

# Assignment Project Exam Help

<https://eduassistpro.github.io>

Add WeChat University of Leeds edu\_assist\_pr

Lecture 16: GPU memory t

## Previous lectures

# Assignment Project Exam Help

In the last lecture we saw how to program a simple GPU application:

- 
- 
- 

<https://eduassistpro.github.io>

computation.

- Performed by many **work items**
- Arranged into **work groups** of parallel computation.
  - Can arrange work items in 1, 2 or 3 dimensions (=  $n$ -dimensional range).
- **Copy** the result back from the device to the host.

## Today's lecture

# Assignment Project Exam Help

In the vector addition kernel, all of the arguments had the prefix `__global__`:

```
1 __kernel
2 void ker
3 {
4     // ke
5 }
```

<https://eduassistpro.github.io>

Today we will see what this means

- The different **memory types** available
- How and when to use them.
- Performance issues related to **register overflow**.

Add WeChat [edu\\_assist\\_pro](#)

## GPU memory

GPUs are designed as high throughput devices.

- Many threads (100's, 1000's, ...) execute simultaneously.
- ... s, i.e.

<https://eduassistpro.github.io>

To maximise throughput, GPUs have multiple

- Architecture varies greatly between
- Performance would ideally be optimised for which the code may be deployed.

---

<sup>1</sup>Although modern GPUs increasingly also have memory caches.

## Shared virtual, or ‘unified’, memory

In this module we treat CPU and GPU memory as separate

Assignment Project Exam Help

- Typical of early GPU architectures.

Increase  
program



- Integrated GPUs may share **physi**

- Supported from OpenCL 2.0 (CUDA 4.0)

Add WeChat edu\_assist\_pr

We will not consider unified memory in this module.

---

<sup>1</sup>See e.g. Wilt, *The CUDA Handbook* (Addison-Wesley, 2013); Han and Sharma, *Learn CUDA Programming* (Packt, 2019).

## Memory coalescing

When copying from e.g. global to local memory, **adjacent** threads will often access **adjacent** memory locations.

GPUs detect and optimise for this by **memory coalescing**:

- <https://eduassistpro.github.io>
- nearby memory locations wherever possible
- For 2D/3D data can use a technique known as **coalescing**

You are not expected to optimise your code for memory coalescence in this module.

---

<sup>1</sup>Rauber and Rünger, *Parallel Programming* (Springer, 2012). Tiling is also used to optimise cache access in CPUs.

## Memory types

GPUs typically contain 4 different types of memory:

**Global** Accessible to **all** work items in **all** work groups.

**L**

**Pr**

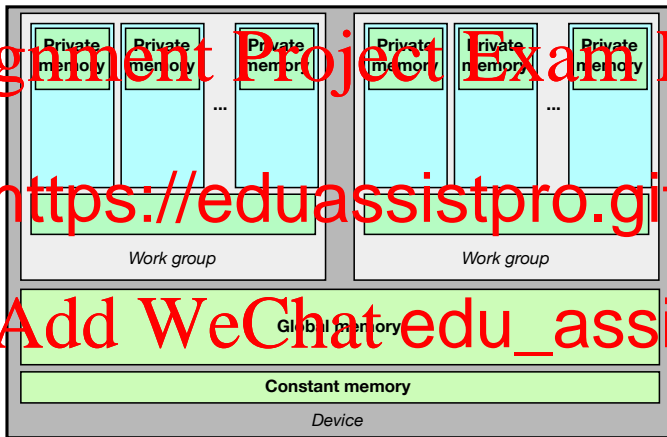
**Constant** Global memory optimised for **constant** work items. Not available on all GPU

These are **disjoint** - it is **not** allowed to c  
another.

---

<sup>1</sup>In CUDA and Nvidia devices, *local* memory is referred to as *shared*.

# GPU memory types<sup>1</sup>



<sup>1</sup>After Kaeli *et al.*, *Heterogeneous computing with OpenCL 2.0* (Morgan-Kaufman, 2015).



## Analogy with CPU memory

Notice there is a (loose) analogy with CPU memory types:

CPU	GPU	Similarity
Shared memory	Shared memory	(CPU); <b>work group</b> (GPU).
Private memory	Private memory	(GPU); <b>process</b>

To send data between work groups, must use global memory (the host); a form of **communication**.

Cannot **directly** send data between the local memory of different work groups.

- The analogy with distributed memory CPUs breaks down.

## Allocating device memory

# Assignment Project Exam Help

When allocating buffer memory on a device, it is placed in **global** memory.

For example, to allocate a buffer of size `N` for a device with capacity `capa`:

```
1 cl_mem device_array= clCreateBuffer(context,  
    CL_MEM_READ_ONLY, N*sizeof(float),...);
```

Add WeChat edu\_assist\_pr

Note that `CL_MEM_READ_ONLY` does **not** make it constant memory.

- Still allocated in the device's **global** memory.

