



Leibniz-Rechenzentrum
der Bayerischen Akademie der Wissenschaften



DEEP
LEARNING
INSTITUTE



DEEP
LEARNING
INSTITUTE

Fundamentals of Accelerated Computing with CUDA C/C++

Dr. Momme Allalen LRZ | 07.11.2023

Fundamentals of Accelerated Computing with CUDA C/C++



- You learn the basics of **CUDA C/C++** by:
 - Accelerating CPU-only applications to run their latent parallelism on GPUs.
 - Utilizing essential **CUDA memory** management techniques to optimize accelerated applications
 - Exposing accelerated application potential for concurrency and exploiting it with **CUDA streams**
 - Leveraging command line and visual profiling to guide and check your work.
- Upon completion, you'll be able to accelerate and optimize existing C/C++ CPU-only applications using the most essential **CUDA tools** and techniques. You'll understand an iterative style of CUDA development that will allow you to ship accelerated applications fast.

Tentative Agenda



DEEP
LEARNING
INSTITUTE



10:00-10:15 Intro@**CUDA**

10:15-12:00 Accelerating Applications with **CUDA C/C++**

12:00-13:00 Lunch break

13:00-14:20 Managing Accelerated Application Memory
with **CUDA** Unified Memory and **nsys**

14:20-14:30 Coffee break

14:30-15:45 Asynchronous Streaming and Visual Profiling for
Accelerated Applications with **CUDA C/C++**

15:45-16:00 Q&A, Final Remarks

Workshop Webpage



DEEP
LEARNING
INSTITUTE



- **Lecture material will be made available under:**
 - <https://tinyurl.com/hdli3w23>
- **Access CUDA C/C++ Code :**
 - See the **Chat Window**

Training Setup



- To get started, follow these steps:
- Create an NVIDIA Developer account at <http://courses.nvidia.com/join> Select "Log in with my NVIDIA Account" and then "Create Account".
- If you use your own laptop, make sure that WebSockets works for you:
Test your Laptop at <http://websocketstest.com>
 - Under ENVIRONMENT, confirm that "WebSockets" is checked yes.
 - Under WEBSOCKETS (PORT 80]. confirm that "Data Receive", "Send", and "Echo Test" are checked yes.
 - If there are issues with WebSockets, try updating your browser.
We recommend Chrome, Firefox, or Safari for an optimal performance.
- Visit <http://courses.nvidia.com/dli-event> and enter the event code provided by the instructor.
- You're ready to get started.

And now



DEEP
LEARNING
INSTITUTE



Enjoy the course !

Why do we need to program for GPU?

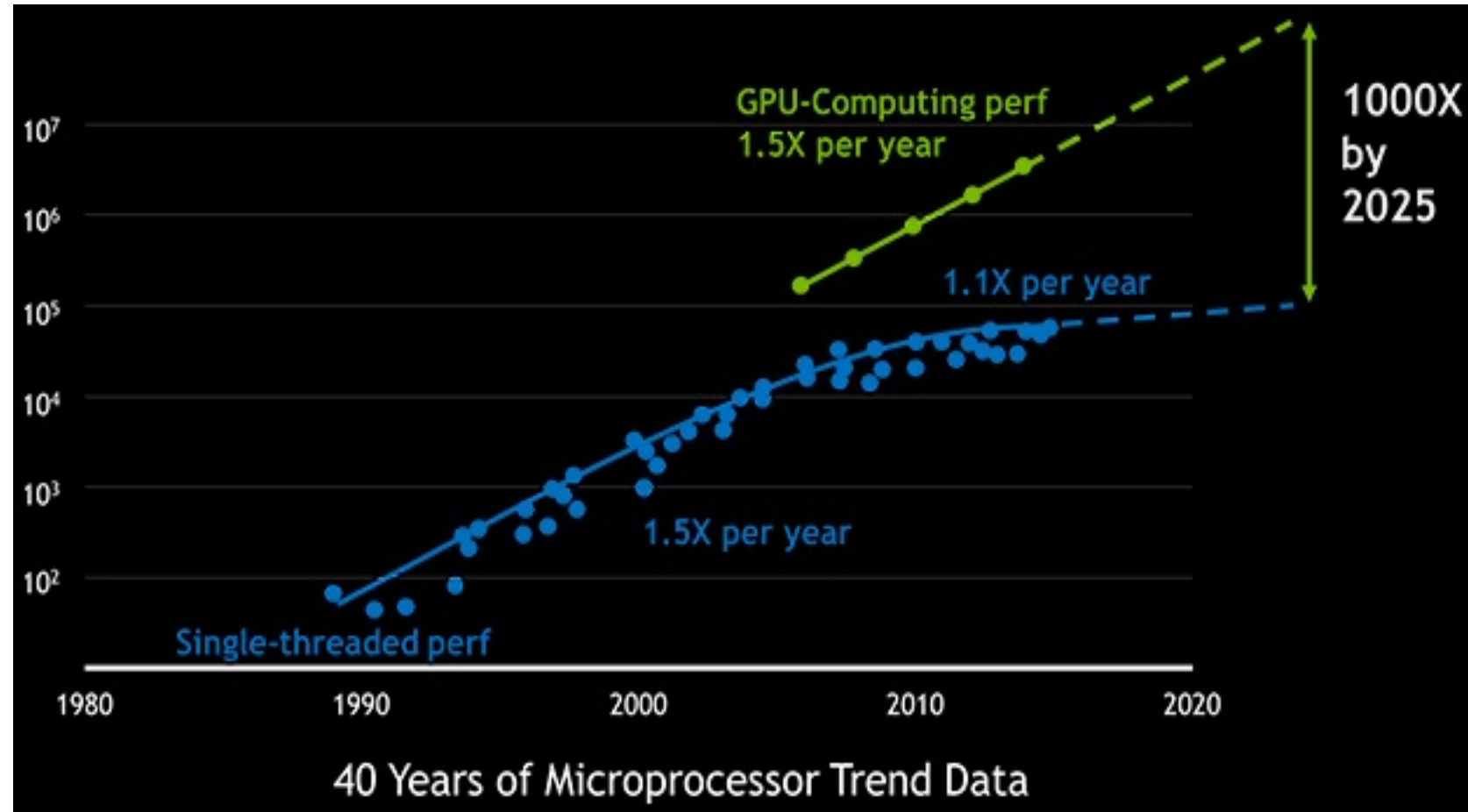


Moore's law is dead !

The long-held notion that the processing power of computers increases exponentially every couple of years has hit its limit....

The free lunch is over..

Future is parallel !



Why do we need to program for GPU?

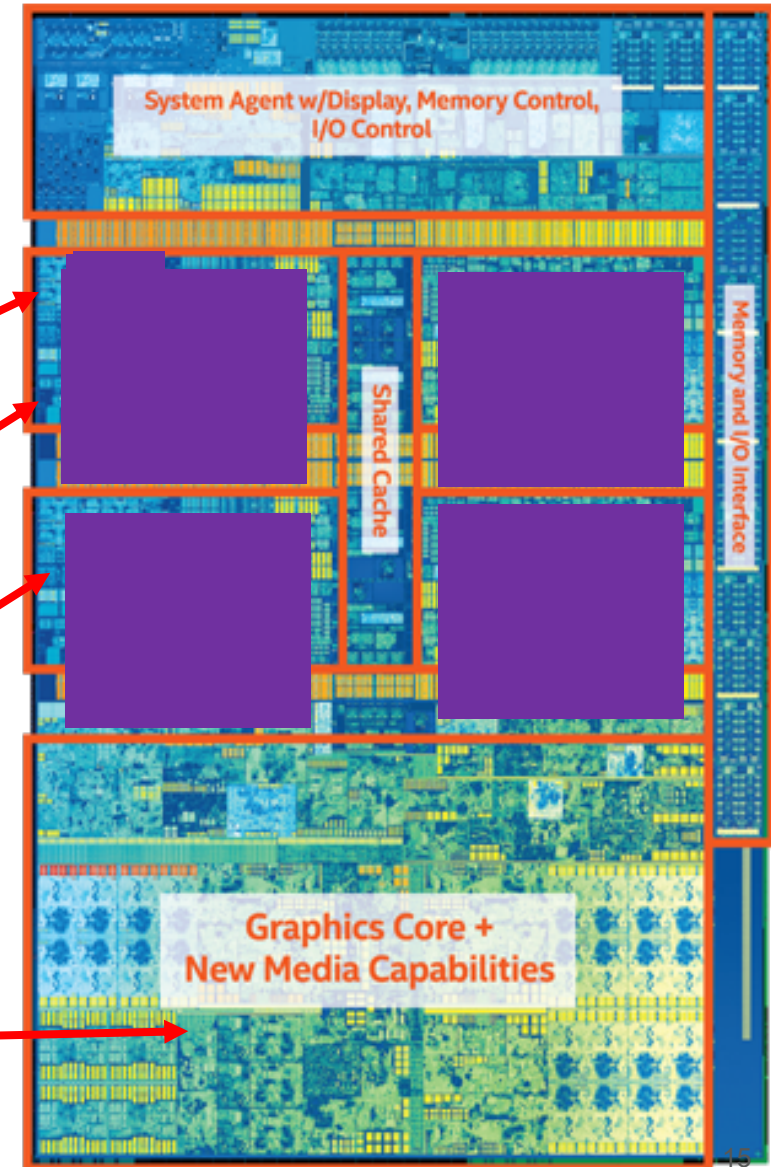
- Typical example Intel chip: Core i7 Gen Kaby Lake processors

- 4*CPU cores
- With hyperthreading
- Each with 8-wide AVX instructions
- GPU with 1280 processing elements

- Programming on chip:

- Serial C/C++ .. Code only takes advantage of a very small amount of the available resources of the chip.
- Using vectorisation allows you to fully utilise the resources of a single hyper-thread
- Using multi-threading allows you to fully utilise all CPU cores

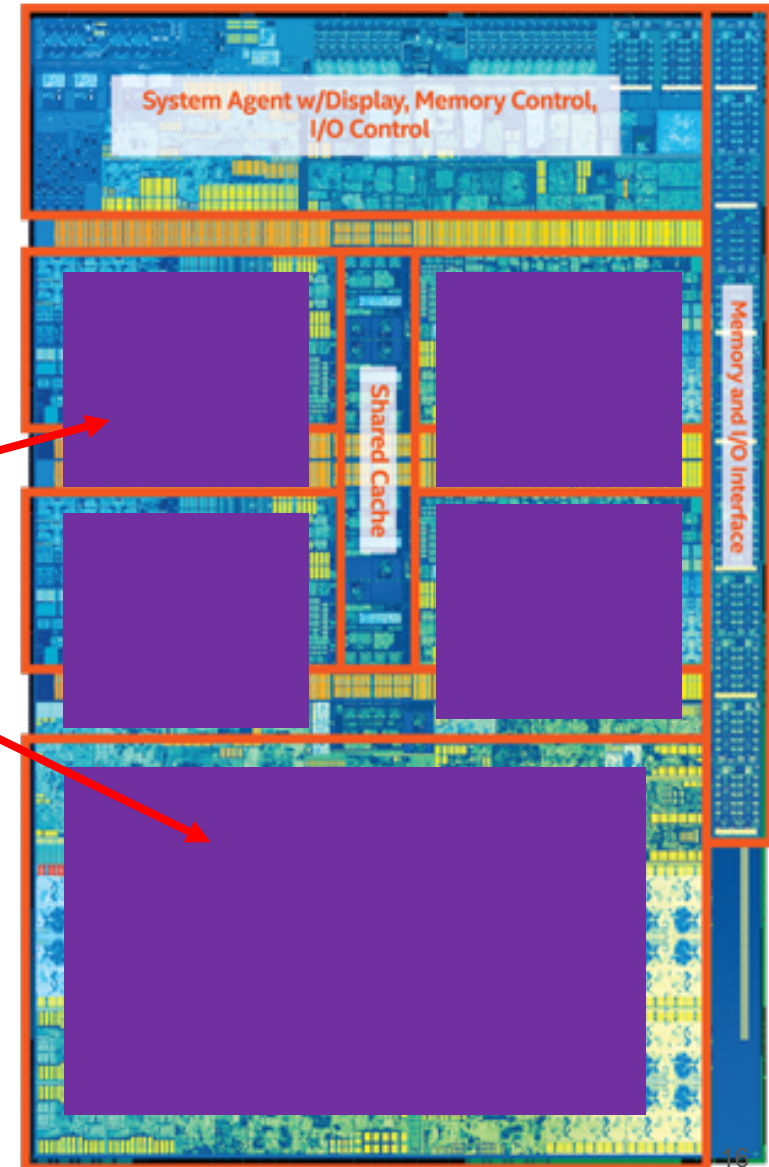
GPU need to be used ?



Why do we need to program for GPU?



