Assignment 1

Exercise 1 - Reflections on GPU-accelerated Computing

1

- CPUs are faster and better at handling complex instructions
- GPUs have a lot more cores which means they can work on a lot of tasks in parallel
- CPUs usually have higher clock speeds and there fore can execute a single instruction faster than a GPU

2

Rank	Name	GPU model	Rpeak / Power (TFlops/kW)
1	Frontier	AMD Instinct MI250X	73.95
2	Fugaku	Fujitsu A64FX	14.78
3	LUMI	AMD Instinct MI250X	51.36
4	Leonardo	Nvidia Ampere A100	32.14
5	Summit	Nvidia Tesla V100 / Volta GV100	14.66
6	Sierra	Nvidia Tesla V100 / Volta GV100	12.64
7	Sunway TaihuLight	-	6.05
8	Perlmutter	Nvidia Ampere A100	27.04
9	Selene	Nvidia Ampere A100	23.81
10	Tianhe-2A	Matrix-2000	3.30

- 9 out of 10 have a GPU
- Out of those 5 are by Nvidia 2 by AMD and one each from Fujitsu and Matrix
- source: https://www.top500.org/lists/top500/2023/06/

Exercise 2 - Query Nvidia GPU Compute Capability

!./deviceQuery

→ ./deviceQuery Starting...

CUDA Device Query (Runtime API) version (CUDART static linking)

Detected 1 CUDA Capable device(s)

Device 0: "Tesla T4" CUDA Driver Version / Runtime Version 12.0 / 11.8 7.5 CUDA Capability Major/Minor version number: 15102 MBytes (15835398144 bytes) Total amount of global memory: (040) Multiprocessors, (064) CUDA Cores/MP: 2560 CUDA Cores GPU Max Clock rate: 1590 MHz (1.59 GHz) Memory Clock rate: 5001 Mhz Memory Bus Width: L2 Cache Size: 256-bit 4194304 bytes 1D=(131072), 2D=(131072, 65536), 3D=(16384, 16384, 16384) 1D=(32768), 2048 layers Maximum Texture Dimension Size (x,y,z) Maximum Layered 1D Texture Size, (num) layers Maximum Layered 2D Texture Size, (num) layers 2D=(32768, 32768), 2048 layers Total amount of constant memory: 65536 bytes Total amount of shared memory per block: 49152 bytes Total shared memory per multiprocessor: 65536 bytes Total number of registers available per block: 65536 Warp size: 32 Maximum number of threads per multiprocessor: 1024 Maximum number of threads per block: 1024 Max dimension size of a thread block (x,y,z): (1024, 1024, 64) (x,y,z): (2147483647, 65535, 65535) Max dimension size of a grid size 2147483647 bytes Maximum memory pitch: Texture alignment: 512 bytes Concurrent copy and kernel execution: Yes with 3 copy engine(s) Run time limit on kernels: No Integrated GPU sharing Host Memory: No Support host page-locked memory mapping: Yes Alignment requirement for Surfaces: Yes Device has ECC support: Enabled Device supports Unified Addressing (UVA): Yes Device supports Managed Memory: Yes Device supports Compute Preemption: Yes Supports Cooperative Kernel Launch: Yes

Compute Mode: < Default (multiple host threads can use ::cudaSetDevice() with device simultaneously) >

Yes 0 / 0 / 4

deviceQuery, CUDA Driver = CUDART, CUDA Driver Version = 12.0, CUDA Runtime Version = 11.8, NumDevs = 1 Result = PASS

The compute capability is 7.5.

• Memory Clock rate: 5001 Mhz

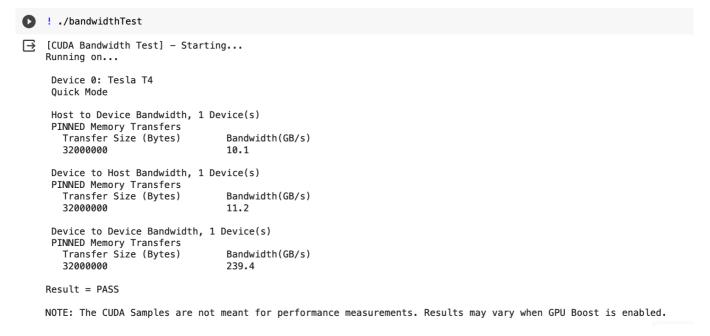
Supports MultiDevice Co-op Kernel Launch:

Device PCI Domain ID / Bus ID / location ID:

• Bus width: 256-bit

• DDR: 2

Memory bandwith: 5001 * 256 / 8 / 1024 * 2 = 312 GB/s



The bandiwth from the test is lower than the value calculated.

