

A MODERN MULTI-CORE PROCESSOR

* Forms of llism & understanding latency + bandwidth

1. Quick review

* single instruction stream = single list of instr. to execute

* superscalar execution is max 4x.

* speedup blocked by communication, shared resources, work distribution (load balancing).

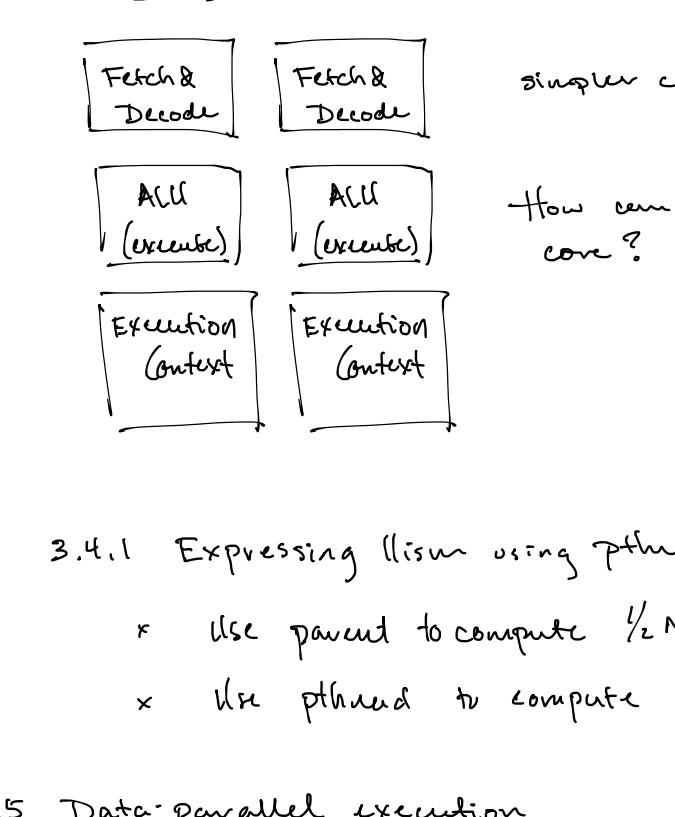
2. Overview

- * two concepts about ll execution
- * two challenges of accessing memory

3. Parallel Execution

- * compute $\sin(x)$, for vector x , via Taylor expansion.

