

Assignment Project Exam Help

XJCO3221 Parallel Computation

<https://powcoder.com>

Peter Jimack

University of Leeds
Add WeChat powcoder

Lecture 17: Synchronisation

Previous lectures

Many of the previous lectures has mentioned **parallel synchronisation** in some form. However, there are many ways to synchronise:

- **locks** in shared memory systems [lectures 6 and 7].
- **Synchronisation barrier** at each level of a binary tree **reduction** [Lecture 11].
- **Blocking communication**, which affords a form of synchronisation, in distributed memory systems [Lecture 12].
- ...

Also recall that GPU's have multiple memory types, some of which can be viewed as *shared* (`--global`), and some which can be viewed as *distributed* (`--local`) [Lecture 16].

Today's lecture

Assignment Project Exam Help

Today's lecture will look at **synchronisation** on a GPU.

- How to synchronise **within** a work group.
- How to synchronise **between** work groups.

<https://powcoder.com>

We will also see how the SIMD cores can potentially **reduce** or **improve** performance.

Add WeChat powcoder

- Threads within a **subgroup** are **automatically synchronised**.
- Threads performing different calculations can lead to **divergence** and reduced performance.

Reminder: Scalar product

Assignment Project Exam Help

As an example, we will use the **vector product** between two n -vectors **a** and **b**, as in Lecture 11.

Written mathematically as (*indexing starting from 1*):

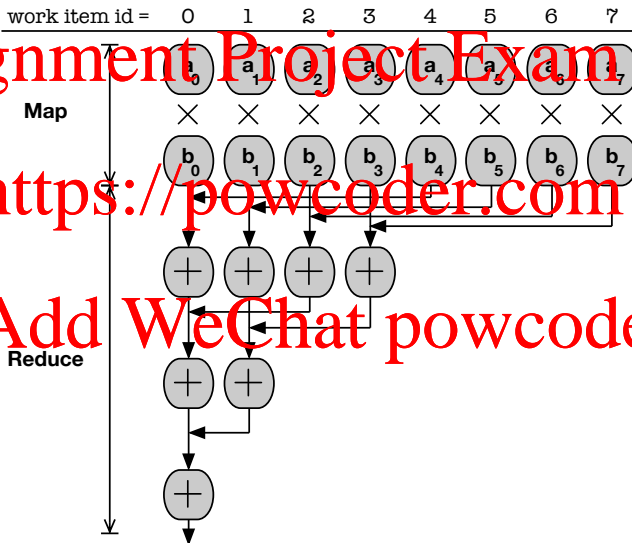
$$\mathbf{a} \cdot \mathbf{b} = \sum_{i=1}^n a_i b_i = a_1 b_1 + a_2 b_2 + \dots + a_n b_n$$

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

In serial CPU code (*indexing starting from 0*):

```
1 float dot = 0.0f;  
2 for( i=0; i<n; i++ ) dot += a[i] * b[i];
```

MapReduce pattern for $n = 8$



Reduction in local memory

First, consider n equal to or less than the work group size.

Use **local memory** for the intermediate quantities [Lecture 16].

- **Faster** than global memory; allows **communication** within a work group.

<https://powcoder.com>

Each work item copies $a[i]*b[i]$ to local memory first.

- Reduce using the binary tree pattern on the previous slide¹.
- Each addition performed by the work item with the lower i.d.
- Final result in work item with i.d. = 0 copied to the answer (in global memory).

¹Divide-by-two implemented by **compound bitwise right shift** operator '>>='.

Kernel code

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

Assignment Project Exam Help

```
1 __kernel
2 void reduceNoSync( __global float *device_a, __global
    float *device_b, __global float *dot, __local
    float *scratch )
3 {
4     int stride;
5     id      = get_local_id  (0),
6     groupSize = get_local_size(0);  //=work group
7
8     scratch[id] = device_a[id] * device_b[id];
9
10    for( stride=groupSize/2; stride>0; stride>>=1 )
11        if( id < stride )
12            scratch[id] += scratch[id+stride];
13
14    if(id==0) *dot = scratch[0];
15 }
```

<https://powcoder.com>

Add WeChat powcoder

Calling C-code

```
1 // float array of size 1 on device.  
2 cl_mem device_dot = clCreateBuffer(...);  
3  
4 ... // Set kernel arguments 0, 1 and 2.  
5 clSetKernelArg(kernel, 3, N*sizeof(float), NULL);  
6 // NULL => local memory of given size  
7  
8 // Add to the command queue.  
9 size_t indexSpaceSize[1]={N}, workGroupSize[1]={N};  
10 clEnqueueNDRangeKernel(queue, kernel, 1, NULL,  
    indexSpaceSize, workGroupSize, 0, NULL, NULL);  
11  
12 // Get the result back to host float 'dot'.  
13 float dot;  
14 clEnqueueReadBuffer(queue, device_dot, CL_TRUE, 0, sizeof(  
    float), &dot, 0, NULL, NULL);
```

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Barriers

Without synchronisation, this reduction is *not* guaranteed to work on *all* systems.

