



NVIDIA.COM

Introduction to utilizing the power of Nvidia GPUs

Agenda

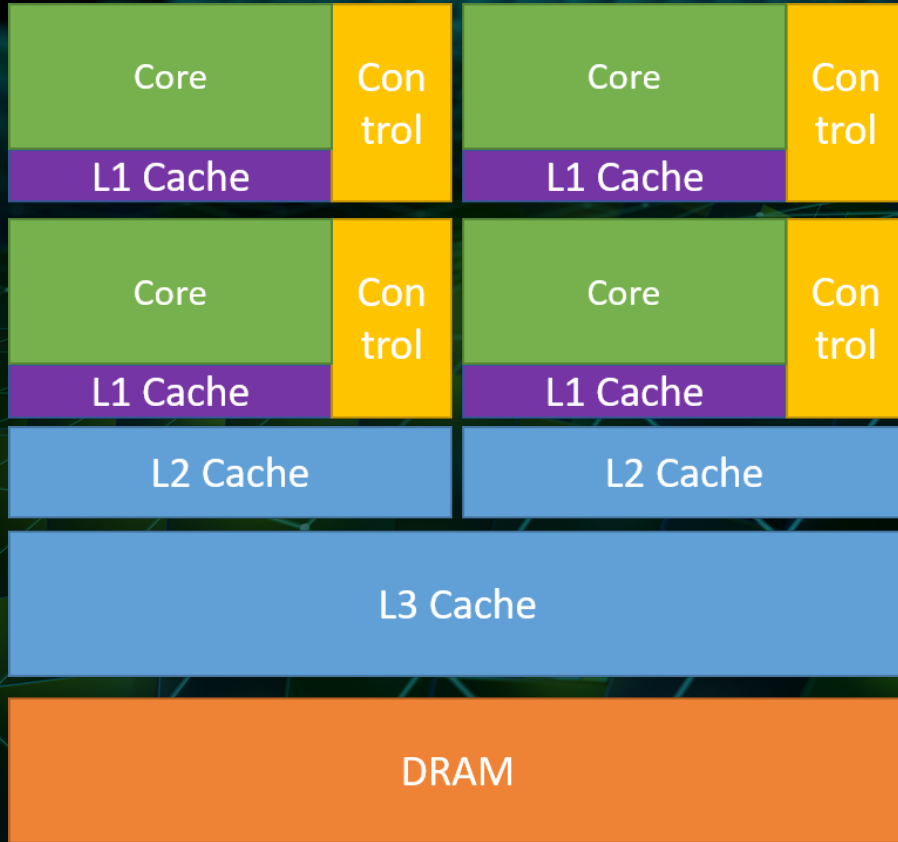
1. Introduction
 - 1.1 The Benefits of using GPUs
 - 1.2 Performance Demo
 - 1.3 Cuda, short overview
2. Programming Model
3. Implementation
 - 3.1 Kernels
 - 3.2 Thread hierarchy
 - 3.3 Heterogeneous programming
4. Usage
 - 4.1 Installation
 - 4.2 IDEs & Tools
5. Demo