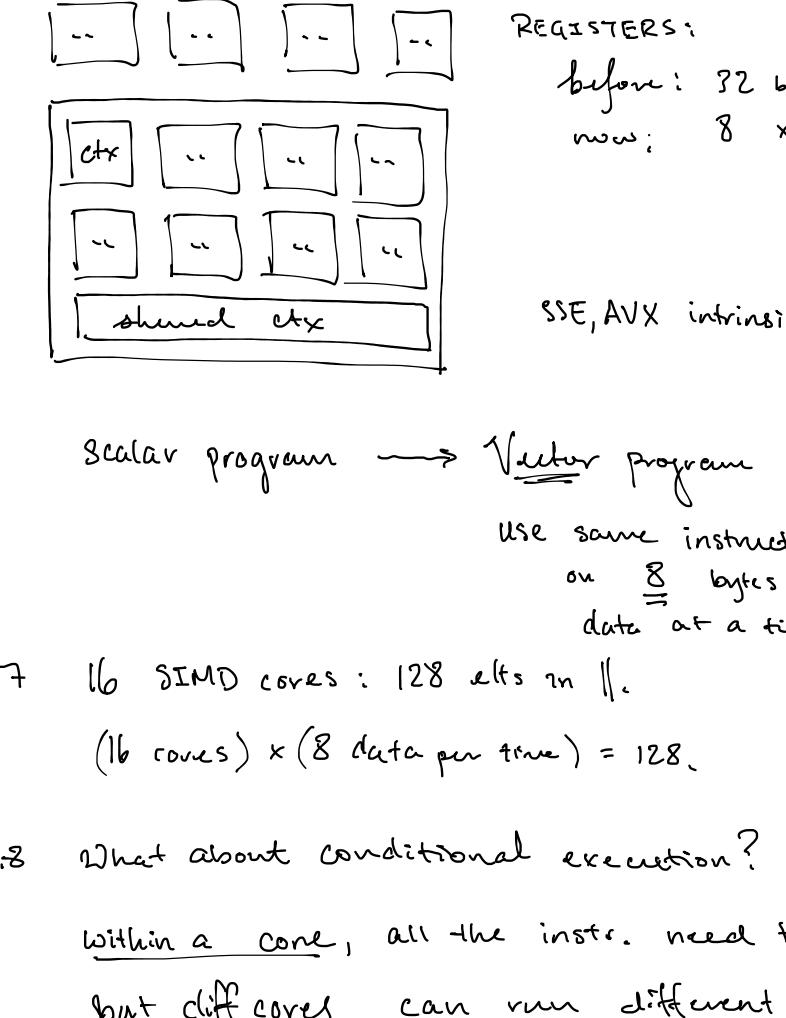
3.1 Compile



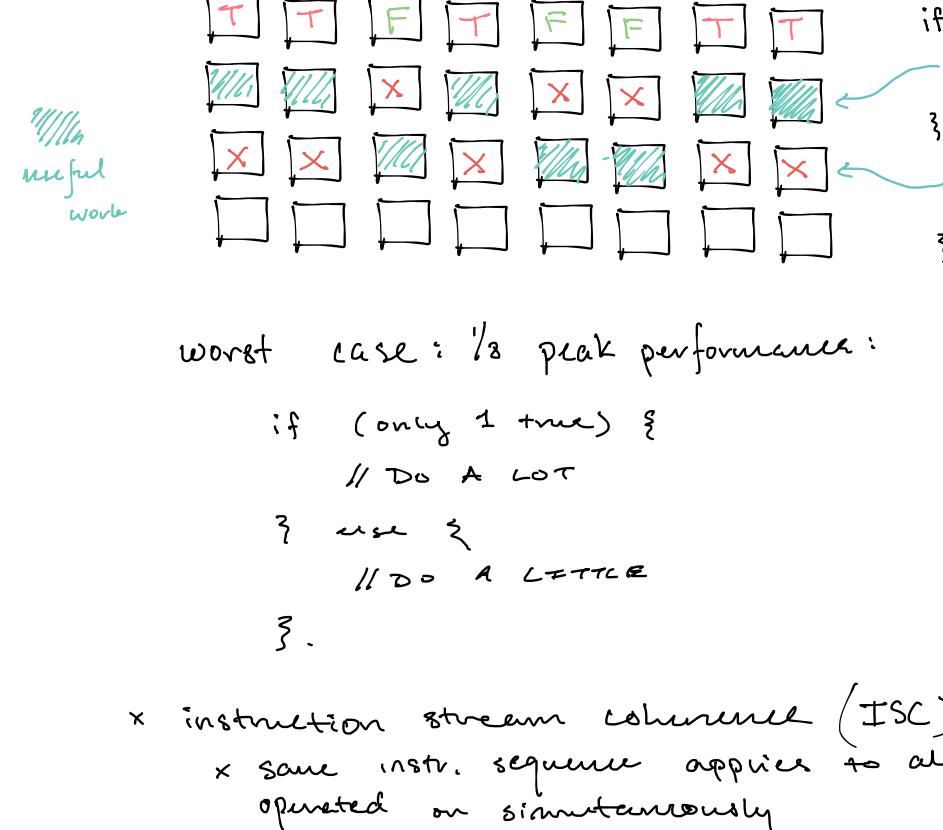
* instruction stream \leftrightarrow thread. (ILP)

\hookrightarrow not how many ALUs you have.

3.2 Processor: Pre-Multi-core Era



3.3 Two Cores



3.4.1 Expressing llism using pthreeds.

- * Use pthreeds to compute $\frac{1}{2}N$.
- * Use pthreeds to compute second $\frac{1}{2}N$.

