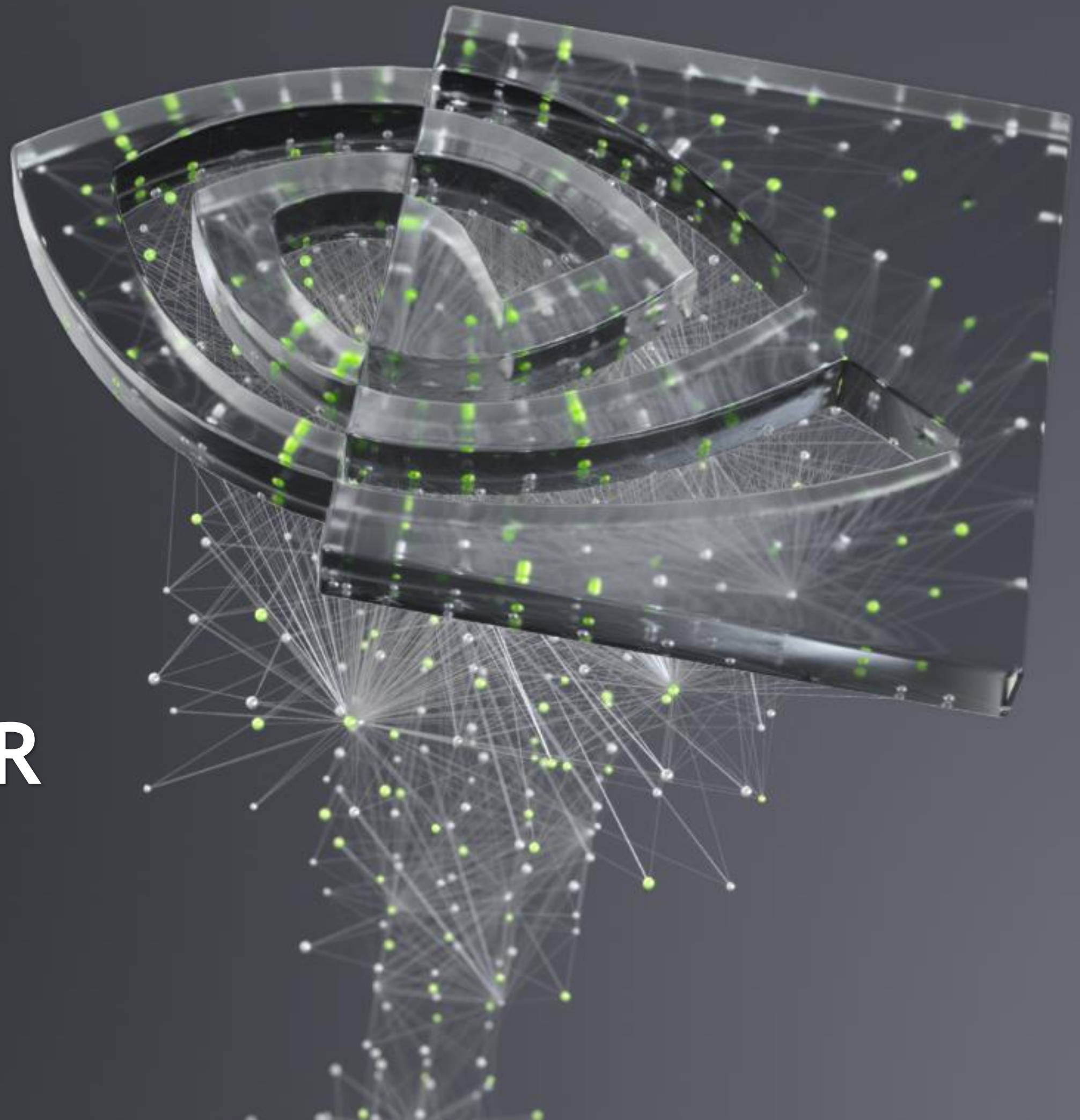# WHAT THE PROFILER IS TELLING YOU

Mozhgan Kabiri chimeh, Sept 2020

# BEFORE YOU START
## Steps to enlightenment

- Know your application

  - What does it compute? How is it parallelized? What final performance is expected?

- Know your hardware

  - What are the target machines and how many? Machine-specific optimizations okay?

- Know your tools

  - Strengths and weaknesses of each tool? Learn how to use them.

- Know your process

  - Performance optimization is a constant learning process.

Outline

1. Overview of the tools
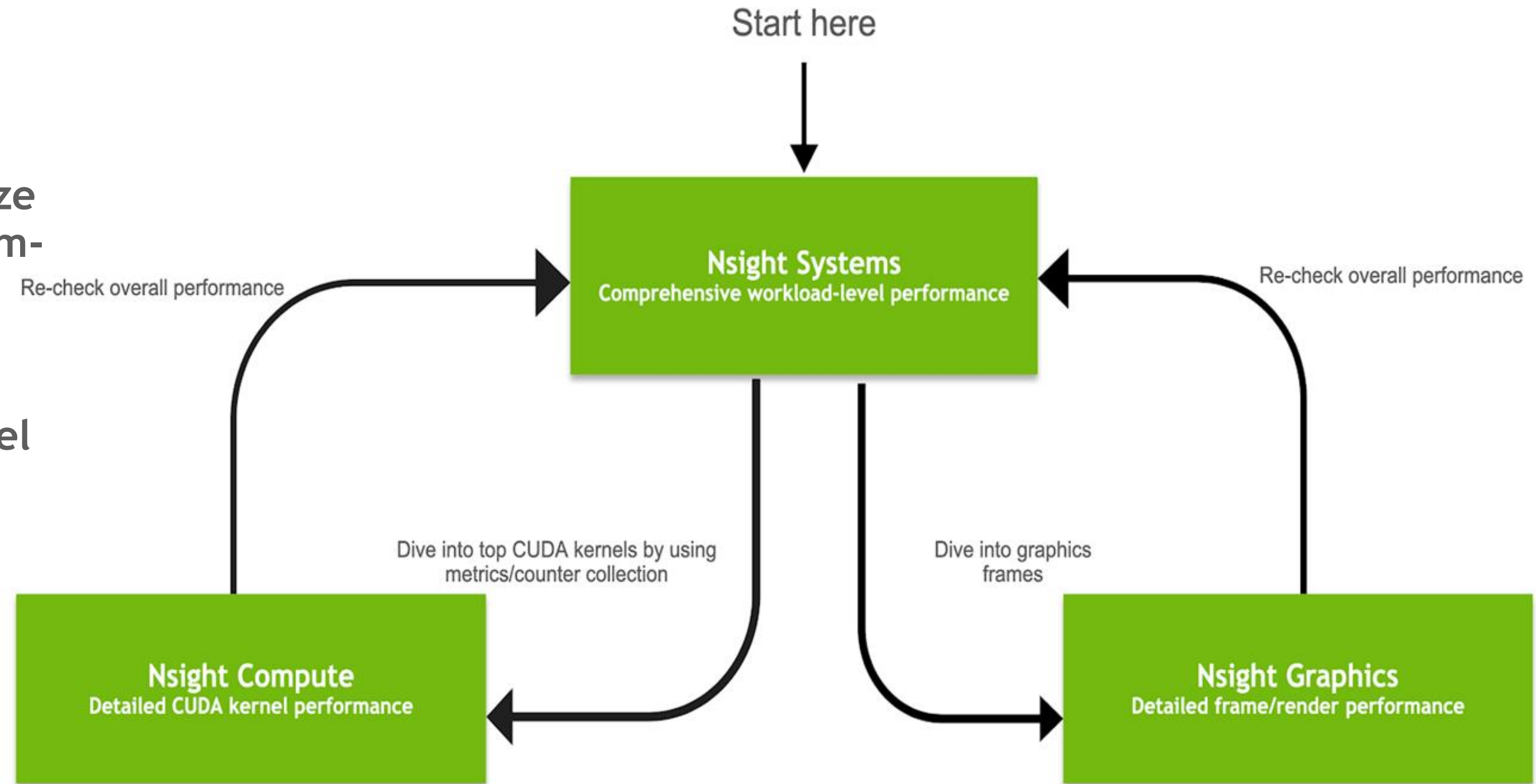
2. Demo

3. Example application

4. Optimization

NVIDIA.

