currently.

## Ideal Cache Model



Cache "hit": CPU needs item from cache

Cache "miss": CPU needs item not in cache, item then loaded into cache for future use.

Idea: loaded item replaces item that will not be used for *longest time in future*.

⇒ Analyze algorithms by counting Cache misses  $\mathcal{O}$ .

Idealized model but easier to analyze and reveals basic ideas needed to understand memory-efficient algorithms. Gets within a constant of optimal!

## Cache performance of naive Mat. Mult.

Notation:  $f = \Theta(g)$  if constants  $\alpha, \beta > 0$  s.t.

$$\alpha g(x) \leq f(x) \leq \beta g(x) \text{ as } x \rightarrow +\infty.$$

" $f$  Asymptotically proportional to  $g$ "

$$\begin{array}{c} C \\ \left[ \begin{array}{c} \cdot \\ \vdots \\ \vdots \end{array} \right] \\ i\text{th entry} \end{array} = \begin{array}{c} A \\ \left[ \begin{array}{c} \text{\textit{i}th row} \\ \hline \vdots \end{array} \right] \\ \text{outer loop over rows of A} \end{array} \left[ \begin{array}{c} B \\ \vdots \\ \text{\textit{i}th column} \\ \vdots \end{array} \right] \\ \text{inner loop over columns of B} \end{array}$$

Cache size  $Z$

Cache misses  $Q(m, Z) = m(m^2 - Z)$

↑ loop through all entries of  $B$  for every row of  $A$ .

$$Q(m, Z) = \Theta(m^3)$$

$\Rightarrow$  makes no use of Cache asymptotically

## Cache-Aware Matrix Mult.

Idea: Load in blocks of  $A, B$  and do all ops with these blocks before loading new data into cache.

$$\begin{array}{c} C \\ \left[ \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \right] = \begin{array}{c} A \\ \left[ \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \right] \end{array} \begin{array}{c} B \\ \left[ \begin{array}{|c|c|c|} \hline & & \\ \hline \end{array} \right] \end{array} \left( \frac{n}{b} \right)^2 \text{ blocks}$$

$$C_{21} = A_{21}B_{11} + A_{22}B_{21} + A_{23}B_{31}$$

$$\text{block size} = b \Rightarrow 3b^2 = Z$$

$\uparrow$   
 store  $C_{21}, A_{2i}, B_{j1}$

$$\begin{aligned} \text{Now, } Q(n, Z) &= \underbrace{\Theta\left(\underbrace{2b^2}_{\text{block size}} \underbrace{\left(\frac{n}{b}\right)^3}_{\text{\# of blocks}}\right)}_{\text{block size}} \\ &= \Theta\left(\frac{n^3}{\sqrt{Z}}\right) \quad \text{as } n \rightarrow \infty \end{aligned}$$

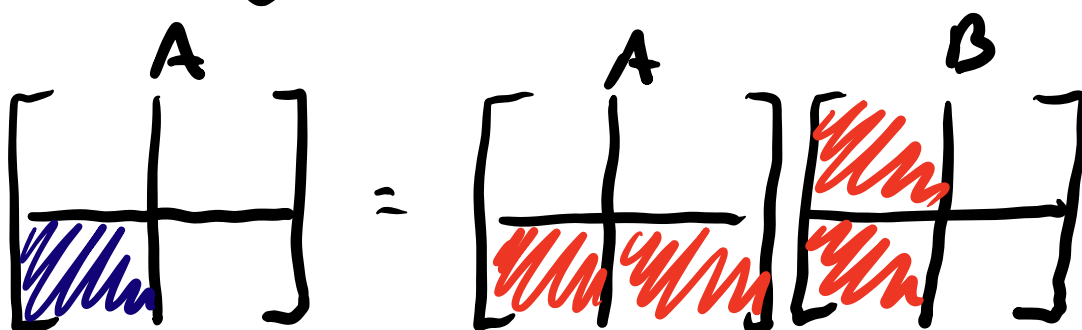
$\Rightarrow$  benefit from cache

$\Rightarrow$  Actually need to do this hierarchically.

## Cache - Oblivious

→ Divide recursively

$$Q(n) = \begin{cases} 8Q(\frac{n}{2}) & \text{otherwise} \\ 3n^2 & 3n^2 \leq 2 \end{cases}$$



$$C_{21} = A_{21} B_{11} + A_{22} B_{21}$$

Call recursively until n sub. small

for  $n^2 \gg 2$

$$Q(n) = 8Q(\frac{n}{2})$$

$$= 8 \cdot 8 \cdot Q(\frac{n}{4})$$

$$= 8^2$$

$$= \underbrace{8 \dots 8}_{k \text{ times}} \left[ 3 \left( \frac{m}{2^{2^k}} \right)^2 \right] \quad \frac{1}{4} < f \leq 1$$

$$\Rightarrow = (2^1)^k 3 \left( \frac{m}{2^{2^k}} \right)^2 = 2^k 3m^2$$

$$2^k := m \sqrt{\frac{3}{f2}}$$

$$= \Theta \left( \frac{m^3}{\sqrt{f2}} \right)$$

$\Rightarrow$  Same asymptotic cache benefit,  
but w/out explicit knowledge of cache size!