GPU Teaching Kit

Accelerated Computing

ILLINOIS
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN
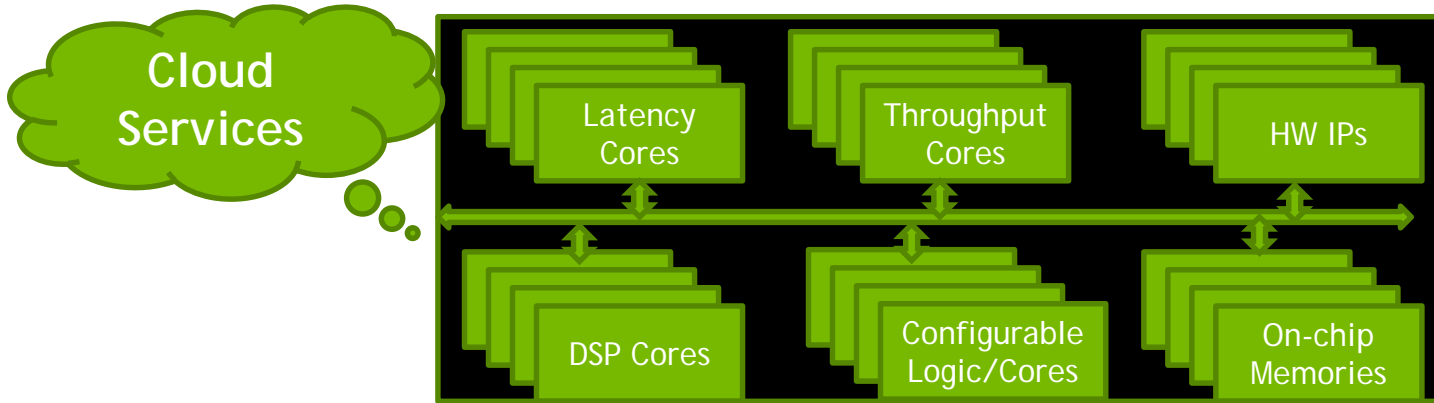
Lecture 1.2 – Course Introduction

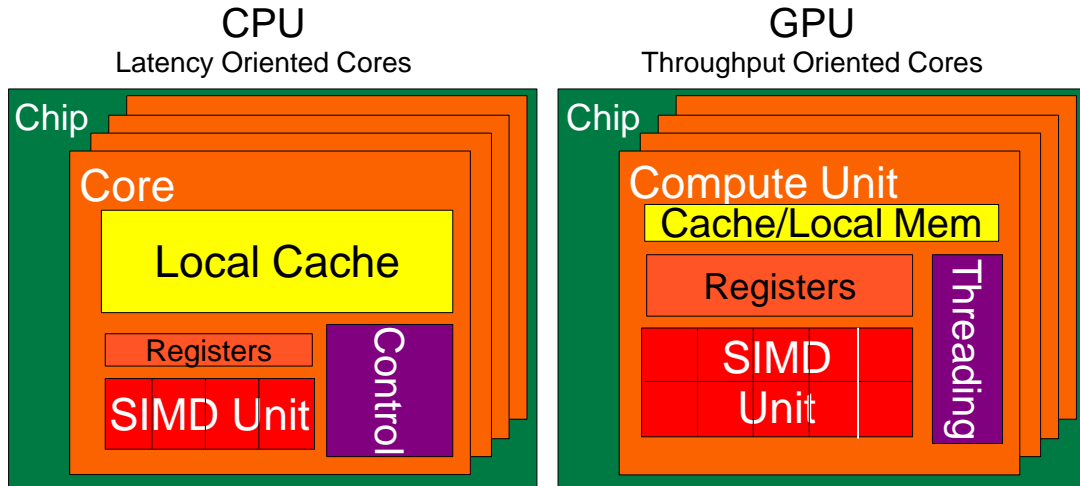Introduction to Heterogeneous Parallel Computing

# Objectives

- To learn the major differences between latency devices (CPU cores) and throughput devices (GPU cores)
- To understand why winning applications increasingly use both types of devices

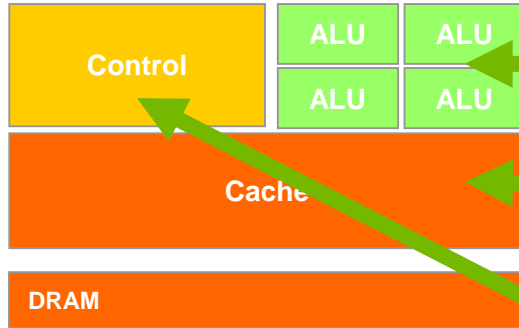# Heterogeneous Parallel Computing

– Use the best match for the job (heterogeneity in mobile SOC)

