



Custom Memory Cube (CMC)



RC Class Presentation

Gongyu Wang
Ph.D candidate

UF UNIVERSITY of
FLORIDA



THE GEORGE
WASHINGTON
UNIVERSITY
WASHINGTON, DC

Virginia
Tech

VIRGINIA POLYTECHNIC INSTITUTE
AND STATE UNIVERSITY

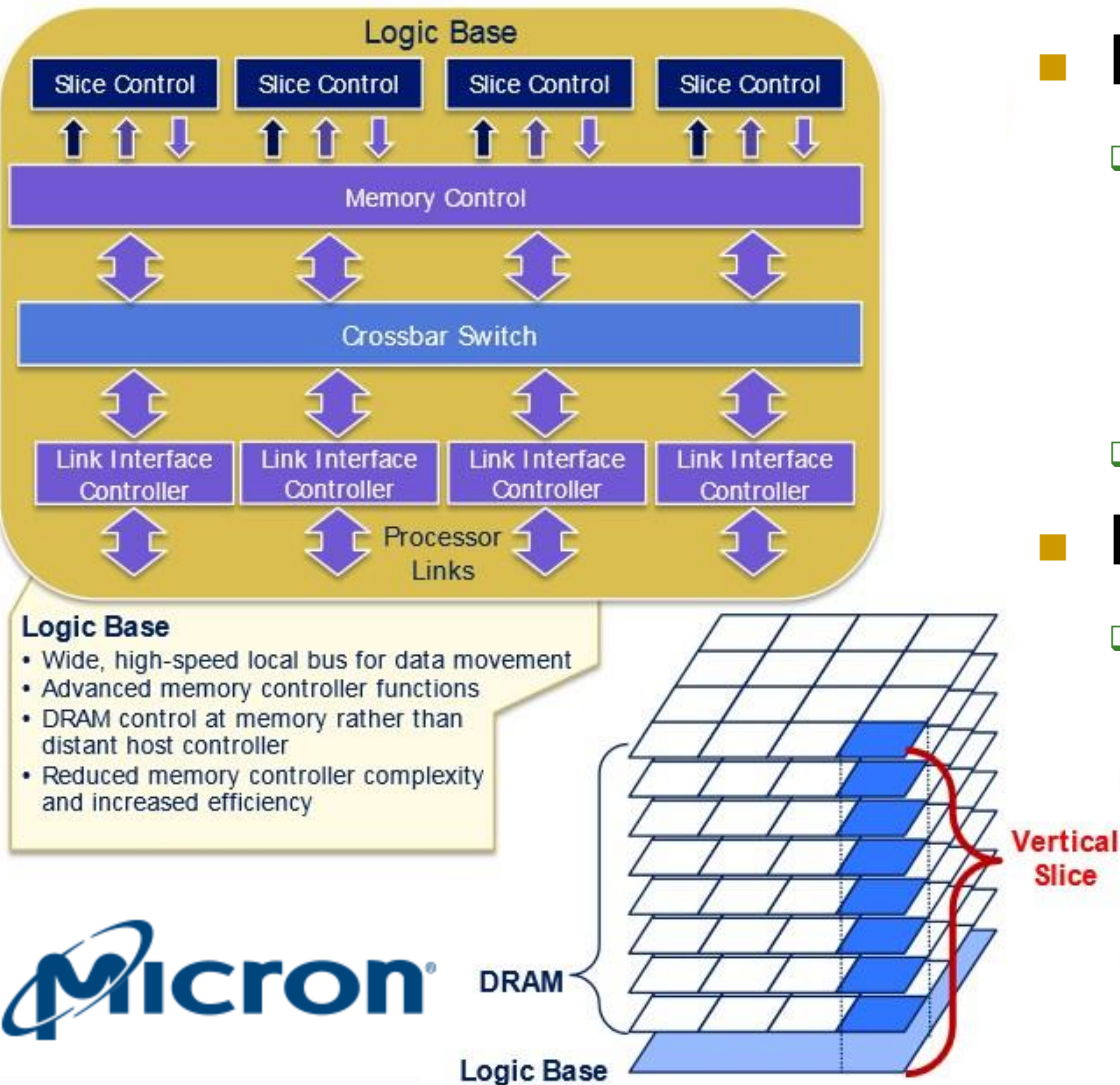
BYU

BRIGHAM YOUNG
UNIVERSITY

Agenda

- Introduction: Hybrid memory cube (HMC)
- Motivation and goals for CMC
- Approach to CMC
 - Two phases
 - FPGA/HMC board – Merlin
- Programming Merlin
 - Hybrid-threading toolset
- Potential RC class course projects

Introduction



- Hybrid memory cube
 - High-performance RAM interface for through-silicon vias (TSV)-based stacked DRAM memory
 - Logic base (layer)
- High bandwidth
 - ~10% energy per bit of current DDR memories

