



GPU Teaching Kit
Accelerated Computing



Lecture 2.1 - Introduction to CUDA C

CUDA C vs. Thrust vs. CUDA Libraries

Objective

- To learn the main venues and developer resources for GPU computing
 - Where CUDA C fits in the big picture

3 Ways to Accelerate Applications

Applications

Libraries

Compiler
Directives

Programming
Languages

Easy to use
Most Performance

Easy to use
Portable code

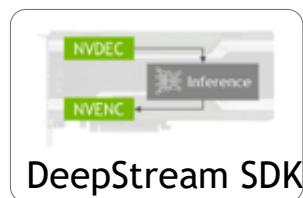
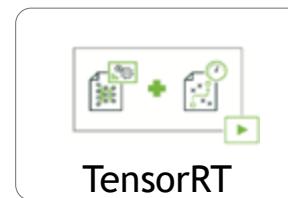
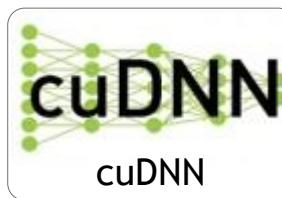
Most Performance
Most Flexibility

Libraries: Easy, High-Quality Acceleration

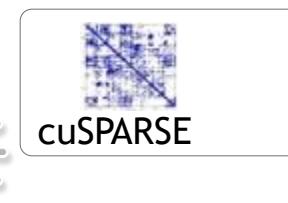
- **Ease of use:** Using libraries enables GPU acceleration without in-depth knowledge of GPU programming
- **“Drop-in”:** Many GPU-accelerated libraries follow standard APIs, thus enabling acceleration with minimal code changes
- **Quality:** Libraries offer high-quality implementations of functions encountered in a broad range of applications

NVIDIA GPU Accelerated Libraries

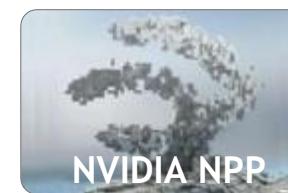
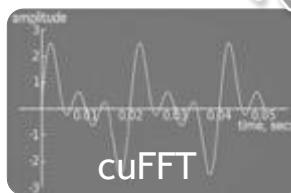
DEEP LEARNING



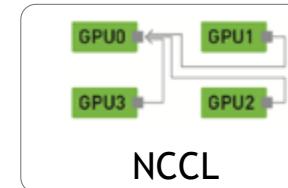
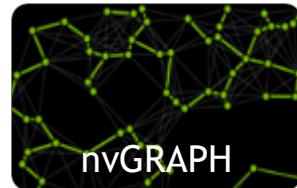
LINEAR ALGEBRA



SIGNAL, IMAGE, VIDEO



PARALLEL ALGORITHMS



Vector Addition in Thrust

```
#include <thrust/device_vector.h>
#include <thrust/copy.h>

int main(void) {
    size_t inputLength = 500;
    thrust::host_vector<float> hostInput1(inputLength);
    thrust::host_vector<float> hostInput2(inputLength);
    thrust::device_vector<float> deviceInput1(inputLength);
    thrust::device_vector<float> deviceInput2(inputLength);
    thrust::device_vector<float> deviceOutput(inputLength);

    thrust::copy(hostInput1.begin(), hostInput1.end(), deviceInput1.begin());
    thrust::copy(hostInput2.begin(), hostInput2.end(), deviceInput2.begin());

    thrust::transform(deviceInput1.begin(), deviceInput1.end(),
                    deviceInput2.begin(), deviceOutput.begin(),
                    thrust::plus<float>());
}
```

Compiler Directives: Easy, Portable Acceleration

- **Ease of use:** Compiler takes care of details of parallelism management and data movement
- **Portable:** The code is generic, not specific to any type of hardware and can be deployed into multiple languages
- **Uncertain:** Performance of code can vary across compiler versions

OpenACC

- Compiler directives for C, C++, and FORTRAN

```
#pragma acc parallel loop
copyin(input1[0:inputLength],input2[0:inputLength]),
    copyout(output[0:inputLength])
for(i = 0; i < inputLength; ++i) {
    output[i] = input1[i] + input2[i];
}
```

Programming Languages: Most Performance and Flexible Acceleration

- **Performance:** Programmer has best control of parallelism and data movement
- **Flexible:** The computation does not need to fit into a limited set of library patterns or directive types
- **Verbose:** The programmer often needs to express more details

GPU Programming Languages

Numerical analytics ➤

MATLAB, Mathematica, LabVIEW

Python ➤

PyCUDA, Numba

Fortran ➤

CUDA Fortran, OpenACC

C ➤

CUDA C, OpenACC

C++ ➤

CUDA C++, Thrust

C# ➤

Hybridizer

CUDA - C

Applications

Libraries

Compiler
Directives

Programming
Languages

Easy to use
Most Performance

Easy to use
Portable code

Most Performance
Most Flexibility



GPU Teaching Kit

Accelerated Computing



The GPU Teaching Kit is licensed by NVIDIA and the University of Illinois under the [Creative Commons Attribution-NonCommercial 4.0 International License](#).



NVIDIA®

GPU Teaching Kit
Accelerated Computing



ILLINOIS
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

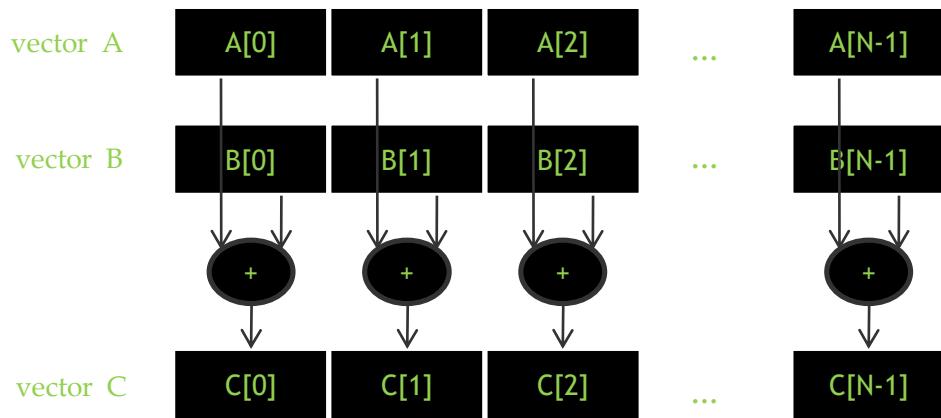
Lecture 2.2 - Introduction to CUDA C

Memory Allocation and Data Movement API Functions

Objective

- To learn the basic API functions in CUDA host code
 - Device Memory Allocation
 - Host-Device Data Transfer

Data Parallelism - Vector Addition Example



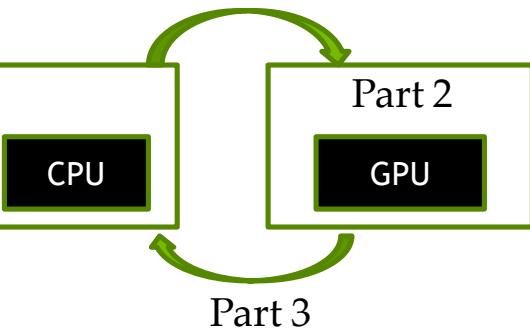
Vector Addition – Traditional C Code

```
// Compute vector sum C = A + B
void vecAdd(float *h_A, float *h_B, float *h_C, int n)
{
    int i;
    for (i = 0; i<n; i++) h_C[i] = h_A[i] + h_B[i];
}

int main()
{
    // Memory allocation for h_A, h_B, and h_C
    // I/O to read h_A and h_B, N elements
    ...
    vecAdd(h_A, h_B, h_C, N);
}
```

Heterogeneous Computing vecAdd CUDA Host Code

Part 1

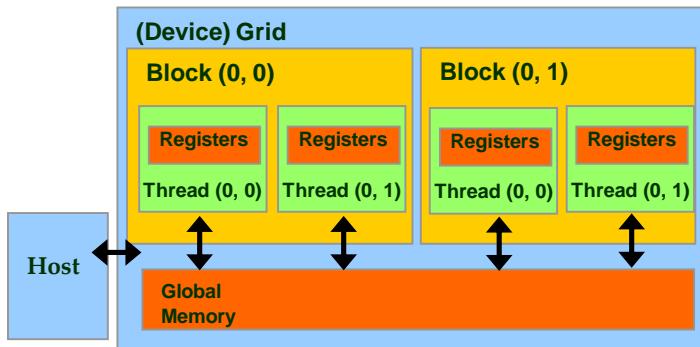


```
#include <cuda.h>
void vecAdd(float *h_A, float *h_B, float *h_C, int n)
{
    int size = n * sizeof(float);
    float *d_A, *d_B, *d_C;
    // Part 1
    // Allocate device memory for A, B, and C
    // copy A and B to device memory

    // Part 2
    // Kernel launch code – the device performs the actual vector addition

    // Part 3
    // copy C from the device memory
    // Free device vectors
}
```

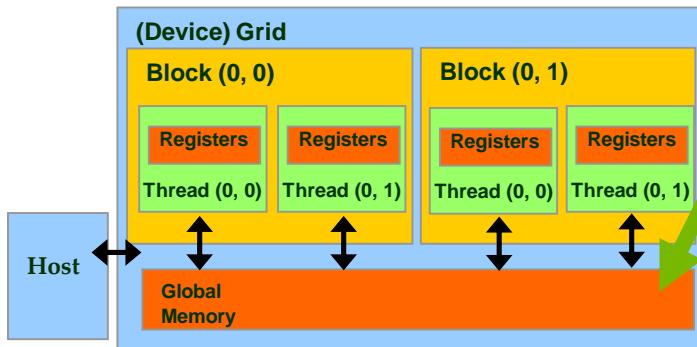
Partial Overview of CUDA Memories



- Device code can:
 - R/W per-thread **registers**
 - R/W all-shared **global memory**
- Host code can
 - Transfer data to/from per grid **global memory**

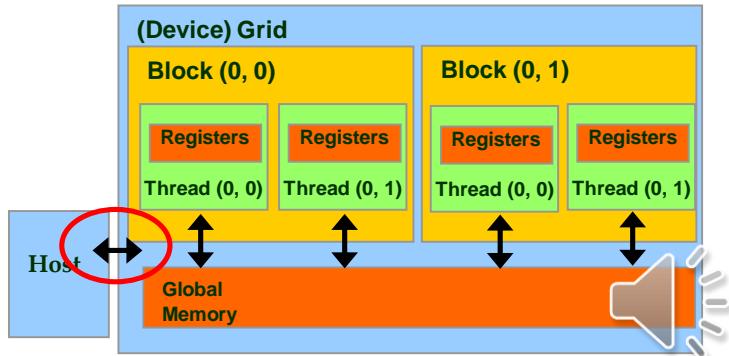
We will cover more memory types and more sophisticated memory models later.

CUDA Device Memory Management API functions



- `cudaMalloc()`
 - Allocates an object in the device global memory
 - Two parameters
 - **Address of a pointer** to the allocated object
 - **Size of allocated object** in terms of bytes
- `cudaFree()`
 - Frees object from device global memory
 - One parameter
 - **Pointer** to freed object

Host-Device Data Transfer API functions



- `cudaMemcpy()`
 - memory data transfer
 - Requires four parameters
 - Pointer to destination
 - Pointer to source
 - Number of bytes copied
 - Type/Direction of transfer
- Transfer to device is synchronous with respect to the host

Vector Addition, Explicit Memory Management

... Allocate *h_A*, *h_B*, *h_C* ...

```
void vecAdd(float *h_A, float *h_B, float *h_C, int n)
{
    int size = n * sizeof(float); float *d_A, *d_B, *d_C;
```

```
cudaMalloc((void **) &d_A, size);
cudaMalloc((void **) &d_B, size);
cudaMalloc((void **) &d_C, size);
```

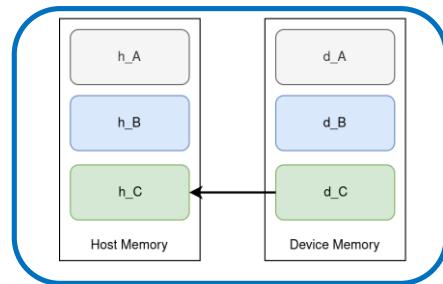
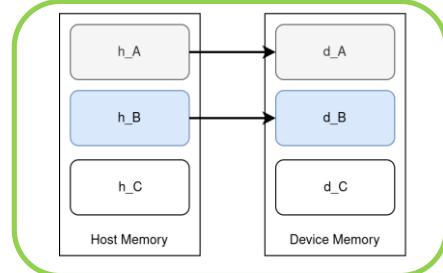
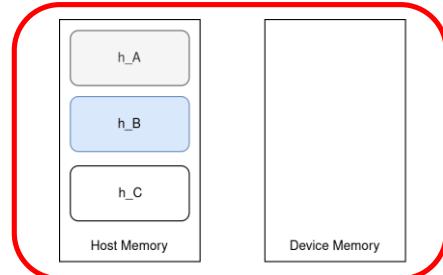


```
cudaMemcpy(d_A, h_A, size, cudaMemcpyHostToDevice);
cudaMemcpy(d_B, h_B, size, cudaMemcpyHostToDevice);
```

// Kernel invocation code – to be shown later

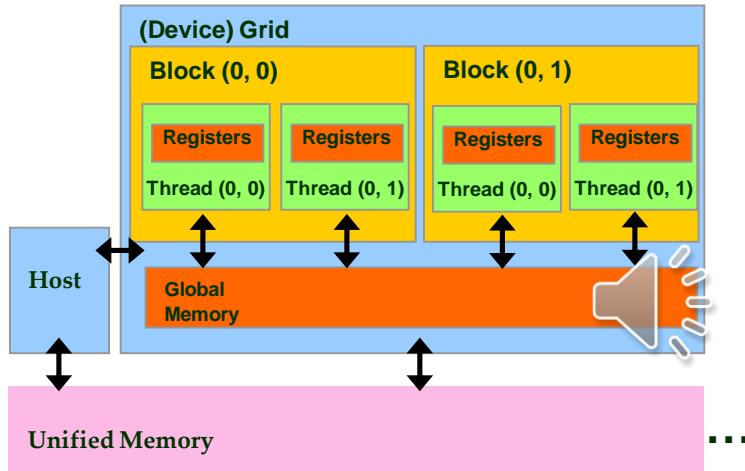
```
cudaMemcpy(h_C, d_C, size, cudaMemcpyDeviceToHost);
cudaFree(d_A); cudaFree(d_B); cudaFree(d_C);
}
```

... Free *h_A*, *h_B*, *h_C* ...



