

## GPU Computing and Accelerators: Part I





(b) Memory bandwidth

Figure: Throughput comparison of Multicore CPUs and CUDA enabled GPUs (taken from CUDA C Programming Guide)



CSC

Architecture	GFLOPS	GFLOPS/Watt	Utilization
Core i7-960	96	1.14	95%
Nvidia <sup>®</sup> GTX280	410	2.6	66%
Cell	200	5.0	88%
Nvidia <sup>®</sup> GTX480	940	5.4	70%
TI C66x DSP	74	7.4	57%

Table: Power efficient comparison of Multicore CPUs and accelerator chips (taken from Conference Poster by F. Igual and M. Ali)





Figure: Schematic of a general parallel system



## Memory Hierarchy with Accelerators

Graphics Processing Units (GPUs)



Figure: Memory configuration of a CUDA Device (taken from CUDA C Programming Guide)



#### Memory Hierarchy with Accelerators

Field Programmable Gate Arrays (FPGAs)



Figure: Comparison of CPUs and FPGA execution models.



# GPU Computing and Accelerators: Part II



**CUDA** is two things at the same time:

## 1. platform model

for the hardware implementation of general purpose graphics processing units made by the  $\mathsf{Nvidia}^{(\!R\!)}$  Corporation.

## 2. programming model

realizing the software implementation and scheduling of tasks of the parallel programs on the above hardware.



#### **Definition (thread)**

A *thread*, or more precisely *GPU-thread* is the smallest unit of data and instructions to be executed in a parallel CUDA program.



**Definition (thread)** 

A *thread*, or more precisely *GPU-thread* is the smallest unit of data and instructions to be executed in a parallel CUDA program.

In contrast to CPU-threads a task switch between GPU-threads is usually almost for free due to the special CUDA architecture.



Definition (thread)

**Basic Definitions** 

A *thread*, or more precisely *GPU-thread* is the smallest unit of data and instructions to be executed in a parallel CUDA program.

In contrast to CPU-threads a task switch between GPU-threads is usually almost for free due to the special CUDA architecture.

#### Definition (warp)

The CUDA hardware consists of streaming multi-processors that are executing several threads simultaneously. The GPU-threads are therefore grouped in so called *warps* of threads per multi-processor.



Definition (thread)

**Basic Definitions** 

A *thread*, or more precisely *GPU-thread* is the smallest unit of data and instructions to be executed in a parallel CUDA program.

In contrast to CPU-threads a task switch between GPU-threads is usually almost for free due to the special CUDA architecture.

#### Definition (warp)

The CUDA hardware consists of streaming multi-processors that are executing several threads simultaneously. The GPU-threads are therefore grouped in so called *warps* of threads per multi-processor.

The number of threads in a warp may depend on the hardware. One finds mostly 32 threads per warp which in turn is the smallest number of tasks executed in SIMD style.



A *block* is a larger group of threads that can contain 64-512 threads.



A *block* is a larger group of threads that can contain 64-512 threads.

Ideally it contains a multiple of 32 threads so it can be split optimally into warps by the CUDA environment for scheduling.



A *block* is a larger group of threads that can contain 64-512 threads.

Ideally it contains a multiple of 32 threads so it can be split optimally into warps by the CUDA environment for scheduling.

Definition (grid)

The actual work to be performed by a program or algorithm is distributed to a one or two dimensional *grid* of blocks.



A *block* is a larger group of threads that can contain 64-512 threads.

Ideally it contains a multiple of 32 threads so it can be split optimally into warps by the CUDA environment for scheduling.

#### Definition (grid)

The actual work to be performed by a program or algorithm is distributed to a one or two dimensional *grid* of blocks.

The grid represents the largest freedom in design that the developer has.



**Basic Definitions** 



Figure: Grids of Thread Blocks (taken from CUDA C programming guide)



The central notions to understand data management in a CUDA program are those of *host* and *device*. Here *host* refers to the computer that hosts the GPU. Especially the CPU and memory of the host are relevant. The *device* then is the GPU installed on the host system.



The central notions to understand data management in a CUDA program are those of *host* and *device*. Here *host* refers to the computer that hosts the GPU. Especially the CPU and memory of the host are relevant. The *device* then is the GPU installed on the host system.

In case multiple GPUs are installed on a single host system with multiple CPUs, each GPU is connected to a single CPU representing a single NUMA node of the host system.



The central notions to understand data management in a CUDA program are those of *host* and *device*. Here *host* refers to the computer that hosts the GPU. Especially the CPU and memory of the host are relevant. The *device* then is the GPU installed on the host system.

In case multiple GPUs are installed on a single host system with multiple CPUs, each GPU is connected to a single CPU representing a single NUMA node of the host system.

The host CPU controls the execution of the program. However host and device may execute their tasks asynchronously. When not specified differently data transfers between them serve as implicit synchronization points.



**Basic Definitions** 

#### **Definition (kernel)**

The *kernel* is the core element of a CUDA parallel program. It represents the function that specifies the work a certain thread in a block on a grid has to execute.