## Read and write buffer flags

For OpenGL the buffer type should be specified as read-only, write-only, or read-write, by using one of the following flags:

CL\_...  
CL\_...  
CL\_MEM\_READ\_WRITE Both read and write allowed.

- Hint to allow **optimisations** by the r...
- The default is CL\_MEM\_READ\_WRITE.
- These refer to the **device** accessibility, *i.e.* **from inside kernels**, *not* from the host.

## Memory type 1. Global memory

Assignment Project Exam Help

Global memory

wer

com

<https://eduassistpro.github.io>

This is the memory type we have used already

- Convenient from a programming perspective
- Generally poor performance, although typically still faster than host-device communication.

## Using global memory in OpenCL<sup>1</sup>

Allocate global memory using `clCreateBuffer()`:

```
1 cl_mem device_buffer = clCreateBuffer(context, flags,  
    size, ...);
```

In the kernel:

```
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 }
```

---

<sup>1</sup>In CUDA: `clCreateBuffer()` → `cudaMalloc()`; no `__global` specifier.

## Memory type 2. Local memory

Assignment Project Exam Help

Local memory

A

tween

<https://eduassistpro.github.io>

Typically used as a **scratch space** for c  
than one work item in a group.

- Use for intermediate calculations.
- Place final answer in global memory.

In practice, this also requires **synchronisation** which is next lecture's topic, so will see an example of local memory then.

## Static local arrays

To use local memory in a kernel, simply place `--local` before the variable<sup>1</sup>.

```
1 __kernel
2 void ker
3 {
4     --
5     // Ca
6     ...
7     // Place final answer in device_array.
8 }
```

However, this is **static allocation** - the size is known **when the kernel is built**.

---

<sup>1</sup>In CUDA: `--local` → `--shared--`.

## Dynamic local arrays

Assignment Project Exam Help

To create dynamic local arrays, declare it as a kernel argument with the specifier `__local`:

```
1 __ker
2 void ker
3 {
4     // Ca
5     ...
6     // Place final answer in device_a.
7 }
```

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pro

Then when setting kernel arguments, specify the size **but set the pointer to NULL**:

```
1 clSetKernelArg(kernel, 0, N*sizeof(int), NULL);
```



## Memory type 3. Private memory

# Assignment Project Exam Help

Private memory

<https://eduassistpro.github.io>

In practice

hardware as **registers**:

- Small amount of memory that can be accessed
- Faster access than even local memory.
- **Automatic storage duration**, *i.e.* deallocated at the end of the kernel (or enclosing code block).

## Using private memory

# Assignment Project Exam Help

There is no specifier for private memory (i.e. no `__private`).



For in

<https://eduassistpro.github.io>

```
1 __kernel
2 void kernel(__global float *array)
3 {
4     int gid = get_global_id(0);
5     ...
6 }
```

Add WeChat edu\_assist\_pro

...the variable `gid` is treated as private memory.

## Private kernel arguments

# Assignment Project Exam Help

Kernel arguments without a specifier are also treated as private.

In this e

```
1 __kernel
2 void kernel
3 {
4     // Calculations involving N.
5     ...
6 }
```

<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pro

The corresponding call to `setKernelArg()` would be:

```
1 int N=...;
2 clSetKernelArg(kernel, 1, sizeof(int), &N);
```

## Register overflow

Code on Minerva: `registerOverflow.c`, `registerOverflow.cl`, `helper.h`

# Assignment Project Exam Help

Devices have a **fixed** amount of register memory. What happens if this is exceeded is device-dependent, but is usually one of:

①

②

<https://eduassistpro.github.io>

Either mechanism reduces performance.

Add WeChat `edu_assist_pr`

### Guidance

Kernels should be **small** (in the sense of low register usage) to limit the risk of register overflow.

## Memory type 4. Constant memory

### Constant memory

Accessible by all work items and work groups, but read only (by a kernel)

GPU  
read fr

- Known as **constant** memory.
- Can still be written to by the host.
- Originally to accelerate the mapping of ('texture' memory<sup>1</sup>).
- Typically much smaller than global memory, even if it exists.

---

<sup>1</sup>CUDA treats texture memory separately to other constant memory.

## Using constant memory<sup>1</sup>

For kernel arguments, use the `__constant` specifier:

```
1 __kernel
2 void ker
3 { ... }
```

Initial  
argu

Variables within kernels can also be `__`  
specified at compile time:

```
1 __constant float pi = 3.1415926;
```

---

<sup>1</sup>In CUDA: Device data specified `__constant__`, and copy from host using `cudaMemcpyToSymbol()`.

## Summary and next lecture

# Assignment Project Exam Help

Today we have look at the different **memory types** in GPUs:



<https://eduassistpro.github.io>

Add WeChat edu\_assist\_pro

The main use of **local** memory is to coordi between work items in a group. We will see an example next time when we look at **synchronisation** in GPUs.