# NVIDIA NSIGHT FAMILY

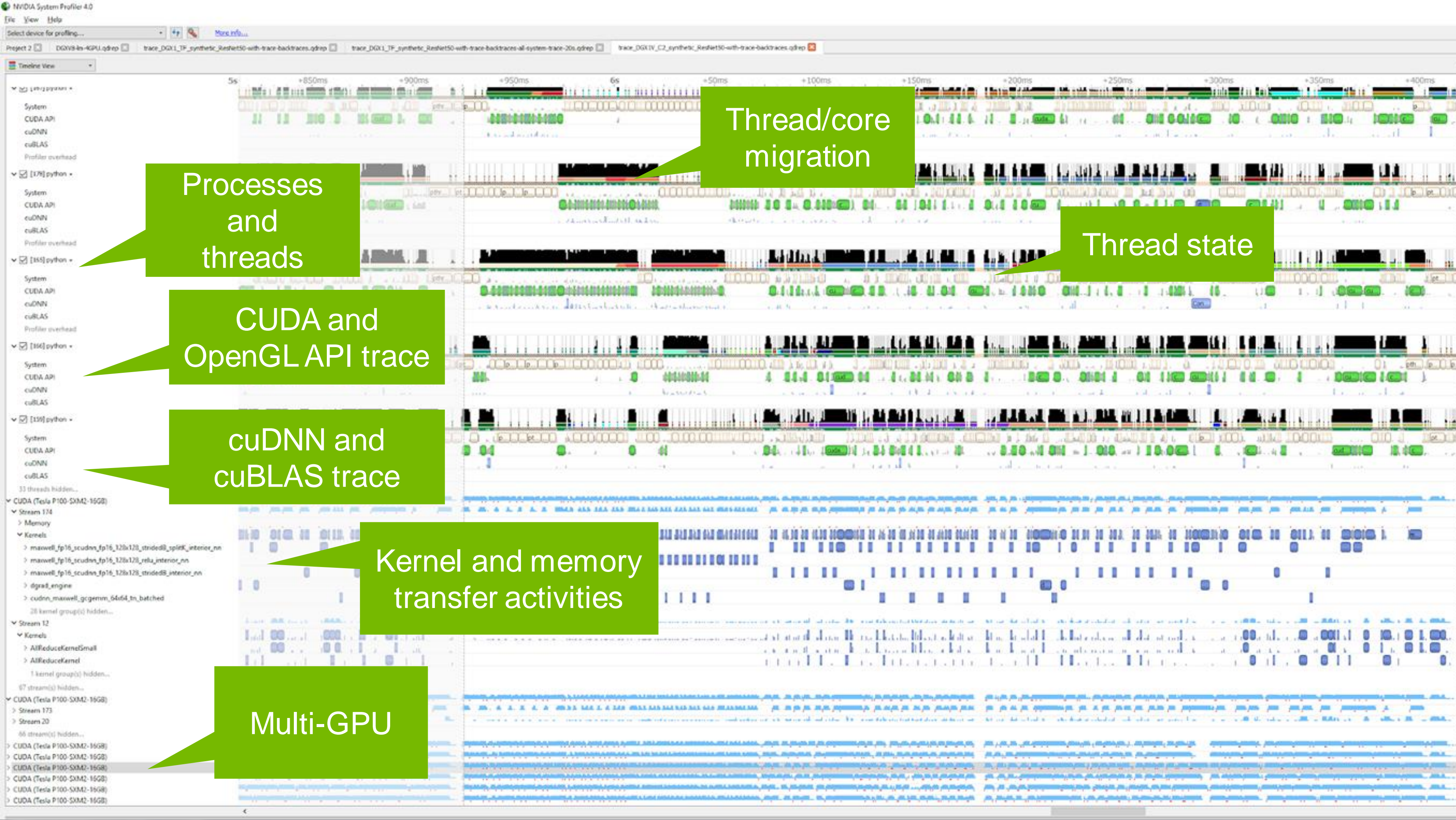# Nsight Product Family

## Workflow

Start here

**Nsight Systems** - Analyze application algorithm system-wide

**Nsight Compute** - Debug/optimize CUDA kernel

**Nsight Graphics** - Debug/optimize graphics workloads

**Nsight Systems**
Comprehensive workload-level performance

Re-check overall performance

Re-check overall performance

Dive into top CUDA kernels by using metrics/counter collection

Dive into graphics frames

**Nsight Compute**
Detailed CUDA kernel performance

**Nsight Graphics**
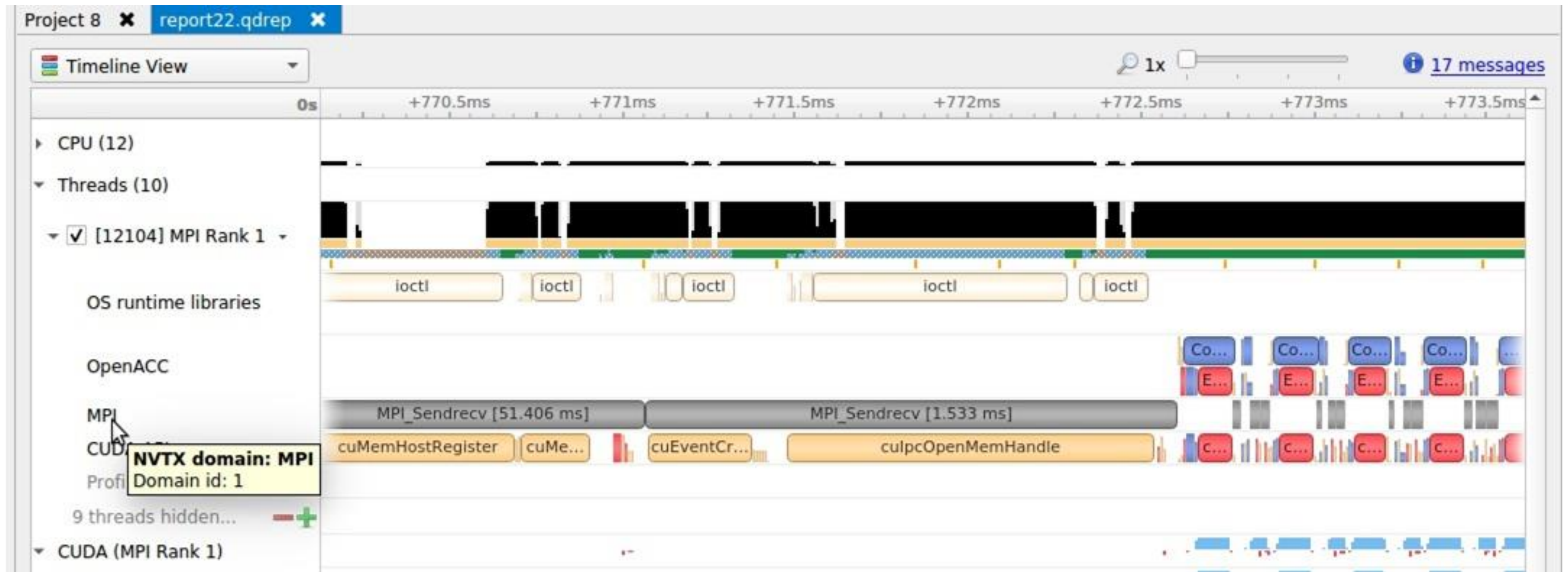Detailed frame/render performance

NVIDIA.

DEMO

# USER ANNOTATIONS APIS FOR CPU & GPU
# NVTX, OPENGL, VULKAN, AND DIRECT3D PERFORMANCE MARKERS

EXAMPLE:  VISUAL MOLECULAR DYNAMICS (VMD) ALGORITHMS VISUALIZED WITH NVTX ON CPU
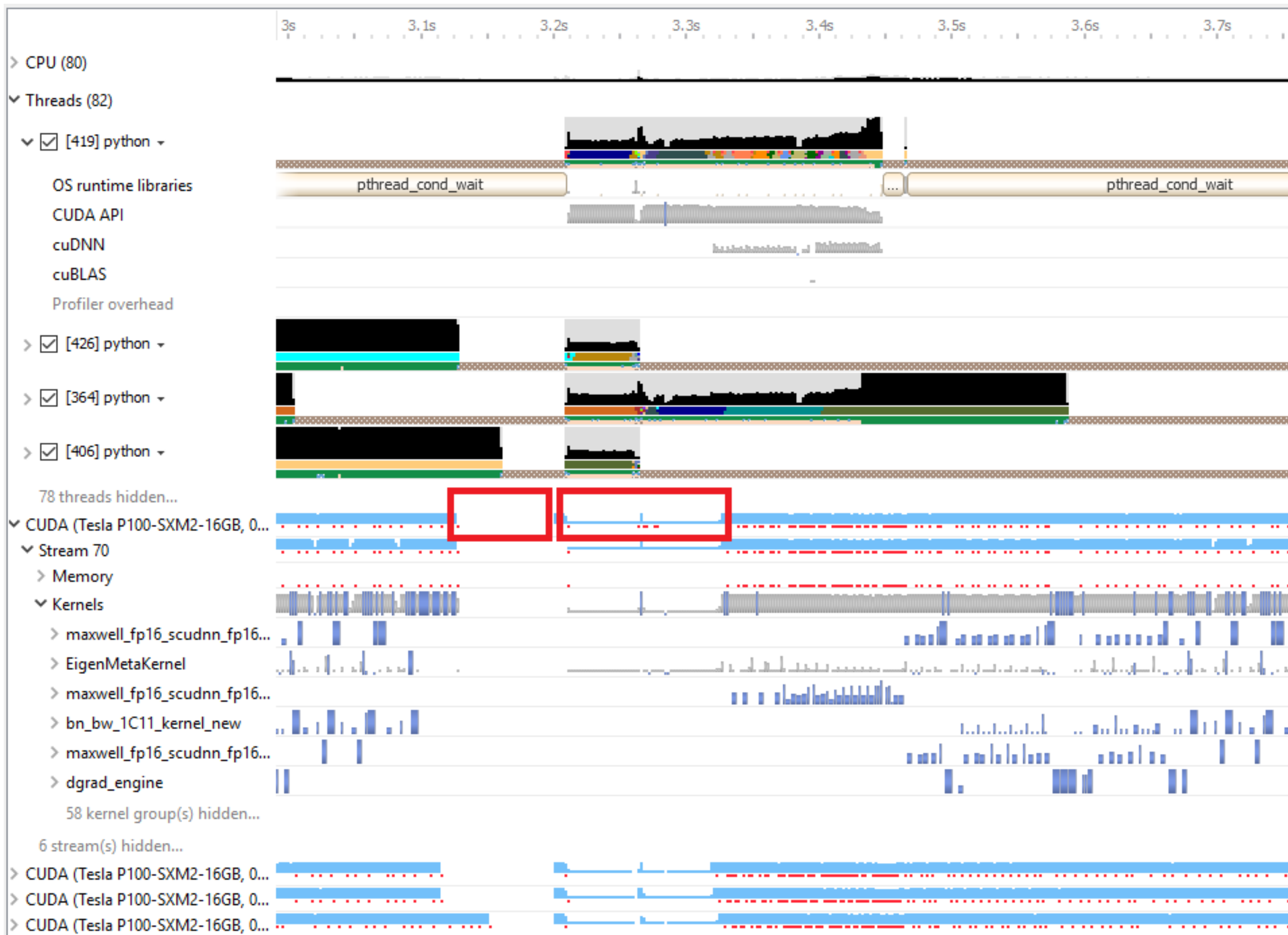
# MPI & OPENACC TRACE