Exercise 3 - Rodinia CUDA benchmarks and Comparison with CPU

I've carried out experiments the heartwall and k-means tests. Changing the Makefile was not necessary however I removed the unnecessary file complitations to simplify the process and shorten the compilcation times. The two selected tasks are imaging and data mining tasks which can be done in parallel. Looking at the results in the below figures it can be seen that for both tests the execution time is more than twice as fast on a GPU compared to a CPU. While both executions can be parallelised the OMP code has been run using a single thread while the CUDA code uses 256 threads per block for the imaging task. This is means that it can do many more calculations at the same time.

```
\frac{\checkmark}{0s} [109] ! time ./heartwall ../../data/heartwall/test.avi 10
       WG size of kernel = 256
       frame progress: 0 1 2 3 4 5 6 7 8 9
                0m0.435s
       user
               0m0.136s
       sys
               0m0.249s
[110] ! nvprof ./heartwall ../../data/heartwall/test.avi 10
       WG size of kernel = 256
       ==45201== NVPROF is profiling process 45201, command: ./heartwall ../../data/heartwall/test.avi 10
       frame progress: 0 1 2 3 4 5 6 7 8 9
       ==45201== Profiling application: ./heartwall ../../../data/heartwall/test.avi 10
       ==45201== Profiling result:
                                                   lls Avg Min Max
10 12.225ms 13.024us 13.590ms
                   Type
                          Time(%)
                                       Time
                                                 Calls
                                                                                       Name
        GPU activities:
                           97.96% 122.25ms
                                                                                       kernel(void)
                                                   26 97.669us
                                                                  1.3120us
                                                                                       [CUDA memcpy HtoD]
[CUDA memcpy DtoH]
                                   2.5394ms
                            2.03%
                                                                            273.37us
                            0.01%
                                   10.784us
                                                     4 2.6960us
                                                                  2.5920us
                                                                            2.8160us
             API calls:
                                   231.58ms
                                                   623 371.73us
                                                                  2.0820us
                                                                            228.13ms
                                                                                       cudaMalloc
                           70.43%
                           27.93%
                                                   18 5.1015ms
                                   91.827ms
                                                                  4.2070us
                                                                            13.590ms
                                                                                       cudaMemcpy
                            1.30%
                                   4.2592ms
                                                   623
                                                        6.8360us
                                                                  2.6220us
                                                                            148.27us
                                                                                       cudaFree
                            0.24%
                                   792.45us
                                                   12
                                                        66.037us
                                                                  10.028us
                                                                            81.935us
                                                                                       cudaMemcpyToSymbol
                            0.05%
                                   177.59us
                                                    10
                                                       17.759us
                                                                  13.920us
                                                                            34.008us
                                                                                       cudaLaunchKernel
                            0.04%
                                   121.13us
                                                   101
                                                        1.1990us
                                                                     133ns
                                                                            49.039us
                                                                                       cuDeviceGetAttribute
                            0.01%
                                   24.702us
                                                    1
                                                        24.702us
                                                                  24.702us
                                                                            24.702us
                                                                                       cuDeviceGetName
                            0.00%
                                   6.4520us
                                                     1
                                                        6.4520us
                                                                  6.4520us
                                                                            6.4520us
                                                                                       cuDeviceGetPCIBusId
                                                        1.3040us
                            0.00%
                                   2.6080us
                                                                            2.4430us
                                                     2
                                                                     165ns
                                                                                       cuDeviceGet
                            0.00%
                                   1.9790us
                                                     3
                                                           659ns
                                                                     224ns
                                                                            1.5300us
                                                                                       cuDeviceGetCount
                            0.00%
                                                           495ns
                                                                     495ns
                                                                                495ns
                                      495ns
                                                                                       cuModuleGetLoadingMode
                                                     1
                            0.00%
                                      471ns
                                                           471ns
                                                                     471ns
                                                                                471ns
                                                                                       cuDeviceTotalMem
                                                     1
                            0.00%
                                      232ns
                                                     1
                                                           232ns
                                                                     232ns
                                                                                232ns
                                                                                       cuDeviceGetUuid
[111] %cd /content/drive/MyDrive/DD2360/rodinia_3.1/bin/linux/omp
       /content/drive/MyDrive/DD2360/rodinia_3.1/bin/linux/omp
(64) ! chmod +x heartwall
112] ! time ./heartwall ../../../data/heartwall/test.avi 10 1
       num of threads: 1 frame progress: 0 1 2 3 4 5 6 7 8 9
       real