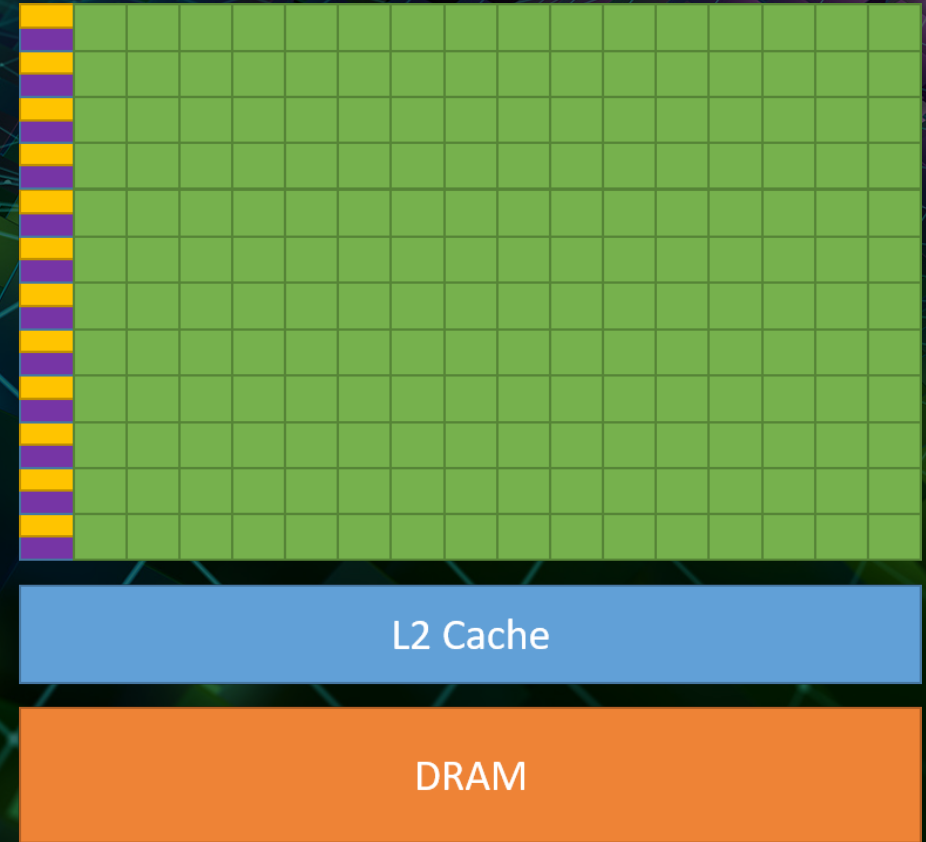
NVIDIA.COM

Introduction

The Benefits of using GPUs



CPU



GPU

Performance Demo

Results:

- Matrix size: 3440 x 3440
 - C++: 22.7181 seconds
 - Cuda: 0.563215 seconds
- => 3934% Performance increase!

Cuda – Short Overview

- Introduced in 2006
- General purpose parallel computing platform
- Only available for Nvidia GPUs
- Use Cases:
 - Machine Learning
 - Computational Chemistry
 - Bioinformatics
 - Computational Fluid Dynamics
 - Data Science
 - Weather and Climate

GPU Computing Applications

Libraries and Middleware





cuDNN TensorRT	cuFFT cuBLAS cuRAND cuSPARSE	CUDA MAGMA	Thrust NPP	VSIP SVM OpenCurrent	PhysX OptiX iRay	MATLAB Mathematica
-------------------	---------------------------------------	---------------	---------------	----------------------------	------------------------	-----------------------

Programming Languages

C	C++	Fortran	Java Python Wrappers	DirectCompute	Directives (e.g. OpenACC)
---	-----	---------	----------------------------	---------------	------------------------------