# CORRELATION
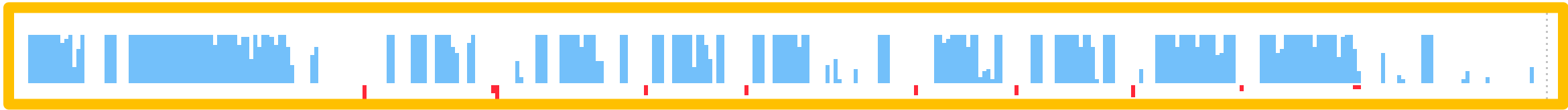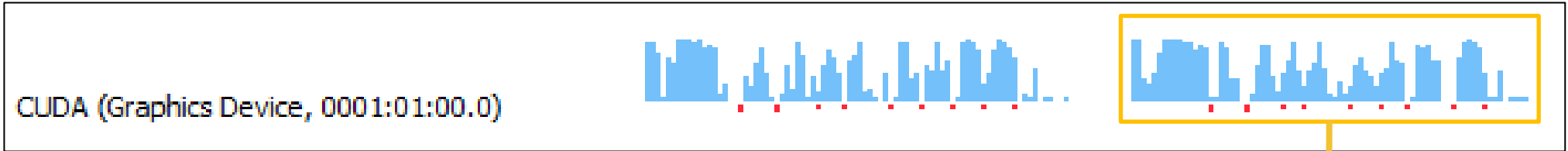
GPU IDLE AND LOW UTILIZATION LEVEL OF DETAIL

# GPU UTILIZATION BASED ON PERCENTAGE TIME COVERAGE



ZOOMING IN REVEALS GAPS WHERE THERE WERE VALLEYS

EXAMPLE APPLICATION

# A SAMPLE OF A FLUID SIMULATION

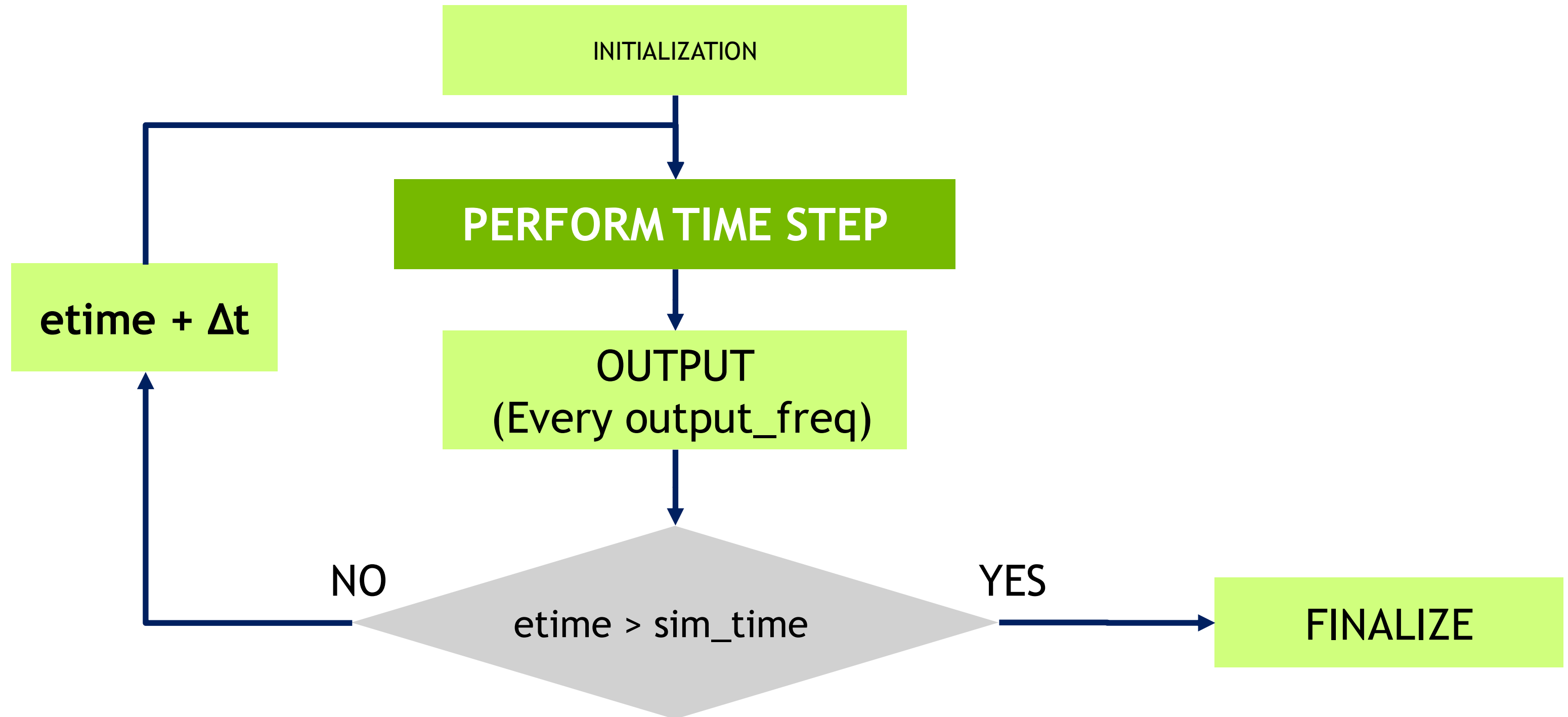## In the context of atmosphere and weather simulation[1]

A narrow jet of fast and slightly cold wind is injected into a balanced, neutral atmosphere at rest from the left domain near the model top.
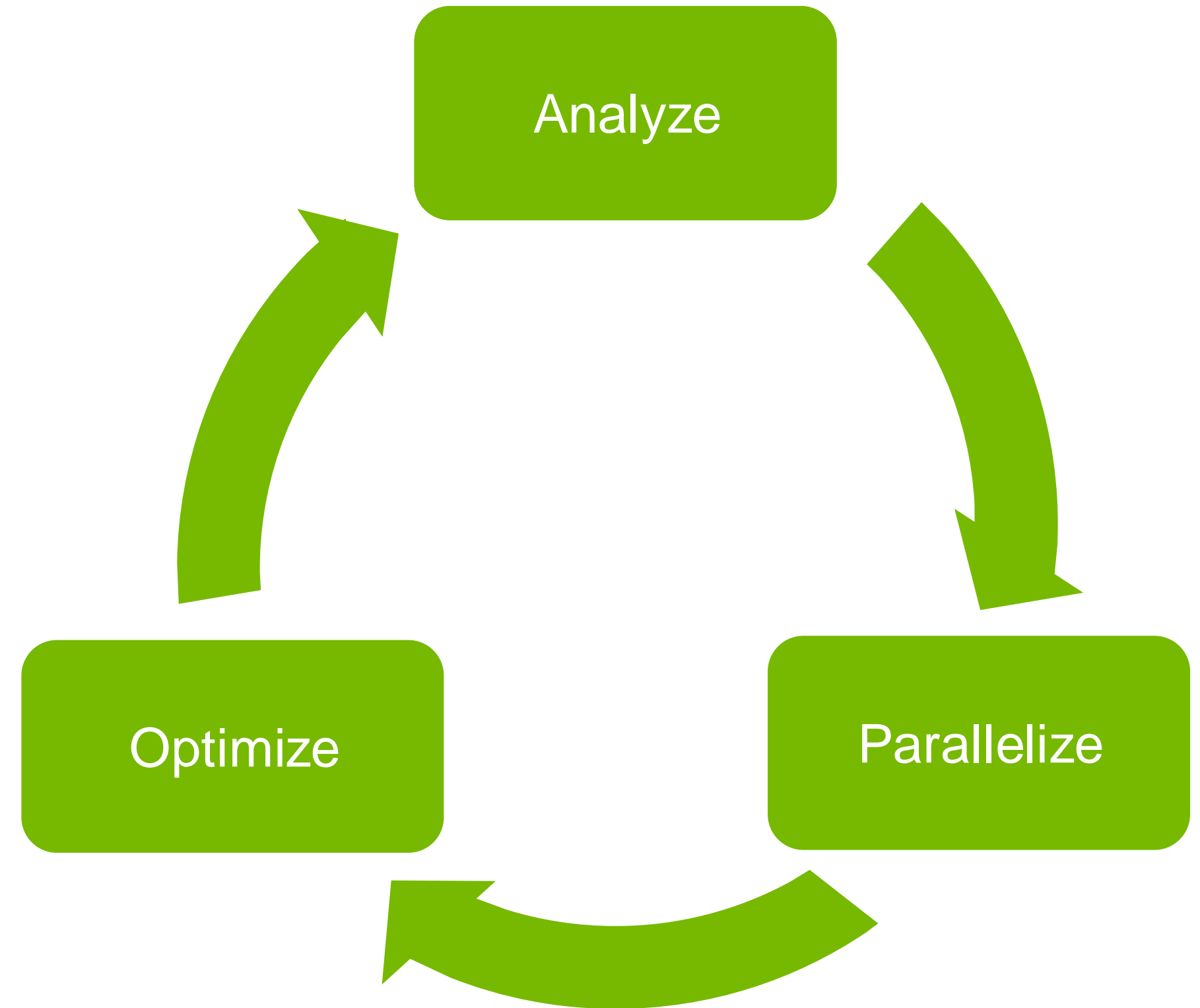


