



GPU Teaching Kit

Accelerated Computing



西南石油大学 计算机科学学院

SCHOOL OF COMPUTER SCIENCE, SOUTHWEST PETROLEUM UNIVERSITY



# Lecture 1.1 – Course Introduction

Course Introduction and Overview

# Course Goals

- **Learn how to program heterogeneous (异构) parallel computing systems and achieve**
  - High performance and energy-efficiency
  - Functionality(功能性) and maintainability(可维护性)
  - Scalability(可扩展性) across future generations
  - Portability(可移植性) across vendor devices
- **Technical subjects**
  - Parallel programming API, tools and techniques
  - Principles and patterns of parallel algorithms
  - Processor architecture features and constraints

# Course Content

Module 1 Course Introduction	<ul style="list-style-type: none"><li>• Course Introduction and Overview</li><li>• Introduction to Heterogeneous Parallel Computing</li><li>• Portability and Scalability in Heterogeneous Parallel Computing</li></ul>
Module 2 Introduction to CUDA C	<ul style="list-style-type: none"><li>• CUDA C vs. CUDA Libs vs. OpenACC</li><li>• Memory Allocation and Data Movement API Functions</li><li>• Data Parallelism and Threads</li><li>• Introduction to CUDA Toolkit</li></ul>
Module 3 CUDA Parallelism Model	<ul style="list-style-type: none"><li>• Kernel-Based SPMD Parallel Programming</li><li>• Multidimensional Kernel Configuration</li><li>• Color-to-Greyscale Image Processing Example</li><li>• Blur Image Processing Example</li></ul>
Module 4 Memory Model and Locality	<ul style="list-style-type: none"><li>• CUDA Memories</li><li>• Tiled Matrix Multiplication</li><li>• Tiled Matrix Multiplication Kernel</li><li>• Handling Boundary Conditions in Tiling</li><li>• Tiled Kernel for Arbitrary Matrix Dimensions</li></ul>
Module 5 Kernel-based Parallel Programming	<ul style="list-style-type: none"><li>• Histogram (Sort) Example</li><li>• Basic Matrix-Matrix Multiplication Example</li><li>• Thread Scheduling</li><li>• Control Divergence</li></ul>

# Course Content

Module 6 Performance Considerations: Memory	<ul style="list-style-type: none"><li>• DRAM Bandwidth</li><li>• Memory Coalescing in CUDA</li></ul>
Module 7 Atomic Operations	<ul style="list-style-type: none"><li>• Atomic Operations</li></ul>
Module 8 Parallel Computation Patterns (Part 1)	<ul style="list-style-type: none"><li>• Convolution</li><li>• Tiled Convolution</li><li>• 2D Tiled Convolution Kernel</li></ul>
Module 9 Parallel Computation Patterns (Part 2)	<ul style="list-style-type: none"><li>• Tiled Convolution Analysis</li><li>• Data Reuse in Tiled Convolution</li></ul>
Module 10 Performance Considerations: Parallel Computation Patterns	<ul style="list-style-type: none"><li>• Reduction</li><li>• Basic Reduction Kernel</li><li>• Improved Reduction Kernel</li></ul>
Module 11 Parallel Computation Patterns (Part 3)	<ul style="list-style-type: none"><li>• Scan (Parallel Prefix Sum)</li><li>• Work-Inefficient Parallel Scan Kernel</li><li>• Work-Efficient Parallel Scan Kernel</li><li>• More on Parallel Scan</li></ul>

# Course Content

Module 12 Performance Considerations: Scan Applications	<ul style="list-style-type: none"><li>• Scan Applications: Per-thread Output Variable Allocation</li><li>• Scan Applications: Radix Sort</li><li>• Performance Considerations (Histogram (Atomics) Example)</li><li>• Performance Considerations (Histogram (Scan) Example)</li></ul>
Module 13 Advanced CUDA Memory Model	<ul style="list-style-type: none"><li>• Advanced CUDA Memory Model</li><li>• Constant Memory</li><li>• Texture Memory</li></ul>
Module 14 Sparse Matrix Computation	<ul style="list-style-type: none"><li>• Parallel SpMV Using CSR</li><li>• Padding and Transposition</li><li>• Using a Hybrid Approach to Regulate Padding</li><li>• Sorting and Partitioning for Regularization</li></ul>
Module 15 Merge Sort	<ul style="list-style-type: none"><li>• Sequential Merge Algorithm</li><li>• Parallelization Approach</li><li>• Co-Rank Function Implementation</li><li>• Basic Parallel Merge Kernel</li><li>• Tiled Merge Kernel</li><li>• Circular-Buffer Merge Kernel</li></ul>
Module 16 Graph Search	<ul style="list-style-type: none"><li>• Breadth-First Search</li><li>• Sequential BFS Function</li><li>• Parallel BFS Function</li><li>• Optimizations</li></ul>
Module 17 CUDA Python Using Numba	<ul style="list-style-type: none"><li>• CUDA Python using Numba</li></ul>

# Teaching Arrangement

- **32 Lecture Hours**
- **16 Laboratory Hours (Important)**

# Grade Components

- **In-class performance (70%)**
  - Final written exam (70%)
  - Regular performance (30%)
    - Attendance (20%)
    - In-class quizzes (80%)
- **Lab performance (30%)**
  - Nvidia Certifications (50%)
    - 2 NVIDIA certifications
    - 2 corresponding quizzes.
  - 4 assignments (50%)

Students who miss more than three classes will be disqualified from participating in the final exam!!!