CUDA-Enabled NVIDIA GPUs

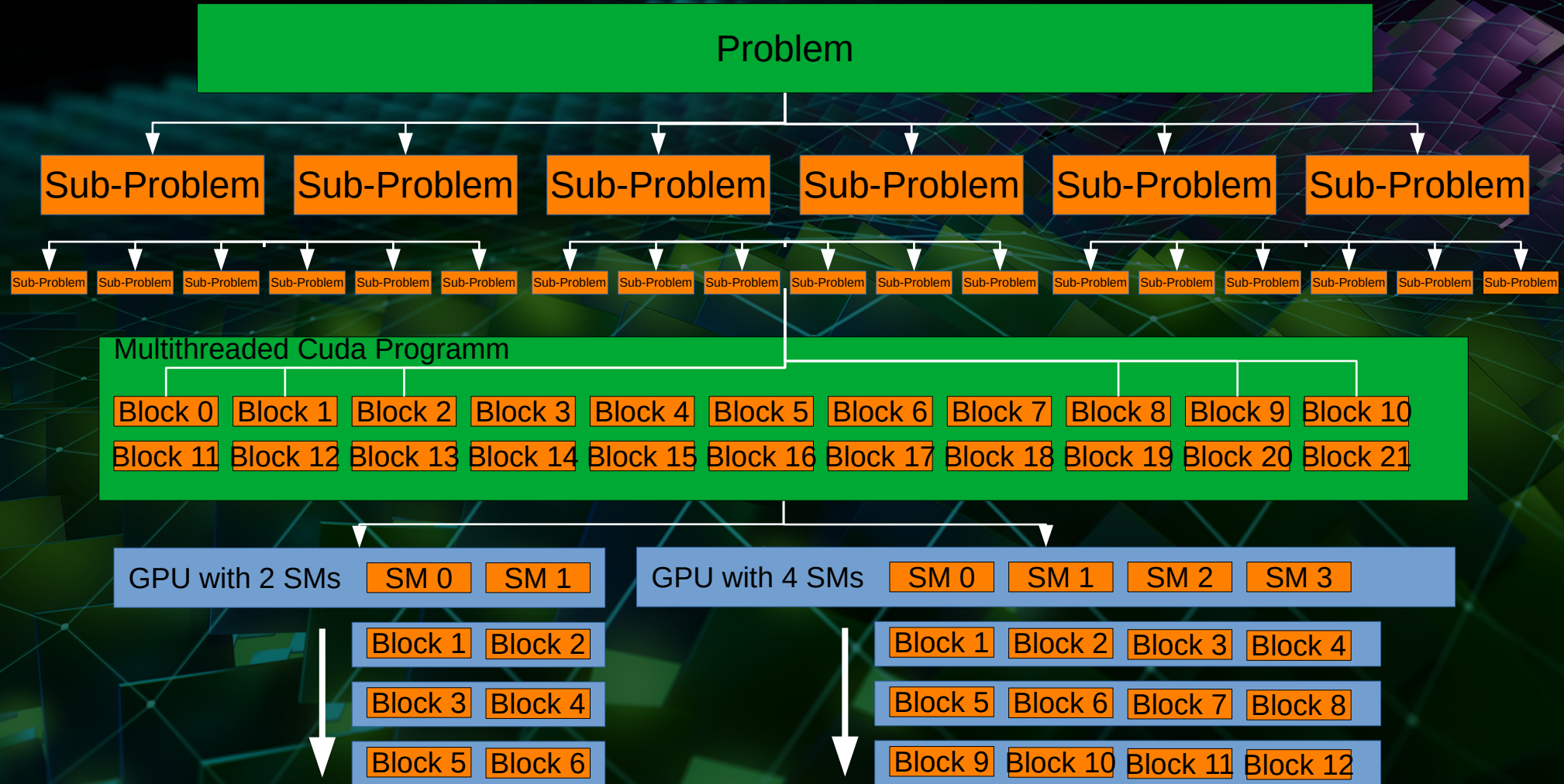
NVIDIA Ampere Architecture (compute capabilities 8.x)				Tesla A Series
NVIDIA Turing Architecture (compute capabilities 7.x)		GeForce 2000 Series	Quadro RTX Series	Tesla T Series
NVIDIA Volta Architecture (compute capabilities 7.x)	DRIVE/JETSON AGX Xavier		Quadro GV Series	Tesla V Series
NVIDIA Pascal Architecture (compute capabilities 6.x)	Tegra X2	GeForce 1000 Series	Quadro P Series	Tesla P Series
	 Embedded	 Consumer Desktop/Laptop	 Professional Workstation	 Data Center



NVIDIA.COM

Programming model

Programming model





NVIDIA.COM

Implementation

Kernels

- Kernel represents smallest possible sub-Problem
- C++ function
- Gets executed N times by N **different** Cuda threads
- Defined by the `__global__` declaration specifier:

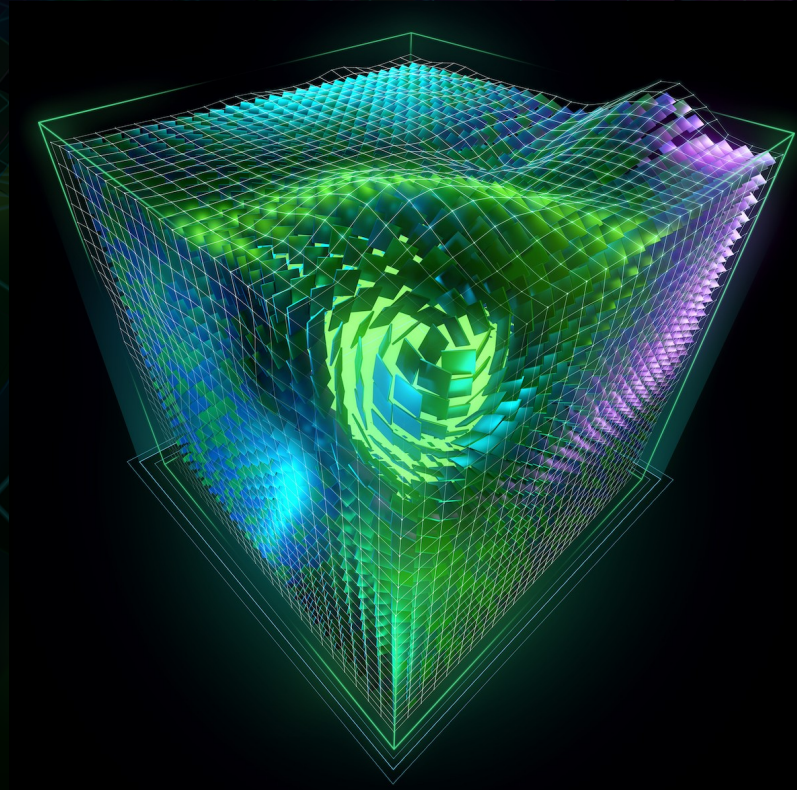
```
__global__ void VecAdd(float* A, float* B, float* C)
{
    int i = threadIdx.x;
    C[i] = A[i] + B[i];
}
```

