

April 1, 2025

Announcements

- Project feed back; soon
- github.com/google/highway

Goals

Review

- Core?
- something that fetches insns.
- Core has a "scratchpad"
- in order
- funny, pseudo-scalar programming model
- a lot of "registers" ?

- latency hiding
 - ↳ "variable-size" register file
 - ↳ SMT
- hardware scheduler

Outline

Introduction

Machine Abstractions

C

OpenCL/CUDA

Convergence, Differences in Machine Mapping

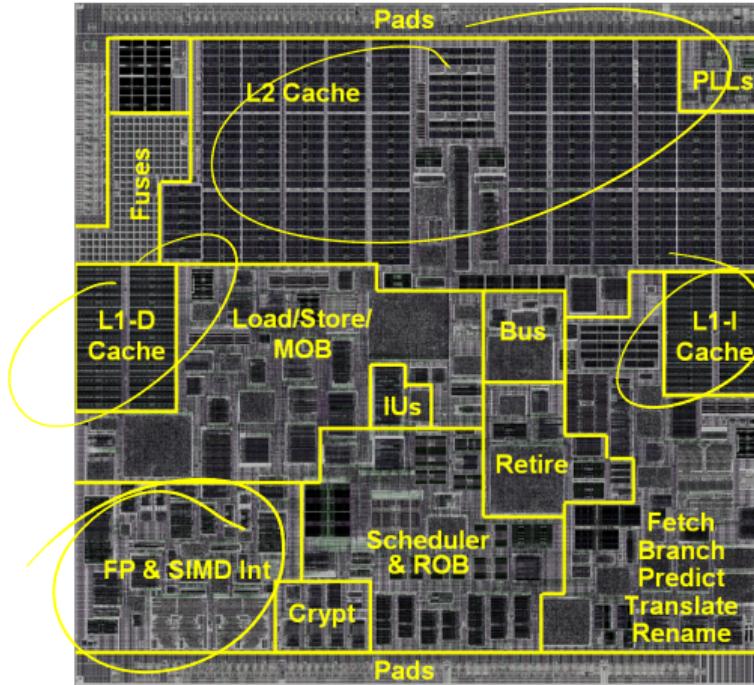
Lower-Level Abstractions: SPIR-V, PTX

Performance: Expectation, Experiment, Observation

Performance-Oriented Languages and Abstractions

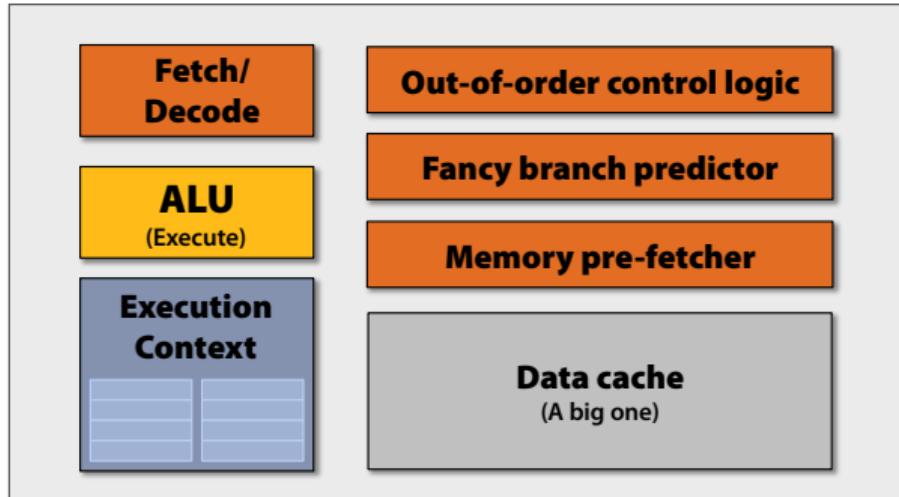
Polyhedral Representation and Transformation

