

Getting Started with CUDA



September 30, 2016

Outline

Is the Device Visible?

Device Drivers and Supporting Utilities

General CUDA Requirements

Hello World

The CUDA Compiler

Profiling CUDA Applications

The Device



- ▶ Before programming in CUDA, and even before installing a driver, need to check if the OS sees an NVIDIA device.
- ▶ On linux this can be as easy as `ls -l /dev/nv*` or use the `lspci` command which provides detailed information about all PCI buses and devices in the system
- ▶ On OS X goto “About This Mac” (beware graphics-switching)



- ▶ CUDA installer comes with a driver but do not blindly rely on this driver version to match your specific GPU.
- ▶ Go to the [NVIDIA drivers download page](#) to lookup current driver version for your device
- ▶ On Windows and Linux systems, the NVIDIA display driver installs the **nvidia-smi** command line utility. The [NVIDIA System Management Interface](#) is intended to aid in the management and monitoring of NVIDIA GPU devices.
- ▶ The **nvidia-smi** command line utility is the go-to tool for querying installed drivers, device information, device state, device performance (e.g. power, temp, clock speed), etc

General requirements to use CUDA on your system:

- ▶ a CUDA-capable GPU device connected via PCI
- ▶ gcc or Clang compiler and toolchain
- ▶ the [NVIDIA CUDA Toolkit](#)

Links to detailed installation guides for each OS below:

- ▶ [Windows](#)
- ▶ [Linux](#)
- ▶ [OSX](#)

```
1 #include<stdio.h>
2
3 // GPU Kernel definition
4 __global__ void sayHello(void){
5     printf("Hello World from the GPU!\n");
6 }
7
8 int main(void){
9
10     //launch kernel
11     sayHello<<<1,5>>>();
12
13     //wait for the kernel to finish
14     cudaDeviceSynchronize();
15
16     //that's all
17     return 0;
18 }
```

Listing 1: Hello world in CUDA (runtime API)

- ▶ For managing the device and organizing threads, CUDA consists of a higher-level *runtime API* and a low-level *driver API*.
- ▶ Each function of the runtime API calls are broken down into more basic driver API operations.
- ▶ While the driver API does offer more explicit control over how the GPU device is used but does not provide any performance gains over the runtime API and is more difficult to program.
- ▶ Most (all?) modern CUDA applications and libraries (e.g. cuBLAS, cuFFT) are built on the runtime API.
- ▶ Note driver API calls start with “cu” while runtime API calls start with “cuda” (e.g. `cudaDeviceSynchronize()`)

- ▶ The CUDA **nvcc** compiler is based on the widely used LLVM open source compiler infrastructure.
- ▶ CUDA programs consist of a mixture of *host* code for the CPU and *device* code for the GPU.
- ▶ The **nvcc** compiler separates the device code from the host code during the compilation process.
- ▶ The host code is standard C and is compiled with C compilers while the CUDA C device code is compiled by **nvcc**.
- ▶ During the link stage, CUDA runtime libraries (libcudart.so.4) are added for kernel procedure calls and explicit GPU device manipulation.
- ▶ Many additional details available via the [nvcc documentation](#)

- ▶ PTX which stands for “Parallel Thread eXecution” is the *intermediate representation* of the compiled GPU code that can be further compiled into native GPU microcode.
- ▶ This is the mechanism that enables CUDA applications to be “future-proof” against instruction set innovations by NVIDIA.
- ▶ The PTX code is typically compiled into hardware specific microcode in an on-demand fashion (**JIT**ted) by the CUDA driver. This online compilation process happens automatically when running CUDART applications compiled with the `--fatbin` option (default).
- ▶ The PTX code can be manually compiled into microcode using the PTX assembler **ptxas**. The resulting CUDA binary microcode is called a “cubin” (pronounced like “Cuban”).
- ▶ Cubin files can be disassembled with `cuobjdump` using the option `--dump-sass`. See the [docs](#) for more information on CUDA binary utilities.

- ▶ Both `.cubin` microcode and PTX representations of each kernel are included in the **nvcc** “fatbin” executable.
- ▶ If the executable is run on hardware which does not support any of the `.cubin` representations, the driver compiles the PTX version.
- ▶ Since PTX compilation can be time consuming, the driver caches these compiled kernels on disk for repeated invocation.
- ▶ Note that PTX code can be generated at runtime and compiled explicitly by the driver by calling `cuModuleLoadEx()`.

- ▶ As a compiler driver, **nvcc** does nothing more than set up a build environment and spawn a combination of native tools (e.g. the C compiler installed on the system) and CUDA specific command-line tools (e.g. **ptxas**)
- ▶ To compile CUDA program simply invoke **nvcc** `myprog.cu`
- ▶ Use the `--verbose` option to view the build process or `--dryrun` option to generate the build commands without actually executing them.
- ▶ There are many options for guiding the code generation. In particular the `--gpu-architecture` option for specifying which PTX version to emit and the `--gpu-code` option for specifying which version of Streaming Multiprocessor (SM) microcode (`.cubin`) to produce (`sm_1[0123]`, `sm_2[01]`, `sm_3[05]`).

- ▶ There are three key header files when programming in CUDA:
 - ▶ `cuda.h` defining types and host functions for the CUDA *driver* API.
 - ▶ `cuda_runtime_api.h` which defines types and host functions and types for the CUDA *runtime* API.
 - ▶ `cuda_runtime.h` contains a superset of definitions including everything from `cuda_runtime_api.h`, as well as built-in type definitions, function overlays for the CUDA language extensions, and device intrinsic functions.
- ▶ Notice that when compiling with **nvcc** the appropriate CUDA headers are included automatically.