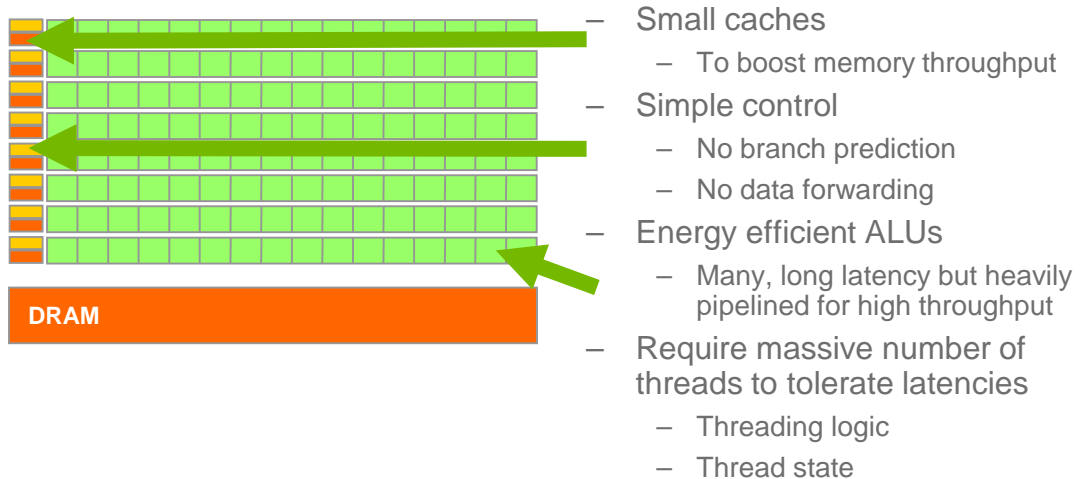# CPU and GPU are designed very differently



CPU
Latency Oriented Cores

GPU
Throughput Oriented Cores

Chip

Core

Local Cache

Registers

SIMD Unit

Control

Chip

Compute Unit

Cache/Local Mem

Registers

SIMD Unit

Threading

# CPUs: Latency Oriented Design



- – Powerful ALU
  - – Reduced operation latency
- – Large caches
  - – Convert long latency memory accesses to short latency cache accesses
- – Sophisticated control
  - – Branch prediction for reduced branch latency
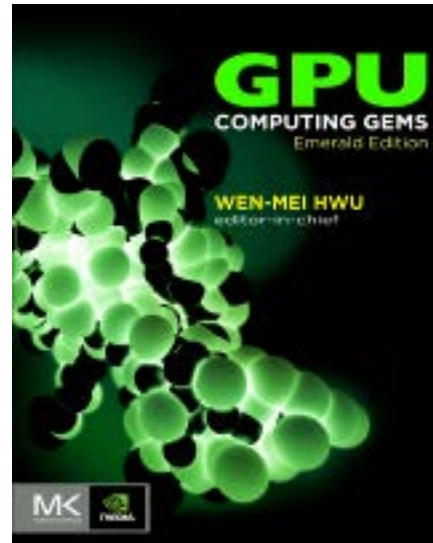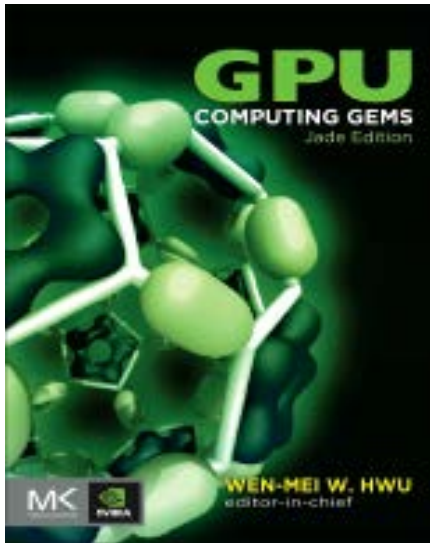  - – Data forwarding for reduced data latency

# GPUs: Throughput Oriented Design



**DRAM**

– Small caches
  – To boost memory throughput
– Simple control
  – No branch prediction
  – No data forwarding
– Energy efficient ALUs
  – Many, long latency but heavily pipelined for high throughput
– Require massive number of threads to tolerate latencies
  – Threading logic
  – Thread state

NVIDIA

# Winning Applications Use Both CPU and GPU

– CPUs for sequential parts where latency matters
  – CPUs can be 10X+ faster than GPUs for sequential code

– GPUs for parallel parts where throughput wins
  – GPUs can be 10X+ faster than CPUs for parallel code

NVIDIA   ILLINOIS

# GPU computing reading resources



90 articles in two volumes

# Heterogeneous Parallel Computing in Many Disciplines

| | | | | |
|---|---|---|---|---|
| Financial Analysis | Scientific Simulation | Engineering Simulation | Data Intensive Analytics | Medical Imaging |
| Digital Audio Processing | Digital Video Processing | Computer Vision | Biomedical Informatics | Electronic Design Automation |
| | | | Statistical Modeling | Numerical Methods |
| | | | Ray Tracing Rendering | Interactive Physics |

NVIDIA   ILLINOIS

GPU Teaching Kit

Accelerated Computing