- Using heterogeneous programming allows you to dispatch and fully utilise the entire chip



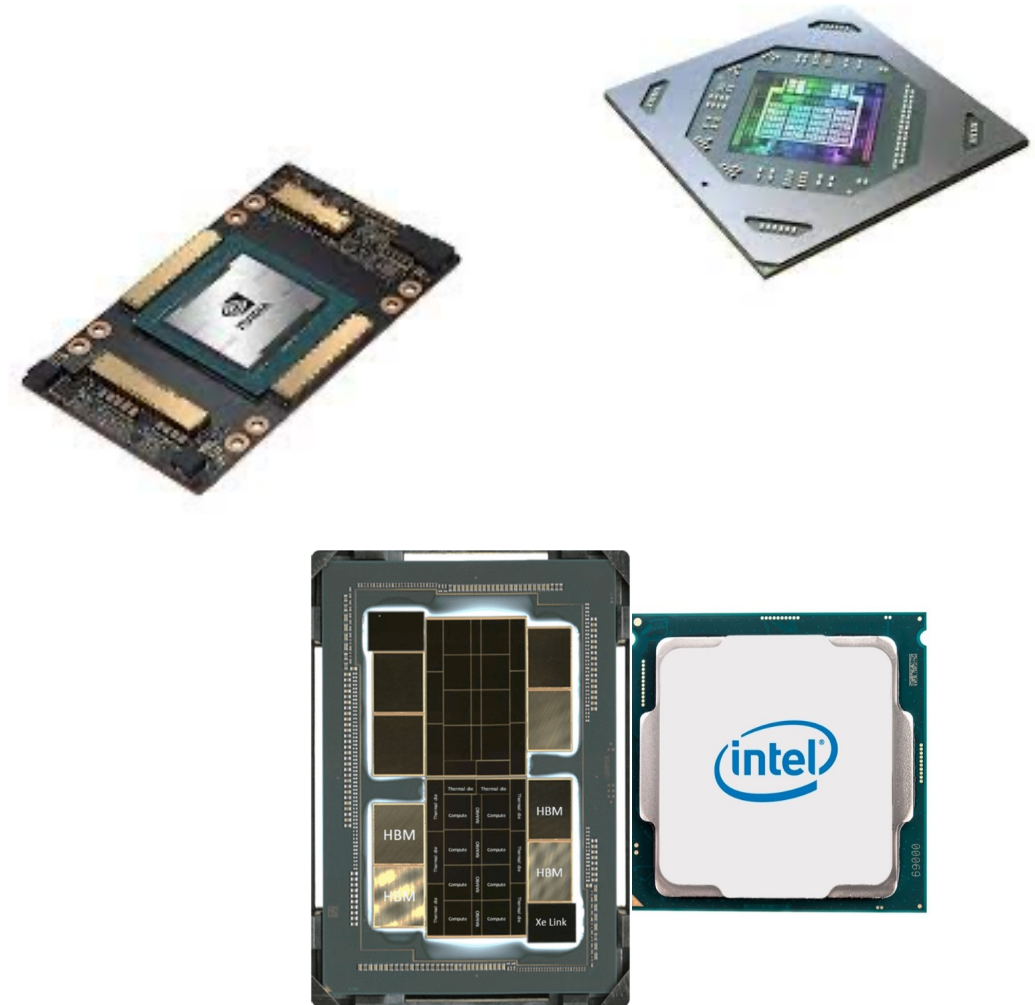
Why do we need to program for GPU?



GPU programming:

- Limited only to a specific domain
- Separate source solutions
- Verbose low Levels APIs

- **oneAPI & DPC++**
- **CUDA C/C++**
- **HIP**
- **SYCL**
- **OpenCL**
- **Kokkos**
- **HPX ...**



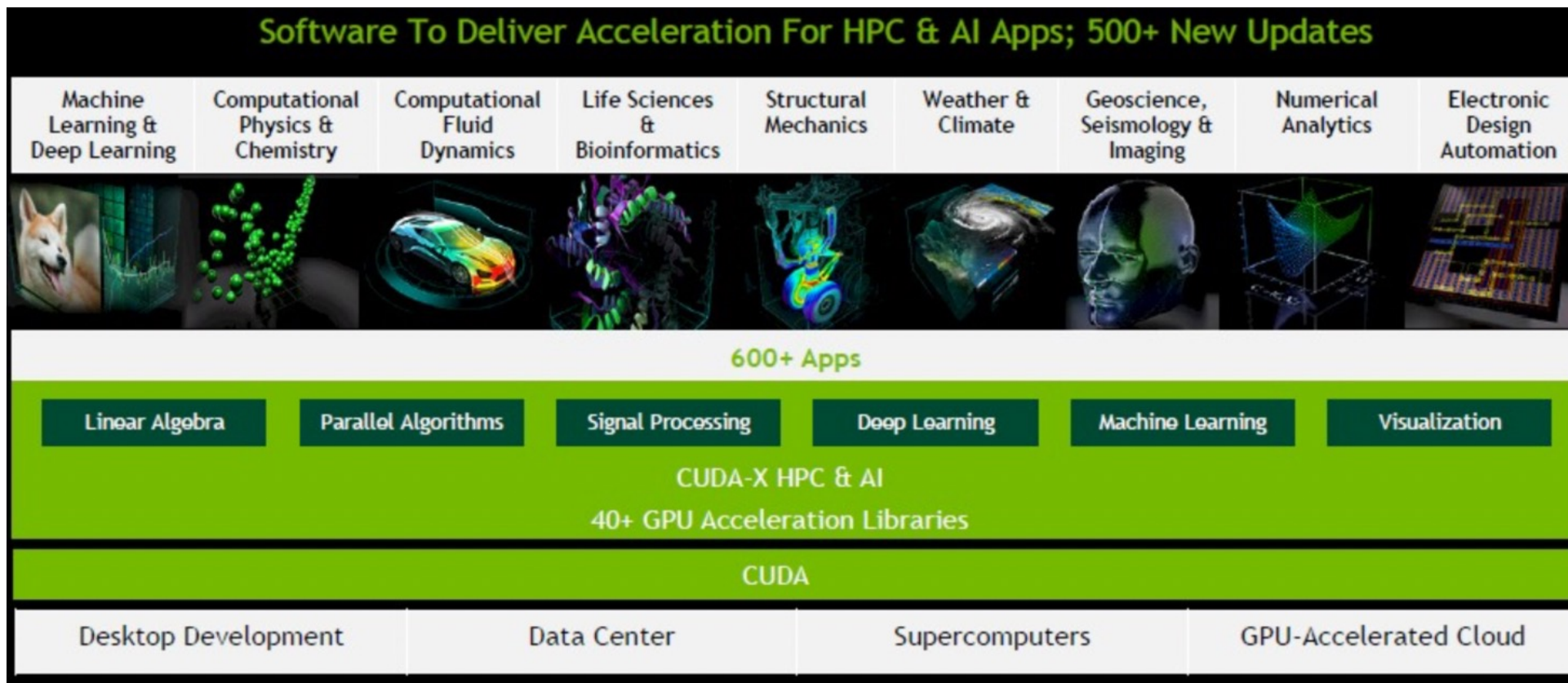
Why do we need GPUs on HPC?



DEEP
LEARNING
INSTITUTE



- Increase in parallelism
- Today almost a **similar amount of efforts** on using CPUs vs GPUs by real applications
- GPUs well-suited to deep learning.



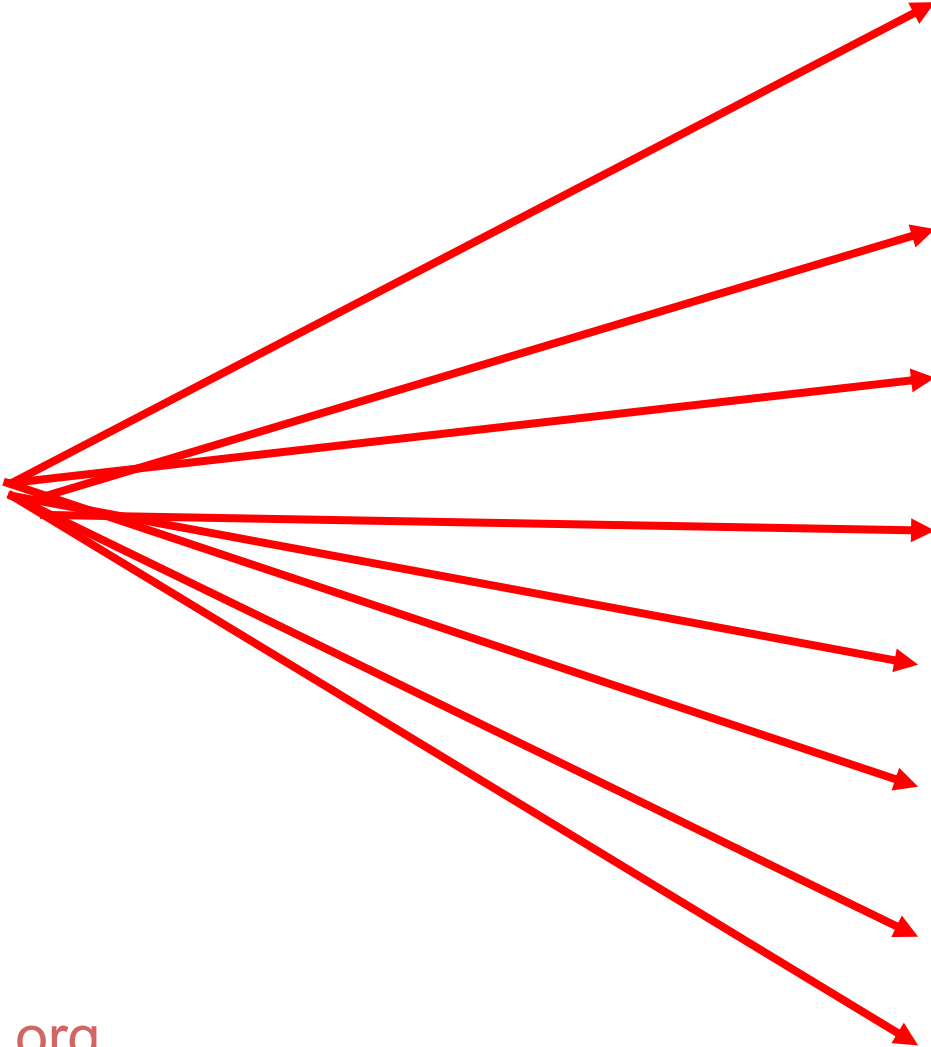
NVIDIA Software uses CUDA

Why do we need “GPU accelerators” on HPC?

GPU-accelerated systems



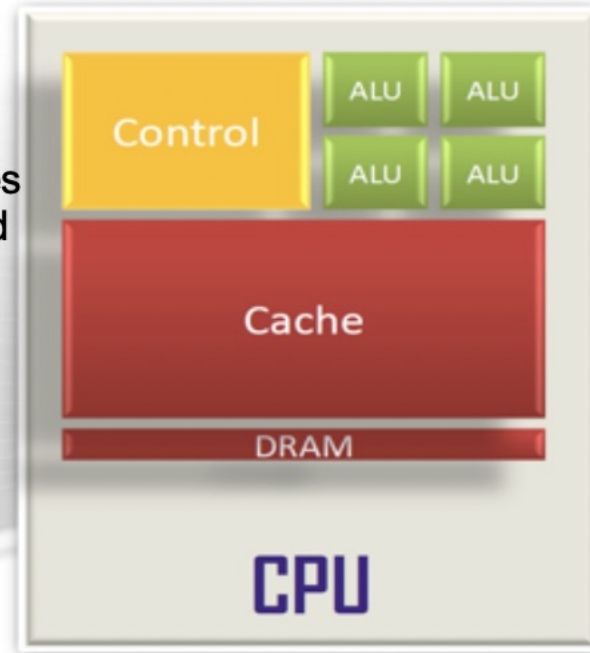
www.top500.org



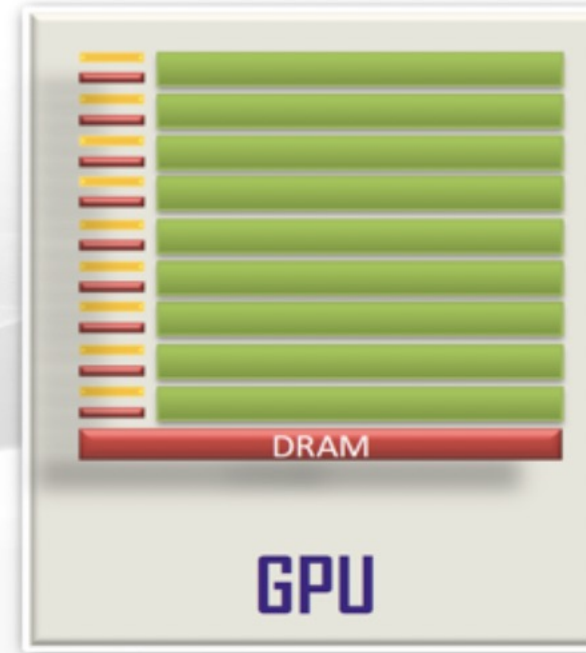
| Rank | System | Cores | Rmax (PFlop/s) | Rpeak (PFlop/s) | Power (kW) |
|------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------|----------------|-----------------|------------|
| 1 | Frontier - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, HPE DOE/SC/Oak Ridge National Laboratory United States | 8,699,904 | 1,194.00 | 1,679.82 | 22,703 |
| 2 | Supercomputer Fugaku - Supercomputer Fugaku, A64FX 48C 2.2GHz, Tofu interconnect D, Fujitsu RIKEN Center for Computational Science Japan | 7,630,848 | 442.01 | 537.21 | 29,899 |
| 3 | LUMI - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, HPE EuroHPC/CSC Finland | 2,220,288 | 309.10 | 428.70 | 6,016 |
| 4 | Leonardo - BullSequana XH2000, Xeon Platinum 8358 32C 2.6GHz, NVIDIA A100 SXM4 64 GB, Quad-rail NVIDIA HDR100 Infiniband, Atos EuroHPC/CINECA Italy | 1,824,768 | 238.70 | 304.47 | 7,404 |
| 5 | Summit - IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband, IBM DOE/SC/Oak Ridge National Laboratory United States | 2,414,592 | 148.60 | 200.79 | 10,096 |
| 6 | Sierra - IBM Power System AC922, IBM POWER9 22C 3.1GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband, IBM / NVIDIA / Mellanox DOE/NNSA/LLNL United States | 1,572,480 | 94.64 | 125.71 | 7,438 |
| 7 | Sunway TaihuLight - Sunway MPP, Sunway SW26010 260C 1.45GHz, Sunway, NRCC National Supercomputing Center in Wuxi China | 10,649,600 | 93.01 | 125.44 | 15,371 |
| 8 | Perlmutter - HPE Cray EX235n, AMD EPYC 7763 64C 2.45GHz, NVIDIA A100 SXM4 40 GB, Slingshot-10, HPE DOE/SC/LBNL/NERSC United States | 761,856 | 70.87 | 93.75 | 2,589 |
| 9 | Selene - NVIDIA DGX A100, AMD EPYC 7742 64C 2.25GHz, NVIDIA A100, Mellanox HDR Infiniband, Nvidia NVIDIA Corporation United States | 555,520 | 63.46 | 79.22 | 2,646 |
| 10 | Tianhe-2A - TH-IVB-FEP Cluster, Intel Xeon E5-2692v2 12C 2.2GHz, TH Express-2, Matrix-2000, NUDT National Super Computer Center in Guangzhou China | 4,981,760 | 61.44 | 100.68 | 18,482 |

