Overview
Vector addition on a GPU
Work items and work groups
Summary and next lecture

Assignment Project Exam Help

https://eduassistpro.github.i

University of Leeds

Add WeChat edu_assist_pr

Lecture 15: GPU threads and k

Overview
Vector addition on a GPU
Work items and work groups
Summary and next lecture

Previous lecture
Today's lecture

# Previous lecture

In the last lecture we started looking at General Purpose GPU prog

- https://eduassistpro.github.i

- Thread scheduling is performed **in hardware**.

- Programmable using **OpenCL** (th
  others

- Device discovery performed at run time (*cf.* the
  `displayDevices.c` example).

Assignment Project Exam Help

Add WeChat edu_assist_pr

Overview
Vector addition on a GPU
Work items and work groups
Summary and next lecture

Previous lecture
Today's lecture

## Today's lecture

Assignment Project Exam Help

Today we will see how to perform vector addition on a GPU:

- https://eduassistpro.github.i

- **Work items** are the basic unit of concurr
- Arranged into **work groups** for **s** Add WeChat edu_assist_pr
- How to set the work group size.

Overview
Vector addition on a GPU
Work items and work groups
Summary and next lecture

Communication between host and device
GPU kernels
Copying data between device and host

## Vector addition
Code on Minerva: `vectorAddition.c`, `vectorAddition.cl` and `helper.h`

Assignment Project Exam Help

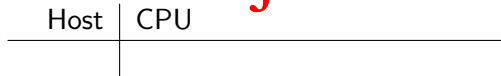Once again use **vector addition** as our first worked example:

https://eduassistpro.github.i

In seri

```
1  for ( i=0; i<N; i++ )
2    c[i] = a[i] + b[i];
```

Add WeChat edu_assist_pr

where vectors **a**, **b** and **c** all have *N* el
mathematical and computer indexing differ by o

This is a **map**/**data parallel** problem with no **data dependencies**.

Overview
**Vector addition on a GPU**
Work items and work groups
Summary and next lecture

Communication between host and device
GPU kernels
Copying data between device and host

## Host and device

The CPU is the **host**, and the GPU is the **device**:

| Host | CPU |
|------|-----|

Assu

If the initial data is only accessible to the CPU, must          to
the GPU to perform the calculations, then
to the CPU.

- This requires **explicit communication**, somewhat similar to
  the distributed memory model.

---

[1]Some modern GPUs support **unified memory** — see next lecture.

Overview
**Vector addition on a GPU**
Work items and work groups
Summary and next lecture

Communication between host and device
GPU kernels
Copying data between device and host

## Typical program structure

Assignment Project Exam Help

1. **Send** problem

   t https://eduassistpro.github.i

2. calculations on
   the device.

3. **Return** results
   from device to
   the host.



(CPU)

evice

https://eduassistpro.github.i

Add WeChat edu_assist_pr

Overview
**Vector addition on a GPU**
Work items and work groups
Summary and next lecture

Communication between host and device
GPU kernels
Copying data between device and host

## Contexts and command queues

Recall from last lecture that to use OpenCL we first need to:

1. Identify the **platform** and a suitable **device**.

2. [text obscured] **eue**.

The [text obscured]

```
1  cl_device_id device;
2  cl_context context = simpleOpenContext_GPU(&device);
3
4  cl_int status;
5  cl_command_queue queue = clCreateCommandQueue(context,
      device,0,&status);
6  ...     // Use the GPU.
7  clReleaseCommandQueue(queue);
8  clReleaseContext(context);
```

Overview
Vector addition on a GPU
Work items and work groups
Summary and next lecture

Communication between host and device
GPU kernels
Copying data between device and host

# Device memory allocation

Suppose arrays a, b and c initialised on the **host**:

```
1  float *host_a = (float*) malloc(N*sizeof(float));
```

Simil

Can a
**host a**

```
1  cl_mem device_a = clCreateBuffer(
2    context,
3    CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,  // Flags
4    N*sizeof(float),           // Size in bytes.
5    host_a,                    // Copy from this host array.
6    &status,                   // Error status.
7  );
```

Similar for device_b, device_c.

Overview
Vector addition on a GPU
Work items and work groups
Summary and next lecture

Communication between host and device
GPU kernels
Copying data between device and host

# clCreateBuffer() usage

Assignment Project Exam Help

- The context has been initialised for the GPU.
- The flag CL_MEM_READ_ONLY refers to how the **device** accesses

https://eduassistpro.github.i

- The flag CL_MEM_COPY_HOST_PTR
  existing **host** array (the $4^{th}$ argu

Add WeChat edu_assist_pr

- For device_c, where no host data (yet) ex
  CL_MEM_WRITE_ONLY and the $4^{th}$ argument is NULL.