0s          300s                              1000s        **Simulated Time**

[1]. https://github.com/mrnorman/miniWeather

# OUTER LOOP
## The main() function



INITIALIZATION

PERFORM TIME STEP

OUTPUT
(Every output_freq)

etime + Δt

NO    etime > sim_time    YES

FINALIZE

# DEVELOPMENT CYCLE

- **Analyze** your code to determine most likely places needing parallelization or optimization.

- **Parallelize** your code by starting with the most time consuming parts and check for correctness.

- **Optimize** your code to improve observed speed-up from parallelization.

Analyze

Parallelize

Optimize

# TOOLS WE WILL USE: NSIGHT SUITE
## Application-wide profiling (Systems), Kernel-level profiling (Compute)

Instrument with NVIDIA Tools Extension (NVTX):
Automatic or manual

    Create (nested) ranges, define macros

    Compiler instrumentation

Tracing: CUDA API calls, NVTX trace

Sampling, hardware counters

NVTX primer: https://devblogs.nvidia.com/parallelforall/cuda-pro-tip-generate-custom-application-profile-timelines-nvtx/

**Nsight Systems**

**Nsight Compute**

# A FIRST (I)NSIGHT

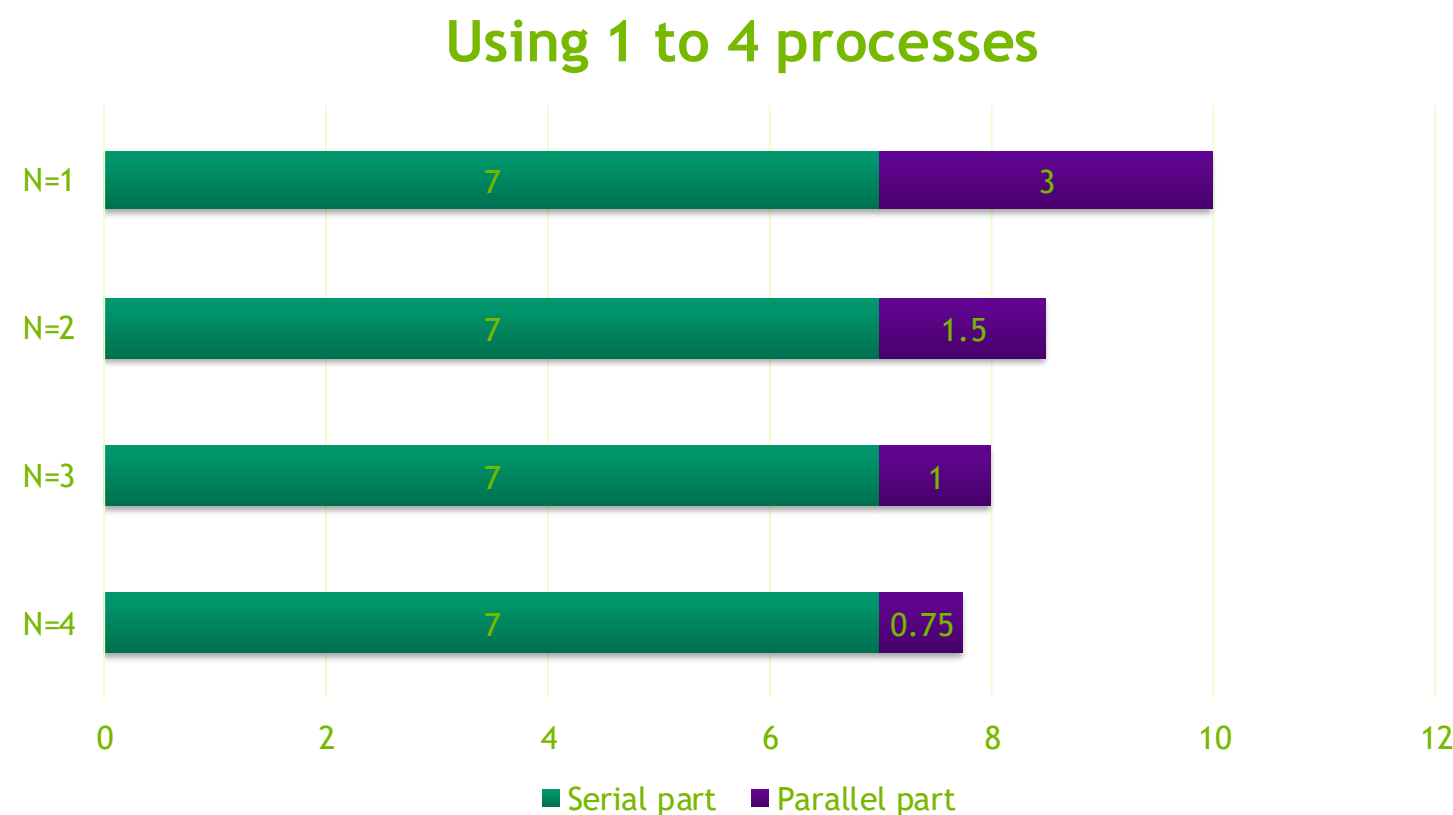## Maximum achievable speedup: Amdahl's law

Amdahl's law states overall speedup $s$ given the parallel fraction $p$ of code and number of processes $N$

$$s = \frac{1}{1 - p + \frac{p}{N}} < \frac{1}{1 - p}$$

Limited by serial fraction, even for $N \to \infty$

Example for $p = 30\%$

Also valid for per-method speedups

**Using 1 to 4 processes**



N=1 — 7, 3

N=2 — 7, 1.5

N=3 — 7, 1

N=4 — 7, 0.75

■ Serial part  ■ Parallel part

# A FIRST (I)NSIGHT
## Recording an application timeline

1) Recording, via the GUI (not shown here) or via command line

```
nsys profile -t nvtx,openacc --stats=true --force-overwrite true -o
my_report ./myapp
```

- `profile` - start a profiling session
- `-t`: Selects the APIs to be traced (nvtx and openacc in this example)
- `--stats`: if true, it generates summary of statistics after the collection
- `--force-overwrite`: if true, it overwrites the existing generated report
- `-o` – name for the intermediate result file, created at the end of the collection (`.qdrep` filename)

```
nsys --help or nsys [specific command] --help
```

2) Inspect results: Open the report file in the GUI

Also possible to get details on command line (documentation)

See also https://docs.nvidia.com/nsight-systems/, "Profiling from the CLI on Linux Devices"

# A FIRST (I)NSIGHT

## Timeline overview in Nsight Systems GUI



Application already ported to GPU – basic guidelines followed

S7122: CUDA Optimization Tips, Tricks and Techniques (2017)

