Physical locality ← transfer cost
Temporal locality ← re-use
Correlation locality

Metrics:

hit rate
composite performance metrics — access time
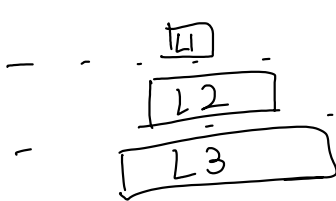
$$\text{arithmetic intensity} = \frac{\#ops}{\text{data transferred}}$$

$$\text{"Locality"} = \frac{\#\text{words accessed locally}}{\#\text{words accessed for computation}}$$

$$= \frac{\text{total} - \text{transfers}}{\#ops \times \text{operands/operation}} \qquad \left(\frac{\#bytes}{\#bytes}\right)$$

⇒ #ops ⇒ ISA dependent ⇒ keep it at 'C' level like ADD, MUL

What does "locally" mean?



⇒ Put a line in hierarchy; above that is local.

⇒ Depends on the point of view.

Correlation locality

→ Hit rate may include correlation locality
(Paper — Ask Mattan)

EXAMPLE: MAT_MUL



N [ A ] X [ B ] = [ C ]     Look at the hit rates.
       N

Counting misses w.r.t. a cache.

for i
  for j
    for k

$$C_{ij} = A_{ik} B_{kj}$$

First Time: (2K+1) locations needed to store entire row of A, entire column of B & $C_{ij}$.

$^{(k)}$ ... $^{(k)}$ ... $^{(1)}$

Now move to next column of B.

Cache Size Z

Parameters: Fully Associative
    Cache Size = Z
      Line Size = 1
      Replacement Policy = LRU
      Write-back
      Allocation : R/W  (anything we touch, put into Cache)
      Prefetcher : None.

$\Rightarrow$   $Z > 2N+1$     (keep entire row of A)

<span style="color:red">$\Rightarrow$ How can we get N+1 with different replacement policy ??</span>

Any $\overset{ideal}{\wedge}$ replacement policy is within 2X of LRU.

Total no. of Misses:   $N \cdot N \cdot (2N+1) = 2N^3 + N^2 \leftarrow$ transferred.

    total $= N \cdot N \cdot N \cdot (4)$

           $\longrightarrow$ (old C, new C, A, B)

$$\frac{4N^3 - (2N^3 + N^2)}{4N^3} \to \frac{1}{2} \qquad \text{Not very good.}$$

If we have Cache Line Size = L

no. of misses = $\dfrac{N+1}{L} + N$

# of transfers = $\left( \dfrac{N^3}{1} + \dfrac{N^3}{L} + \dfrac{N^2}{L} \right) \cdot L$

Things get worse because 'L' comes into picture., since B has no cache line locality.

Change iteration order

      For i
        (for j) $\longrightarrow$ loop interchange        For i
                                        for k

$$\text{For } k$$

$$C_{ijp} = A_{ik} \times B_{kj}$$

$$\text{for } j$$

$$C_{igt} = A_{ik} \times B_{kj}$$

→ Now we have more misses on C

⇒ Compilers might do this optimization automatically.
(loop inter-change)

— no fast_math

## if $Z > 2N+1$:

now instead of $2N^3 + N^2 = N^2(2N+1)$
we get better locality. Locality $= N^2(N+2)$

## if $Z > 3N^2$:  all matrices  (Cache fits all matrices)

$$\text{locality} = \frac{-3N^2 + 4N^3}{4N^3} \longrightarrow \text{going towards } 1$$

— — — — — — — — — —

Take a block out of full matrices:



$C_{00} = A_{00}B_{00} + A_{01}B_{10} + A_{02}B_{20}$
↓
$b \times b$ sub-matrix

Did it improve locality?

Choose b such that $Z > 3b^2$

$$b < \sqrt{\frac{Z}{3}}$$

$$b^2 + \underset{\substack{\downarrow \\ \text{Bring } C}}{} \; \underset{\substack{\downarrow \\ \text{Bring} \\ A_{00} \, B_{00}}}{2b^2} + \underset{\substack{\downarrow \\ A_{01} \, B_{10}}}{2b^2} + 2b^2$$

$$= b^2 + \sum_{0}^{N/b} 2b^2 \qquad = b^2 + \frac{2N}{b} \cdot b^2$$

$$= \left(b^2 + 2Nb\right) \cdot \left(\frac{N}{b}\right)^2 \quad = \quad N^2 + \frac{2N^3}{b} = N^2 + \frac{2\sqrt{3}N^3}{\sqrt{Z}}$$

$\Downarrow$
Total no. of subblocks

Earlier: $\quad 2N^3 + N^2$

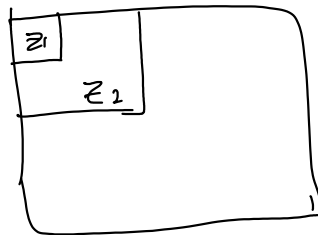Now $\quad : \quad \dfrac{2\sqrt{3}\,N^3}{\sqrt{Z}} + N^2$

$\longrightarrow$ improved by $\sqrt{Z}$ times

Now locality $= \dfrac{4N^3 - \left(\frac{3.4N^3}{\sqrt{Z}} + N^2\right)}{N^3}$

Hit Rate improves significantly by re-ordering.

$\sim$ $\rightarrow$ This Technique is called Blocking (Tiling).

One more Cache: $\qquad Z_1 \quad Z_2$
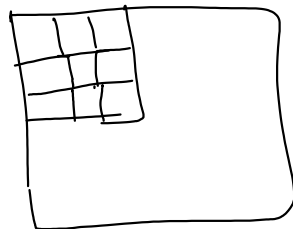


To gain Locality in L2,
we'll incur misses in L1
but less ~~memory~~ load from memory.

$\rightarrow$ optimize for L1 or L2 ??

within L2 Block, optimize for L1.

Nesting Tiling



for $ii$
  for $jj$
    $C_{ii,jj} = 0$
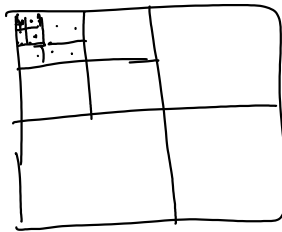    for $kk$
      for $i = b \cdot ii$
        for $j = b \cdot jj$

How to come up with an order - oblivious of $Z$.

Base $\rightarrow$ $2 \times 2$ / $4 \times 4$

Start from top:

This is divide & conq algo.

$C_{00} = A_{00} \cdot B_{00} + A_{01} \cdot B_{10}$

$C_{01} = A_{00} \cdot B_{01} + A_{01} \cdot B_{11}$

$C_{10} = A_{10} \cdot B_{00} + A_{11} \cdot B_{10}$

$C_{11} = A_{10} \cdot B_{01} + A_{11} \cdot B_{11}$

$\rightarrow$ Blend between Cache-aware & Cache-obliviousness;

$\rightarrow$ Stop recursion when sub-block $>$ thresh.

$r$ - recursion depth

$C$ - read $2^{r \sim \log(N)}$ times

no. of misses

$Q(N) < 8 \cdot Q\left(\dfrac{N}{2}\right) \Rightarrow \Omega\left(\dfrac{N^3}{Z} + N^2\right)$

$\Rightarrow$ Labs: Prefetcher ON $\rightarrow$ will change the performance

PIN 2
Performance Counters.