GPU vs CPU Architecture

- * Small number of large cores
- * More control structures and less processing units
- * Optimised for latency which requires quite a lot of power



General purpose architecture



Massively data parallel

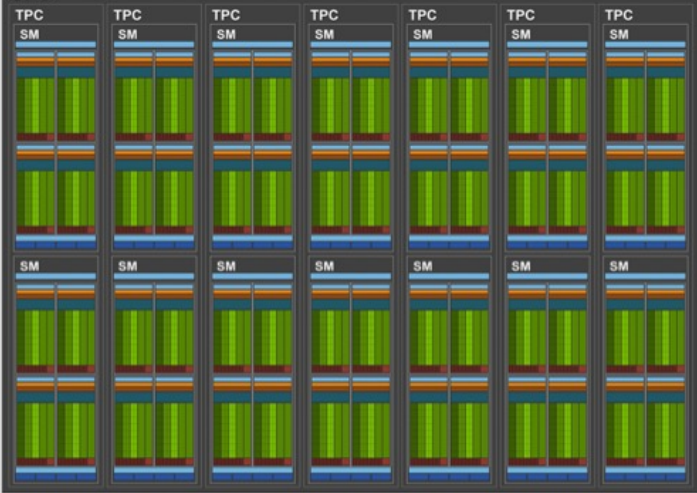
- * Large number of small cores
- * Less control structured and more processing units
- * Less flexible program model
- * There're more restrictions but Requires a lot less power

• GPU devotes more transistors data processing rather than data caching and flow control. Same problem executed on many data elements in parallel.

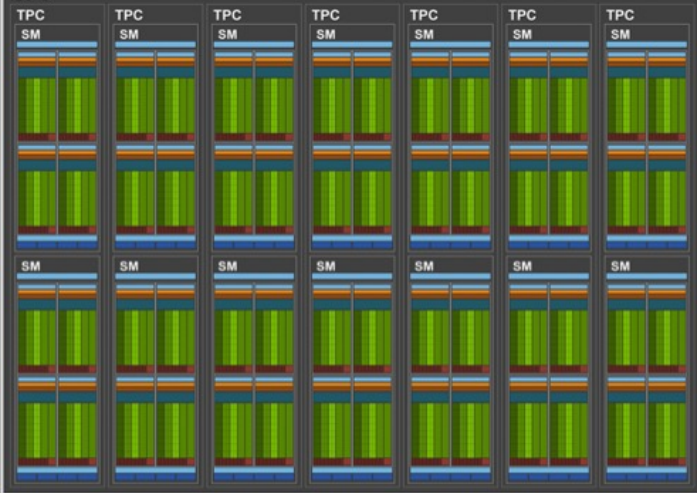
PCI Express 3.0 Host Interface

GigaThread Engine

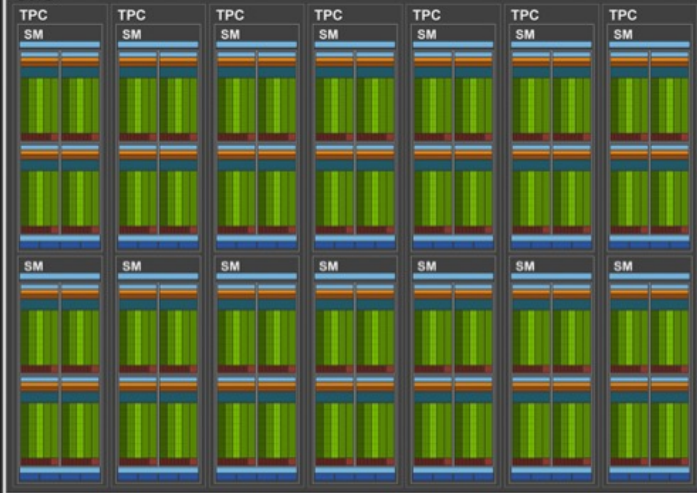
GPC



GPC



GPC



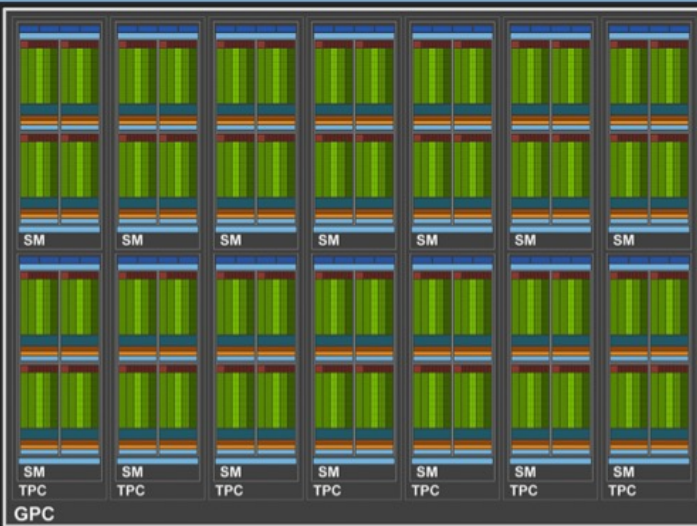
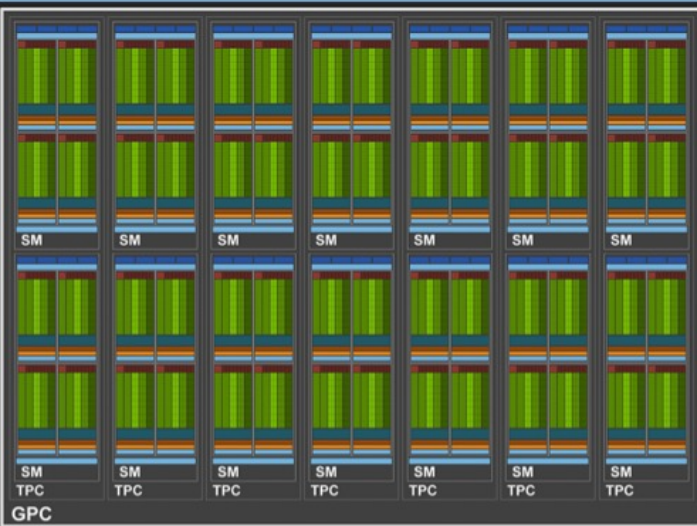
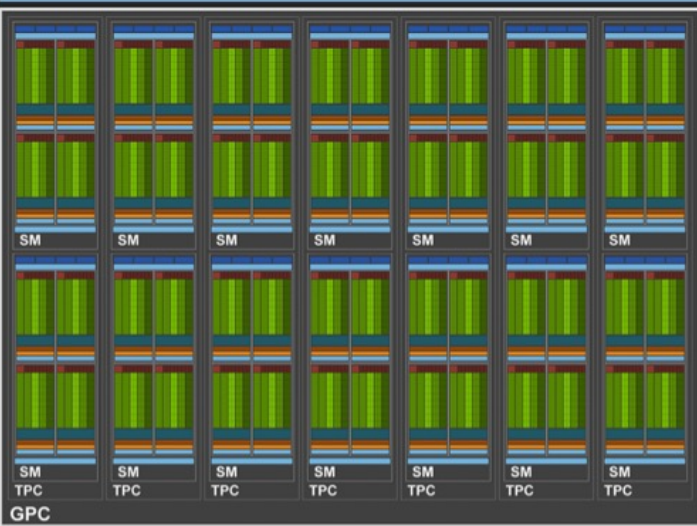
Memory Controller

Memory Controller

Memory Controller

Memory Controller

L2 Cache



High-Speed Hub

NVLink

NVLink

NVLink

NVLink

NVLink

NVLink

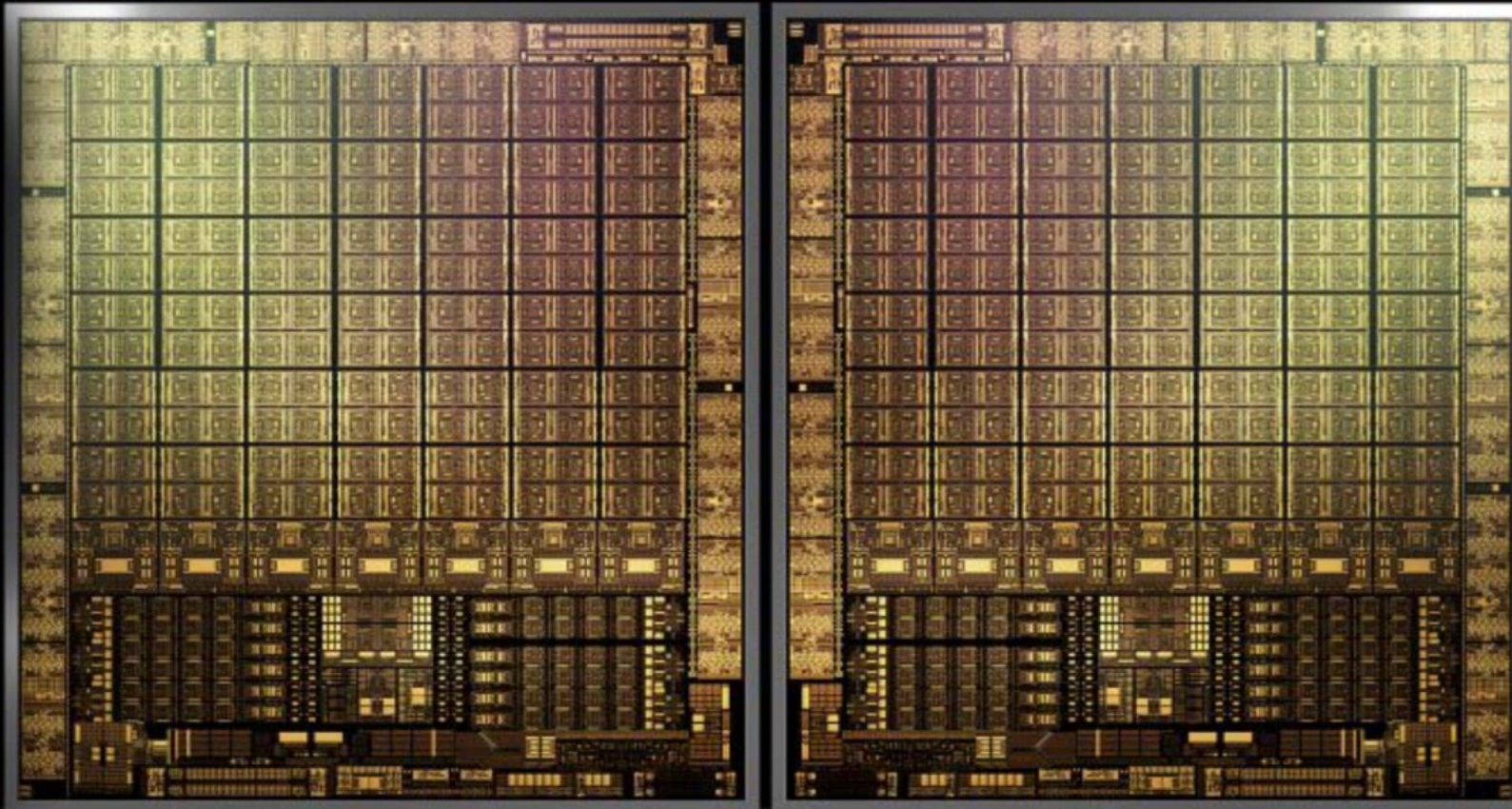
HBM2

HBM2

HBM2

HBM2





- Hopper GPU (H100) with over 80 Billion Transistors on an 814 mm²
- 89 GB memory
- First support PCIe gen5 and utilize the HBM3 enabling 3TB/s
- 30 Tflops of peak FP64, 60Tflops with FP64 tensor-core or 32 FP performance

What and Why CUDA C/C++ ?