- status is set to CL_SUCCESS if the operation was successful,
  otherwise some other error code.

Overview
Vector addition on a GPU
Work items and work groups
Summary and next lecture

Communication between host and device
GPU kernels
Copying data between device and host

## GPU kernel

Assignment Project Exam Help

> **Definition**
>
> **Kernels** are functions that execute on the device.

https://eduassistpro.github.i

Use st

Add WeChat edu_assist_pr

```
1  __kernel
2  void vectorAdd(__global float *a, __global float *b,
       __global float *c)
3  {
4    int gid = get_global_id(0);
5    c[gid] = a[gid] + b[gid];
6  }
```

Overview
Vector addition on a GPU
Work items and work groups
Summary and next lecture

Communication between host and device
GPU kernels
Copying data between device and host

## OpenCL kernels

- All kernels¹ preceded with __kernel.

- 

- __

- get_global_id() returns the (globa
  - For this problem it is the index of the vector.
  - See later.

---
¹CUDA kernels are preceded __global__ (if they are callable by the host).

Overview
**Vector addition on a GPU**
Work items and work groups
Summary and next lecture

Communication between host and device
**GPU kernels**
Copying data between device and host

## Building a kernel

Assignment Project Exam Help

- OpenCL kernels are compiled **at run time** (of the C code).
  - Allows optimisation for the device that executes it.

Requ https://eduassistpro.github.i

1. Start with the program as a char* **string** (typically read from file ending in .cl).

2. Create the **program** for the context v edu_assist_pr clCreateProgramWithSource()

3. **Build** (compile and link) using clBuildProgram().

4. Create a **kernel** using clCreateKernel().

Overview
Vector addition on a GPU
Work items and work groups
Summary and next lecture

Communication between host and device
GPU kernels
Copying data between device and host

## Building a kernel with `helper.h`

To simplify this process, the file `helper.h` contains the routine `compileKernelFromFile()`.

For th

```
1  cl_k
2    "v
3    "vectorAdd",        // Name of function.
4    context,            // Same as before.
5    device              // Same as before.
6  );
7
8  ... // Use kernel.
9
10 clReleaseKernel(kernel);
```

It also includes some basic error handling.

Overview
Vector addition on a GPU
Work items and work groups
Summary and next lecture

Communication between host and device
GPU kernels
Copying data between device and host

## Setting kernel arguments

Each kernel argument must be set by using clSetKernelArg().

```
1  statu  = c
2    ke
3    0,              // The argument number.
4    sizeof(cl_mem), // The size of the argument.
5    &device_a       // The value.
6  );
```

This is repeated for argument 1 ($\rightarrow$device_b) and argument 2 ($\rightarrow$device_c) for the vector addition example.

Overview
**Vector addition on a GPU**
Work items and work groups
Summary and next lecture

Communication between host and device
**GPU kernels**
Copying data between device and host

# Starting a kernel in OpenCL[1]

To start a kernel, you place it on the **command queue** using
`clEnqueueNDRangeKernel()`:

```
1  // Will cover this later.
2  size_  i
3  index
4  work
5
6  // Place the kernel onto the command queue.
7  status = clEnqueueNDRangeKernel(queue,kernel,1,NULL
        indexSpaceSize,workGroupSize,0,NULL,NULL);
```

There are many arguments; we will cover some lat

Note that `size_t` is an **unsigned integer**.

---

[1]In CUDA: `kernel<<<workGroupSize,indexSpaceSize>>>(...)`.

Overview
Vector addition on a GPU
Work items and work groups
Summary and next lecture

Communication between host and device
GPU kernels
Copying data between device and host

# Copying data between device and host[1]

To get the result (device_c) back to the host (host_c), enqueue a **read buffer** command:

```
1  status = clEnqueueReadBuffer(
2    qu                // The command queue
3    de                // Devic
4    CL                // Bloc
5    0,                // Offset; must be zero.
6    N*sizeof(float),  // Data size.
7    host_c,           // Host memory.
8    0, NULL, NULL     // Event
9  );
```

Note this is a **blocking** communication call - **it will not return until the copy has finished** — like MPI_Send()/MPI_Recv().

---

[1]In CUDA: cudaMemcpy(...,cudaMemcpyDeviceToHost).

Overview
**Vector addition on a GPU**
Work items and work groups
Summary and next lecture

Communication between host and device
GPU kernels
**Copying data between device and host**

# Copying data from host to device[1]

If we had not used CL_MEM_COPY_HOST_PTR earlier, we would need two calls to clEnqueueWriteBuffer():

```
1  statu  = c

2  stat    = c
```

- Copies **from** host **to** device.
- CL_FALSE used for **non-blocking**
- The device memory **always** come
  argument list.

