

# Module 2 Lab

## Device Query

*GPU Teaching Kit – Accelerated Computing*

### OBJECTIVE

The purpose of this lab is to introduce the student to the CUDA hardware resources along with their capabilities. The student is not expected to understand the details of the code, but should understand the process of compiling and running code that will be used in subsequent modules.

### INSTRUCTIONS

The code provided queries the GPU hardware on the system. Do not concentrate on the API calls, but on functions starting with `wb`.

The `wbLog` function logs results, specifically we log the following hardware features:

- GPU card's name
- GPU computation capabilities
- Maximum number of block dimensions
- Maximum number of grid dimensions
- Maximum size of GPU memory
- Amount of constant and share memory
- Warp size

The details are explained in the first and second modules of the teaching kit.

### LOCAL SETUP INSTRUCTIONS

The most recent version of source code for this lab along with the build-scripts can be found on the [Bitbucket repository](#). A description on how to use the [CMake](#) tool in along with how to build the labs for local development found in the [README](#) document in the root of the repository.

The executable generated as a result of compiling the lab can be run using the following command:

```
./DeviceQuery_Template
```

## QUESTIONS

(1) What is the compute capability of the NVIDIA Fermi architecture?

ANSWER: **2.1.**

(2) What are the maximum block dimensions for GPUs with 3.0 compute capability?

ANSWER: **1024 in X, 64 in Y and Z.**

(3) Suppose you are launching a one dimensional grid and block. If the hardware's maximum grid dimension is 65535 and the maximum block dimension is 512, what is the maximum number threads can be launched on the GPU?

ANSWER:  **$65535 * 512 = 33553920$**

(4) Under what conditions might a programmer choose not want to launch the maximum number of threads?

ANSWER: **There may not be enough data to effectively use the maximum number of threads.**

(5) What can limit a program from launching the maximum number of threads on a GPU?

ANSWER: **A program may be limited by shared memory or registers used across all threads.**

(6) What is shared memory?

ANSWER: **Fast read-write memory shared between all threads in a block.**

(7) What is global memory?

ANSWER: **Read-write memory that can be accessed by any thread.**

(8) What is constant memory?

ANSWER: **Fast read-only memory for values that do not change during kernel execution.**

(9) What does warp size signify on a GPU?

ANSWER: **The number of threads that execute in lockstep.**

(10) Is double precision supported on GPUs with 1.3 compute capability?

ANSWER: **Double precision is supported starting with CC 1.3.**

## CODE TEMPLATE

The following code is suggested as a starting point for students. The code handles the import and export as well as the checking of the solution. Students are expected to insert their code in the sections demarcated with `//@@`. Students are expected to keep the other code unchanged. The tutorial page describes the functionality of the `wb*` methods.

```

1  #include <wb.h>
2
3  /// The purpose of this code is to become familiar with the submission
4  /// process. Do not worry if you do not understand all the details of
5  /// the code.
6
7  int main(int argc, char **argv) {
8      int deviceCount;
9
10     wbArg_read(argc, argv);
11
12     cudaGetDeviceCount(&deviceCount);
13
14     wbTime_start(GPU, "Getting GPU Data."); /// start a timer
15
16     for (int dev = 0; dev < deviceCount; dev++) {
17         cudaDeviceProp deviceProp;
18
19         cudaGetDeviceProperties(&deviceProp, dev);
20
21         if (dev == 0) {
22             if (deviceProp.major == 9999 && deviceProp.minor == 9999) {
23                 wbLog	TRACE, "No CUDA GPU has been detected");
24                 return -1;
25             } else if (deviceCount == 1) {
26                 /// WbLog is a provided logging API (similar to Log4J).
27                 /// The logging function wbLog takes a level which is either
28                 /// OFF, FATAL, ERROR, WARN, INFO, DEBUG, or TRACE and a
29                 /// message to be printed.
30                 wbLog	TRACE, "There is 1 device supporting CUDA");
31             } else {
32                 wbLog	TRACE, "There are ", deviceCount,
33                     " devices supporting CUDA");
34             }
35         }
36
37         wbLog	TRACE, "Device ", dev, " name: ", deviceProp.name);
38         wbLog	TRACE, " Computational Capabilities: ", deviceProp.major, ".",
39             deviceProp.minor);
40         wbLog	TRACE, " Maximum global memory size: ",
41             deviceProp.totalGlobalMem);
42         wbLog	TRACE, " Maximum constant memory size: ",
43             deviceProp.totalConstMem);
44         wbLog	TRACE, " Maximum shared memory size per block: ",
45             deviceProp.sharedMemPerBlock);
46         wbLog	TRACE, " Maximum block dimensions: ",
47             deviceProp.maxThreadsDim[0], " x ", deviceProp.maxThreadsDim[1],
48             " x ", deviceProp.maxThreadsDim[2]);
49         wbLog	TRACE, " Maximum grid dimensions: ", deviceProp.maxGridSize[0],
50             " x ", deviceProp.maxGridSize[1], " x ",
51             deviceProp.maxGridSize[2]);
52         wbLog	TRACE, " Warp size: ", deviceProp.warpSize);
53     }