	Bandwidth	Power	W / GBs
DDR3-1333	10.66 GB/s	5.52 W	518 x
DDR4-2666	21.34 GB/s	6.60 W	309 x
HMC (4 DRAMs)	128.00 GB/s	11.08 W	87 x



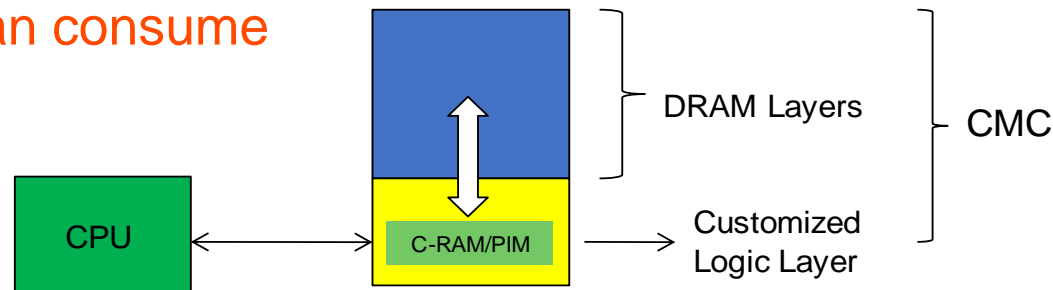
Custom Memory Cube (CMC)

■ Motivation

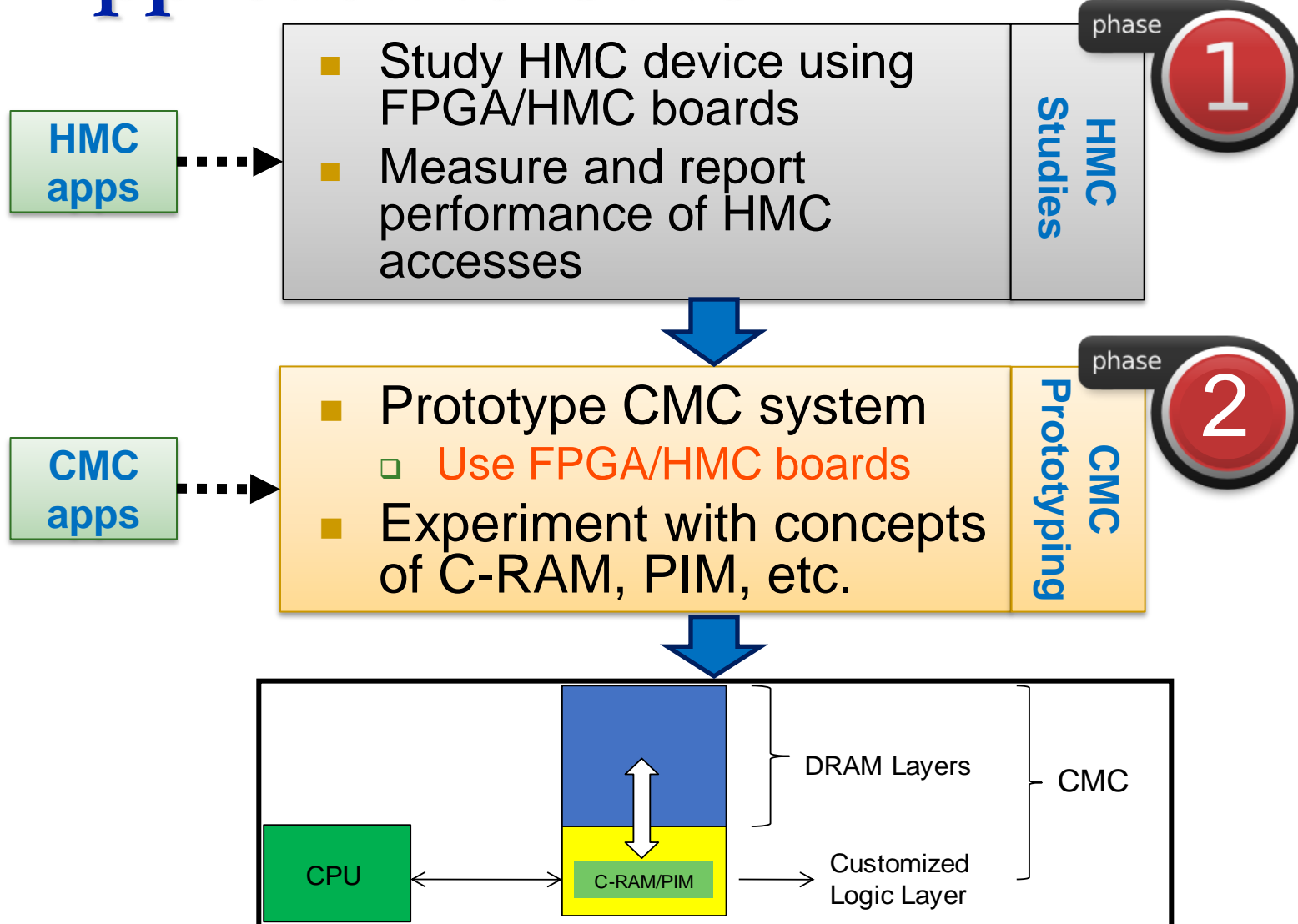
- Increasing need for fast and low-power parallel processing
- Memory accesses of data-intensive apps consumes considerable time and power
- HMC presents much higher memory bandwidth than current processors can consume