Chip Real Estate



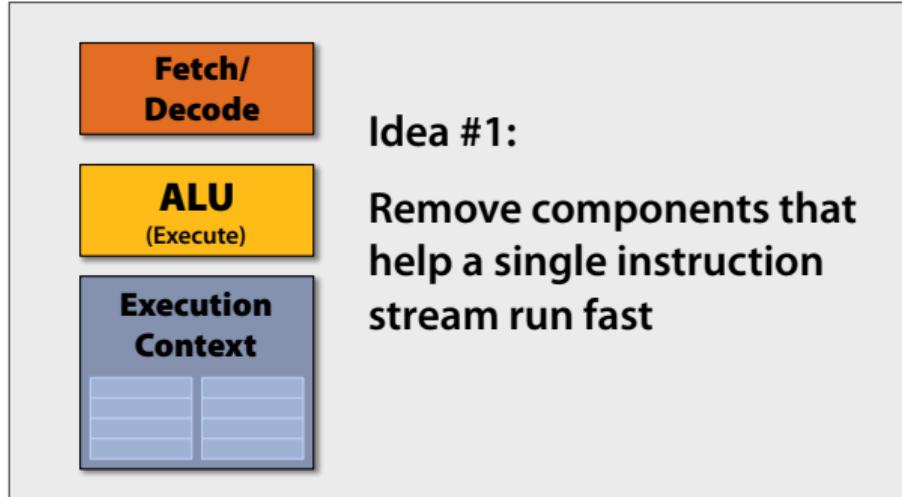
Die floorplan: VIA Isaiah (2008).
65 nm, 4 SP ops at a time, 1 MiB L2.

“CPU-style” Cores



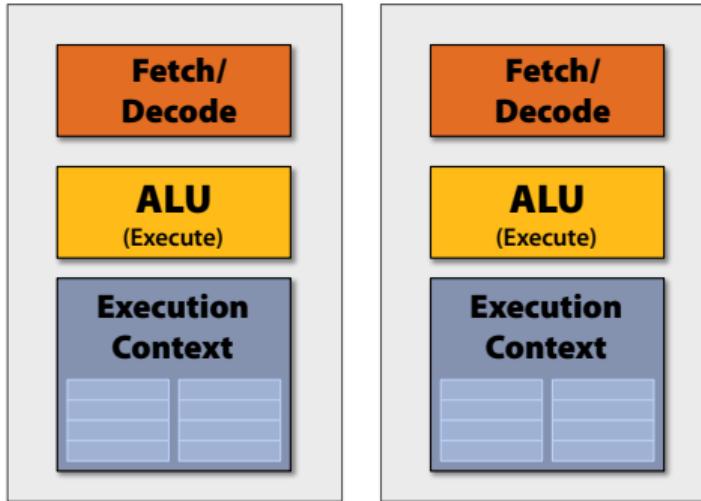
[Fatahalian ‘08]

Slimming down



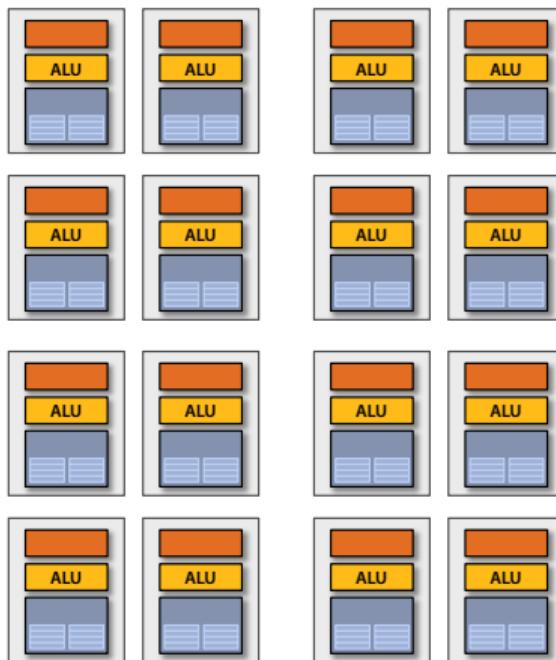
[Fatahalian '08]