Recall that **barriers** are points in code that no processing unit can leave until all units reach it [v.f. Lecture 11]

- `#pragma omp barrier` in OpenMP.
- `MPI_Barrier()` in MPI.

In OpenCL¹, a barrier **within a work group** is implemented as:

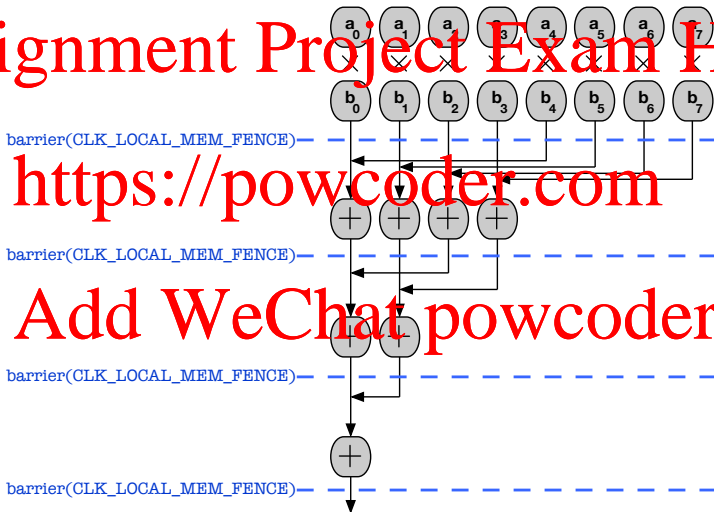
```
1 barrier(CLK_LOCAL_MEM_FENCE);
```

¹In CUDA: `__syncthreads()` synchronises within a **thread block**=work group.

Reduction with synchronisation

```
1 void reduceWithSync(...) // Same arguments
2 {
3     int id=..., groupSize=..., stride; // As before.
4
5     scratch[id] = device_a[id] * device_b[id];
6     barrier(CLK_LOCAL_MEM_FENCE) // Sync.
7
8     for( stride=groupSize/2; stride>0; stride>=1 )
9     {
10         if( id < stride )
11             scratch[id] += scratch[id+stride];
12
13         barrier(CLK_LOCAL_MEM_FENCE); // Sync.
14     }
15
16     if(id==0) *dot = scratch[0];
17 }
```

Reduction with barrier(CLK_LOCAL_MEM_FENCE)



Problems larger than a single work group?

If we could synchronise **between** work groups, could use the same method as before.

- 1 Make device vectors and scratch **global**.
- 2 Replace **local** barriers with **global** barriers.

However, **no such global barrier exists**.

GPUs cannot synchronise between work groups/thread blocks²

Add WeChat powcoder

¹`barrier(CLK_GLOBAL_MEM_FENCE)` *does* exist, but refers to accesses to global memory; it still only synchronises *within* a work group.

²Some modern GPUs support **cooperative groups** that allow synchronisation across multiple thread blocks; e.g. CUDA 9.0.

Warning!

Assignment Project Exam Help

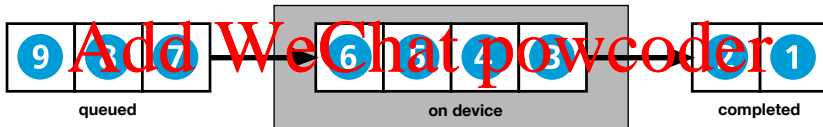
You might see claims that it is possible to synchronise globally on any GPU by constantly **polling** a global memory location.

- i.e. work items constantly read/write to synchronise.

This may work, but **only for small problems**.

<https://powcoder.com>

If there are too many work groups for the device, it **queues** them:



If they are not all on the device at the same time, it is **impossible to synchronise within one kernel** using this method.

Solution: Multiple kernels

The solution is to break the kernel at the barrier point into multiple kernels called consecutively:



This way kernel 1 completes before kernel 2 starts.

Reduction across work groups

Assignment Project Exam Help

- 1 Repeatedly call kernel that reduces an array of **partial sums** until less than maximum work group size.
- 2 Final kernel call to reduce these partial sums.

<https://powcoder.com>

It is simpler (although less efficient) to use the CPU:

- 1 Each work group inserts its partial sum into a global array.
- 2 Final summation performed on the host

Add WeChat powcoder

This is conceptually similar to an MPI program performing final calculations on rank 0.

¹Wilt, *The CUDA handbook* (Addison-Wesley, 2013).

Subgroups (warp, wavefront, etc.)

Assignment Project Exam Help

Recall that GPUs are based on SIMD cores.

- Each core contains **multiple hardware threads** that perform the **same operation**.

<https://powcoder.com>

In OpenCL, the number of **work items** (i.e. **threads**) simultaneously on a single SIMD core is known as a **subgroup**.

- Smaller** than a work group.