# LOOKING CLOSER
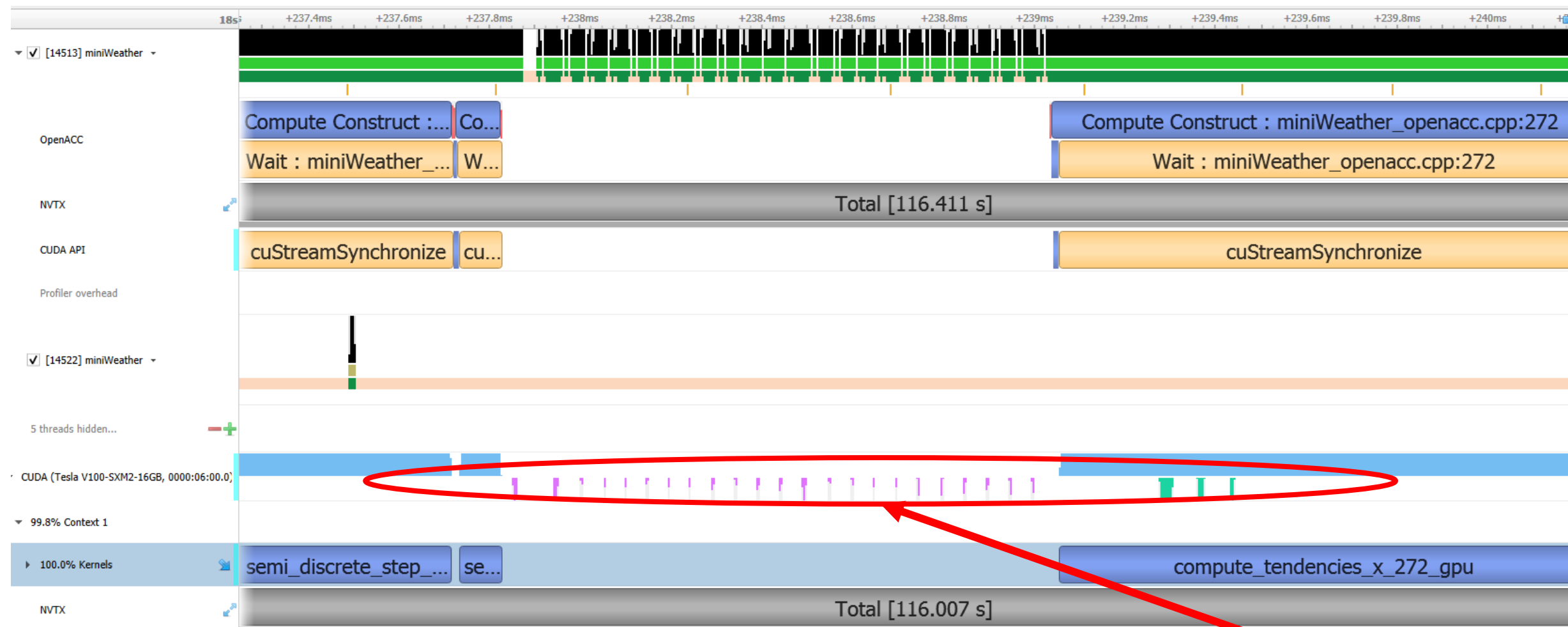


Zooming in and using Events View for NVTX

Useful for other rows, e.g. CUDA API

Hierarchy of ranges, use to locate on timeline

# LOOKING CLOSER



Data transfer

Learn more about Unified Memory:

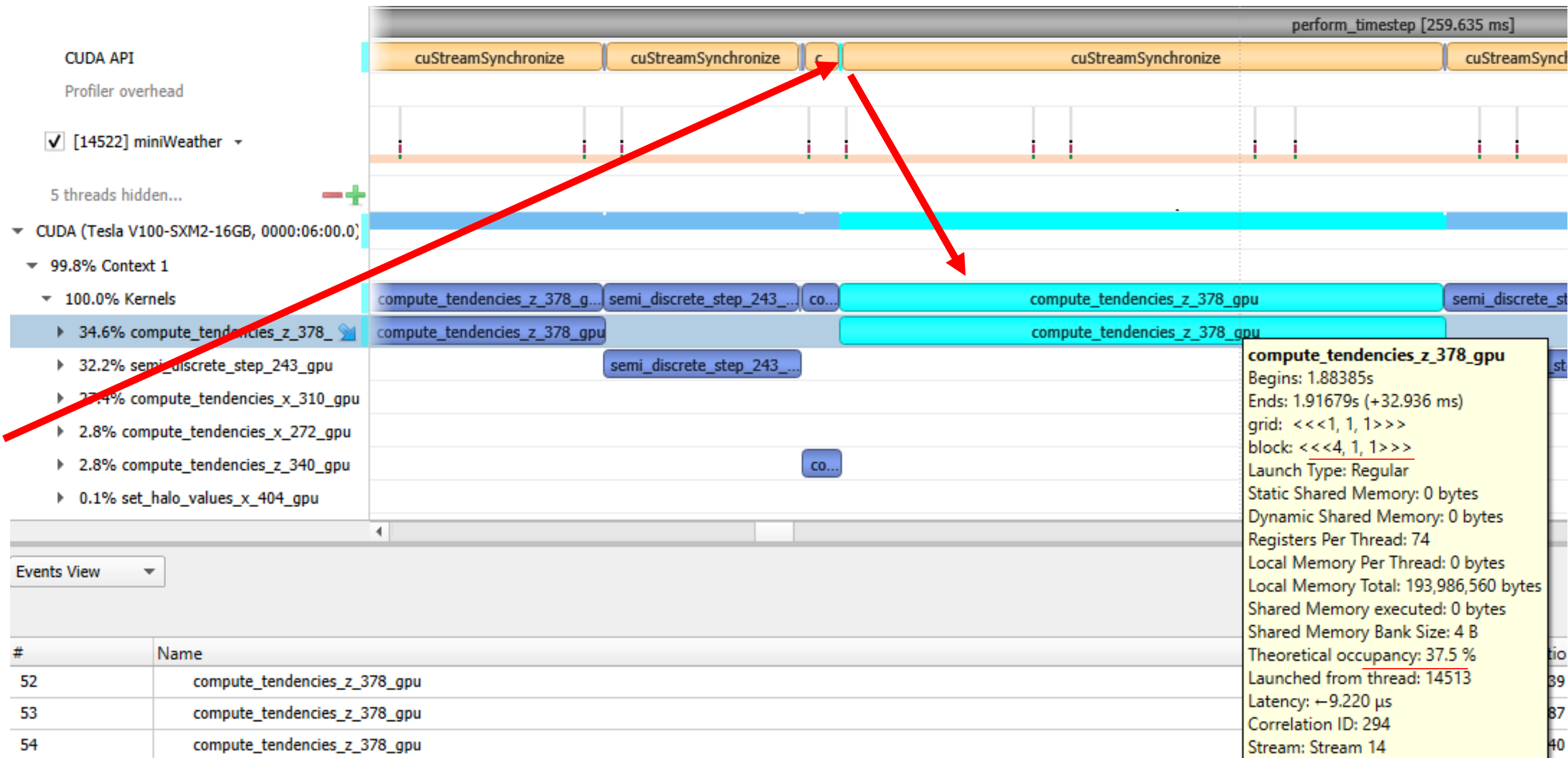S8430: Everything You Need to Know About Unified Memory (2018)

S9727: Memory Management on Modern GPU Architectures (2019)