■ Goals

- HMC as high-bandwidth memory
 - Exploit performance and power advantages
- CMC: Develop a customizable logic layer for HMC
- Offload processing tasks to CMC
 - Concepts of C-RAM & PIM (e.g., SIMD, ARM64, & RISC-V in logic layer)



Approach to CMC



Phase 1

phase

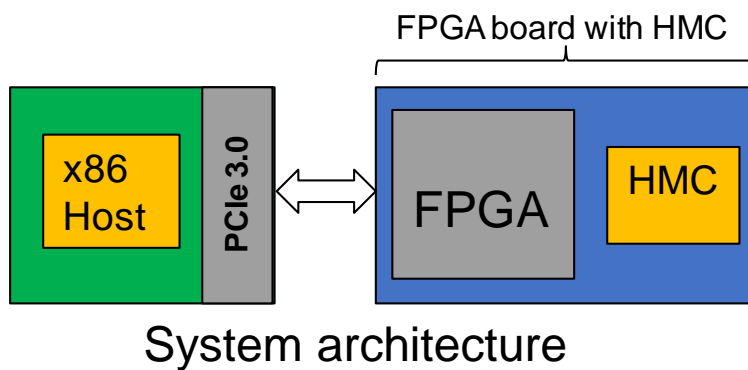
1

■ HMC studies

- Test and study HMC's amenability in selected app. domains
 - Amenability: bandwidth, latency, power, etc.
 - App. Domains: sort, search, cryptography, FFT, etc.

■ Approach

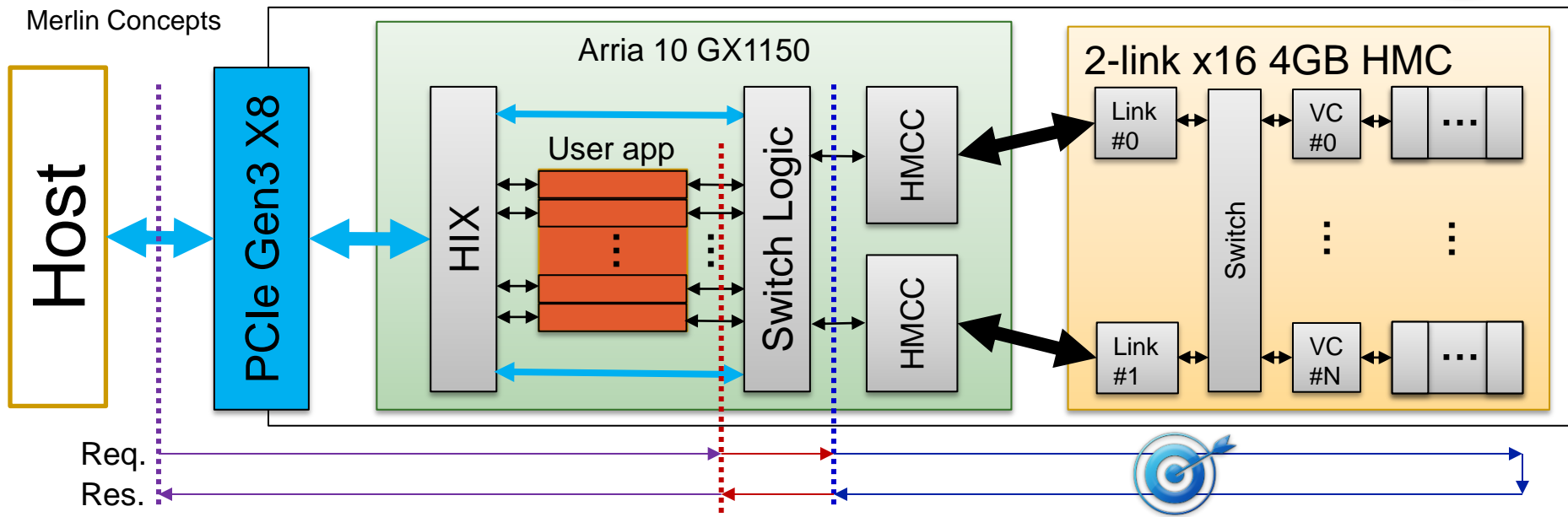
- Develop test applications on FPGA/HMC board
- Measure report access time of HMC in FPGA and report results



Phase 1 with Merlin Board

phase

1



- Goal: test various user apps and report access time of HMC
 - **Wanted latency:** Round-trip latency at inputs of HMCC
 - **Unwanted latencies:** PCIe, HIX, switch logic
 - **Ability to set parameters of HMC access:** packet size, etc.
- Approaches to measuring route-trip latency of HMC access
 - Direct measurement OR subtract unwanted

