Overview
Host/device latency hiding
Task graphs
Summary and next lecture

University of Leeds

Lecture 19: Task parall

Overview
Host/device latency hiding
Task graphs
Summary and next lecture

Previous lectures
Today's lecture

## Previous lectures

So much of this module have considered **loop parallel** problems:

- Same operation applied to multiple data sets.

In Le

- independent **tasks** to all other units.
- These **worker** processes performed result back to the main process.

We referred to this as **task parallelism** since the emphasis was on parallelising **tasks** rather than the data.

Overview
Host/device latency hiding
Task graphs
Summary and next lecture

Previous lectures
Today's lecture

# Today's lecture

Today we will look at task parallelism in more detail:

- How GPU **command queues** or **streams** can permit:
  -
  -
  -

- **Task graphs** that are derived from thes

- The **work span model** that estima                    **d
  up** from a task graph.

Firstly, we will see how to time an OpenCL program using an **event**.

Overview
Host/device latency hiding
Task graphs
Summary and next lecture

Queues and events in OpenCL
Overlapping host and device computation
Overlapping computation with communication

# Timing kernels in OpenCL
Code on Minerva: `timedReduction.c`, `timedReduction.cl`, `helper.h`

For **profiling** purposes we often want to **time** how long a kernel take

Timi https://eduassistpro.github.i

1. Ensure the **command queue** supports profiling.
2. Declare an **event**, and attach to the kern enqueued.

   Add WeChat edu_assist_pr

3. Extract the time taken **once the ker**

Some of previous code examples already do this.

Overview
Host/device latency hiding
Task algorithms
Summary and next lecture

Queues and events in OpenCL
Overlapping host and device computation
Overlapping computation with communication

```
1  // Ensure queue supports profiling.
2  cl_command_queue queue = clCreateCommandQueue
3       (context,device,CL_QUEUE_PROFILING_ENABLE,&status);
4
5  // OpenCL event.
6  cl_ev
7
8  // Enque
9  statu  = c
10
11 // ... (once the kernel has finished)
12 cl_ulong start,end;
13 clGetEventProfilingInfo(timer,
14     CL_PROFILING_COMMAND_START,sizeof(cl_ulong),&start
15     ,NULL);
14 clGetEventProfilingInfo(timer,CL_PROFILING_COMMAND_END
        ,sizeof(cl_ulong),&end,NULL);
15 printf("Time: %g ms\n",1e-6*(cl_double)(end-start));
```

Overview
Host/device latency hiding
Task glilies
Summary and next lecture

Queues and events in OpenCL
Overlapping host and device computation
Overlapping computation with communication

# Blocking communication

Recall that when we copy the data from device to host at the end
of the calculation, we typically use a **blocking** call:

```
1  clEnqueueReadBuffer ( queue , device_dot , CL_TRUE ,...);
```

- 

- Similar to `MPI_Recv()` *[cf. Lecture 9]*.

Replacing `CL_TRUE` with `CL_FALSE` ma
copy command[1]:

- Will return 'immediately,' **before** the copy is complete.

- Similar to `MPI_Irecv()` *[cf. Lecture 12]*.

---

[1]In CUDA: Use `cudaMemcpyAsync()` rather than `cudaMemcpy()`.

Overview
Host/device latency hiding
Task graphs
Summary and next lecture

Queues and events in OpenCL
Overlapping host and device computation
Overlapping computation with communication

## Potential consequences of non-blocking

# Assignment Project Exam Help

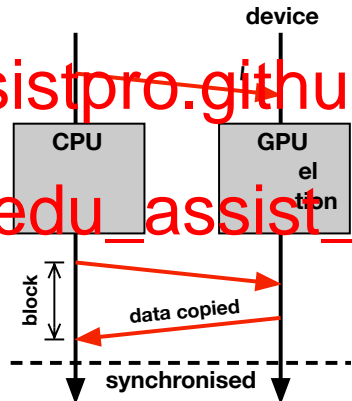For this example, using a non-blocking copy can mean:

1

# https://eduassistpro.github.i

2

# Add WeChat edu_assist_pr

Note that the read did not start until the kernel had fin

- It was enqueued on the **same command queue**.