Unified Memory

- `cudaMallocManaged(
void** ptr, size_t size)`



- Single memory space for all CPUs/GPUs
 - Maintain single copy of data
 - CUDA-managed data
 - On-demand page migration
 - Compatible with `cudaMalloc()`, `cudaFree()`
 - Can be optimized
 - `cudaMemAdvise()`, `cudaMemPrefetchAsync()`,
`cudaMemcpyAsync()`

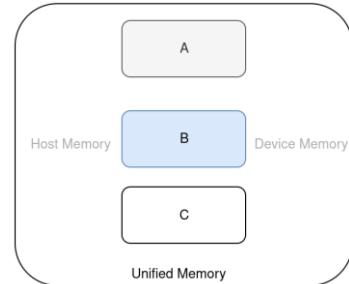
Vector Addition, Unified Memory

```
float *A, *B, *C  
cudaMallocManaged(&A, n * sizeof(float));  
cudaMallocManaged(&B, n * sizeof(float));  
cudaMallocManaged(&C, n * sizeof(float));
```

```
// Initialize A, B
```

```
void vecAdd(float *A, float *B, float *C, int n)  
{  
    // Kernel invocation code – to be shown later  
}
```

```
cudaFree(A);  
cudaFree(B);  
cudaFree(C);
```



In Practice, Check for API Errors in Host Code

```
cudaError_t err = cudaMalloc((void **) &d_A, size);

if (err != cudaSuccess) {
    printf("%s in %s at line %d\n", cudaGetErrorString(err), __FILE__,
           __LINE__);
    exit(EXIT_FAILURE);
}
```



GPU Teaching Kit

Accelerated Computing



The GPU Teaching Kit is licensed by NVIDIA and the University of Illinois under the [Creative Commons Attribution-NonCommercial 4.0 International License](#).



NVIDIA®

GPU Teaching Kit
Accelerated Computing



ILLINOIS
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

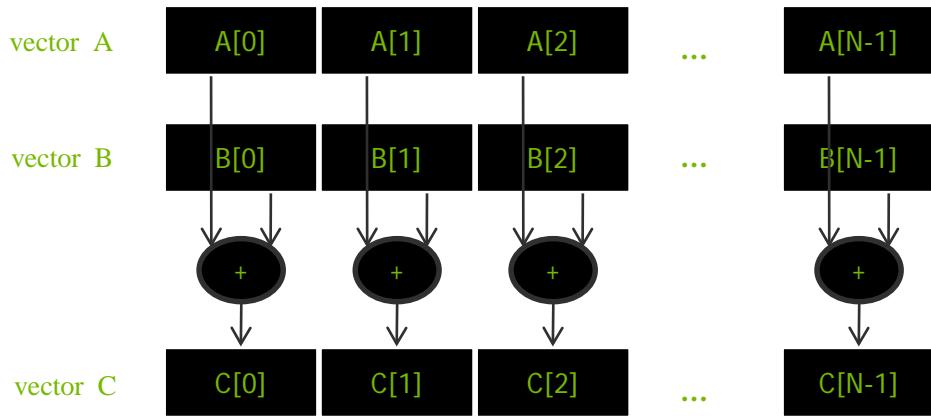
Lecture 2.3 – Introduction to CUDA C

Threads and Kernel Functions

Objective

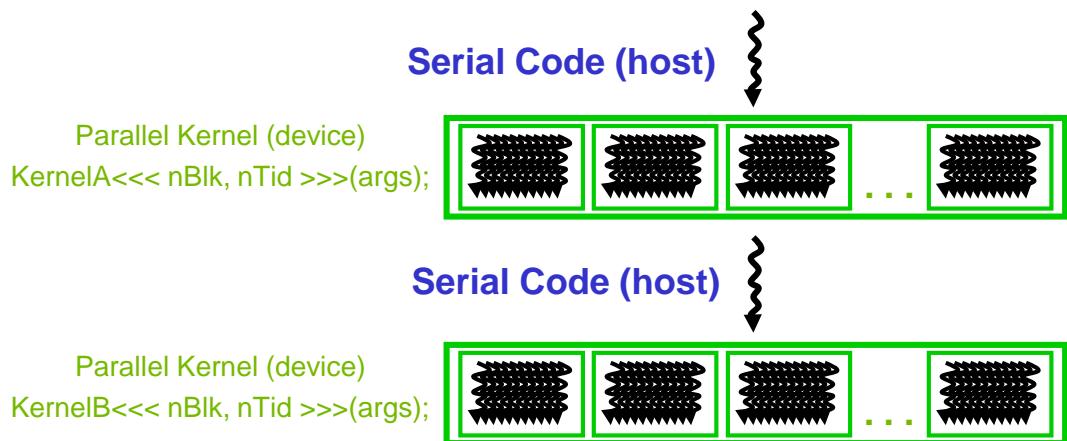
- To learn about CUDA threads, the main mechanism for exploiting of data parallelism
 - Hierarchical thread organization
 - Launching parallel execution
 - Thread index to data index mapping

Data Parallelism - Vector Addition Example

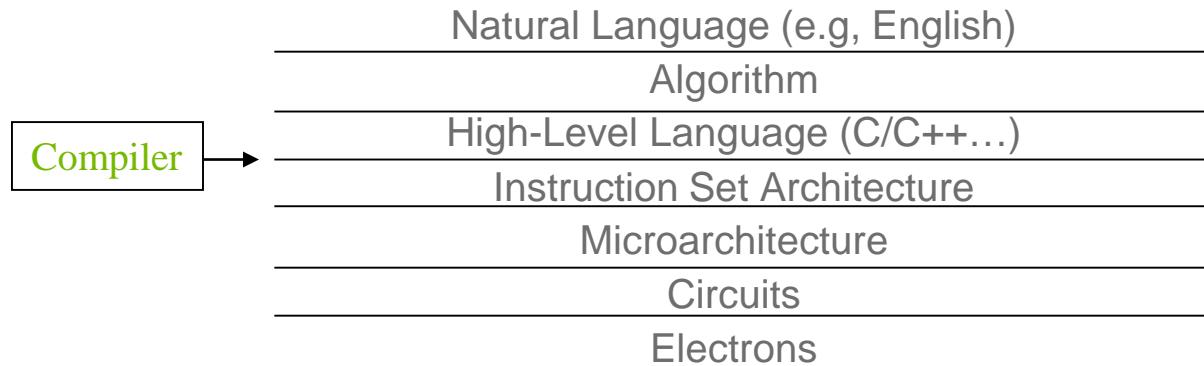


CUDA Execution Model

- Heterogeneous host (CPU) + device (GPU) application C program
 - Serial parts in **host** C code
 - Parallel parts in **device** SPMD kernel code



From Natural Language to Electrons



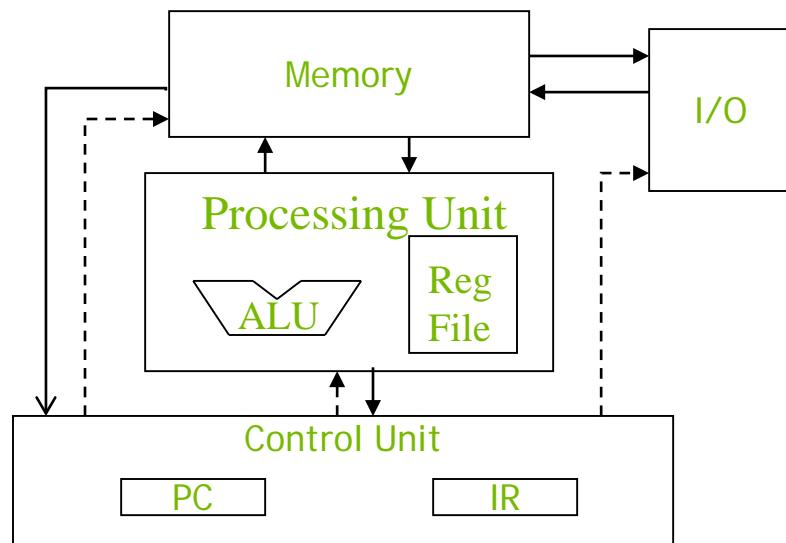
©Yale Patt and Sanjay Patel, *From bits and bytes to gates and beyond*

A program at the ISA level

- A program is a set of instructions stored in memory that can be read, interpreted, and executed by the hardware.
 - Both CPUs and GPUs are designed based on (different) instruction sets
- Program instructions operate on data stored in memory and/or registers.

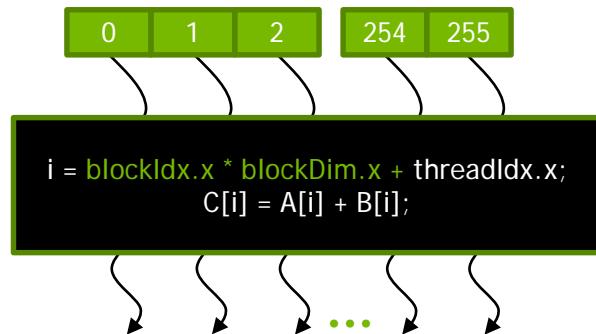
A Thread as a Von-Neumann Processor

A thread is a “virtualized” or
“abstracted”
Von-Neumann Processor

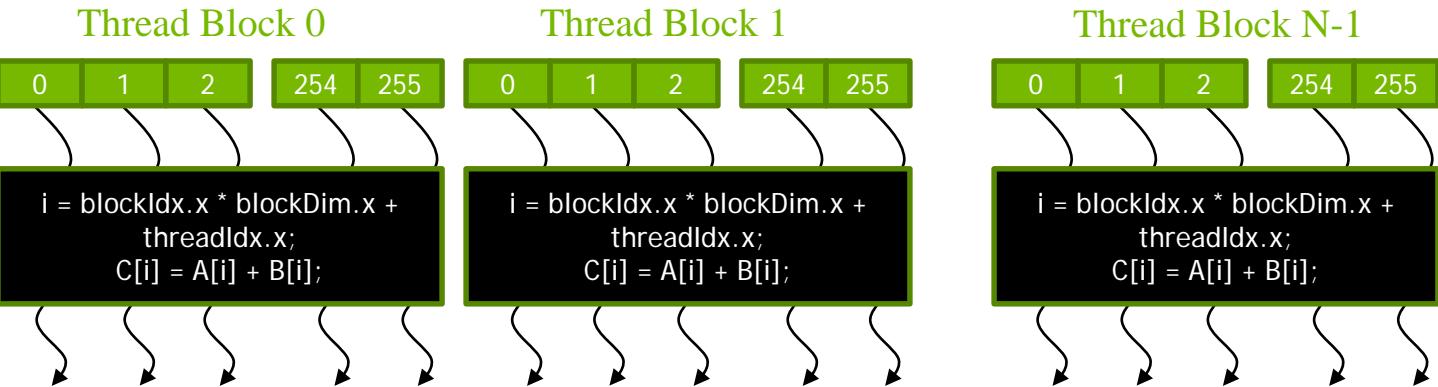


Arrays of Parallel Threads

- A CUDA kernel is executed by a **grid** (array) of threads
 - All threads in a grid run the same kernel code (Single Program Multiple Data)
 - Each thread has indexes that it uses to compute memory addresses and make control decisions



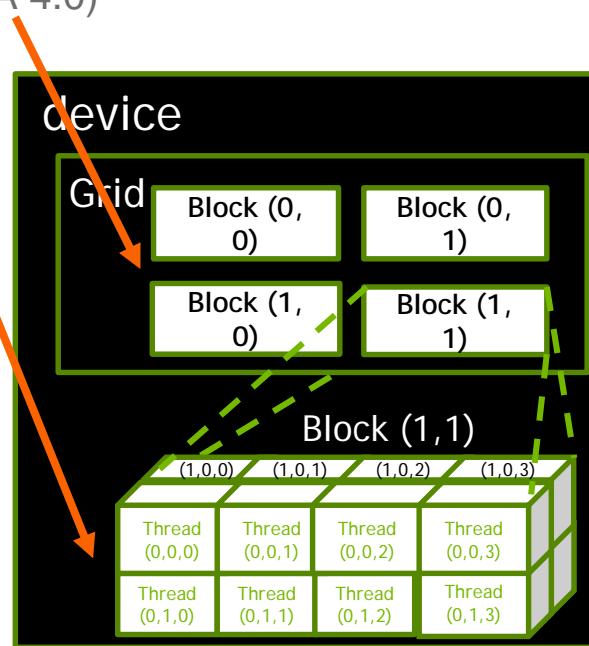
Thread Blocks: Scalable Cooperation



- Divide thread array into multiple blocks
 - Threads within a block cooperate via **shared memory, atomic operations** and **barrier synchronization**
 - Threads in different blocks do not interact

blockIdx and threadIdx

- Each thread uses indices to decide what data to work on
 - blockIdx: 1D, 2D, or 3D (CUDA 4.0)
 - threadIdx: 1D, 2D, or 3D
- Simplifies memory addressing when processing multidimensional data
 - Image processing
 - Solving PDEs on volumes
 - ...





GPU Teaching Kit

Accelerated Computing



The GPU Teaching Kit is licensed by NVIDIA and the University of Illinois under the [Creative Commons Attribution-NonCommercial 4.0 International License](#).



nVIDIA®

GPU Teaching Kit
Accelerated Computing



ILLINOIS

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

Lecture 2.4 – Introduction to CUDA C

Introduction to the CUDA Toolkit

Objective

- To become familiar with some valuable tools and resources from the CUDA Toolkit
 - Compiler flags
 - Debuggers
 - Profilers

GPU Programming Languages

Numerical analytics ➤ MATLAB, Mathematica, LabVIEW

Python ➤ PyCUDA, Numba

Fortran ➤ CUDA Fortran, OpenACC

C ➤ CUDA C, OpenACC

C++ ➤ CUDA C++, Thrust

C# ➤ Hybridizer



CUDA - C

Applications

Libraries

Easy to use
Most Performance

Compiler
Directives

Easy to use
Portable code

Programming
Languages

Most Performance
Most Flexibility

NVCC Compiler

- NVIDIA provides a CUDA-C compiler
 - nvcc
- NVCC compiles device code then forwards code on to the host compiler (e.g. g++)
- Can be used to compile & link host only applications

Example 1: Hello World

```
#include <cstdio>

int main() {
    printf("Hello World!\n");
    return 0;
}
```

Instructions:

1. Build and run the hello world code
2. Modify Makefile to use nvcc
instead of g++
3. Rebuild and run

CUDA Example 1: Hello World

```
#include <cstdio>

__global__ void mykernel(void) {
}

int main(void) {
    mykernel<<<1,1>>>();
    printf("Hello World!\n");
    return 0;
}
```

Instructions:

1. Add kernel and kernel launch to **main.cc**
2. Try to build

CUDA Example 1: Build Considerations

- Build failed
 - Nvcc only parses .cu files for CUDA
- Fixes:
 - Rename main.cc to main.cu
 - nvcc –x cu
 - Treat all input files as .cu files

Instructions:

1. Rename main.cc to main.cu
2. Rebuild and Run

Hello World! with Device Code

```
#include <cstdio>

__global__ void mykernel(void) {
}

int main(void) {
    mykernel<<<1,1>>>();
    printf("Hello World!\n");
    return 0;
}
```

Output:

```
$ nvcc main.cu
$ ./a.out
Hello World!
```

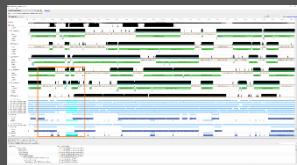
- mykernel (does nothing, somewhat anticlimactic!)

