# **CSCI 4060U Test #3**

Started: Apr 10 at 10:01a.m.

# **Quiz Instructions**

This test is take-home and you will have 48 hours to complete it starting at 9:40am on Monday, April 10th. On this test you will be given questions with answers written by ChatGPT3. You will also be given a marking rubric. For each question/answer pair you need to mark the question and provide written feedback. Your feedback should explain your marking. For correct parts of the answer explain why the marks were earned. For incorrect parts of the answer explain what was wrong and what the correct answer should be. Good luck!

Question 1 12 pts

(11 marks) Consider the OpenCL memory and device models.

- A. **(4) What are the four types of device memory?** [1 mark for each device memory type]
- B. (3) Is it possible to access device memory from the host? Explain. [1 mark for answering yes/no. 2/2 marks for a correct explanation with full detail (including method calls from host), 0/2 marks for an incorrect explanation and partial marks for partially correct]
- C. **(4)** What is the difference between a compute unit and a processing element? [4 marks for a correct explanation, 1.5 marks for defining each correctly but not explicitly listing the difference.]

#### **ChatGPT Answer:**

Α.

OpenCL defines four types of device memory that can be used for storing data during the execution of a kernel:

Global memory: This is the largest memory space available on a device and is used for storing data that needs to be accessed by all work-items in a kernel. Global memory is persistent and can be read and written by the host and all work-items.

Constant memory: This memory space is read-only and is used for storing data that does not change during the execution of a kernel. Constant memory is often used to store lookup tables and other data that is used repeatedly by a kernel.

Local memory: Each work-group in a kernel has its own local memory, which is used for storing data that is shared among work-items in the same work-group. Local memory is faster than global memory, but it is limited in size.

Private memory: This is memory that is used by a single work-item in a kernel. Private memory is often used to store temporary variables and is the fastest type of memory available on a device.

It's important to note that the different types of memory have different characteristics in terms of size, speed, and accessibility, and choosing the right type of memory for a particular task can have a significant impact on the performance of an OpenCL kernel.

В.

Yes, it is possible to access device memory from the host in OpenCL, but it requires the use of specific functions and mechanisms provided by the OpenCL API.

C.

In OpenCL, a compute unit and a processing element are both hardware components that are used to perform computations on a device. However, there are some differences between the two.

A compute unit is a term used to describe a group of processing elements that share a common instruction stream and can execute a subset of the total number of workitems in a kernel. A compute unit typically includes a local data store and a set of processing elements, and it is responsible for scheduling and executing work-items within its assigned subset of the kernel. A device can have multiple compute units, each of which can execute work-items in parallel.

A processing element, on the other hand, is a term used to describe an individual unit of computation within a compute unit. A processing element typically includes an ALU (Arithmetic Logic Unit), a register file, and access to the local data store. Processing elements are responsible for executing instructions and performing arithmetic and logic operations on data.

A) (4/4) The 4 types of device memory are global, constant, local, private. They can be declared by doing the following: (\_\_global, \_\_constant, \_\_local, \_\_private)

D) (2/2) Compatibility of Device Province In the Company of the Co

B) (2/3). Correct answer for yes/no. Device memory can be accessed by the host. -1 mark for not mentioning correct function. We can Call the clCreateBuffer() function with the CL\_MEM\_ALLOC\_HOST\_PTR flag which will instruct OpenCL to allocate memory on the host and return a pointer to that memory.

C) (4/4). Correct Answer. Each device is organized as a group of compute units, and each compute unit is organized as a group of processing elements. Compute units typically include shared resources such as memory and registers. A processing element is responsible for for executing a single instruction on a

p





129 words | </> /





## **Question 2**

4 pts

(5 marks) Consider the following statement and explain why it is true or false:

"In OpenCL programs, host data types and device data types are identical."

[2 marks for true/false with no explanation. 3/3 marks for explanation in full, 0/3 marks for no explanation or wrong explanation, partial marks otherwise.]

#### **ChatGPT Answer:**

The statement "In OpenCL programs, host data types and device data types are identical" is false.