# Textbooks, Reference Books, and Resources

## – Textbooks

- 大卫·B. 柯克(David B. Kirk),胡文美(Wen-meí W. Hwu)著
- 大规模并行处理器程序设计 (英文版·原书第3版)
- 机械工业出版社, 2020年

## – Reference Books

- 高性能计算系列丛书·CUDA并行政程序设计: GPU编程指南  
[CUDA Programming: A Developer's Guide to Parallel Computing with GPUs], 机械工业出版社, 2014年
- Professional Cuda C Programming, Wiley, 2014年

## – Resources

- <https://docs.nvidia.com/cuda/doc/index.html>





GPU Teaching Kit  
Accelerated Computing



西南石油大学 计算机科学学院  
SCHOOL OF COMPUTER SCIENCE, SOUTHWEST PETROLEUM UNIVERSITY



# Lecture 1.2 – Course Introduction

Introduction to Heterogeneous Parallel Computing

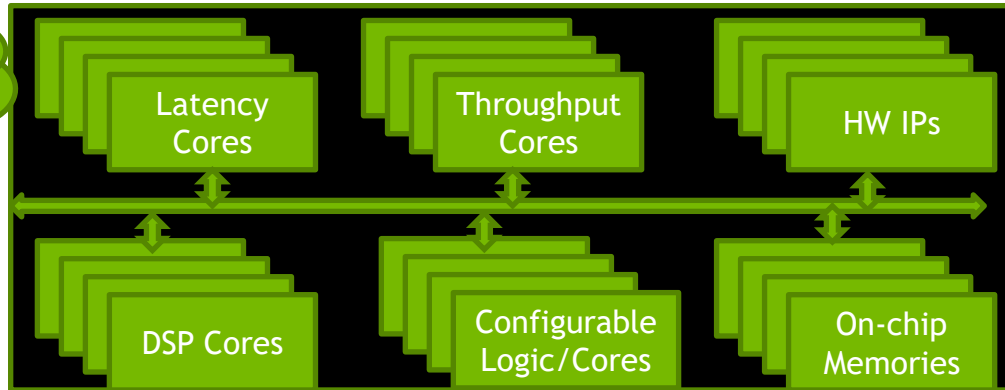
# Objectives

- 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

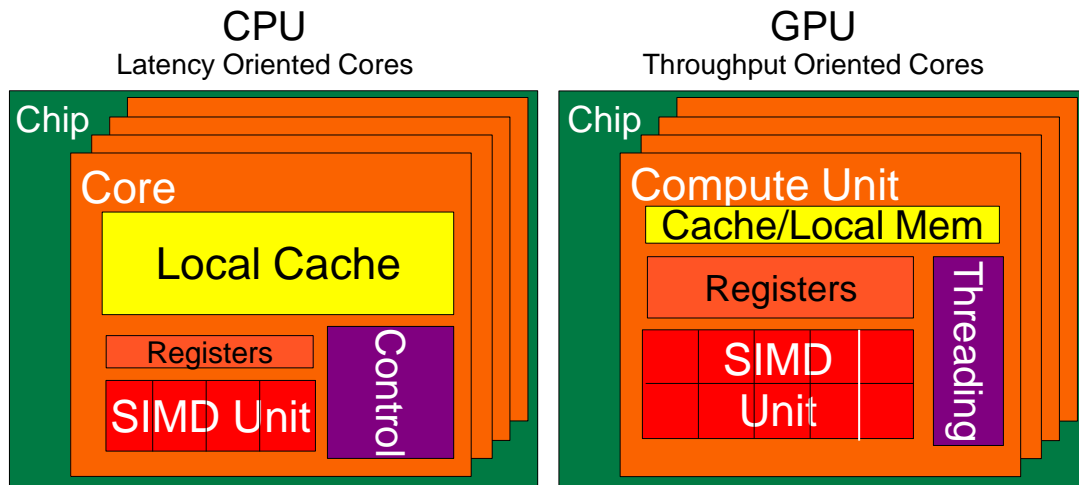
# Heterogeneous Parallel Computing

- Use the best match for the job (heterogeneity in mobile SOC(System on a Chip))