Developer Tools - Debuggers

Nsight



Nsight
Systems



CUDA-GDB



CUDA
MEMCHECK



NVIDIA Provided

arm
FORGE

TotalView®

3rd Party

<https://developer.nvidia.com/debugging-solutions>

Compiler Flags

- Remember there are two compilers being used
 - NVCC: Device code
 - Host Compiler: C/C++ code
- NVCC supports some host compiler flags
 - If flag is unsupported, use `-Xcompiler` to forward to host
 - e.g. `-Xcompiler -fopenmp`
- Debugging Flags
 - `-g`: Include host debugging symbols
 - `-G`: Include device debugging symbols
 - `-lineinfo`: Include line information with symbols

CUDA-MEMCHECK

- Memory debugging tool
 - No recompilation necessary
%> cuda-memcheck ./exe
- Can detect the following errors
 - Memory leaks
 - Memory errors (OOB, misaligned access, illegal instruction, etc)
 - Race conditions
 - Illegal Barriers
 - Uninitialized Memory
- For line numbers use the following compiler flags:
 - -Xcompiler -rdynamic -lineinfo

<http://docs.nvidia.com/cuda/cuda-memcheck>

Example 2: CUDA-MEMCHECK

Instructions:

1. Build & Run Example 2
Output should be the numbers 0-9
Do you get the correct results?
2. Run with cuda-memcheck
%> cuda-memcheck ./a.out
3. Add nvcc flags “-Xcompiler -rdynamic -lineinfo”
4. Rebuild & Run with cuda-memcheck
5. Fix the illegal write

<http://docs.nvidia.com/cuda/cuda-memcheck>

CUDA-GDB

- cuda-gdb is an extension of GDB
 - Provides seamless debugging of CUDA and CPU code
- Works on Linux and Macintosh
 - For a Windows debugger use NVIDIA Nsight Eclipse Edition or Visual Studio Edition

<http://docs.nvidia.com/cuda/cuda-gdb>

Example 3: cuda-gdb

Instructions:

1. Run exercise 3 in cuda-gdb

```
%> cuda-gdb --args ./a.out
```

2. Run a few cuda-gdb commands:

```
(cuda-gdb) b main           //set break point at main
(cuda-gdb) r                 //run application
(cuda-gdb) l                 //print line context
(cuda-gdb) b foo              //break at kernel foo
(cuda-gdb) c                 //continue
(cuda-gdb) cuda thread       //print current thread
(cuda-gdb) cuda thread 10    //switch to thread 10
(cuda-gdb) cuda block        //print current block
(cuda-gdb) cuda block 1      //switch to block 1
(cuda-gdb) d                 //delete all break points
(cuda-gdb) set cuda memcheck on //turn on cuda memcheck
(cuda-gdb) r                 //run from the beginning
```

3. Fix Bug

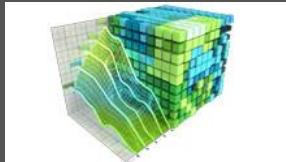
<http://docs.nvidia.com/cuda/cuda-gdb>

Developer Tools - Profilers

NSIGHT



NVVP

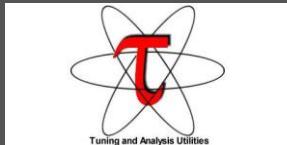


NVPROF

```
--29561: Profiling result:  
Time(%) Time Calls Avg Min Max Name  
49.88% 866.69ms 594758 1.7170us 1.5040us 2.0160us void th  
int, thrust::detail::device_generate_functor<thrust::detail::fill_<br>_t, thrust::detail::device_generate_functor<thrust::detail::fill_<br>_t, thrust::detail::device_generate_functor<thrust::detail::fill_<br>_t>::operator()>::operator()  
17.07% 296.68ms 200 1.4830us 1.2040us 1.7253ms kerComp  
2.98% 51.81ms 200 259.89us 246.97us 264.81us kerMake  
1.10% 20.09ms 200 1.0000us 999.99us 1.0000us [CUBLA m  
0.53% 16.19ms 200 89.991us 71.840us 90.751us kerColV  
0.73% 12.636ns 400 31.589us 14.720us 50.432us [CUBLA m  
0.69% 12.075ns 200 60.376us 59.680us 62.384us kerRowA  
0.53% 11.975ns 200 54.960us 54.960us 54.960us kerRowB  
0.32% 5.6524ns 200 22.559us 22.559us 33.152us [CUBLA m  
0.12% 2.1342ns 1 2.1342ns 2.1342ns 2.1342ns void th
```

NVIDIA Provided

TAU



VampirTrace



3rd Party

<https://developer.nvidia.com/performance-analysis-tools>

NVPROF

Command Line Profiler

- Compute time in each kernel
- Compute memory transfer time
- Collect metrics and events
- Support complex process hierarchy's
- Collect profiles for NVIDIA Visual Profiler
- No need to recompile

Example 4: nvprof

Instructions:

1. Collect profile information for the matrix add example

```
%> nvprof ./a.out
```

2. How much faster is add_v2 than add_v1?

3. View available metrics

```
%> nvprof --query-metrics
```

4. View global load/store efficiency

```
%> nvprof --metrics
```

```
gld_efficiency,gst_efficiency ./a.out
```

5. Store a timeline to load in NVVP

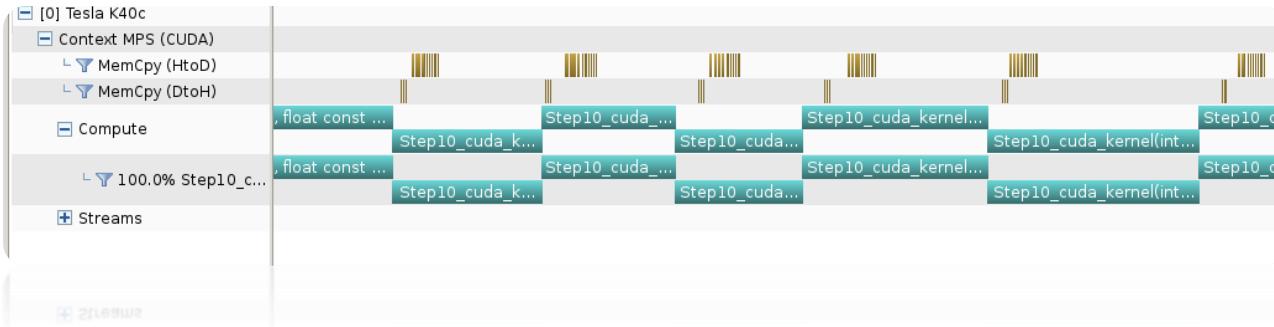
```
%> nvprof -o profile.timeline ./a.out
```

6. Store analysis metrics to load in NVVP

```
%> nvprof -o profile.metrics --analysis-metrics  
./a.out
```

NVIDIA's Visual Profiler (NVVP)

Timeline



Guided System

1. CUDA Application Analysis

2. Performance-Critical Kernels

3. Compute, Bandwidth, or Latency Bound

The first step in analyzing an individual kernel is to determine if the performance of the kernel is bounded by computation, memory bandwidth, or instruction/memory latency. The results at right indicate that the performance of kernel "Step10_cuda_kernel" is most likely limited by compute.

Perform Compute Analysis

The most likely bottleneck for performance for this kernel is compute so you should first perform compute analysis to determine how it is limiting performance.

Perform Latency Analysis

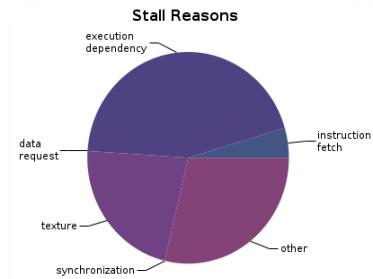
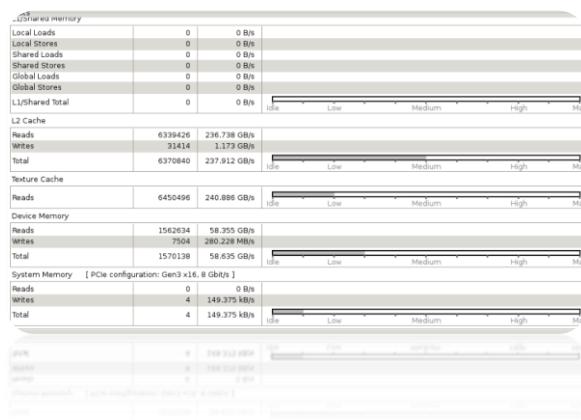
Perform Memory Bandwidth Analysis

Instruction and memory latency and memory bandwidth are likely not the primary performance bottlenecks for this kernel, but you may still want to perform those analyses.

Rerun Analysis

If you modify the kernel you need to rerun your application to update this analysis.

Analysis



Example 4: NVVP

Instructions:

1. Import nvprof profile into NVVP

Launch nvvp

Click File/ Import/ Nvprof/ Next/ Single process/ Next / Browse

Select profile.timeline

Add Metrics to timeline

Click on 2nd Browse

Select profile.metrics

Click Finish

2. Explore Timeline

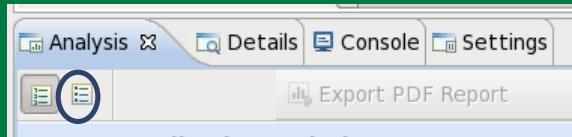
Control + mouse drag in timeline to zoom in

Control + mouse drag in measure bar (on top) to measure time

Example 4: NVVP

Instructions:

1. Click on a kernel
2. On Analysis tab click on the unguided analysis



2. Click Analyze All

Explore metrics and properties

What differences do you see between the two kernels?

Note:

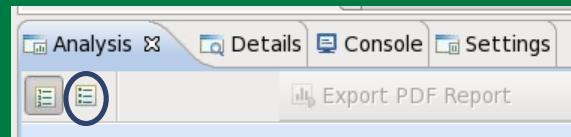
If kernel order is non-deterministic you can only load the timeline or the metrics but not both.

If you load just metrics the timeline looks odd but metrics are correct.

Example 4: NVVP

Let's now generate the same data within NVVP

1. Click File / New Session / Browse
Select Example 4/a.out
Click Next / Finish



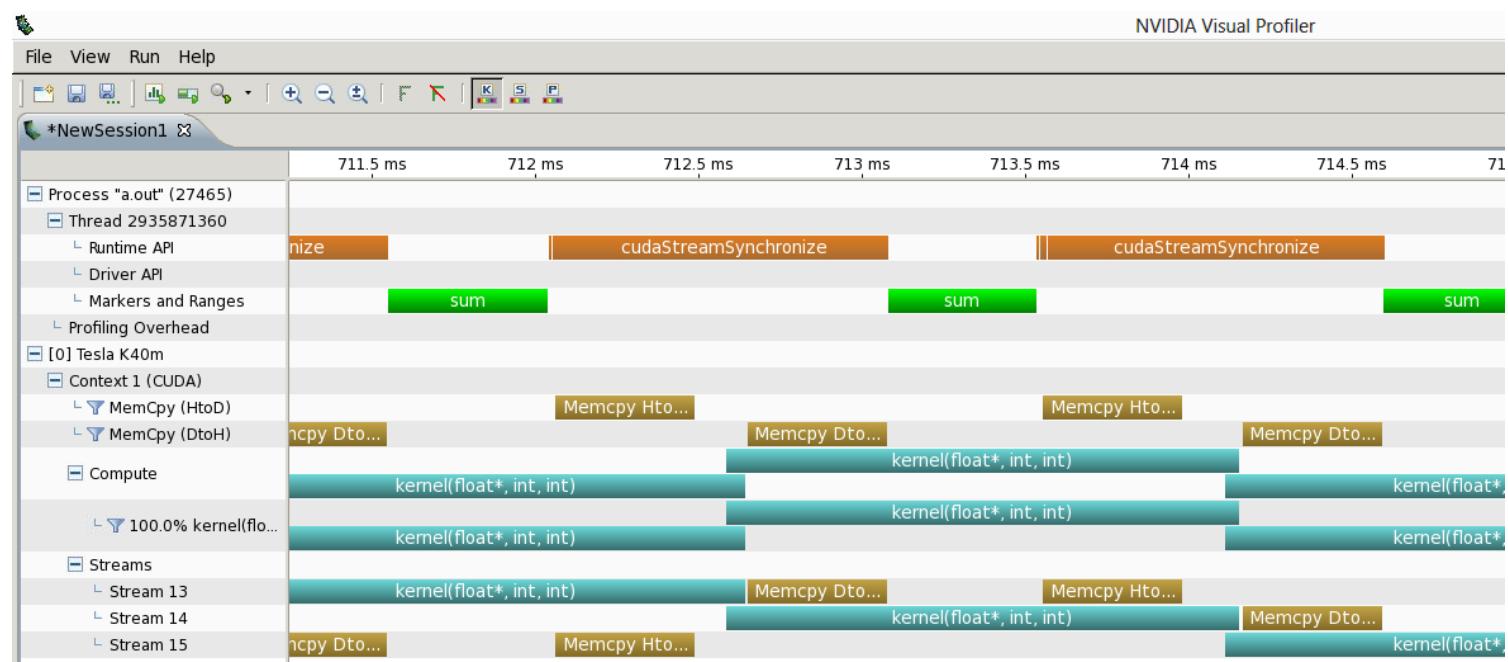
2. Click on a kernel
Select Unguided Analysis
Click Analyze All

NVTX

- Our current tools only profile API calls on the host
 - What if we want to understand better what the host is doing?
- The NVTX library allows us to annotate profiles with ranges
 - Add: #include <nvToolsExt.h>
 - Link with: -lNvToolsExt
- Mark the start of a range
 - nvtxRangePushA("description");
- Mark the end of a range
 - nvtxRangePop();
- Ranges are allowed to overlap

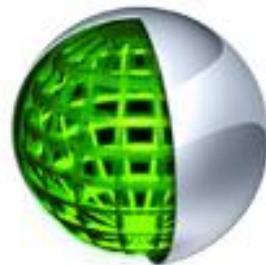
<http://devblogs.nvidia.com/parallelforall/cuda-pro-tip-generate-custom-application-profile-timelines-nvtx/>

NVTX Profile



NSIGHT

- CUDA enabled Integrated Development Environment
 - Source code editor: syntax highlighting, code refactoring, etc
 - Build Manager
 - Visual Debugger
 - Visual Profiler
- Linux/Macintosh
 - Editor = Eclipse
 - Debugger = cuda-gdb with a visual wrapper
 - Profiler = NVVP
- Windows
 - Integrates directly into Visual Studio
 - Profiler is NSIGHT VSE