In OpenCL, the host refers to the computer's CPU, while the device refers to the device's processor, such as a GPU or an FPGA. These processors may have different architectures and data types. For example, the host may use a floating-point format that is different from the floating-point format used by the device. Additionally, the host and device may have different byte orders, so data that is transmitted between the two may need to be converted.

OpenCL provides mechanisms for transferring data between the host and device and for converting data between different formats and byte orders. The OpenCL API provides functions for allocating memory on the device, transferring data between the host and device, and synchronizing the execution of kernels. It also provides data types, such as cl float and cl double, which are specific to the device's architecture and can be used to ensure compatibility between the host and device.

In summary, host data types and device data types are not identical in OpenCL, but OpenCL provides mechanisms for managing the differences between the two and ensuring compatibility between the host and device.



- +2 for specifying that the statement is false
- +3 for the correct explanation. Host and device types may have different representations for data, and could also differ in the size of that data.

Total: 5/5

р





**Question 3** 8 pts

## (4 marks) Graphics card performance.

- A. (2) What is the main *performance* issue when running OpenCL device code on a graphics card compared with running it on an accelerated processing unit (APU)? [2 marks for correct main performance issue, 0 marks for an incorrect performance issue, partial marks for a performance issue that is not the main one]
- B. (2) Is there a hardware solution that can alleviate this issue? [2 marks for correct explanation, 0 marks for incorrect, partial otherwise]

#### **ChatGPT Answer:**

Α.

One of the main performance issues when running OpenCL device code on a graphics card (GPU) is memory bandwidth. Graphics cards are optimized for data-parallel processing, which means that they are designed to execute the same instruction on a large set of data in parallel. However, in order to achieve high performance, the GPU needs to be able to access data quickly.

Memory bandwidth refers to the rate at which data can be read from or written to memory. GPUs typically have high memory bandwidth to support the large amount of data that is processed in parallel. However, if the GPU is not able to access data quickly enough, it can become a bottleneck that limits performance.

There are several factors that can affect memory bandwidth on a GPU, including memory latency, memory size, and memory access patterns. Memory latency refers to the time it takes for the GPU to retrieve data from memory, and can be affected by factors such as the memory clock speed and the memory bus width. Memory size refers to the amount of memory available on the GPU, and can affect performance if the GPU needs to access more memory than is available on the device. Memory access patterns refer to the way in which data is accessed in memory, and can affect performance if the GPU needs to access data that is not stored in a contiguous block of memory.

To maximize performance when running OpenCL device code on a GPU, it is important to optimize memory access patterns, minimize memory latency, and use data structures and algorithms that are optimized for data-parallel processing. Additionally, it may be necessary to adjust the work-group size and other parameters to achieve the best performance for a particular device and workload.

В.

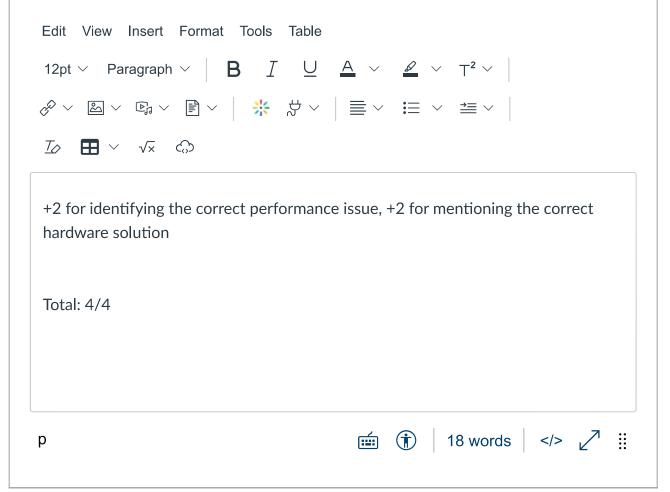
Yes, there are hardware solutions that can help alleviate the memory bandwidth issue when running OpenCL device code on a graphics card.

One such solution is the use of high-bandwidth memory (HBM). HBM is a type of memory that is designed to provide high bandwidth and low latency access to memory. HBM is stacked directly on top of the GPU, allowing for shorter interconnects and faster data transfer between the GPU and memory.