- Execution configuration syntax:

```
// Kernel invocation with N threads
VecAdd<<<1, N>>>>(A, B, C);
```

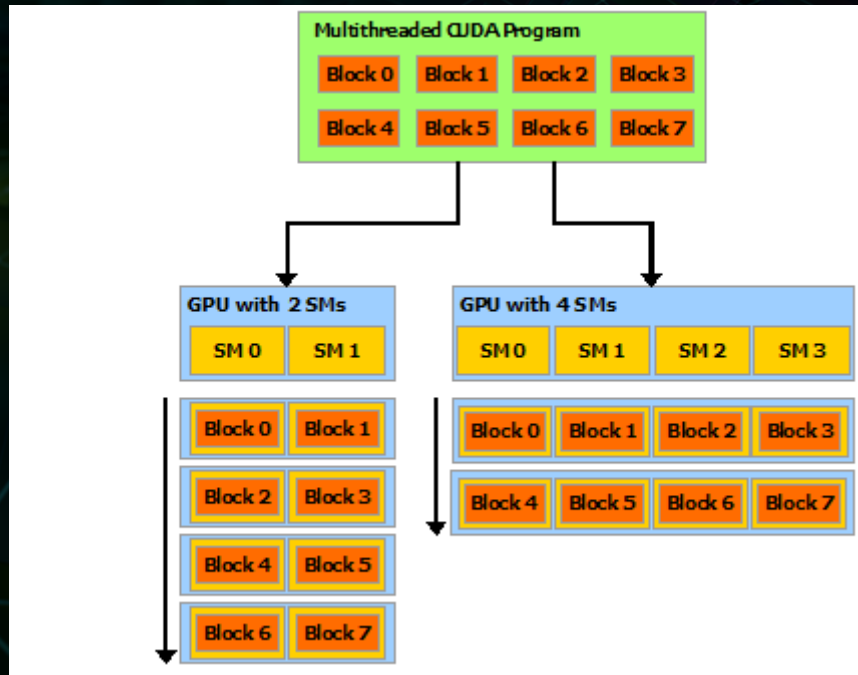

Thread hierarchy

- Each thread gets a unique threadID
 - accessed by build-in variables:
`int i = threadIdx.x`
- ThreadID is a 3 component Vector
 - Thread identification using up to 3-dimensional indices
 - Threads are arranged in blocks of threads
 - Thread Blocks



Thread hierarchy

- Number of Threads per Block is set by user



→ Calculate blockSize and number of Blocks!