---

[1]In CUDA: cudaMemcpy(...,cudaMemcpyHostToDevice).

Overview
Vector addition on a GPU
**Work items and work groups**
Summary and next lecture

**Work item hierarchy**
Specifying work groups and NDRange
What group size to use?

# Work items

**Definition**

The **work item** is the unit of concurrent execution. It usually

https://eduassistpro.github.i

As thr

is (essentially) **no overhead** in launchin

- No problem **oversubscribing**, i.e.
  there are physical cores.

**Normally issue as many threads as the problem requires**.

Overview
Vector addition on a GPU
**Work items and work groups**
Summary and next lecture

**Work item hierarchy**
Specifying work groups and NDRange
What group size to use?

# Work item hierarchy

To remain scalable, the hardware does not allow communication (including synchronisation) between *all* threads at once.

Inste

- ...
- Communication (including synchronisation) only possible **within** a work group.
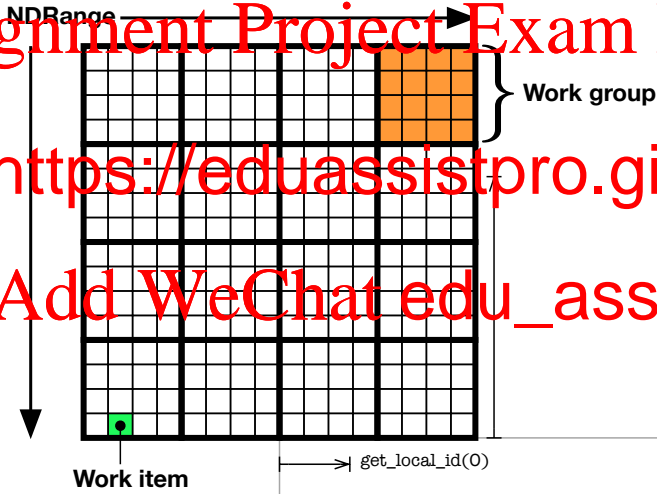
The full range of all threads is called *n*-dimensional range[2].

---

[1] *Threads* and *thread blocks* in CUDA.
[2] *Grid* in CUDA.

Overview
Vector addition on a GPU
**Work items and work groups**
Summary and next lecture

**Work item hierarchy**
Specifying work groups and NDRange
What group size to use?

# Hierarchy of work items: 2D example



**NDRange**

**Work group**

**Work item**

get_local_id(0)

Overview
Vector addition on a GPU
Work items and work groups
Summary and next lecture

Work item hierarchy
Specifying work groups and NDRange
What group size to use?

# Specifying the *n*-dimensional range NDRange

The NDRange must be 1, 2 or 3 dimensions.

A 2-dimensional example:

```
1  size_    g
2  size    w
3
4  statu  = c
      globalSize , workGroupSize ,0 , NULL , NULL );
```

- Launches X*Y kernels in total (one per w
- In work groups of 8*16.

OpenCL 2.0 allows X and Y to be arbitrary, but in earlier versions they must be multiples of the work group size (8 and 16 here).

Overview
Vector addition on a GPU
**Work items and work groups**
Summary and next lecture

Work item hierarchy
**Specifying work groups and NDRange**
What group size to use?

Once in a kernel, can get the **global** indices using
get_global_id(). For this 2D example:

```
1  get_global_id(0);  // Varies from 0 to X-1 inc.
2  get_global_id(1);  // Varies from 0 to Y-1 inc.
```

Similarly
get_

```
1  get_local_id(0);  // Varies from 0 to 7 inc.
2  get_local_id(1);  // Varies from 0 to 15 inc.
```

Can also get the number of work items in a group or in the
NDRange using get_local_size() and get_global_size():

```
1  get_local_size (1);  // Returns 16.
2  get_global_size(0);  // Returns X.
```

Overview
Vector addition on a GPU
Work items and work groups
Summary and next lecture

Work item hierarchy
Specifying work groups and NDRange
What group size to use?

# What group size to use?

Devices have a **maximum work group size** they can support.
This can be determined at run time as follows:

```
1  size_t maxWorkItems;
2  clGet
```

Note

Other factors may suggest using work group size l
maximum

- We will look at one of these next time.

Passing `NULL` as the work group argument lets OpenCL try to
determine a suitable size **automatically**.

Overview
Vector addition on a GPU
Work items and work groups
Summary and next lecture

Summary and next lecture

## Summary and next lecture

Today we have looked at a complete GPGPU solution:

- 
- 
- 
- Group into **work groups**, within whi
  possible.

Next time we will look at the different memory types on a GPU.