```

```

54
55     wbTime_stop(GPU, "Getting GPU Data."); /// stop the timer
56
57     return 0;
58 }

```

## CODE SOLUTION

The following is a possible implementation of the lab. This solution is intended for use only by the teaching staff and should not be distributed to students.

```

1  #include <wb.h>
2
3  /// The purpose of this code is to become familiar with the submission
4  /// process. Do not worry if you do not understand all the details of
5  /// the code.
6
7  int main(int argc, char **argv) {
8      int deviceCount;
9
10     wbArg_read(argc, argv);
11
12     cudaGetDeviceCount(&deviceCount);
13
14     wbTime_start(GPU, "Getting GPU Data."); /// start a timer
15
16     for (int dev = 0; dev < deviceCount; dev++) {
17         cudaDeviceProp deviceProp;
18
19         cudaGetDeviceProperties(&deviceProp, dev);
20
21         if (dev == 0) {
22             if (deviceProp.major == 9999 && deviceProp.minor == 9999) {
23                 wbLog	TRACE, "No CUDA GPU has been detected";
24                 return -1;
25             } else if (deviceCount == 1) {
26                 /// WbLog is a provided logging API (similar to Log4J).
27                 /// The logging function wbLog takes a level which is either
28                 /// OFF, FATAL, ERROR, WARN, INFO, DEBUG, or TRACE and a
29                 /// message to be printed.
30                 wbLog	TRACE, "There is 1 device supporting CUDA";
31             } else {
32                 wbLog	TRACE, "There are ", deviceCount,
33                     " devices supporting CUDA";
34             }
35         }
36
37         wbLog	TRACE, "Device ", dev, " name: ", deviceProp.name);
38         wbLog	TRACE, " Computational Capabilities: ", deviceProp.major, ".",
39             deviceProp.minor);
40         wbLog	TRACE, " Maximum global memory size: ",

```

```

41         deviceProp.totalGlobalMem);
42     wbLog(TRACE, " Maximum constant memory size: ",
43           deviceProp.totalConstMem);
44     wbLog(TRACE, " Maximum shared memory size per block: ",
45           deviceProp.sharedMemPerBlock);
46     wbLog(TRACE, " Maximum block dimensions: ",
47           deviceProp.maxThreadsDim[0], " x ", deviceProp.maxThreadsDim[1],
48           " x ", deviceProp.maxThreadsDim[2]);
49     wbLog(TRACE, " Maximum grid dimensions: ", deviceProp.maxGridSize[0],
50           " x ", deviceProp.maxGridSize[1], " x ",
51           deviceProp.maxGridSize[2]);
52     wbLog(TRACE, " Warp size: ", deviceProp.warpSize);
53 }
54
55 wbTime_stop(GPU, "Getting GPU Data."); //@@ stop the timer
56
57 return 0;
58 }

```