

A large teal geometric shape, resembling a stylized 'A' or a series of nested triangles, occupies the left and bottom-left portions of the slide. A smaller, separate teal parallelogram is positioned above the main shape, to the right of the center.

THE AMD gem5 APU SIMULATOR: MODELING GPUS USING THE MACHINE ISA

TONY GUTIERREZ, SOORAJ PUTHOOR, TUAN TA*, MATT SINCLAIR,
AND BRAD BECKMANN
AMD RESEARCH, *CORNELL
JUNE 2, 2018

OBJECTIVES AND SCOPE



▲ Objectives

- Introduce the Radeon Open Compute Platform (ROCm)
- AMD's Graphics Core Next (GCN) architecture and GCN3 ISA
- Describe the gem5-based APU simulator

▲ Scope

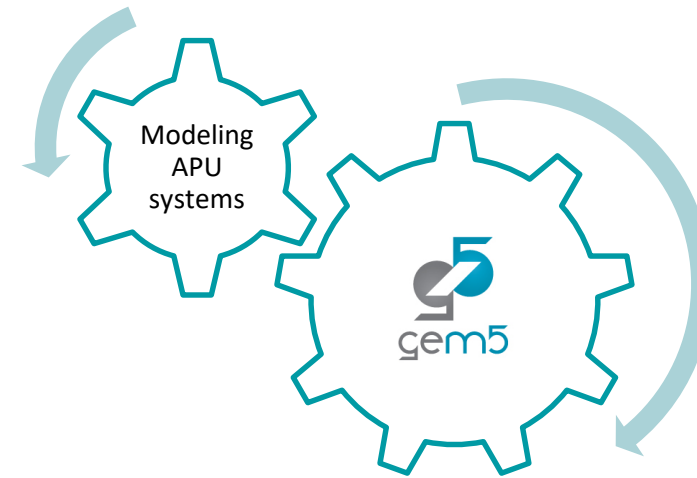
- Emphasis on the GPU side of the simulator
 - APU (CPU+GPU) model, not discrete GPU
 - Covers GPU arch, GCN3 ISA, and HW-SW interfaces

▲ Why are we releasing our code?

- Encourage AMD-relevant research
- Modeling ISA and real system stack is important [1]
- Enhance academic collaborations
 - Enable intern candidates to get experience before arriving
 - Enable interns to take their experience back to school

▲ Acknowledgement

- AMD Research's gem5 team



[1] Gutierrez et al. *Lost in Abstraction: Pitfalls of Analyzing GPUs at the Intermediate Language Level*. HPCA, 2018.

QUICK SURVEY



- ▲ Who is in our audience?
 - Graduate students
 - Faculty members
 - Working for government research labs
 - Working for industry
- ▲ Have you written an GPU program?
 - CUDA, OpenCL™, HIP, HC, C++ AMP, other languages
- ▲ Have you used these simulators?
 - GPGPU-Sim
 - Multi2Sim
 - gem5
 - Our HSAIL-based APU model
- ▲ Are you familiar with our HPCA 2018 paper?
 - *Lost in Abstraction: Pitfalls of Analyzing GPUs at the Intermediate Language Level*

OUTLINE



Topic	Presenter	Time
Background	Tony	8:00 – 8:15
ROCm, GCN3 ISA, and GPU Arch	Tony	8:15 – 9:15
HSA Implementation in gem5	Sooraj	9:15 – 10:00
Break		10:00 – 10:30
Ruby and GPU Protocol Tester	Tuan	10:30 – 11:15
Demo and Workloads	Matt	11:15 – 11:50
Summary and Questions	All	11:50 – 12:00

- ▲ Overview of gem5
 - Source tree
- ▲ GPU terminology and system overview
- ▲ HSA standard and building blocks
 - Coherent shared virtual memory
 - User-level queues
 - Signals
 - etc.

OVERVIEW OF gem5



- ▲ Open-source modular platform for system architecture research
 - Integration of M5 (Univ. of Michigan) and GEMS (Univ. of Wisconsin)
 - Actively used in academia and industry
- ▲ Discrete-event simulation platform with numerous models
 - CPU models at various performance/accuracy trade-off points
 - Multiple ISAs: x86, ARM, Alpha, Power, SPARC, MIPS
 - Two memory system models: Ruby and “classic” (M5)
 - Including caches, DRAM controllers, interconnect, coherence protocols, etc.
 - I/O devices: disk, Ethernet, video, etc.
 - Full system or app-only (system-call emulation)
- ▲ Cycle-level modeling (not “cycle accurate”)
 - Accurate enough to capture first-order performance effects
 - Flexible enough to allow prototyping new ideas reasonably quickly
- ▲ See <http://www.gem5.org>
- ▲ More information available from Jason Lowe-Power’s tutorial
 - <http://learning.gem5.org/tutorial/>

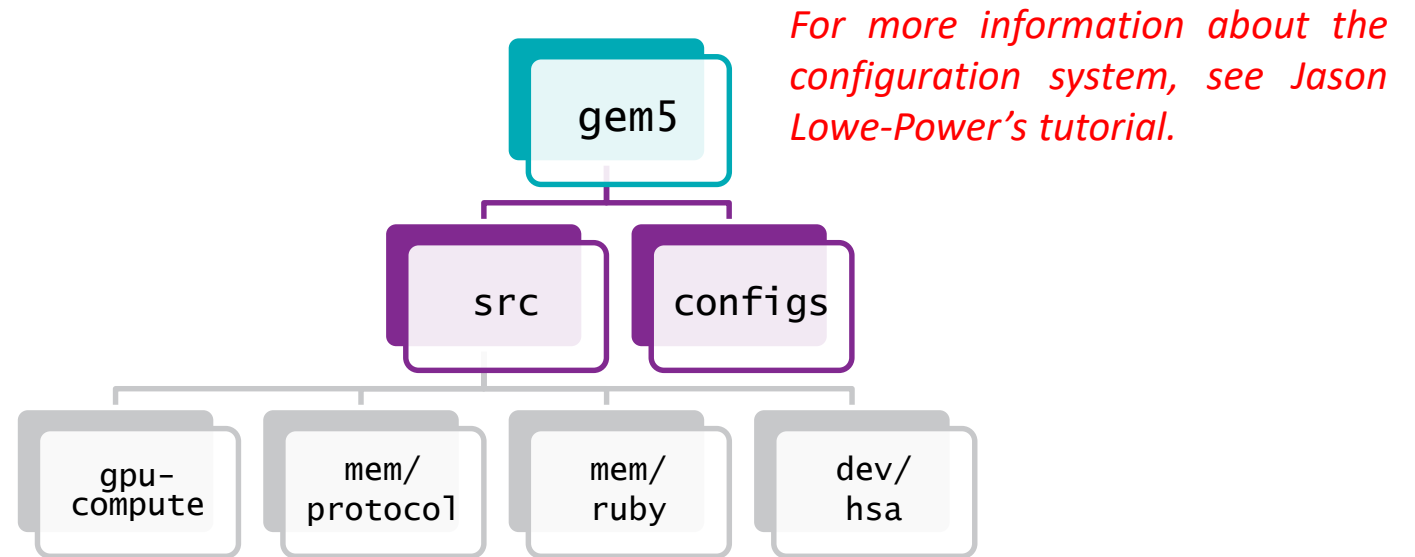


APU SIMULATOR CODE ORGANIZATION



▲ Gem5 ← top-level directory

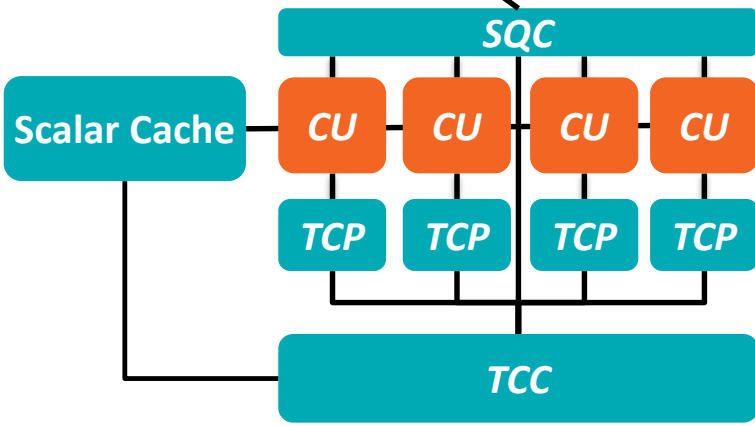
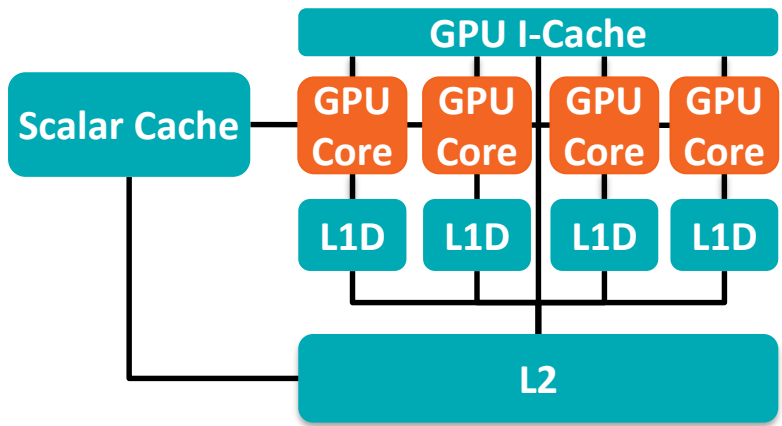
- src/
 - gpu-compute/ ← GPU core model
 - mem/protocol/ ← APU memory model
 - mem/ruby/ ← APU memory model
 - dev/hsa/ ← HSA device models
- configs/
 - example/ ← apu_se.py sample script
 - ruby/ ← APU protocol configs



For the remainder of this talk, files without a directory prefix are located in **src/gpu-compute/**

SQC: Sequencer Cache (shared L1 instruction)

CU: Compute Unit (SM in NVIDIA terminology)



TCP: Texture Cache per Pipe (private L1 data)

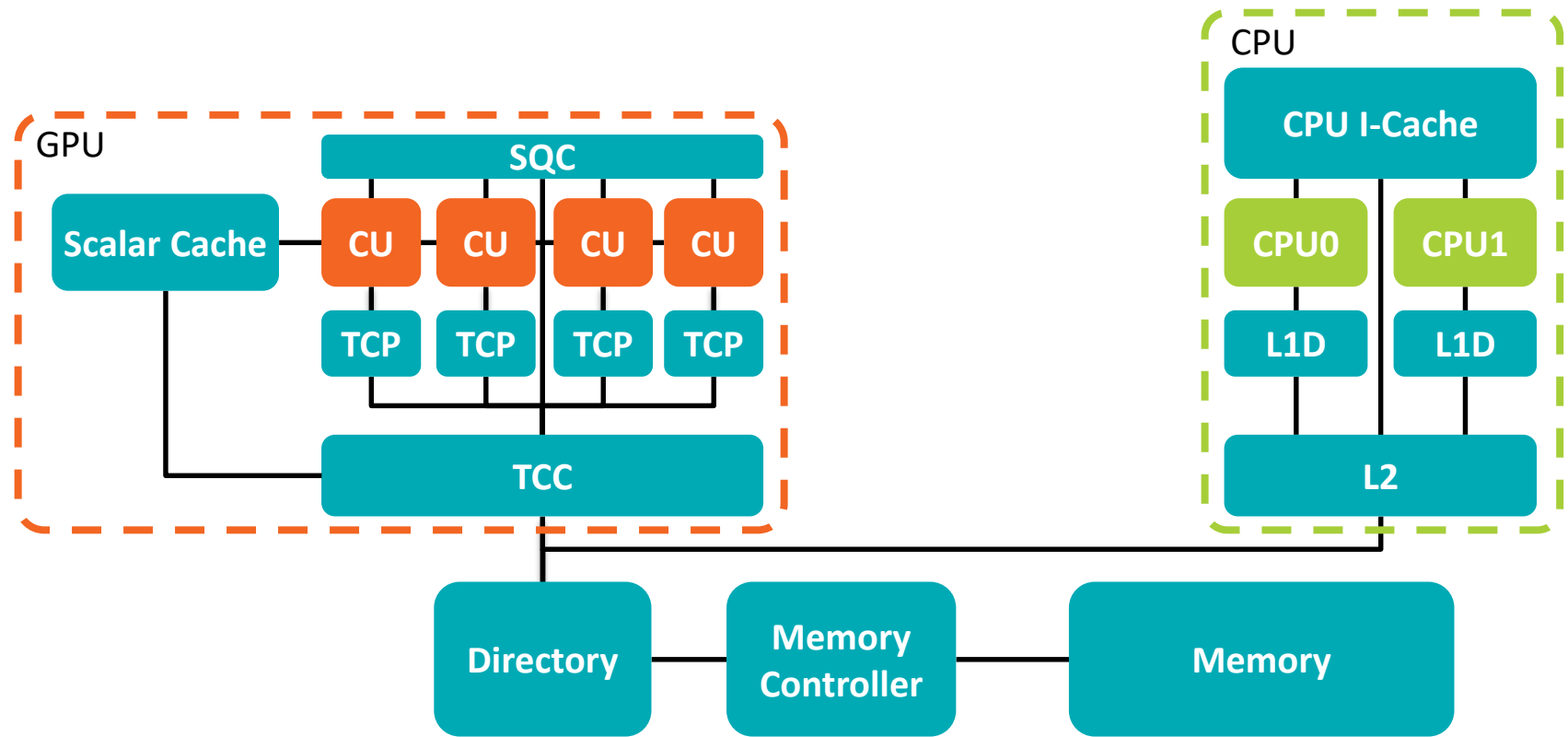
TCC: Texture Cache per Channel (shared L2)

AMD terminology

Not shown (per GPU core):
LDS: Local Data Share (Shared memory in NVIDIA terminology, sometimes called “scratch pad”)

EXAMPLE APU SYSTEM

GPU + CPU CORE-PAIR WITH A SHARED DIRECTORY

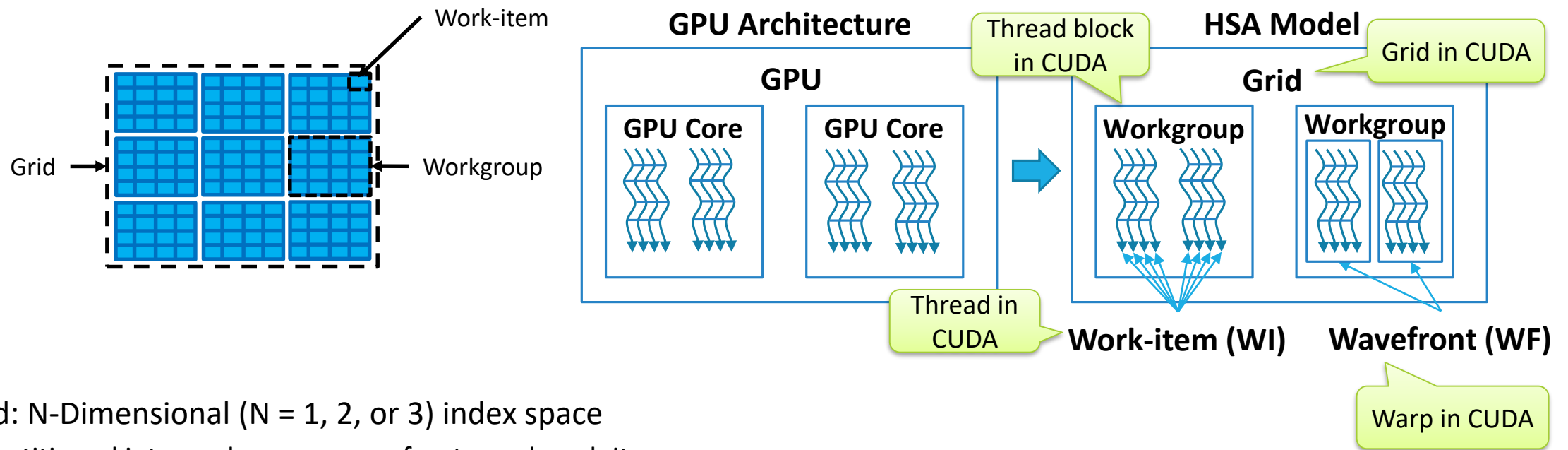


AMD TERMINOLOGY IN A NUTSHELL



▲ Heterogeneous Systems Architecture (HSA) programming abstraction

- Standard for heterogeneous compute – supported by AMD hardware
- Light abstractions of parallel physical hardware
- Captures basic HSA and OpenCL constructs, plus much more



▲ Grid: N-Dimensional (N = 1, 2, or 3) index space

- Partitioned into workgroups, wavefronts, and work-items

SPECIFICATION BUILDING BLOCKS



HSA Hardware Building Blocks

- ▲ Shared Virtual Memory
 - Single address space
 - Coherent
 - Pageable
 - Fast access from all components
 - Can share pointers
- ▲ Architected User-Level Queues
- ▲ Signals
- ▲ Platform Atomics
- ▲ Defined Memory Model
- ▲ Context Switching

Open-Source

HSA Platform System Arch Specification



Industry standard, architected requirements for how devices share memory and communicate with each other

HSA Software Building Blocks

- ▲ HSA Runtime
 - *Implemented by the ROCm runtime*
 - Create queues
 - Allocate memory
 - Device discovery
- ▲ Multiple high-level compilers
 - CLANG/LLVM
 - C++, HIP, OpenMP, OpenACC, Python
- ▲ GCN3 Instruction Set Architecture
 - Kernel state
 - ISA encodings
 - Program flow control

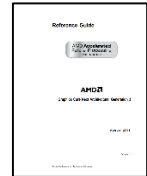
Open-Source

HSA System Runtime Specification



Open-Source

GCN3 ISA Specification



Industry specifications to enable existing programming languages to target the GPU

<http://hsafoundation.com>

<http://github.com/HSAFoundation>

HSA Hardware Building Blocks

- ▲ Shared virtual memory
 - Single address space
 - Coherent
 - Fast access from all components
 - Can share pointers
 - Pageable
- ▲ Architected user-level queues
 - Via architected queuing language (AQL)
- ▲ Signals
- ▲ Platform atomics
- ▲ Defined memory model
 - Acquire and release semantics as implemented by the compiler
 - Merging functional and timing models
- ▲ Context switching

HSA Software Building Blocks

- ▲ Radeon Open Compute platform (ROCm)
 - AMD's implementation of HSA principles
 - Create queues
 - Device discovery
 - AQL support
 - Allocate memory
- ▲ Machine ISA
 - GCN3
- ▲ Heterogeneous Compute Compiler (HCC)
 - CLANG/LLVM – direct to GCN3 ISA
 - C++, C++ AMP, HIP, OpenMP, OpenACC, Python

Legend

Included in this release

Work-in-progress / may be released

Longer term work

OUTLINE



Topic	Presenter	Time
Background	Tony	8:00 – 8:15
ROCm, GCN3 ISA, and GPU Arch	Tony	8:15 – 9:15
HSA Implementation in gem5	Sooraj	9:15 – 10:00
Break		10:00 – 10:30
Ruby and GPU Protocol Tester	Tuan	10:30 – 11:15
Demo and Workloads	Matt	11:15 – 11:50
Summary and Questions	All	11:50 – 12:00

HW-SW INTERFACES



- ▲ ROCm – high-level SW stack
- ▲ HW-SW interfaces
- ▲ Kernel launch flow
- ▲ GCN3 ISA overview

SW STACK AND HIGH-LEVEL SIMULATION FLOW



- Clang front end and LLVM-based backend
 - Direct to ISA
 - Multi-ISA binary (x86 + GCN3)

- HCC libraries
- Runtime layer – **ROCr**
- Thunk (user space driver) – **ROCr-t**
- Kernel fusion driver (KFD) – **ROCr-k**

- Command processor (CP) HW aids in implementing HSA standard
- Rich application binary interface (ABI) shader, [h

- Runtime ELF loaders for GCN3 binary

```
gpu_command_processor.[hh|cc]
```



DETAILED VIEW OF KERNEL LAUNCH

GPU FRONTEND AND HW-SW INTERFACE



- ▲ User space SW talks to GPU via `ioctl()`
 - HCC/ROCr/ROCr are off-the-shelf ROCm
 - ROCK is emulated in `gem5`
 - Handles `ioctl` commands

▲ CP frontend

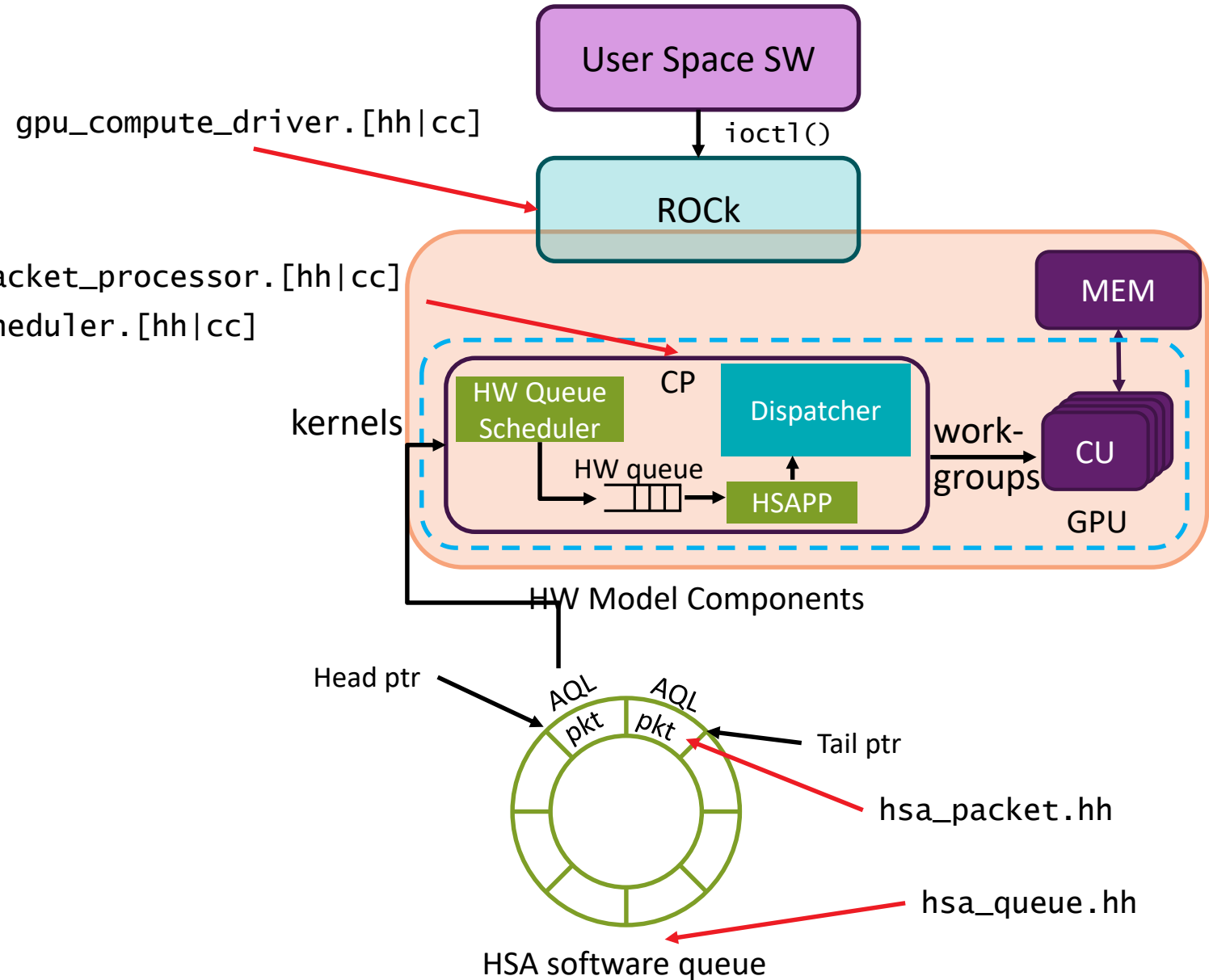
- Two primary components:
 - HSA packet processor (HSAPP)
 - Workgroup dispatcher

▲ Runtime creates soft HSA queues

- HSAPP maps them to hardware queues
- HSAPP schedules active queues

▲ Runtime creates and enqueues AQL packets

- Packets include:
 - Kernel resource requirements
 - Kernel size
 - Kernel code object pointer
 - More...

