

Assignment Project Exam Help

XJCO3221 Parallel Computation

<https://powcoder.com>

Peter Jimack

University of Leeds
Add WeChat powcoder

Lecture 16: GPU memory types

Previous lectures

Assignment Project Exam Help

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

- **Allocate** memory on the **device** (the GPU).
- **Copy** host (CPU) data to the device.
- Build and execute a **kernel** on the device that performs the computation.
 - Performed by many **work items** (*threads*).
 - Arranged into **work groups** of programmable size.
 - Can arrange work items in 1, 2 or 3 dimensions; the **NDRange** (= *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 kernel(__global float *a, __global float *b, ...)
3 {
4     // Kernel body.
5 }
```

Today we will see what this means:

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

GPU memory

GPUs are designed as **high throughput** devices.

- **Many** threads (100's, 1000's, ...) execute simultaneously.
- By contrast, CPUs are **latency reducing** architectures, *i.e.* fast memory access by use of caches¹, instruction level parallelism [cf. Lectures 1, 2], etc.

To maximise throughput, GPUs have multiple **memory types**:

- Architecture varies greatly between, and within, vendors.
- Performance would ideally be optimised for **each** GPU on which the code may be deployed.

¹Although modern GPUs increasingly also have memory caches.

Shared virtual, or 'unified', memory

In this module we treat CPU and GPU memory as **separate**

- Typical of early GPU architectures.

Increasingly, CPU and GPU memory are presented to the programmer as **unified**¹

- API decides whether CPU or GPU memory is allocated.
- Integrated GPUs may share **physical** memory with the CPU.
- Supported from OpenCL 2.0 (CUDA 4.0).

We will not consider unified memory in this module.

¹See 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**:

- **Coalesces** multiple small memory accesses into a single large memory access — much faster.
- To exploit this, programmers can ensure nearby threads access nearby memory locations wherever possible.
- For 2D/3D data can use a technique known as **tiling**¹.

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

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

Memory types

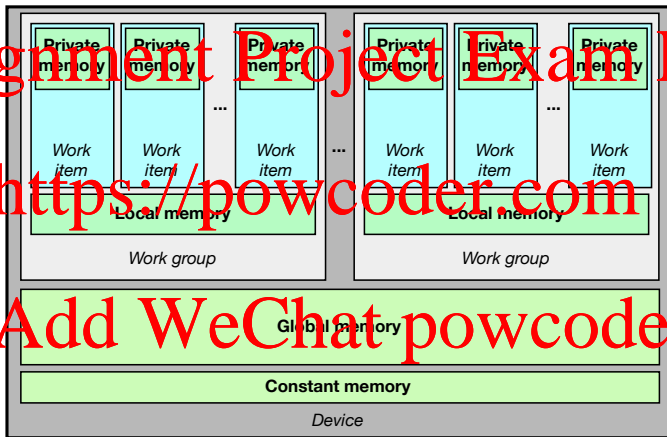
GPUs typically contain 4 different types of memory.

- Global** Accessible to **all** work items in **all** work groups.
- Local**¹ Shared by work items within **one work group only**.
- Private** Accessible to a **single work item only**.
- Constant** Global memory optimised for **read-only** operations. Not available on all GPUs.

These are **disjoint** - it is **not** allowed to cast one address space to another.

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

GPU memory types¹



¹After Kaeli *et al.*, *Heterogeneous computing with OpenCL 2.0* (Morgan-Kauffman, 2015).

Analogy with CPU memory

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

CPU	GPU	Similarity
Shared memory	Global memory	Accessible by all processing units [threads (CPU); work items (GPU)]
Distributed memory	Local memory	Only accessible to one process (CPU); work group (GPU).

To send data between work groups, must use global memory (or 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 global buffer (called device array) capable of storing N floats:

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

Add WeChat powcoder

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_MEM_READ_ONLY` The buffer should only be read from.
`CL_MEM_WRITE_ONLY` The buffer should only be written to.
`CL_MEM_READ_WRITE` Both read and write allowed.

- Hint to allow **optimisations** by the runtime system.
- 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

Accessible to all processing units, but bandwidth is **slower** compared to the other memory types. Typically cached in modern GPUs, but not in older devices.

<https://powcoder.com>

This is the memory type we have used already (cf. Lecture 15).

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

Add WeChat powcoder

Using global memory in OpenCL¹

Allocate global memory using `clCreateBuffer()`:

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

In the kernel, prepend `__global` before the pointer(s):

```
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 }
```

¹In CUDA: `clCreateBuffer()` → `cudaMalloc()`; no `__global` specifier.

Memory type 2. Local memory

Assignment Project Exam Help

Local memory

Accessible to all work items in a work group, but *not* between groups. Much faster than global memory.

<https://powcoder.com>

Typically used as a **scratch space** for calculations involving more 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 variable¹.

```
1 __kernel
2 void kernel(__global int *device_array)
3 {
4     __local int temp[128];
5     // Calculations involving temp.
6     ...
7     // Place final answer in device_array.
8 }
```

However, this is **static allocation** - the size of the array must be known **when the kernel is built**.

¹In 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 __kernel
2 void kernel(__local int *temp, __global int *device_a)
3 {
4     // Calculations involving temp.
5     ...
6     // Place final answer in device_a.
7 }
```

<https://powcoder.com>
Add WeChat powcoder

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

Only accessible to each work item. *Very fast access.*

<https://powcoder.com>

In practice, private memory is almost always implemented in hardware as **registers**:

- Small amount of memory that can be accessed quickly.
- 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`).

- Variables declared within a kernel **default** to private.

<https://powcoder.com>

For instance, in this kernel ...

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

...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 example, N is treated as a private variable:

```
1 __kernel
2 void kernel(__global float *array, int N)
3 {
4     // Calculations involving N.
5     ...
6 }
```

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 Private memory will 'spill over' into **global memory**.
- 2 Fewer work groups will be launched simultaneously, resulting in an **under-utilisation** of available processing units.

Either mechanism reduces performance.

Add WeChat powcoder

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). Faster than global memory. Not available on all GPUs.

GPUs often have memory that is global in scope, but can only be read from within kernels:

- Known as **constant** memory.
- Can still be written to **by the host**.
- Originally to accelerate the mapping of **textures** to polygons ('texture' memory¹).
- Typically much smaller than global memory, even if it exists.

¹CUDA treats texture memory separately to other constant memory.

Using constant memory¹

For kernel arguments, use the `__constant` specifier:

```
1 __kernel
2 void kernel(__constant float *a,...)
3 { ... }
```

Initialise the array (*i.e.* copy from the host) and set the kernel argument as normal.

Variables within kernels can also be `__constant`, but must be specified at compile time.

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

¹In 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:

- **Global** (slow), **local** (faster), **private** (very fast).
- Possibly also **constant** (faster than global)
- Private memory can **overflow**, resulting in performance loss.

<https://powcoder.com>

Add WeChat powcoder

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