Another solution is the use of cache memory. GPUs typically have multiple levels of cache memory, which can be used to reduce the number of memory accesses and improve memory access patterns. By storing frequently accessed data in cache memory, the GPU can reduce the amount of time it spends waiting for data to be retrieved from memory.

Finally, newer generations of GPUs are incorporating more advanced memory controllers that can dynamically adjust memory clock speeds and access patterns based on the workload. This can help to maximize memory bandwidth and reduce memory latency for a wide range of workloads.

In summary, there are hardware solutions, such as HBM, cache memory, and advanced memory controllers, that can help alleviate the memory bandwidth issue when running OpenCL device code on a graphics card.



Question 4 4 pts

(5 marks) What is the basic structure of an OpenCL host program? [1 mark for each correct part of the structure]

**ChatGPT Answer:** 

The basic structure of an OpenCL host program typically involves the following steps:

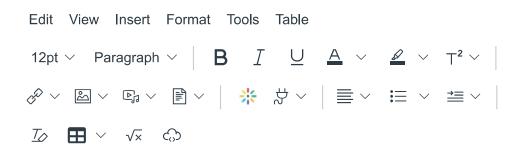
Discovering and selecting a compute device: This involves identifying the available OpenCL platforms, selecting a specific platform, and choosing a compute device (such as a GPU or CPU) on that platform to perform the computations.

Creating an OpenCL context: A context is an object that encapsulates the state of an OpenCL execution environment, including the selected compute device, memory objects, and command queues. It is created by passing the selected device(s) to the clCreateContext() function.

Creating memory objects: Memory objects are used to store data that will be processed by the compute device(s). In OpenCL, there are two types of memory objects: buffers and images. Buffers are simple one-dimensional arrays, while images are more complex data structures that can represent two- or three-dimensional arrays or other types of data.

Building and compiling OpenCL programs: An OpenCL program consists of one or more compute kernels, which are small programs that execute on the compute device(s). Kernels are written in a variant of the C programming language, called OpenCL C. The OpenCL program is compiled and built using the clCreateProgramWithSource() and clBuildProgram() functions.

Executing the kernels on the compute device(s): Once the OpenCL program is built and the memory objects are created, the kernels can be executed on the compute device(s) by enqueuing commands in a command queue. Commands are enqueued using the clEnqueueNDRangeKernel() function, which specifies the kernel to be executed, the work dimensions, the work group size, and other parameters.



(5/5) correctly Identifies the basic structure of host program

Define the platform

Create and Build the program

Setup memory objects

Dofina the karnal

Denne the kerner

Submit the commands (transfer memory objects and execute kernels)

р





(†) 34 words </>





**Question 5** 4 pts

## [3 marks] What's the BIG win in using OpenCL? Answer by considering traditional loops like:

```
void multiply(const int n, const float *a, const float *b, float *c) {
  int i;
  for (i=0; i < n; i++) {
    c[i] = a[i] * b[i];
  }
}
```

[3 marks for correct explanation, 0 marks for incorrect, partial otherwise]

#### ChatGPT Answer:

The big win in using OpenCL is the ability to harness the computing power of multiple processors on a single machine, or even across multiple machines. Traditional loops like the one provided can only use a single processor to perform the computations, which can limit the speed at which the computations are performed. However, with OpenCL, the same computations can be performed in parallel on multiple processing units (e.g. CPUs, GPUs, etc.), leading to much faster execution times. This is particularly useful for large-scale computations that involve a large amount of data and require significant computational power.

```
Edit View Insert Format Tools Table
12pt \vee Paragraph \vee | B I \underline{\cup} \underline{A} \vee \underline{\mathscr{D}} \vee \top^2 \vee
& ∨ ≧ ∨ ₽ ∨ ₽ ∨ | ★ ♡ ∨ | ≡ ∨ ≒ ∨ |
```

(2.5/3) -0.5 for not mentioning the SIMD architecture, It uses kernels to take advantage of the single instruction multiple data architecture to take advantage of the fact that every piece of data will have the same instruction applied to it in the loop. Since every element being calculated in the loop is having the same instruction being applied to them (multiplication) this is the perfect example that could take advantage of a SIMD architecture.