Overview
Host/device latency hiding
Task graphs
Summary and next lecture

Queues and events in OpenCL
Overlapping host and device computation
Overlapping computation with communication

# Overlapping host and device computation

The queuing model means we can perform calculations on the host (CPU) and device (GPU) **simultaneously**.

```
1  // Enqueue k...
2  clEnqu...
3
4  // Perform useful operations
5  // on the host.
6  ...
7
8  // Blocking copy device->host
9  clEnqueueReadBuffer(...);
10
11 // Device and host in sync.
```

Overview
Host/device latency hiding
Task graphs
Summary and next lecture

Queues and events in OpenCL
Overlapping host and device computation
**Overlapping computation with communication**

# Overlapping computation with communication

Recall from Lecture 12 that we can reduce **latency** by overlapping **computation** with **communication**:
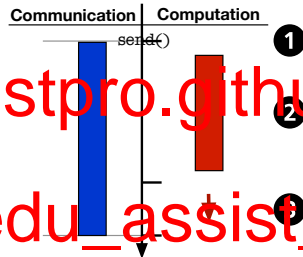


1. [text obscured]

2. Perform **calculations**.

3. **Synchronise** using MPI_Wait().

> Similar benefits can be achieved on a GPU using multiple
> **command queues** (OpenCL) / **streams** (CUDA).

Overview
Host/device latency hiding
Task graphs
Summary and next lecture

Queues and events in OpenCL
Overlapping host and device computation
**Overlapping computation with communication**

# Multiple command queues: Example

Consider the following problem:

- 
- 
- 

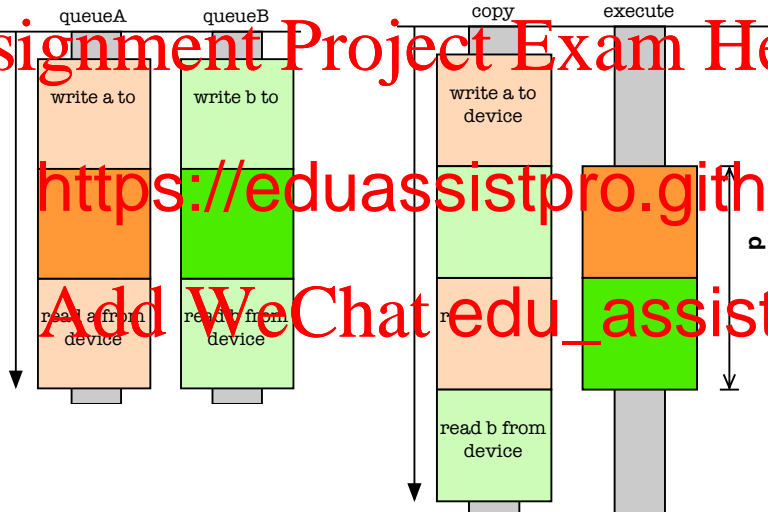Suppose our device supports **asynchro**
**simultaneous data transfer and kernel exec**

- **Not guaranteed**, although common in modern GPUs.
- May require device to have **direct access** to host memory.

Overview
Host/device latency hiding
Task graphs
Summary and next lecture

Queues and events in OpenCL
Overlapping host and device computation
Overlapping computation with communication

# OpenCL with two command queues: Outline

```
 1  // Initialise two command queues.
 2  cl_command_queue queueA = clCreateCommandQueue(...);
 3  cl_command_queue queueB = clCreateCommandQueue(...);
 4
 5  // Enque
 6  clEnq
 7  clEnq
 8
 9  // Enqueue both kernels.
10  clEnqueueNDRangeKernel(queueA,kernelA,...);
11  clEnqueueNDRangeKernel(queueB,kernelB,...);
12
13  // Enqueue data transfer device->host (blocking).
14  clEnqueueReadBuffer(queueA,...,CL_TRUE,...);
15  clEnqueueReadBuffer(queueB,...,CL_TRUE,...);
16
17  ... // Process results; clear up.
```