               0m11.396s
               0m11.193s
       user
               0m0.051s
       svs
```

```
[102] ! time ./kmeans -i /content/drive/MyDrive/DD2360/rod/rodinia_3.1/data/kmeans/819200.txt
          I/O completed
          Number of objects: 819200
Number of features: 34
iterated 2 times
Number of Iteration: 1
           real
                      0m2.108s
[107] ! nvprof ./kmeans -i /content/drive/MyDrive/DD2360/rod/rodinia_3.1/data/kmeans/819200.txt
           ==44770== NVPROF is profiling process 44770, command: ./kmeans -i /content/drive/MyDrive/DD2360/rod/rodinia_3.1/data/kmeans/819200.txt
          I/O completed
          Number of objects: 819200
Number of features: 34
iterated 2 times
           Number of Iteration: 1
          ==44770== Profiling application: ./kmeans -i /content/drive/MyDrive/DD2360/rod/rodinia_3.1/data/kmeans/819200.txt
==44770== Profiling result:
    Type Time(%) Time Calls Avg Min Max Name
                                                                               Avg
5.0637ms
5.6269ms
2.2615ms
            Type GPU activities:
                                      69.40% 25.318ms
15.42% 5.6269ms
12.40% 4.5230ms
                                                                                              1.3750us
5.6269ms
2.2563ms
                                                                                                             24.247ms
5.6269ms
2.2668ms
                                                                                                                             Name
[CUDA memcpy HtoD]
invert_mapping(float*, float*, int, int)
kmeansPoint(float*, int, int, int*, float*, float*, int*)
                                        2.77%
                                                  1.0115ms
212.60ms
                                                                           2 505.74us
4 53.149ms
                                                                                              408.25us
                                                                                                              603.23us
                                                                                                                             [CUDA memcpy DtoH] cudaMalloc
                   API calls:
                                      84.06%
                                                                                               72.386us
                                                                                                              212.22ms
                                                                              53.149ms
4.4185ms
2.2651ms
856.69us
1.1321ms
1.1560us
31.171us
5.0180us
25.208us
2.3290us
1.5660us
                                       12.23%
1.79%
1.36%
                                                  30.929ms
4.5301ms
3.4268ms
1.1321ms
                                                                                                                             cudaMemcpy
cudaThreadSynchronize
                                                                                               71.935us
                                                                                                              24.445ms
                                                                                               2.2598ms
231.07us
1.1321ms
                                                                                                              2.2703ms
1.1591ms
                                                                                                                             cudaFree
cuDeviceGetPCIBusId
                                        0.45%
                                                                                                              1.1321ms
                                                  116.79us
                                        0.05%
                                                                        101
                                                                                                   149ns
                                                                                                              48.870us
                                                                                                                             cuDeviceGetAttribute
                                                  93.515us
30.113us
25.208us
6.9870us
                                                                                                              34.190us
17.137us
25.208us
2.6780us
                                                                                               28.102us
                                        0.04%
                                                                                                                             cudal aunchKernel
                                                                                               1.3350us
25.208us
2.0910us
                                        0.01%
                                                                                                                             cudaBindTexture
                                        0.01%
0.01%
0.00%
0.00%
                                                                                                                             cuDeviceGetName
cudaSetDevice
                                                  3.1930us
                                                                                1.5960us
                                                                                              1.4670us
                                                                                                              1.7260us
                                                                                                                             cudaMemcpyToSymbol
                                        0.00%
                                                  1.4550us
                                                                                     485ns
                                                                                                   195ns
                                                                                                              1.0340us
                                                                                                                             cuDeviceGetCount
                                                                                                   117ns
170ns
604ns
377ns
247ns
                                                                                                                  378ns
767ns
604ns
377ns
                                                                                                                             cudaCreateChannelDesc
cuDeviceGet
cuModuleGetLoadingMode
                                        0.00%
                                                  1.2140us
                                                                                     202ns
                                                       937ns
604ns
377ns
247ns
                                        0.00%
                                                                                     604ns
377ns
                                        0.00%
                                                                                                                             cuDeviceTotalMem
                                                                                                                  247ns
                                                                                                                             cuDeviceGetUuid
/<sub>la</sub> [103] %cd /content/drive/MyDrive/DD2360/rod/rodinia_3.1/bin/linux/omp
          /content/drive/MyDrive/DD2360/rod/rodinia_3.1/bin/linux/omp
[89] !chmod +x kmeans
   [104] ! time ./kmeans -i /content/drive/MyDrive/DD2360/rod/rodinia_3.1/data/kmeans/819200.txt
          I/O completed
           num of threads = 1
          number of Clusters 5
number of Attributes 34
          Time for process: 3.201367
           real
                      0m4.910s
```

Exercise 4 - Run a HelloWorld on AMD GPU

To launch the code on the AMD GPUs in Dardel. First, one needs to get a GPU allocation. Then the executable created can be launched on the GPU using the srun command and specifying which node the program should run on with the -n flag.

```
[aorucu@uan01:~/Private/hw1> srun -n 1 ./HelloWorld
   System minor 0
   System major 9
   agent prop name
   input string:
   GdkknVnqkc

output string:
HelloWorld
Passed!
aorucu@uan01:~/Private/hw1> []
```