Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

University of Leeds

Lecture 14: Introduction to GPGPU p

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Previous lectures
Today's lecture

## Previous lectures

So far we have looked at CPU programming.

- are

- on

- Many **common** parallelism issues (sc
synchronisation; binary tree reduction).

- Also some **unique** to each type (locks a
shared memory; explicit communication for distributed
memory).

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Previous lectures
Today's lecture

# Today's lecture

Today's lecture is the first of 6 on programming **GPU**s (<u>G</u>raphics <u>Proc</u>      _

- _                                    eneral
  _

- GPU **devices** contain multiple **SIMD** units.

- Different **memory types**, some 'shar
  interpreted as 'distributed.'

- Programmable using a variety of C/C++-based languages, notably **OpenCL** and **CUDA**.

Overview
**Anatomy of a GPU**
General purpose programming on a GPU
Summary and next lecture

**Development of GPUs**
Overview of GPU architecture
Books

# Development of GPUs[1]

Early accelerators were driven by graphical operating systems and high-end applications (defense, science and engineering *etc.*).

-

- https://eduassistpro.github.i

Consumer applications employing 3D domin

- First person shooters in mid-90s (Doom, Q
- Add WeChat edu_assist_pr
- 3D graphics accelerators by Nvidia, ATI
- Initially as external **graphics cards**.

---

[1]Sanders and Kandrot, *CUDA By Example* (Addison-Wesley, 2011).

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Development of GPUs
Overview of GPU architecture
Books

# Programmable GPUs

The first **programmable** graphics cards were Nvidia's GeForce3 series (2001).

- **e** vertex

- 

Early **general purpose** applications 'dis **graphical**.

1. Input data converted to pixel **colours**.
2. Pixel shaders performed calculations on this data.
3. Final 'colours' converted back to **numerical data**.

Overview
**Anatomy of a GPU**
General purpose programming on a GPU
Summary and next lecture

**Development of GPUs**
Overview of GPU architecture
Books

# GPGPUs

In 2006 Nvidia released its first GPU with CUDA.

- General calculations *without* converting to/from colours.

Now h

- https://eduassistpro.github.i

-

- Vendors include Nvidia, AMD and Intel.

Originally designed for **data parallel** g Add WeChat edu_assist_pr

- Increasing use of GPUs for *e.g.* **machine learning**[1] and **cryptocurrencies**.

---

[1]Now also have **neural processing units** (NPUs) for machine learning.

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Development of GPUs
Overview of GPU architecture
Books

## Overview of GPU architectures

Design and terminology of GPU hardware differs between vendors.

- Nvidia *different to* AMD *different to* Intel *different to* . . .

Typi

-

-

SIMD processors contain **SIMD functio**                                :

- Each SIMD core contains multiple
- Executes **the same** instruction on m

**Hierarchy:**

Threads $\in$ SIMD Cores $\in$ SIMD Processors $\in$ GPU

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Development of GPUs
Overview of GPU architecture
Books

## SIMD processor

A typical SIMD processor has:

- A **thread scheduler**.

- Multiple **SIMD function**

- 

Not shown but usually present:

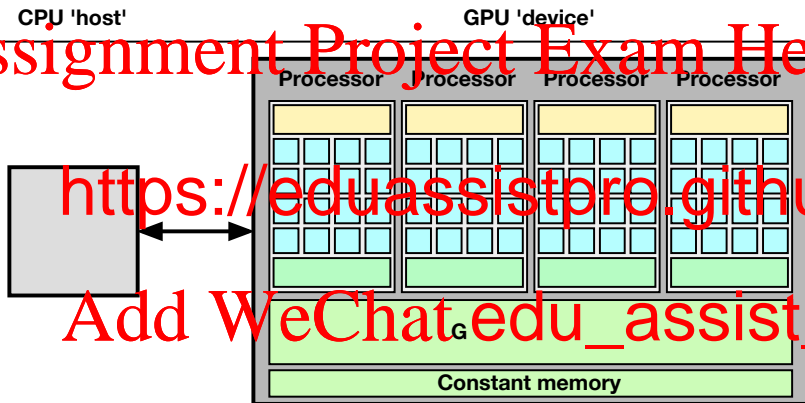- Registers, special floating point units, . . .