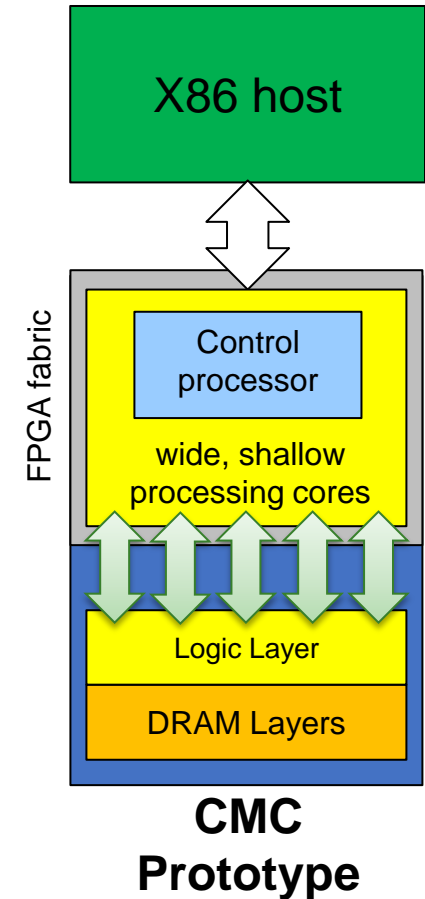


Phase 2

phase

2

- Prototype CMC using FPGA/HMC boards
 - Emulate processing capabilities of CMC on FPGA
- Tasks
 - Define CMC operations based on apps
 - Explore processing capabilities
 - Data processing: “wide” (many-core) and “shallow” (simple) processors; e.g., SIMD
 - Control: e.g., Simplified RISC-V/ARM64
 - Develop demo applications

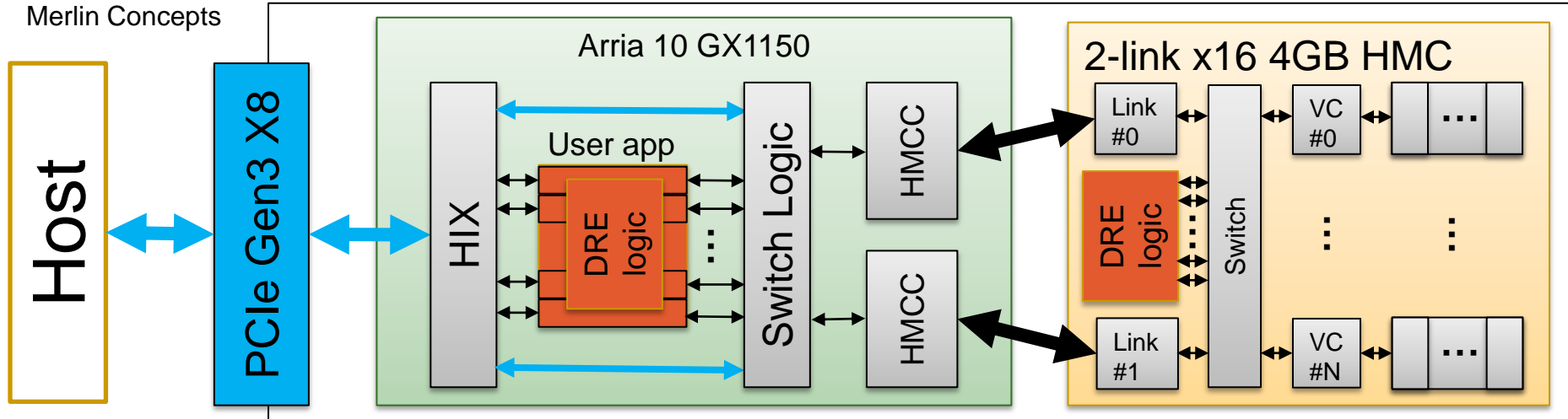


Phase 2 with Merlin Board

phase

2

Merlin Concepts

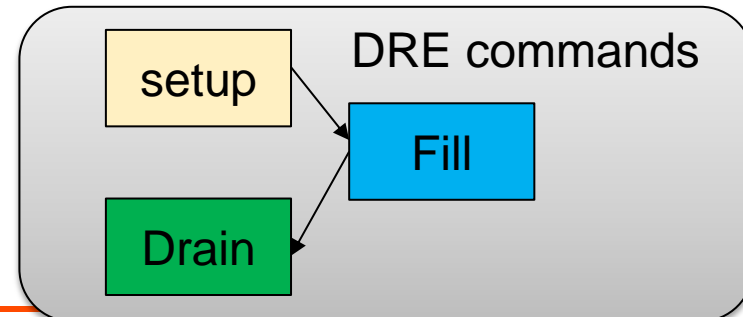


■ Prototype CMC

- ❑ Leverage LLNL's DRE* work for initial prototyping
- ❑ Long term: evaluate CMC based on requirements of other relevant apps