We will see in the course of this Chapter how the thread executing the kernel knows which part of the global problem it has to perform.



We will next introduce the most basic elements of the CUDA C language extension. These consist of two important things.

- 1. **qualifiers** that apply to functions and specify where the function should be executed,
- 2. **launch size specifiers** that control the grid and block sizes that are used to run a kernel.

An extensive API, defining C-style functions and data types to be used in CUDA programs, together with a handful of libraries for several kinds of tasks (e.g., a BLAS implementation) complete the picture.



#### Most Basic Syntax of the CUDA C Extension

GPU Computing Applications												
Libraries and Middleware												
CUFFT CUBLAS CULA CURAND MAGMA CUSPARSE		Thrust NPP		VSIPL SVM OpenCurrent		PhysX OptiX iRay		cuDNN TensorRT	MATLAB Mathemati	ica		
Programming Languages												
c c++			Fortran		j Py Wra	ava thon ppers	DirectCompute		Directives (e.g. OpenACC)			
CUDA-enabled NVIDIA GPUs												
	Pascal Architecture (compute capabilities 6.x)     G       Maxwell Architecture (compute capabilities 5.x)     G		GeForce 1000 Series		Quadro P Series		Tesla P Series					
			GeForce 900 Series		Quadro M Series		Tesla M Series					
Kepler Architecture (compute capabilities 3.x)		GeForce 700 Series GeForce 600 Series		Quadro K Series								
Fermi Architecture (compute capabilities 2.x)		GeForce 500 Series GeForce 400 Series										
			108	Entertainmer	nt		rofessional Graphics		Hide	ettermance mputing		

Figure: The CUDA GPU computing applications framework (taken from CUDA C programming guide)



 \_\_global\_\_ This qualifier applies to a function and is used to indicate that it in fact represents a kernel.



- \_\_global\_\_ This qualifier applies to a function and is used to indicate that it in fact represents a kernel.
- \_\_device\_\_ The qualifier that specifies functions that should be run on the device, but are not kernels. It can be useful for subtasks called in a kernel. It also applies to variables determining them to reside on the device.



- \_\_global\_\_ This qualifier applies to a function and is used to indicate that it in fact represents a kernel.
- \_\_device\_\_ The qualifier that specifies functions that should be run on the device, but are not kernels. It can be useful for subtasks called in a kernel. It also applies to variables determining them to reside on the device.
- \_\_host\_\_\_ Being basically redundant this qualifier can be used to explicitly state that a function is to be executed on the host. It is therefore optional.



- \_\_global\_\_ This qualifier applies to a function and is used to indicate that it in fact represents a kernel.
- \_\_device\_\_ The qualifier that specifies functions that should be run on the device, but are not kernels. It can be useful for subtasks called in a kernel. It also applies to variables determining them to reside on the device.
- \_\_host\_\_\_ Being basically redundant this qualifier can be used to explicitly state that a function is to be executed on the host. It is therefore optional.
- \_\_shared\_\_applies to a variable declaring that it should reside in the shared memory of a streaming multiprocessor



- \_\_global\_\_ This qualifier applies to a function and is used to indicate that it in fact represents a kernel.
- \_\_device\_\_ The qualifier that specifies functions that should be run on the device, but are not kernels. It can be useful for subtasks called in a kernel. It also applies to variables determining them to reside on the device.
- Lost\_\_\_Being basically redundant this qualifier can be used to explicitly state that a function is to be executed on the host. It is therefore optional.
- \_\_shared\_\_applies to a variable declaring that it should reside in the shared memory of a streaming multiprocessor
- \_\_constant\_\_applies to a variable specifying the residence in the constant memory.



- \_\_global\_\_ This qualifier applies to a function and is used to indicate that it in fact represents a kernel.
- \_\_device\_\_ The qualifier that specifies functions that should be run on the device, but are not kernels. It can be useful for subtasks called in a kernel. It also applies to variables determining them to reside on the device.
- \_\_host\_\_\_ Being basically redundant this qualifier can be used to explicitly state that a function is to be executed on the host. It is therefore optional.
- \_\_shared\_\_applies to a variable declaring that it should reside in the shared memory of a streaming multiprocessor
- \_\_constant\_\_applies to a variable specifying the residence in the constant memory.

Note that \_\_global\_\_ and \_\_device\_\_ functions are not allowed to be recursive.



The basic launch size specification for a kernel takes the form

```
<<< grid , block size >>>
```

where grid specifies the block distribution and block size indicates the number of threads per block in the grid.

- <<<1, 1>>> launches 1 block with 1 thread
- $\blacksquare$  <<<N, 1>>> launches N blocks with 1 thread each
- $\blacksquare$  <<<1, N>>> launches 1 block with N threads
- $\blacksquare$  <<<N, M>>> launches a 1d grid of N blocks running M threads each



Both the arguments can be two dimensional distributions. CUDA defines special tuple hiding types for these declarations. Using

```
dim3 grid(3,2)
dim3 threads(16,16)
```

one defines a  $3 \times 2$  grid of blocks for running 256 threads arranged in a  $16 \times 16$  local grid. These are then used in the launch specification as

```
<<< grid, threads>>>
```

Launch size specifications are simply appended to the kernel function name upon calling it.