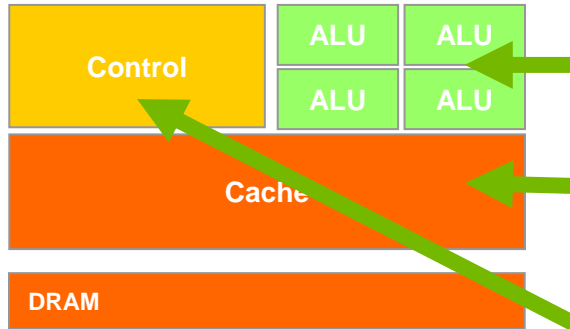
Cloud  
Services



# CPU and GPU are designed very differently

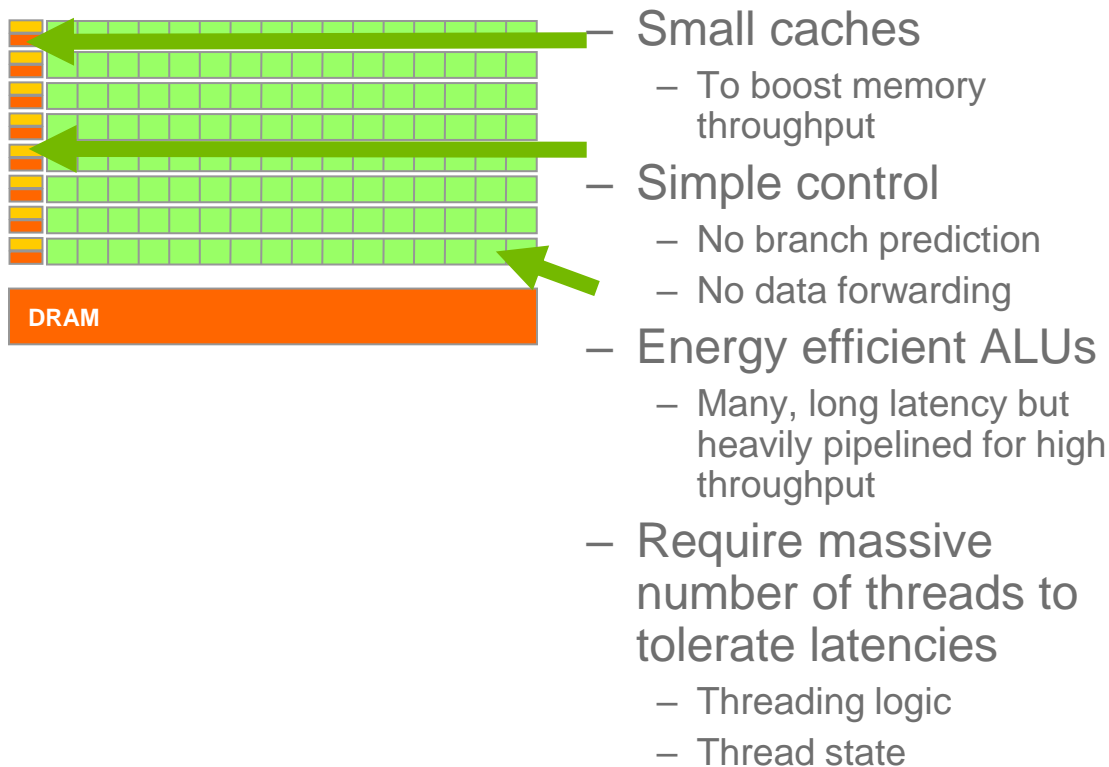


# CPUs: Latency Oriented Design



- 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



# 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

# 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





GPU Teaching Kit  
Accelerated Computing



西南石油大学 计算机科学学院  
SCHOOL OF COMPUTER SCIENCE, SOUTHWEST PETROLEUM UNIVERSITY



# Lecture 1.3 – Course Introduction

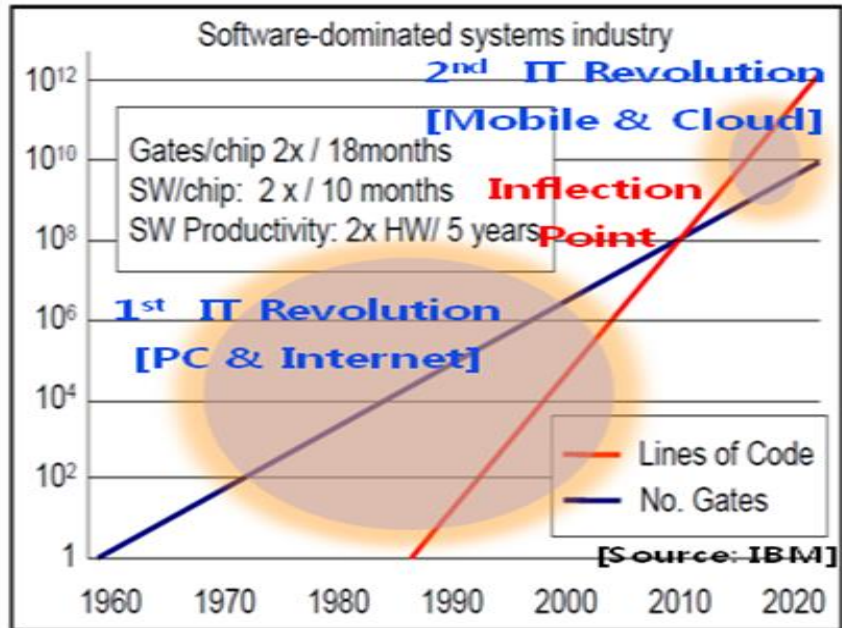
Portability and Scalability in Heterogeneous Parallel Computing