# IDENTIFYING INTERESTING REGIONS
## How to correlate ranges, API and kernel calls

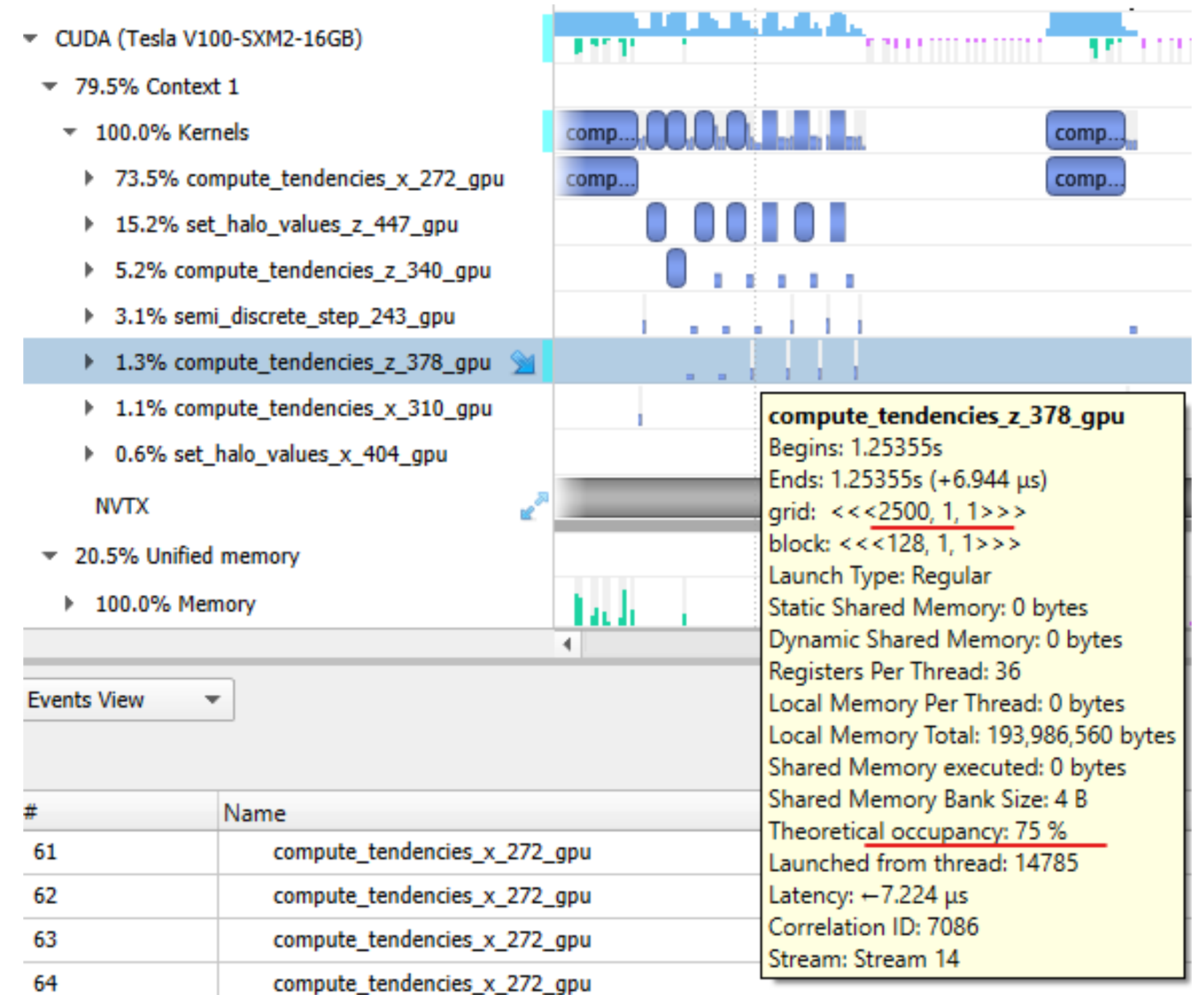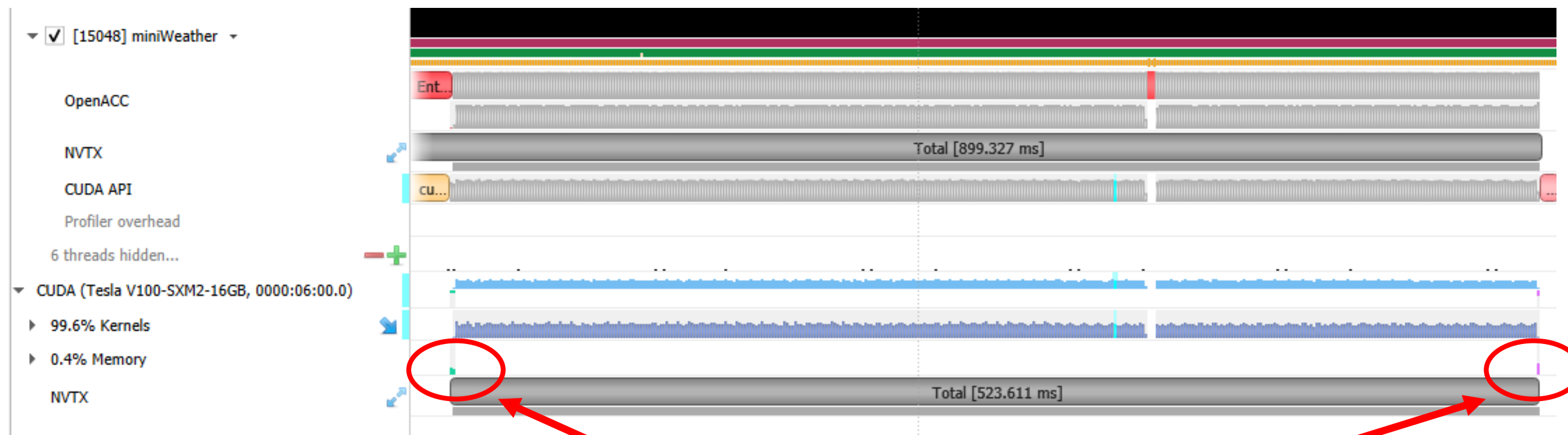Identify components. Mark kernel in CUDA API row, find kernel launch



Finding correlations

# BUILDING A HYPOTHESIS

```
#pragma acc parallel loop private(indt, indf1, indf2)
  for (ll = 0; ll < NUM_VARS; ll++) {
    for (k = 0; k < nz; k++) {
      for (i = 0; i < nx; i++) {
        indt = ll * nz * nx + k * nx + i;
        indf1 = ll * (nz + 1) * (nx + 1) + k * (nx + 1) + i;
        indf2 = ll * (nz + 1) * (nx + 1) + k * (nx + 1) + i + 1;
        tend[indt] = -(flux[indf2] - flux[indf1]) / dx;
} } }
```

**Hypothesis:** small iteration count for outer loop

**Solution:** flatten the loop by collapsing the tightly-nested loops

# LOOKING CLOSER



Data transfers

Learn more about Unified Memory:

S8430: Everything You Need to Know About Unified Memory (2018)

S9727: Memory Management on Modern GPU Architectures (2019)

**Result:**
Baseline = 116 seconds
Current = 900 ms

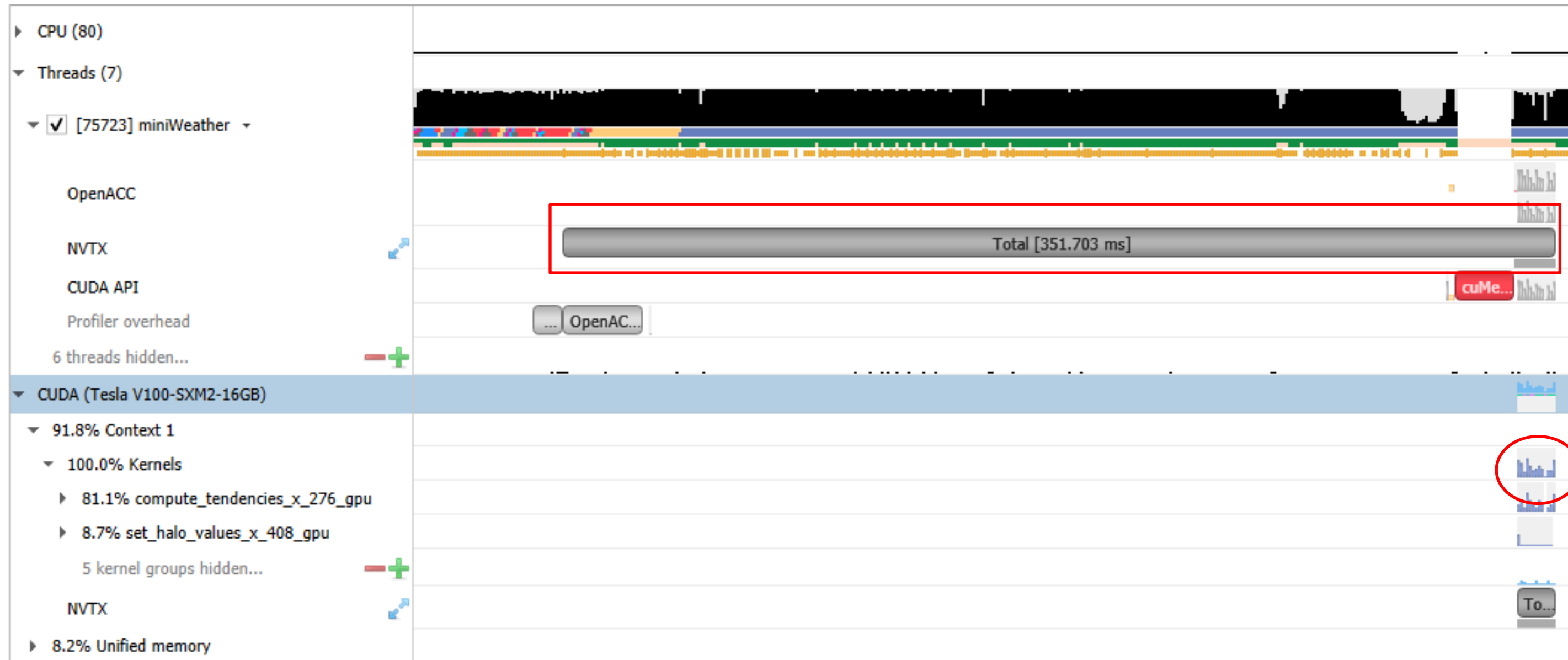# RECAP: HIGH-LEVEL ANALYSIS
## Application timeline with Nsight Systems

Annotated code

Analyzed regions (NVTX)

Identified first optimization target (Wallclock, Amdahl's law)

Correlated with actual kernel launch

**Now**: Look briefly at the Nsight Compute

# EXAMPLE ANALYSIS OF A KERNEL
## Analysis with Nsight Compute



Right-click menu in Nsight Systems,
get command line

Run command line

```
ncu --set full -k compute_tendencies_x_276_gpu -s 4 -c 1 -o myreport ./myapp
```

Important switches for metrics collection, pre-selected sets

Fully customizable, `ncu --help`. Check `--list-metrics` and `--query-metrics`

We use GUI for analysis and load report file

Alternatively, interactively run and analyze directly through GUI

See also https://docs.nvidia.com/nsight-compute/, "Nsight Compute CLI"