CUDA = "Compute Unified Device Architecture"

* Introduced and released in 2006 for the GeForce 8800 *

- GPU = massively data parallel - co-processor

C/C++ plus a few simple extensions

- Compute oriented drivers, language, and tools

Documentations:

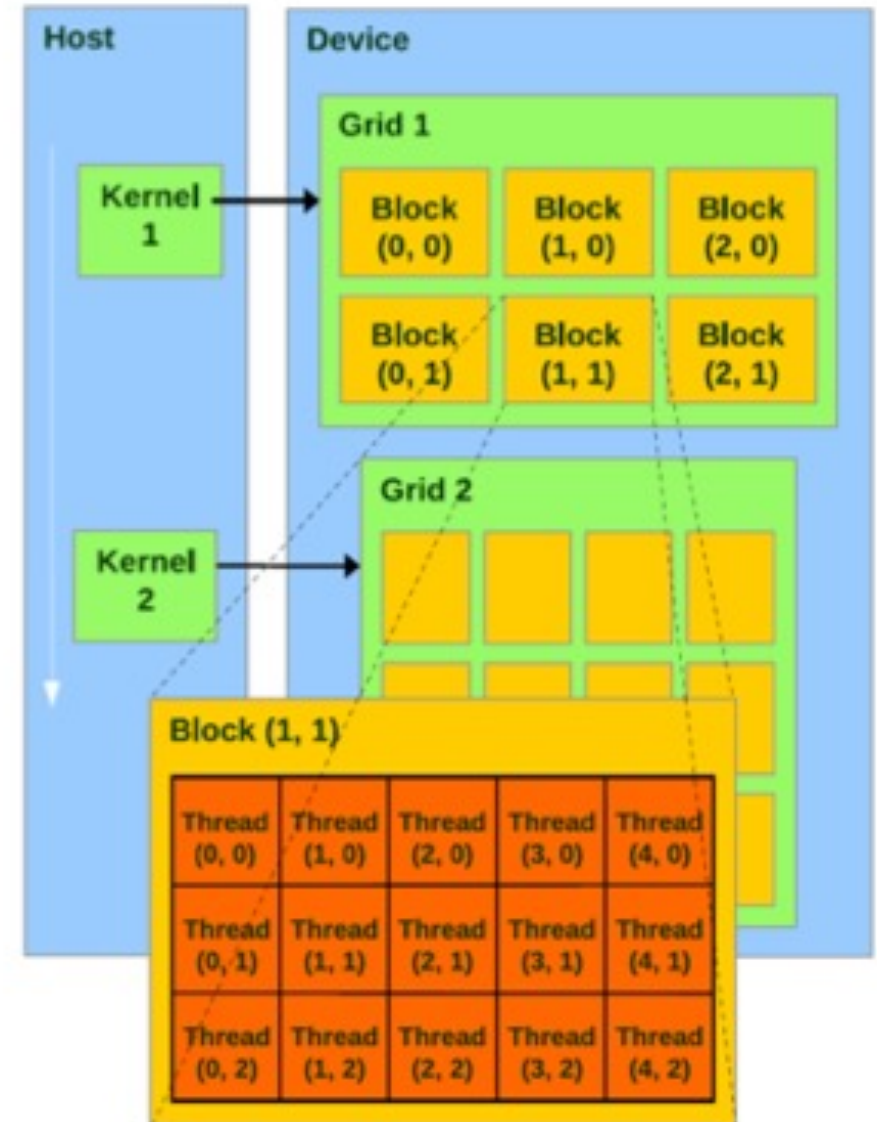
CUDA_C_Programming_Guide.pdf

CUDA_C_Getting_Started.pdf

CUDA_C_Toolkit_Release.pdf

CUDA Programming Model

- A kernel is executed as a grid of thread blocks
- All threads share data memory space
- A thread block is a batch of threads that can cooperate with each other by:
 - Synchronizing their execution
 - Efficiently sharing data through a low latency shared memory
- Two threads from two different blocks cannot cooperate
- Sequential code launches **asynchronously** GPU kernels



Terminology:

Host: The CPU and its memory
(host memory)



Host

Device: The GPU and its memory
(device memory)

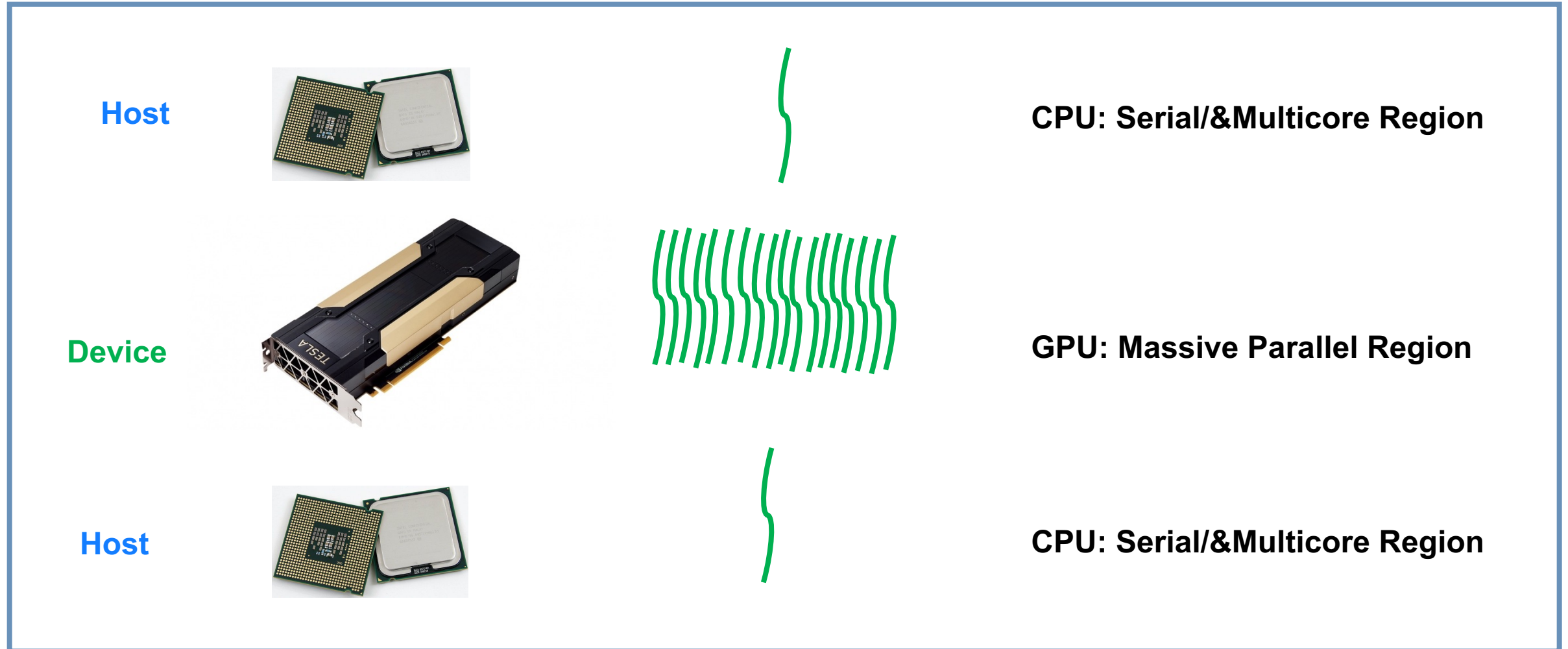


Device

CUDA Devices and Threads Execution Model



DEEP
LEARNING
INSTITUTE



The CPU allocates memory on the GPU
The CPU copies data from CPU to GPU
The CPU launches kernels on the GPU
The CPU copies data to CPU from GPU



NVCC Compiler



- NVIDIA provides a CUDA-C compiler

→ **nvcc**

- NVCC splits your code in 2: **Host** code and **Device** code.
- **Device** code sent to NVIDIA device compiler.

- **nvcc** is capable of linking together both host and device code into a single executable.

- **Convention:** C++ source files containing CUDA syntax are typically given the extension **.cu**.

- For „**.cpp**“ extension use:
`nvcc -x cu -arch=sm_70 -o exe code.cpp`



DEEP
LEARNING
INSTITUTE



Lab1: Accelerating Applications with CUDA C/C++

Dr. Momme Allalen

Leibniz Computing Centre, Munich Germany - www.lrz.de

Deep Learning Certified Instructor, NVIDIA Deep Learning Institute NVIDIA Corporation.

Lab1: Accelerating Applications with CUDA C/C++



Prerequisites

You should already be able to:

- Declare variables, write loops, and use if / else statements in C.
- Define and invoke functions in C.
- Allocate arrays in C.
- No previous CUDA knowledge is required.

Objectives

By the time you complete this lab, you will be able to:

- Write, compile, and run C/C++ programs that both call **CPU functions** and **launch GPU kernels**.
- Control parallel **threadhierarchy** using **execution configuration**.
- Refactor serial loops to execute their iterations in parallel on a **GPU**.
- Allocate and free memory available to both **CPUs** and **GPUs**.
- Handle errors generated by CUDA code.
 - Accelerate **CPU-only applications**.

nvc; nvc++ Compiler



nvc :is a C11 compiler for NVIDIA GPUs and AMD, Intel, OpenPOWER, and Arm CPUs. It invokes the C compiler, assembler, and linker for the target processors with options derived from its command line arguments. **nvc** supports ISO C11, supports GPU programming with OpenACC, and supports multicore CPU programming with OpenACC and OpenMP.

nvc++ : is a C++17 compiler for NVIDIA GPUs and AMD, Intel, OpenPOWER, and Arm CPUs. It invokes the C++ compiler, assembler, and linker for the target processors with options derived from its command line arguments. **nvc++** supports ISO C++17, supports GPU and multicore CPU programming with C++17 parallel algorithms, OpenACC, and OpenMP.

nvfortran : is a Fortran compiler for NVIDIA GPUs and AMD, Intel, OpenPOWER, and Arm CPUs. It invokes the Fortran compiler, assembler, and linker for the target processors with options derived from its command line arguments. **nvfortran** supports ISO Fortran 2003 and many features of ISO Fortran 2008, supports GPU programming with CUDA Fortran, and GPU and multicore CPU programming with ISO Fortran parallel language features, OpenACC, and OpenMP.

nvcc : is the CUDA C and CUDA C++ compiler driver for NVIDIA GPUs. **nvcc** accepts a range of conventional compiler options, such as for defining macros and include/library paths, and for steering the compilation process. **nvcc** produces optimized code for NVIDIA GPUs and drives a supported host compiler for AMD, Intel, OpenPOWER, and Arm CPUs.

Lab2: Managing Accelerated Application Memory with CUDA Unified Memory and nsys

Dr. Momme Allalen

Leibniz Computing Centre, Munich Germany - www.lrz.de

Deep Learning Certified Instructor, NVIDIA Deep Learning Institute NVIDIA Corporation.

Lab2: Managing Accelerated Application Memory with CUDA Unified Memory and nsys



Prerequisites

You should already be able to:

- Write, compile, and run C/C++ programs that both call CPU functions and **launch GPU kernels**.
- Control parallel **thread hierarchy** using **execution configuration**.
- Refactor serial loops to execute their iterations in parallel on a GPU.
- Allocate and free Unified Memory.

Objectives

- By the time you complete this lab, you will be able to:
- Use the Nsight Systems command line tool (**nsys**) to profile accelerated application performance.
 - Leverage and understanding of **Streaming Multiprocessors** to optimize execution configurations.
 - Understand the behavior of **Unified Memory** with regard to page faulting and data migrations.
 - Use **asynchronous memory prefetching** to reduce page faults and data migrations for increased performance.
 - Employ an iterative development cycle to rapidly accelerate and deploy applications.

CUDA® PROFILING TOOLS



nvvc: NVIDIA visual profiler

nvprof: tool to understand and optimize the performance of your CUDA, OpenACC or OpenMP applications,
Application level opportunities

- Overall application performance

 - Overlap CPU and GPU work, identify the bottlenecks (CPU or GPU)

- Overall GPU utilization and efficiency

 - Overlap compute and memory copies

 - Utilize compute and copy engines effectively.

Kernel level opportunities

 - Use memory bandwidth efficiently

 - Use compute resources efficiently

 - Hide instruction and memory latency

There are more features, example for Dependency Analysis

Command: **nvprof** --dependency-analysis --cpu-thread-tracing on ./executable_cuda



Nsight Systems
Nsight Compute

THE NSIGHT SUITE COMPONENTS

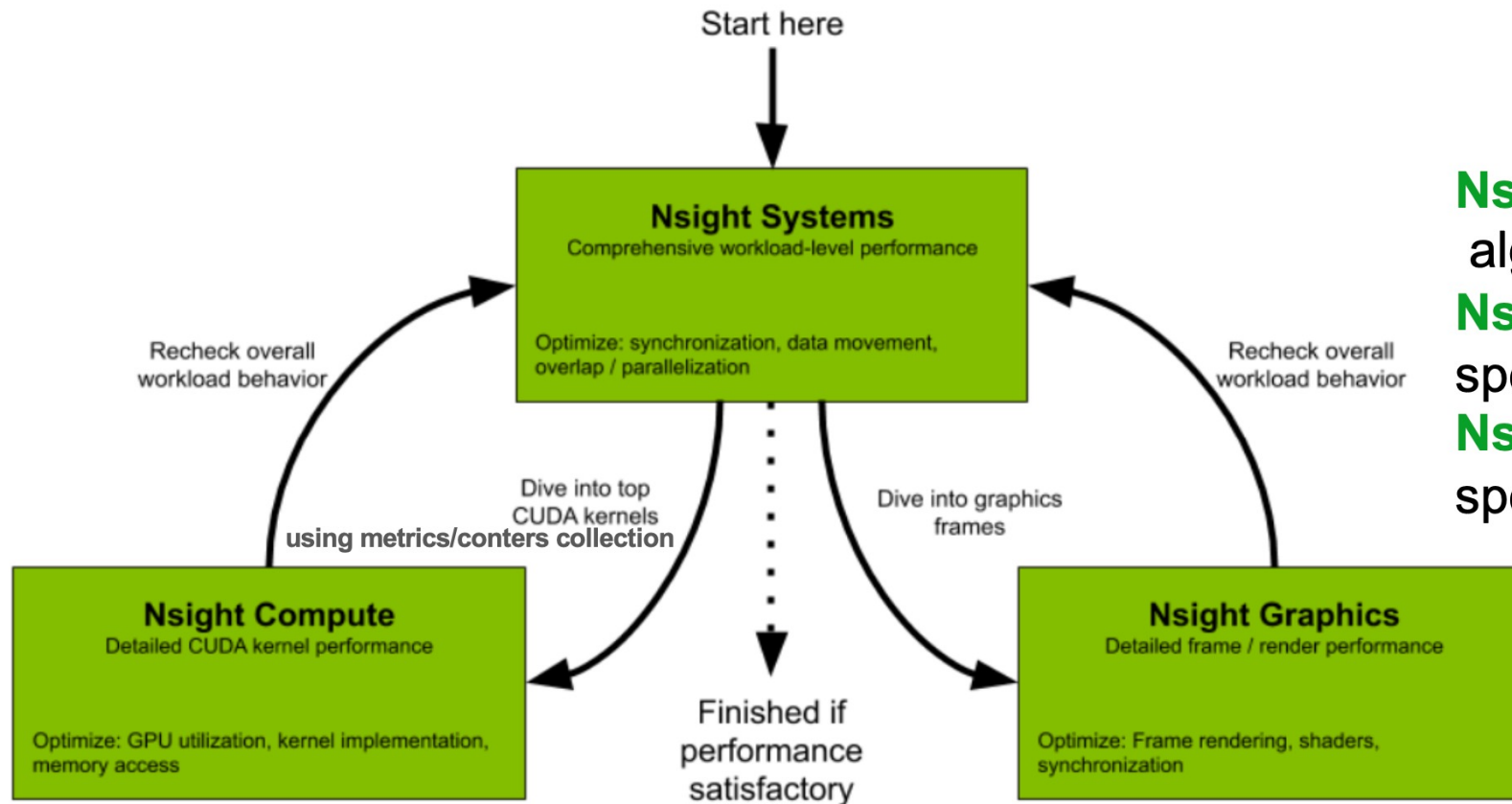


Figure 1. Flowchart describing working with new NVIDIA Nsight tools for performance optimization