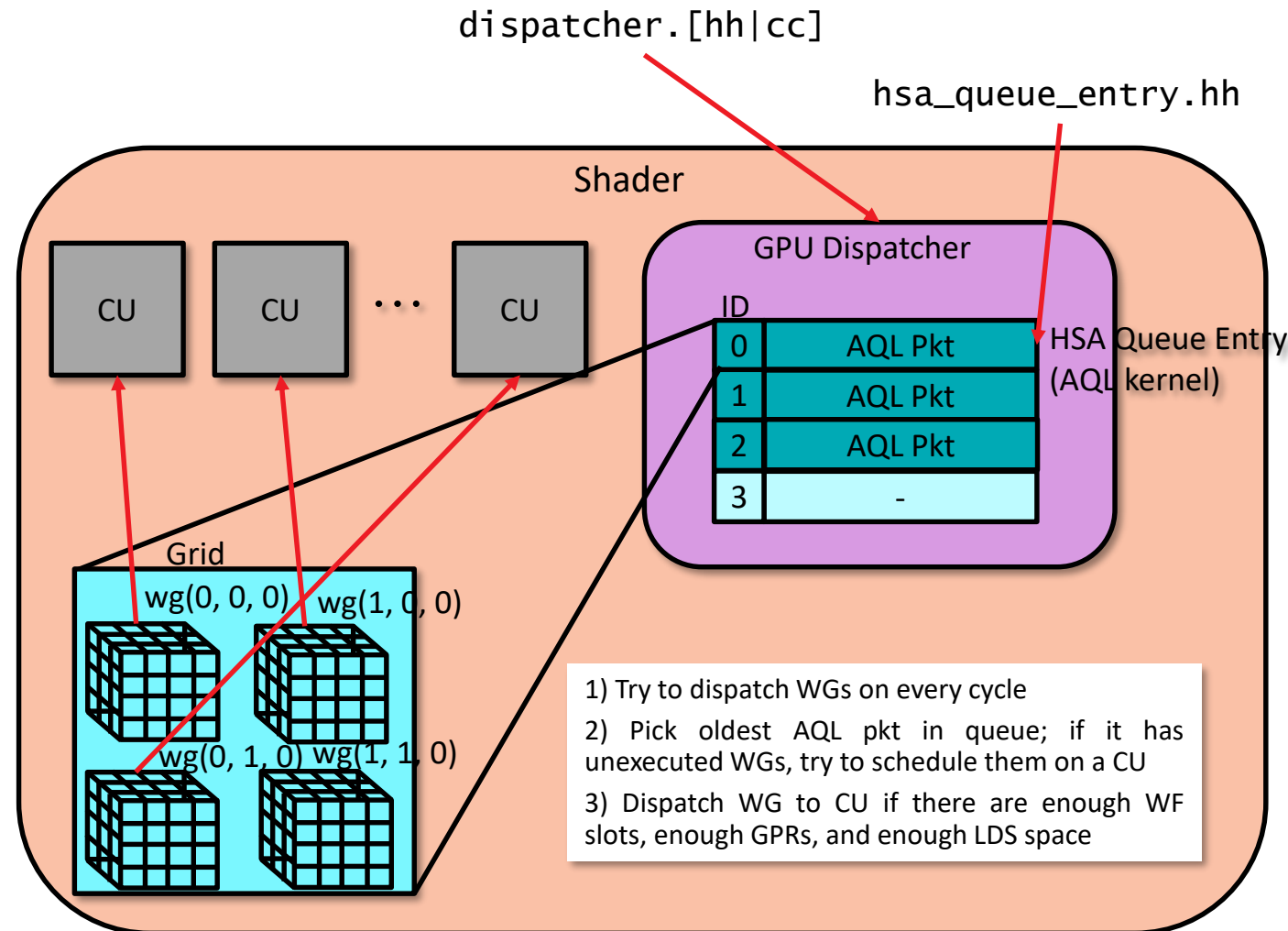


DETAILED VIEW OF KERNEL LAUNCH

DISPATCHER WORKGROUP ASSIGNMENT



- Kernel dispatch is resource limited
 - WGs are scheduled to CUs
- Dispatcher tracks status of in-flight/pending kernels
 - If a WG from a kernel cannot be scheduled, it is enqueued until resources become available
 - When all WGs from a task have completed, the dispatcher frees CU resources and notifies the host



DETAILED VIEW OF KERNEL LAUNCH

GPU ABI INITIALIZATION

▲ ABI

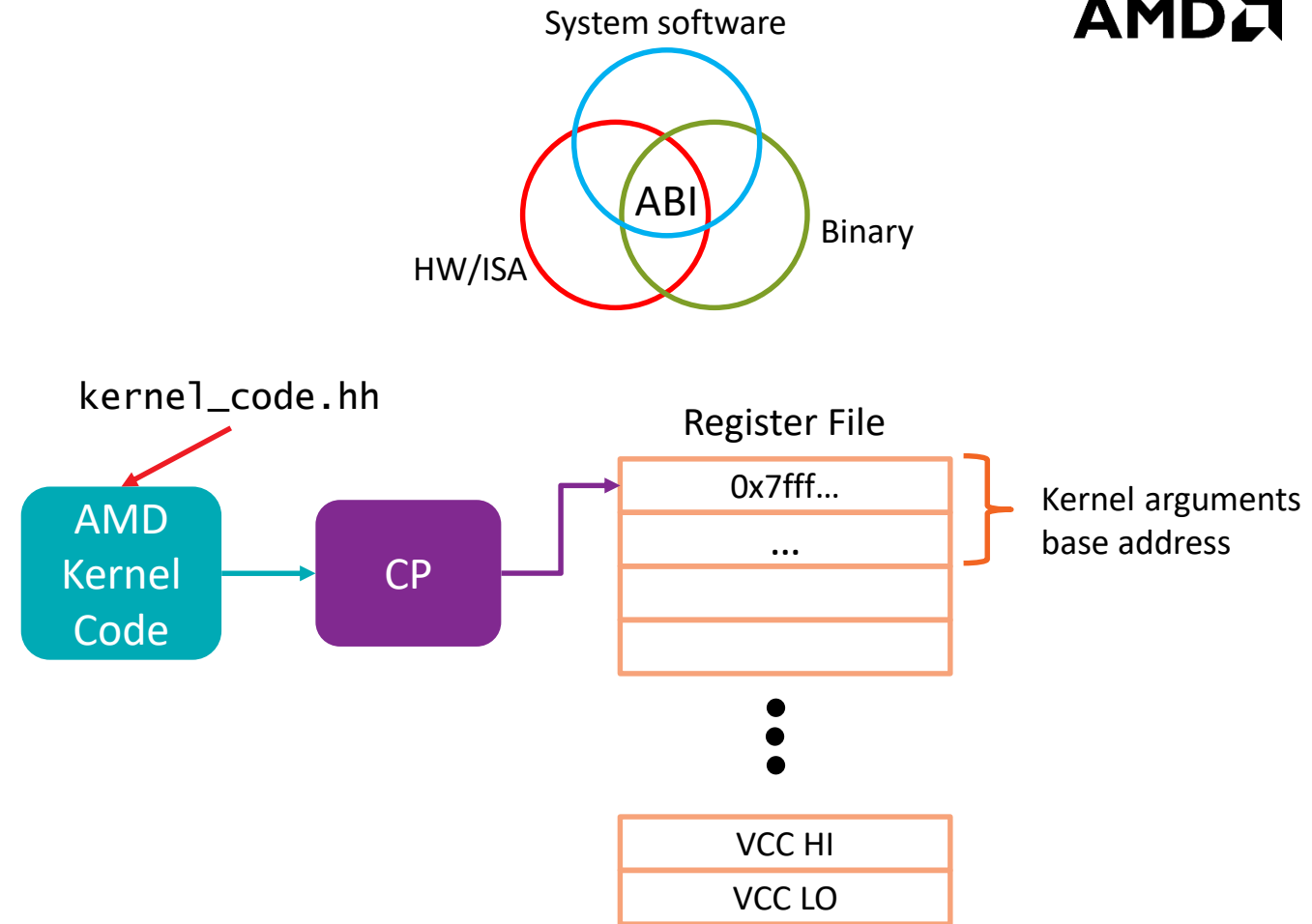
- Interface between the application binary (ELF) and other components (e.g., OS, ISA, and HW)
- Function call convention, location of special values, etc.

▲ ABI state primarily maintained in register file

- Some state maintained in the WF context for simplicity

▲ CP responsible for initializing register file state

- Extracts metadata from code object (GPU ELF binary)
- Initializes both scalar and vector registers
- Kernel argument base address, vector condition codes (VCC), etc.



GCN3 ISA application loading a kernel argument

```
# Load 2nd 64b kernel arg into s[2:3]
s_load_dwordx2 s[2:3], s[0:1], 0x08
```

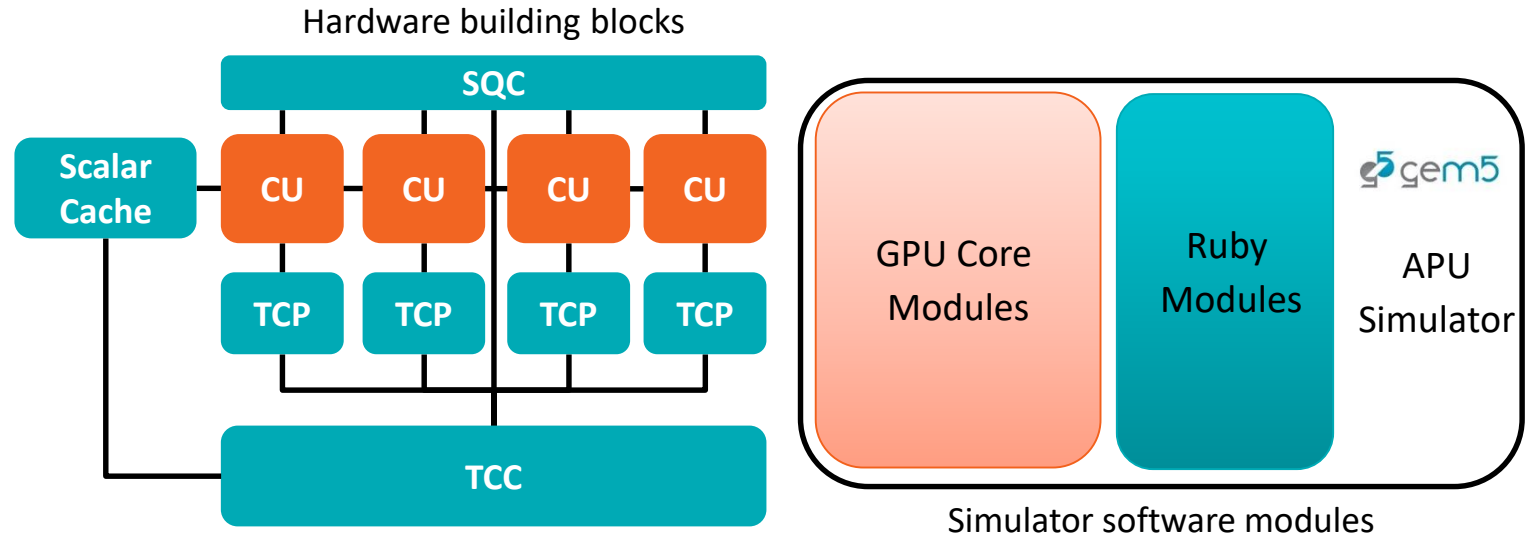
- ▲ High-level gem5 and Ruby core models
- ▲ gem5 ISA/HW separation
 - Object oriented design
 - Modular, extensible...
 - Possible to support multiple-ISAs
- ▲ GCN microarchitecture
- ▲ gem5's conceptual pipeline and timing flow
- ▲ GCN3 ISA description
- ▲ gem5 compute unit implementation
 - Pipeline breakdown

GPU CORE MODULES

GPU CORE MODULES VS. RUBY MODULES



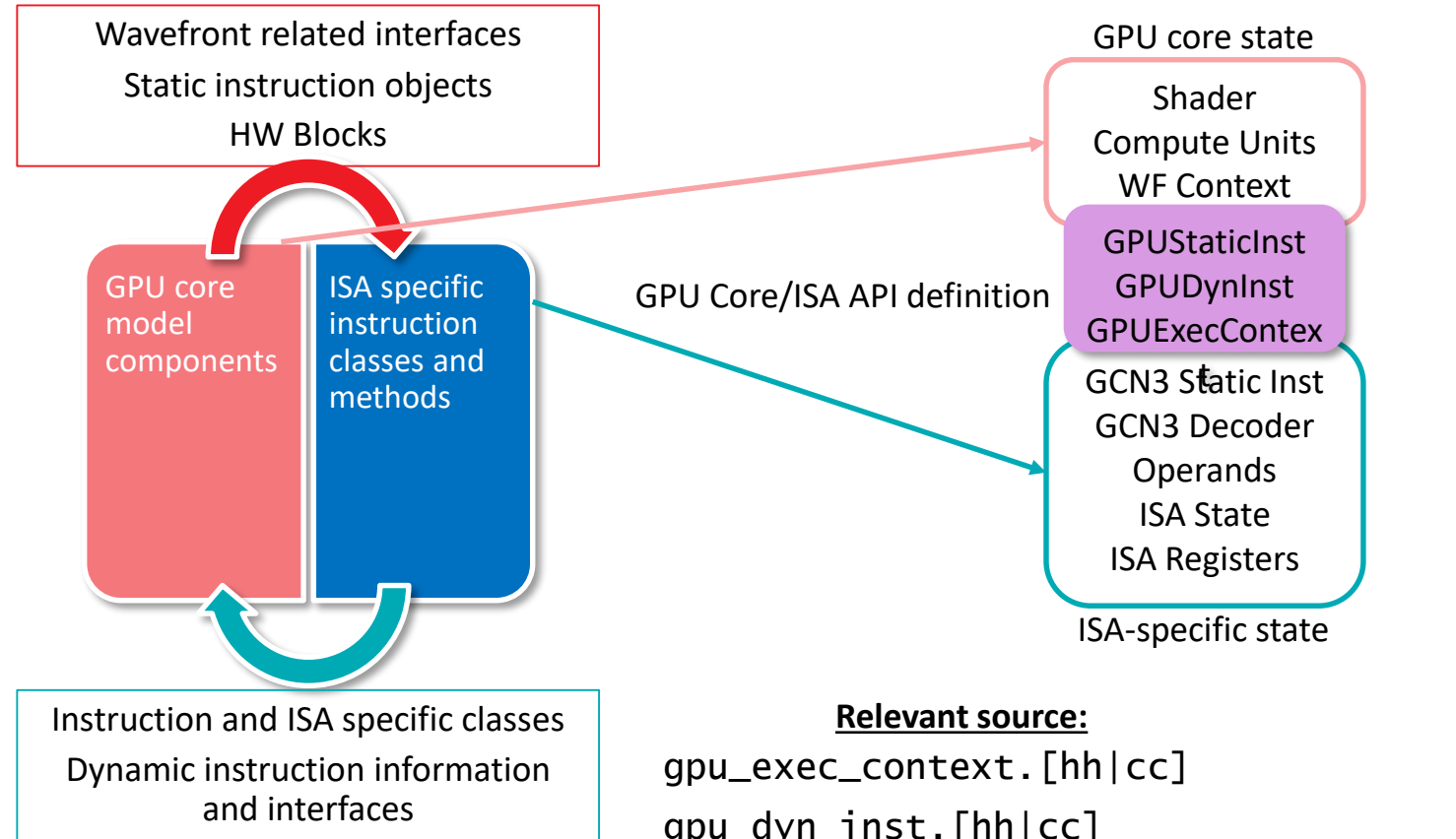
- GPU core is the compute unit
 - Resources inside GPU Core
 - Instruction buffering, Registers, Vector ALUs
 - Resources outside GPU Core
 - TCP, TCC, SQC (Ruby based)
- Shader: object containing all GPU core models
 - The top-level view of the GPU in gem5
 - Contains multiple CUs
 - With other miscellaneous components



ISA DESCRIPTION/MICROARCHITECTURE SEPARATION



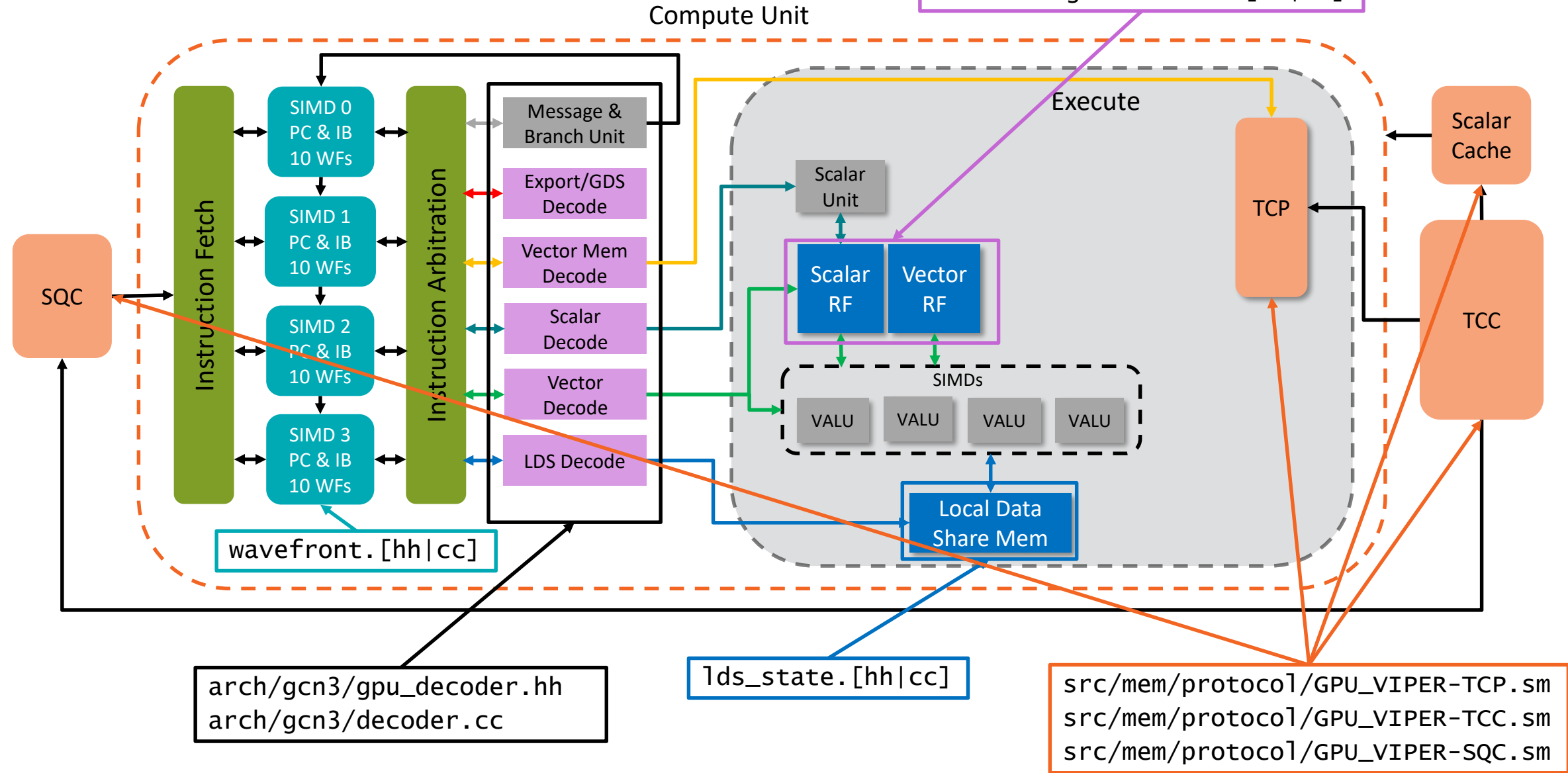
- ▲ GPUStaticInst & GPUDynInst
 - Architecture-specific code src/arch/
 - Base instruction classes
 - ISA decoder
 - ISA state, etc.
 - Define API for instruction execution
 - e.g., execute() – perform instruction execution
 - Implemented by ISA-specific instruction classes
- ▲ GPUExecContext & GPUISA
 - Define API for accessing ISA state
- ▲ Object-oriented design
 - Base classes define the API
 - Model uses base class pointers
 - Configuration system instantiates ISA-specific objects



Relevant source:

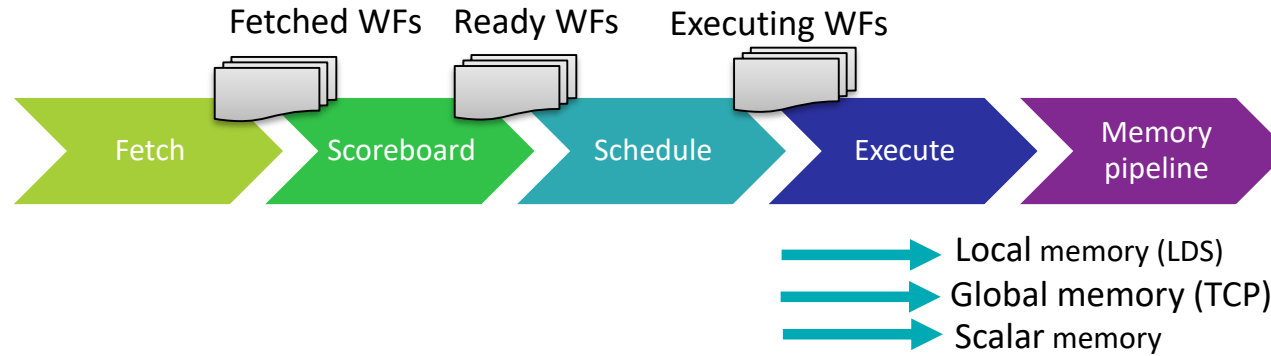
```
gpu_exec_context.[hh|cc]  
gpu_dyn_inst.[hh|cc]  
src/arch/gcn3/gpu_static_inst.hh|cc]  
src/arch/gcn3/gpu_isa.hh,  
src/arch/gcn3/isa.cc  
src/arch/gcn3/operand.hh  
src/arch/gcn3/registers.[hh|cc]
```

GPU CORE BASED ON GCN ARCHITECTURE



GPU CORE TIMING

CONCEPTUAL TIMING STAGES



▲ Execute-in-execute philosophy

▲ Pipeline stages

- Fetch: fetch for dispatched WFs - `fetch_stage.[hh|cc]` and `fetch_unit.[hh|cc]`
- Scoreboard: Check which WFs are ready - `scoreboard_check_stage.[hh|cc]`
- Schedule: Select a WF from the ready pool - `schedule_stage.[hh|cc]`
- Execute: Run WF on execution resource - `exec_stage.[hh|cc]`
- Memory pipeline: Execute LDS/global memory operation
 - `local_memory_pipeline.[hh|cc]`
 - `global_memory_pipeline.[hh|cc]`
 - `scalar_memory_pipeline.[hh|cc]`

GCN3 GPU ISA




arch/gcn3/insts/op_encodings.[hh|cc]
arch/gcn3/insts/instructions.[hh|cc]

▲ Vector and scalar instructions

- Single instruction stream
- Not really “SIMT”
- Divergence handled by scalar unit
 - Can directly modify execution mask
 - Jump over basic blocks when EXEC = 0

▲ Instructions broken down by OP type

- Op types map to different functional units in CU
- The CU can issue one instruction to each unit in the same cycle
 - Export/GDS not supported in gem5



Op Type	Functional Unit	Usage
SOPP	Branch/Message Unit	Branching, NOPs, Barriers, waitcnts, messaging
SOPC/SOPK/SOP1/SOP2	Scalar ALU	General scalar computation/divergence handling
SMEM	Scalar Memory	Scalar memory access, cache maintenance
VOPC/VOP1/VOP2/VOP3	SIMD Unit	General SIMD computation
DS	LDS	Private scratch pad memory
MUBUF	Vector Memory	Accessing vector memory, cache maintenance
FLAT	LDS/Vector Memory	Accessing vector memory, may resolve to LDS or system memory
VINTRP/MTBUF/MIMG/EXP	Varies	Primarily used by graphics. Not currently modeled in gem5; however, the infrastructure to do so is present.

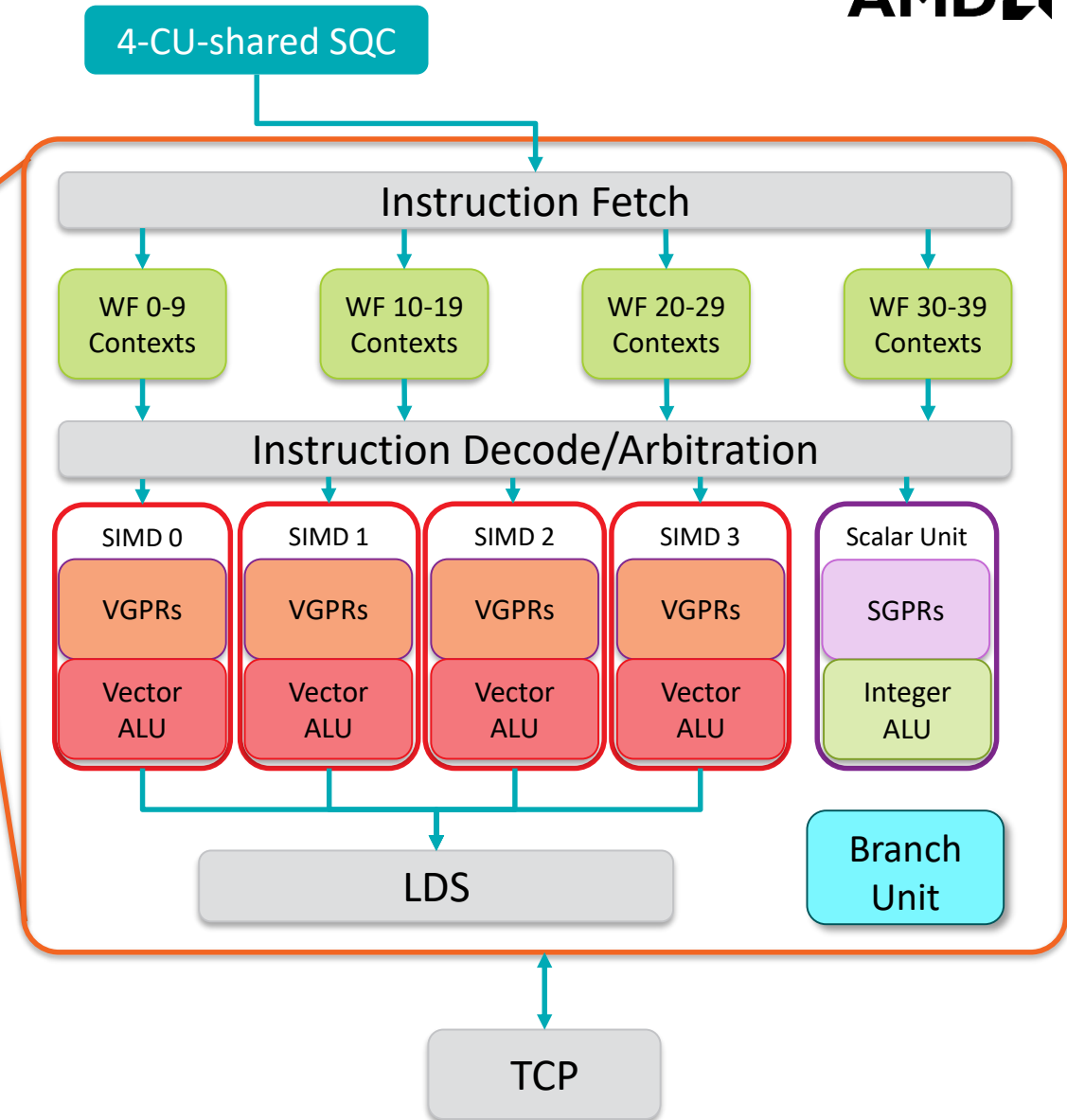
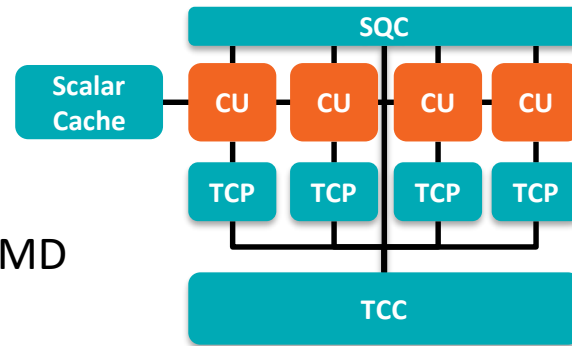
Full GCN3 spec available at: <https://gpuopen.com/compute-product/amd-gcn3-isa-architecture-manual/>

GPU CORE MODULE INTERNALS

SHARED VS. PRIVATE STRUCTURES

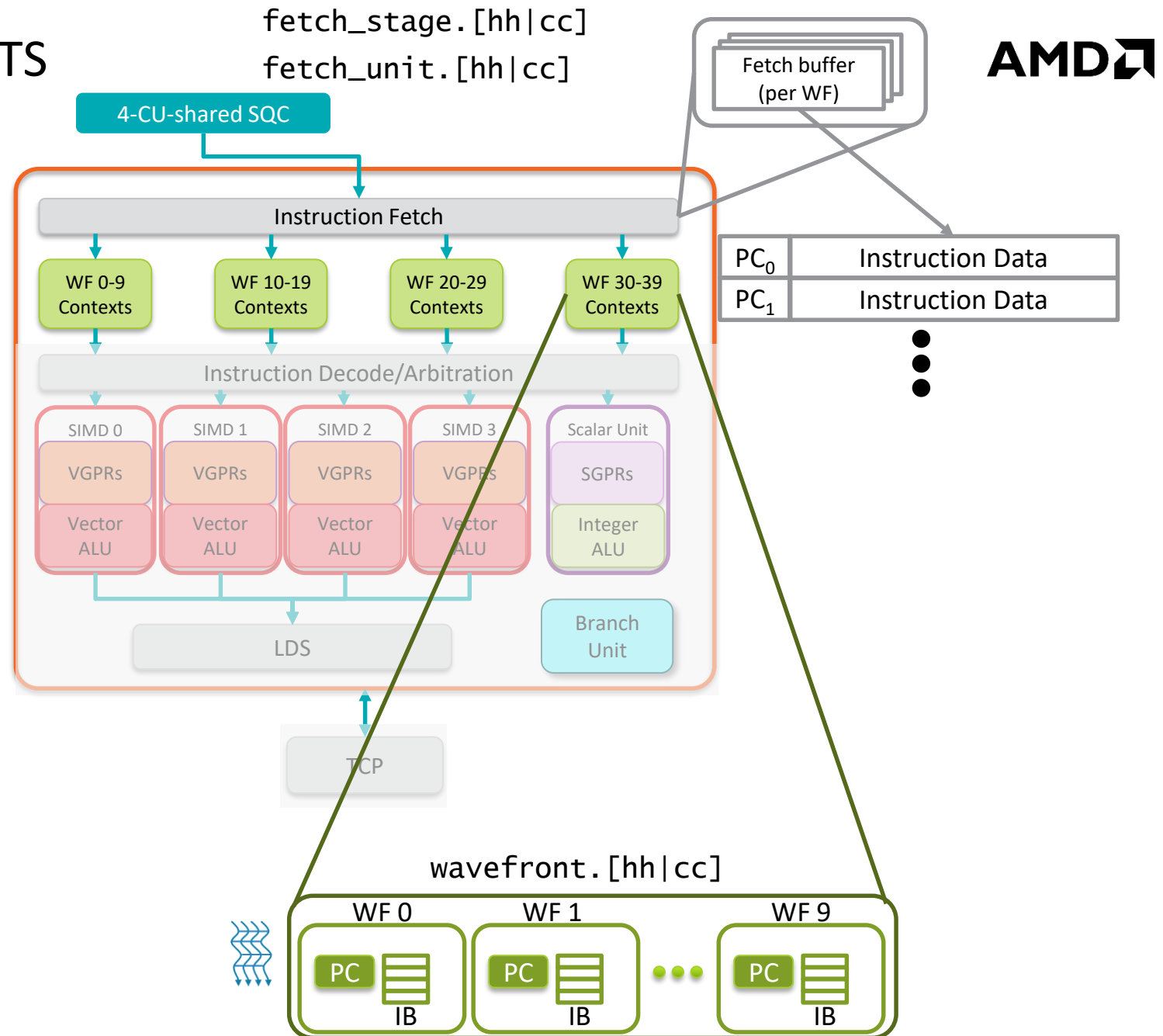
▲ Compute unit

- Four 16-wide SIMD units
- SIMD hosts WFs
- **Private** resources to each SIMD
 - Instruction buffering
 - Registers
 - Vector ALUs
- **Shared** resources
 - Fetch and decode
 - TCP
 - LDS

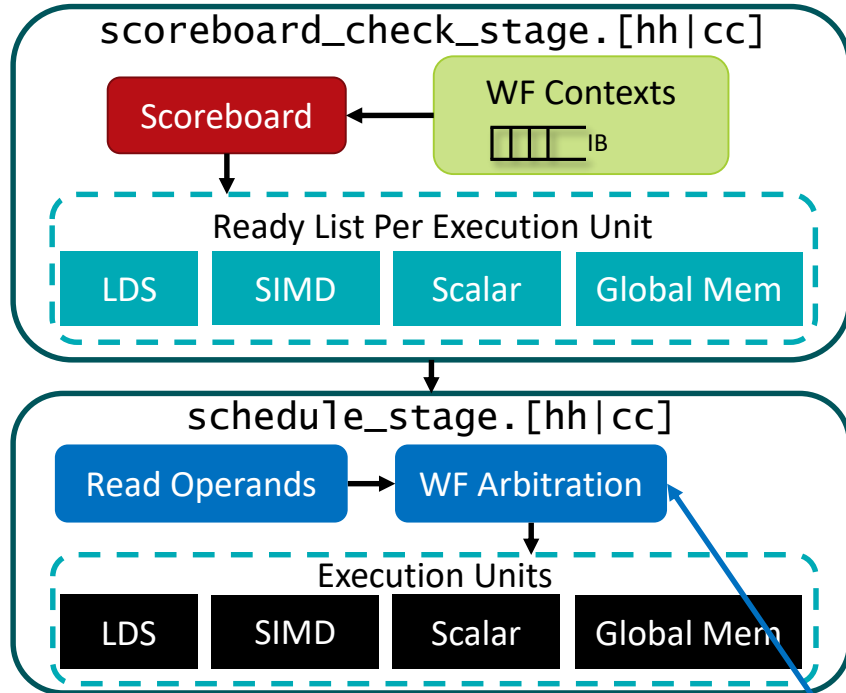


FETCH AND WAVEFRONT CONTEXTS

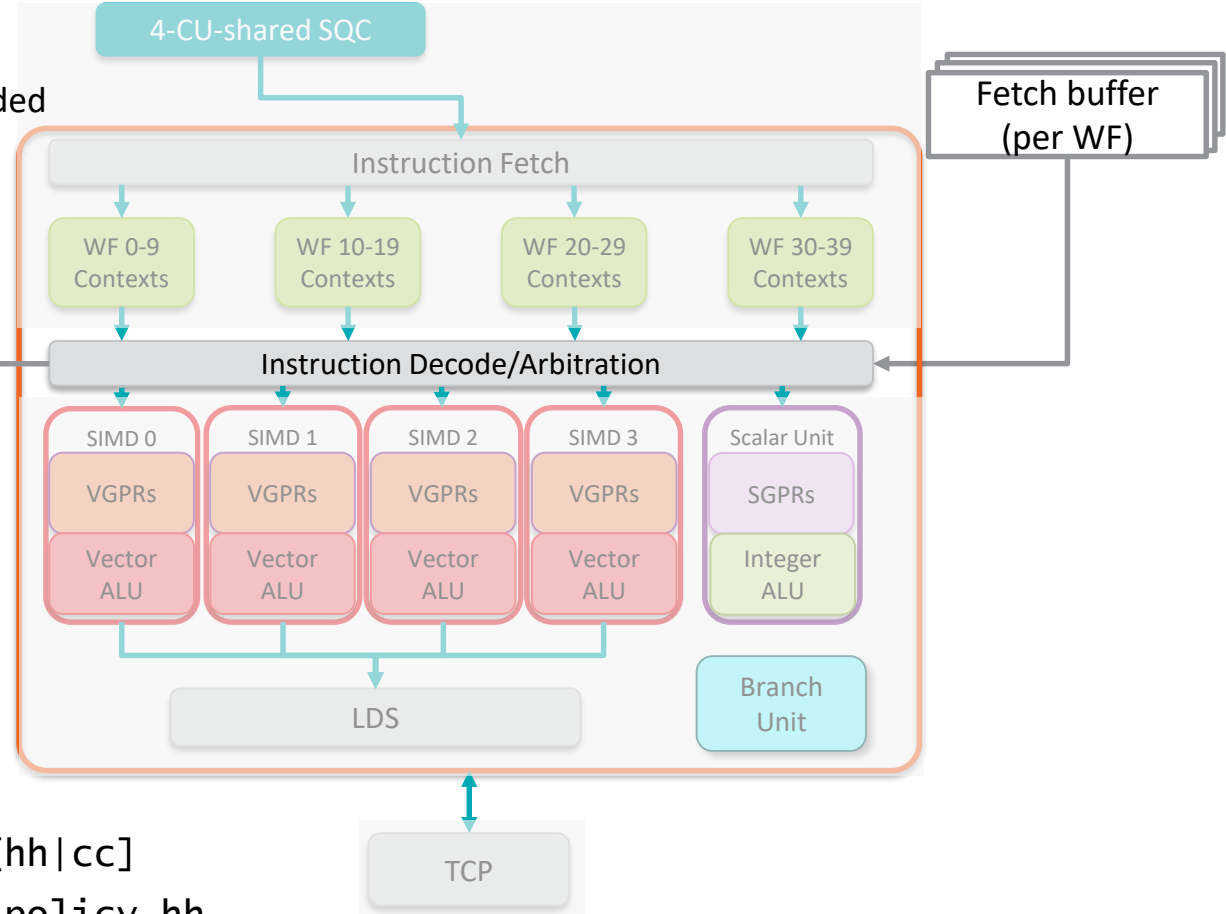
- ▲ SQC shared by 4 CUs
 - Number of SQCs and CUs are configurable
- ▲ Fetch
 - Shared and arbitrated between SIMDs in a CU
 - Fetch to each SIMD unit
 - Buffers fetched cache lines per WF
- ▲ WF Contexts
 - 10 WFs per SIMD, 40 per CU
 - PC and decoded instruction buffers (IB)
 - Register file and LDS allocation
 - ¼ of WF executes each cycle
 - 4 cycles needed to fully execute single SIMD instruction



DECODE AND ISSUE



Instruction data decoded into instruction buffer



- ▲ Instructions are decoded out of fetch buffers
- ▲ Instruction arbitration
 - Can issue to each functional unit each cycle
 - Finds ready WFs
 - Scheduling policy dictates which WFs have priority
 - Oldest first, easy to add others

scheduler.[hh|cc]
scheduling_policy.hh

REGISTER FILES

vector_register_file.[hh|cc]
scalar_register_file.[hh|cc]
simple_pool_manager.[hh|cc]
static_register_manager_policy.[hh|cc]

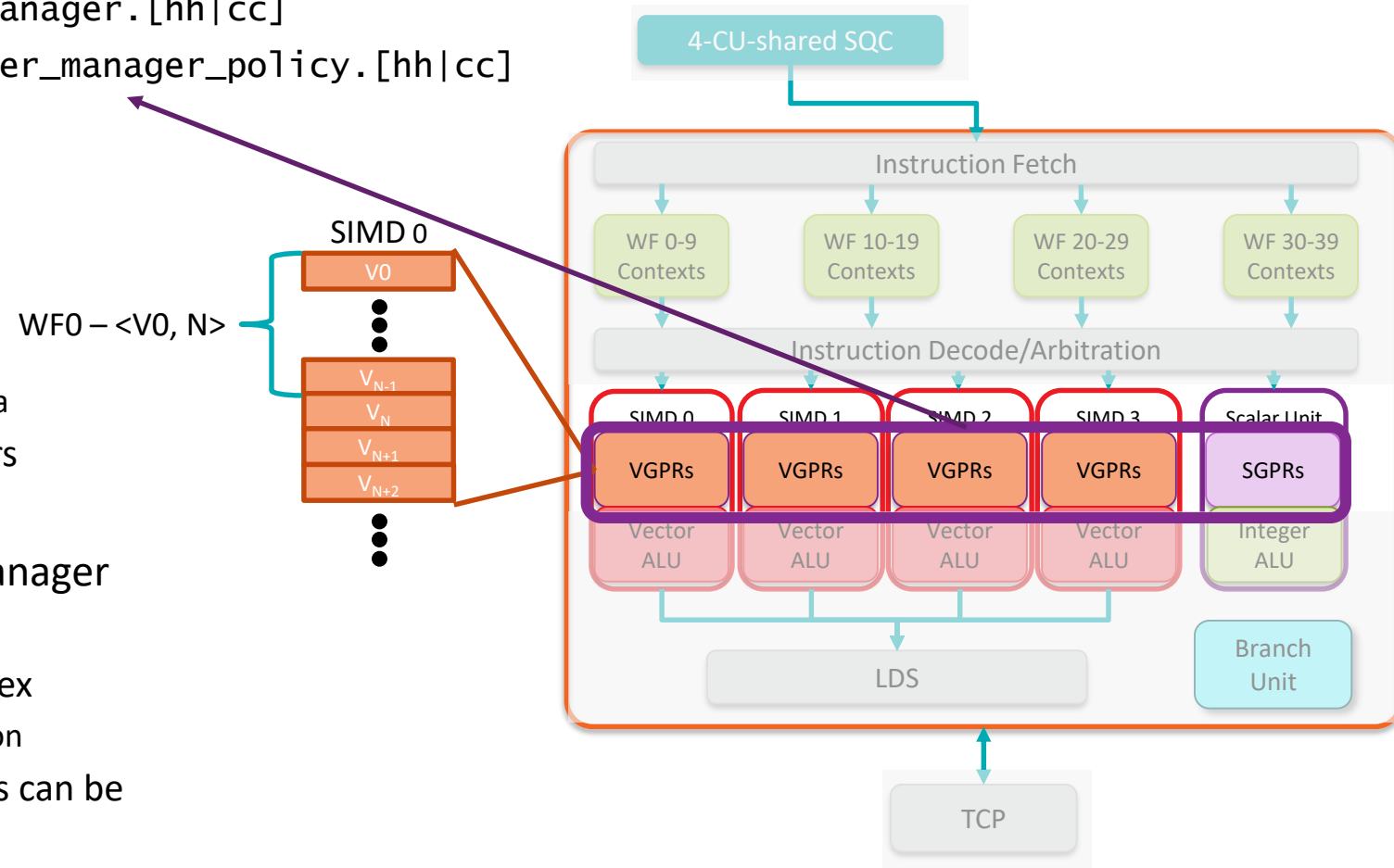


General Purpose Registers (GPRs)