- ▶ As of CUDA 5.0, **nvprof** is available to help collect timeline information from the application CPU and GPU activity (e.g. kernel execution, memory transfers, and API calls).
- ▶ The Visual Profiler, **nvvp**, displays a timeline of your application's activity on both the CPU and GPU so that you can identify opportunities for performance improvement.
- ▶ In addition, **nvvp** will analyze the application to detect potential performance bottlenecks and provide recommendations to eliminate or reduce those bottlenecks.
- ▶ Both **nvprof** and **nvvp** are powerful tools to help understand where time is being spent in an application. In many GPU based workloads it is important to understand the compute to communication ratio. That is, number of instructions per byte accessed. Most HPC workloads are bound by memory bandwidth.

- ▶ **CUDA-GDB** is the NVIDIA tool for debugging CUDA kernels running on Linux and Mac operating systems.
- ▶ CUDA-GDB is an extension to the x86-64 port of GDB, the GNU Project debugger.
- ▶ The CUDA debugger is designed to allow simultaneous debugging of both GPU and CPU code within the same application as well as supports debugging of both the CUDA driver API and/or the CUDA runtime API.
- ▶ In order to use debugging tools, compile with **nvcc -g -G** which embeds debugging information for both host and device code as well as disable optimizations so that program state can be inspected during execution.

- ▶ **CUDA-MEMCHECK** is a functional correctness checking suite included in the CUDA toolkit.
- ▶ In short, `memcheck` is a memory access error and leak detection tool.
- ▶ The `memcheck` tool can precisely detect and report out of bounds and misaligned memory accesses to global, local, shared and global atomic instructions in CUDA applications
- ▶ To use the `memcheck` tool on Linux, compile with **`nvcc -Xcompiler -rdynamic -lineinfo`**
- ▶ Similarly, on Windows compile with **`nvcc -Xcompiler /Zi -lineinfo`**
- ▶ When compiling with these additional options, **`nvcc`** generates executables that will contain sufficient metadata for `memcheck` to display helpful messages but maintain performance characteristics of the original application.

- ▶ The **runtime API** provides a variety of functions for managing GPU devices and querying their associated information.
- ▶ To figure out how many devices are visible to the host use `cudaGetDeviceCount()`
- ▶ To switch between GPU devices use the `cudaSetDevice()`
- ▶ Each GPU device properties can be queried using `cudaGetDeviceProperties()` which returns a struct of type **cudaDeviceProp** containing various information about the device.


```
1 #include <stdio.h>
2 #include <cuda_runtime.h>
3
4 int main(void){
5     int deviceCount;
6     cudaDeviceProp deviceProp;
7     cudaGetDeviceCount(&deviceCount);
8     cudaGetDeviceProperties(&deviceProp,0);
9     printf("There are %d gpu devices\n",deviceCount);
10    printf("Device %s has %f GB of global memory\n",
11           deviceProp.name,
12           deviceProp.totalGlobalMem/pow(1024.0,3)
13    );
14 }
```

Listing 2: Getting information about devices in CUDA (runtime API)

- ▶ Similarly the driver API has a variety of low-level functions for [device attributes](#).
- ▶ Before a device can be queried for attributes, the device must be initialized with `cuInit(int devId)`
- ▶ Most driver API calls require handle to the particular device which is provided by `cuDeviceGet()`
- ▶ There are explicit functions for device count, device name, and total global memory.
- ▶ All other device attributes are queried via the `cuDeviceGetAttribute()` function and a set of attribute macros.

```
1 #include<stdio.h>
2 #include<cuda.h>
3
4 int main(void){
5     int deviceCount;
6     char deviceName[256];
7     CUdevice device;
8     size_t szMem; int szProc;
9     cuInit(0);
10    cuDeviceGetCount(&deviceCount);
11    cuDeviceGet(&device,0);
12    cuDeviceGetName(deviceName,255,device);
13    cuDeviceTotalMem(&szMem,device);
14    cuDeviceGetAttribute(&szProc,
15        CU_DEVICE_ATTRIBUTE_MULTIPROCESSOR_COUNT,device);
16    printf("There are %d devices detected\n",deviceCount);
17    printf("Device %s has %f GB of global memory\n",
18        deviceName,szMem/pow(1024.0,3));
19    printf("Device multiprocessor count: %d\n",szProc);
20 }
```

Listing 3: Getting information about devices in CUDA (driver API)

- ▶ The CUDA Toolkit is prepackaged with a variety of CUDA [sample codes](#) which cover both [runtime](#) and [driver](#) APIs.
- ▶ These samples cover everything from basic bandwidth tests to using features such as Zero-Copy Memory, Asynchronous Data Transfers, Unified Virtual Addressing, Peer-to-Peer Communication, Concurrent Kernels, sharing data between CUDA and Direct3D/OpenGL graphics APIs, using CUDA with MPI and OpenMP, Image Processing, Video encode/decode, CFD, FDTD, and more.
- ▶ Typically CUDA code samples are selected for installation when installing the CUDA Toolkit.
- ▶ The code samples can be installed after toolkit installation using the `cuda-install-samples-X.X.sh` script located in the root CUDA directory.