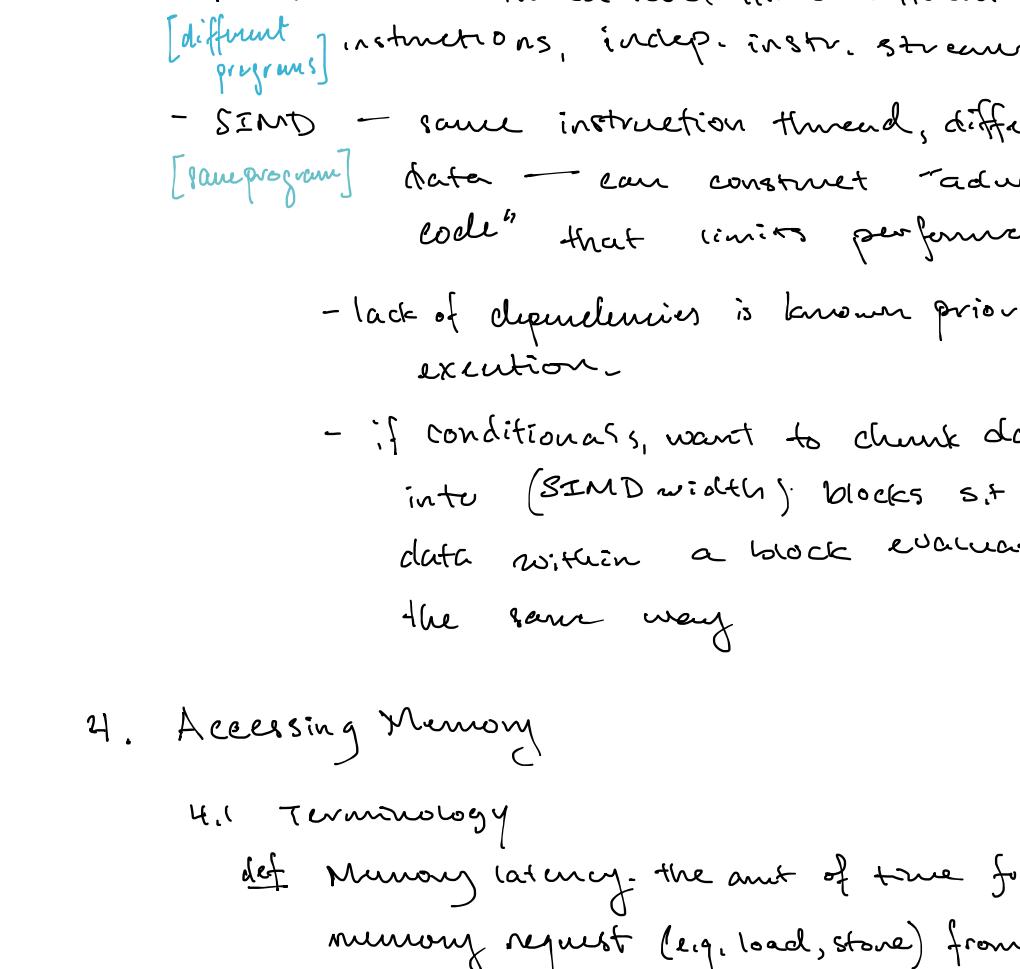
3.5 Data-parallel execution

- * "forall" \Rightarrow all iterations are data-independent
- launch a pthreed per core.

* this is better b/c it is [# core independent]

* also tells that all iterations do EXACTLY the same thing

3.6 SIMD Processing



Scalar program \rightarrow Vector program \leftarrow still a single instruction stream, but each inst. does 8 things at once.

3.7 16 SIMD cores: 128 elts in ll.