- Vector registers (VGPR) partitioned per SIMD
 - Configurable size
 - Because each SIMD executes independent WF
- 32-bit wide
 - Combine adjacent VGPRs for 64-bit or 128-bit data
- Each WF also has scalar general purpose registers (SGPRs)

Register Allocation Done by a Simple Pool Manager

- *Only allows one WG at a time*
- Statically mapped virtual → physical register index
 - $\langle \text{base}, \text{limit} \rangle$ pair of registers specify GPR allocation
- Modular design – more advanced pool managers can be swapped into the VRF seamlessly
- Simple timing model with constant delay

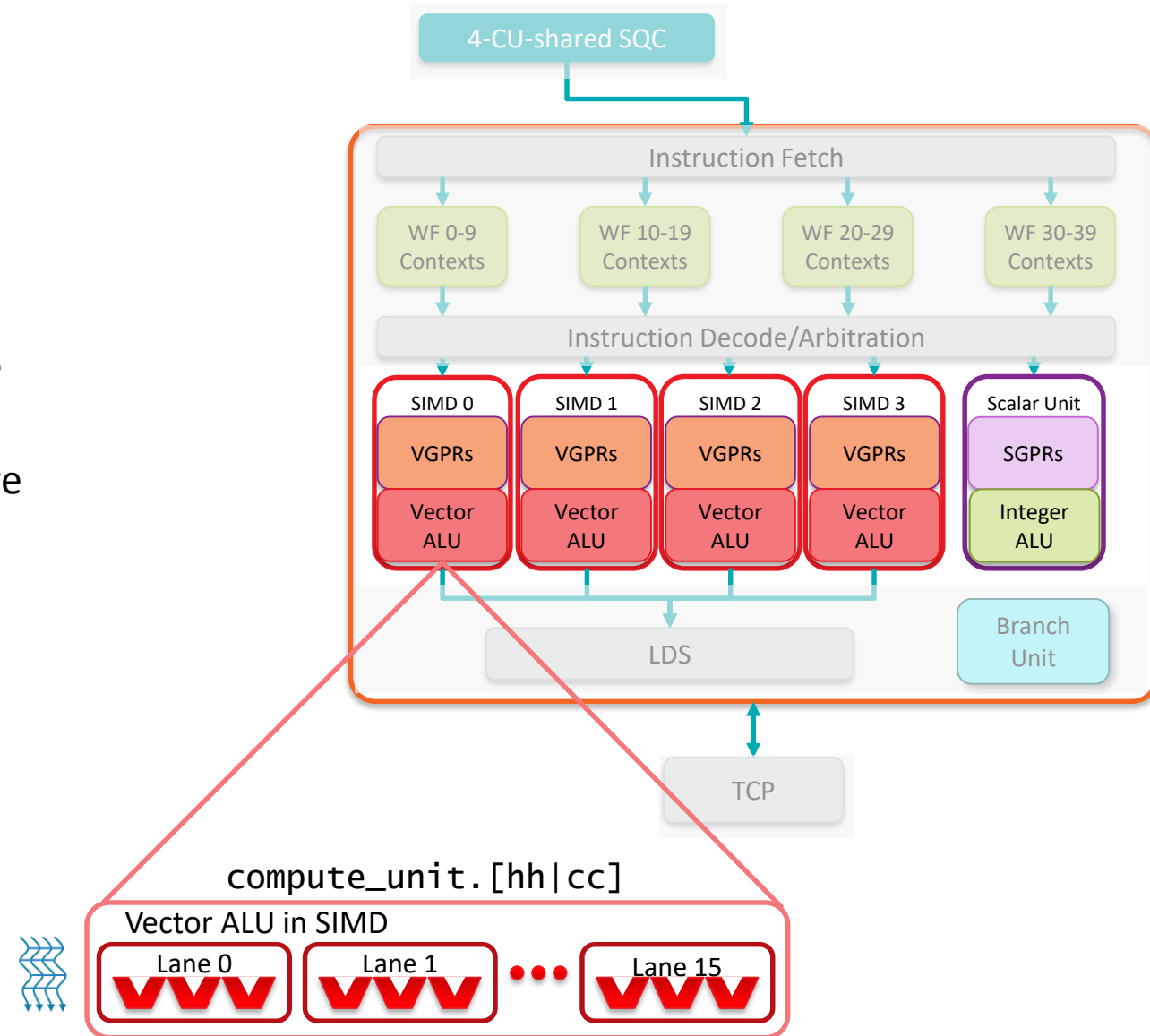


This is an area where gem5 user contributions would be extremely valuable.

VECTOR ALUs



- ▲ 16-lane vector pipeline per SIMD
 - Each lane has a set of functional units
 - One work-item per lane
- ▲ 4 cycles to execute a WF for all 64 work-items
 - In gem5, 64 work-items are executed in one tick and ticks are multiplied by 4
- ▲ SIMD execution may take longer if work-items in WF have dissimilar behaviors
 - Example 1: Branch (or spatial) divergence
 - Branches executed through predication
 - When control flow diverges, all lanes take all paths
 - Example 2: Memory (or temporal) divergence
 - Longer access latency by one work-item stalls entire WF



GPU CORE TIMING

HANDLING MEMORY INSTRUCTIONS

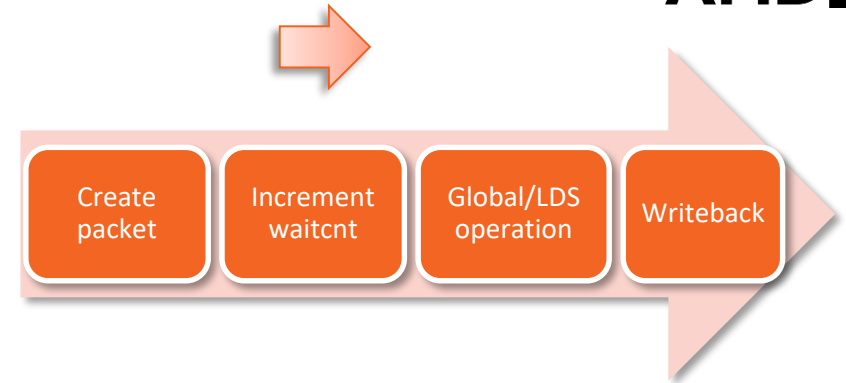
```
gpu_dyn_inst.[hh|cc]  
arch/gcn3/instructions.[hh|cc]  
arch/gcn3/insts/op_encodings[hh|cc]
```

- ▲ Memory instructions generate memory requests
 - Part of GPU instruction definition (ISA-specific)
- ▲ Three phases
 - *execute()*
 - Read operands, calculate address, increment wait count, and issue to appropriate memory pipe
 - *initiateAcc()*
 - Issue request to memory system
 - *completeAcc()*
 - For loads write back data. Stores do nothing.
- ▲ New machine ISAs can use this capability to support their own memory instructions
- ▲ Individual stages contribute to the memory instruction timing
 - Additionally memory end timing handled by Ruby and memory technology parameters
- ▲ Memory dependencies are preserved using ***waitcnts***

Example GCN3 code:

```
flat_load_dword v4, v[4:5] ①  
flat_load_dword v16, v[8:9]  
flat_load_dword v23, v[14:15]  
flat_load_dword v10, v[10:11]  
s_waitcnt          cnt(3) ②  
v_ashrrev          v5, 31, v4 ③
```

GPU dynamic memory instruction



```
local_memory_pipeline.[hh|cc]  
global_memory_pipeline.[hh|cc]  
scalar_memory_pipeline.[hh|cc]
```

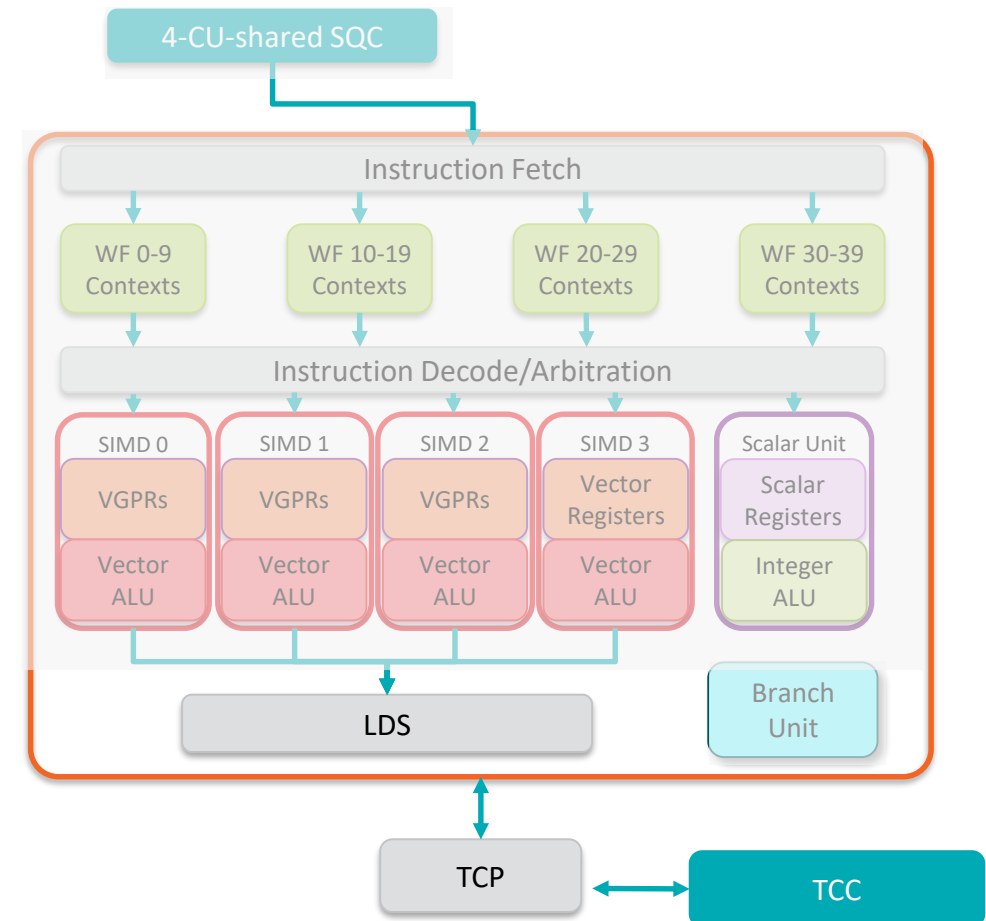
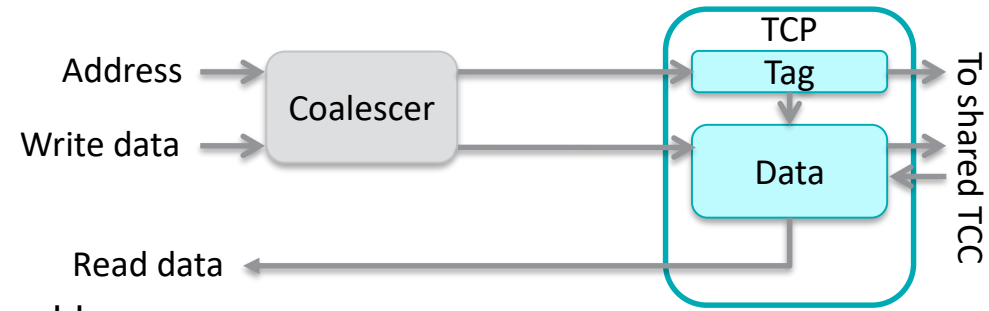
1. flat_load writes v4
2. Waitcnt specifies wait count value must be ≤ 3
3. Arithmetic shift later reads from v4
4. Waitcnt waits until at least #1 is finished

VECTOR MEMORY EXECUTION



- ▲ In gem5:
 - Address calculation: `arch/gcn3/insts/op_encodings.hh`
 - Address coalescing
 - `mem/ruby/system/GPUCoalescer.hh|cc`
 - `mem/ruby/system/VIPERCoalescer.hh|cc`
 - `mem/protocol/GPU_VIPER-TCP.sm`
 - `mem/protocol/GPU_VIPER-TCC.sm`

- ▲ LDS
 - User-managed address space
 - Scratchpad for WFs in workgroup
 - Used for data sharing and synchronization within workgroup
 - Cleared when workgroup completes
 - In gem5, functional model with a pointer per workgroup



CONTROL FLOW DIVERGENCE

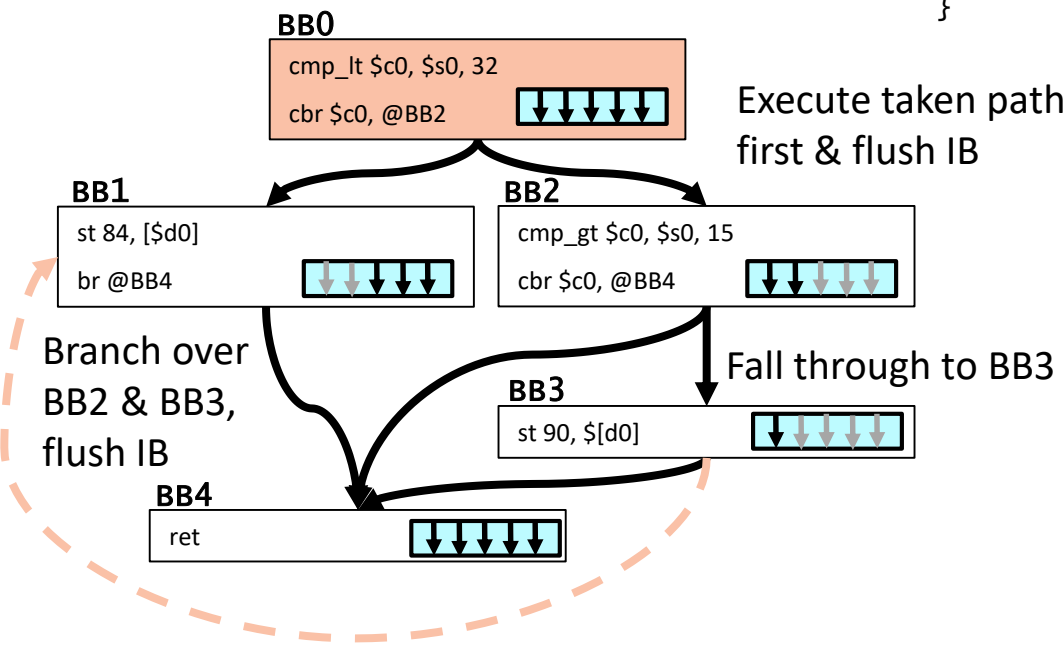
SIMT VS. VECTOR EXECUTION MODEL

Source code:

```
if (i > 31) {  
    *x = 84;  
} else if (i < 16) {  
    *x = 90;  
}
```



HSAIL



Instruction buffer



Reconvergence point reached, HW initiated jump to divergent path

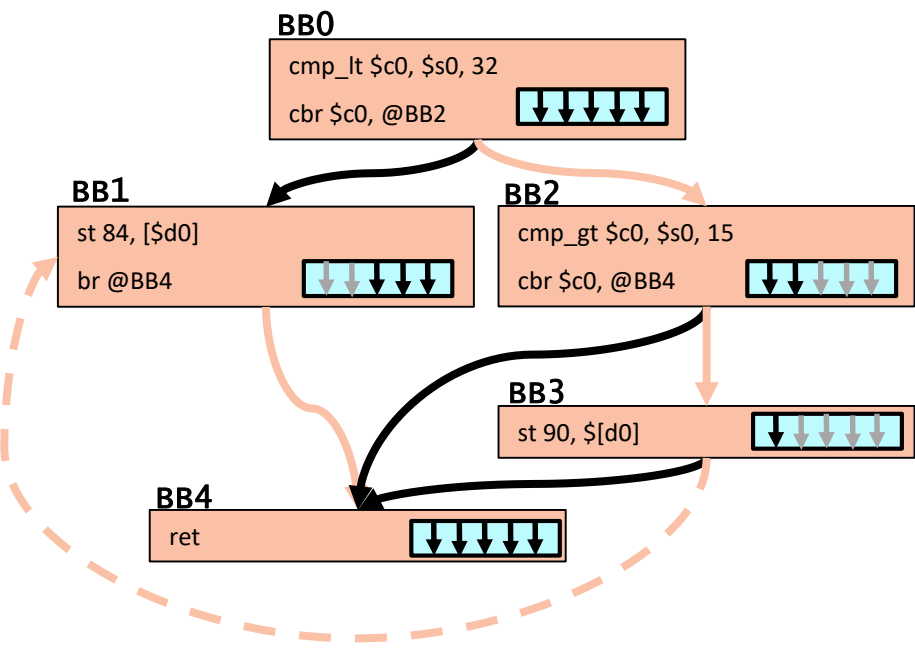
CONTROL FLOW DIVERGENCE

SIMT VS. VECTOR EXECUTION MODEL

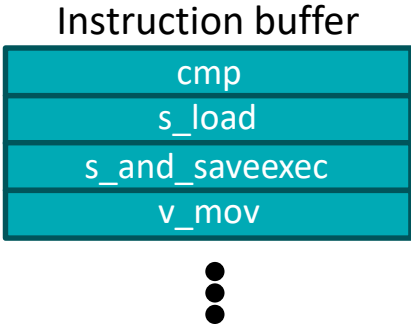
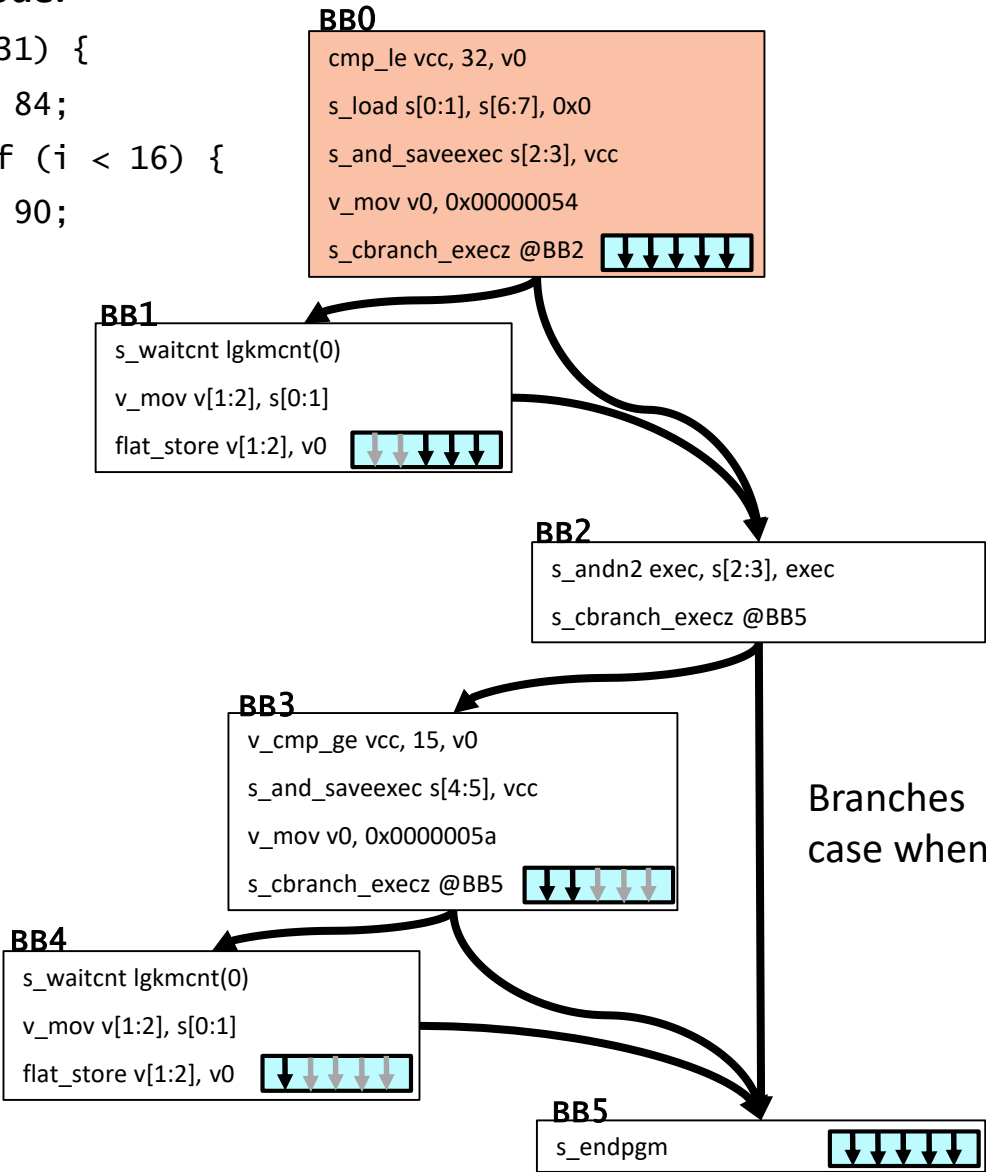
Source code:

```
if (i > 31) {  
    *x = 84;  
} else if (i < 16) {  
    *x = 90;  
}
```

HSAIL



GCN3



Branches are optimizations for case when EXEC = 0 for a BB

OUTLINE



Topic	Presenter	Time
Background	Tony	8:00 – 8:15
ROCm, GCN3 ISA, and GPU Arch	Tony	8:15 – 9:15
HSA Implementation in gem5	Sooraj	9:15 – 10:00
Break		10:00 – 10:30
Ruby and GPU Protocol Tester	Tuan	10:30 – 11:15
Demo and Workloads	Matt	11:15 – 11:50
Summary and Questions	All	11:50 – 12:00

WHAT IS HSA?