Nsight Systems – Analyze application algorithm system wide

Nsight Compute – Debug/&optimize specific CUDA kernels

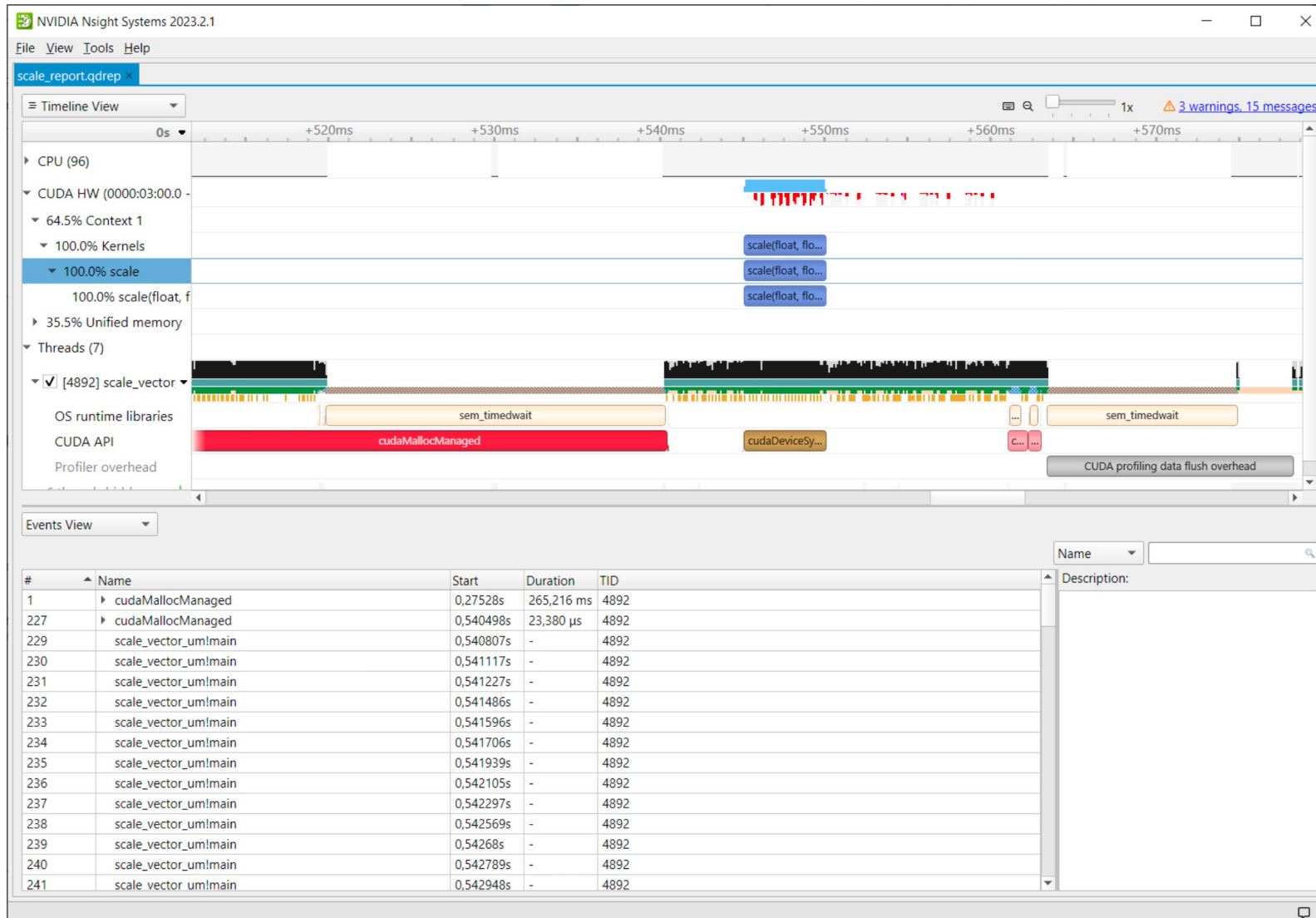
Nsight Graphics – Debug/&optimize specific graphics workloads

nvprof replaced with **nsys –profile....**

<https://developer.nvidia.com/nsight-systems>

Nsight Systems GUI

Main timeline view, Events View



- Nsight Systems is an extremely low overhead profiling Tool across any number of CPUs and GPUs.
- Nsight System is your first stop on your profiling workflow, inspect your algorithm timing and GPU interaction and identify a large number of opportunities for optimization.

NSIGHT SYSTEMS



- Provides users with a more complete view of how their codes balance workload across multiple CPUs and GPUs
- Locate optimization opportunities, helps and allows to identify issues such as:
 - *GPU starvation*
 - *Insufficient CPU parallelisation or pipelining*
 - *Unexpectedly expensive CPU or GPU algorithm*
 - *Unnecessary GPU synchronization*
- The tool uses low overhead tracing and sampling techniques to collect process and thread activity and visualize millions of events on a very fast GUI timeline
- Correlates that data across CPU cores and GPU streams, allowing users to investigate bottlenecks.
- Multi-platform: Linux & Windows, x86-64, Tegra, Power, MacOSX (host only)

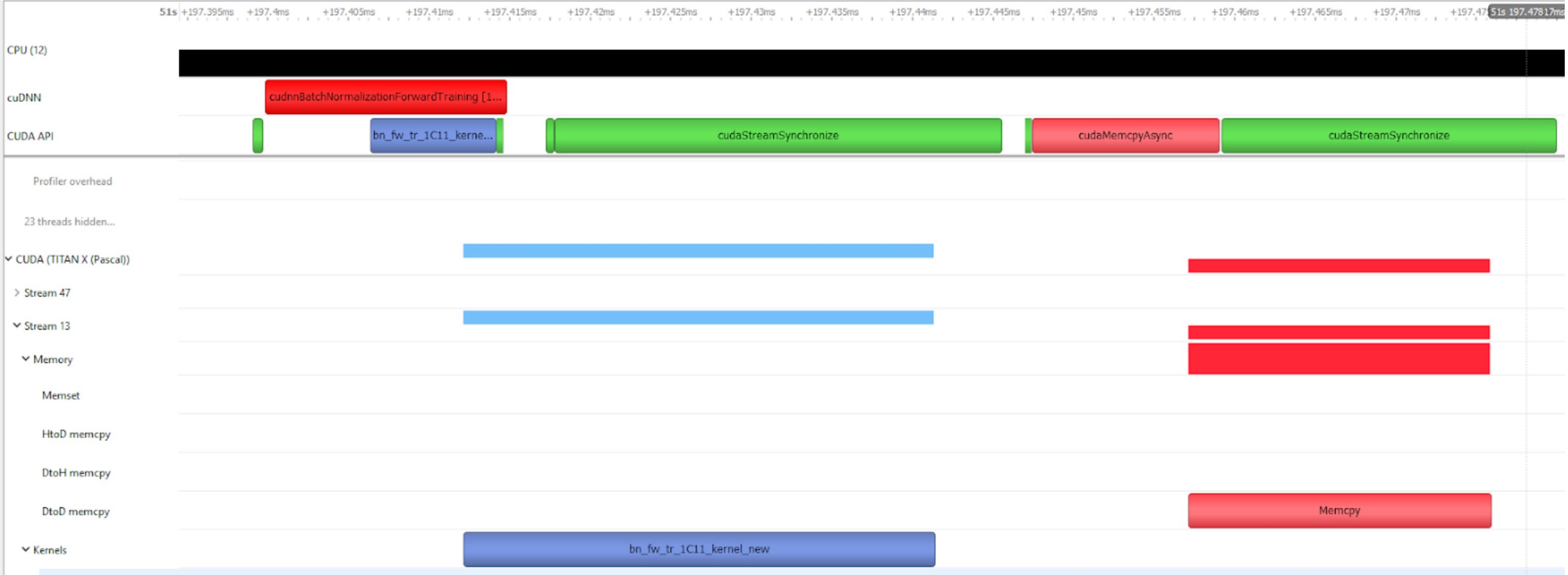
<https://developer.nvidia.com/nsight-systems>

GPU Starvation Investigations



<https://developer.nvidia.com/nsight-systems>

Unnecessary GPU Synchronisation Calls



<https://developer.nvidia.com/nsight-systems>

NVIDIA NSIGHT SYSTEMS



- Support: **MPI**, **OpenACC**, **OpenMP**
- Complex data mining capabilities, enables to go beyond basic statistics.
- Support multiple simultaneous sessions.
- **MPI trace** feature enables to analyse when the threads are busy or blocked in long-running functions of the **MPI** standard, available on **OpenMPI**, **MPICH** and **NVShmem**.
- **OpenACC** trace enables to see where code has been offload and parallelized onto the GPU, which helps you to analyse the activities executing on the CPUs and GPUs in parallel.
- Tracing **OpenMP** code is available for compilers supporting **OpenMP5** and **OMPT interface**. This capability enables tracing of the parallel regions of code that are distributed either across multiple threads or to the GPU.
- Provides support for **CUDA** graphs. To understand the execution of the source of **CUDA** kernels and execution of **CUDA** graphs, kernels can be correlated back through the graph launch, instantiation, and all the way back to the code creation, to identify the origin of the kernel execution on the GPU.

Command Line Options nsys



| Command | Description |
|----------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| profile | A fully formed profiling description requiring and accepting no further input. The command switch options used (see below table) determine when the collection starts, stops, what collectors are used (e.g. API trace, IP sampling, etc.), what processes are monitored, etc. |
| start | Start a collection in interactive mode. The start command can be executed before or after a launch command. |
| stop | Stop a collection that was started in interactive mode. When executed, all active collections stop, the CLI process terminates but the application continues running. |
| cancel | Cancels an existing collection started in interactive mode. All data already collected in the current collection is discarded. |
| launch | In interactive mode, launches an application in an environment that supports the requested options. The launch command can be executed before or after a start command. |
| shutdown | Disconnects the CLI process from the launched application and forces the CLI process to exit. If a collection is pending or active, it is cancelled |
| export | Generates an export file from an existing .nsys-rep file. For more information about the exported formats see the /documentation/nsys-exporter directory in your Nsight Systems installation directory. |
| stats | Post process existing Nsight Systems result, either in .nsys-rep or SQLite format, to generate statistical information. |
| analyze | Post process existing Nsight Systems result, either in .nsys-rep or SQLite format, to generate expert systems report. |
| status | Reports on the status of a CLI-based collection or the suitability of the profiling environment. |
| sessions | Gives information about all sessions running on the system. |

<https://docs.nvidia.com/nsight-systems/UserGuide/index.html>

NSIGHT COMPUTE (ncu)



Interactive CUDA Kernel profiler

- Targeted metric sections for various performance aspects (Debug/Profile)
- API debugging via a user interface command line tool
- Very high freq. GPU perf counter, customizable data collection and presentation (tables, charts ..)
- Python-based rules for guided analysis (or postprocessing)
- Provides a customizable and data-driven user interface and metric collection and can be extended with analysis scripts for post-processing results.