Example 4: NSIGHT

Let's import an existing Makefile project into NSIGHT

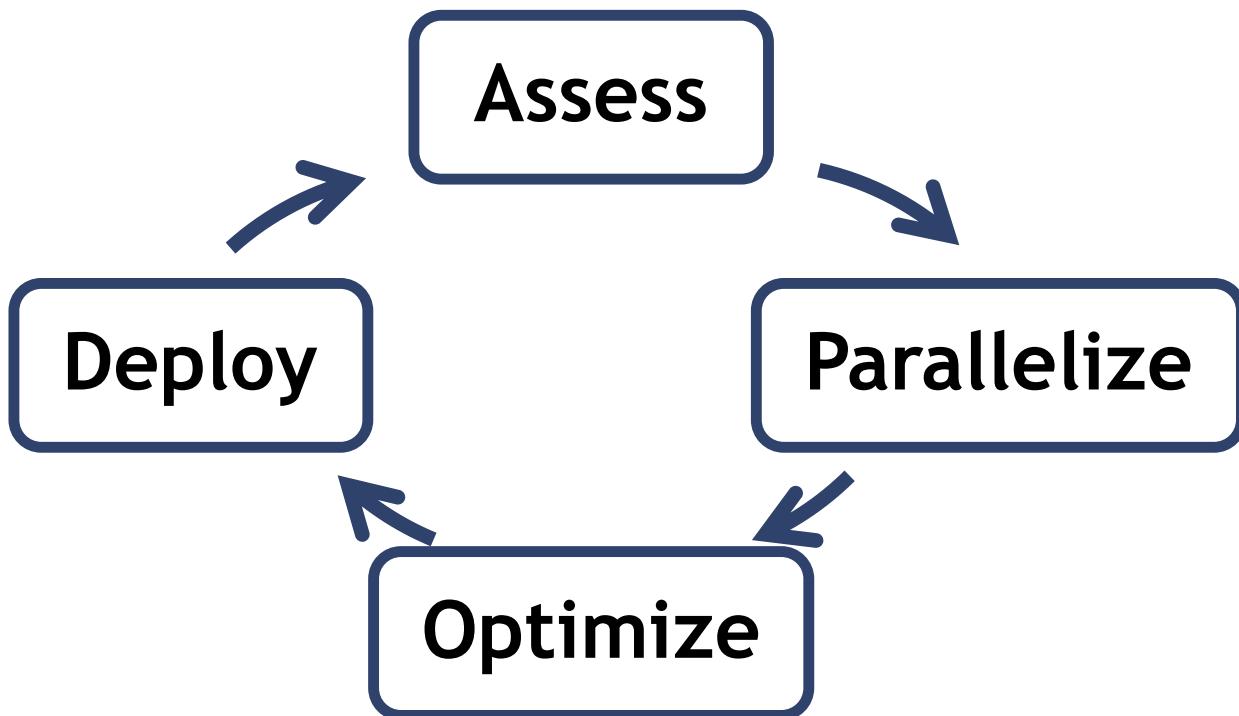
Instructions:

1. Run nsight
Select default workspace
2. Click File / New / Makefile Project With Existing CodeTest
3. Enter Project Name and select the Example15 directory
4. Click Finish
5. Right Click On Project / Properties / Run Settings / New / C++ Application
6. Browse for Example 4/a.out
7. In Project Explorer double click on main.cu and explore source
8. Click on the build icon
9. Click on the run icon
10. Click on the profile icon

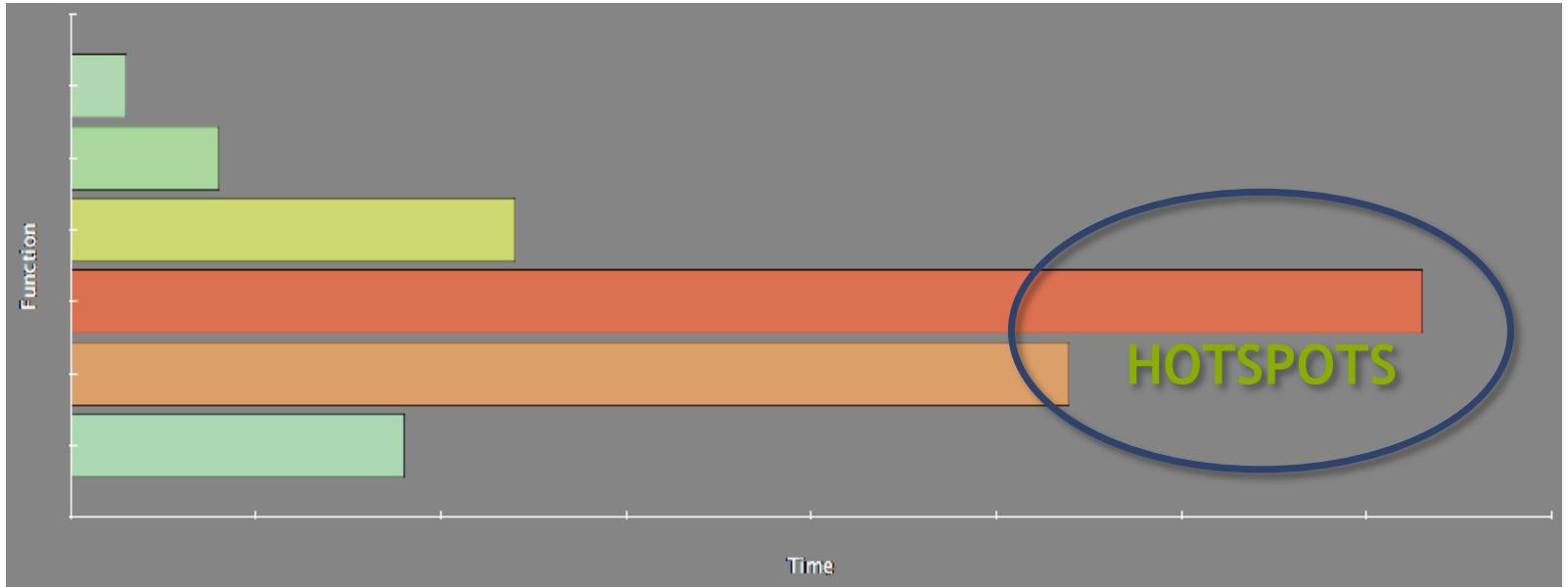
Profiler Summary

- Many profile tools are available
- NVIDIA Provided
 - NVPROF: Command Line
 - NVVP: Visual profiler
 - NSIGHT: IDE (Visual Studio and Eclipse)
- 3rd Party
 - TAU
 - VAMPIR

Optimization



Assess



- Profile the code, find the hotspot(s)
- Focus your attention where it will give the most benefit

Parallelize

Applications

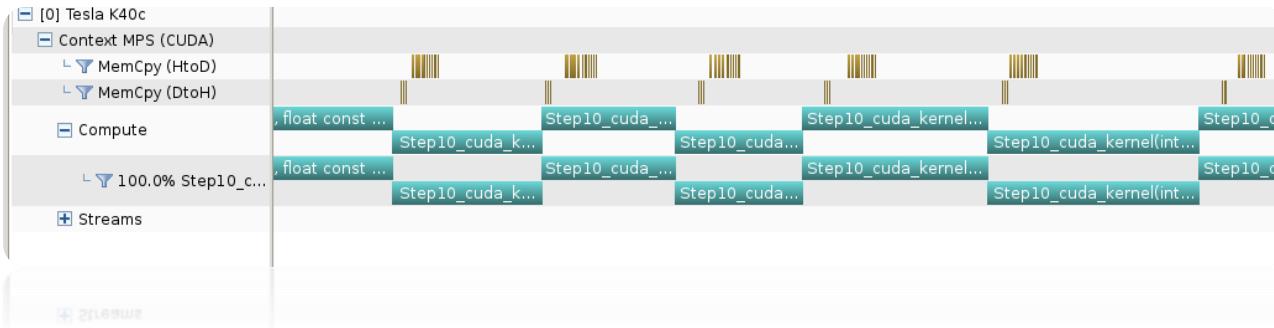
Libraries

Compiler
Directives

Programming
Languages

Optimize

Timeline



Guided System

1. CUDA Application Analysis

2. Performance-Critical Kernels

3. Compute, Bandwidth, or Latency Bound

The first step in analyzing an individual kernel is to determine if the performance of the kernel is bounded by computation, memory bandwidth, or instruction/memory latency. The results at right indicate that the performance of kernel "Step10_cuda_kernel" is most likely limited by compute.

Perform Compute Analysis

The most likely bottleneck for performance for this kernel is compute so you should first perform compute analysis to determine how it is limiting performance.

Perform Latency Analysis

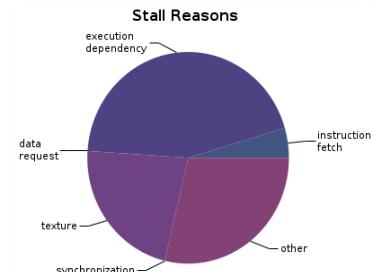
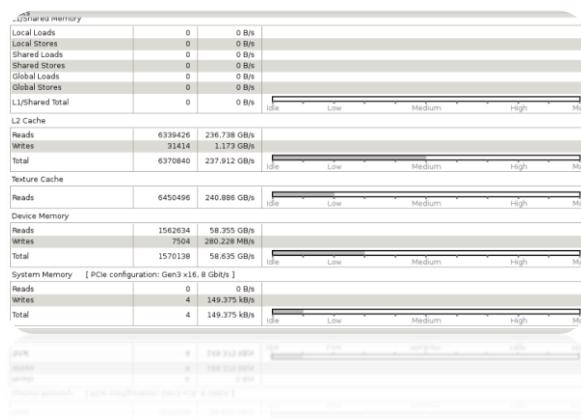
Perform Memory Bandwidth Analysis

Instruction and memory latency and memory bandwidth are likely not the primary performance bottlenecks for this kernel, but you may still want to perform those analyses.

Rerun Analysis

If you modify the kernel you need to rerun your application to update this analysis.

Analysis



Bottleneck Analysis

- Don't assume an optimization was wrong
- Verify if it was wrong with the profiler

129 GB/s → 84 GB/s

L1/Shared Memory		
Local Loads	0	0 B/s
Local Stores	0	0 B/s
Shared Loads	2097152	1,351.979 GB/s
Shared Stores	131072	84.499 GB/s
Global Loads	131072	42.249 GB/s
Global Stores	131072	42.249 GB/s
Atomic	0	0 B/s
L1/Shared Total	2490368	1,520.977 GB/s



gpuTranspose_kernel(int, int, float const *, float*)	
Start	547.303 ms (5)
End	547.716 ms (5)
Duration	413.872 µs
Grid Size	[64,64,1]
Block Size	[32,32,1]
Registers/Thread	10
Shared Memory/Block	4 KiB
Efficiency	
Global Load Efficiency	100%
Global Store Efficiency	100%
Shared Efficiency	5.9%
Warp Execution Efficiency	100%
Non-Predicated Warp Execution Efficiency	97.1%
Occupancy	
Achieved	86.7%
Theoretical	100%
Shared Memory Configuration	
Shared Memory Requested	48 KiB
Shared Memory Executed	48 KiB

⚠ Shared Memory Alignment and Access Pattern

Memory bandwidth is used most efficiently when each shared memory load and store has proper alignment and access pattern.

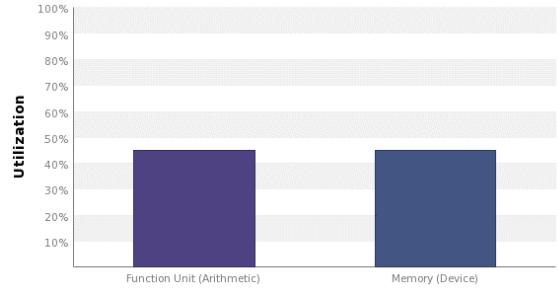
Optimization: Select each entry below to open the source code to a shared load or store within the kernel with an inefficient alignment or access pattern. For each access pattern of the memory access.

▼ Line / File main.cu - /home/jluitjens/code/CudaHandsOn/Example19

49 Shared Load Transactions/Access = 16, Ideal Transactions/Access = 1 [2097152 transactions for 131072 total executions]

Performance Analysis

gpuTranspose_kernel(int, int, float const *, float *, float const *, float const *, float const *, float const *)	
Start	770.0671
End	770.3241
Duration	256.7140
Grid Size	[64,64,1]
Block Size	[32,32,1]
Registers/Thread	10
Shared Memory/Block	4.125 KiB
▼ Efficiency	
Global Load Efficiency	100%
Global Store Efficiency	100%
Shared Efficiency	⚠ 50%
Warp Execution Efficiency	100%
Non-Predicated Warp Execution Efficiency	97.1%
▼ Occupancy	
Achieved	87.7%
Theoretical	100%
▼ Shared Memory Configuration	
Shared Memory Requested	48 KiB
Shared Memory Executed	48 KiB



84 GB/s → 137 GB/s

L1/Shared Memory			
Local Loads	0	0 B/s	
Local Stores	0	0 B/s	
Shared Loads	131072	138.433 GB/s	
Shared Stores	131720	139.118 GB/s	
Global Loads	131072	69.217 GB/s	
Global Stores	131072	69.217 GB/s	
Atomic	0	0 B/s	
L1/Shared Total	524936	415.984 GB/s	Idle Low Medium
L2 Cache			
L1 Reads	524288	69.217 GB/s	
L1 Writes	524288	69.217 GB/s	
Texture Reads	0	0 B/s	
Atomic	0	0 B/s	
Noncoherent Reads	0	0 B/s	
Total	1048576	138.433 GB/s	Idle Low Medium
Texture Cache			
Reads	0	0 B/s	Idle Low Medium
Device Memory			
Reads	524968	69.306 GB/s	
Writes	524289	69.217 GB/s	
Total	1049257	138.523 GB/s	Idle Low Medium



GPU Teaching Kit



The GPU Teaching Kit is licensed by NVIDIA and the University of Illinois under the [Creative Commons Attribution-NonCommercial 4.0 International License](#).



nVIDIA®

GPU Teaching Kit

Accelerated Computing



ILLINOIS

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

Lecture 2.5 – Nsight Compute and Nsight Systems

Introduction to the CUDA Toolkit

Objective

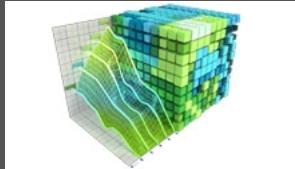
- To become familiar with Nsight Systems and Compute

Profiling Tools

nvprof

```
==29561== Profiling result:
Time(%):   Time    Calls  Avg     Min     Max  Name
49.88% 866.09ms 504758  1.7170us 1.5940us 2.0160us void th
  time, thrust::detail::device_generate_functor<thrust::detail::fill_
  15.85% 650.06ms 504758  1.7410us 1.7410us 1.7410us void th
  t, thrust::detail::device_generate_functor<thrust::detail::fill_fu
  17.07% 296.69ms 280  1.4830ns 1.2840ns 1.7230ns kerComp
  2.98% 51.819ms 280  23.0190us 24.970us 264.830us kernakel
  1.30% 57.375ms 504758  1.6030us 1.5960us 1.6100us kerTriv
  6.93% 16.198ms 280  80.991us 71.840us 96.751us kerColV
  6.73% 12.636ms 400  31.5890us 14.720us 56.432us [CUDA m
  6.69% 12.075ms 280  60.7760us 59.680us 62.364us kerRhoA
  6.65% 11.975ms 280  60.7760us 59.680us 62.364us kerRhoB
  6.32% 5.5520ns 280  27.781us 22.559us 33.152us [CUDA m
  6.12% 2.1342ns   1  2.1342ns 2.1342ns 2.1342ns void th
```