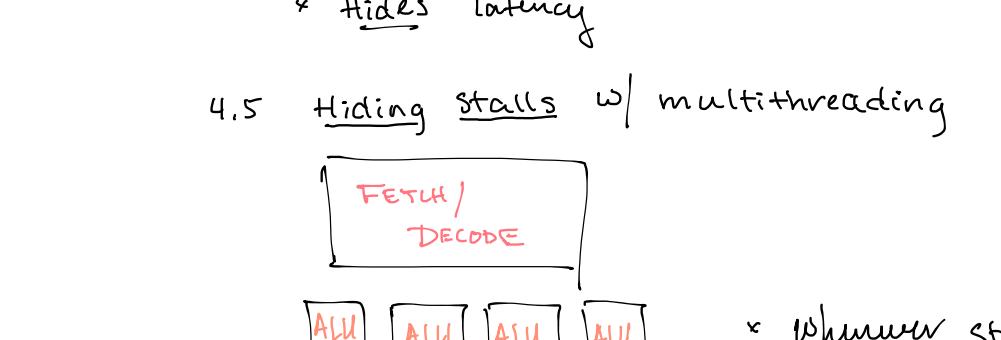
(16 cores) \times (8 data per core) = 128.

3.8 What about conditional execution?

Within a core, all the instr. need to be the same.

But diff cores can run different instr. streams.

Continue through, but don't execute certain ALUs in the SIMD core (or mask out unwanted).



worst case: $\frac{1}{3}$ peak performance:

```
if (only 1 true) {  
    // Do A LOT  
}  
? else {  
    // Do A LITTLE  
}
```

* instruction stream coherence (ISC)

* same instr. sequence applies to all elts operated on simultaneously

* necessary for efficient use of SIMD,

* NOT necessary for efficient parallelization across cores

* "divergent" execution (lack of ISC)

3.9 Modern CPUs, SIMD

* SSE 128 bit ops: 4×32 bits (4 wide float vec).

* AVX 256 bit ops: 8×32 bits (8 wide float vec)

* llism conveyed using ll language semantics.

* llism explicitly requested by programmer by intrinsics.

* def "explicit SIMD" — SIMD generated at compile time.

3.10 SIMD execution on GPUs.

* "Implicit SIMD"

* compiler generates scalar binary.

* N instances of program always run together on the processor

* GPU, SIMD width ranges from 8 to 32.

\hookrightarrow divergence would mean $\frac{1}{32}$ of peak performance!

3.11 Summary

* several forms of ll execution

- Multicore \rightarrow thread-level llism: different programs

- SIMD — same instruction thread, different data \rightarrow can construct "adversarial code" that ruins performance.

- lack of dependencies is known prior to execution

- if conditionals, want to chunk data into (SIMD width) blocks s.t. all data within a block evaluates the same way

4. Accessing Memory

4.1 Terminology

def memory latency: the amt of time for a memory request (e.g. load, store) from a processor to be serviced by the memory system

def memory bandwidth: the rate at which a memory system can provide data to a processor

\hookrightarrow a rate (bytes/sec)

$$\text{delay} = \text{latency} + \frac{\text{size of request}}{\text{bandwidth}}$$

4.2 Stalls

* a processor "stalls" when the next instr. can't be run b/c of dependence on previous instruction

* memory access time is measure of latency

* accessing memory is major source of stalls.

4.3 Caches \downarrow [reduce latency]

4.4 Prefetching

* hides latency

4.5 Hiding stalls w/ multithreading

aka a different thread

\hookrightarrow run diff set of instructions

\hookrightarrow [Aka how to use resources effectively!]

\hookrightarrow tradeoff more L1# per thread, or more threads?

\hookrightarrow If instruction stream per core, running concurrently (not simultaneously)

SIMULTANEOUS = # cores

CONCURRENT = # cores \times # CTXs per core

* Nvidia GTX 480:

15 cores, 32 wide SIMD per core, 48 execution CTXs per core.

\hookrightarrow 15 diff instruction streams

\hookrightarrow 32 per instruction

\hookrightarrow 48 diff instr. streams.

\hookrightarrow SIMULTANEOUS

\hookrightarrow ratio of $\frac{\text{memory}}{\text{memory}}$ is important!