<https://docs.nvidia.com/nsight-compute/NsightCompute/index.html>

NVIDIA NSIGHT COMPUTE

Important Features



DEEP
LEARNING
INSTITUTE



- Result comparison across one or multiple reports within the tool
- Graphical profile report
- Interactive kernel profiler and API debugger: debugging CPU and GPU simultaneously and capable of handling thousands of simultaneous threads.
- Fast data collection
- GUI and command line interface
- Fully customizable reports and analysis rules

<https://developer.nvidia.com/nsight-systems>

Nsight Compute Feature Spotlight in CUDA Toolkit 11 and A100

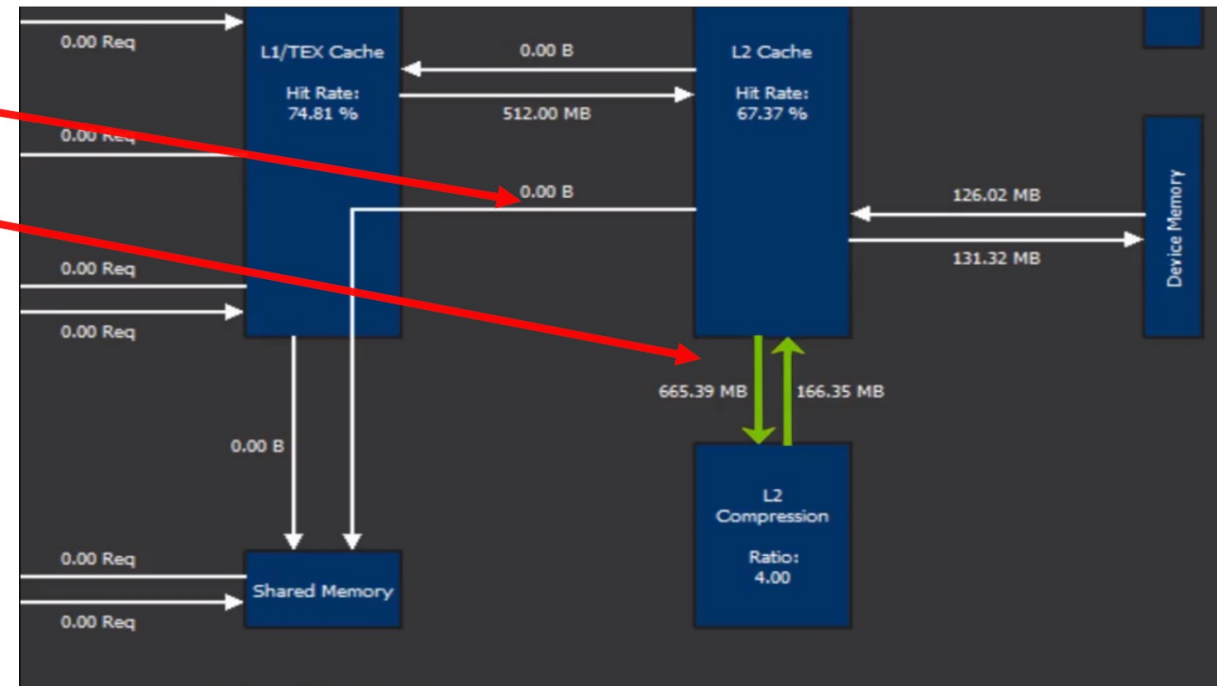
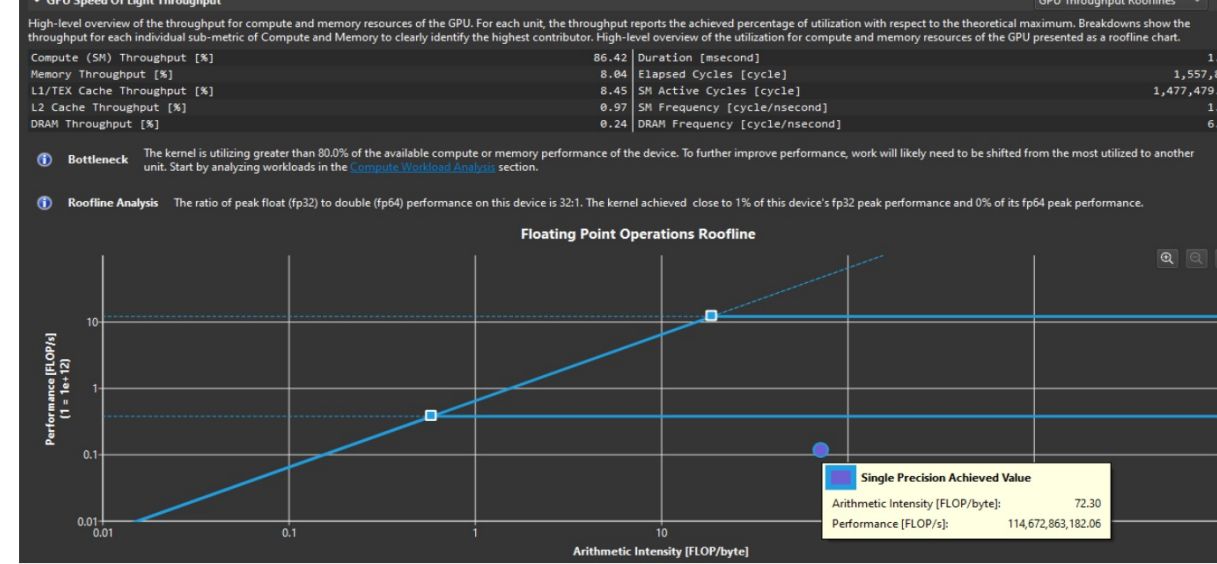
- **Roofline Analysis**

Arithmetic intensity = Compute/Memory
FLOPS = Floating Points Ops/Second

- **Asynchronous copy**

- **Sparse Data Compression**

Shows the amount of data compressed through this feature and the compression ratio, helps on kernels with bandwidth or cache issues.



Docs/product: <https://developer.nvidia.com/nsight-compute>

NVIDIA® Tools Extension SDK (NVTX)



- C-based Application Programming Interface (API) for annotating events, code ranges, and resources in your applications
- Codes which integrate NVTX can use NVIDIA Nsight, Tegra System Profiler, and Visual Profiler to capture and visualize these events and ranges.

```
[allalen1@jwlogin22 v2]$ ncu -h | grep nvtx
--nvtx                Enable NVTX support.
--nvtx-include arg    Adds include statement to the NVTX filter, which allows selecting kernels to
--nvtx-exclude arg    Adds exclude statement to the NVTX filter, which allows selecting kernels to
--print-nvtx-rename arg (=none) Select how NVTX should be used for renaming:
                        per-nvtx
Usage of --nvtx-include and --nvtx-exclude:
ncu --nvtx --nvtx-include "Domain A@Range A"
ncu --nvtx --nvtx-exclude "Range A]"
ncu --nvtx --nvtx-include "Range A" --nvtx-exclude "Range B"
```

<https://docs.nvidia.com/nsight-visual-studio-edition/nvtx/index.html>

NVIDIA® Tools Extension SDK (NVTX)



```
#include <nvToolsExt.h>
#include <sys/syscall.h>
#include <unistd.h>
```

```
static void wait(int seconds) {
    nvtxRangePush(__FUNCTION__);
    nvtxMark("Waiting...");
    sleep(seconds);
    nvtxRangePop();
}
```

```
int main(void) {
    nvtxNameOsThread(syscall(SYS_gettid), "Main Thread");
    nvtxRangePush(__FUNCTION__);
    wait(1);
    nvtxRangePop();
}
```

nsys profile -t nvtx --stats=true ...

Or for Julia code:

**nsys profile -t nvtx,cuda -o output_file.qdrep
julia --project=../.. script.jl**

<https://docs.nvidia.com/nsight-visual-studio-edition/2020.1/nvtx/index.html>

A First (I)Nsight

Recording with the CLI

- Use the command line
 - `srun nsys profile --trace=cuda,nvtx,mpi --force-override=true --output=my_report.%q{SLURM_PROCID} \ ./jacobi -niter 10`
- Inspect results: Open the report file in the GUI
 - Also possible to get details on command line
 - Either add `--stats` to profile command line, or: `nsys stats --help`
- Runs set of reports on command line, customizable (**sqlite** + **Python**):
 - Useful to check validity of profile, identify important kernels

Running [.../reports/**gpukernsum.py** jacobi_metrics_more-nvtx.0.**sqlite**]...

| Time(%) | Total Time (ns) | Instances | Avg (ns) | Med (ns) | Min (ns) | Max (ns) | StdDev (ns) | Name |
|---------|-----------------|-----------|-----------|-----------|----------|----------|-------------|-----------------------|
| 99.9 | 36750359 | 20 | 1837518.0 | 1838466.5 | 622945 | 3055044 | 1245121.7 | void jacobi_kernel |
| 0.1 | 22816 | 2 | 11408.0 | 11408.0 | 7520 | 15296 | 5498.5 | initialize_boundaries |

System-level Profiling with Nsight Systems

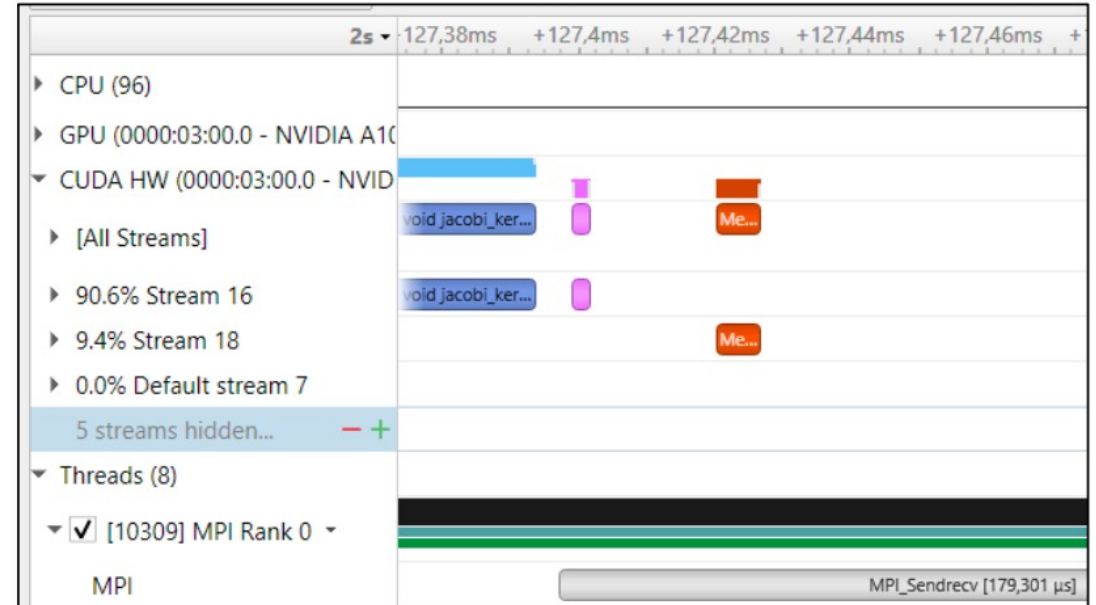
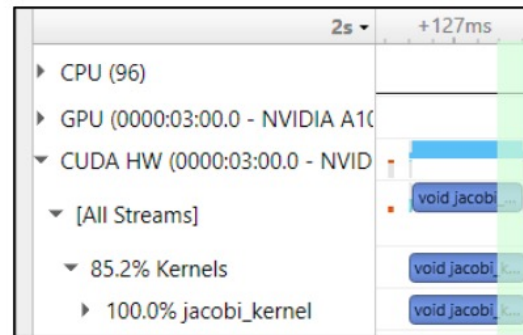
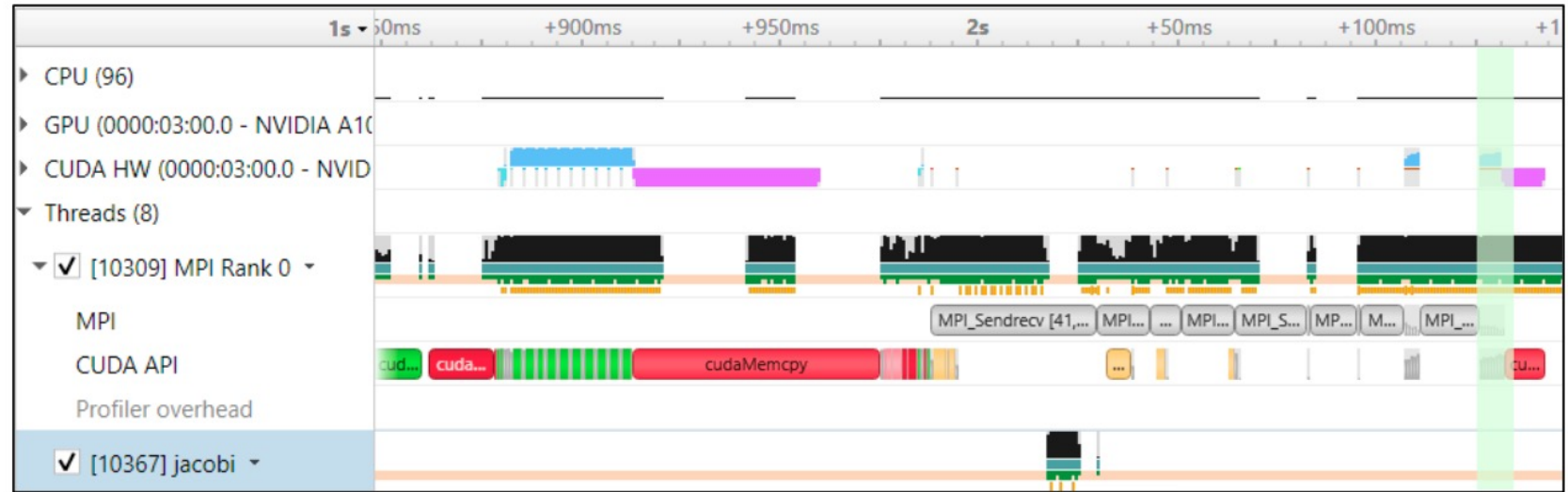
- Global timeline view
 - CUDA HW: streams, kernels, memory
- Different traces, e.g. CUDA, MPI
 - correlations API <-> HW
- Stack samples
 - bottom-up, top-down for CPU code
- GPU metrics
- Events View
 - Expert Systems
- looks at single process (tree)
 - correlate multi-process reports in single timeline

The screenshot displays the NVIDIA Nsight Systems 2021.4.1 interface. The main window shows a 'Timeline View' of system activity. The left sidebar contains a tree view of processes and threads, with 'MPI' and 'CUDA API' circled in red. The main timeline shows various GPU streams and MPI ranks. A vertical blue line marks a specific time point. Below the timeline, the 'Events View' is visible, showing a table of events with columns for #, Name, Start, Duration, TID, GPU, Context, and Description. The event at index 5 is highlighted in blue.

| # | Name | Start | Duration | TID | GPU | Context | Description |
|---|-------------------|----------|----------|-----|-------|-----------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| 4 | Memset | 1,88258s | 3,200 μs | - | GPU 0 | Stream 13 | |
| 5 | void jacobi_ke... | 1,88259s | 3,056 ms | - | GPU 0 | Stream 13 | void jacobi_kernel<(int)32, (int)32>(float *, const float *, float *, int, int, int, bool) Begins: 1,88259s Ends: 1,88565s (+3,056 ms) grid: <<<512, 512, 1>>> |
| 6 | Memcpy DtoD | 1,88565s | 5,024 μs | - | GPU 0 | Stream 14 | |
| 7 | Memcpy DtoH | 1,88565s | 4,864 μs | - | GPU 0 | Stream 13 | |

Discovering Optimization Potential

- Using Jacobi solver example*
- Spot kernels – lots of whitespace
 - Which part is „bad“?
 - Enhance!
- MPI calls
 - Memory copies
 - We know: This is CUDA-aware MPI
- Even without knowing source, insight
- Too complicated for repeated/reliable usage
 - How to simplify navigating and comparing reports?