NVVP

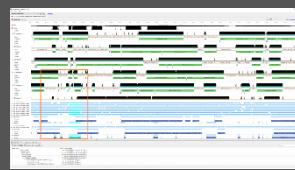


Phasing out

**Nsight
Compute**



**Nsight
Systems**



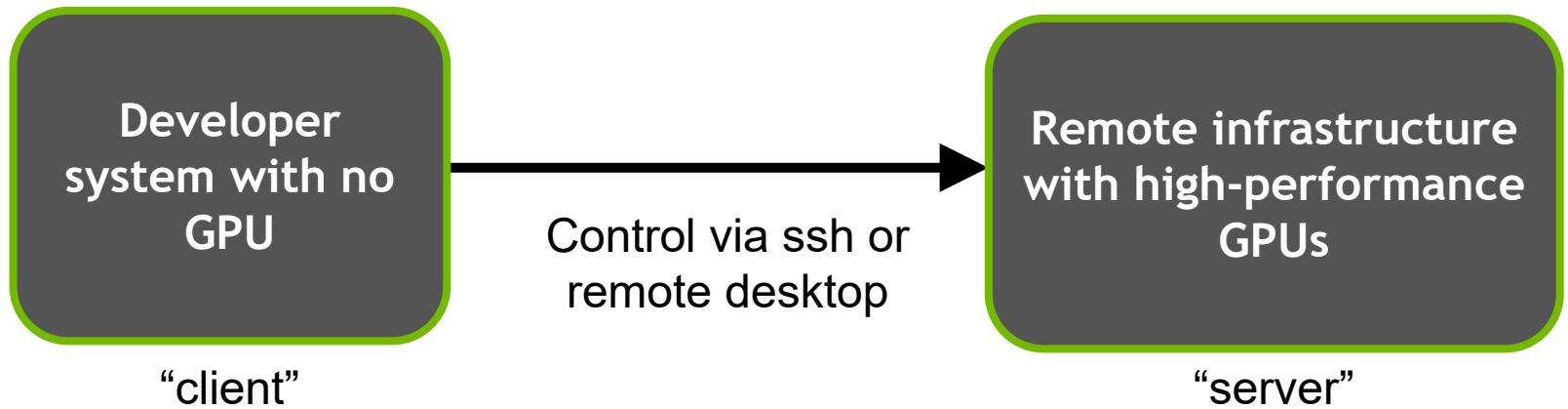
Current

See lecture 2-4 for an overview of all tools

Nsight Compute & Nsight Systems

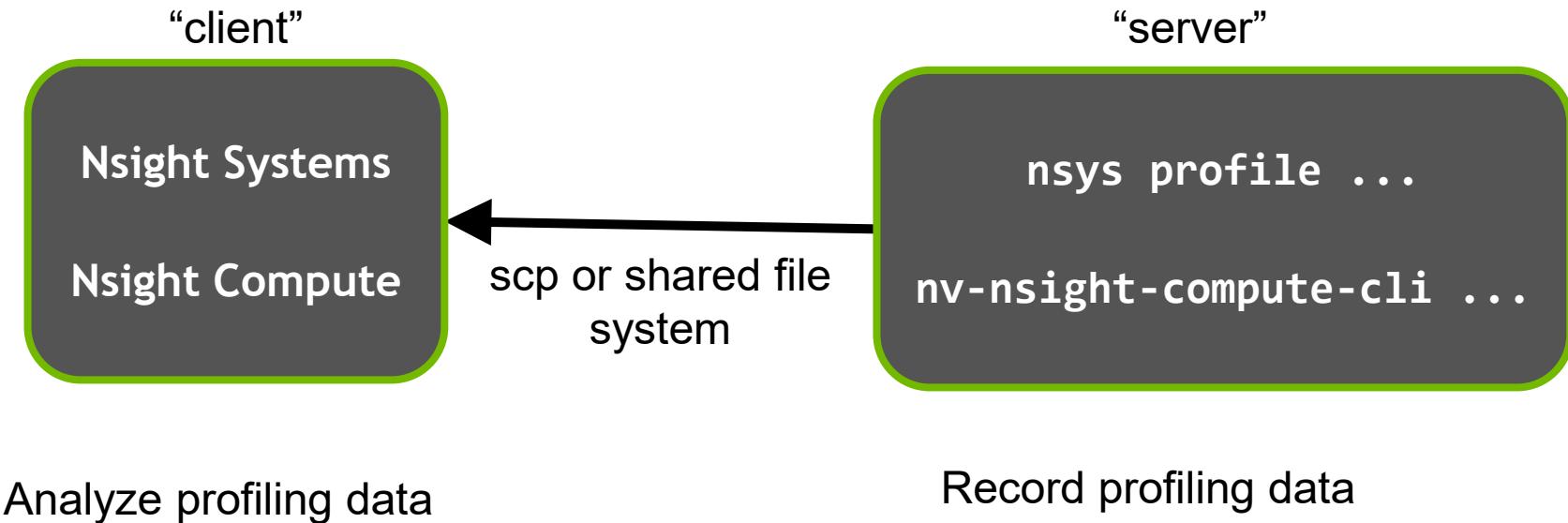
- Command Line and Interactive Profilers
- Bundled with CUDA Toolkit
- Newer standalone downloads available on NVIDIA website
- Nsight Systems
 - “Feeds and speeds:” getting data into the GPU and profiling GPU utilization
- Nsight Compute
 - Kernel-level profiling

A Common GPU Development Model



Two-Phase Profiling

- (interactive profiling also supported)



Before Profiling

- Options to improve your profiling experience
 - Host code annotations with Nvidia Tools Extension Library
 - Correctness & cuda-memcheck
 - Compilation flags
 - Ensure Nsight system environment is correct

Before Profiling: Host Code Annotations

#include <nvToolsExt.h> and link with -lNvToolsExt

Will show up as a named span in the Nsight System GUI

Useful for marking parts of the code for later reference.

```
nvtxRangePush("sleeping");
sleep(100);
nvtxRangePop();
```

Before Profiling: cuda-memcheck

- Certain kinds of errors cause CUDA programs to complete, but crash under profiling
- Check your program with cuda-memcheck if code behaves incorrectly under profiling

```
cuda-memcheck ./my-program ...
```

Before Profiling: Compilation Flags

- Compile device code with optimizations
 - Optimizations dramatically change performance
 - Remove “-G” device-debug flag from nvcc
- Compile device code with line information
 - Minimal information included in binary to map PTX/SASS to source code
 - Add “-lineinfo” flag to nvcc

```
nvcc -G main.cu → nvcc -lineinfo main.cu
```

Before Profiling: Nsight Systems Configuration

Nsight System uses various system hooks to accomplish profiling.

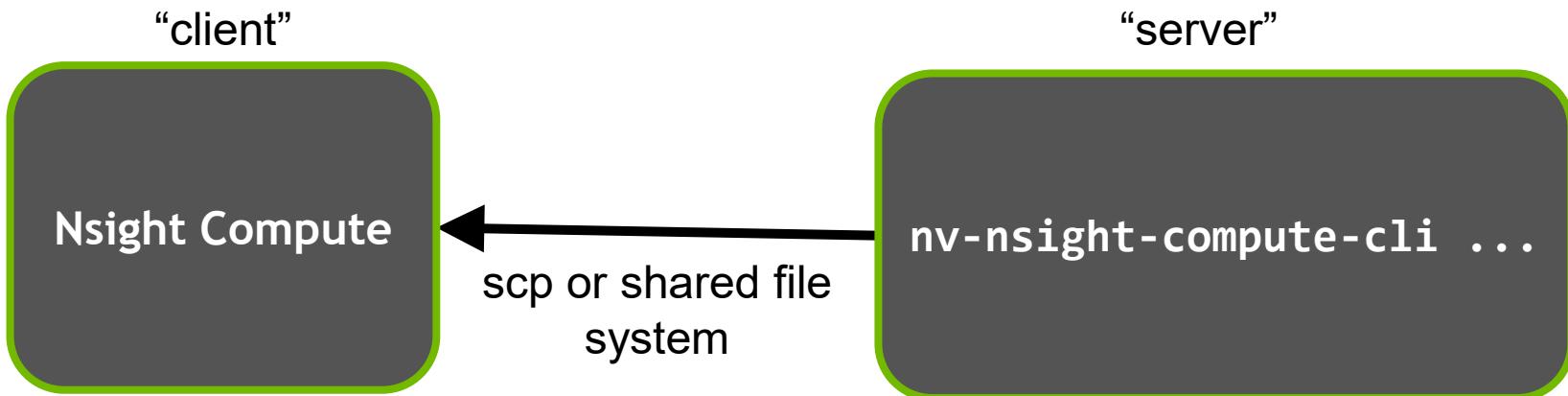
Some errors would reduce the amount or accuracy of gathered info, some will make system profiling impossible. An example of a GOOD output: (check with nsys status -e)

```
$ nsys status -e
>> Sampling Environment Check
>> Linux Kernel Paranoid Level = 2: OK
>> Linux Distribution = Ubuntu
>> Linux Kernel Version = 4.16.15-41615: OK
>> Linux perf_event_open syscall available: OK
>> Sampling trigger event available: OK
>> Intel(c) Last Branch Record support: Available
>> Sampling Environment: OK
```

Consult the Nsight Systems documentation or your system administrator to correct any issues.

Nsight Compute

- Record and analyze detailed kernel performance metrics
- Two interfaces:
 - GUI (nv-nsight-cu)
 - CLI (nv-nsight-cu-cli)
- Directly consuming 1000 metrics is challenging, we use the GUI to help
- Use a two-part record-then-analyze process



Performance Counters

- Device has many performance counters to record detailed information
 - Made available as “metrics”.
- Nsight Compute helps you interpret these

```
$ nv-nsight-cu-cli --devices 0 --query-metrics

lts__t_sectors_srcunit_l1_op_atom_dot_cas
l1tex__data_pipe_lsu_wavefronts_mem_shared_cmd_write
lts__t_sectors_srcunit_l1_aperture_sysmem_op_read
lts__t_requests_op_red_lookup_hit
lts__t_sectors_equiv_l1tagmiss_pipe_tex_mem_texture_op_ld
l1tex__t_bytes_pipe_tex_lookup_miss
l1tex__texin_requests_mem_texture
l1tex__t_bytes_pipe_lsu_mem_local_op_ld_lookup_miss
l1tex__t_bytes_pipe_tex_mem_surface_op_red_lookup_miss
...
```

Record Kernel Traces

- Recording may be for the whole execution, or usually restricted to a single launch of a kernel

```
$ nv-nsight-cu-cli      \ <nsight compute cli>
--kernel-id ::mygemm:6 \ <6th launch of “mygemm” kernel>
--section “.*”      \ <output file name>
-o output_file      \ <executable to run>
executable
```

Open Nsight Compute

- We will import the recorded file

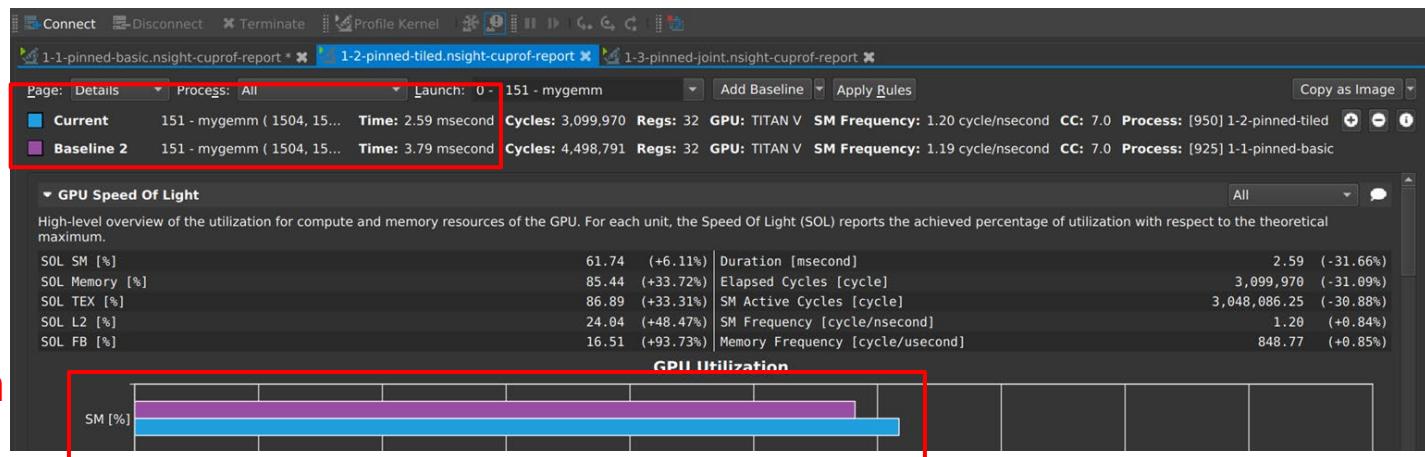
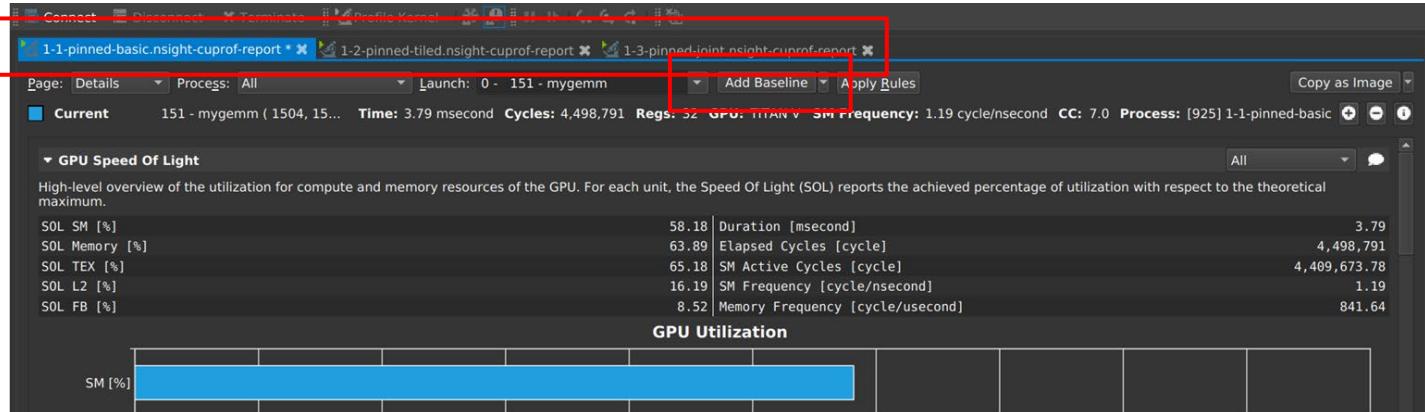
```
File > Open File ... > output_file.nsight-cuprof-report
```

- We can open multiple files in multiple tabs, if desired

First Look

- Can compare multiple open codes with baseline button
 - For comparing effect of optimizations

