

# *CUDA SDK — Getting Started*

These instructions are for setting up and getting started with CUDA programming in the LG04 school lab on Linux systems.

All machines in the lab have a GeForce RTX 2060. Please note that these machines do not, and should be modified to, use those GPUs to drive their monitors. That would interfere with using the GPUs for general purpose parallel programming.

## **1. Environment**

To set the necessary environment variables (on the School computers only) run:

```
module load cuda
```

When you load the `cuda` module, it will display two directories. The second, the examples directory, is where the computer officers have downloaded and compiled a large set of sample programs that come with the CUDA SDK. The binaries can be found below the `bin` sub-directory. Run the sample program `deviceQuery`. Save the results as you will be referring to it often as you develop CUDA programs.

Here is the output from `deviceQuery` for the GeForce GTX 960:

```
Device 0: "GeForce RTX 2060"
  CUDA Driver Version / Runtime Version      10.1 / 10.1
  CUDA Capability Major/Minor version number: 7.5
  Total amount of global memory:             5905 MBytes (6191382528 bytes)
  (30) Multiprocessors, ( 64) CUDA Cores/MP: 1920 CUDA Cores
  GPU Max Clock rate:                       1680 MHz (1.68 GHz)
  Memory Clock rate:                        7001 Mhz
  Memory Bus Width:                         192-bit
  L2 Cache Size:                            3145728 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 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:                     Disabled
  Device supports Unified Addressing (UVA):   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 / 1 / 0
  Compute Mode:
    < Default (multiple host threads can use ::cudaSetDevice() with device simultaneously
```

Here is the output from deviceQuery for a previous generation GeForce GTX 960, which is referred to in the slide decks:

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

Detected 1 CUDA Capable device(s)

Device 0: "GeForce GTX 960"

|  |  |
|--|--|
| CUDA Driver Version / Runtime Version          | 8.0 / 7.0  |
| CUDA Capability Major/Minor version number:    | 5.2  |
| Total amount of global memory:                 | 1996 MBytes (2092957696 bytes)                       |
| ( 8) Multiprocessors, (128) CUDA Cores/MP:     | 1024 CUDA Cores                                      |
| GPU Max Clock rate:                            | 1291 MHz (1.29 GHz)                                  |
| Memory Clock rate:                             | 3600 Mhz   |
| Memory Bus Width:                              | 128-bit  |
| L2 Cache Size:                                 | 1048576 bytes  |
| Maximum Texture Dimension Size (x,y,z)         | 1D=(65536), 2D=(65536, 65536), 3D=(4096, 4096, 4096) |
| Maximum Layered 1D Texture Size, (num) layers  | 1D=(16384), 2048 layers                              |
| Maximum Layered 2D Texture Size, (num) layers  | 2D=(16384, 16384), 2048 layers                       |
| Total amount of constant memory:               | 65536 bytes  |
| Total amount of shared memory per block:       | 49152 bytes  |
| Total number of registers available per block: | 65536  |
| Warp size:                                     | 32   |
| Maximum number of threads per multiprocessor:  | 2048   |
| 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 2 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:                        | Disabled   |
| Device supports Unified Addressing (UVA):      | Yes  |
| Device PCI Domain ID / Bus ID / location ID:   | 0 / 1 / 0  |
| Compute Mode:                                  |  |

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

deviceQuery, CUDA Driver = CUDART, CUDA Driver Version = 8.0, CUDA Runtime Version = 7.0, 1 Device(s)  
Result = PASS

Here is the output from deviceQuery for a different GPU: a GeForce GTX 660. Compare the above values to those for this one:

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

Detected 1 CUDA Capable device(s)

Device 0: "GeForce GTX 660"

```
CUDA Driver Version / Runtime Version      7.5 / 7.5
CUDA Capability Major/Minor version number: 3.0
Total amount of global memory:             2047 MBytes (2146762752 bytes)
( 5) Multiprocessors, (192) CUDA Cores/MP: 960 CUDA Cores
GPU Max Clock rate:                       1098 MHz (1.10 GHz)
Memory Clock rate:                        3004 Mhz
Memory Bus Width:                         192-bit
L2 Cache Size:                            393216 bytes
Maximum Texture Dimension Size (x,y,z)    1D=(65536), 2D=(65536, 65536), 3D=(4096, 4096, 4096)
Maximum Layered 1D Texture Size, (num) layers 1D=(16384), 2048 layers
Maximum Layered 2D Texture Size, (num) layers 2D=(16384, 16384), 2048 layers
Total amount of constant memory:           65536 bytes
Total amount of shared memory per block:   49152 bytes
Total number of registers available per block: 65536
Warp size:                                 32
Maximum number of threads per multiprocessor: 2048
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 1 copy engine(s)
Run time limit on kernels:                Yes
Integrated GPU sharing Host Memory:        No
Support host page-locked memory mapping:   Yes
Alignment requirement for Surfaces:        Yes
Device has ECC support:                   Disabled
Device supports Unified Addressing (UVA):   Yes
Device PCI Domain ID / Bus ID / location ID: 0 / 1 / 0
Compute Mode:
    < Default (multiple host threads can use ::cudaSetDevice() with device simultaneously) >
```

```
eviceQuery, CUDA Driver = CUDART, CUDA Driver Version = 7.5, CUDA Runtime Version = 7.5,
    NumDevs = 1, Device0 = GeForce GTX 660
Result = PASS
```