*See <https://github.com/NVIDIA/multi-gpu-programming-models/>

Adding NVTX

Simple range-based API

- `#include "nvtx3/nvToolsExt.h"`
 - NVTX v3 is header-only, needs just `-ldl`
 - C++ and Python APIs
- Fortran: [NVHPC compilers include module](#)
 - Just use `nvtx` and `-lnvhpcwrapnvtx`
 - Other compilers: See blog posts linked below
- Definitely: Include `PUSH/POP` macros (see links below)
`PUSH_RANGE(name, color_idx)`
- Sprinkle them strategically through code
 - Use hierarchically: Nest ranges
- Not shown: Advanced usage (domains, ...)
- Similar range-based annotations exist for other tools
 - e.g. [SCOREP_USER_REGION_BEGIN](#)

```
int main(int argc, char** argv) {
    PUSH_RANGE("main", 0)
    PUSH_RANGE("init", 1)
    do_initialization();
    POP_RANGE
    /* ... */
    PUSH_RANGE("computation", 2)
    jacobi_kernel<<< /* ... */, compute_stream>>>(...);
    cudaStreamSynchronize(compute_stream);
    POP_RANGE
    /* ... */
    POP_RANGE
}
```

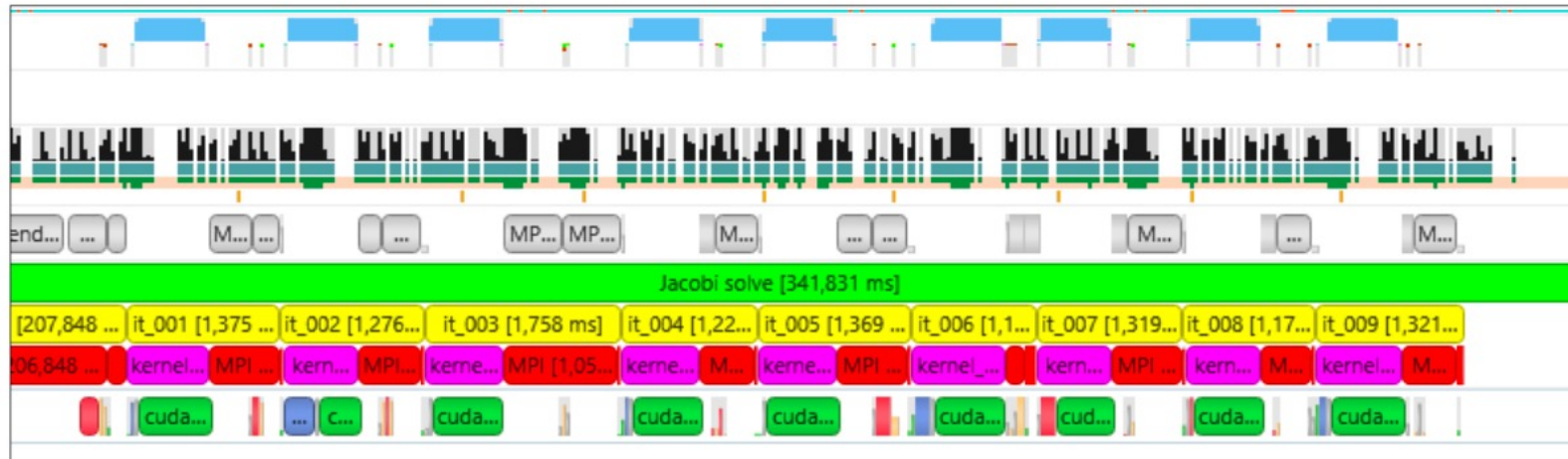
<https://github.com/NVIDIA/NVTX> and <https://nvidia.github.io/NVTX/#how-do-i-use-nvtx-in-my-code>

<https://developer.nvidia.com/blog/cuda-pro-tip-generate-custom-application-profile-timelines-nvtx/>
<https://developer.nvidia.com/blog/customize-cuda-fortran-profiling-nvtx/>

Minimizing Profile Size

Shorter time, smaller files = quicker progress

- Only profile what you need – all profilers have some overhead
 - Example: Event that occurs after long-running setup phase
- Bonus: lower number of events leads to smaller file size
- Add to nsys command line:
 - `--capture-range=nvtx --nvtx-capture=any_nvtx_marker_name \`
`--env-var=NSYS_NVTX_PROFILER_REGISTER_ONLY=0 --kill none`
 - Use [NVTX registered strings](#) for best performance
- Alternatively: `cudaProfilerStart()` and `-Stop()`
 - `--capture-range=cudaProfilerApi`

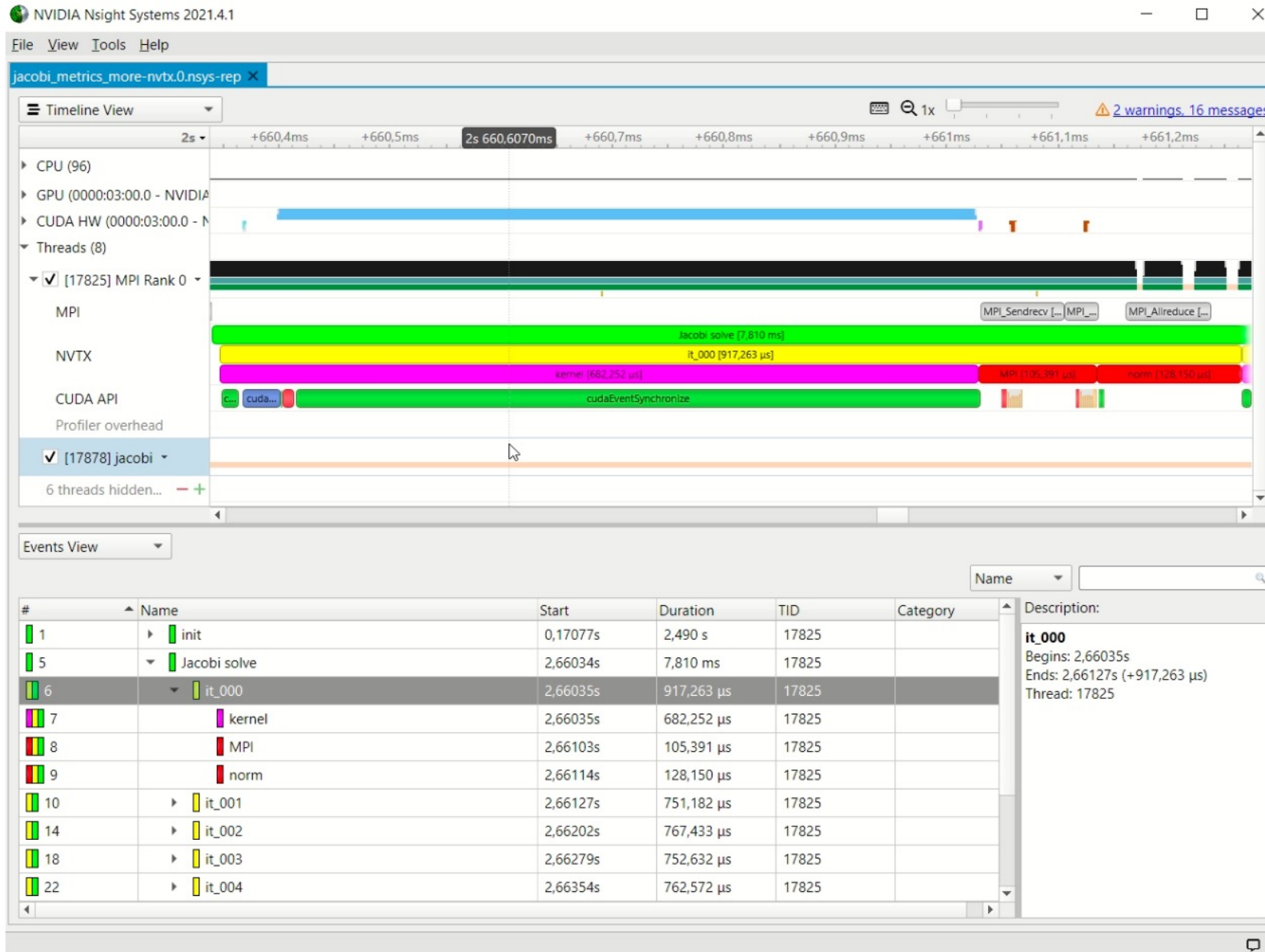


Nsight Systems Workflow with NVTX

Repeating the analysis



DEEP
LEARNING
INSTITUTE



GPU Metrics in Nsight Systems

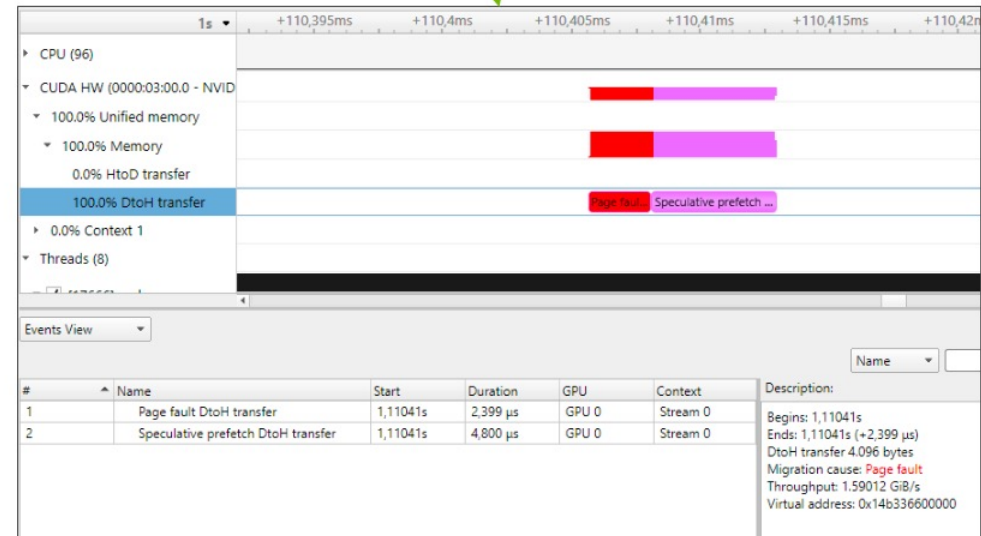
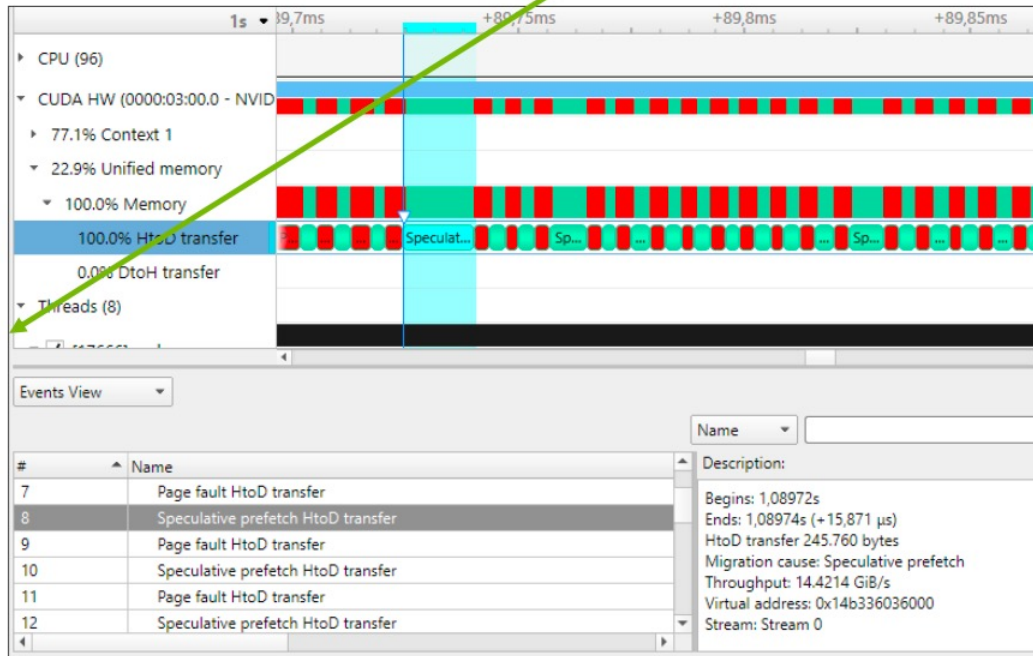
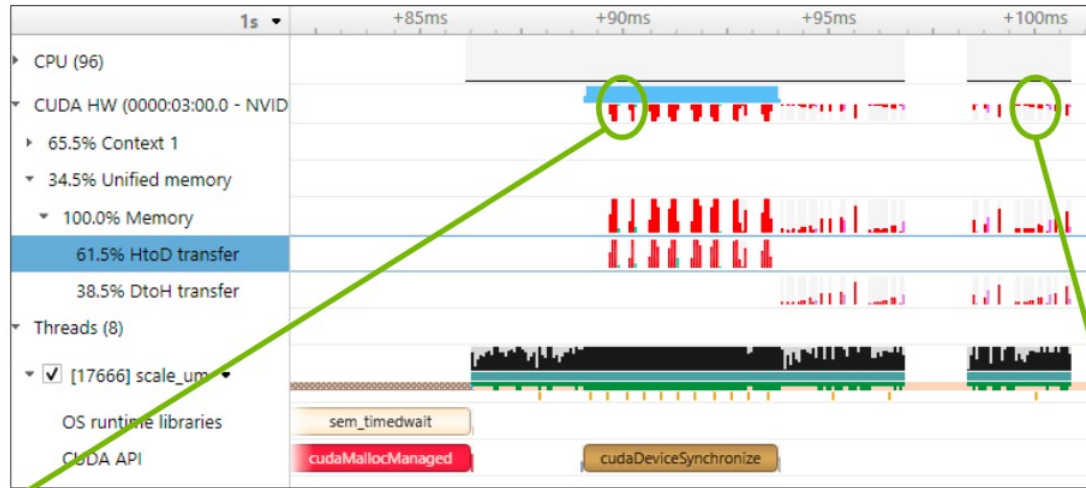
...and other traces you can activate

