

## Task 2

1. The screenshot of the output from running deviceQuery test in /1\_Uilities.

```
! ./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
  CUDA Capability Major/Minor version number: 7.5
  Total amount of global memory:             15102 MBytes (15835398144 bytes)
  (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:                         256-bit
  L2 Cache Size:                            4194304 bytes
  Maximum Texture Dimension Size (x,y,z)    1D=(131072), 2D=(131072, 65536), 3D=(16384, 16384, 16384)
  Maximum Layered 1D Texture Size, (num) layers 1D=(32768), 2048 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)
  Max dimension size of a grid size (x,y,z): (2147483647, 65535, 65535)
  Maximum memory pitch:                     2147483647 bytes
  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
  Supports MultiDevice Co-op Kernel Launch:  Yes
  Device PCI Domain ID / Bus ID / location ID: 0 / 0 / 4
  Compute Mode:
    < Default (multiple host threads can use ::cudaSetDevice() with device simultaneously) >

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

2. What is the Compute Capability of your GPU device?

7.5

3. The screenshot of the output from running bandwidthTest test in /1\_Uilities.

```
!./bandwidthTest
```

```
[CUDA Bandwidth Test] - Starting...  
Running on...
```

```
Device 0: Tesla T4  
Quick Mode
```

```
Host to Device Bandwidth, 1 Device(s)  
PINNED Memory Transfers  
Transfer Size (Bytes)      Bandwidth(GB/s)  
32000000                  11.4
```

```
Device to Host Bandwidth, 1 Device(s)  
PINNED Memory Transfers  
Transfer Size (Bytes)      Bandwidth(GB/s)  
32000000                  10.4
```

```
Device to Device Bandwidth, 1 Device(s)  
PINNED Memory Transfers  
Transfer Size (Bytes)      Bandwidth(GB/s)  
32000000                  239.4
```

```
Result = PASS
```

NOTE: The CUDA Samples are not meant for performance measurements. Results may vary when GPU Boost is enabled.

4. How will you calculate the GPU memory bandwidth (in GB/s) using the output from deviceQuery? (Hint: memory bandwidth is typically determined by clock rate and bus width, and check what double data rate (DDR) may impact the bandwidth). Are they consistent with your results from bandwidthTest?

With DDR :  $2 * 256 * 5001 * 10^6 = 320\text{GB/s}$

Without DDR :  $160\text{GB/s}$

Both figures have a about  $100\text{ GB/s}$  difference in bandwidth with the Device to Device Bandwidth from the test