# Objectives

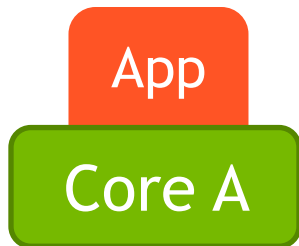
- To understand the importance and nature of **scalability** and **portability** in parallel programming

# Software Dominates System Cost

- SW lines per chip increases at 2x/10 months
- HW gates per chip increases at 2x/18 months
- Future systems must minimize software redevelopment



# Keys to Software Cost Control



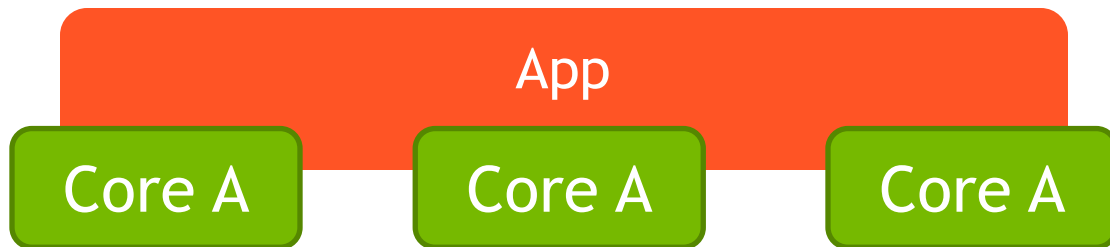
– Scalability

# Keys to Software Cost Control



- Scalability
  - The same application runs efficiently on new generations of cores

# Keys to Software Cost Control



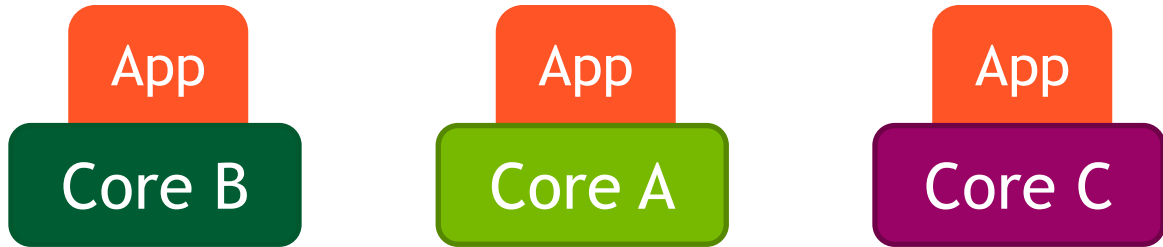
## – Scalability

- The same application runs efficiently on new generations of cores
- **The same application runs efficiently on more of the same cores**

# More on Scalability

- Performance growth with HW generations
  - Increasing number of compute units (cores)
  - Increasing number of threads
  - Increasing vector length
  - Increasing pipeline depth
  - Increasing DRAM burst size
  - Increasing number of DRAM channels
  - Increasing data movement latency

# Keys to Software Cost Control



- Scalability
- **Portability**
  - The same application runs efficiently on different types of cores



# Keys to Software Cost Control



- Scalability
- Portability
  - The same application runs efficiently on different types of cores
  - The same application runs efficiently on systems with different organizations and interfaces

# More on Portability

- Portability across many different HW types
  - Across ISAs (Instruction Set Architectures) - X86 vs. ARM, etc.
  - Latency oriented CPUs vs. throughput oriented GPUs
  - Across parallelism models - SIMD vs. threading
  - Across memory models - Shared memory vs. distributed memory



# GPU Teaching Kit

Accelerated Computing



西南石油大学 计算机科学学院

SCHOOL OF COMPUTER SCIENCE, SOUTHWEST PETROLEUM UNIVERSITY



Some content in the PPT is sourced from The GPU Teaching Kit.  
The GPU Teaching Kit is licensed by NVIDIA and the University of Illinois under  
the **Creative Commons Attribution-NonCommercial 4.0 International License**.