



GPU Teaching Kit  
Accelerated Computing



# Lecture 1.1 – Course Introduction

Course Introduction and Overview

# Course Goals

- Learn how to program heterogeneous parallel computing systems and achieve
  - High performance and energy-efficiency
  - Functionality and maintainability
  - Scalability across future generations
  - Portability across vendor devices
- Technical subjects
  - Parallel programming API, tools and techniques
  - Principles and patterns of parallel algorithms
  - Processor architecture features and constraints

# People

- Wen-mei Hwu (University of Illinois)
- David Kirk (NVIDIA)
- Joe Bungo (NVIDIA)
- Mark Ebersole (NVIDIA)
- Abdul Dakkak (University of Illinois)
- Izzat El Hajj (University of Illinois)
- Andy Schuh (University of Illinois)
- John Stratton (Colgate College)
- Isaac Gelado (NVIDIA)
- John Stone (University of Illinois)
- Javier Cabezas (NVIDIA)
- Michael Garland (NVIDIA)

# 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 Floating Point Considerations	<ul style="list-style-type: none"><li>• Floating Point Precision Considerations</li><li>• Numerical Stability</li></ul>
Module 15 GPU as part of the PC Architecture	<ul style="list-style-type: none"><li>• GPU as part of the PC Architecture</li></ul>
Module 16 Efficient Host-Device Data Transfer	<ul style="list-style-type: none"><li>• Data Movement API vs. Unified Memory</li><li>• Pinned Host Memory</li><li>• Task Parallelism/CUDA Streams</li><li>• Overlapping Transfer with Computation</li></ul>
Module 17 Application Case Study: Advanced MRI Reconstruction	<ul style="list-style-type: none"><li>• Advanced MRI Reconstruction</li></ul>
Module 18 Application Case Study: Electrostatic Potential Calculation	<ul style="list-style-type: none"><li>• Electrostatic Potential Calculation (Part 1)</li><li>• Electrostatic Potential Calculation (part 2)</li></ul>

# Course Content

Module 19 Computational Thinking For Parallel Programming	<ul style="list-style-type: none"><li>• Computational Thinking for Parallel Programming</li></ul>
Module 20 Related Programming Models: MPI	<ul style="list-style-type: none"><li>• Joint MPI-CUDA Programming</li><li>• Joint MPI-CUDA Programming (Vector Addition - Main Function)</li><li>• Joint MPI-CUDA Programming (Message Passing and Barrier) (Data Server and Compute Processes)</li><li>• Joint MPI-CUDA Programming (Adding CUDA)</li><li>• Joint MPI-CUDA Programming (Halo Data Exchange)</li></ul>
Module 21 CUDA Python Using Numba	<ul style="list-style-type: none"><li>• CUDA Python using Numba</li></ul>
Module 22 Related Programming Models: OpenCL	<ul style="list-style-type: none"><li>• OpenCL Data Parallelism Model</li><li>• OpenCL Device Architecture</li><li>• OpenCL Host Code (Part 1)</li><li>• OpenCL Host Code (Part 2)</li></ul>
Module 23 Related Programming Models: OpenACC	<ul style="list-style-type: none"><li>• Introduction to OpenACC</li><li>• OpenACC Subtleties</li></ul>
Module 24 Related Programming Models: OpenGL	<ul style="list-style-type: none"><li>• OpenGL and CUDA Interoperability</li></ul>

# Course Content

Module 25 Dynamic Parallelism	<ul style="list-style-type: none"><li>• Effective use of Dynamic Parallelism</li><li>• Advanced Architectural Features: Hyper-Q</li></ul>
Module 26 Multi-GPU	<ul style="list-style-type: none"><li>• Multi-GPU</li></ul>
Module 27 Using CUDA Libraries	<ul style="list-style-type: none"><li>• Example Applications Using Libraries: CUBLAS</li><li>• Example Applications Using Libraries: CUFFT</li><li>• Example Applications Using Libraries: CUSOLVER</li></ul>
Module 28 Advanced Thrust	<ul style="list-style-type: none"><li>• Advanced Thrust</li></ul>
Module 29 Other GPU Development Platforms: QwickLABS	<ul style="list-style-type: none"><li>• Other GPU Development Platforms: QwickLABS</li></ul>
Where to Find Support	





# GPU Teaching Kit

Accelerated Computing



The GPU Teaching Kit is licensed by NVIDIA and the University of Illinois under the [Creative Commons Attribution-NonCommercial 4.0 International License](https://creativecommons.org/licenses/by-nc/4.0/).