Thread hierarchy

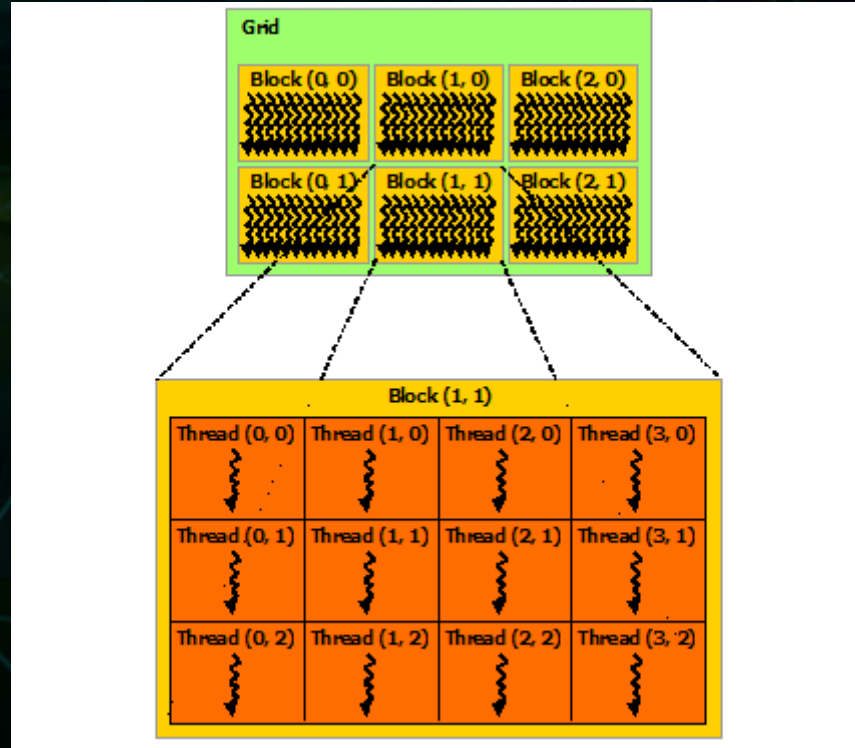
RTX 2070super:

- 40 Streaming Multiprocessors (Sms)
- 2560 Cuda Cores / Shading Units

- $\text{BlockSize} = \text{CudaCores} / \text{Sms}$
 $\text{blockSize} = 2560 / 40 = 64$
- $\text{NumberOfBlocks} = (\text{N} + \text{blockSize} - 1) / \text{blockSize}$
N = total number of Threads needed to solve the Problem
 $\text{numBlocks} = (10000 + 64 - 1) / 64 \approx 158$

Thread hierarchy

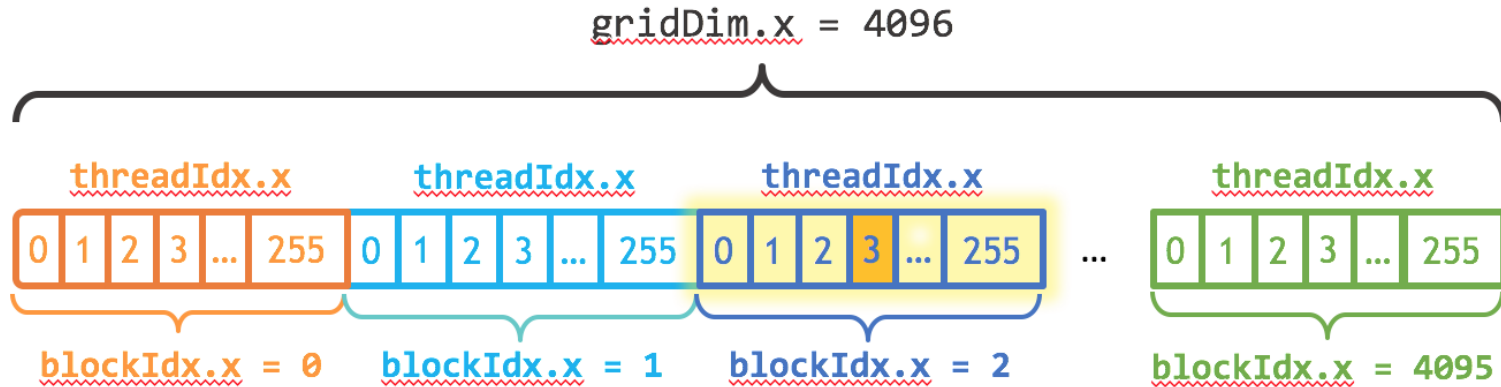
- Blocks get organized in a Grid



- Configuration specified by the Execution Configuration Syntax

Thread hierarchy

- BlockDim = number of threads in current block
- BlockIdx = ID of current block
- GridDim = number of threads in current Grid