More Space: Double the Number of Cores



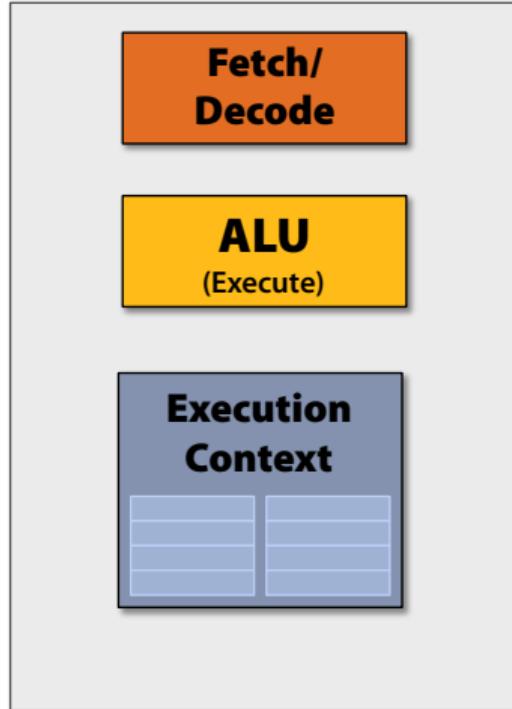
[Fatahalian '08]

Even more



[Fatahalian '08]

SIMD

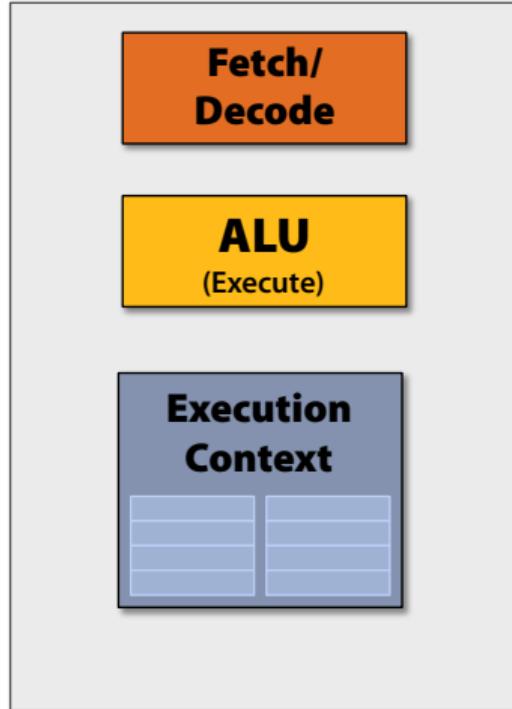


Idea #2: SIMD

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

[Fatahalian '08]

SIMD

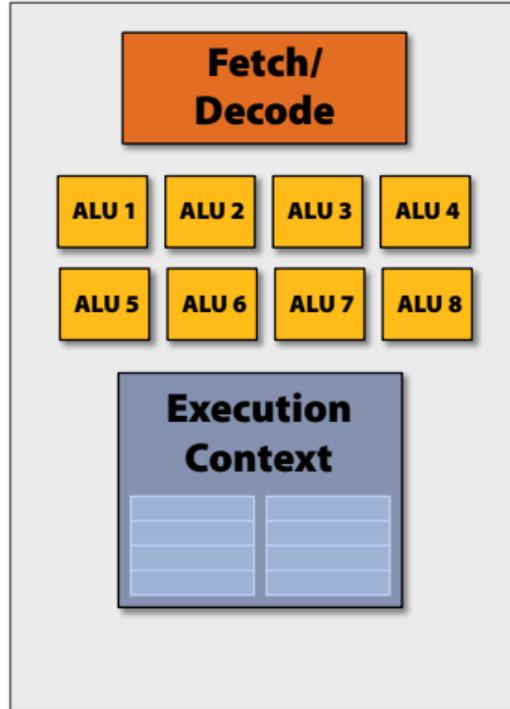


Idea #2: SIMD

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

[Fatahalian '08]