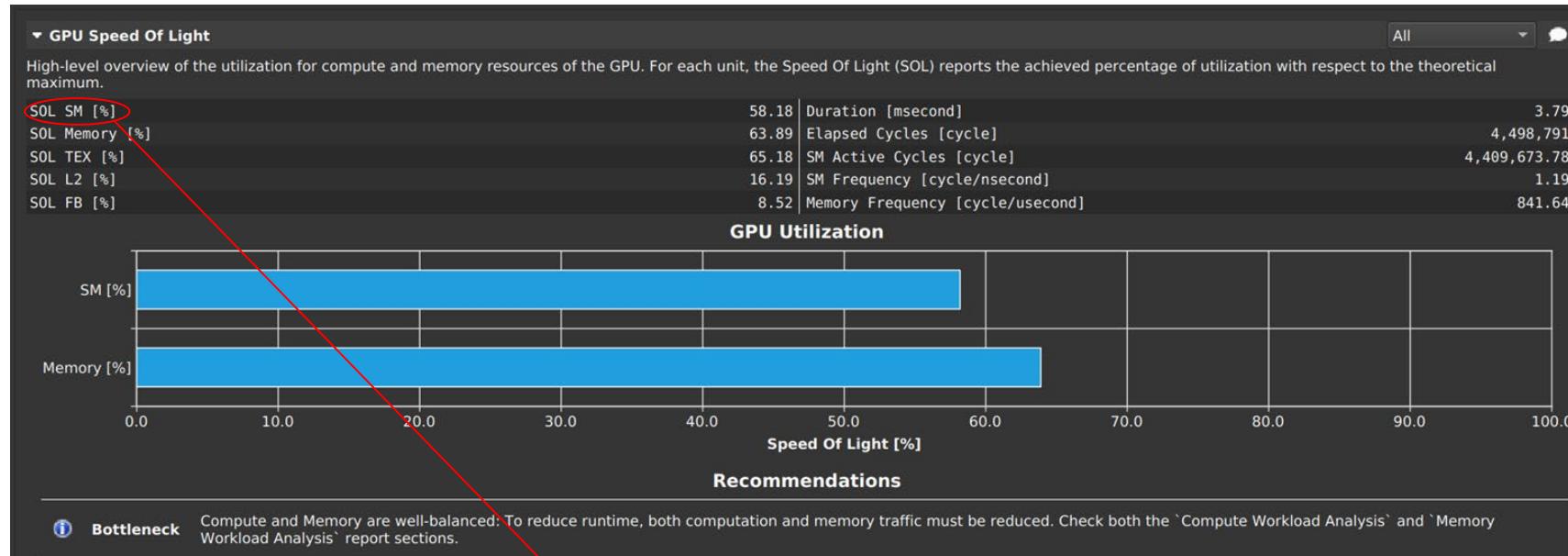
Tabs and baseline button



Next tab now has comparison

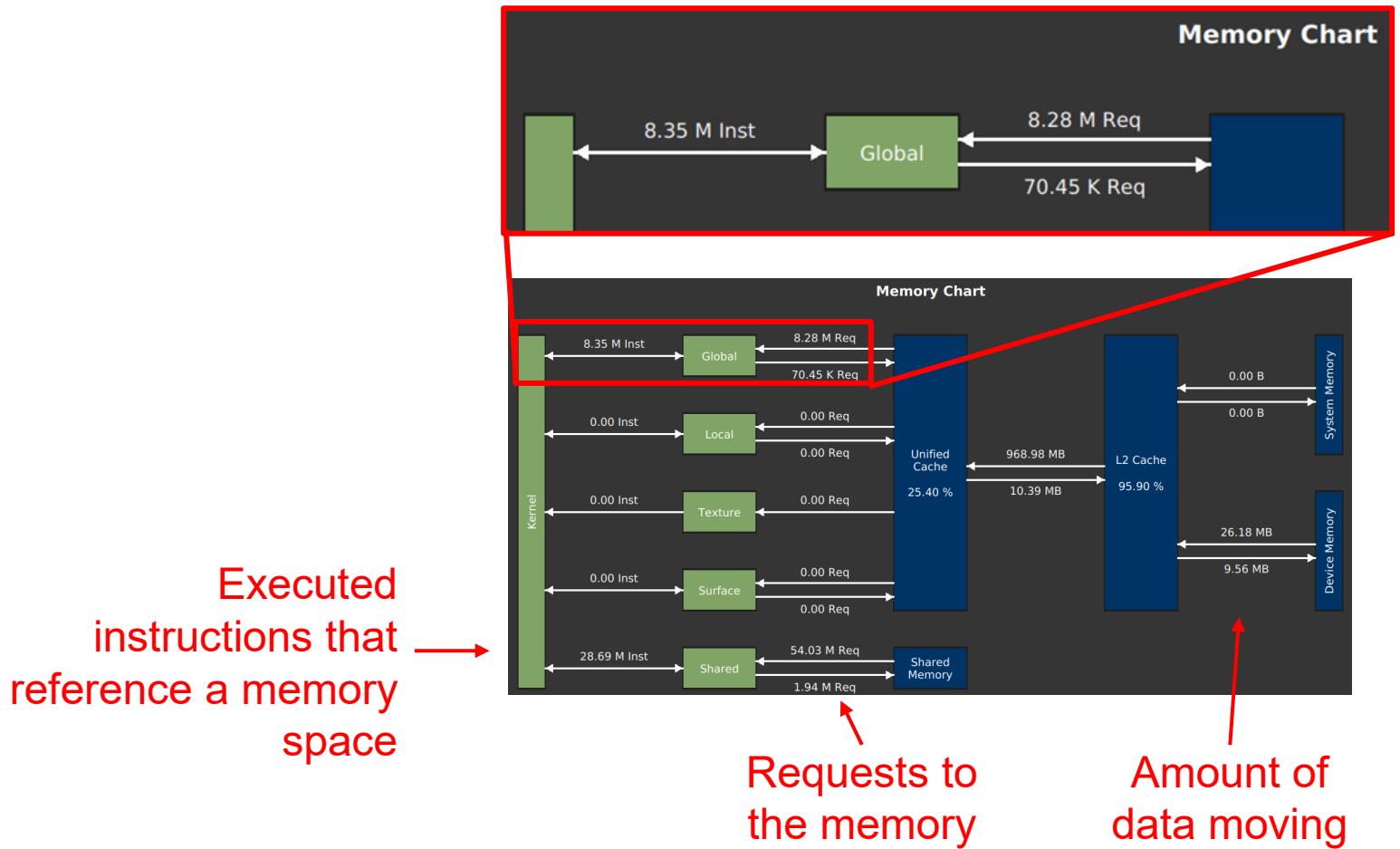
GPU Speed of Light

- Utilization compared to theoretical maximums



Tip: mouse over each to see the associated metric

Workload Memory Analysis Chart

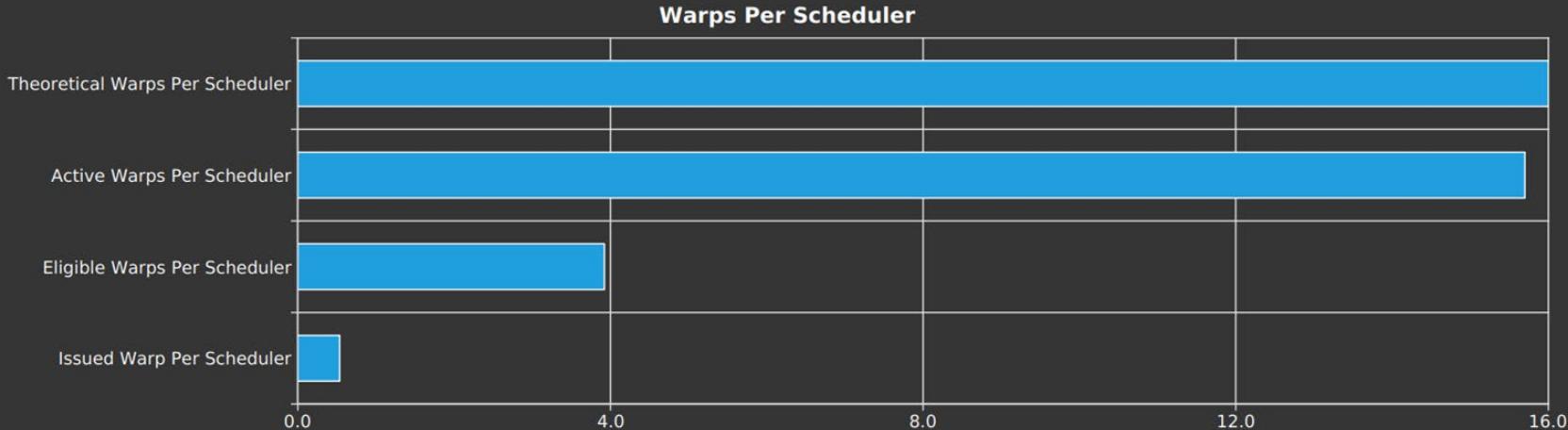


Scheduler Statistics

- Activity of the scheduler issuing instructions.
 - Maximum possible warps per scheduler
 - Warps that have not exited a kernel
 - Warps that have execution dependencies satisfied
 - Warps actually issued

Summary of the activity of the schedulers issuing instructions. Each scheduler maintains a pool of warps that it can issue instructions for. The upper bound of warps in the pool (Theoretical Warps) is limited by the launch configuration. On every cycle each scheduler checks the state of the allocated warps in the pool (Active Warps). Active warps that are not stalled (Eligible Warps) are ready to issue their next instruction. From the set of eligible warps the scheduler selects a single warp from which to issue one or more instructions (Issued Warp). On cycles with no eligible warps, the issue slot is skipped and no instruction is issued. Having many skipped issue slots indicates poor latency hiding.

Active Warps Per Scheduler [warp]	15.70	Instructions Per Active Issue Slot [inst/cycle]	1
Eligible Warps Per Scheduler [warp]	3.92	No Eligible [%]	46.25
Issued Warp Per Scheduler	0.54	One or More Eligible [%]	53.75



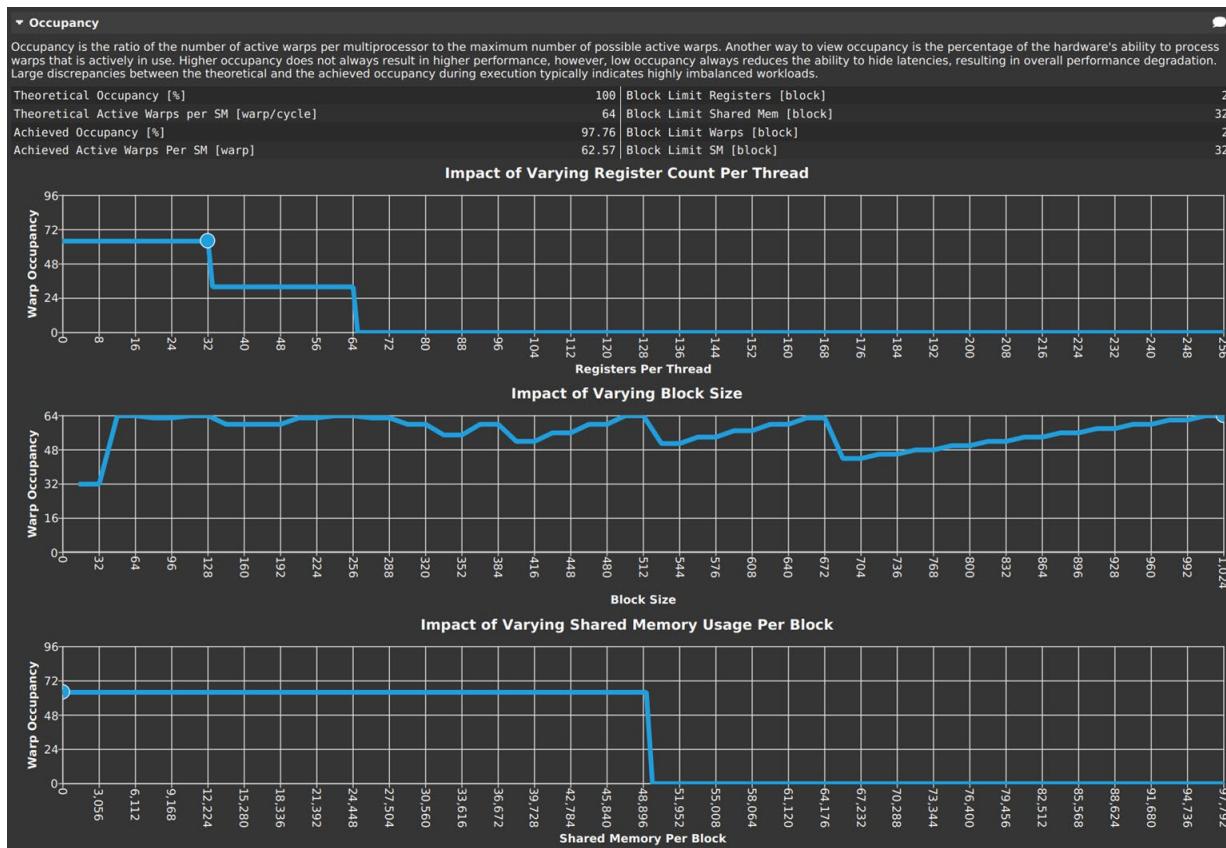
Warp State Statistics

- How much time warps spend in each state



Occupancy

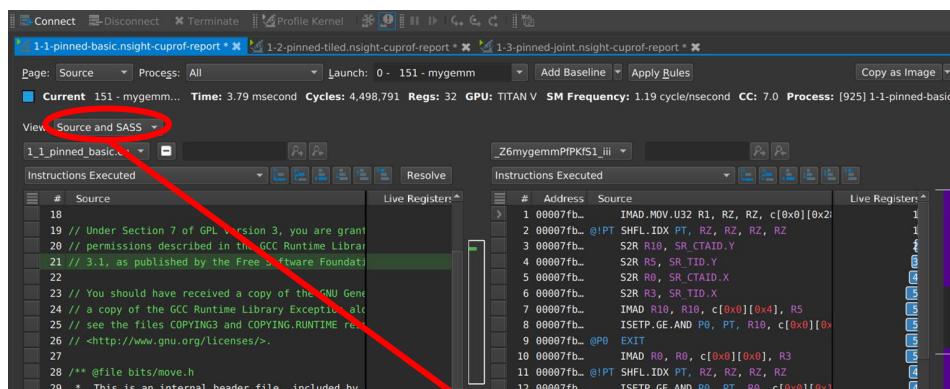
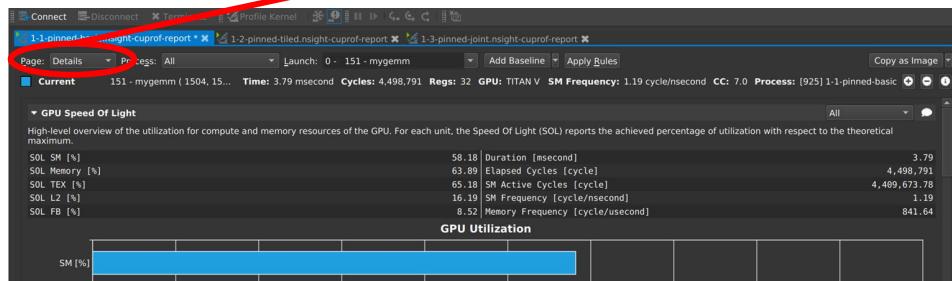
- How many warps are active compared to the maximum possible
 - Achieved: true number of active warps as average
 - Charts show how kernel resources and grid dimension affect occupancy