Total: 2.5/3

P

Total: 2.5/3

Total: 2.5/3

Question 6 8 pts

## (4 marks) Is the below example also legal in OpenCL C?

```
__global int* x;
__private int* y;
x = y;
```

## If yes, explain why. If no, explain how the code could be fixed.

[2 marks for identifying if code is correct/incorrect. [2 marks for correct explanation, 0 marks for incorrect, partial otherwise]

#### **ChatGPT Answer:**

No, the given example is not legal in OpenCL C.

In OpenCL C, the \_\_global and \_\_private storage qualifiers are used to define variables that are allocated in global memory or private memory, respectively, on the device.

In the given example, \_\_global int\* x declares a pointer to a global memory buffer on the device, whereas \_\_private int\* y declares a pointer to a private memory buffer on the device.

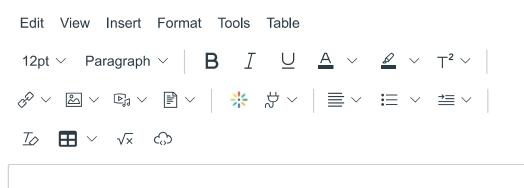
The line x = y; is not legal in OpenCL C because it attempts to assign a private memory buffer pointer to a global memory buffer pointer, which have different storage qualifiers and point to different memory spaces on the device.

Therefore, the types of the two pointers must match when they are assigned to each other in OpenCL C.

To fix this, you would need to allocate memory for both x and y within the same address space, and use appropriate OpenCL memory allocation functions. For

example, you could allocate x using clCreateBuffer with the CL\_MEM\_READ\_WRITE flag to allocate memory in the global memory space, and y using clLocalAlloc or cl mem alloc local to allocate memory in the private memory space.

After allocating the memory, you could copy the contents of y to x using OpenCL memory transfer functions such as clEnqueueWriteBuffer and clEnqueueReadBuffer.



(4/4) +2 marks for specifying that the mentioned operation is not a legal operation in OpenCL. +2 marks for specifying the correct reason being that they are in different types of memory.

p

(i) 36 words </> √/ ∷

# Question 7 8 pts

## (8 marks) Matrix multiplication example

### A. (5) Consider the code

```
void seq_matrix_mult(
int m_dim, int k_dim, int n_dim,
float *a, float *b, float *ab);
```

### Write an optimal OpenCL kernel for this code.

[2 marks for header and parameters, 1 mark for local variables, 2 marks for computation loop]

B. (3) If m\_dim=1000, k\_dim=500 and n\_dim=50, how many processing elements are ideal? [1.5 marks for correct answer, 1.5 marks for explanation (1.5 correct, 0 incorrect and 1 otherwise).]

### **ChatGPT Answers:**

Here's an OpenCL kernel for matrix multiplication that can be optimized further based on the specific hardware:

```
__kernel void matrix_mult(
    const int m_dim, const int k_dim, const int n_dim,
    __global float* a, __global float* b, __global float* ab)
{
    const int i = get_global_id(0);
    const int j = get_global_id(1);
    float sum = 0.0f;
    for (int k = 0; k < k_dim; k++)
    {
        sum += a[i * k_dim + k] * b[k * n_dim + j];
    }
    ab[i * n_dim + j] = sum;
}</pre>
```

This kernel uses global memory for accessing the matrices a, b, and ab. The kernel is launched with 2D global work size where the first dimension is m\_dim and the second dimension is n\_dim. Each work-item computes a single element of the output matrix ab.

To optimize this kernel, we can consider the following techniques:

Use local memory to cache input matrices a and b for reuse within a work group.

Use vector data types for input matrices a and b to take advantage of SIMD units.

Use kernel fusion to combine matrix multiplication with other operations to reduce memory bandwidth and improve performance.

Use the appropriate work-group size based on the hardware architecture to maximize utilization of compute units and memory hierarchy.

