



# Introduction to CUDA

Hadi Salimi and Nima Ghaemian Distributed Systems Lab. School of Computer Engineering, Iran University of Science and Technology, hsalimi@iust.ac.ir and nima@comp.iust.ac.ir

## Traditional CPUs



- Traditional CPUs:
  - Performance increases
  - Cost reductions
  - Brought GFLOPS to the desktop
  - Brought hundreds of GFLOPS to clusters
- This increase has slowed since 2003 due to power consumption issues that limited the increase of the clock frequency.
- Thus processor vendors have switched to multicore and many-core models.

# Traditionally Programs



- Traditionally, the vast majority of software applications are written as sequential programs.
- A sequential program will only run on one of the processor cores.
- To use the whole power of the cores of the new multi-core processors, parallel programs are needed.

#### **GPUs and CPUs Performance Gap**



# 1 teraflop on your desktop

In June 2008, NVIDIA introduced the GT200 chip, which delivers almost 1 teraflop (1,000 gigaflops) of single precision and almost 100 gigaflops of double precision performance.

#### GPGPU & CUDA

- What is GPGPU?
  - General-Purpose computing on a Graphics
    Processing Unit
  - Using graphic hardware for non-graphic computations
- What is CUDA?
  - Compute Unified Device Architecture
  - Software architecture for managing data-parallel programming

# CPU vs. GPU

- CPU
  - Fast caches
  - Branching adaptability
  - High performance
- GPU
  - Multiple ALUs
  - Fast onboard memory
  - High throughput on parallel tasks
    - Executes program on each fragment/vertex
- CPUs are great for task parallelism
- GPUs are great for data parallelism

# CPU vs. GPU (Cont.)



- The design of a CPU is optimized for sequential code performance.
- The general philosophy for GPU design is to optimize for the execution of massive number of threads.
- Graphics chips have been operating at approximately 10x the bandwidth of contemporaneously available CPU chips.



# Major problem



The major problem with traditional parallel processing applications:

Applications that can be run on a processor with a small market place will not have a large customer base!

The G80 family of CUDA-capable processors and its successors have shipped almost 100 million units to date.

#### IEEE Floating-Point Standard



- IEEE Floating-Point Standard makes it possible to have predictable results across processors from different vendors.
- Support for it was not widespread in early GPUs, but this has changed for the GeForce 8 series.
- The GPUs floating-point arithmetic units are primarily single precision today but we have already seen many applications where single precision floating is sufficient.

# Traditionally GPGPU



- Until 2006, graphics chips were very difficult to use because programmers had to use the equivalent of graphic API to access the processor cores, (by OpenGL or direct3D techniques)
- This technique was called GPGPU, for General Purpose Programming using a Graphics Processing Unit.

## CUDA



- But everything changed in 2007 with the release of CUDA.
- NVIDIA actually devoted silicon area to facilitate the ease of parallel programming, so this does not represent software changes alone; additional hardware was added to the chip.

#### Architecture of a CUDA-capable GPU







- SP: Streaming Processor
- Each SP has a multiply-add (MAD) unit, and an additional multiply (MUL) unit, all running at 1.35 GHz.
- SMs: highly threaded Streaming Multiprocessors
  - Each SM has 8 streaming processors (SPs)
- Building block: A pair of SMs





#### Speed up



- A good implementation on a GPU can achieve more than 100 times (100x) of speedup over a CPU.
- If the application includes "data parallelism," it's a simple task to achieve a 10x speedup with just a few hours of work.







*It depends on the portion of the application that can be parallelized.* 

## Example 1



- if 40% of the execution time is in the parallel portion, and we have a 50X speedup, how much speed-up can be expected from this application?
- Answer:

It will reduce the application execution to 60.8% that means 1.6X speedup.



# Web Resources



# CUDA ZONE: http://www.nvidia.com/object/cuda\_education.html

- David Kirk and Wen-mei W. Hwu, Lecture Notes of Programming Massively Parallel Processors, University of Illinois, Urbana-Champaign, 2009
- Rob Farber, "CUDA, Supercomputing for the Masses", *Dr. Dobb's Journal, can be found at:*

http://www.ddj.com/hpc-high-performancecomputing/

# Acknowledgment

Special Thanks to Mr. Ebrahim Khademi for his technical comments.

