## Announcements:
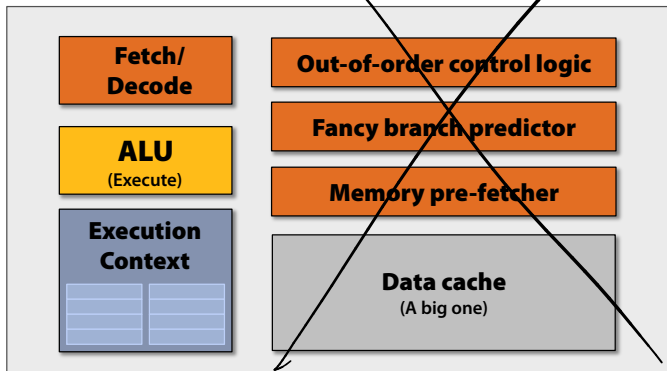
- Project submission logistics
- GPU-focused HW3: tonight / tomorrow

# "CPU-style" Cores
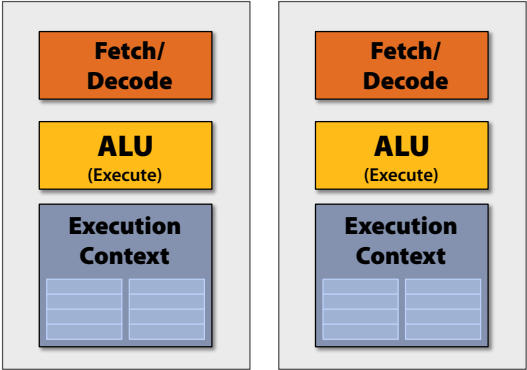


**Fetch/ Decode**

**ALU** (Execute)

**Execution Context**

**Out-of-order control logic**

**Fancy branch predictor**

**Memory pre-fetcher**

**Data cache** (A big one)

[Fatahalian '08]

# Slimming down



Fetch/
Decode

ALU
(Execute)

Execution
Context

Idea #1:

Remove components that
help a single instruction
stream run fast

[Fatahalian '08]

# More Space: Double the Number of Cores



[Fatahalian '08]

# Even more



[Fatahalian '08]

# SIMD



**Fetch/ Decode**

**ALU** (Execute)

**Execution Context**

Idea #2: SIMD

Amortize cost/complexity of managing an instruction stream across many ALUs
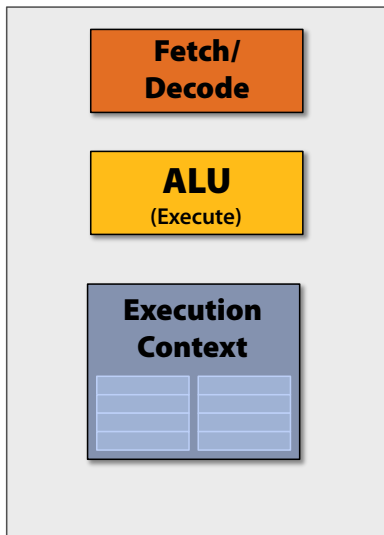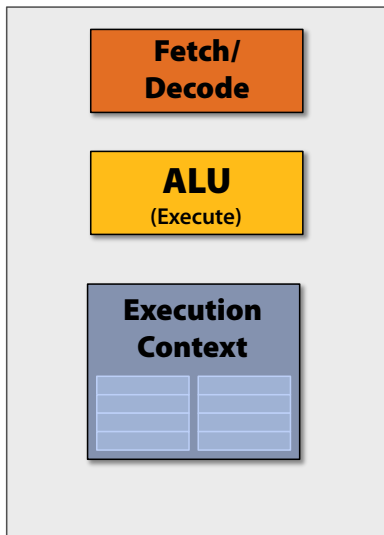
[Fatahalian '08]

# SIMD



**Idea #2: SIMD**

Amortize cost/complexity of managing an instruction stream across many ALUs

[Fatahalian '08]

# SIMD



**Fetch/ Decode**

| ALU 1 | ALU 2 | ALU 3 | ALU 4 |

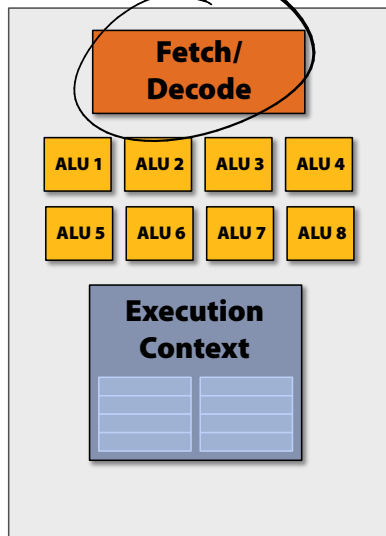| ALU 5 | ALU 6 | ALU 7 | ALU 8 |

**Execution Context**
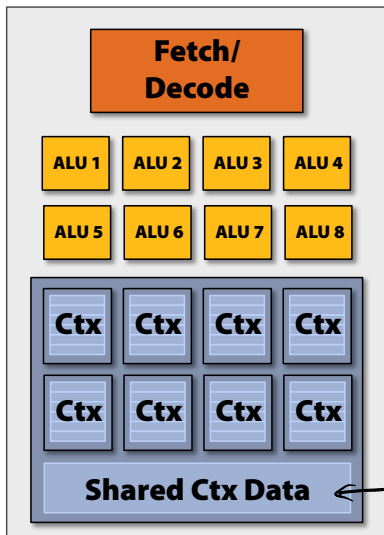
[Fatahalian '08]

Idea #2: SIMD

Amortize cost/complexity of managing an instruction stream across many ALUs

# SIMD



**Idea #2: SIMD**

Amortize cost/complexity of managing an instruction stream across many ALUs

Diagram labels: Fetch/Decode, ALU 1, ALU 2, ALU 3, ALU 4, ALU 5, ALU 6, ALU 7, ALU 8, Ctx (×8), Shared Ctx Data ← scratch pad

[Fatahalian '08]

# Latency Hiding

▶ Latency (mem, pipe) hurts non-OOO cores

▶ Do *something* while waiting
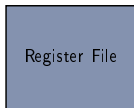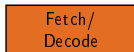
What is the unit in which work gets scheduled on a GPU?