■ Tasks

- ❑ Study, port, and modify DRE for our need
 - Port to Arria 10 FPGA on Merlin board
- ❑ Prototype CMC architecture in FPGA
- ❑ Develop demo applications



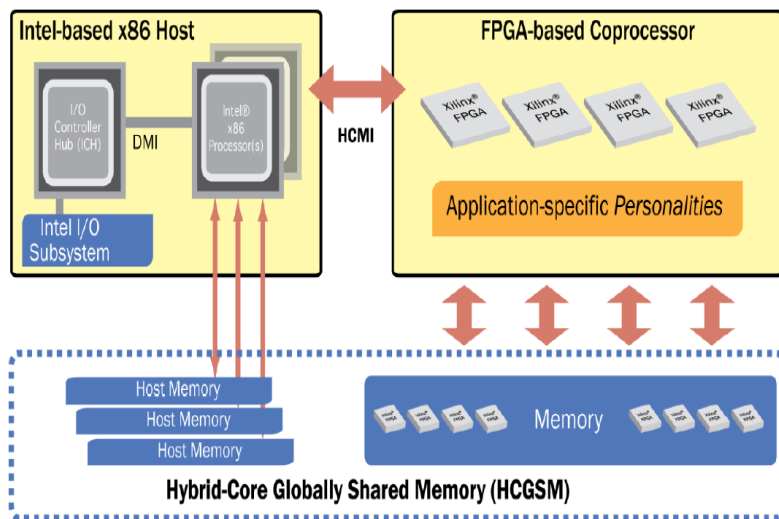
*DRE: Data Reordering/Rearrangement Engine

Agenda

- Introduction: Hybrid memory cube (HMC)
- Motivation and goals for CMC
- Approach to CMC
 - Two phases
 - FPGA/HMC board – Merlin
- Programming Merlin
 - Hybrid-threading toolset
- Potential RC class course projects

Convey HT Background

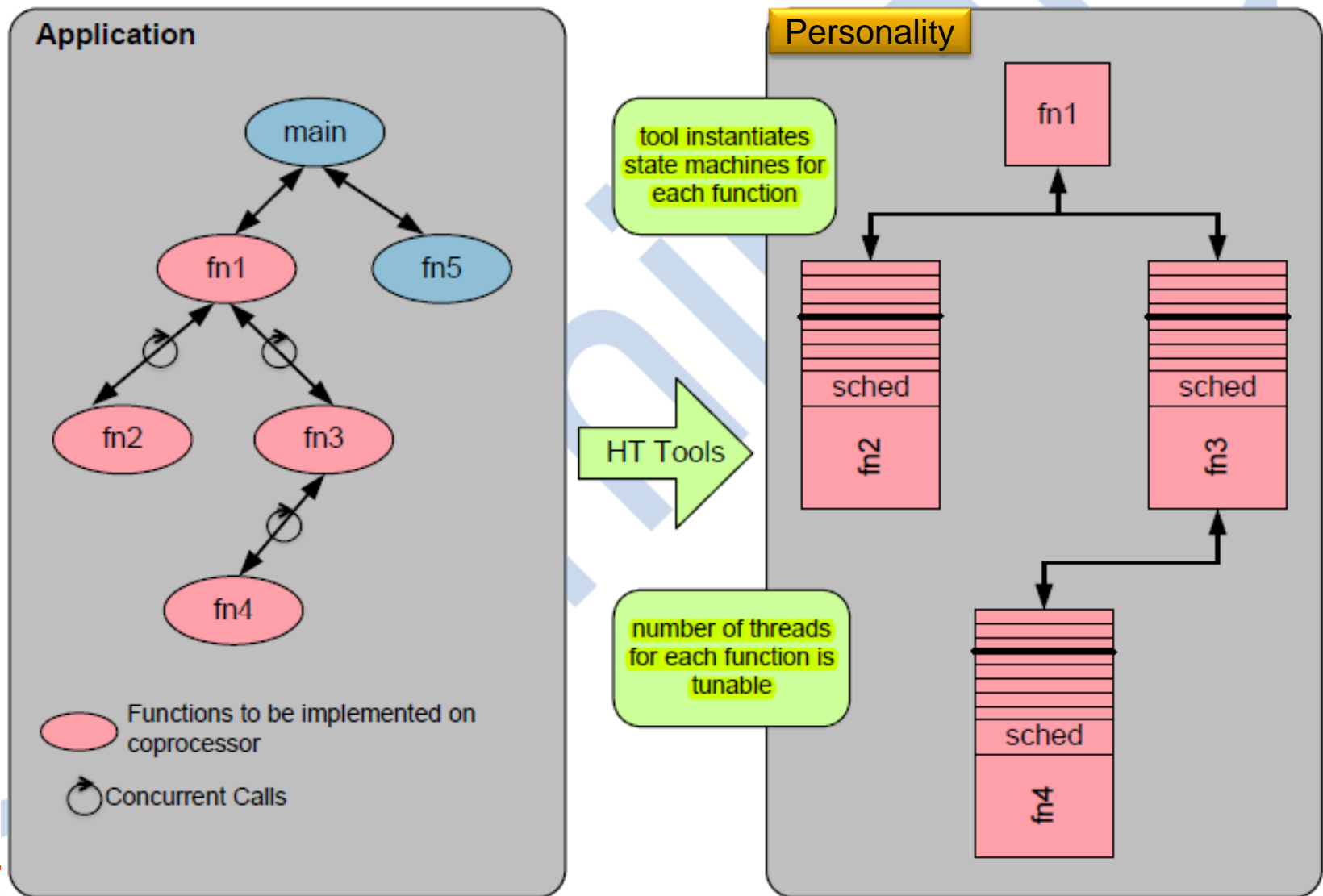
- Convey computers
 - ❑ Startup since 2006, based @ Texas
 - ❑ Hybrid-core computing platforms
 - ❑ Joined CHREC in 2013
 - Via F3, focusing on bioinformatics app acceleration
- Convey system philosophy – hybrid-core computing