# NSIGHT COMPUTE OVERVIEW



switch between report pages

Kernel

Sections

Current    171 - compute_tendencies_x_277_gpu ( 896,  1,  1)  **Time:** 11.94 usecond  **Cycles:** 9,838  **Regs:** 72  **GPU:** TITAN V  **SM Frequency:** 824.15 cycle/usecond  **CC:** 7.0  **Process:** [4570] miniWeather_cpp

▼ GPU Speed Of Light ⚠    All

High-level overview of the utilization for compute and memory resources of the GPU. For each unit, the Speed Of Light (SOL) reports the achieved percentage of utilization with respect to the theoretical maximum.

| | | |
|---|---|---|
| SOL SM [%] | 0.98 | Duration [usecond] 11.94 |
| SOL Memory [%] | 1.50 | Elapsed Cycles [cycle] 9,838 |
| SOL TEX [%] | 2.11 | SM Active Cycles [cycle] 1,059.54 |
| SOL L2 [%] | 1.25 | SM Frequency [cycle/usecond] 824.15 |
| SOL FB [%] | 1.50 | Memory Frequency [cycle/usecond] 579.09 |

**GPU Utilization**

SM [%]
Memory [%]
0.0  10.0  20.0  30.0  40.0  50.0  60.0  70.0  80.0  90.0  100.0
Speed Of Light [%]

**SOL SM Breakdown**

| | |
|---|---|
| SOL SM: Pipe Shared Cycles Active [%] | 0.98 |
| SOL SM: Pipe Fp64 Cycles Active [%] | 0.98 |
| SOL SM: Issue Active [%] | 0.67 |
| SOL SM: Inst Executed [%] | 0.66 |
| SOL SM: Inst Executed Pipe Cbu Pred On Any [%] | 0.36 |
| SOL SM: Pipe Fma Cycles Active [%] | 0.28 |
| SOL SM: Pipe Alu Cycles Active [%] | 0.27 |
| SOL SM: Inst Executed Pipe Xu [%] | 0.19 |
| SOL SM: Inst Executed Pipe Lsu [%] | 0.10 |
| SOL SM: Mio Pq Read Cycles Active [%] | 0.07 |
| SOL SM: Mio Pq Write Cycles Active [%] | 0.07 |
| SOL SM: Mio2rf Writeback Active [%] | 0.07 |
| SOL SM: Mio Inst Issued [%] | 0.05 |
| SOL SM: Inst Executed Pipe Adu [%] | 0.01 |
| SOL SM: Inst Executed Pipe Fp16 [%] | 0 |
| SOL SM: Inst Executed Pipe Ipa [%] | 0 |
| SOL SM: Inst Executed Pipe Tex [%] | 0 |
| SOL IDC: Request Cycles Active [%] | 0 |
| SOL SM: Pipe Tensor Cycles Active [%] | 0 |

**SOL Memory Breakdown**

| | |
|---|---|
| SOL GPU: Dram Throughput [%] | 1.50 |
| SOL L2: T Sectors [%] | 1.25 |
| SOL L2: Lts2xbar Cycles Active [%] | 1.21 |
| SOL L2: Xbar2lts Cycles Active [%] | 0.59 |
| SOL L2: D Sectors [%] | 0.40 |
| SOL L1: Data Pipe Lsu Wavefronts [%] | 0.23 |
| SOL L2: D Sectors Fill Device [%] | 0.20 |
| SOL L1: M L1tex2xbar Req Cycles Active [%] | 0.20 |
| SOL L2: T Tag Requests [%] | 0.19 |
| SOL L1: Lsu Writeback Active [%] | 0.14 |
| SOL L1: M Xbar2l1tex Read Sectors [%] | 0.12 |
| SOL L1: Lsuin Requests [%] | 0.10 |
| SOL L1: Data Bank Reads [%] | 0.07 |
| SOL L1: Texin Sm2tex Req Cycles Active [%] | 0.03 |
| SOL L1: F Wavefronts [%] | 0.03 |
| SOL L1: Data Bank Writes [%] | 0.03 |
| SOL L2: D Atomic Input Cycles Active [%] | 0 |
| SOL L2: D Sectors Fill Sysmem [%] | 0 |
| SOL L1: Tex Writeback Active [%] | 0 |
| SOL L1: Data Pipe Tex Wavefronts [%] | 0 |

2

**Recommendations**

⚠ **Bottleneck**   [Warning] This kernel grid is too small to fill the available resources on this device. Look at `Launch Statistics` for more details.

▼ Launch Statistics ⚠

Summary of the configuration used to launch the kernel. The launch configuration defines the size of the kernel grid, the division of the grid into blocks, and the GPU resources needed to execute the kernel. Choosing an efficient launch configuration maximizes device utilization.

| | | |
|---|---|---|
| Grid Size | 7 | Registers Per Thread [register/thread] 72 |
| Block Size | 128 | Static Shared Memory Per Block [byte/block] 0 |
| Threads [thread] | 896 | Dynamic Shared Memory Per Block [byte/block] 0 |
| Waves Per SM | 0.01 | Shared Memory Configuration Size [byte] 0 |

3

**Recommendations**

⚠ **Launch Configuration**   [Warning] The grid for this launch is configured to execute only 7 blocks, which is less than the GPU's 80 multiprocessors. This can underutilize some multiprocessors. If you do not intend to execute this kernel concurrently with other workloads, consider reducing the block size to have at least one block per multiprocessor or increase the size of the grid to fully utilize the available hardware resources.

▶ Occupancy

Occupancy is the ratio of the number of active warps per multiprocessor to the maximum number of possible active warps. Another way to view occupancy is the percentage of the hardware's ability to process warps that is actively in use. Higher occupancy does not always result in higher performance, however, low occupancy always reduces the ability to hide latencies, resulting in overall performance degradation. Large discrepancies between the theoretical and the achieved occupancy during

30  ⬢ NVIDIA.

Page: Details | Process: All | Launch: compute_tendencies_x_277_gpu | **Add Baseline** | **1** | Rules | Copy as Image

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **Current** | 171 - compute_tendencies_x_277_gpu (80256, 1, 1) | **Time:** 23.30 usecond | **Cycles:** 21,789 | **Regs:** 72 | **GPU:** TITAN V | **SM Frequency:** 932.74 cycle/usecond | **CC:** 7.0 | **Process:** [4606] miniWeather_cpp |
| **Baseline 10** | 171 - compute_tendencies_x_277_gpu ( 896, 1, 1) | **Time:** 11.94 usecond | **Cycles:** 9,838 | **Regs:** 72 | **GPU:** TITAN V | **SM Frequency:** 824.15 cycle/usecond | **CC:** 7.0 | **Process:** [4570] miniWeather_cpp |

## ▼ GPU Speed Of Light ⚠

All

High-level overview of the utilization for compute and memory resources of the GPU. For each unit, the Speed Of Light (SOL) reports the achieved percentage of utilization with respect to the theoretical maximum.

| | | | |
|---|---|---|---|
| SOL SM [%] | 42.98 (+4,265.33%) | Duration [usecond] | 23.30 (+95.17%) |
| SOL Memory [%] | 42.50 (+2,732.85%) | Elapsed Cycles [cycle] | 21,789 (+121.48%) |
| SOL TEX [%] | 8.18 (+288.10%) | SM Active Cycles [cycle] | 23,515.61 (+2,119.42%) |
| SOL L2 [%] | 17.69 (+1,312.64%) | SM Frequency [cycle/usecond] | 932.74 (+13.18%) |
| SOL FB [%] | 42.50 (+2,732.85%) | Memory Frequency [cycle/usecond] | 664.84 (+14.81%) |

**2**

### GPU Utilization

(Bar chart: SM [%] and Memory [%] vs Speed Of Light [%], axis 0.0 to 100.0)

### SOL SM Breakdown

| | |
|---|---|
| SOL SM: Pipe Shared Cycles Active [%] | 42.98 (+4,265.33%) |
| SOL SM: Pipe Fp64 Cycles Active [%] | 42.98 (+4,265.33%) |
| SOL SM: Issue Active [%] | 28.67 (+4,176.52%) |
| SOL SM: Inst Executed [%] | 28.53 (+4,215.35%) |
| SOL SM: Inst Executed Pipe Cbu Pred On Any [%] | 14.18 (+3,888.23%) |
| SOL SM: Pipe Fma Cycles Active [%] | 11.90 (+4,215.21%) |
| SOL SM: Pipe Alu Cycles Active [%] | 11.61 (+4,218.01%) |
| SOL SM: Inst Executed Pipe Xu [%] | 8.08 (+4,229.83%) |
| SOL SM: Inst Executed Pipe Lsu [%] | 4.47 (+4,233.24%) |
| SOL SM: Mio Pq Read Cycles Active [%] | 3.01 (+4,334.69%) |
| SOL SM: Mio2rf Writeback Active [%] | 2.93 (+4,354.76%) |
| SOL SM: Mio Pq Write Cycles Active [%] | 2.88 (+4,265.33%) |
| SOL SM: Mio Inst Issued [%] | 2.31 (+4,223.97%) |
| SOL SM: Inst Executed Pipe Adu [%] | 0.29 (+3,955.14%) |
| SOL SM: Inst Executed Pipe Fp16 [%] | 0 (+0.00%) |
| SOL SM: Inst Executed Pipe Ipa [%] | 0 (+0.00%) |
| SOL SM: Inst Executed Pipe Tex [%] | 0 (+0.00%) |
| SOL IDC: Request Cycles Active [%] | 0 (+0.00%) |
| SOL SM: Pipe Tensor Cycles Active [%] | 0 (+0.00%) |