For the moment, pay particular attention to the values for :

- CUDA Capability Major/Minor version number
- Total amount of global memory
- Number of Multiprocessors
- Number of Cores
- Total amount of shared memory per block
- Warp size
- Maximum number of threads per multiprocessor
- Maximum number of threads per block
- Maximum dimension size of a thread block
- Maximum dimension size of a grid size

Investigate the meaning of these values on the Web.

## 2. The nsight IDE

nsight is a version of the Eclipse IDE modified by NVidia for CUDA project development. Like Eclipse, nsight uses a workspace directory (by default ~/cuda-workspace) and, when you create a project, offers to store your project files there. I recommend that you do **NOT** do so, but instead that you put your project files somewhere outside the workspace directory. A good approach is to make a directory called “cuda” somewhere in your school GIT repository (you can make a new repository if you wish to keep it separate from your other repositories).

Occasionally, NSIGHT can misbehave and the only way to recover may be to delete the ~/cuda-workspace directory. By not putting your project files in the ~/cuda-workspace directory you can delete it without losing your projects. Also it means you can easily work at home on your own machine sharing your project files through GIT. Note: it is possible to install the CUDA SDK on your home machine and be able to edit and compile your programs even if you

don't have a CUDA capable GPU — in this case, of course, you won't be able to run the programs locally.

First we will import one of the sample applications and compile and run it:

1. Start `nsight`
2. Select `File|New|Cuda C/C++ Project`
3. Set the project name to “`deviceQuery`”, untick “`use default location`”, set the location to a new directory “`DeviceQuery`” under your chosen `cuda` directory, select `Executable|Import CUDA Sample` and click next
4. Select “`deviceQuery`” and click next
5. This page lets you select the version of CUDA code to generate. If you have a CUDA capable GPU on the machine it will detect it and automatically set it correctly. If not you must set it manually. The “`CUDA Capability Major/Minor version number`” from the output of `deviceQuery` tells you what you should set this to to work on that GPU. Click next
6. It is possible to run `nsight` locally while compiling for a different architecture and running the result code remotely. In this case we don't wish to do that so click next
7. click finish
8. The “`deviceQuery.cpp`” file will open. Click on the hammer button in the toolbar to build the project.
9. Click on the dropdown arrow beside the green run button in the toolbar and select `Run Configurations...`
10. Double click on `C/C++ Application`. So long as you had first built the project, this should set all the parameters correctly. Click Run

### 3. Importing existing project (e.g. from GIT)

Assuming you have kept your project files in a directory called `cuda`, away from your `~/cuda-workspace` directory, then, in `nsight`, select `File|Import...|General|Existing Projects into Workspace`, select your `cuda` directory in the `Select root directory` field and click next. You can then select any projects not previously imported to import. Note that you should not try to copy your `~/cuda-workspace` directory between your home and school machines.

### 4. Installing the CUDA SDK on your own machine

I will assume you are using a (fairly) recent version of Ubuntu. Do not use the Ubuntu CUDA packages. If you have them installed remove them. With one exception follow the instructions on: <https://www.quantstart.com/articles/Installing-Nvidia-cuda-on-Ubuntu-14-04-for-Linux-GPU-Computing>

The exception is to get the latest version of the CUDA SDK available that matches your system. The website <https://developer.nvidia.com/cuda-downloads> helps you to choose the correct one.

**Be very careful installing CUDA on a machine that does not have an NVidia GPU:** it can override the video driver with an NVidia one that will not work on your machine. Check which video driver you are running before installing CUDA and make sure it is still installed afterwards. You can check which video driver you are running with the command:

```
sudo lshw -c video
```

On a non-Nvidia machine, this should have an output like this:

```
*-display
  description: VGA compatible controller
  product: Core Processor Integrated Graphics Controller
  vendor: Intel Corporation
  physical id: 2
  bus info: pci@0000:00:02.0
  version: 02
  width: 64 bits
  clock: 33MHz
  capabilities: msi pm vga_controller bus_master cap_list rom
  configuration: driver=i915 latency=0
  resources: irq:27 memory:fe000000-fe3fffff memory:d0000000-dfffffff ioport:f160 (size=64K)
```

The important bit of information is the `driver=i915`. You can find out more about that driver with:

```
modinfo i915
```

If your correct driver is not installed, switch it back. you should be able to do this by running:

```
software-properties-gtk
```

and switching to the “Additional Drivers” tab and choosing the correct driver.

In the worst case an easy fix is to simply remove the NVidia driver. On Ubuntu:

```
apt remove nvidia-3*
```

This should leave CUDA and `nsight` installed but the `nvidia` driver inactive