Overview
Host/device latency hiding
Task graphs
Summary and next lecture

Queues and events in OpenCL
Overlapping host and device computation
**Overlapping computation with communication**

**Program logic**

**On the device**

Overview
Host/device latency hiding
Task graphs
Summary and next lecture

Events in queues and streams
Task dependencies as directed acyclic graphs
Work-span model
Task graphs as a programming pattern

# Events in queues and streams
Code on Minerva: `taskGraph.c`, `taskGraph.cl`, `helper.h`

Earlier we saw how an **event** can be used as a timer.

```
1  cl_ev
```

In ge https://eduassistpro.github.i

**kern**

The last arguments on enqueue commands are:

```
1  clEnqueue...(...,numWait,waitEvents,event);
```

| cl_uint numWait | Number of events to wait for. |
|---|---|
| cl_event *waitEvents | List of events to wait for. |
| cl_event *event | Used to identify when this operation completes. |

Overview
Host/device latency hiding
Task graphs
Summary and next lecture

Events in queues and streams
Task dependencies as directed acyclic graphs
Work-span model
Task graphs as a programming pattern

## Example (fragment)

Link together reads, writes and kernels **on multiple queues** — not necessary when using a single queue.

```
1  cl_event writeEvent, kernelEvent, readEvent;
2
3  // Non-b
4  clEnq
5
6  // Enqueue kernel.
7  clEnqueueNDRangeKernel(..,1,&writeEvent,&kerne
8
9  // Non-blocking read device->host.
10 clEnqueueReadBuffer(..,1,&kernelEvent,&readEve
11
12 // Synchronise (wait for read to complete).
13 clWaitForEvents(1,&readEvent);  // Sim. to MPI_Wait().
```
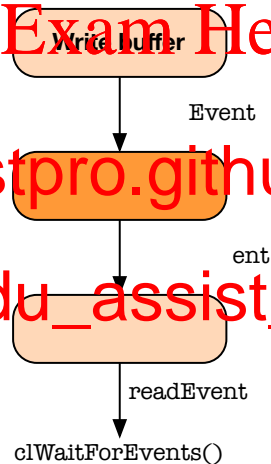
Overview
Host/device latency hiding
Task graphs
Summary and next lecture

Events in queues and streams
Task dependencies as directed acyclic graphs
Work-span model
Task graphs as a programming pattern

# Task graphs

**Events** are used to link two tasks when the first must complete before the sec

Simp

- **Directed, acyclic graph**.
- Nodes are **tasks**.
- Edges denote **dependencies**.
- **Direction** denotes which task must complete before the other begins.

**Write buffer**

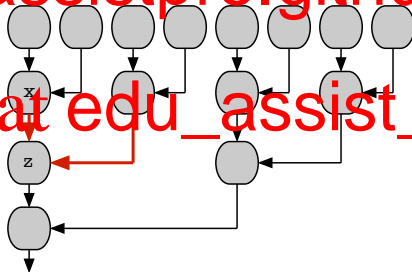Event

ent

readEvent

clWaitForEvents()

Overview
Host/device latency hiding
**Task graphs**
Summary and next lecture

Events in queues and streams
Task dependencies as directed acyclic graphs
Work-span model
Task graphs as a programming pattern

## Earlier task graphs

In fact we have already seen examples of **task graphs**.

In **bi**
must

Exa

Must know $x$ and $y$ before
we can calculate $z = x \otimes y$.
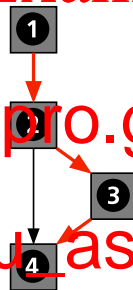
*(These two dependencies
highlighted in the diagram).*

Overview
Host/device latency hiding
Task graphs
Summary and next lecture

Events in queues and streams
Task dependencies as directed acyclic graphs
Work-span model
Task graphs as a programming pattern

## Satisfying dependencies

However, many processing units, dependencies **must be satisfied**.



**Invalid**                    **Valid**

This means you must always take the 'longest path'.