- Valuable low-overhead insight into HW usage:
 - SM instructions
 - DRAM Bandwidth, PCIe Bandwidth (GPUDirect)
- Also: Memory usage, Page Faults (higher overhead)
 - CUDA Programming guide: [Unified Memory Programming](#)
- Can save kernel-level profiling effort!
- `nsys profile`
 - `--gpu-metrics-device=0`
 - `--cuda-memory-usage=true`
 - `--cuda-um-cpu-page-faults=true`
 - `--cuda-um-gpu-page-faults=true`
 - `./app`



Unified Memory movement

Observing transfers in Nsight Systems



NSIGHT SYSTEMS



DEEP
LEARNING
INSTITUTE



System-wide application algorithm tuning

Multi-process tree support

Locate optimization opportunities

Visualize millions of events on a very fast GUI timeline

Or gaps of unused CPU and GPU time

Balance your workload across multiple CPUs and GPUs

CPU algorithms, utilization, and thread state

GPU streams, kernels, memory transfers, etc

Multi-platform: Linux & Windows, x86-64, Te g r a, Power, MacOSX(host only)

GPUs: Volta, Turing

Docs/product: <https://developer.nvidia.com/nsight-systems>

CUDA Kernel profiler

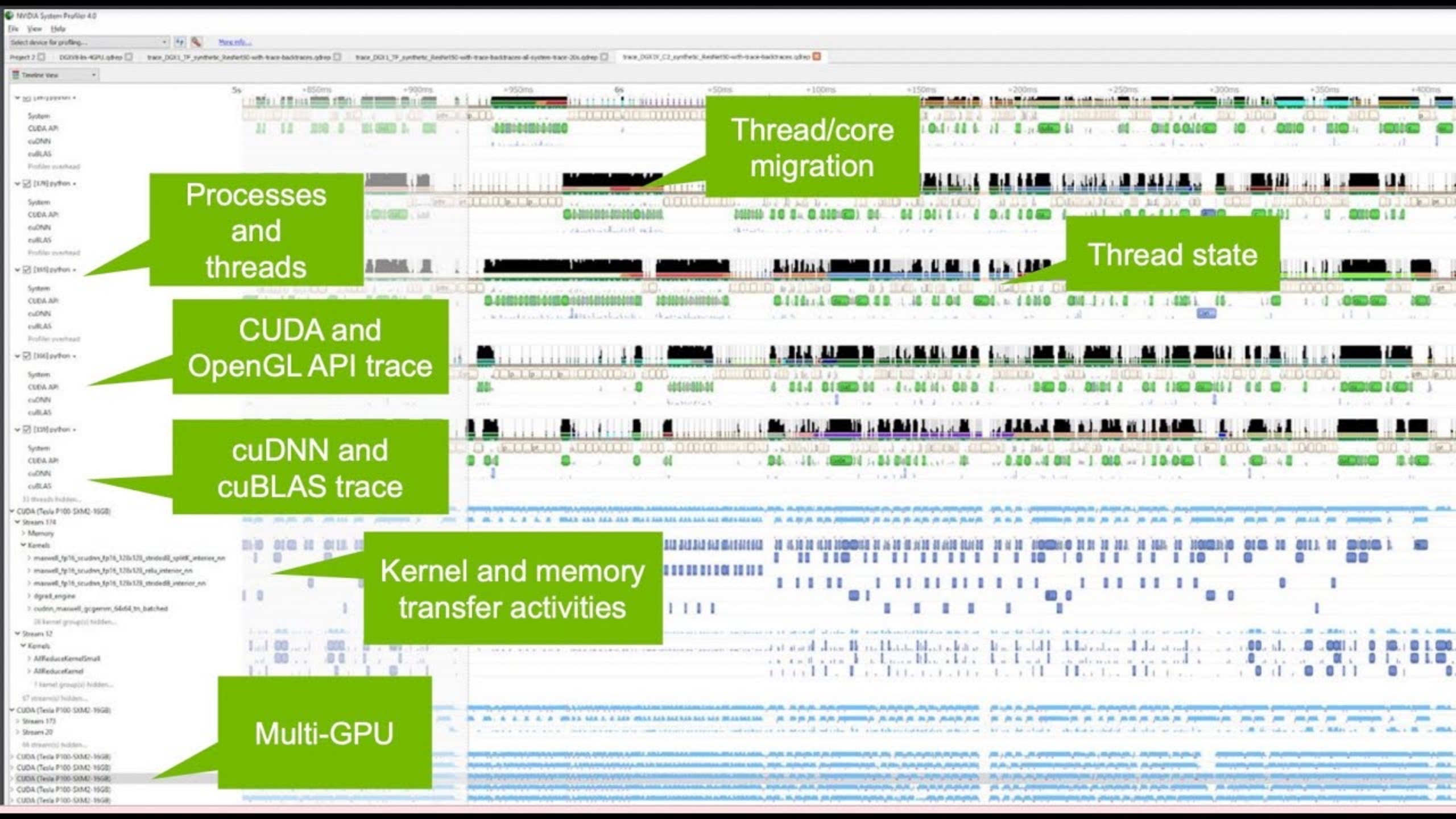
Targeted metric sections for various performance aspects (Debug/&Profile)

Very high freq GPU perf counter, customizable data collection and presentation (tables, charts ...)

Python-based rules for guided analysis (or postprocessing)

GPUs: Volta, Turing, Amper...

Docs/product: <https://developer.nvidia.com/nsight-systems>



NVIDIA Tools Extension API Library (NVTX)



The NVIDIA Tools Extension SDK (NVTX) is a C-based Application Programming Interface (API) for annotating events, code ranges, and resources in your applications. Applications which integrate NVTX can use NVIDIA Nsight VSE to capture and visualize these events and ranges.

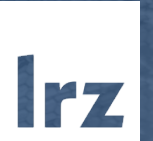
```
void Wait(int waitMilliseconds)
{
    nvtxNameOsThread("MAIN");
    nvtxRangePush(__FUNCTION__);
    nvtxMark(>"Waiting...");
    Sleep(waitMilliseconds);
    nvtxRangePop();
}
int main(void)
{
    nvtxNameOsThread("MAIN");
    nvtxRangePush(__FUNCTION__);
    Wait();
    nvtxRangePop();
}
```

`nsys profile -t nvtx --stats=true ...`

<https://docs.nvidia.com/nsight-visual-studio-edition/2020.1/nvtx/index.html>



DEEP
LEARNING
INSTITUTE



Lab3: Asynchronous Streaming, and Visual Profiling with CUDA C/C++

Dr. Momme Allalen

Leibniz Computing Centre, Munich Germany - www.lrz.de

Deep Learning Certified Instructor, NVIDIA Deep Learning Institute NVIDIA Corporation.

Lab2: Managing Accelerated Application Memory with CUDA Unified Memory and nvprof



Prerequisites

To get the most out of this lab you should already be able to:

- Write, compile, and run C/C++ programs that both call CPU functions and launch GPU kernels.
- Control parallel thread hierarchy using execution configuration.
- Refactor serial loops to execute their iterations in parallel on a GPU.
- Allocate and free CUDA Unified Memory.
- Understand the behavior of Unified Memory with regard to page faulting and data migrations.
- Use asynchronous memory prefetching to reduce page faults and data migrations.

Objectives

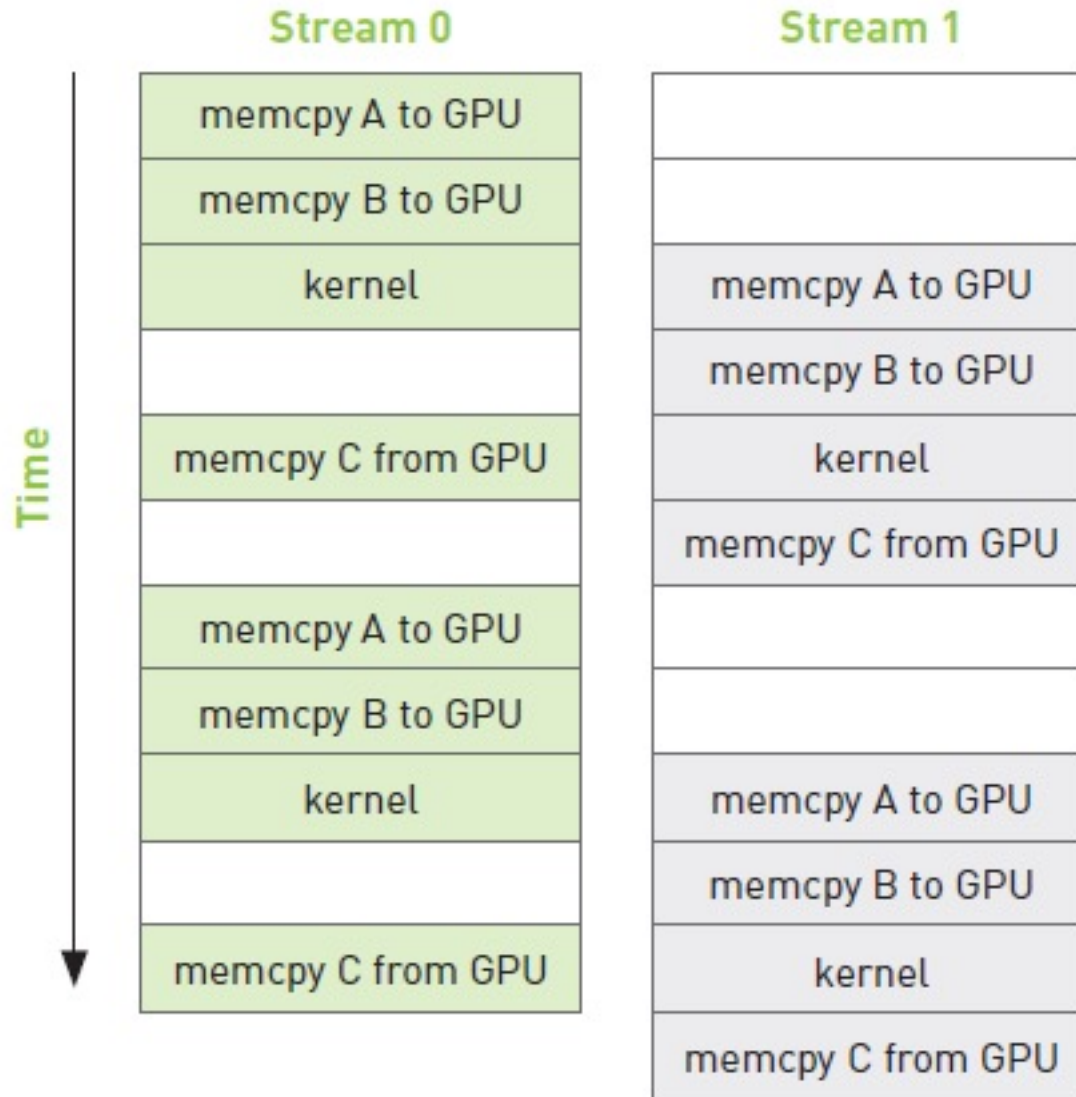
By the time you complete this lab you will be able to:

- Use the **Nsight Systems** to visually profile the timeline of GPU-accelerated CUDA applications.
- Use **Nsight Systems** to identify, and exploit, optimization opportunities in GPU-accelerated CUDA applications.
- Utilize **CUDA streams** for concurrent kernel execution in accelerated applications.
- **(Optional Advanced Content)** Use manual memory allocation, including allocating pinned memory, in order to asynchronously transfer data in concurrent CUDA streams.

Multiple Streams



Overlap copy
with kernel



Multiple Streams



```
for (int i=0; i<FULL_SIZE; i+= N*2) {
// copy the locked memory to the device, async
cudaMemcpyAsync (dev_a0, host_a+i, N * sizeof(int), cudaMemcpyHostToDevice, stream0);
cudaMemcpyAsync (dev_b0, host_b+i, N * sizeof(int), cudaMemcpyHostToDevice, stream0);

kernel<<<N/256,256,0,stream0>>>( dev_a0, dev_b0, dev_c0 );

// copy the data from device to locked memory
cudaMemcpyAsync (host_c+i, dev_c0, N * sizeof(int), cudaMemcpyDeviceToHost, stream0);
// copy the locked memory to the device, async
cudaMemcpyAsync (dev_a1,host_a+i+N, N * sizeof(int), cudaMemcpyHostToDevice, stream1);
cudaMemcpyAsync (dev_b1,host_b+i+N, N * sizeof(int), cudaMemcpyHostToDevice, stream1);

kernel<<<N/256,256,0,stream1>>>( dev_a1, dev_b1, dev_c1 );

// copy the data from device to locked memory
cudaMemcpyAsync (host_c+i+N,dev_c1, N * sizeof(int), cudaMemcpyDeviceToHost, stream1);
}
```



DEEP
LEARNING
INSTITUTE



THANK YOU

Instructor: Dr. Momme Allalen
www.nvidia.com/dli