Instruction Hotspots

- Show metrics by source line, PTX instruction, or SASS instruction

Switch page to “Source”



“Source and PTX” (usually) or “Source and SASS”

Instruction Hotspots

- Source file may not load if profiling was recorded on another system

The screenshot shows the NVIDIA Nsight Compute interface. At the top, there are three tabs: '1-1-pinned-basic.nsight-cuprof-report', '1-2-pinned-tiled.nsight-cuprof-report', and '1-3-pinned-joint.nsight-cuprof-report'. Below the tabs, the 'Current' tab is selected, showing details: Process: All, Launch: 0 - 151 - mygemm, Time: 3.79 msecond, Cycles: 4,498,791, Rgs: 32, GPU: TITAN V, SM Frequency: 1.19 cycle/nsecond, CC: 7.0, Process: [925] 1-1-pinned-basic. The 'View' dropdown is set to 'Source and SASS'. The left pane, titled 'Instructions Executed', displays the source code of '1_1_pinned_basic.cu':

```
# Source
18
19 // Under Section 7 of GPL version 3, you are grant
20 // permissions described in the GCC Runtime Library
21 // 3.1, as published by the Free Software Foundat
22
23 // You should have received a copy of the GNU Gene
24 // a copy of the GCC Runtime Library Exception alc
25 // see the files COPYING3 and COPYING.RUNTIME resp
26 // <http://www.gnu.org/licenses/>.
27
28 /** @file bits/move.h
29 * This is an internal header file, included by
```

The right pane, titled 'Instructions Executed', shows the assembly instructions for the current process:

#	Address	Source	Live Registers
1	00007fb...	IMAD.MOV.U32 R1, RZ, RZ, c[0x0][0x2]	1
2	00007fb...	@!PT SHFL.IDX PT, RZ, RZ, RZ, RZ	1
3	00007fb...	S2R R10, SR_CTAID.Y	4
4	00007fb...	S2R R5, SR_TID.Y	5
5	00007fb...	S2R R0, SR_CTAID.X	4
6	00007fb...	S2R R3, SR_TID.X	5
7	00007fb...	IMAD R10, R10, c[0x0][0x4], R5	5
8	00007fb...	ISETP.GE.AND P0, PT, R10, c[0x0][0x4]	5
9	00007fb...	@P0 EXIT	5
10	00007fb...	IMAD R0, R0, c[0x0][0x8], R3	5
11	00007fb...	@!PT SHFL.IDX PT, RZ, RZ, RZ, RZ	4
12	00007fb...	ISETP.GE.AND P0, PT, R0, c[0x0][0x4]	4

Click “resolve” and find your local copy of the code that was compiled or run remotely.

Instruction Hotspots

Click a line to highlight lines from other side...

#	Source	Live Registers	Sampling Data (All)	Sampling Data (Not Issued)	Instructions Executed
1	1_pinned_basic.cu:		0	0	
2	#include <algorithm>		0	0	
3			0	0	
4	#include <nvToolsExt.h>		0	0	
5			0	0	
6	#include <argparse/argparse.hpp>		0	0	
7			0	0	
8	#include "common.hpp"		0	0	
9			0	0	
10	/* NOTE: A and C are column major, B is row major		0	0	
11	*/		0	0	
12	<u>_global_ void mygemm(float * __restrict__ c, //</u>		0	0	
13	const float *a, //		0	0	
14	const float *b, //		0	0	
15	const int M, const int N,		0	0	
16			0	0	
17	#define A(i, j) A[(i) + (j)*M]		0	0	
18	#define B(i, j) b[(i)*N + (j)]		0	0	
19	#define C(i, j) C[(i) + (j)*M]		0	0	
20			0	0	
21	int gidx = blockDim.x * blockIdx.x + threadIdx.x;	3	111	50	
22	int gidy = blockDim.y * blockIdx.y + threadIdx.y;	5	39	15	
23		0	0	0	
24	for (int i = gidy; i < M; i += gridDim.y * blo	5	61	35	
25	for (int j = gidx; j < N; j += gridDim.x * b	10	63	18	
26	float acc = 0;	0	0	0	
27	for (int k = 0; k < k; ++k) {	26	171,378	53,73	
28	acc += A(i, k) * B(k, j);	10	454	331	
29	}	0	0	0	
30	C(i, j) = acc;	10	454	331	
31	}	0	0	0	
32	}	0	0	0	
33		0	0	0	
34	#undef A	0	0	0	
35	#undef B	0	0	0	
36	#undef C	0	0	0	
37	}	0	0	0	
38		356	333	0	

Program counter spends most of its time on instructions from this line.
Mouse over for breakdown.

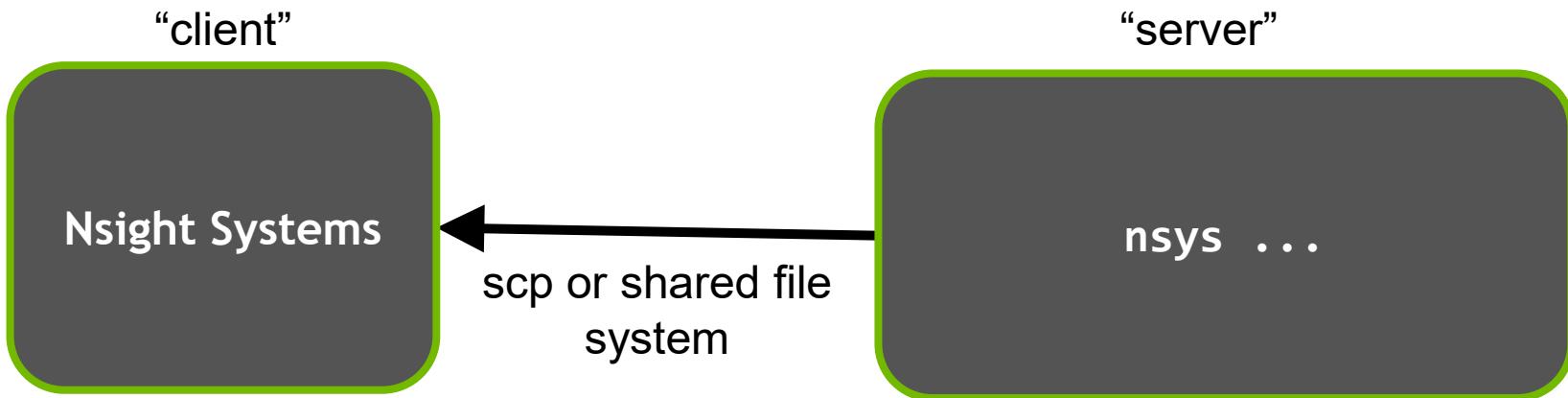
...corresponding PTX/SASS lines over here.

#	Source	Live Registers	Sampling Data (All)	Sampling Data (Not Issued)	Instructions Executed	
144	B80.13:	28	3,564	478	32,682,061	
145	mul.wide.s32 %rd21, %rs4, 4;	28	0	0	26,103,655	
146	add.s64 %rd22, %rd1, %rd21;	28	12,047	0	26,103,655	
147	cvt.a.to.global.u64 %rd23, %rd5;	29	22,495	0	26,103,655	
148	mult.wide.s32 %rd24, %rs3, 4;	28	0	0	26,103,655	
149	add.s64 %rd25, %rd3, %rd24;	28	13,128	5,166	26,103,655	
150	fma.rn.f32 %f22, %f21, %f20, %f35;	29	0	0	26,103,655	
151	ld.global.f32 %f21, [%rd21];	28	139,966	0	26,103,655	
152	shl.b32 %r45, %r24, 2;	13	9	0	139,966	
153	cvt.s64.s32 %rd26, %r45;	14	26	2	139,966	
154	add.s64 %rd27, %rd22, %rd26;	27	1,560	138	52,287,318	
155	shl.b32 %r46, %r25, 2;	13	1	0	139,966	
156		0	0	0		
157	cvt.s64.s32 %rd26, %r45;	15	4	0	139,966	
158	add.s64 %rd27, %rd22, %rd26;	27	2,854	412	52,287,318	
159	ld.global.f32 %f23, [%rd23];	28	9,326	3,169	26,103,655	
160	fma.rn.f32 %f25, %f24, %f23, %f22;	28	9,598	3,304	26,103,655	
161	add.s64 %rd31, %rd29, %rd28;	28	1,888	192	52,287,318	
162	ld.global.f32 %f26, [%rd31];	28	0	0		
163	ld.global.f32 %f27, [%rd32];	28	10,914	3,614	26,103,655	
164	ld.global.f32 %f28, [%rd33];	27	7	7,437	2,094	26,103,655
165	fma.rn.f32 %f25, %f24, %f23, %f22;	28	0	0		
166	add.s64 %rd31, %rd29, %rd28;	28	0	0		
167	add.s64 %rd31, %rd29, %rd28;	28	11,190	3,986	26,103,655	
168		27	10,914	3,614	26,103,655	
169	ld.global.f32 %f26, [%rd31];	27	0	0		
170	ld.global.f32 %f27, [%rd32];	28	0	0		
171	fma.rn.f32 %f28, %f27, %f26, %f25;	27	0	0		
172	add.s64 %rd32, %rd30, %rd26;	28	2,214	561	52,287,318	
173	add.s64 %rd33, %rd31, %rd28;	28	3,305	561	52,287,318	
174		28	0	0		
175	ld.global.f32 %f29, [%rd33];	27	19,188	6,837	26,103,655	
176	ld.global.f32 %f30, [%rd34];	27	9,524	3,664	26,103,655	
177	fma.rn.f32 %f35, %f30, %f29, %f28;	28	4,301	1,146	26,103,655	
178		28	0	0		
179	add.s64 %r54, %r54, %r45;	24	244	0	26,033,676	
180	add.s32 %r53, %r53, %r46;	26	232	0	26,033,676	
181		28	0	0		

Sometimes, stalls can show in a following instruction that depends on a previous one

Nsight Systems

- Record and analyze system utilization information
- Two interfaces:
 - GUI (nsight-sys)
 - CLI (nsys)
- Use a two-part record-then-analyze process