> **Note:**
>
> Thread scheduling is performed **in hardware**.

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Development of GPUs
Overview of GPU architecture
Books

# CPU with a single GPU

**CPU 'host'**　　　　　　　　　　　**GPU 'device'**



- The **data bus** between CPU and GPU is **very slow**.
- Faster for **integrated GPU**s.

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Development of GPUs
Overview of GPU architecture
Books

## SIMD *versus* SIMT

Assignment Project Exam Help

Nvidia refer to their architectures as **SIMT** rather than SIMD.

- <u>S</u>ingle <u>I</u>nstruction <u>M</u>ultiple **Threads**.

- 

https://eduassistpro.github.i

- 
  **simultaneously**.

- Therefore 'in between' SIMD and MIMD.

Add WeChat edu_assist_pr

Will look at this more closely in Lecture 17, where we will see how it can be detrimental to performance.

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Development of GPUs
Overview of GPU architecture
Books

# Books

McCool et al. [lecture 1] includes some OpenCL, but does not address GPUs specifically. Books for GPU programming include:

- li,

  https://eduassistpro.github.i

- **CUDA by example**, *Sanders and K* ey 2011).
  - Slightly old, but a gentle introduction.
  - Only considers CUDA, whereas we will use OpenCL, but may still be useful.

You do not need any of these books for this module!

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Available languages
Installing and building OpenCL
Platforms, devices and contexts
'Hello world'

# GPU programming languages 1. CUDA

Assignment Project Exam Help

The first language for GPGPU programming was Nvidia's **CUDA**.

https://eduassistpro.github.i

- **Only works on CUDA**-enabled device .

Add WeChat edu_assist_pr

As the first GPGPU language it has much docume

Therefore we will reference CUDA concepts and terminology quite frequently, often in footnotes.

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Available languages
Installing and building OpenCL
Platforms, devices and contexts
'Hello world'

# GPU programming languages 2. OpenCL

Assignment Project Exam Help

- Currently the main alternative to CUDA is **OpenCL** (2008).
  - Stands for **Open Computing Language**.
  -
    rs
    https://eduassistpro.github.i
  -
  - Can also run on CPUs, FPGAs (=F _ _ e
    Arrays...
    Add WeChat edu_assist_pr
  - C/C++ based.
  - Similar programming model to CUDA.
  - OpenCL 3.0 released Sept. 2020.

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Available languages
Installing and building OpenCL
Platforms, devices and contexts
'Hello world'

## Directive based programming abstractions

**OpenACC** (2011):

- Open ACCelerator, originally intended for **accelerators**.
-
- https://eduassistpro.github.i

**OpenMP**:

- GPU support from version 4.0 onwards.
- Usual #pragma omp directives, wit

Both give **portable** code, but both require some understanding of the hierarchical nature of GPU hardware to produce reasonable performance.

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Available languages
Installing and building OpenCL
Platforms, devices and contexts
'Hello world'

## Installing OpenCL

Already installed on cloud-hpc1.leeds.ac.uk (and most Macs).

Othe
archi

Nvidia: `https://developer.nvidia.com`

Intel: `https://software.intel.com/en-`

AMD: `https://www.amd.com/en` and search for OpenCL.

Assignment Project Exam Help

https://eduassistpro.github.i

Add WeChat edu_assist_pr

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Available languages
Installing and building OpenCL
Platforms, devices and contexts
'Hello world'

## OpenCL header file

All OpenCL programs need to include a header file

Sinc
UNI
have t

```
1  #ifde   _
2  #include <OpenCL/opencl.h>
3  #else
4  #include <CL/cl.h>
5  #endif
```

Note that the coursework will be marked on a system similar to
`cloud-hpc1.leeds.ac.uk`, so it **must** run on that system.