Heterogeneous System Architecture

Processor design that makes it easy to harness the entire computing power of GPUs for faster and more power-efficient devices, including personal computers, tablets, smartphones, and servers



Bringing GPU performance to a wide variety of applications

KEY FEATURES OF HSA



hUMA

Heterogeneous Unified Memory Architecture

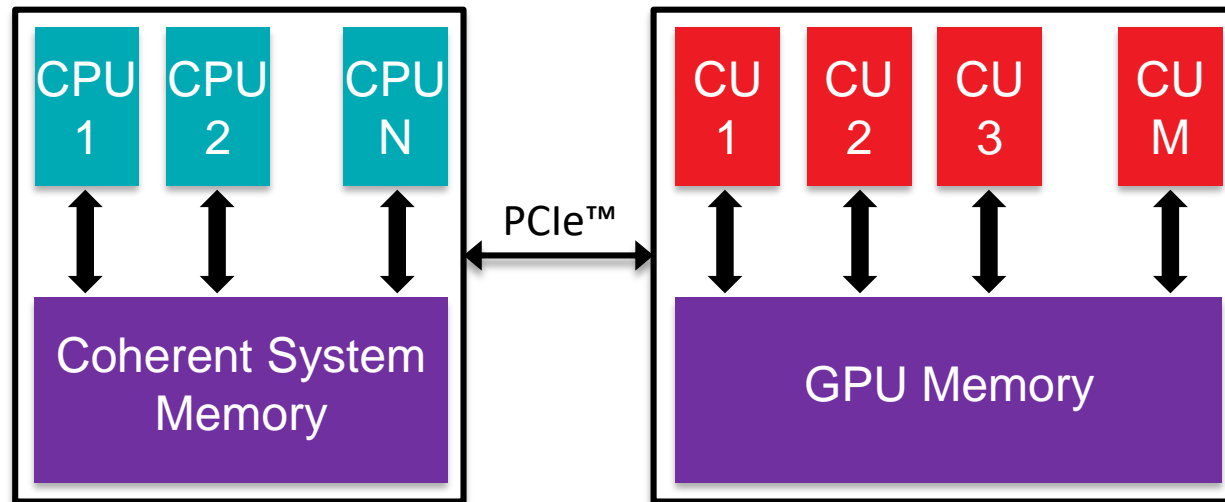
hComm

Heterogeneous Communication via Signals and Atomics

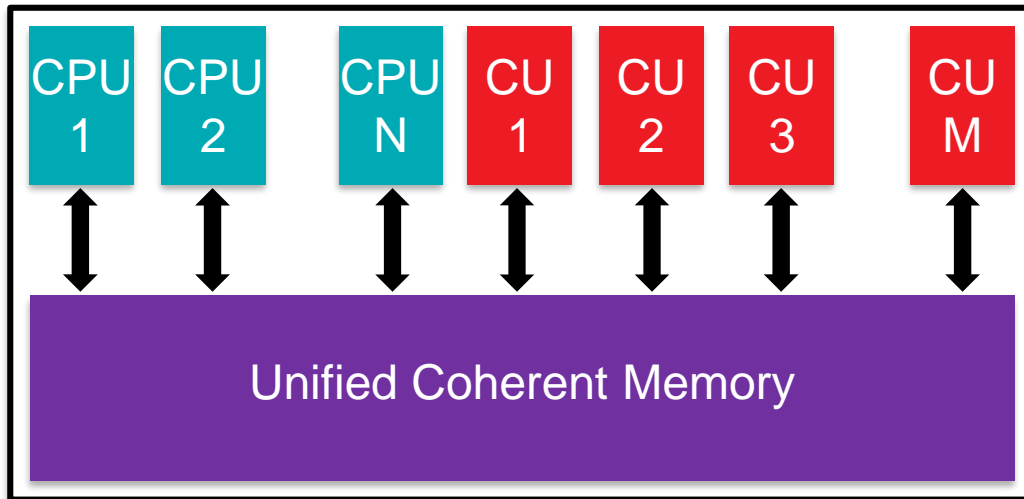
hQ

Heterogeneous Queuing

TRADITIONAL DISCRETE GPU

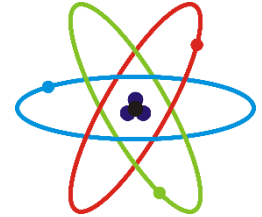


- ▲ Separate memory
- ▲ Separate addr space
 - No pointer-based data structures
- ▲ Explicit data copying
 - High latency
 - Low bandwidth
- ▲ Need lots of compute on GPU to amortize copy overhead
- ▲ Very limited GPU memory capacity



- ▲ Unified address space
 - GPU uses user virtual addresses
 - Fully coherent
- ▲ No explicit copying
 - Data movement on demand
- ▲ Pointer-based data structures shared across CPU & GPU
- ▲ Pageable virtual addresses
 - No GPU capacity constraints

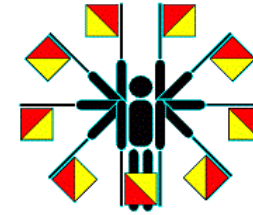
- ▲ Atomic memory updates fundamental to efficient thread synchronization
 - Implement primitives like mutexes, semaphores, histograms, ..., previously only implemented on CPU
- ▲ HSA supports 32bit or 64bit values for atomic ops
 - CAS, SWAP, add, increment, sub, decrement, ...
 - and other common arithmetic and logic atomic ops
- ▲ On PCI-Express system, atomics map to PCI-E atomics
- ▲ HSA specifies a well-defined “SC for HRF” memory model
 - A variant of “Release Consistency” model
 - Acquire: pull latest data (to me)
 - Release: push latest data (to others)
 - Compatible with C++11, Java, OpenCL, and .NET memory models
 - Details: “HSA Platform System Arch Specification”, <http://hsafoundation.com>



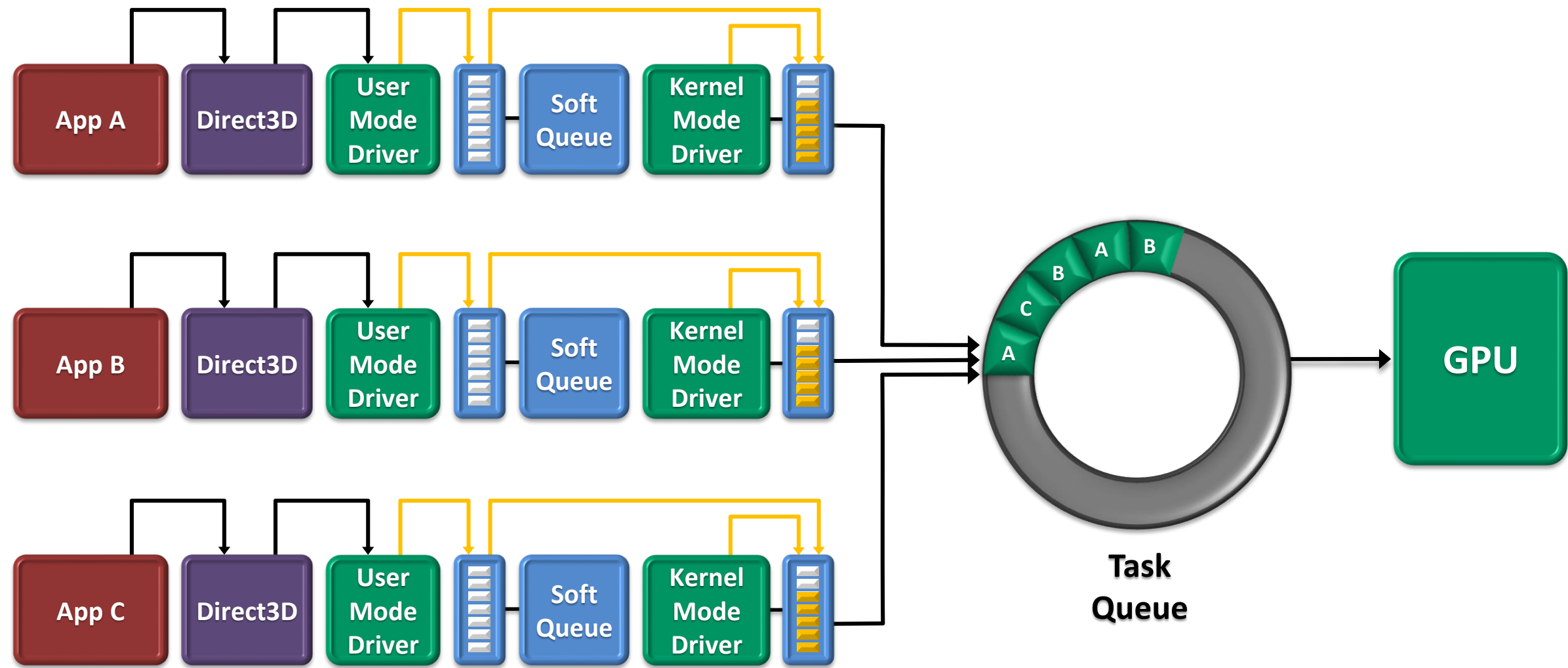
THE HSA SIGNALS INFRASTRUCTURE



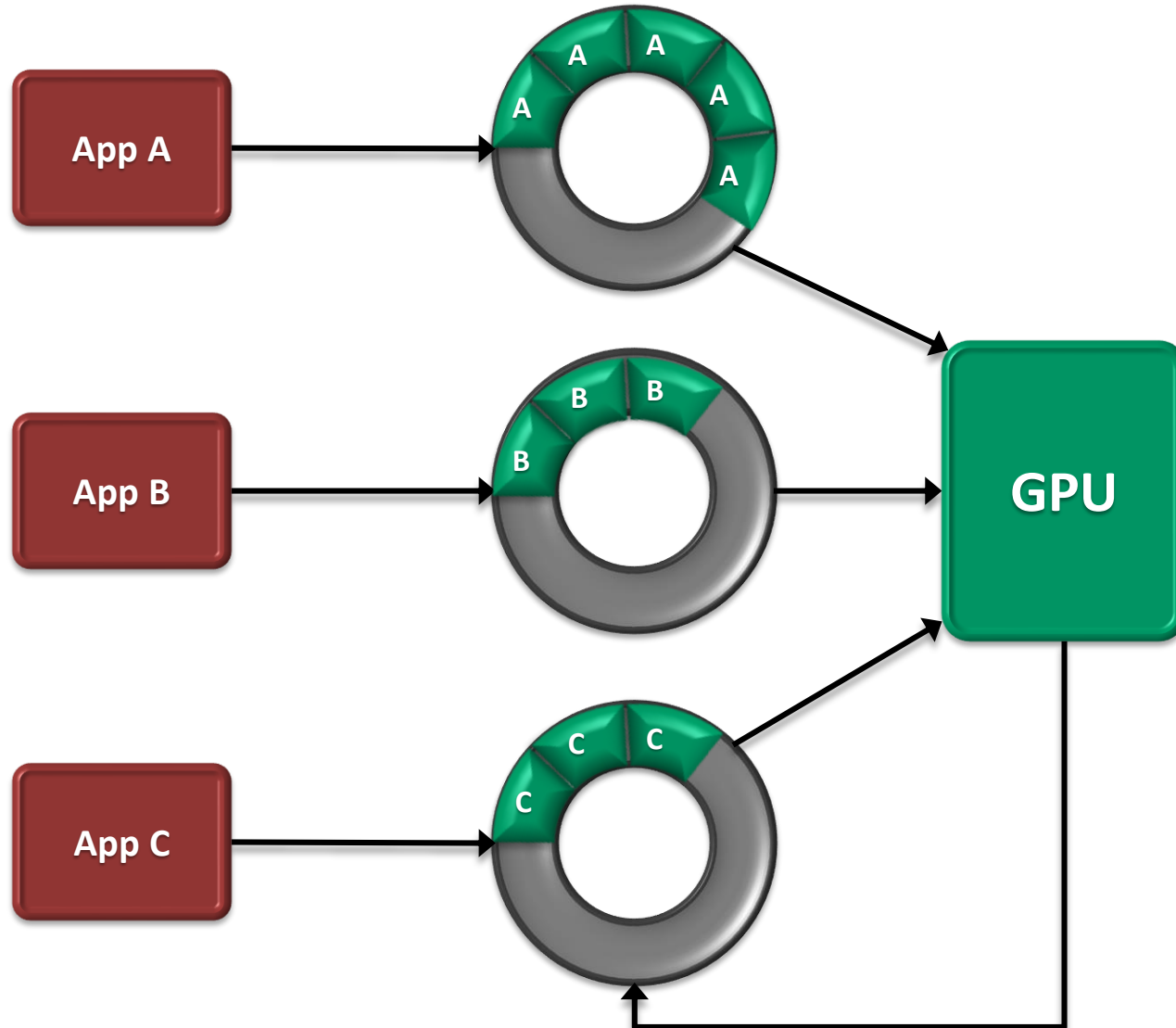
- ▲ Hardware-assisted signaling and synchronization primitives
 - Memory semantics, equivalent to atomics
 - e.g., 32bit or 64bit value, content updated atomically
 - Threads can wait on a value
 - Power-efficient synchronization between CPU and GPU threads
- ▲ Allows one-to-one, many-to-one, and one-to-many signaling
 - Used by system software, runtime, and application SW
 - Infrastructure to build higher-level synchronization primitives like mutexes, semaphores, etc.
- ▲ Updating the value of a signal is equivalent to sending the signal
 - Release semantics: push data to others
- ▲ Waiting on a signal is also permitted
 - Via a wait instruction or via a runtime API
 - Acquire semantics: pull data from others



TRADITIONAL COMMAND AND DISPATCH FLOW



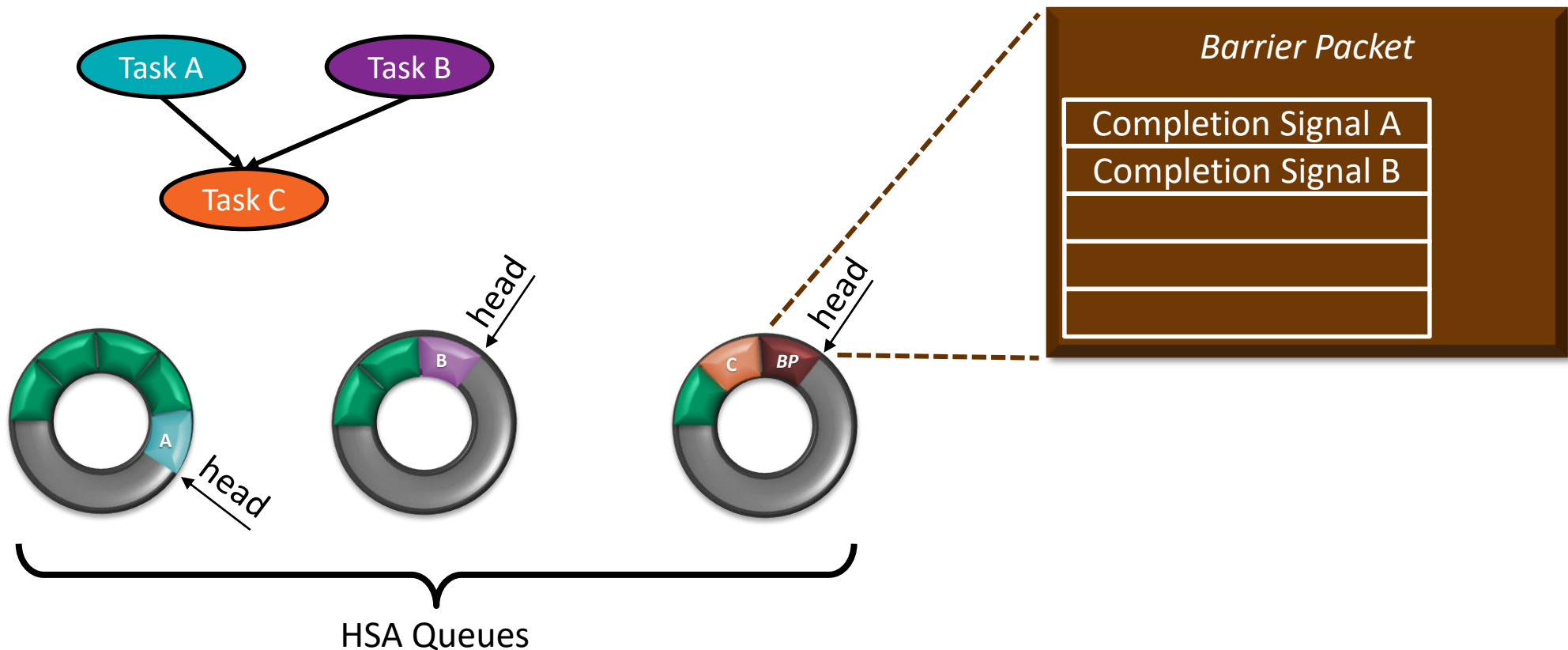
hQ COMMAND AND DISPATCH FLOW



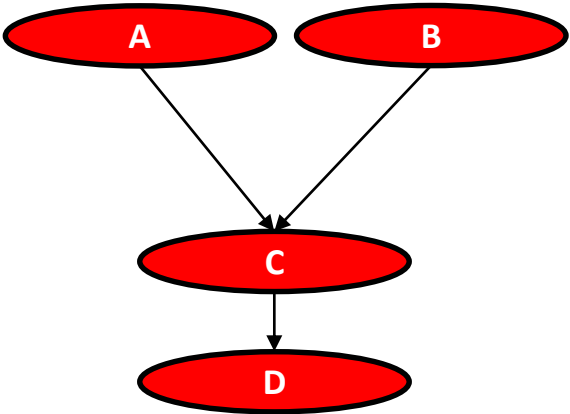
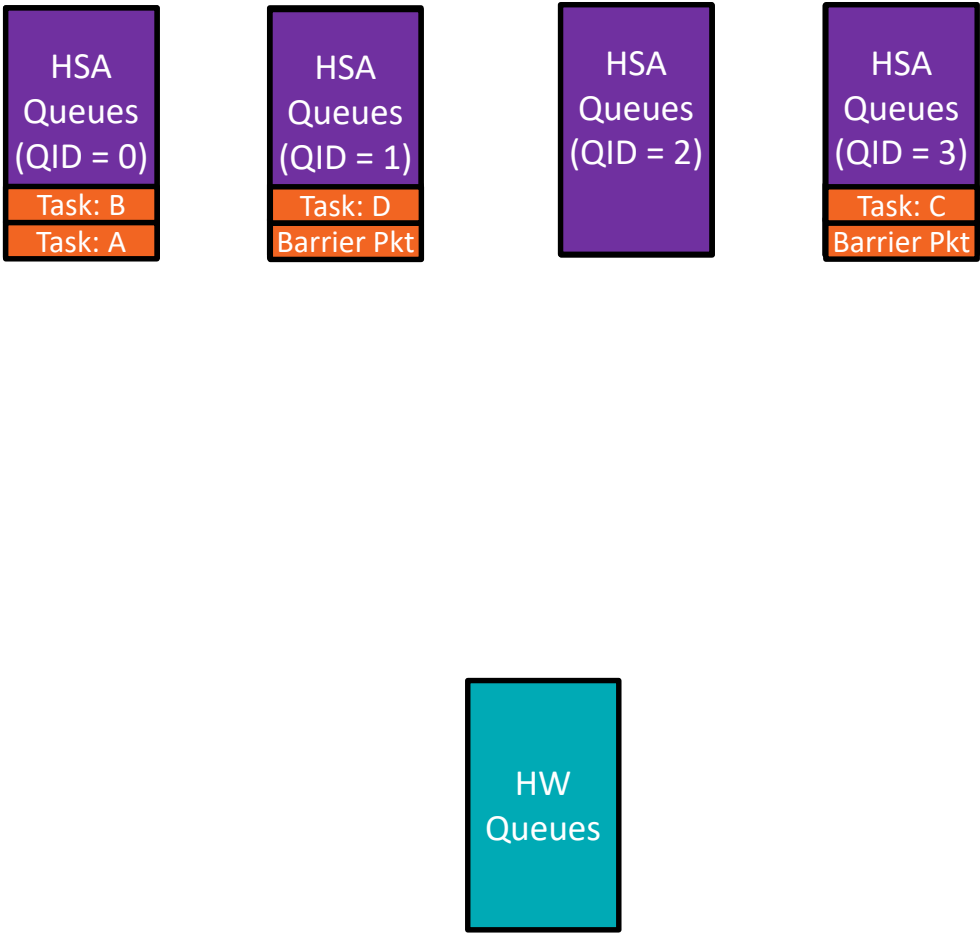
- ▲ User-mode application talks directly to the hardware
 - HSA Architected Queuing Language (AQL) defines vendor-independent format
 - No system call
 - No kernel driver involvement
- ▲ Hardware scheduling
- ▲ Greatly reduced dispatch overhead
 - less overhead to amortize
 - profitable to offload smaller tasks
- ▲ Device enqueue: GPU kernels can self-enqueue additional tasks (dynamic parallelism)

NATIVE SUPPORT FOR DATA-DEPENDENT TASKS

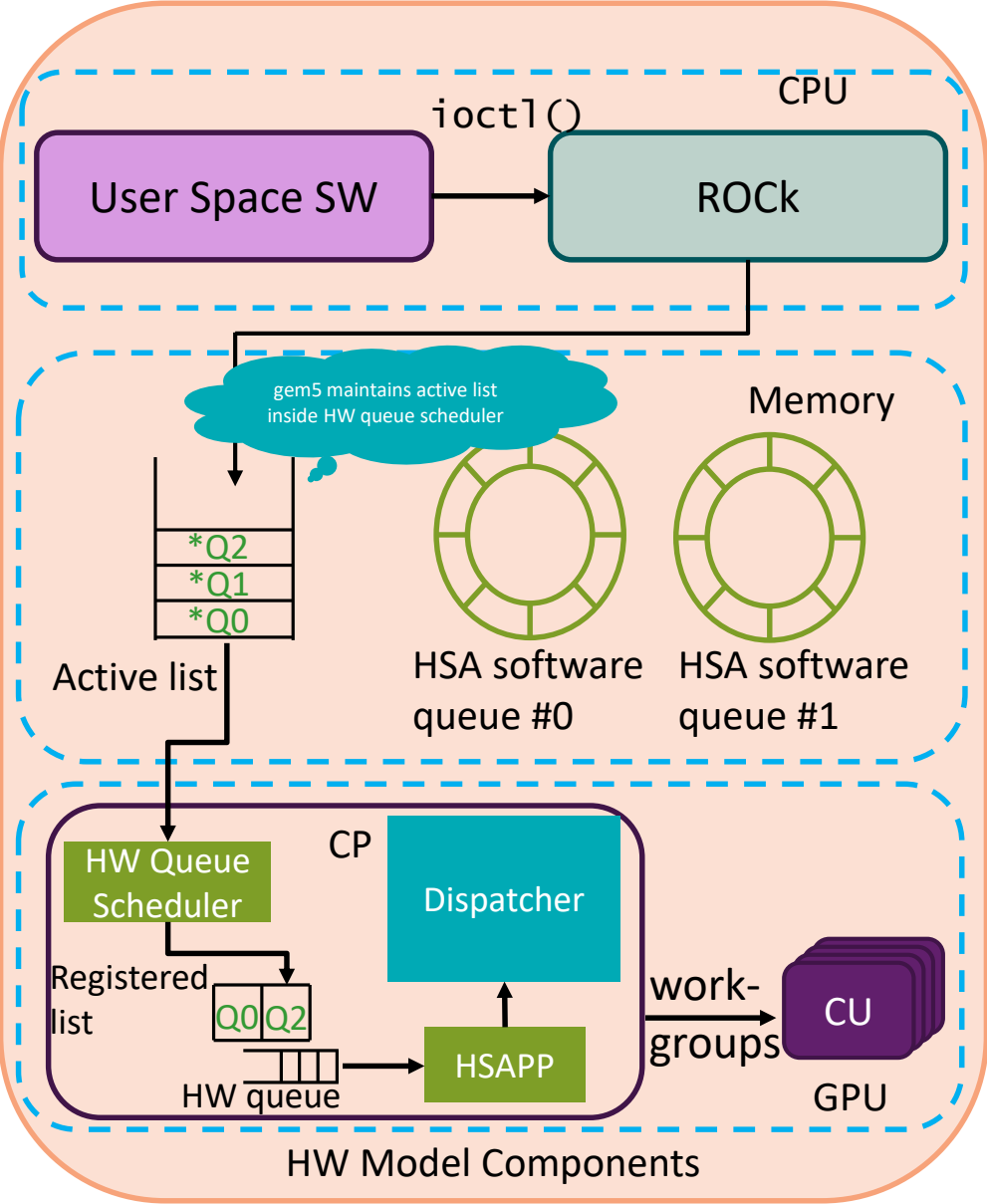
EXPOSE DIRECTED ACYCLIC GRAPHS (DAG) TO HARDWARE



GPU hardware is responsible for managing and scheduling tasks (No Operating System!)