Vector
Nvidia: warp / wavefront ✓
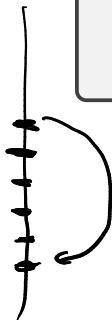
How can we keep busy?

- SMT
- ILP

Change in architectural picture?

Before:

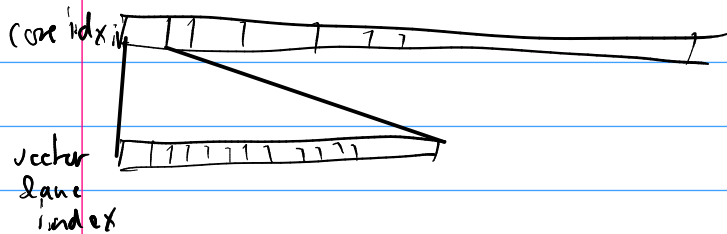Fetch/
Decode

Register File

Scratchpad/L1

After:

more
state space

# GPUs: Core Architecture Ideas

Three core ideas:

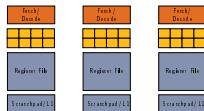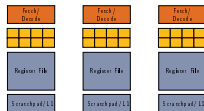- Simpler / many of cores
- SIMD
- latency hiding through concurrency

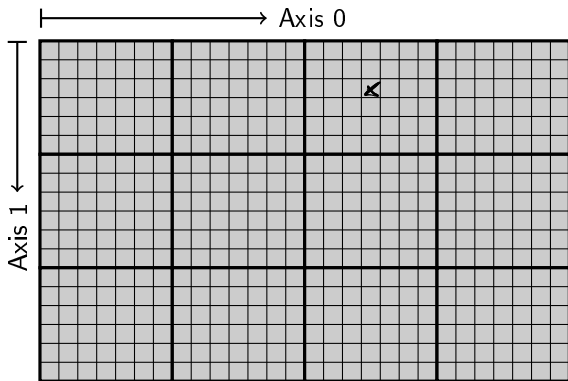core idx in 

vector
lane
index

$(core\ idx,\ vector\ lane\ idx)$

$x = vec\ \%\ 17$

$y = vec\ //\ 17$

‘SIMT’

# Wrangling the Grid



Axis 0

Axis 1

Handwritten annotations:
blockIdx.xyz
gridDim: 3
1:1

Dim 0: 4
1: 1

threadIdx.xyz

- ▶ get_local_id(axis)?/size(axis)?
- ▶ get_group_id(axis)?/num_groups(axis)?
- ▶ get_global_id(axis)?/size(axis)?
axis=0,1,2,...

# Demo CL code

Demo: machabstr/Hello GPU

# 'SIMT' and Branches



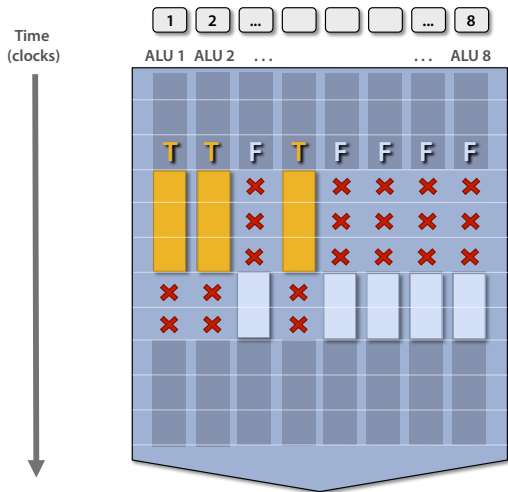[Fatahalian '08]

# GPU Abstraction: Core Model Ideas

How do these aspects show up in the model?

- ▶ View concrete counts as an implementation detail
  - ▶ SIMD lane
  - ▶ Core
  - ▶ Scheduling slot
- ▶ Program as if there are infinitely many of them
- ▶ Hardware division is expensive
  Make $n$D grids part of the model to avoid it
- ▶ Design the model to expose *extremely* fine-grain concurrency
  (e.g. between loop iterations!)
- ▶ Draw from the same pool of concurrency to hide latency

# GPU Program 'Scopes'

| Hardware | CL adjective | OpenCL | CUDA |
| --- | --- | --- | --- |
| SIMD lane | private | Work Item | Thread |
| SIMD Vector | — | Subgroup | Warp |
| Core | local | Workgroup | Thread Block |
| Processor | global | NDRange | Grid |

# GPU: Communication

What forms of communication exist at each scope?

- SIMD lanes; Vector shuffles
- Scratchpad + barrier, atomics + memfences
- across workgroups; atomics + memfences

Can we just do locking like we might do on a CPU?

no: indep. fw progress required

# GPU Programming Model: Commentary

- "Vector" / "Warp" / "Wavefront"
  - Important hardware granularity
  - Poorly/very implicitly represented
- What is the impact of reconvergence?