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Available languages
Installing and building OpenCL
Platforms, devices and contexts
'Hello world'

## Compiling and running

We use the CUDA nvcc compiler on cloud-hpc1 leeds.a....uk:

```
1  nvcc -lOpenCL -o <executable> <source>.c
```

Note t

**Exe**

To execute on a GPU it will be necessary to use the batc
(see next slide). However, it is also possible to run an O
code on the login nodes CPU by launching as any
executable:

```
1  ./<executable> [any command line arguments]
```

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Available languages
Installing and building OpenCL
Platforms, devices and contexts
'Hello world'

## Running on GPU via batch jobs

The batch node of cloud-hpc1.leeds.ac.uk may be configured with a Tesla T4 GPU.

Henc
follo

- https://eduassistpro.github.i

- Create a job submission script as outlined below;

- Submit ot the batch queue using

Here is a typcal batch script to run 'g

```
#!/bin/bash
#SBATCH --partition=gpu --gres=gpu:t4:1
./gpu-example
```

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Available languages
Installing and building OpenCL
Platforms, devices and contexts
'Hello world'

## Compiling and running: Macs

**Compiling:**
Use the OpenCL framework:

```
1  gcc -Wal
```

- https://eduassistpro.github.i

- If you see deprecation **errors**, try
gcc.

Add WeChat edu_assist_pr

**Executing:**
Launch as any normal executable

```
1  ./<executable> [any command line arguments]
```

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Available languages
Installing and building OpenCL
Platforms, devices and contexts
'Hello world'

## Platforms, devices and contexts

Since OpenCL runs on many different devices by many different vendors, it can be quite laborious to initialise.

Need

| Platform | ... |
|----------|-----|
| **Device** | Belongs to a platform; may b |

Need to **initialise**:

| Context | Coordinates inter |
| | **device** (*e.g.* a GPU). One per device. |
| **Command queue** | To request action by a device. Normally one per device, but can have more *[Lecture 19]*. |

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Available languages
Installing and building OpenCL
Platforms, devices and contexts
'Hello world'

## Initialisation code

Most code for this module will come with helper.h which
contains two useful routines:

simp

- https://eduassistpro.github.i

compileKernelFromFile()

- Compiles an OpenCL **kernel** to be e
  Will cover this next lecture.

Add WeChat edu_assist_pr

You don't need to understand how these routines work, but are
welcome to take a look.

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Available languages
Installing and building OpenCL
Platforms, devices and contexts
'Hello world'

# Using `simpleOpenContext_GPU()`

```
1  #include "helper.h"    // Also includes OpenCL.
2  int main() {
3    // Get context and device for a GPU.
4    cl_device_id device;
5    cl
6
7    // Op
8    cl
9    cl_command_queue queue = clCreateCommandQueue(
       context,device_0,&status);
10
11   ...  // Use the GPU through 'queue'.
12
13   // At end of program.
14   clReleaseCommandQueue(queue);
15   clReleaseContext(context);
16 }
```

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Available languages
Installing and building OpenCL
Platforms, devices and contexts
'Hello world'

# 'Hello world' in OpenCL
## Code on Minerva: `displayDevices.c`

Assignment Project Exam Help

Since most GPU's cannot print in the normal sense, there is no simple 'Hello World' program.

Inste https://eduassistpro.github.i
help

- Loops through **all** platforms and dev
- Lists all **OpenCL-compatible** de Add WeChat edu_assist_pr
- Also a list of extensions; *e.g.* `cl_ _` supports double precision floating point arithmetic.
- In the output, a **compute unit** is a SIMD processor or streaming multiprocessor.

Overview
Anatomy of a GPU
General purpose programming on a GPU
Summary and next lecture

Summary and next lecture

## Summary and next lecture

Assignment Project Exam Help

Today we have started looking at GPU programming:

- Overview of GPU architectures.
-
- https://eduassistpro.github.i
-

using the functions:

- clGetPlatformIDs
- clGetDeviceIDs

Add WeChat edu_assist_pr

Next time we will implement a "real" program in OpenCL: **vector addition**.