QUEUE SCHEDULING HARDWARE



Component	Source file
driver	gpu-compute/gpu_compute_driver.cc dev/hsa/kfd_ioctl.h
hardware scheduler	dev/hsa/hw_scheduler.[hh cc]
packet processor	dev/hsa/hsa_packet_processor.[hh cc]
dispatcher	gpu-compute/dispatcher.[hh cc]

DOORBELL PAGES AND EVENT PAGES



1. `HSADriver::mmap(mmap args)`

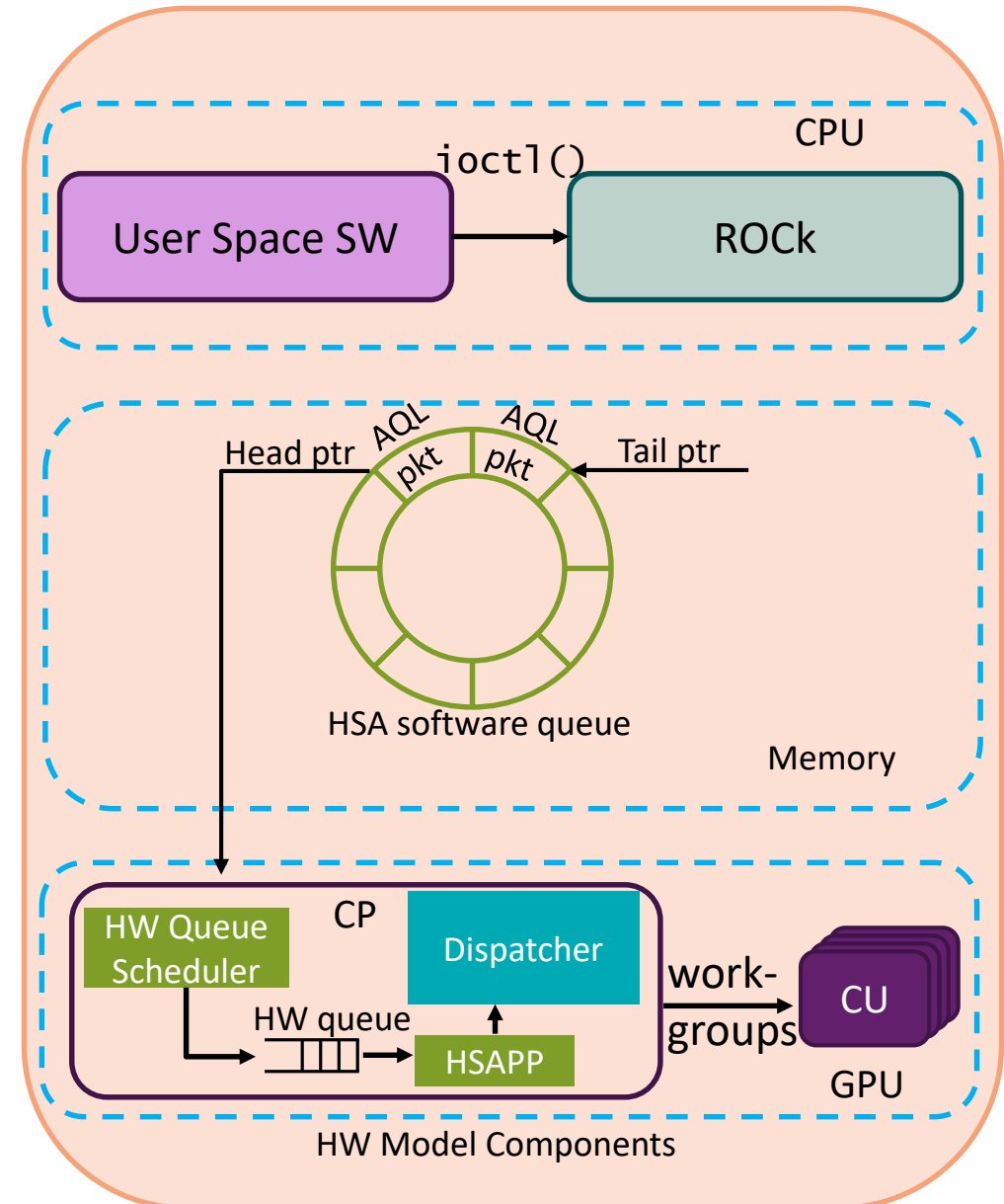
- runtime calls mmap on the driver
- mmap offset distinguishes event page vs doorbell page mmap

2. driver allocates doorbell page and returns the page address

- For doorbells -> Driver maps the page address to PIO address in the PT
- For events -> Event pages are user pages

3. runtime allocates doorbell address for each queue based on *queue ID*.

- Event pages are currently unused; model relies on functional interfaces for event notification.



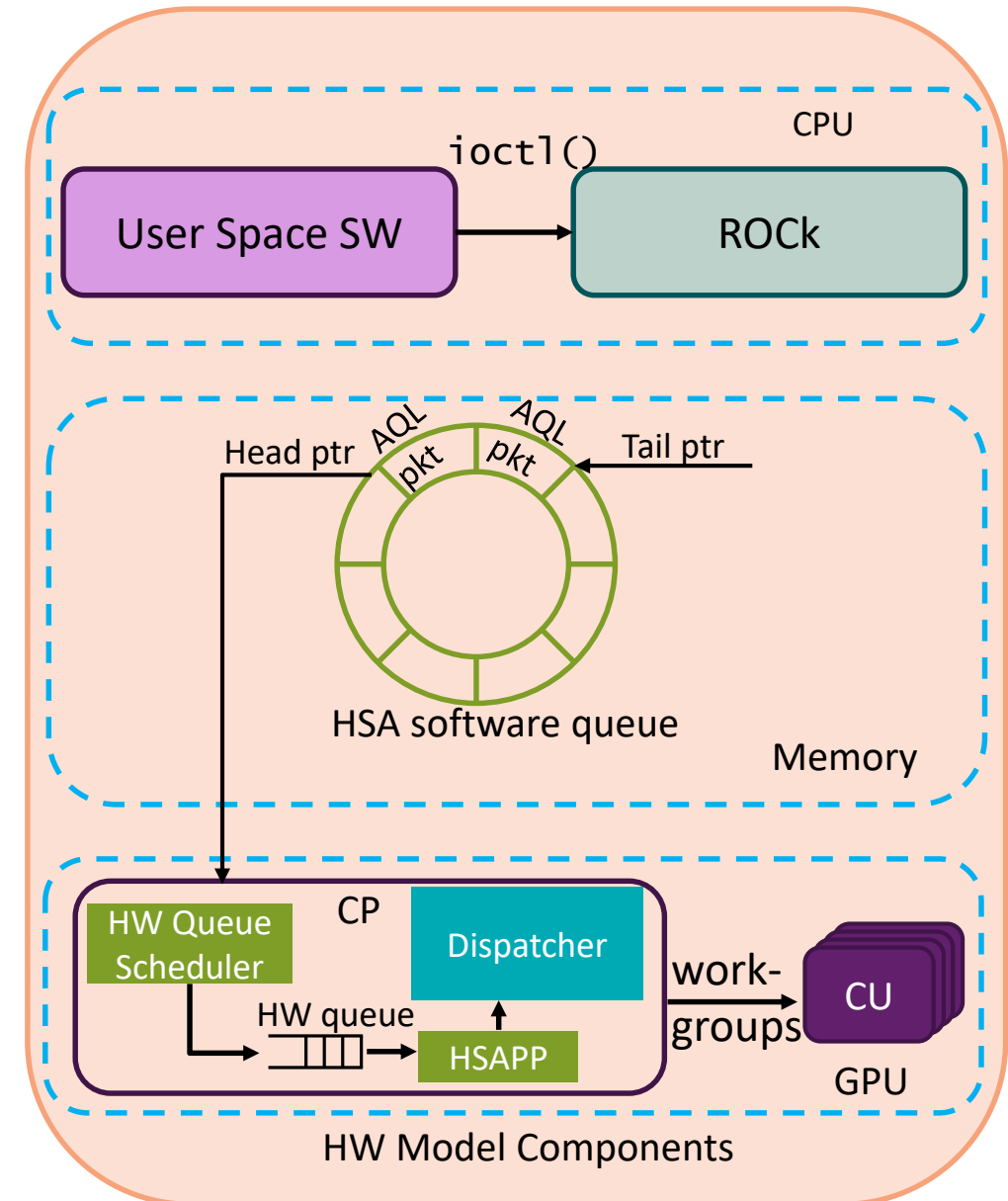
HSA QUEUE CREATION



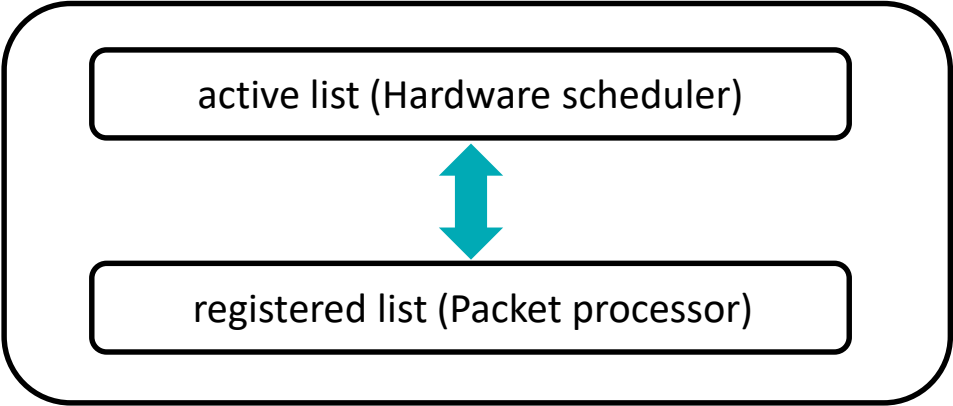
1. `hsa_queue_create(some args)`

2. `GPUComputeDriver::ioctl(tc, AMDKFD_IOC_CREATE_QUEUE)`

3. `HWScheduler::registerNewQueue(some args)`

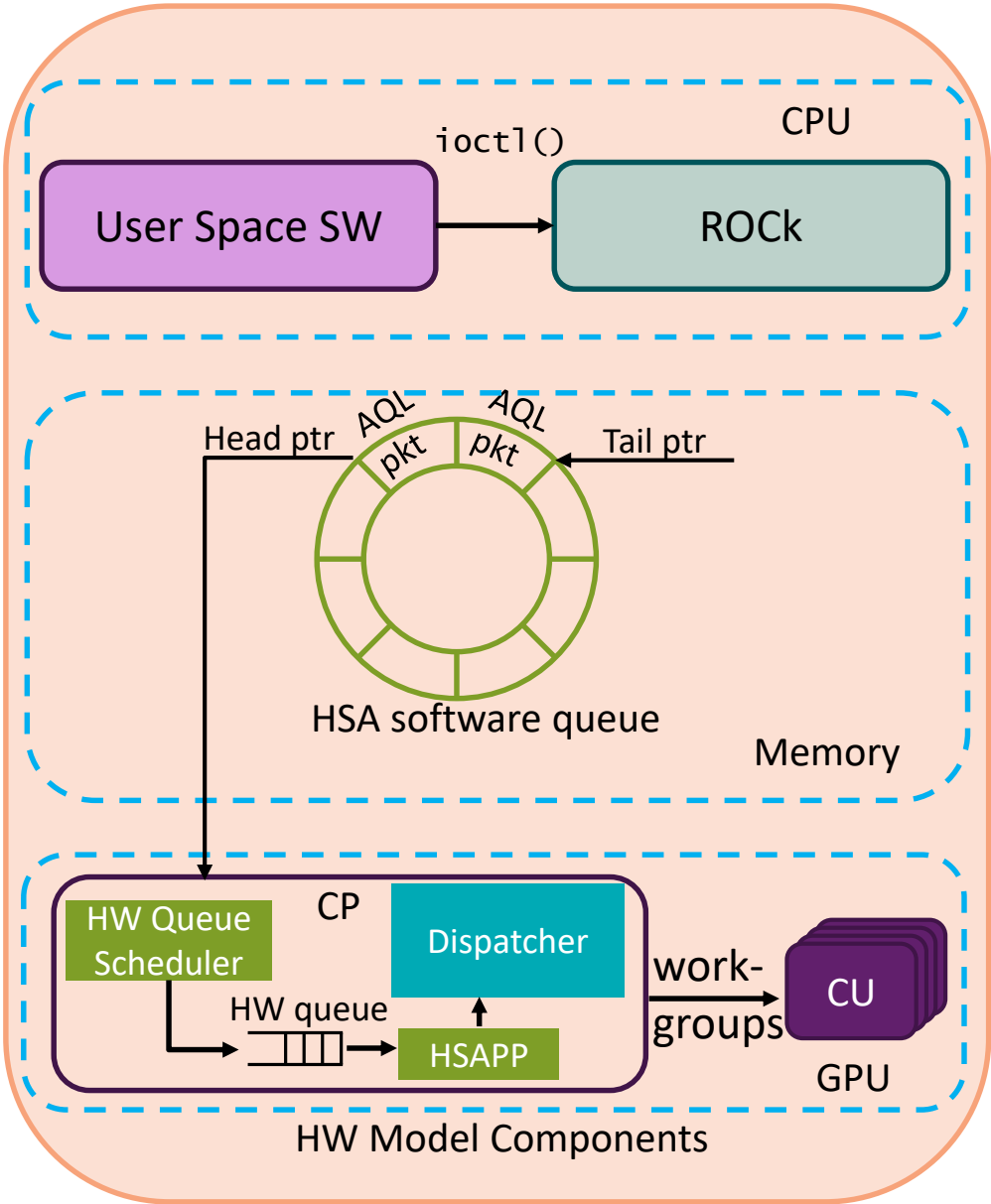


QUEUE SWAPPING LOGIC



```
// wakeup every wakeupDelay ticks
1. HWScheduler::wakeup()

2. HWScheduler::contextSwitchQ ()
```



HSA SIGNAL CREATION



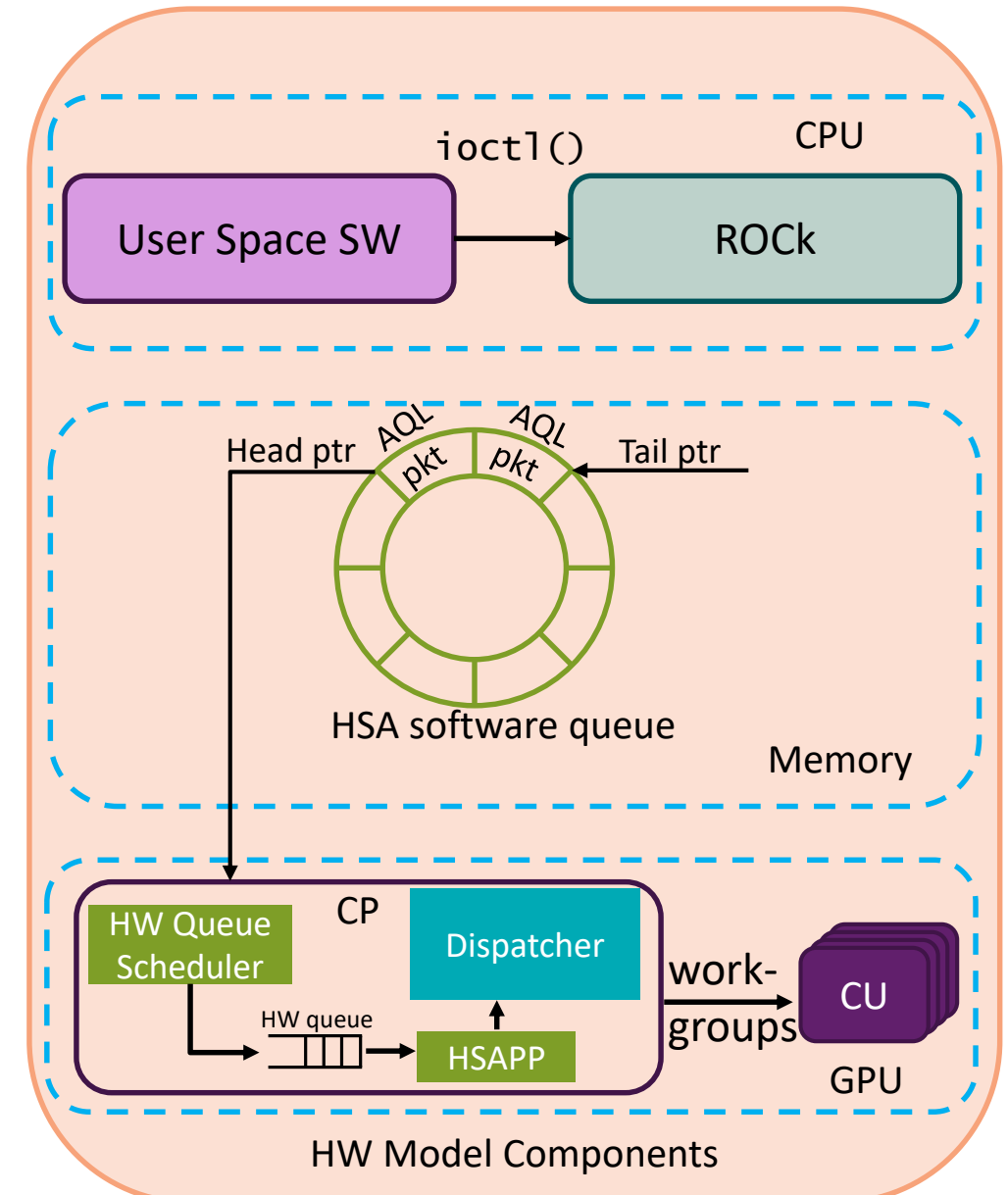
1. `hsa_signal_create(some args)`

2. `GPUComputeDriver::ioctl(tc, AMDKFD_IOC_CREATE_EVENT)`

3. `hsa_signal_wait_scacquire(some args)`

4. `GPUComputeDriver::ioctl(tc, AMDKFD_IOC_WAIT_EVENTS)`

5. `GPUComputeDriver::ioctl(tc, AMDKFD_IOC_SET_EVENT) //or`
5. `signalWakeupEvent(event ID)`



HSA DOORBELLS

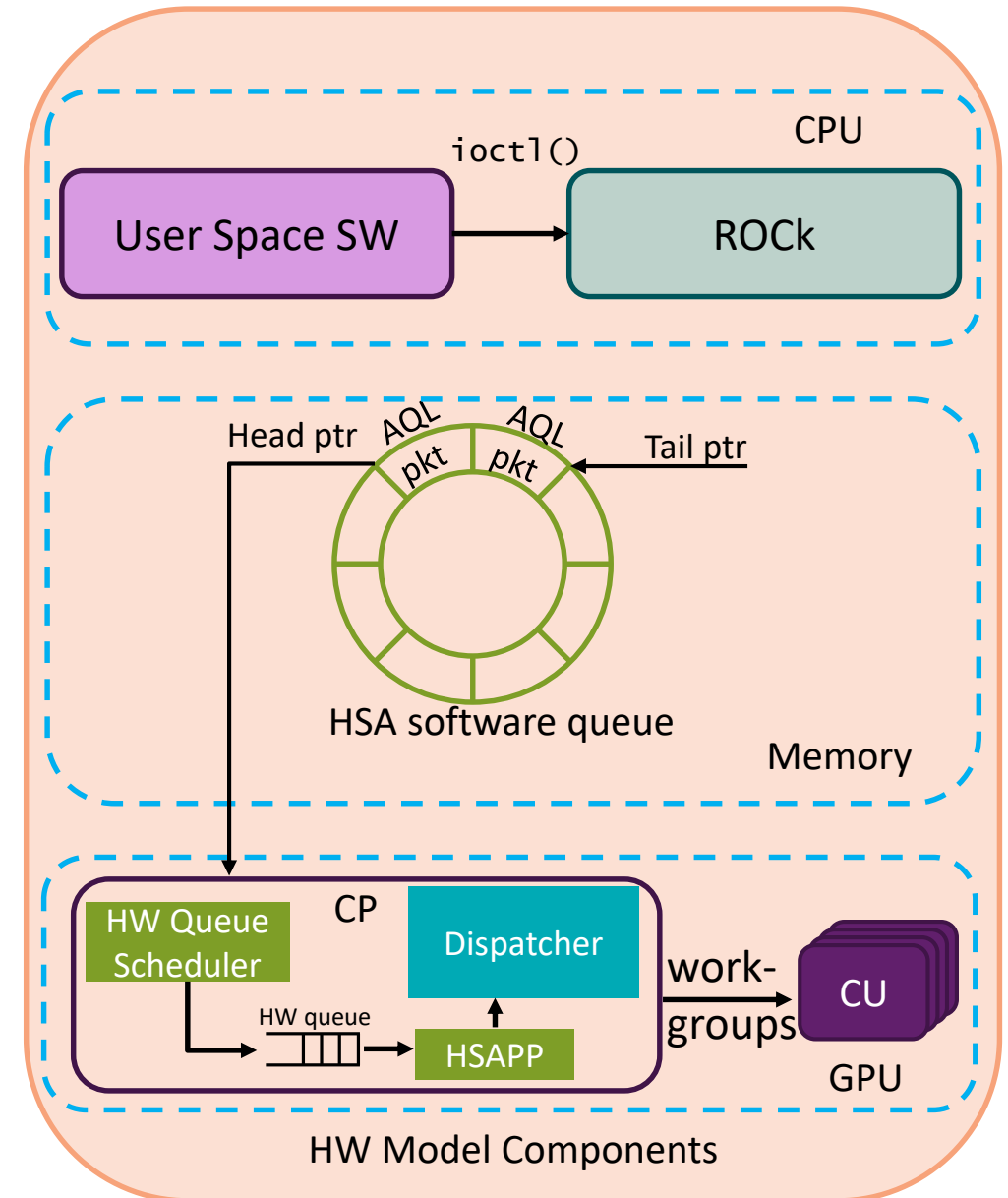


```
1. hsa_signal_store_screlease(queue->doorbell_signal, packet_id);
```

```
2. HSAPacketProcessor::write(some args)
```

```
3. HWScheduler::write(Addr db_addr, uint32_t doorbell_reg)
```

```
4. HSAPacketProcessor::getCommandsFromHost(some args)
```



OUTLINE



Topic	Presenter	Time
Background	Tony	8:00 – 8:15
ROCm, GCN3 ISA, and GPU Arch	Tony	8:15 – 9:15
HSA Implementation in gem5	Sooraj	9:15 – 10:00
Break		10:00 – 10:30
Ruby and GPU Protocol Tester	Tuan	10:30 – 11:15
Demo and Workloads	Matt	11:15 – 11:50
Summary and Questions	All	11:50 – 12:00

OUTLINE



Topic	Presenter	Time
Background	Tony	8:00 – 8:15
ROCm, GCN3 ISA, and GPU Arch	Tony	8:15 – 9:15
HSA Implementation in gem5	Sooraj	9:15 – 10:00
Break		10:00 – 10:30
Ruby and GPU Protocol Tester	Tuan	10:30 – 11:15
Demo and Workloads	Matt	11:15 – 11:50
Summary and Questions	All	11:50 – 12:00

RUBY MEMORY CONTRIBUTIONS

OUTLINE



- ▲ Ruby Background
- ▲ CU – Memory Interface
- ▲ GPU VIPER Protocol
- ▲ GPU SLICC Protocol Tester

RUBY MEMORY CONTRIBUTIONS

OUTLINE



- ▲ Ruby Background
- ▲ CU – Memory Interface
- ▲ GPU VIPER Protocol
- ▲ GPU SLICC Protocol Tester

▲ Flexible Memory System

- Rich configuration
 - Simulate combination of caches, coherence, interconnect, etc.
- Rapid prototyping
 - Domain-Specific Language (SLICC) for coherence protocols
 - Modular components

▲ Detailed statistics

- Latency distributions for requests
- Generated state transitions, network utilization, etc.