The following examples are taken from the "CUDA by Example" book.

```
#include "../common/book.h"
__global__ void kernel( void ) {
    int main( void ) {
        kernel<<<1,1>>>();
        printf( "Hello, _World!\n" );
        return 0;
}
```



Introductory Examples

```
#include "../common/book.h"
__global__ void add( int a, int b, int *c ) {
    *c = a + b;
int main( void ) {
   int c;
    int *dev c;
    HANDLE ERROR ( cudaMalloc( (void **) & dev c, sizeof(int) ) );
    add<<<1,1>>>(2,7, dev c);
    HANDLE_ERROR ( cudaMemcpy ( &c, dev_c, sizeof(int),
                               cudaMemcpyDeviceToHost ) );
    printf( "2_+_7_=_%d\n", c );
    HANDLE ERROR ( cudaFree ( dev c ) );
    return 0;
```



Introductory Examples

```
#include ".../common/book.h"
 device int addem( int a, int b ) {
    return a + b;
 global void add( int a, int b, int *c ) {
    *c = addem(a, b);
int main( void ) {
    int c;
    int *dev c;
    HANDLE ERROR ( cudaMalloc( (void **) & dev c, sizeof(int) ) );
    add<<<1,1>>>(2,7, dev c);
    HANDLE ERROR ( cudaMemcpy ( &c, dev c, sizeof(int),
                              cudaMemcpyDeviceToHost ) );
    printf( "2, +, 7, =, %d\n", c );
    HANDLE ERROR ( cudaFree ( dev c ) );
    return 0;
```



In order to be able to compile the previous examples, one needs to check a few prerequisites:

- NVIDIA<sup>®</sup> device drivers and hardware,
- NVIDIA<sup>®</sup> CUDA toolkit installation,
- compiler for the host code.



Basic information on CUDA in general can be found at http://www.nvidia.com/cuda.



Basic information on CUDA in general can be found at http://www.nvidia.com/cuda. The Toolkit and all the information on the included accelerated libraries and developer tools can be found at https://developer.nvidia.com/cuda-toolkit.



Basic information on CUDA in general can be found at http://www.nvidia.com/cuda. The Toolkit and all the information on the included accelerated libraries and developer tools can be found at https://developer.nvidia.com/cuda-toolkit.

Regarding the hardware, basically every NVIDIA<sup>®</sup>GPU released after the appearance of the GeForce 8800 GTX in 2006 is CUDA enabled. However, one needs to make sure that the OS version, the device driver and CUDA Toolkit version are fitting. Working combinations should be available in the toolkits documentation.



Basic information on CUDA in general can be found at http://www.nvidia.com/cuda. The Toolkit and all the information on the included accelerated libraries and developer tools can be found at https://developer.nvidia.com/cuda-toolkit.

Regarding the hardware, basically every NVIDIA<sup>®</sup> GPU released after the appearance of the GeForce 8800 GTX in 2006 is CUDA enabled. However, one needs to make sure that the OS version, the device driver and CUDA Toolkit version are fitting. Working combinations should be available in the toolkits documentation.

Regarding the compilers  $\mathsf{NVIDIA}^{\textcircled{R}}$  recommends the following

- Microsoft Windows: Visual Studio
- **Linux:** Gnu Compiler Collection (GCC)
- MacOS: GCC as well via Apple's Xcode



We will in the following restrict ourselves to the Linux world again.

Consider our basic "Hello World!" example is stored in a text file called hello\_world.cu. Using the nvcc compiler provided in the CUDA Toolkit we can compile it by

nvcc hello\_world.cu

Since on Linux nvcc uses gcc to compile the host code this will also generate a binary called a.out. As for gcc we can specify the output filename, i.e. name of the resulting executable via

nvcc hello\_world.cu -o hello\_world

The file extension .cu is used to indicate that we have a C file with CUDA C extensions.



Among the further compiler options we meet many old friends:

- $-{\tt c}$  for generating object files of single .c or .cu files
- $-g\,$  for generating debug information in the host code
- $-{\tt pg}\,$  the same for profiling information
  - $\ensuremath{\textsc{o}}$  for specifying the optimization level for the host code
  - -m specify 32 vs 64bit host architecture



Among the further compiler options we meet many old friends:

- $-{\tt c}$  for generating object files of single .c or .cu files
- $-g\,$  for generating debug information in the host code
- -pg the same for profiling information
  - $\ensuremath{\bigcirc}$  for specifying the optimization level for the host code
  - -m specify 32 vs 64bit host architecture

And we have a few more for the device code, e.g.

- $-\ensuremath{\mathsf{G}}$  generates debug information for the device code
- arch specifies the GPU architecture to be assumed, i.e. the compute capabilities of the device (e.g. -arch=sm\_20)