

Lecture 4 : Announcements

- hw time spent
- office hours
- don't forget: talk topic + SSH key
- talk time assignment

Case Study: Streaming Workloads

Q: Estimate expected throughput for saxpy on an architecture with caches. What are the right units?

$$z_i = \alpha x_i + y_i \quad (i = 1, \dots, n)$$



Demo: https://github.com/lcw/stream_ispc

Special Store Instructions

"Non-temporal stores"

At least two aspects to keep apart:

- Non-temporal locality
- Spatial locality



What hardware behavior might result from these aspects?

- Non-temporal: writing past cache, hint at fastevl
- Spatial: don't bring in data, mark rest of line undret

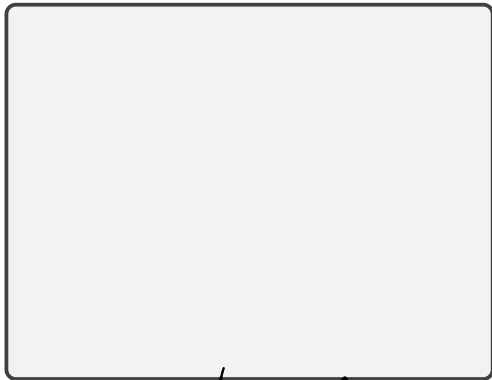
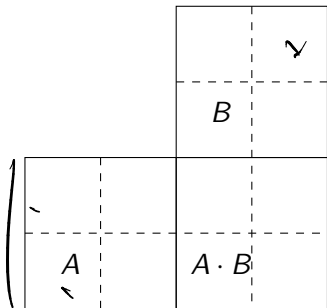
- ▶ Comment on what a compiler can promise on these aspects.
- ▶ Might these 'flags' apply to loads/prefetches?

(see also: [\[McCalpin '18\]](#))

Case study: Matrix-Matrix Mult. ('MMM'): Code Structure

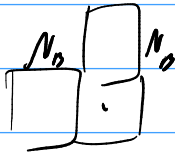
$$O(n^2) - O(n^3)$$

- ▶ How would you structure a high-performance `_MMM`?
- ▶ What are sources of concurrency? ←
- ▶ What should you consider your working set?



concurrency / absence of ord. req. / parallelism: $\downarrow \downarrow$

$$(a+b) \times c + (d+e)$$



Case study: Matrix-Matrix Mult. ('MMM') via Latency

Come up with a simple cost model for MMM in a two-level hierarchy based on latency:

N N_b

Avg. latency per access =
(1 - miss ratio) · cache latency + miss ratio · main mem. lat.

Total accesses: $4N_b^3$

Misses: $3N_b^2$

Miss ratio: $\frac{3N_b^2}{4N_b^3 \cdot \text{cache line size}} = \frac{3}{4N_b \cdot \text{CLS}}$

Case study: Matrix-Matrix Mult. ('MMM') via Bandwidth

- Come up with a cost model for MMM in a two-level hierarchy based on bandwidth:

16×2 FMAs
 $512/32 \uparrow \uparrow \uparrow$ 2 2 2 Fixed Multiplies per FMA-capable exec units

Cycle count for whole calc: $2N^3 / (2 \cdot 32)$

$4N^3 / (N^3/32) = 128$ Floats/cycle

Total mem \leftrightarrow cache data motion;

blocks: (block size) $\cdot 4 = (N/N_B)^3 \cdot N_B^2 \cdot 4$

Required mem bw:

$4N^3 / N_B / (N^3/32) = 128 / N_B$ floats/cycle

total mem. motion / #cycles

Case study: Matrix-Matrix Mult. ('MMM'): Discussion

Discussion: What are the main simplifications in each model?



[\[Yotov et al. '07\]](#)

General Q: How can we analyze cache cost of algorithms in general?