Add WeChat powcoder

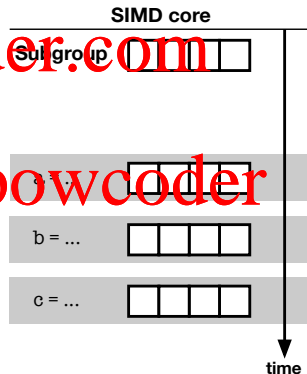
The actual size is vendor specific. For example:

- Nvidia call them **warps**, each of which has 32 threads.
- AMD have 64-thread **wavefronts**.

Lockstep

The SIMD core applies the same operation to all items in the subgroup **simultaneously**. We say it advances in **lockstep**.

```
1  __kernel
2  void kernel(...)
3  {
4      int id = get_global_id(0);
5      float a, b, c;
6
7      a = 4*array[id];
8
9      b = a*a;
10
11     c = b + a;
12 }
```



Reduction with a subgroup

For reduction, this means that once the problem has been reduced to the size of a subgroup, there is no longer any need for explicit synchronisation¹:

```
1 __kernel
2 void reduce(...)
3 {
4     ... // Start as before.
5
6     // Split the loop into two.
7     for(;; stride=group/2; stride>subgroup; stride>>=1) {
8         if(id<stride) scratch[id] += scratch[id+stride];
9         barrier(CLK_LOCAL_MEM_FENCE); // Sync.
10    }
11    // See next slide ...
```

¹Wilt, *The CUDA handbook* (Addison-Wesley, 2013).

Final reduction

Assignment Project Exam Help

For the final reduction, exploit **lockstepping** by removing the synchronisation:

```
1 for (; stride > 0; stride >>= 1)
2 {
3     if (id < stride)
4         scratch[id] += scratch[id + stride];
5
6     // No barrier()
7 }
```

<https://powcoder.com>

Add WeChat powcoder

This avoids any overheads with calling `barrier()` (*i.e.* unnecessarily checking if all threads have reached this point, when we know they must have).

Divergence

Assignment Project Exam Help

Recall that SIMD cores can also be termed **SIMT** (lecture 4)

- Single Instruction stream, Multiple Threads.

This means it is possible to perform different operations with a subgroup, but they **remain in lockstep**.

- Only **one** distinct operation can be performed **at a time**.

Add WeChat powcoder

This can lead to **serialisation** where operations are performed **one after the other**.

- Can lead to a severe performance penalty.

Code that leads to divergence

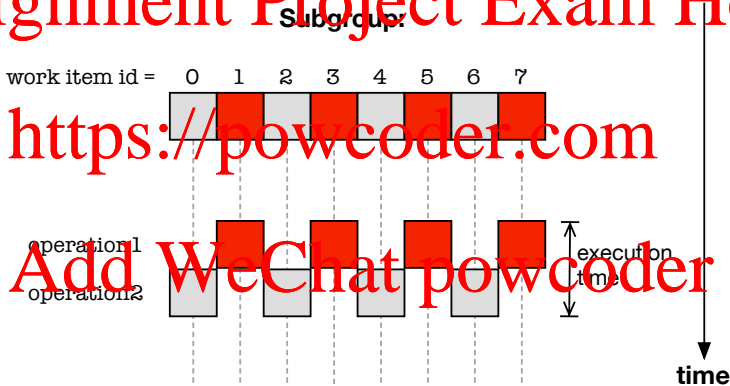
Suppose we want even-numbered work items to perform one operation, and odd-numbered items a different one¹:

```
1 __kernel
2 void kernel(...)
3 {
4     int id = get_global_id(0);
5
6     if (id%2)
7         operation1; // id odd.
8     else
9         operation2; // id even.
10 }
```

¹Recall $i\%2==1$ for i odd, 0 for i even.

In this example the execution time is **double** what was expected:

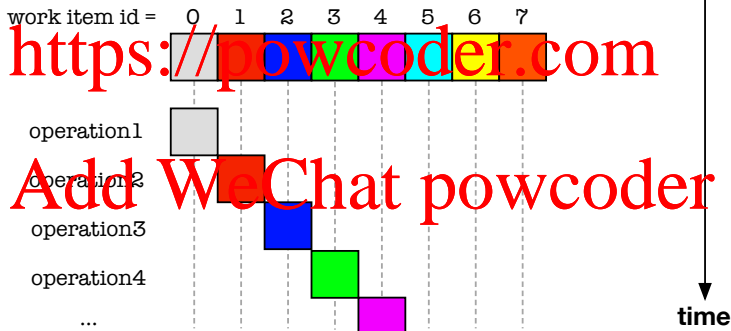
Assignment Project Exam Help



For more operations the execution time increases further, e.g. a switch-case clause where every thread performs a different operation.

Assignment Project Exam Help

Subgroup:



This is true **serialisation**.

Summary and next lecture

Assignment Project Exam Help

Today we have looked at **synchronisation**, focussing on GPUs

- **Barriers** can synchronise within a work group.
- **Cannot** synchronise between work groups within a kernel - must split into separate kernels.
- Work items of threads within a **subgroups** execute in lockstep.
 - No need for explicit synchronisation mechanisms.
 - Can lead to **divergence** and reduced performance.

Next time we will look at **atomic** instructions, continuing what we started in Lecture 6.