

GPGPU Programming

Donato D'Ambrosio

Department of Mathematics and Computer Science
Cubo 30B, University of Calabria, Rende 87036, Italy

Lecturer's email: [donato.dambrosio \[at\] unical.it](mailto:donato.dambrosio@unical.it)

Lecturer's homepage: <http://www.mat.unical.it/~donato>

Course's homepage: <https://www.mat.unical.it/~donato/gpgpu.html>

Course team page: <https://bit.ly/3CU9xiR>

Code to join the team: **3d98qzr**

Academic Year 2021/22

Table of contents

- 1 Introduction
 - Course material
 - Syllabus
 - The Exam
 - The JPDM2 Workstation

Introduction

Introduction

About the Course

The course illustrates the fundamentals concepts and algorithms of **GPGPU programming** by **CUDA**.

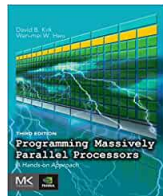


If possible, we will briefly overview OpenCL and MPI, the first representing an open alternative to CUDA, the latter the reference distributed memory programming model.



Course material

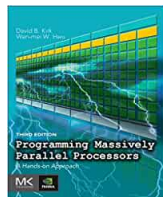
- David B. Kirk, Wen-mei W. Hwu, *Programming Massively Parallel Processors: A Hands-on Approach* (Third Edition). Morgan Kaufmann, 2017
- Online resources:
 - [CUDA Toolkit Documentation](#)
 - [Grid-stride loops](#)
 - [CUDA Unified Memory](#)
 - [CUDA GPU Compute Capability](#)
 - [Query Device and Handle Errors in CUDA](#)
 - [CUDA Occupancy Calculator](#)
 - [Unified L1/Texture Cache in Maxwell GPUs \(see Section 1.4.2.1\)](#)
 - [Constant Memory \(see Section 9.2.6\)](#)
 - [How to Implement Performance Metrics in CUDA](#)
- Further resources will be made available on the [course homepage](#) and on the [course team](#) during the semester.



Syllabus

Chapters 1, 2, 3 of the textbook

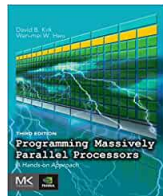
- Introduction and Development of GPU Computing
 - Architecture of GPUs
 - CPU-GPU Architectural Comparison
 - GPGPU and Heterogeneous Computing
- Data Parallelism and CUDA/OpenCL
 - Overview of CUDA (and OpenCL)
 - Programming model
 - A first example of application: Vector Addition
- Data-Parallel Execution Model
 - Thread Organization and Mapping to Data
 - A More Complex Example: Matrix Multiplication
 - Thread Synchronization
 - Thread Scheduling and Latency Tolerance



Syllabus

Chapters 4, 5, 7, 8, 11, 12, 18 of the textbook

- Memory
 - GPGPU Memory Model
 - Reducing Global Memory Traffic
 - A Tiled Matrix Multiplication Kernel
- Improving Performance
 - Warps and Thread Execution
 - Global Memory Bandwidth
 - Using Shared/Local Memory and Registers/Private Memory
- Example of Application
 - Stencil/Convolution
 - Prefix Sum/Merge sort/Graph search
- Programming a Heterogeneous Computing Cluster
 - Message Passing Interface (MPI)



Exam

- The assessment consists of two parts
 - A **written test** (one hour).
 - A **project** (based on CUDA/MPI).
- Assessment criteria of learning
 - The written test aims at verifying the basic knowledge of the course and allows to access to the next phase regarding the evaluation of the project.
 - The project aims to verify the ability to design and develop massive parallel software on GPUs.
- Criteria for measuring learning and assigning the final mark
 - The **partial mark** assigned to each test is between 0 and 30.
 - The **overall mark** is obtained as the average of the two partial outcomes.
 - The **laude** is given in case both tests are considered particularly deserving.

CUDA Workstation @ DeMaCS - UniCAL

- Dual quad-core Intel Xeon CPU E5440 @ 2.83GHz, 16GB of RAM
- Dual **GeForce GTX 980** NVIDIA GPU
 - Compute Capability: 5.2
 - Architecture: Maxwell (please, check-out the [compatibility guide](#))
 - Cores: 2048 (16 SMs, 128 cores per SM) @ 1126-1216 MHz
 - Memory: 4 GB GDDR5 @ 1750MHz
 - Memory Bandwidth 224 GB/s
 - 4.981 TFLOPS FP32 (float) performance
 - 155.6 GFLOPS FP64 (double) (1:32) performance



The JPDM2 Workstation

```
[user00@JPDM2 ~]$ screenfetch
      _
     .o+`
    `ooo/
   `+oooo:
  `+ooooooo:
 -+oooooooo+:
  `/:-:++oooo+:
 `/+////+////////:
  `/+++++////////+:
   `/+++oooooooooooooooo/`
  ./ooosssso++osssssso+`
 .oosssssso-`-`-`-/osssssso+`
 -osssssso.      :ssssssso.
 :osssssso/      ossso+++.
 /osssssssso/    +sssoooo/-
`/osssssso+/:--  -:/+osssso+-
`+sso+:-`        `.-/+oso:
`++:.            `.-/+//
`+`              `.-/+//

[user00@JPDM2 ~]$ nvcc --version
nvcc: NVIDIA (R) Cuda compiler driver
Copyright (c) 2005-2020 NVIDIA Corporation
Built on Wed_Jul_22_19:09:09_PDT_2020
Cuda compilation tools, release 11.0, V11.0.221
Build cuda_11.0_bu.TC445_37.28845127_0

[user00@JPDM2 ~]$
```

The JPDM2 Workstation

- Students can use the workstation to practice CUDA, develop the project for the exam, and for performance assessment.
- The following accounts were created for you students:
user01, user02, ..., use39. Disk quota is 5GB.
- Access the workstation as (\$ is the shell prompt):

```
$ ssh user{01|02|...|39}@160.97.63.93
```

Please, fill the **JPDM2 student accounts.xlsx** in the Files section of the channel **Lectures** to reserve your JPDM2 account.

There you will find the intial passowrd for your account that are invited to change at the first access. Do not annotate the new password in the shared document!

- The accounts will be (hopefully) **ready by October, 11th**.

The JPDM2 Workstation

- The access to the workstation has no limitations; more users can work at the same time.
- A student at a time can submit its program for execution to guarantya proper performance assessment. For this purpose, we use the **Slurm** workload manager. Time limit is fixed to **5 minutes**.



- Compile your program as usual, e.g.:

```
$ nvcc vecAdd.cu -o vecAdd
```

- write a script called `run.sh` or similar:

```
#!/bin/sh
```

```
srun vecAdd
```

- enqueue the job for execution:

```
$ sbatch run.sh
```

- Slurm quick start guide:

<https://slurm.schedmd.com/quickstart.html>