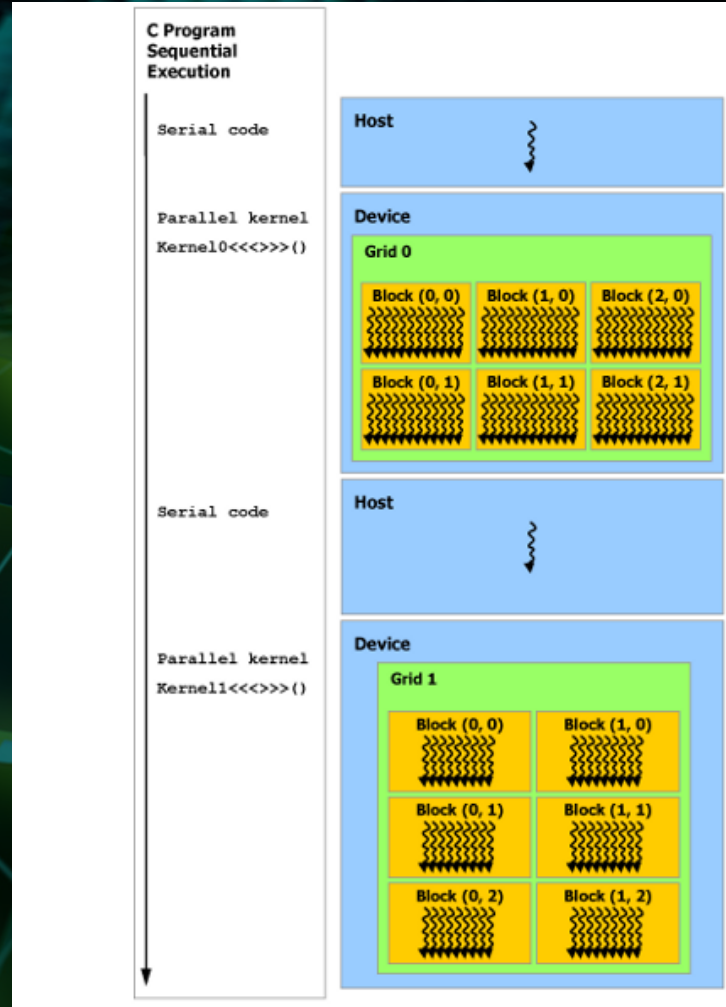
$$\text{index} = \text{blockIdx.x} * \text{blockDim.x} + \text{threadIdx.x}$$

$$\text{index} = (2) * (256) + (3) = 515$$

Heterogeneous Programming

- Cuda Programming Model assumes:
 - Cuda Threads get executed on a physically separate device
 - Both devices maintain their own separate memory Spaces
 - Host memory & device memory
- Memory needs to be allocated on both Devices
- Data needs to be transferred between both devices

Heterogeneous Programming





NVIDIA.COM

Usage

Installation

- Download & Install the CUDA TOOLKIT
<https://developer.nvidia.com/cuda-downloads>
- Prerequisites:
 - Cuda Capable GPU:
<https://developer.nvidia.com/cuda-gpus>
 - You can code and compile Cuda Kernels without a Nvidia GPU
 - Cuda is also available in cloud services:
 - Amazon AWS
 - Microsoft Azure
 - ...

IDEs & Tools

- IDEs:
 - Nsight Studio (eclipse based)
 - Clion (doesn't include debugging functionality)
 - Visual Studio
- Tools:
 - Nvidia Visual Profiler: Performance analysis for Cuda Kernels

IDEs & Tools

The screenshot shows the Microsoft Visual Studio IDE with the OptiX Samples (Debugging) window open. The interface is divided into several panes:

- Source Code:** Displays the C++ source code for the `diffuse()` function in `optiPathTracer.cpp`. The code includes comments and function calls like `normalize`, `faceforward`, and `make_float3`.
- Disassembly:** Shows the PTX and SASS disassembly code for the current function. It includes instructions like `ld.global.nc.v4.u32`, `mul.rn.f32`, and `add.f32`.
- GPU Registers:** Displays the SASS registers and their values. The registers are labeled R0 through R31, and their values are shown in hexadecimal.
- Locals:** Shows the local variables of the current function. The variables are `z1`, `z2`, `p`, `current_prd`, and `result`. Their values and types are listed.
- Call Stack:** Shows the call stack for the current function. The top entry is `[CUDA]000001b3f5928001:closesthit__Z7diffuse_ptb0d6d10b554b9654a_ss_0 Line 207 [0x000001b3f...`.
- Breakpoints:** Shows the breakpoints for the current function. A breakpoint is set on line 207 of `optiPathTracer.cpp`.

The status bar at the bottom shows the current line (Ln 207), column (Col 1), and other information.



NVIDIA.COM

Demo



NVIDIA.COM

More Info

More Info

- Nvidia Cuda Zone:
<https://developer.nvidia.com/cuda-zone>
- Nvidia developer Programm
- Programming Guide:
<https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html>