Features	Gidel	Convey	SRC
Programming	C/Cpp + HDL	C/Cpp/ASM + HT + (HDL)	C/C++
Data Addressing	explicit, separate	explicit, consistent	implicit
Host Connection	PCIe	FSB, QPI*	PCIe

* via PCIe

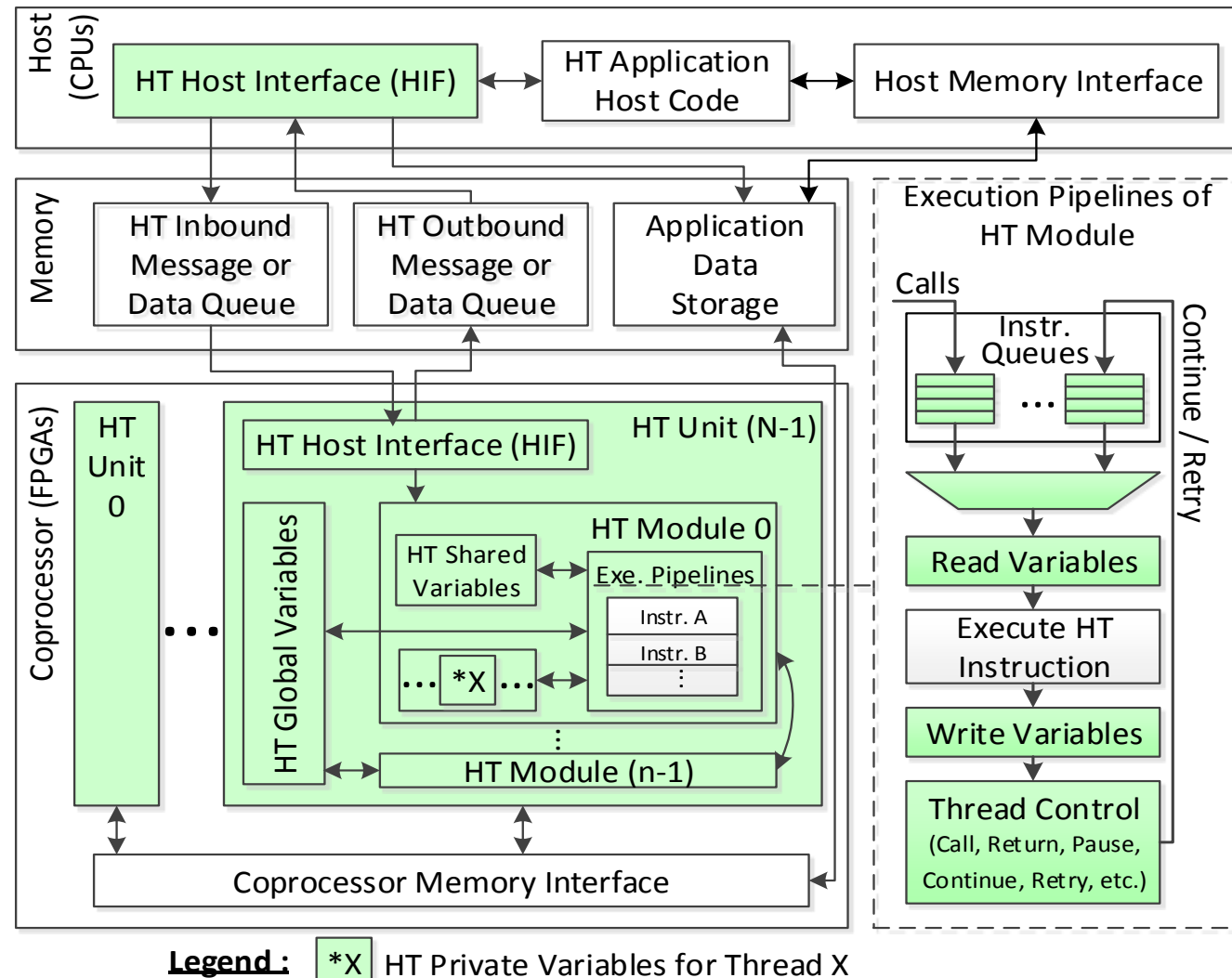
Overview of HT Programming Model



Application "call-graph" structure