### SOL Memory Breakdown

| | |
|---|---|
| SOL GPU: Dram Throughput [%] | 42.50 (+2,732.85%) |
| SOL L2: T Sectors [%] | 17.69 (+1,312.64%) |
| SOL L2: Xbar2lts Cycles Active [%] | 13.81 (+2,223.80%) |
| SOL L2: Lts2xbar Cycles Active [%] | 9.39 (+675.63%) |
| SOL L1: Data Pipe Lsu Wavefronts [%] | 8.85 (+3,799.58%) |
| SOL L2: D Sectors [%] | 8.50 (+1,997.73%) |
| SOL L2: D Sectors Fill Device [%] | 8.18 (+3,964.39%) |
| SOL L1: M L1tex2xbar Req Cycles Active [%] | 7.75 (+3,822.76%) |
| SOL L2: T Tag Requests [%] | 6.05 (+3,011.80%) |
| SOL L1: Lsu Writeback Active [%] | 6.05 (+4,185.62%) |
| SOL L1: M Xbar2l1tex Read Sectors [%] | 4.94 (+4,037.34%) |
| SOL L1: Lsuin Requests [%] | 4.47 (+4,233.24%) |
| SOL L1: Data Bank Reads [%] | 2.90 (+4,318.88%) |
| SOL L1: Data Bank Writes [%] | 1.19 (+4,173.87%) |
| SOL L1: Texin Sm2tex Req Cycles Active [%] | 0.01 (-54.73%) |
| SOL L1: F Wavefronts [%] | 0.01 (-54.73%) |
| SOL L2: D Atomic Input Cycles Active [%] | 0 (+0.00%) |
| SOL L2: D Sectors Fill Sysmem [%] | 0 (+0.00%) |
| SOL L1: Tex Writeback Active [%] | 0 (+0.00%) |
| SOL L1: Data Pipe Tex Wavefronts [%] | 0 (+0.00%) |

**3**

### Recommendations

⚠ **Bottleneck** [Warning] This kernel exhibits low compute throughput and memory bandwidth utilization relative to the peak performance of this device. Achieved compute throughput and/or memory bandwidth below 60.0% of peak typically indicate latency issues. Look at `Scheduler Statistics` and `Warp State Statistics` for potential reasons.

## ▶ Launch Statistics ⚠

Summary of the configuration used to launch the kernel. The launch configuration defines the size of the kernel grid, the division of the grid into blocks, and the GPU resources needed to execute the kernel. Choosing an efficient launch configuration maximizes device utilization.

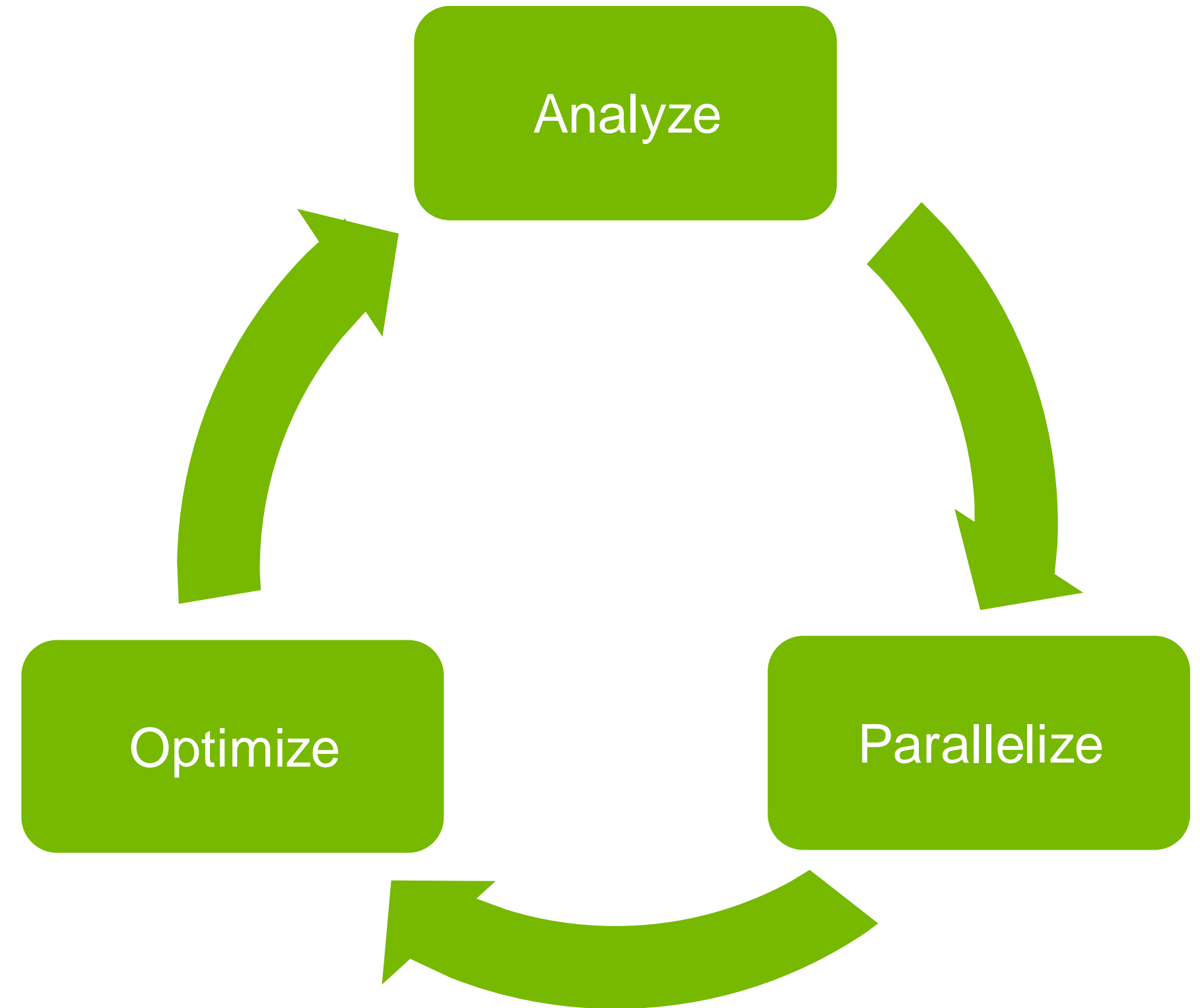| | | | |
|---|---|---|---|
| Grid Size | 627 (+8,857.14%) | Registers Per Thread [register/thread] | 72 (+0.00%) |
| Block Size | 128 (+0.00%) | Static Shared Memory Per Block [byte/block] | 0 (+0.00%) |
| Threads [thread] | 80,256 (+8,857.14%) | Dynamic Shared Memory Per Block [byte/block] | 0 (+0.00%) |
| Waves Per SM | 1.12 (+8,857.14%) | Shared Memory Configuration Size [byte] | 0 (+0.00%) |

## ▶ Occupancy

Occupancy is the ratio of the number of active warps per multiprocessor to the maximum number of possible active warps. Another way to view occupancy is the percentage of the hardware's ability to process warps that is actively in use. Higher occupancy does not always result in higher performance, however, low occupancy always reduces the ability to hide latencies, resulting in overall performance degradation. Large discrepancies between the theoretical and the achieved occupancy during execution typically indicates highly imbalanced workloads.

| | | | |
|---|---|---|---|
| Theoretical Occupancy [%] | 43.75 (+0.00%) | Block Limit Registers [block] | 7 (+0.00%) |
| Theoretical Active Warps per SM [warp/cycle] | 28 (+0.00%) | Block Limit Shared Mem [block] | 32 (+0.00%) |

31

# SUMMARY

## Performance Optimization is a Constant Learning Process

1. Know your application
2. Know your hardware
3. Know your tools
4. Know your process
   1. Identify the Hotspot
   2. Classify the Performance Limiter
   3. Look for indicators

Analyze

Parallelize

Optimize

NVIDIA

# ADDITIONAL REFERENCES

- Documentation

  - https://www.openacc.org/resources

  - https://docs.nvidia.com/cuda/cuda-c-programming-guide/

  - https://docs.nvidia.com/cuda/cuda-c-best-practices-guide/

- GTC 2020 Sessions:

  - What the Profiler is Telling You: How to Get the Most Performance out of Your Hardware [S22141]

    - https://www.nvidia.com/en-us/gtc/on-demand/?search=s22141

- More:

  - NVIDIA Nsight Compute

    - https://vimeo.com/398929189

  - NVIDIA Nsight Systems

    - https://vimeo.com/398838139

THANK YOU!