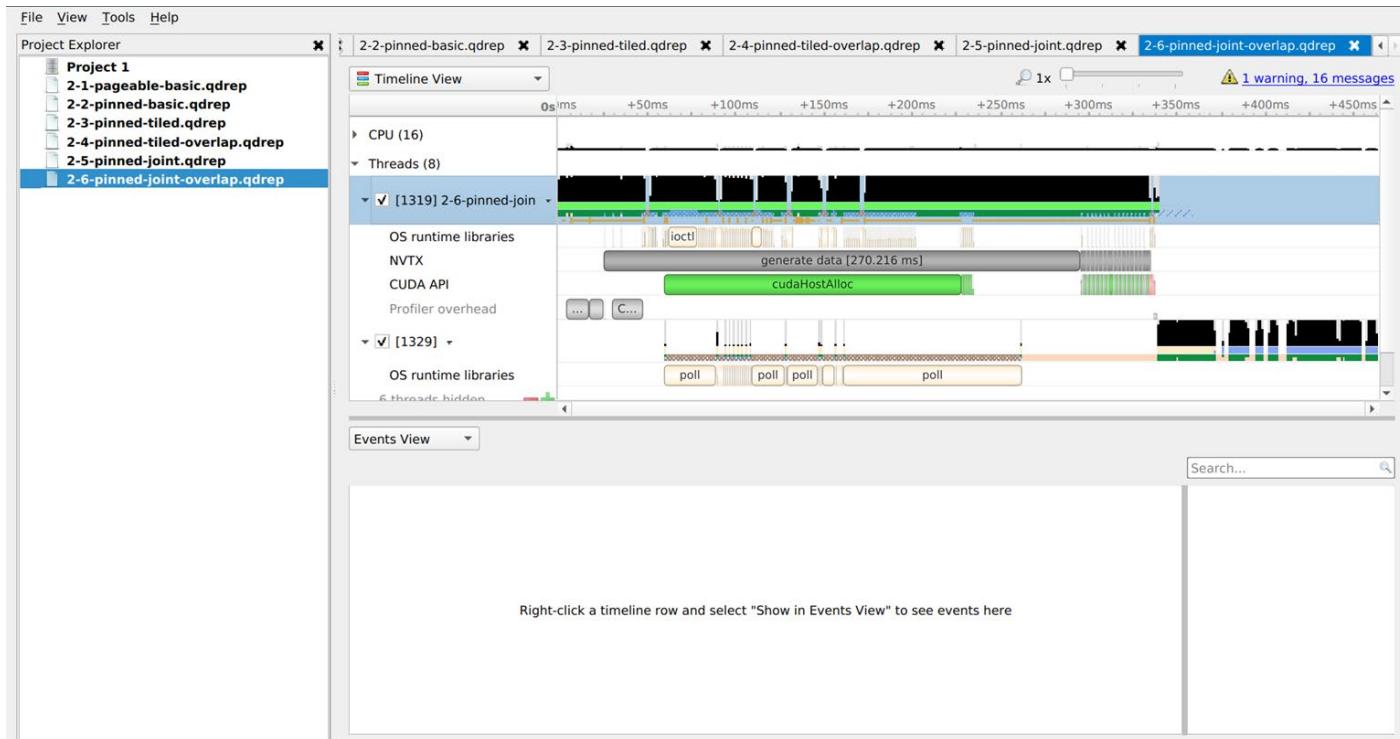
Record Execution Traces

- Record system information for the entire execution

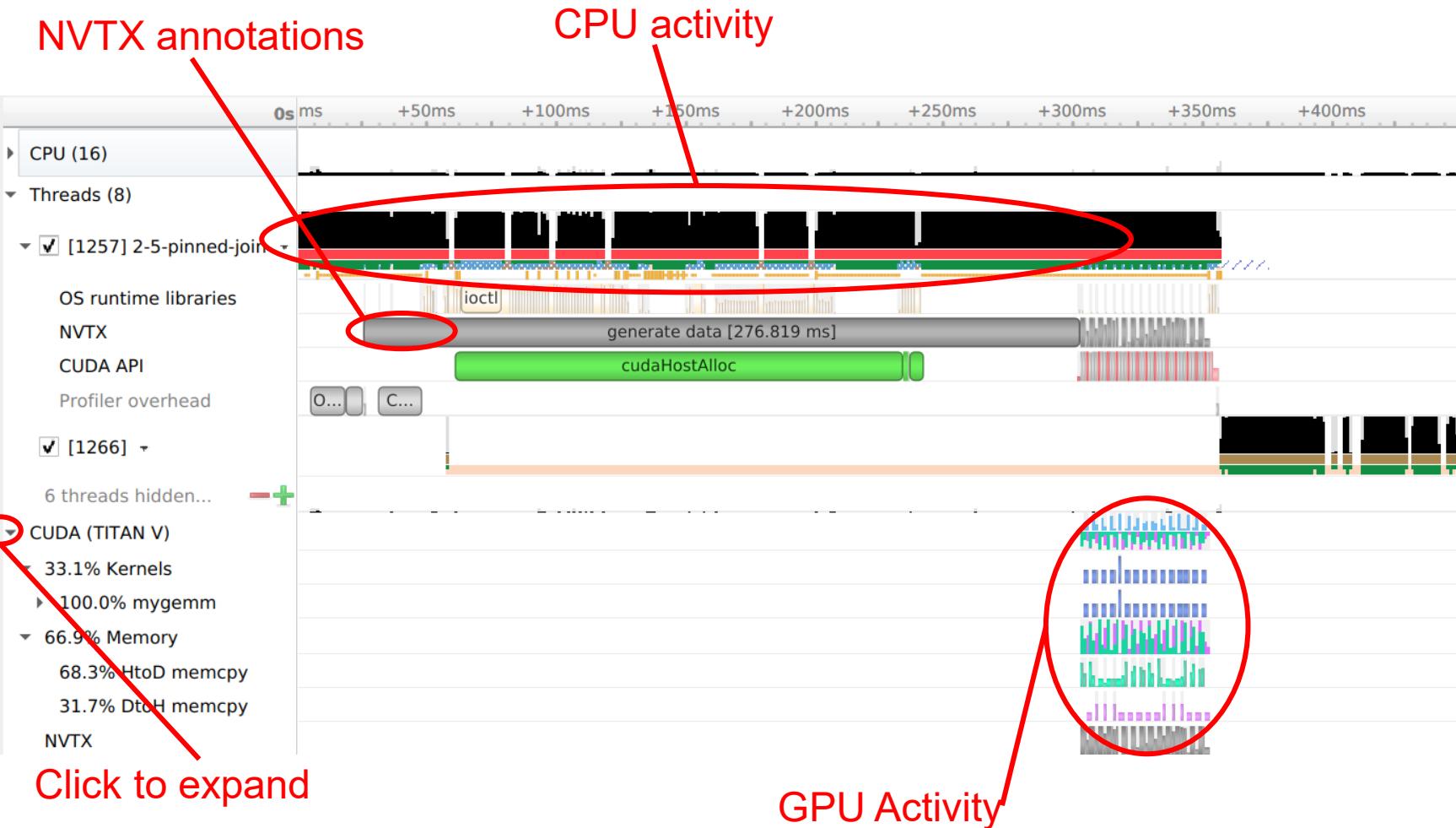
```
$ nsys profile.      \ <nSight Systems CLI>
  -o output_file    \ <record file>
  executable        <executable to run>
```

First Look

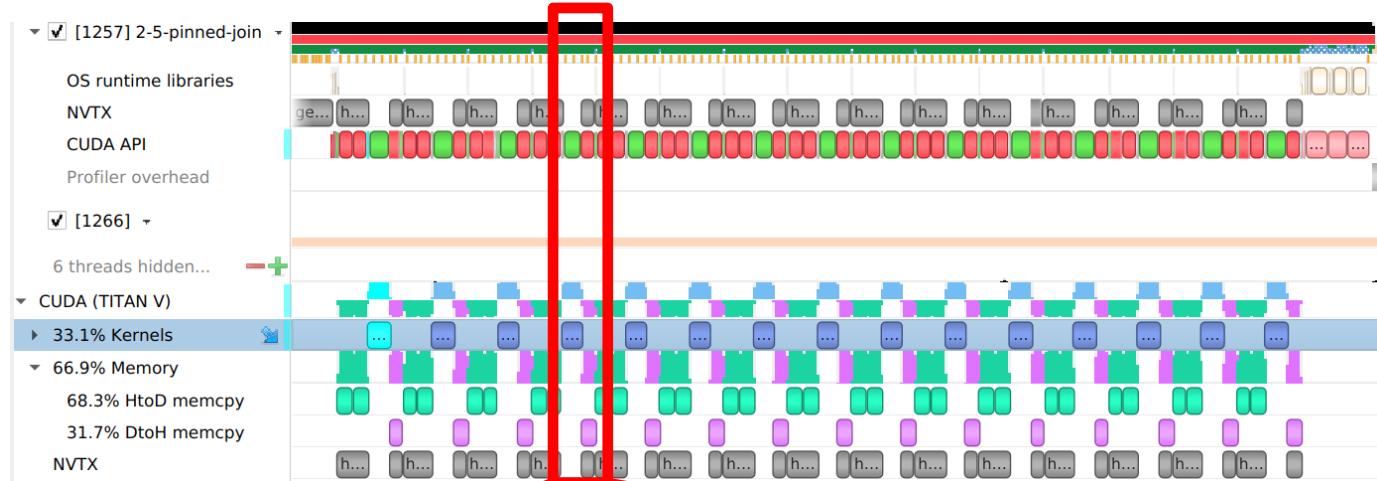
- File > Open > output_file.qdrep
- If multiple files are open, shown in the left pane
- Timeline of OS calls, CUDA calls, NVTX spans, and GPU activity



Timeline View



Click and drag to zoom in



Mouse over spans for more info





GPU Teaching Kit



The GPU Teaching Kit is licensed by NVIDIA and the University of Illinois under the [Creative Commons Attribution-NonCommercial 4.0 International License](#).

GPU Teaching Kit

Accelerated Computing

Lecture 2.6 - Introduction to CUDA C Unified Memory

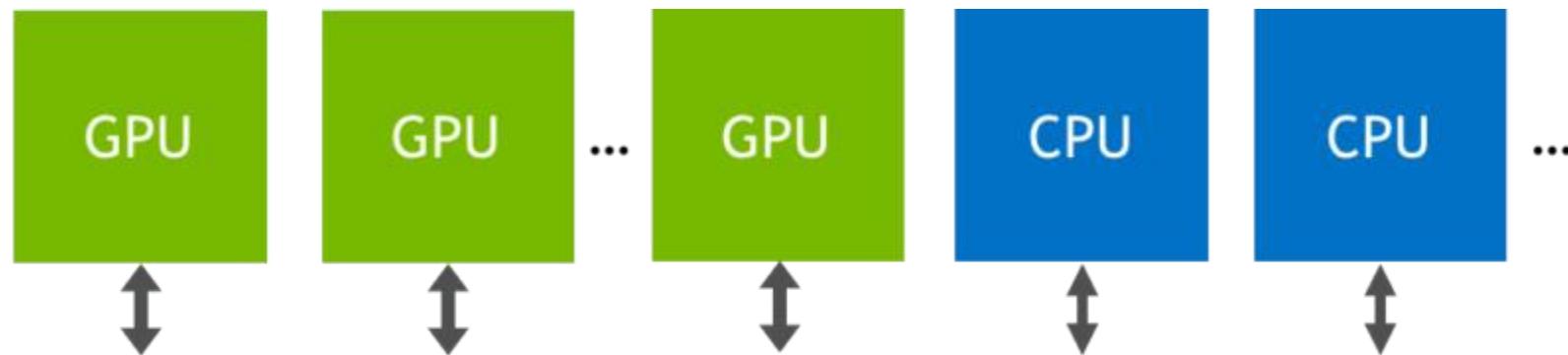


Objective

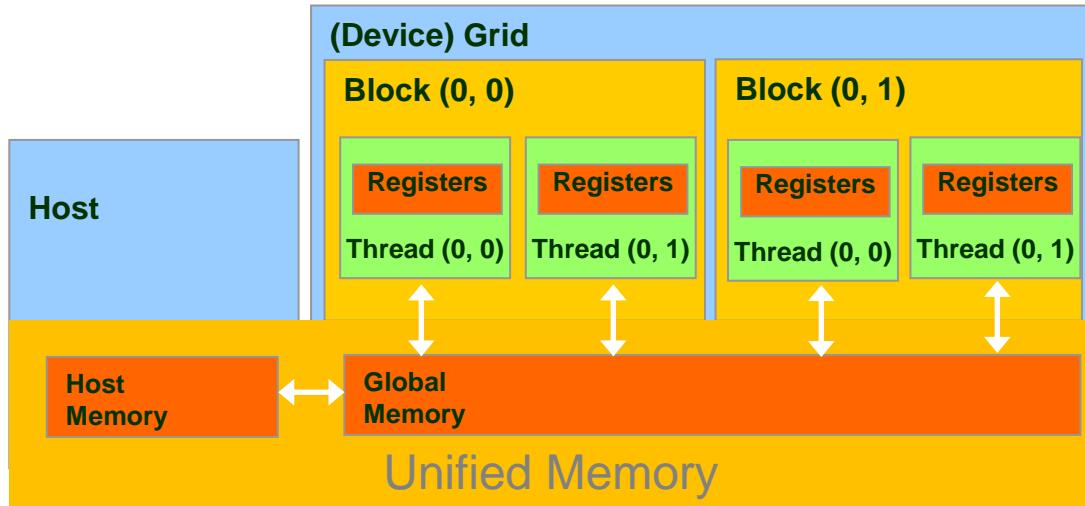
- To learn the basic API functions in CUDA host code for CUDA Unified Memory
 - Unified Memory Allocation
 - Data Transfer in Unified Memory

CUDA Unified Memory (UM)

- Is a single memory address space accessible both from the host and from the device.
- The hardware/software handles automatically the data migration between the host and the device maintaining consistency between them.

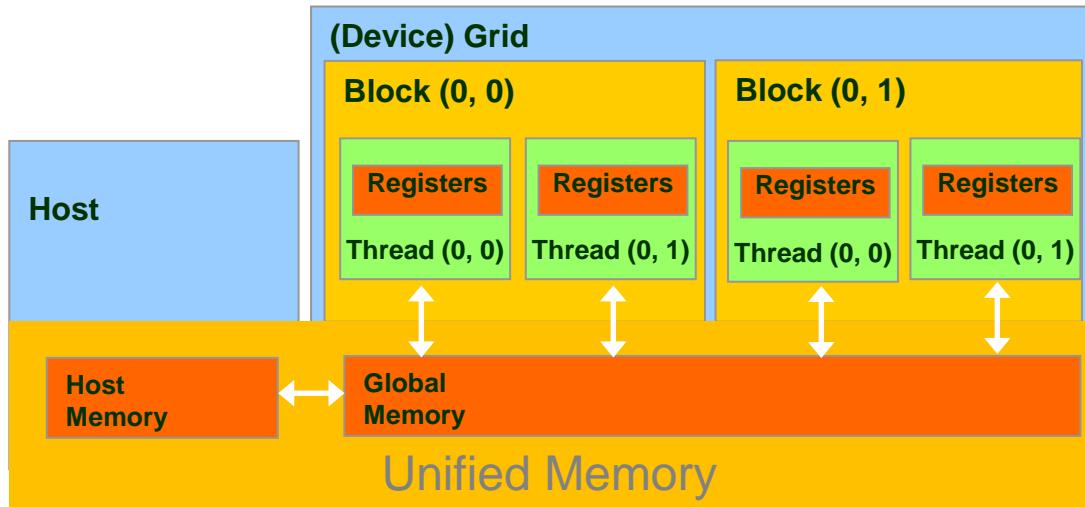


Partial Overview of CUDA Memories



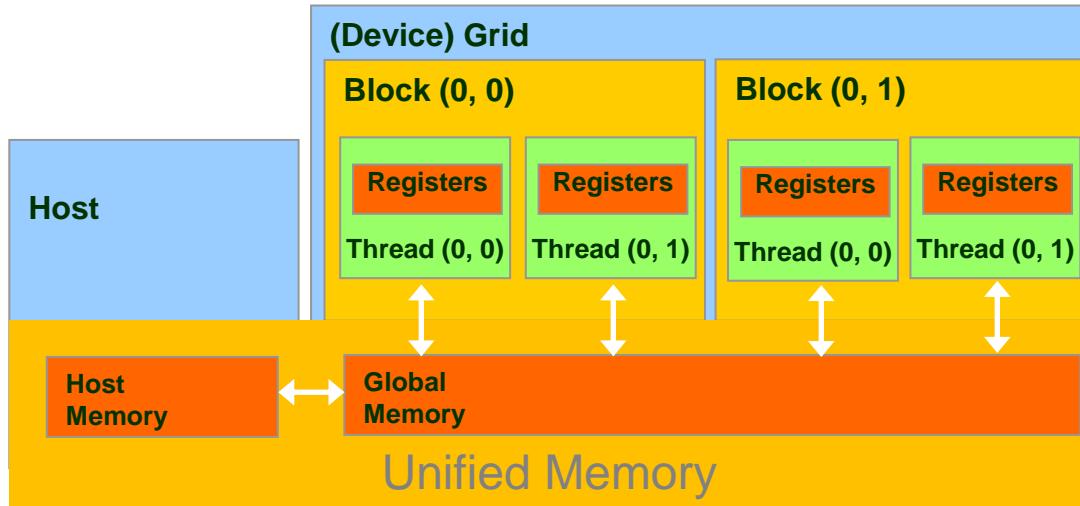
- Device code can:
 - R/W per-thread registers
 - R/W all-shared global memory
 - R/W managed memory (Unified Memory)
- Host code can
 - Transfer data to/from per grid global memory
 - R/W managed memory

Partial Overview of CUDA Memories



- **cudaMallocManaged()**
 - Allocates an object in the Unified Memory address space.
 - Two parameters, with an optional third parameter.
 - Address of a pointer to the allocated object
 - Size of the allocated object in terms of bytes
 - [Optional] Flag indicating if memory can be accessed from any device or stream
- **cudaFree()**
 - Frees object from unified memory.
 - One parameter
 - Pointer to freed object

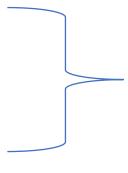
Partial Overview of CUDA Memories



- **cudaMemcpy()**
 - Memory data transfer
 - Requires four parameters
 - Pointer to destination
 - Pointer to source
 - Number of bytes copied
 - Type/Direction of transfer
 - Depending on the transfer type, the driver may decide to use the memory on the host or the device.
 - In Unified Memory this function is utilized to copy data between different arrays, regardless of position.

Putting it all together, vecAdd CUDA host code using Unified Memory

```
int main() {  
  
    float *m_A, float *m_B, float *m_C, int n;  
  
    int size = n * sizeof(float);  
  
    cudaMallocManaged((void**) &m_A, size);  
    cudaMallocManaged((void**) &m_B, size);  
    cudaMallocManaged((void**) &m_C, size);  
  
    // Memory initialization on the Host  
    // Kernel invocation code - to be shown later  
    cudaFree(m_A); cudaFree(m_B); cudaFree(m_C);  
}
```



Allocation of Managed Memory

m_A, m_B gets initialized on the host

The device performs the actual vector addition

CUDA Unified Memory for different architectures

Prior to compute capability 6.x

- There is no specialized hardware units to improve UM efficiency.
- For data migration the full memory block needs to be copied synchronically by the driver.
- No memory oversubscription.

Compute capability 6.x onwards

- There are specialized hardware units managing page faulting.
- Data is migrated on demand, meaning that data gets copied only on page fault.
- Possibility to oversubscribe memory, enabling larger arrays than the device memory size.

GPU Teaching Kit

Accelerated Computing

The GPU Teaching Kit is licensed by NVIDIA under the [Creative Commons Attribution-NonCommercial 4.0 International License..](#)