HT Infrastructure

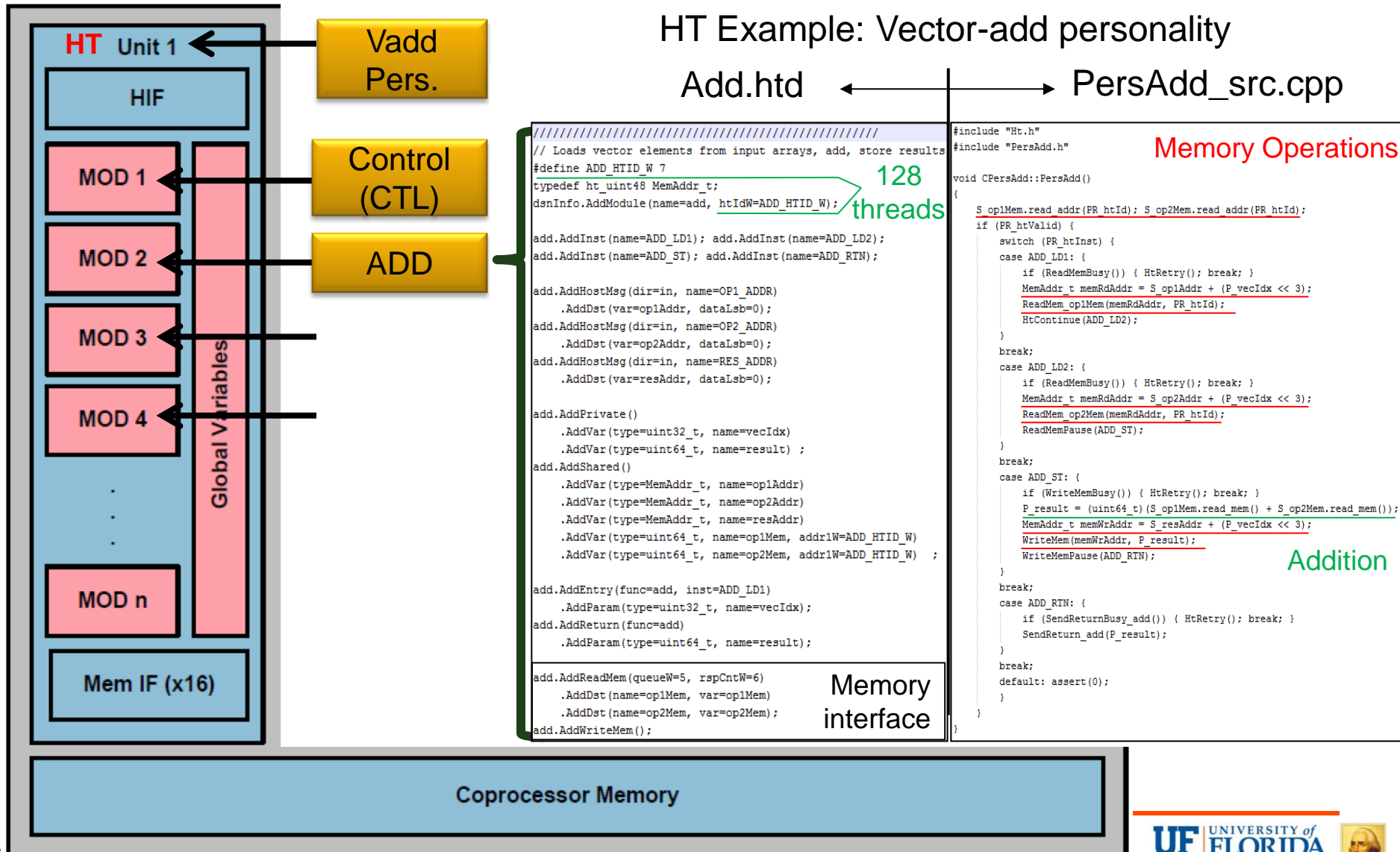
- Host-coproc. strucutre
 - Message/Data interfaces
- HT unit
 - Replication mechanism
- HT module
 - Essential functional blocks
 - Memory hierarchy
 - Can be replicated
 - Messaging interfaces among modules
- Instructions
 - User-defined behaviors of modules