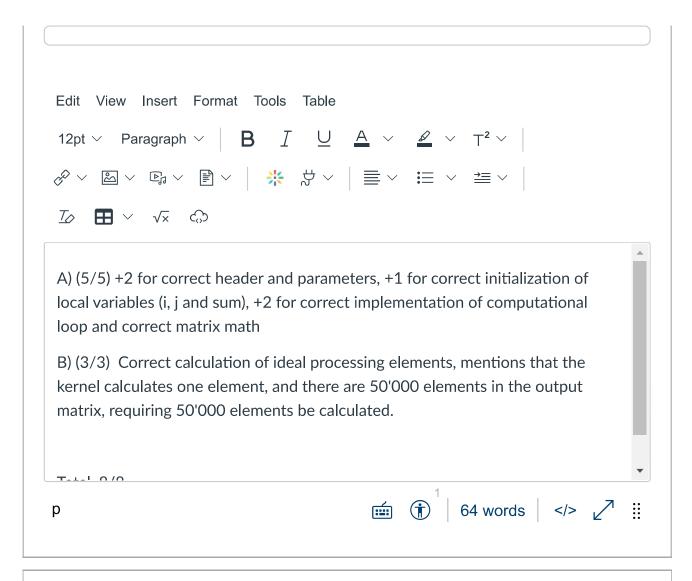
Optimizing the kernel further would require analyzing the hardware and profiling the application to identify bottlenecks and optimize accordingly.

В.

The total number of elements to be computed is m\_dim \* n\_dim = 1000 \* 50 = 50,000.

Each processing element can compute one element of the result array. Therefore, we need at least 50,000 processing elements to perform the computation.

However, the number of processing elements required may be higher depending on other factors such as the available memory bandwidth, cache size, and the specific hardware being used.



Question 8 20 pts

(10 marks) Consider the following statements regarding the benefits of OpenCL over CUDA. Identify each statement as true or false.

- OpenCL has better performance than CUDA on all types of hardware.
- OpenCL allows for heterogeneous computing, while CUDA is only available on NVIDIA GPUs.
- OpenCL has better memory management capabilities than CUDA.
- OpenCL has a more diverse range of supported hardware platforms, while CUDA is limited to NVIDIA GPUs.
- OpenCL can work with CPUs and GPUs simultaneously, while CUDA is limited to GPUs only.
- OpenCL has lower performance compared to CUDA when running on NVIDIA GPUs.
- OpenCL lacks support for advanced features like dynamic parallelism and unified memory.

- OpenCL is open-source, while CUDA is proprietary.
- OpenCL has more flexible memory management than CUDA, allowing for custom memory allocation schemes.
- OpenCL is an open standard, allowing for interoperability between different hardware vendors and platforms.
- OpenCL provides better developer tools and support than CUDA.
- OpenCL has a steeper learning curve than CUDA due to its complex API.

[1 mark per statement for correct true/false]

#### **ChatGPT Answer:**

OpenCL has better performance than CUDA on all types of hardware. (False)

OpenCL allows for heterogeneous computing, while CUDA is only available on NVIDIA GPUs. (True)

OpenCL has better memory management capabilities than CUDA. (False)

OpenCL has a more diverse range of supported hardware platforms, while CUDA is limited to NVIDIA GPUs. (True)

OpenCL can work with CPUs and GPUs simultaneously, while CUDA is limited to GPUs only. (True)

OpenCL has lower performance compared to CUDA when running on NVIDIA GPUs. (False)

OpenCL lacks support for advanced features like dynamic parallelism and unified memory. (False)

OpenCL is open-source, while CUDA is proprietary. (True)

OpenCL has more flexible memory management than CUDA, allowing for custom memory allocation schemes. (True)

OpenCL is an open standard, allowing for interoperability between different hardware vendors and platforms. (True)

OpenCL provides better developer tools and support than CUDA. (False)

OpenCL has a steeper learning curve than CUDA due to its complex API. (False)

Edit View Insert Format Tools Table



- OpenCL has better performance than CUDA on all types of hardware. (False) 1/1
- OpenCL allows for heterogeneous computing, while CUDA is only available on NVIDIA GPUs. (True) 1/1
- OpenCL has better memory management capabilities than CUDA. (False) 1/1
- OpenCL has a more diverse range of supported hardware platforms, while CUDA is limited to NVIDIA GPUs. (True) 1/1
- OpenCL can work with CPUs and GPUs simultaneously, while CUDA is

р





(i) 163 words </>





Quiz saved at 12:04pm

Submit Quiz