▲ Detailed component simulation

- Network (fixed/flexible Garnet pipelines and simple)
- Caches (pluggable replacement policies)
- Memory (shared memory controllers between Classic and Ruby)

▲ gem5 Ruby tutorial: <http://learning.gem5.org/book/part3/index.html>

▲ Our GCN3 GPU model only works with Ruby memory

RUBY MEMORY CONTRIBUTIONS

OUTLINE



- ▲ Ruby Background
- ▲ CU – Memory Interface
- ▲ GPU VIPER Protocol
- ▲ GPU SLICC Protocol Tester

CU - MEMORY INTERFACE

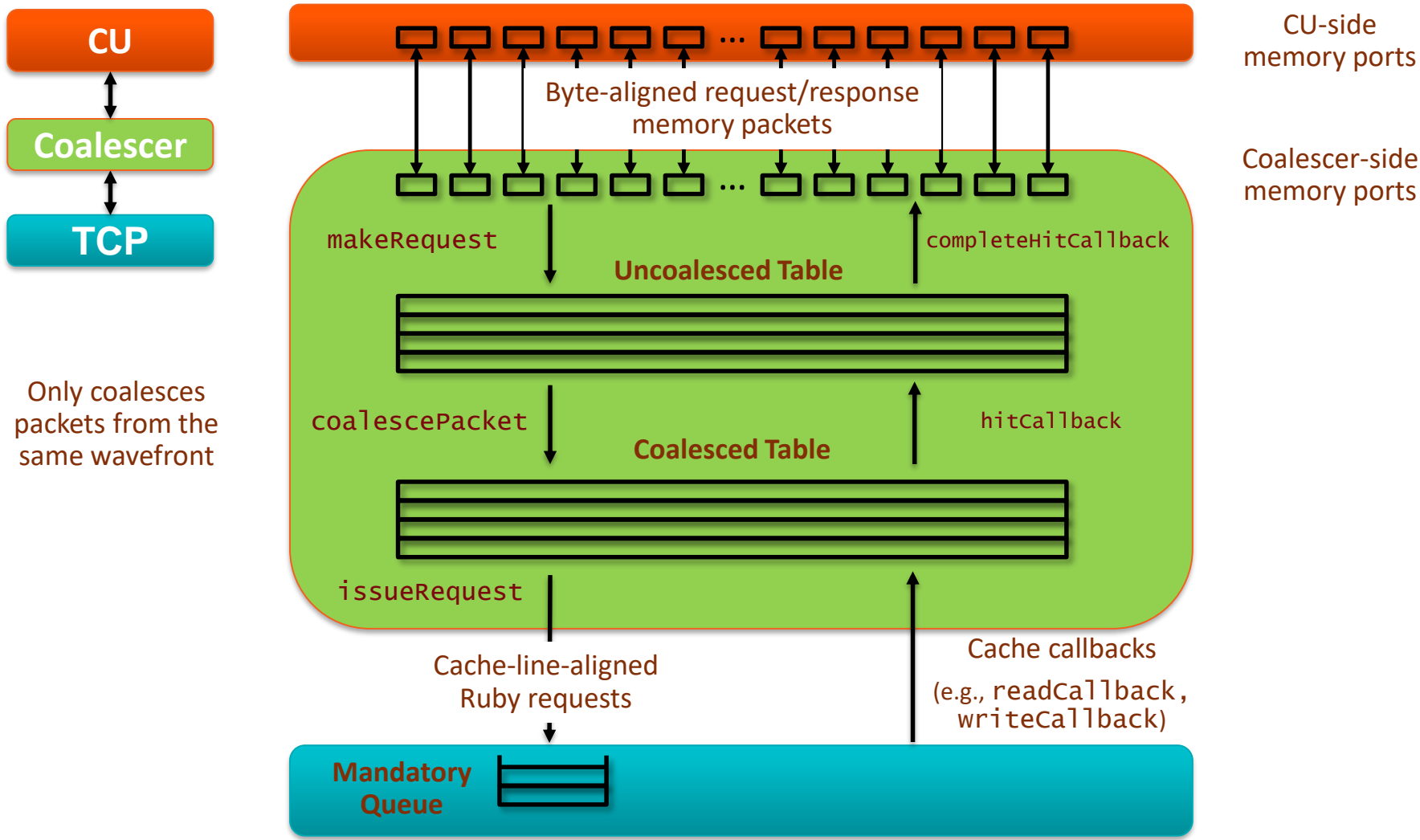
PRIMARY GCN3 MEMORY OPERATIONS



- ▲ Scalar memory operations (to scalar cache)
 - S_LOAD
 - S_STORE
- ▲ Vector memory operations (to TCP)
 - BUFFER_LOAD
 - BUFFER_STORE
 - BUFFER_ATOMIC
- ▲ Cache-wide Synchronization Operation
 - BUFFER_WBINVL1: write back and invalidate TCP cache

CU - MEMORY INTERFACE

GPU MEMORY COALESCING – CRITICAL FOR PERFORMANCE



RUBY MEMORY CONTRIBUTIONS

OUTLINE



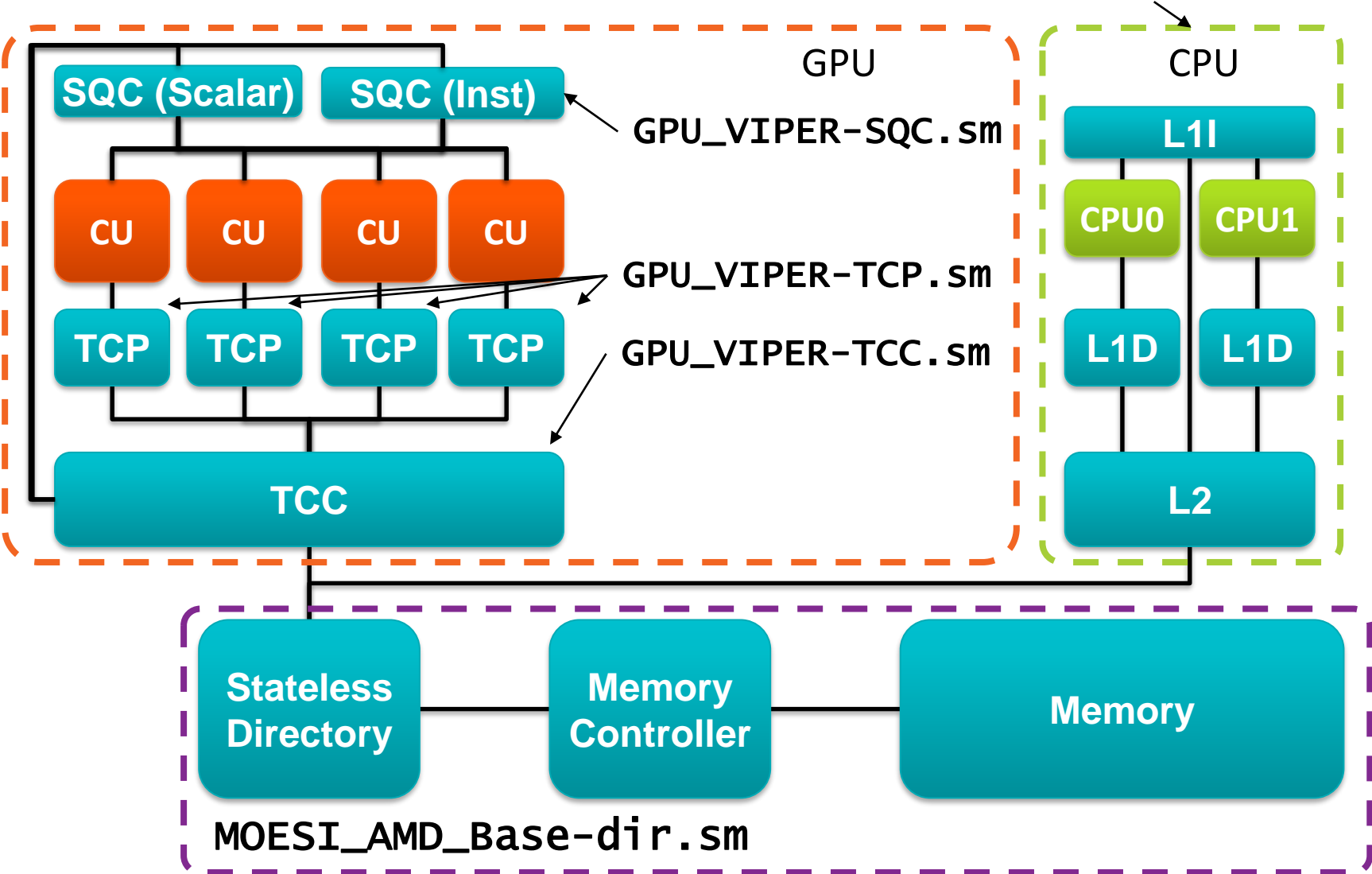
- ▲ Ruby Background
- ▲ CU – Memory Interface
- ▲ GPU VIPER Protocol
- ▲ GPU SLICC Protocol Tester

GPU VIPER PROTOCOL (SEE SRC/MEM/PROTOCOL)

HIGH-LEVEL STRUCTURE



MOESI_AMD_Base-CorePair.sm



GPU VIPER PROTOCOL

WRITE-THROUGH COHERENCE



- ▲ GPU write-through protocol
 - Store & Atomic requests are written through TCP and TCC
 - TCP performs Store & Atomic requests immediately (i.e., no stalling for exclusive permissions)
 - For Store requests, TCP calls back immediately (i.e., `wri teCa llback`) without waiting for ACKs from memory
 - TCP does another callback (i.e., `wri teComp leteCa llback`) when receiving store-complete ACKs from memory
 - Atomic requests are performed in memory
- ▲ Data coherence in TCP
 - TCP-INV requests (e.g., issued by `BUFFER_WBINVL1`) to invalidate entire TCP
- ▲ Data coherence between TCC and CPU caches
 - Data coherence is maintained through a stateless directory (`MOESI_AMD_Base-dir.sm`)
- ▲ Support single coherent address space between CPU and GPU

GPU VIPER PROTOCOL

RELEASE CONSISTENCY SUPPORT



- ▲ Release consistency
 - Hower *et al.* [ASPLOS 2014]
 - HSA System Arch Specification [hsafoundation.com]
- ▲ Acquire fence
 - Issue TCP-INV to invalidate all stale data in TCP
- ▲ Release fence
 - Wait for all outstanding store requests to commit globally (i.e., `writeCompleteCallback`)
 - Wait for all atomic requests (i.e., `atomicCallback`)

RUBY MEMORY CONTRIBUTIONS

OUTLINE



- ▲ Ruby Background
- ▲ CU – Memory Interface
- ▲ GPU VIPER Protocol
- ▲ GPU SLICC Protocol Tester

GPU SLICC PROTOCOL TESTER



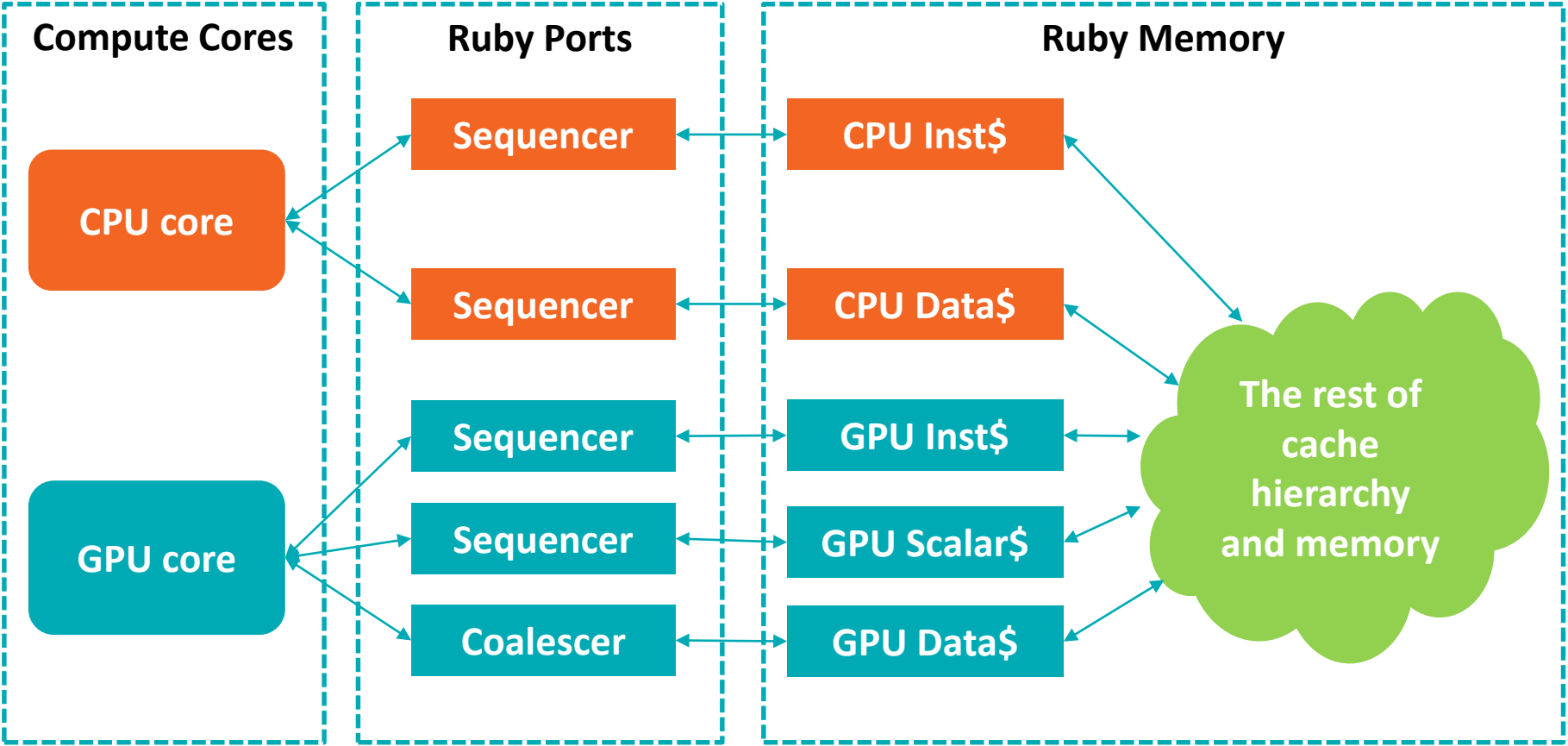
MOTIVATION

- Developing a functionally correct SLICC protocol is challenging
- We need an effective tool to verify a protocol

Desired Features	GPU Protocol Tester
GPU protocol compatibility	✓
Precise validation	✓
Wide bug coverage	✓
Fast bug detection	✓
Intuitive bug report	✓

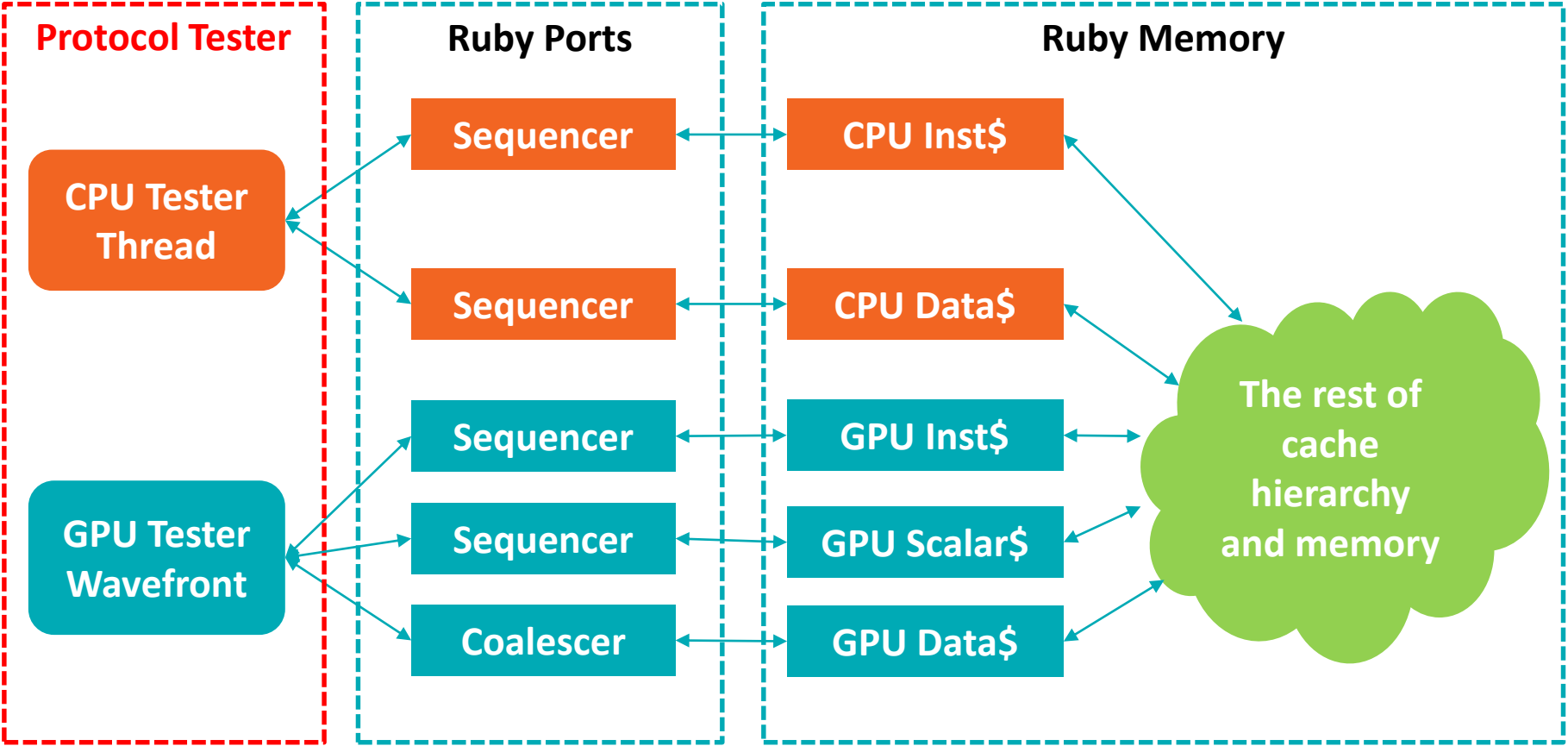
GPU SLICC PROTOCOL TESTER

HIGH-LEVEL STRUCTURE – RUBY INTERFACE



GPU SLICC PROTOCOL TESTER

HIGH-LEVEL STRUCTURE – RUBY INTERFACE



▲ Only support data caches for now

GPU SLICC PROTOCOL TESTER

HIGH-LEVEL STRUCTURE – TESTER THREAD



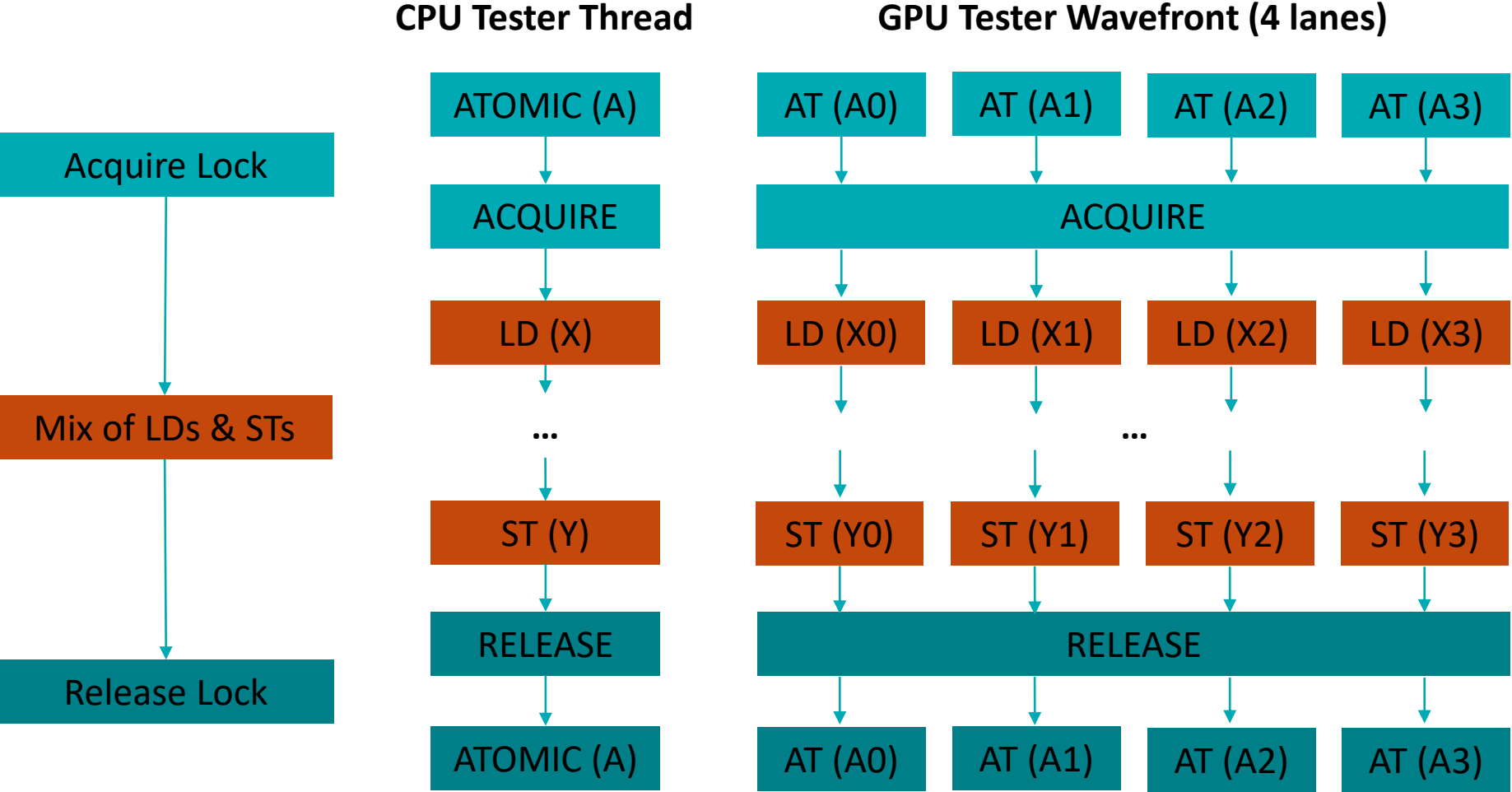
- ▲ Issue sequences (a.k.a., episodes) of memory operations (i.e., Load, Store, Atomic and synchronization operations) with respect to release consistency
- ▲ CPU tester thread
 - Load, Store, and Atomic operations are scalar
 - Acquire and release fence are no-ops
- ▲ GPU Tester Wavefront
 - Load, Store, and Atomic are vector operations (i.e., issued by multiple lanes in a wavefront)
 - Acquire fence issues TCP-INV
 - Release fence waits for all outstanding Store and Atomic operations to complete (i.e., through `writeCompleteCallback` and `atomicCallback`)