User Programmable Parts

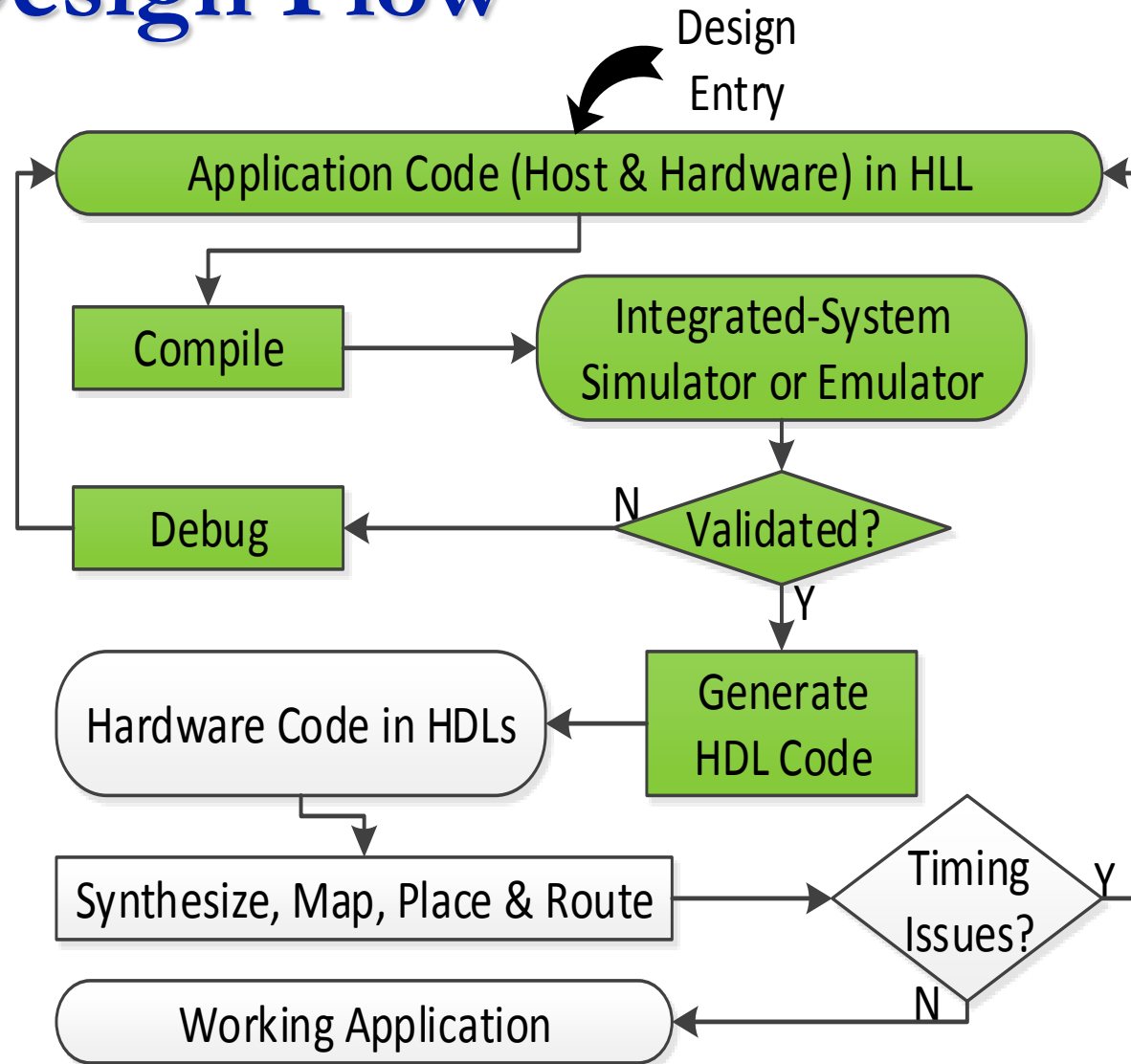
- Hardware description files (.htd)
 - Declare modules, variables, memory interfaces, etc.
- Module instruction files (.cpp)
 - Define behaviors of modules, including memory I/O, messaging interface, computation, etc.
- CTL module example on next slide
 - Responsible for:
 - Accept calls from host
 - Spawn threads onto other HT modules
 - Return to host after all threads returned from other HT modules

HT Code Example



High-level Design Flow

- HT bases on SystemC
 - Cycle-accurate simulation
 - Weed out inter-cycle bugs
- OpenCL
 - Behavioral emulation
 - Tools handle timing issues



Potential Course Project

■ Goals

- ❑ Learn HT programming and develop application for Merlin board using HMC device
- ❑ Get involved with CMC research

■ Tasks

- ❑ Start with simple code examples
 - Hello world, Vadd, Bucket Sort, etc.
- ❑ Pick an application
 - Smith-Waterman, Radix Sort, Improved bucket sort, DRE, etc.
 - Or bring your own application
- ❑ Use SystemC simulation to validate design

■ Bonus task

- ❑ Code review with Gongyu Wang
- ❑ Compile to bit file and run application on Merlin board