SIMD

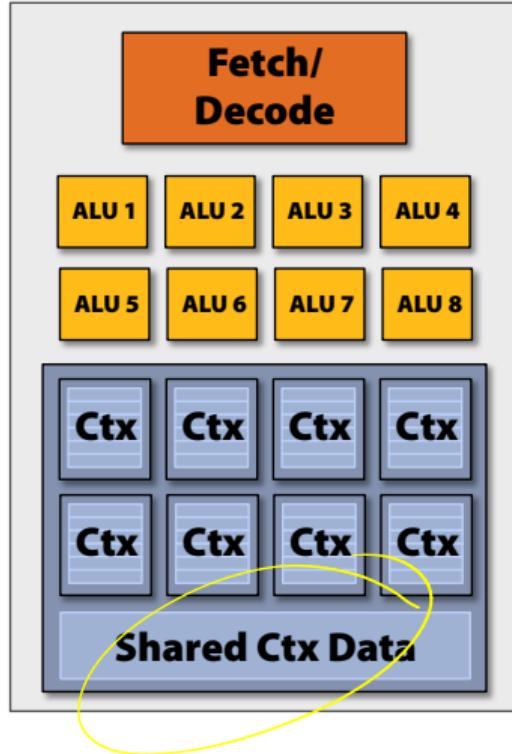


Idea #2: SIMD

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

[Fatahalian '08]

SIMD



[Fatahalian '08]

Idea #2: SIMD

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

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 ; "Subgroup"
Warp (32)
Wavefront (64 → 32) Nr AMG
? (8...64) Intel

How can we keep busy?

By switching to other work

After:

Before:



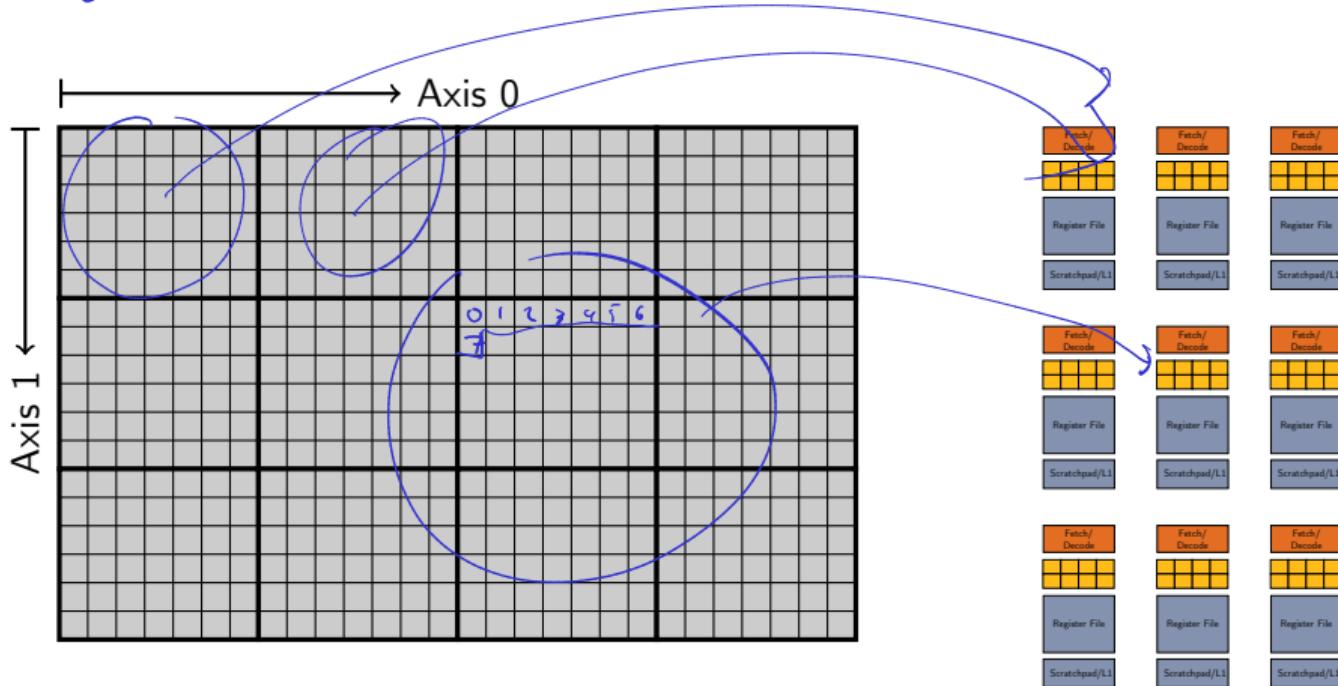
GPUs: Core Architecture Ideas

Three core ideas:

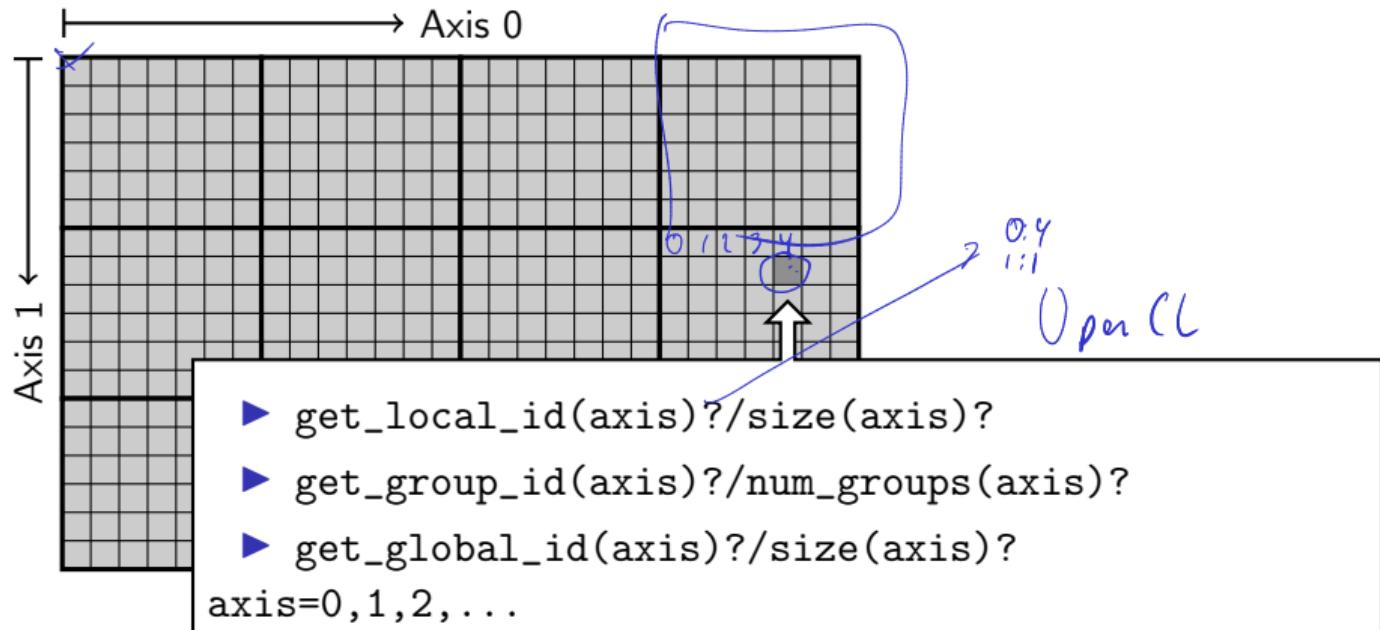
- ▷ Remove stuff for single-thread latency hiding
- ▷ lots of SIMD, lots of cores
- ▷ cover latency with concurrency
 - ▷ SMT ↗?
 - ▷ ILP ↗?