GPU SLICC PROTOCOL TESTER

HIGH-LEVEL STRUCTURE - EPISODE

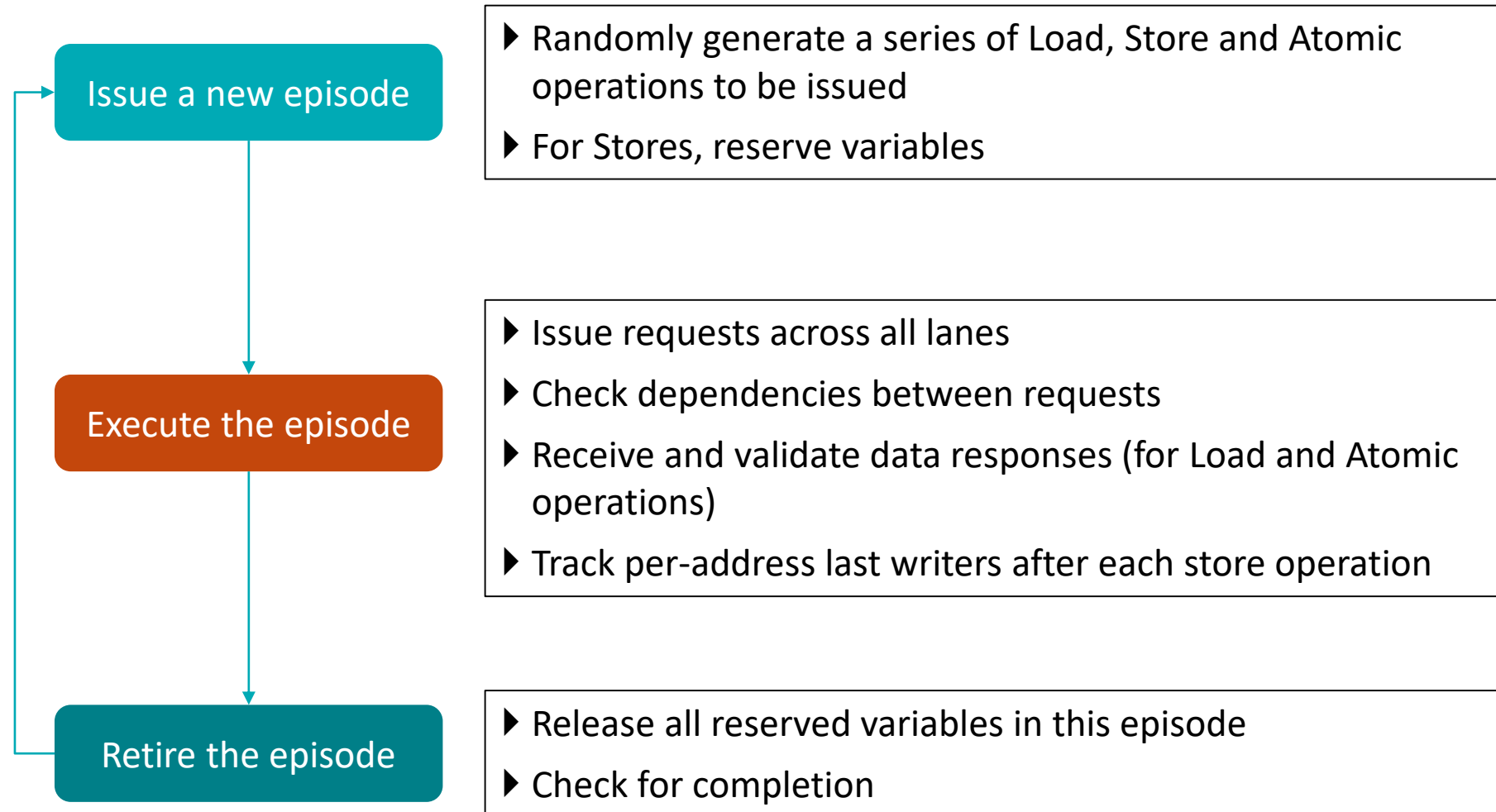


▲ Think of an episode as a critical section



GPU SLICC PROTOCOL TESTER

HIGH-LEVEL STRUCTURE - EXECUTION FLOW PER THREAD

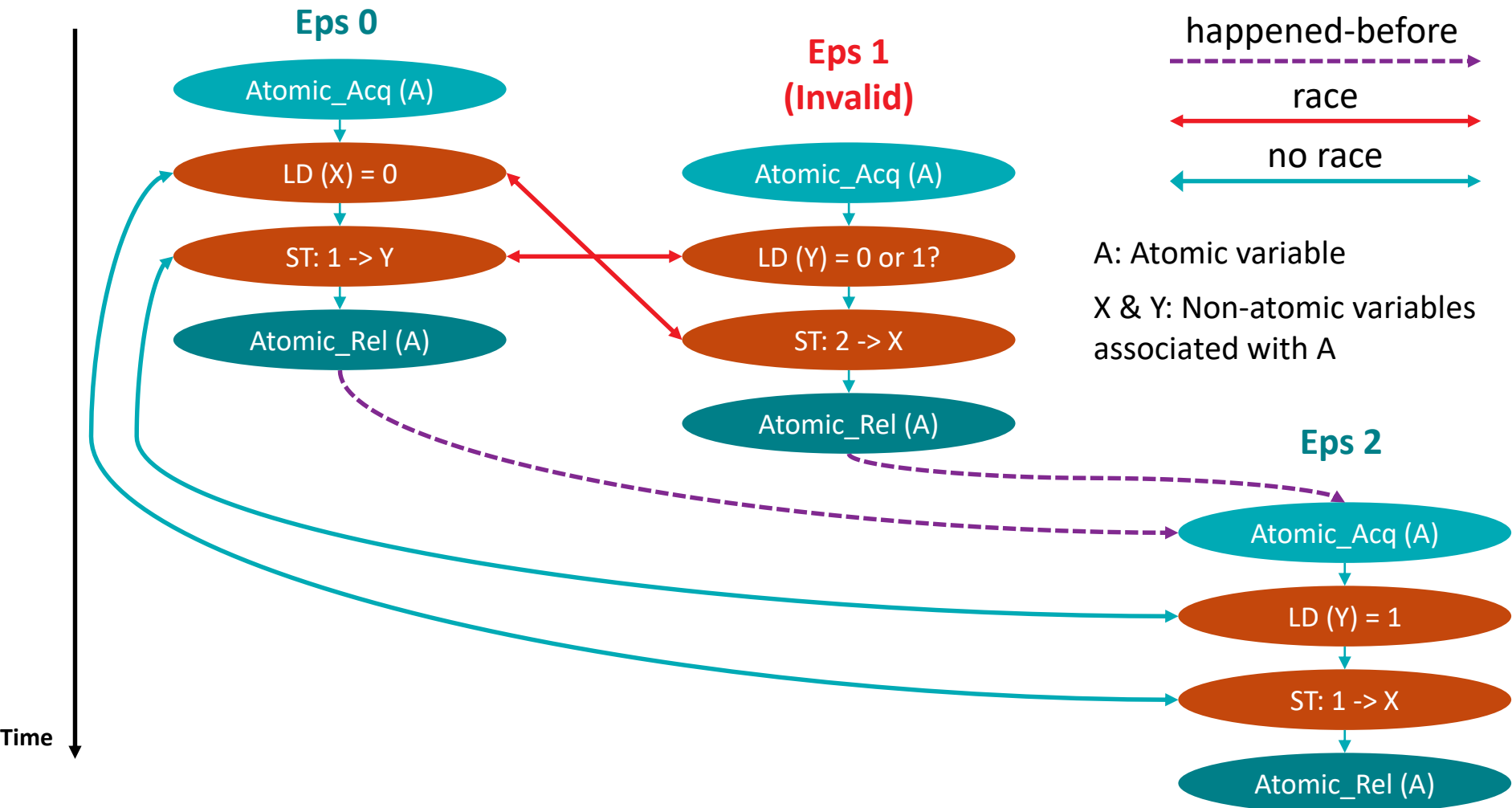


GPU SLICC PROTOCOL TESTER

DATA RACE FREE (DRF) REQUEST STREAM

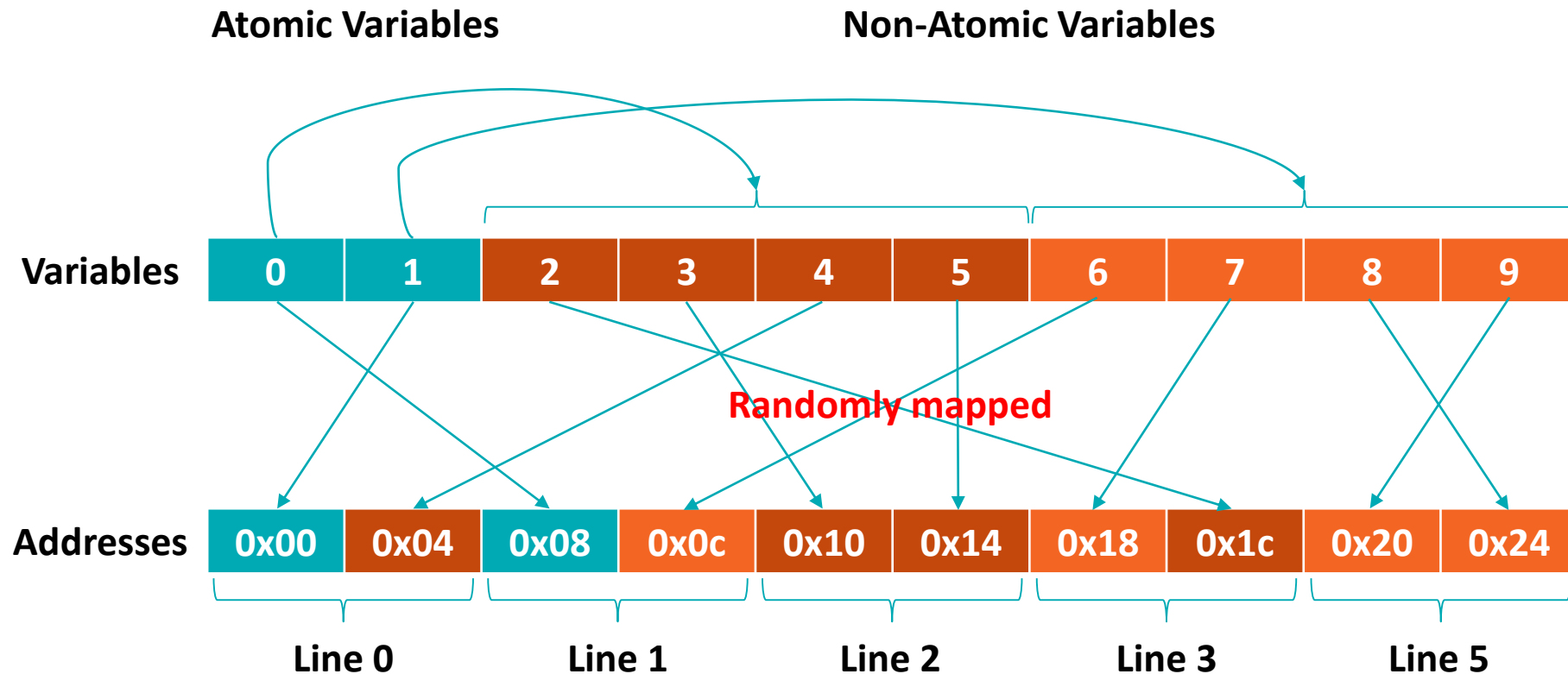


Random GPU Ruby tester does **NOT** cover data-race scenarios



GPU SLICC PROTOCOL TESTER

HIGH LEVEL STRUCTURE – ADDRESS RANGE



- ▲ Random mapping → false sharing
- ▲ An episode can only reserve and issue load/store operations to non-atomic variables associated with the atomic variable it acquires

GPU SLICC PROTOCOL TESTER

VALIDATION



- ▲ Track per-variable last writers
 - Who? (i.e., Lane ID, Wavefront ID, CU ID, and Episode ID)
 - What value was written?
 - When did the store operation happen?
- ▲ For a Load (X), check if return value of X is equal to value written by the last writer to X. The last writer of X could be
 - Either the last ST \rightarrow X in the same episode, or
 - The last ST \rightarrow X in a previous episode
- ▲ For an Atomic (A)
 - Use atomic_inc: increment A by 1 atomically
 - Track how many atomics have been issued so far. Let's say N.
 - Return value must be unique in the range [0..N-1] inclusively

GPU SLICC PROTOCOL TESTER

FAILURE REPORT



▲ Forward progress check

- Report all outstanding requests that have been pending for more than a certain threshold
- Report their target addresses, lane IDs, CU IDs, CPU IDs, and episode IDs

▲ Inconsistency between a Reader and its last Writer (the most common failure!)

```
warn: GpuWavefront(Thread ID = 35, CU ID = 4): Loaded value is not consistent with the last stored value
Thread 35
Episode 727
Lane ID 2
Address [0x52860, line 0x52840]
Loaded value 16
Last writer (Thread ID 35, CU ID 4, Episode ID 727, Value 17, Tick 8454382)
```

▲ Unexpected atomic return

- Report all expected return values
- Report address, lane ID, CU ID, and episode ID

▲ The report really helps trace protocol bugs down **quickly** and **precisely**!

GPU SLICC PROTOCOL TESTER

TESTING OPTIONS AND SCENARIOS



▲ Available configurations

- Cache size
- Address range - # of atomic and normal variables
- System size - # CUs and #CPUs
- Episode length - # of LDs and STs per episode
- Test length - # of episodes

▲ Example testing scenarios

- Small cache size + Wide address range
 - Working set is much larger than the cache size
 - Likely to expose bugs related to cache replacement
- Large cache size + Small address range
 - Likely to expose bugs related to cache coherence (e.g., values are not passed between private caches correctly)
- Long episode length
 - Conflicting accesses in the same episode

Topic	Presenter	Time
Background	Tony	8:00 – 8:15
ROCm, GCN3 ISA, and GPU Arch	Tony	8:15 – 9:15
HSA Implementation in gem5	Sooraj	9:15 – 10:00
Break		10:00 – 10:30
Ruby and GPU Protocol Tester	Tuan	10:30 – 11:15
Demo and Workloads	Matt	11:15 – 11:50
Summary and Questions	All	11:50 – 12:00

HIPIFY'ING CUDA BENCHMARKS



- ▲ HIP code is easy to write if you have experience with CUDA

- ▲ Hipify-perl can automatically convert most CUDA code to HIP code
 - Works well for simple applications (e.g., no library calls)
 - Perl regular expression replacement script
 - Very simple, easy to use (as long as no library calls), fast

- ▲ Useful HIP Information:
 - [Basic Intro to HIP](#)
 - [Basic HIP Tutorial](#)
 - [HIP Porting Guide](#)
 - [Platform-Aware Coding](#)

OUTLINE



Topic	Presenter	Time
Background	Tony	8:00 – 8:15
ROCm, GCN3 ISA, and GPU Arch	Tony	8:15 – 9:15
HSA Implementation in gem5	Sooraj	9:15 – 10:00
Break		10:00 – 10:30
Ruby and GPU Protocol Tester	Tuan	10:30 – 11:15
Demo and Workloads	Matt	11:15 – 11:50
Summary and Questions	All	11:50 – 12:00

COMPARISON TO OTHER GPU SIMULATORS



▲ GPGPU-Sim

- Primarily focused on running Nvidia PTX instructions and CUDA applications
- Functional CPU model, oriented to model discrete GPU systems
- Wisconsin's gem5-gpu added gem5 timing CPU models
 - And a Ruby memory system protocol
- Differences from gem5-GPU:
 - GCN3 instructions and ROCm software stack
 - Unified under the gem5 source control repo

▲ Multi2Sim

- Supports multiple ISAs including AMD Southern Island's Machine ISA
- Limited instruction support
- No transient states in coherence protocol

▲ This is **very different** than the gem5 **NoMAli** emulated GPU device

OBVIOUS IMPROVEMENTS



- ▲ Other GPU ISAs
- ▲ Complete IOMMU model
- ▲ Add graphics functionality
 - Currently compute only
- ▲ Better register model
 - Simple register pool manager: one WG per CU
- ▲ Remove backing store for memory data
 - When running applications, the data from the caches is not used

SUMMARY



- ▲ Covered a very high-level overview of:
 - Introduction to the gem5 APU simulator
 - Mapping between APU system and gem5 APU simulator

- ▲ Topics discussed
 - HSA and GCN Background
 - Compilation and Simulation Flow
 - GPU Core modules
 - GPU memory system model in Ruby
 - GPU protocol tester
 - Comparisons/Improvements

- ▲ Much more detail in the gem5 source code

- ▲ **Please contribute back to this community tool!**

THANK YOU

MANY CONTRIBUTORS OVER THE PAST 8+ YEARS



Alex Dutu	David Roberts	John Alsop	Matt Poremba	Si Li
Ali Jafri	Derek Hower	John Kalamatianos	Matt Sinclair	Sooraj Puthoor
Arka Basu	Dmitri Yudanov	Kishore Punniyamurthy	Michael LeBeane	Steve Reinhardt
Ayse Yilmazer	Eric Van Tassell	Kunal Korgaonkar	Mike Chu	Tanmay Gangwani
Binh Pham	Gagan Sachdev	Lisa Hsu	Myrto Papadopoulou	Tim Rogers
Blake Hechtman	James Wang	Manish Arora	Monir Mozumder	Tony Gutierrez
Brad Beckmann	Jason Power	Marc Orr	Nagesh Lakshminarayana	Tsung Tai Yeh
Brandon Potter	Joel Hestness	Mario Mendez-Lojo	Nilay Vaish	Tushar Krishna
Can Hankendi	Jieming Yin	Mark Leather	Onur Kayiran	Xianwei Zhang
David Hashe	Joe Gross	Mark Wilkening	Srikant Bharadwaj	Yatin Manerkar
		Martin Brown	Shrikanth Ganapathy	Yasuko Eckert

DISCLAIMER & ATTRIBUTION



The information presented in this document is for informational purposes only and may contain technical inaccuracies, omissions and typographical errors.

The information contained herein is subject to change and may be rendered inaccurate for many reasons, including but not limited to product and roadmap changes, component and motherboard version changes, new model and/or product releases, product differences between differing manufacturers, software changes, BIOS flashes, firmware upgrades, or the like. AMD assumes no obligation to update or otherwise correct or revise this information. However, AMD reserves the right to revise this information and to make changes from time to time to the content hereof without obligation of AMD to notify any person of such revisions or changes.

AMD MAKES NO REPRESENTATIONS OR WARRANTIES WITH RESPECT TO THE CONTENTS HEREOF AND ASSUMES NO RESPONSIBILITY FOR ANY INACCURACIES, ERRORS OR OMISSIONS THAT MAY APPEAR IN THIS INFORMATION.

AMD SPECIFICALLY DISCLAIMS ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR ANY PARTICULAR PURPOSE. IN NO EVENT WILL AMD BE LIABLE TO ANY PERSON FOR ANY DIRECT, INDIRECT, SPECIAL OR OTHER CONSEQUENTIAL DAMAGES ARISING FROM THE USE OF ANY INFORMATION CONTAINED HEREIN, EVEN IF AMD IS EXPRESSLY ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

ATTRIBUTION

© 2018 Advanced Micro Devices, Inc. All rights reserved. AMD, the AMD Arrow logo, and combinations thereof are trademarks of Advanced Micro Devices, Inc. Other product names used in this publication are for identification purposes only and may be trademarks of their respective companies. OpenCL is a trademark of Apple Inc. used by permission of Khronos.