Overview
Host/device latency hiding
Task graphs
Summary and next lecture

Events in queues and streams
Task dependencies as directed acyclic graphs
Work-span model
Task graphs as a programming pattern

# Work-span model

Assignment Project Exam Help

Consider the task graph on the right:

- 9
- 1

Assum https://eduassistpro.github.i

The **performance** of a parallel program represented as a task graph can be estimated once the **work** and **span** have been identified.

Add WeChat edu_assist_pr

Overview
Host/device latency hiding
**Task graphs**
Summary and next lecture

Events in queues and streams
Task dependencies as directed acyclic graphs
**Work-span model**
Task graphs as a programming pattern

## Work and span

### Definition

The **work** is the **total** time to complete all tasks.

This t its

### Definition

The **span** is the time taken on an ideal machin

As many tasks in parallel as possible **g**

- The **span** is the **longest path** executed one after the other.
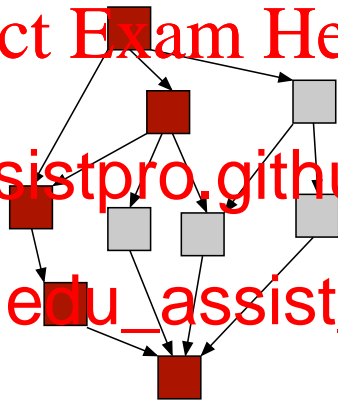- Also called the **critical path**.

Overview                          Events in queues and streams
Host/device latency hiding        Task dependencies as directed acyclic graphs
Task graphs                       Work-span model
Summary and next lecture          Task graphs as a programming pattern

## Span example

In this example, the number of task

**give**

The **span** is therefore 5 tasks.

For uneven sized tasks would be measured in units of seconds/clock cycles/FLOPs
*etc.*



■ = task along the critical path

Overview
Host/device latency hiding
**Task graphs**
Summary and next lecture

Events in queues and streams
Task dependencies as directed acyclic graphs
Work-span model
Task graphs as a programming pattern

# Work-span model

Note that the **work** is just the serial execution time $t_s$.

We have argued that the parallel execution time can never become less than the **span**.

Ther

$$\frac{}{t_p} \quad \frac{}{t_{p=\infty}} \quad \frac{}{}$$

This is the **work-span model**

- Upper limit[1] for $S$ based purely on the t
- $S \leq \frac{9}{5} = 1.8$ for this example.

---

[1]There is also a *lower* bound provided by Brent's lemma. R.P. Brent, *J. Ass. Comp. Mach.* **21**, 201 (1974).

Overview
Host/device latency hiding
**Task graphs**
Summary and next lecture

Events in queues and streams
Task dependencies as directed acyclic graphs
Work-span model
**Task graphs as a programming pattern**

## Superscalar sequences and futures

Some parallel frameworks schedule tasks based on **dependencies specified by the programmer**.

- 
- 
- 

This is sometimes referred to as a **sup**                    [1]

The benefit is that you do not need to

- The runtime system synchronises when n
  the dependencies you provide.

---

[1]McCool *et al.*, *Structured parallel programming* (Morgan-Kauffman, 2012).

Overview
Host/device latency hiding
Task graphs
Summary and next lecture

Summary and next lecture

# Summary of GPGPU programming

| Lec. | Content | Key points |
|---|---|---|
| 14 | GPGPU archi- | SIMD cores; CUDA and OpenCL; start- |
| 1 | Th | el |
| 16 | Memory types | Global, local, private and constant. |
| 17 | Synchronisation | Barriers; breaking up kerne |
| | | ...ing in oc... |
| 18 | Atomics | Global and local atomics; compare-and- |
| | | exchange; lock-free data structures. |
| 19 | Task paral- | GPU queues/streams; events; task |
| | lelism | graphs, work and span. |

Overview
Host/device latency hiding
Task graphs
Summary and next lecture

Summary and next lecture

# Next lecture

This penultimate lecture is the last containing new material.

In the **parallel concept** rather than by architecture.

- Alternative perspective focussing on tran
- Also serves as a useful summary of the modul