Latency hiding comes at the cost
of state space

'SIMT'
D
"thread"



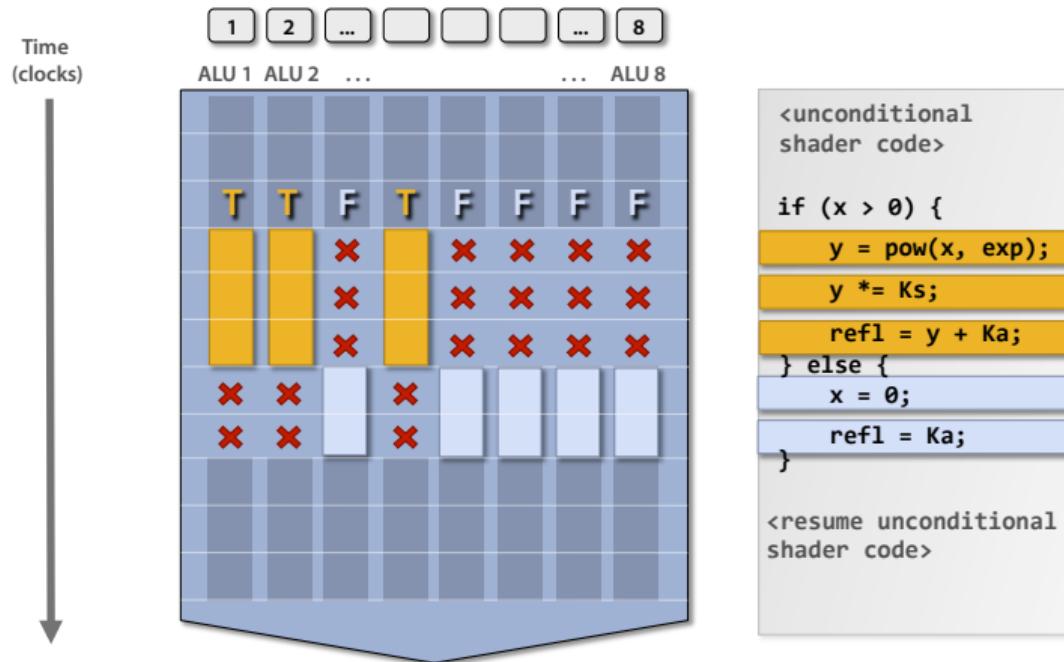
Wrangling the Grid



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?

- ▷ Concrete config (#cores, #lanes) as an impl. detail
 # scheduling slots
- ▷ Program as it is
- ▷ no grid because division is slow

int m, n

$m \% n$

$m \% 8 \rightarrow$ bitwise and

$m \% 7 \rightarrow$ Barrett reduction
Lemire

GPU Address Spaces

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

A Venn diagram with two overlapping circles. The left circle is labeled "OpenCL" and contains the words "Work Item", "Subgroup", "Workgroup", and "NDRange". The right circle is labeled "CUDA" and contains the words "Thread", "Warp", "Thread Block", and "Grid". The overlapping region between the two circles is shaded blue.

GPU: Communication

What forms of communication exist at each scope?

Subgroup:

Shuffles (!)

Workgroup:

Scratchpad + atomic
Global DRAM + atomic

Grid:

Scratchpad + barrier
Global DRAM + atomic

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

- ▷ NO (see Sorenson paper)
- ▷ Independent forward progress (only some GPUs)
needed for intra-group sync.
- ▷ Device partitioning

GPU Programming Model: Commentary

Advantages:

.

Disadvantages: