



GPU PERFORMANCE ANALYSIS SCORE-P AND NVIDIA TOOLS

APRIL 8, 2022 | MICHAEL KNOBLOCH

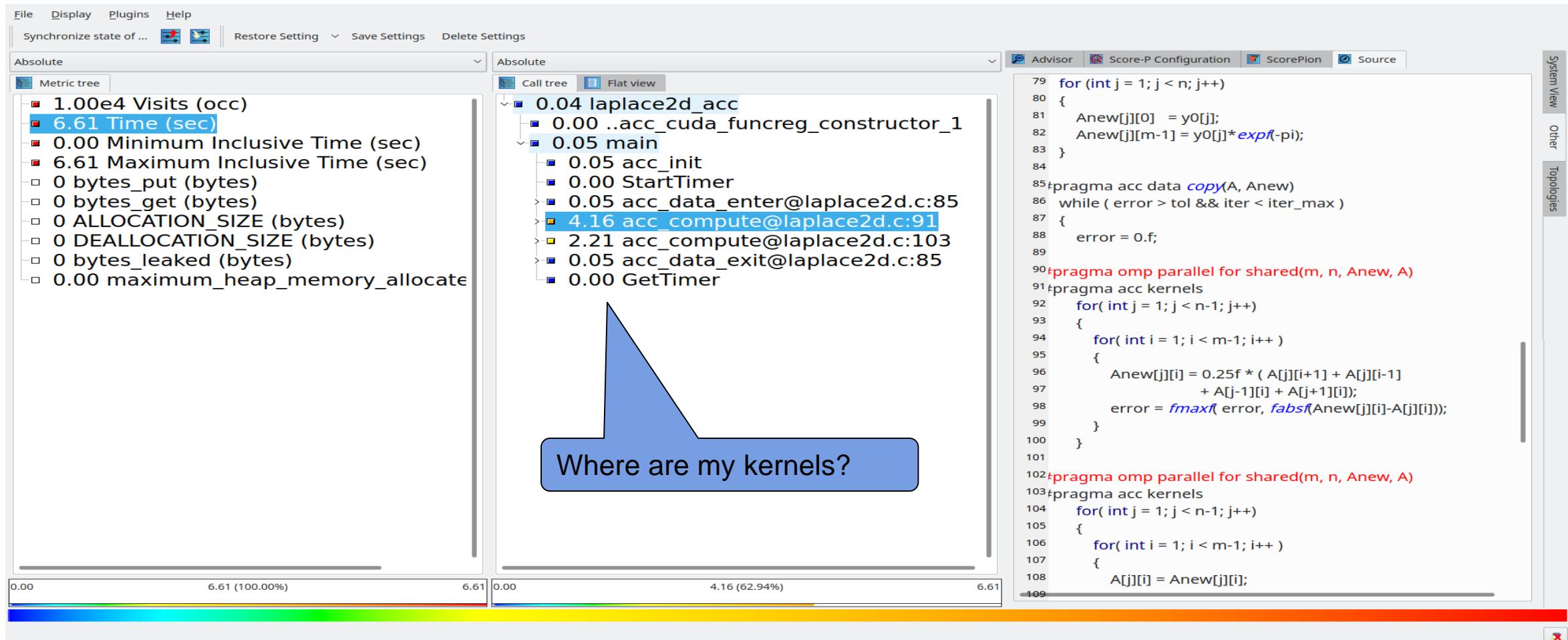
GPU Analysis using the **SCORE-P ECOSYSTEM**

Mitglied der Helmholtz-Gemeinschaft

SCORE-P GPU MEASUREMENTS

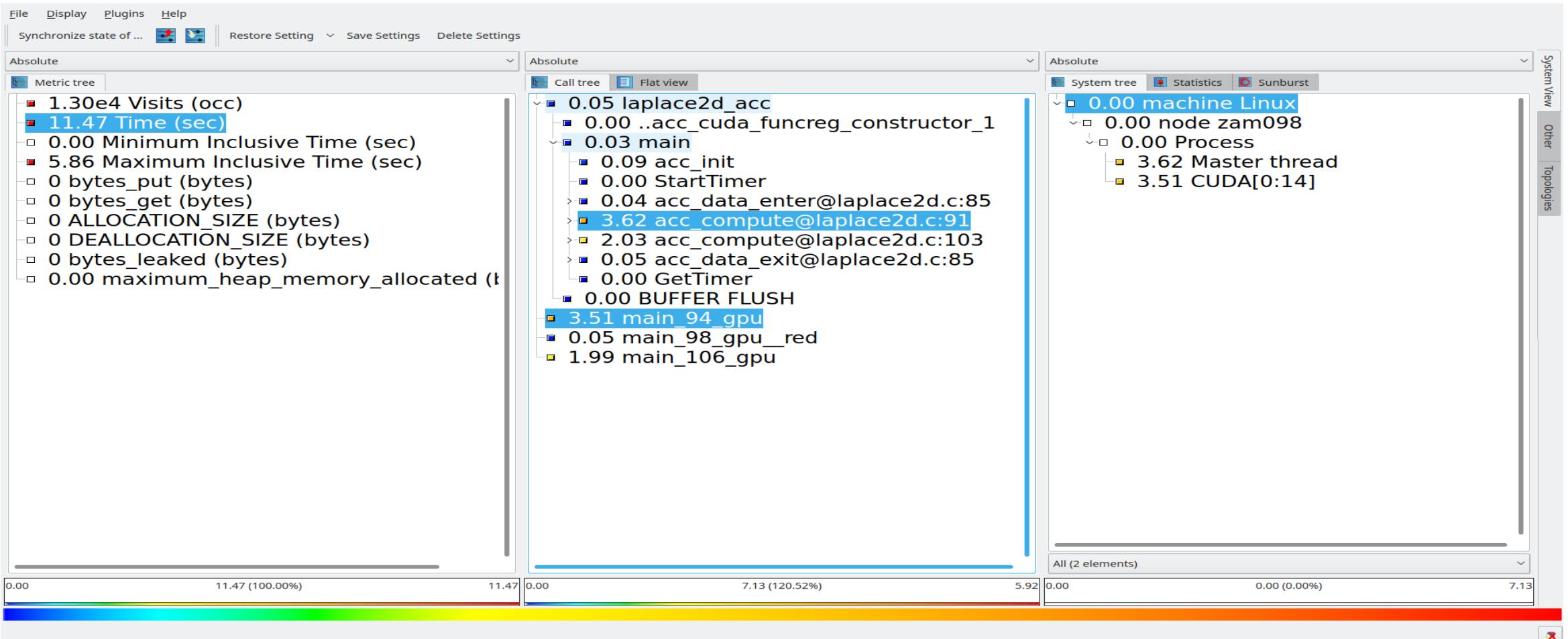
- OpenACC
 - Prefix compiler and linker command with scorep --openacc
 - export ACC_PROFILIB=\$SCOREP_ROOT/lib/libscorep_adapter_openacc_event.so
 - export SCOREP_OPENACC_ENABLE=yes
 - yes refers to: regions, wait, enqueue
 - Full list of options in User Guide
- CUDA
 - Prefix compiler and linker command with scorep --cuda
 - export SCOREP_CUDA_ENABLE=yes
 - yes refers to: runtime, kernel, memcpy
 - Full list of options in User Guide
- OpenCL similar (use SCOREP_OPENCL_ENABLE=yes)

EXAMPLE: OPENACC



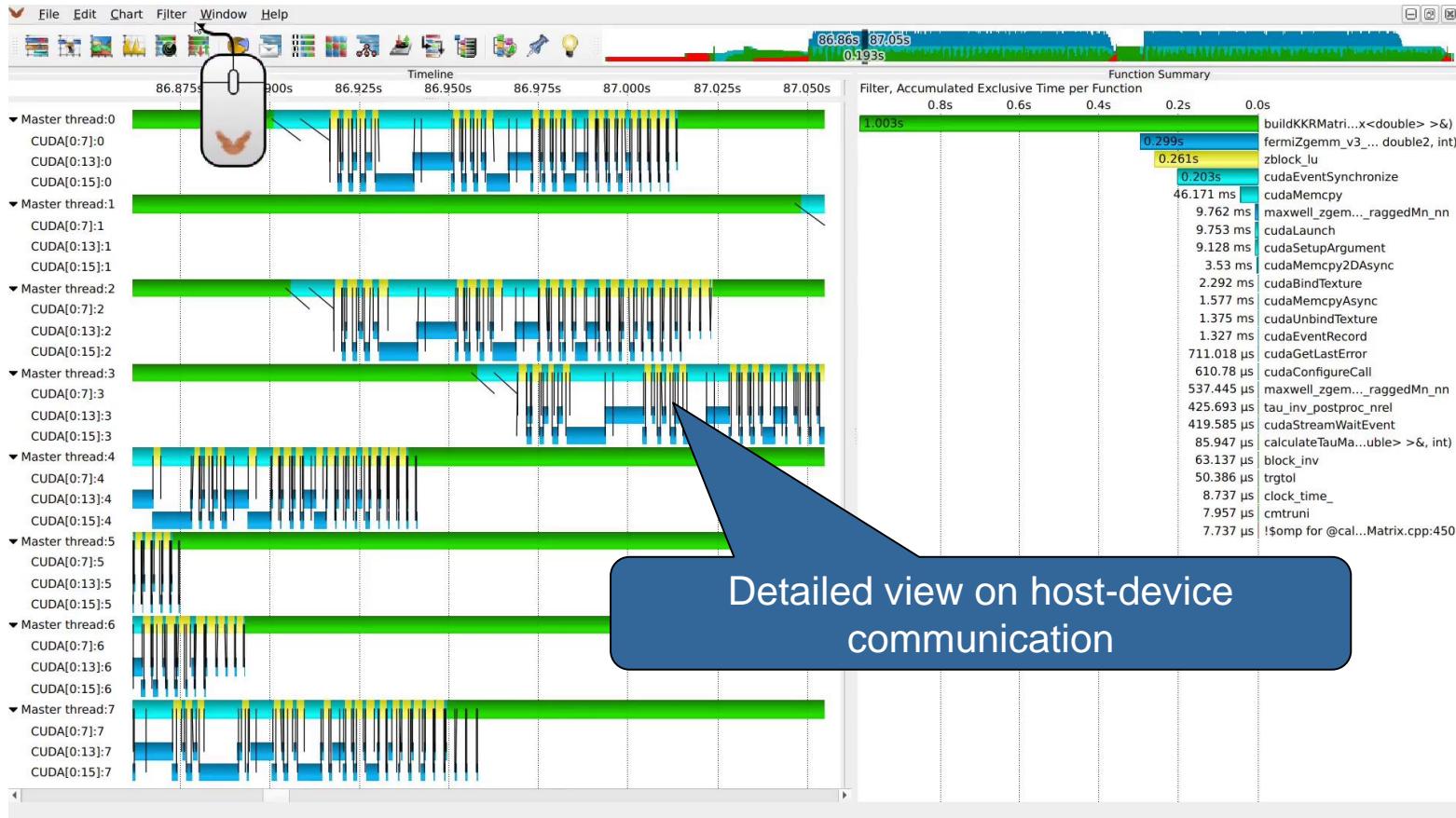
Pure OpenACC measurements contain host-side events only

EXAMPLE: OPENACC + CUDA



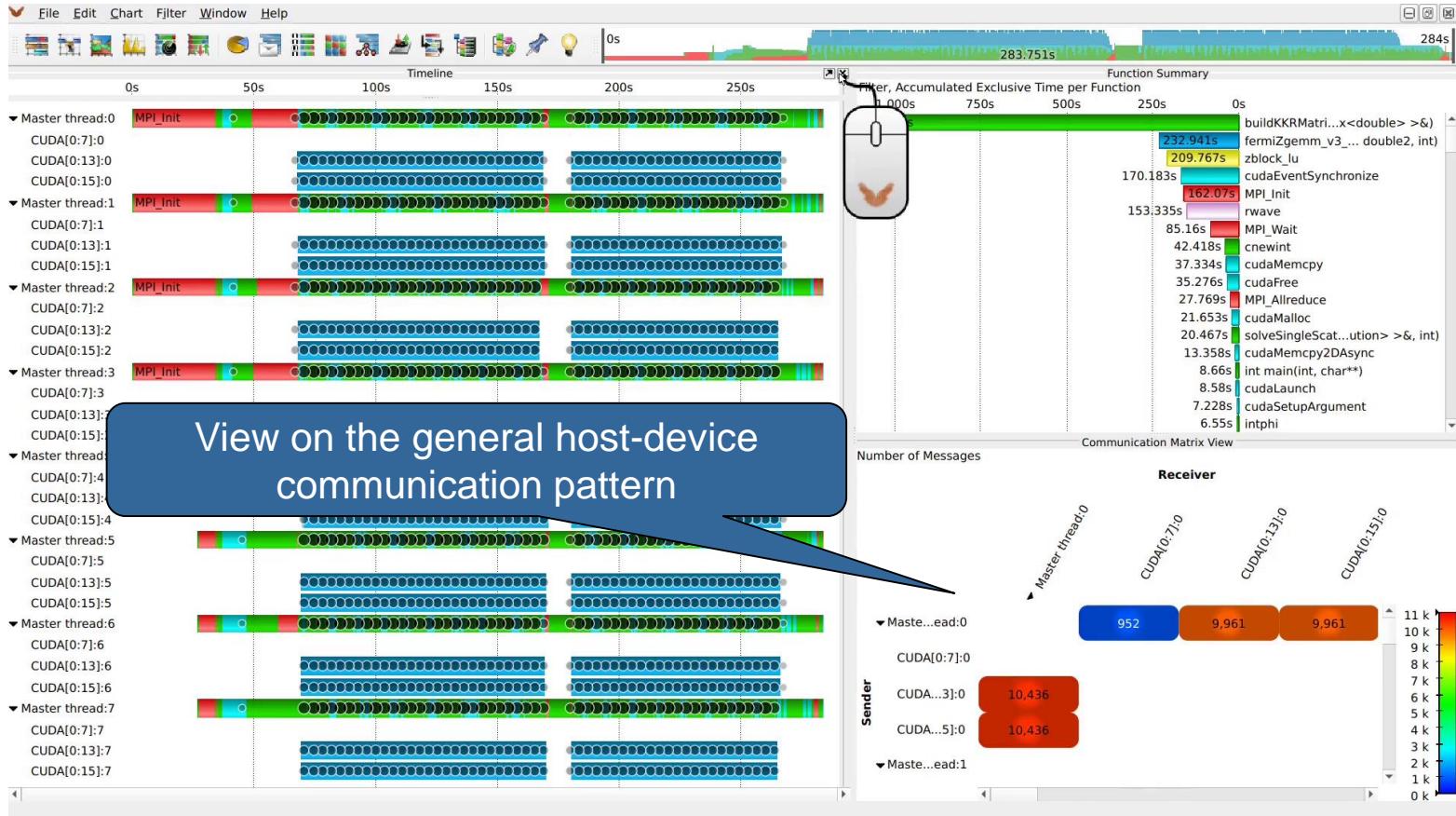
Enabling CUDA also shows kernels on the GPU

VAMPIR CUDA EXAMPLE



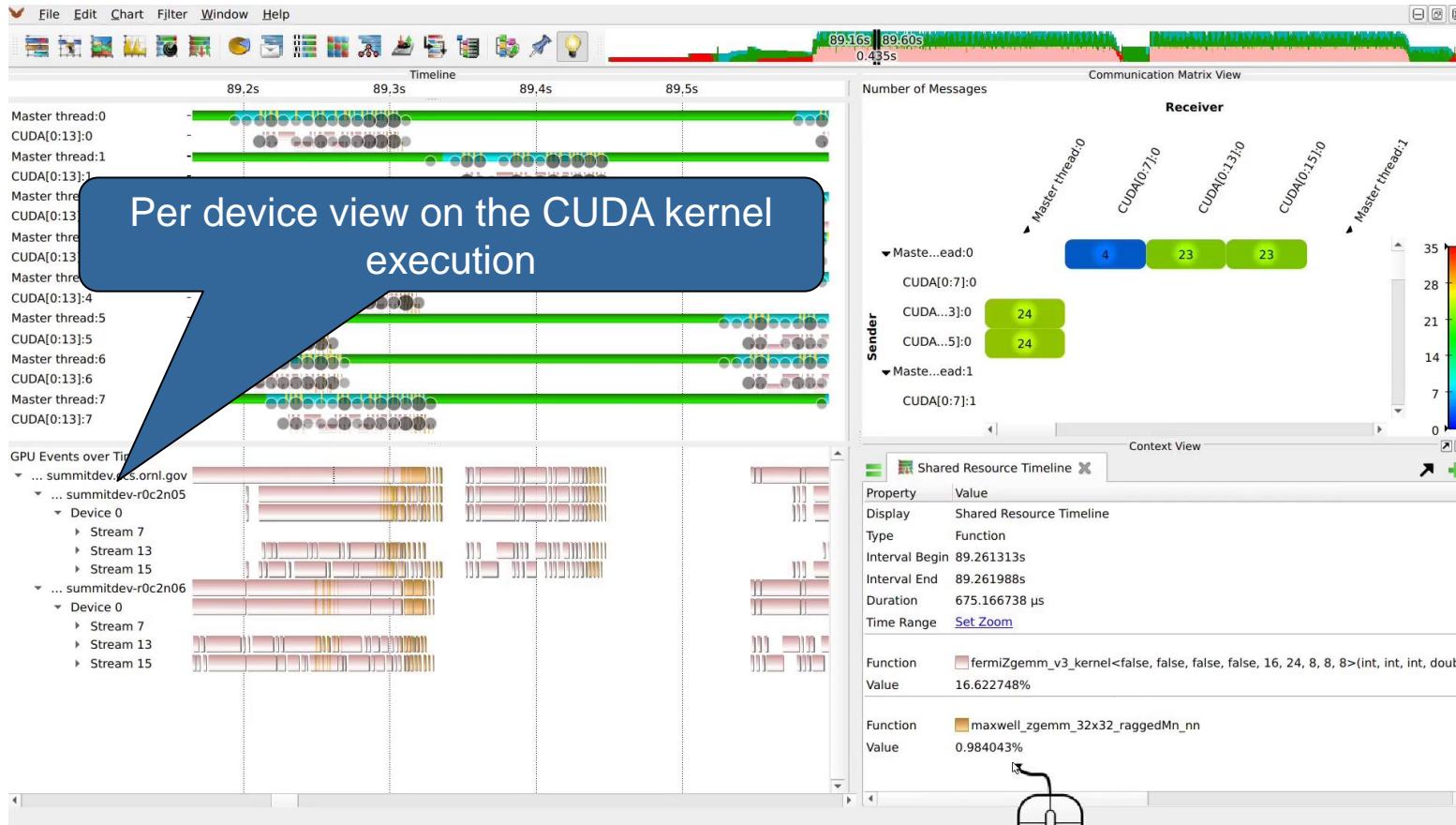
- Material science code LSMS
- CUDA is utilized for heavy computations
- CUDA streams are child's of the owning Process
- Allows an in-depth analysis of host-device communication

VAMPIR CUDA EXAMPLE



- Communication Matrix best for analyzing the general communication pattern
- Expectation: balanced communication, represented by a symmetric matrix
- Problem: communication with stream 7 is different

VAMPIR CUDA EXAMPLE



- Shared Resource Timeline offers a per device view on kernel executions
- Best suited for analyzing multi GPU per node scenarios
- Allows a dedicated analysis of kernel execution patterns
- Yields insights of the actual hardware usage

WHAT ABOUT SCALASCA AND GPUS?

- The Scalasca trace analyzer aborts if it encounters non-CPU locations, e.g. GPU streams
- However, GPU events are not too important for the analysis performed by Scalasca
- To enable a Scalasca analysis of a GPU application, either:
 - set SCOREP_CUDA_ENABLE=0 to disable all GPU events
 - set SCOREP_CUDA_ENABLE=runtime to get host-side events from the runtime

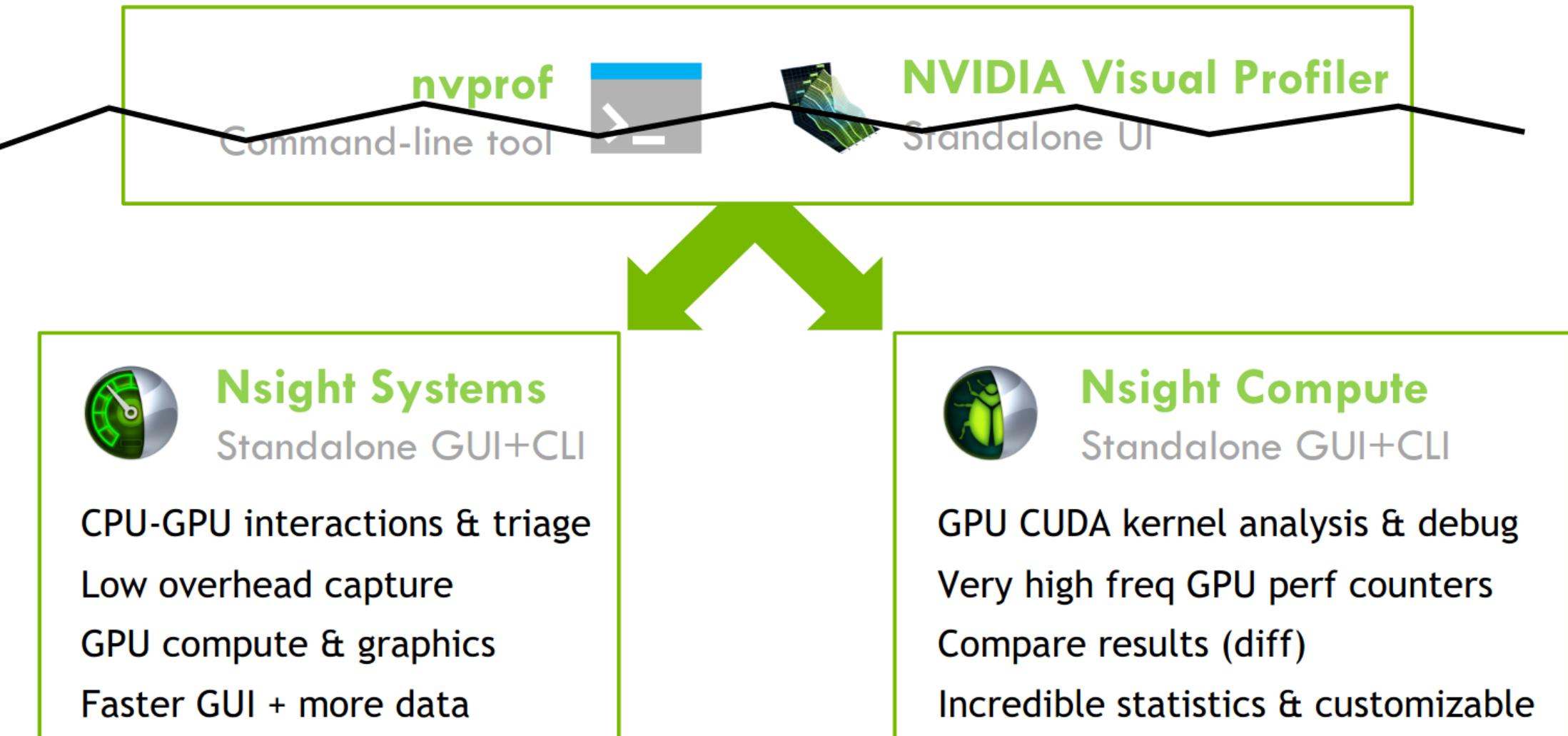
Attention:

- Due to the asynchronous execution of GPU kernels, some advanced Scalasca analysis can be misleading, e.g. Critical-path analysis
- Host side might be idle, while GPU can still be busy

GPU Analysis using **NVIDIA NSIGHT**

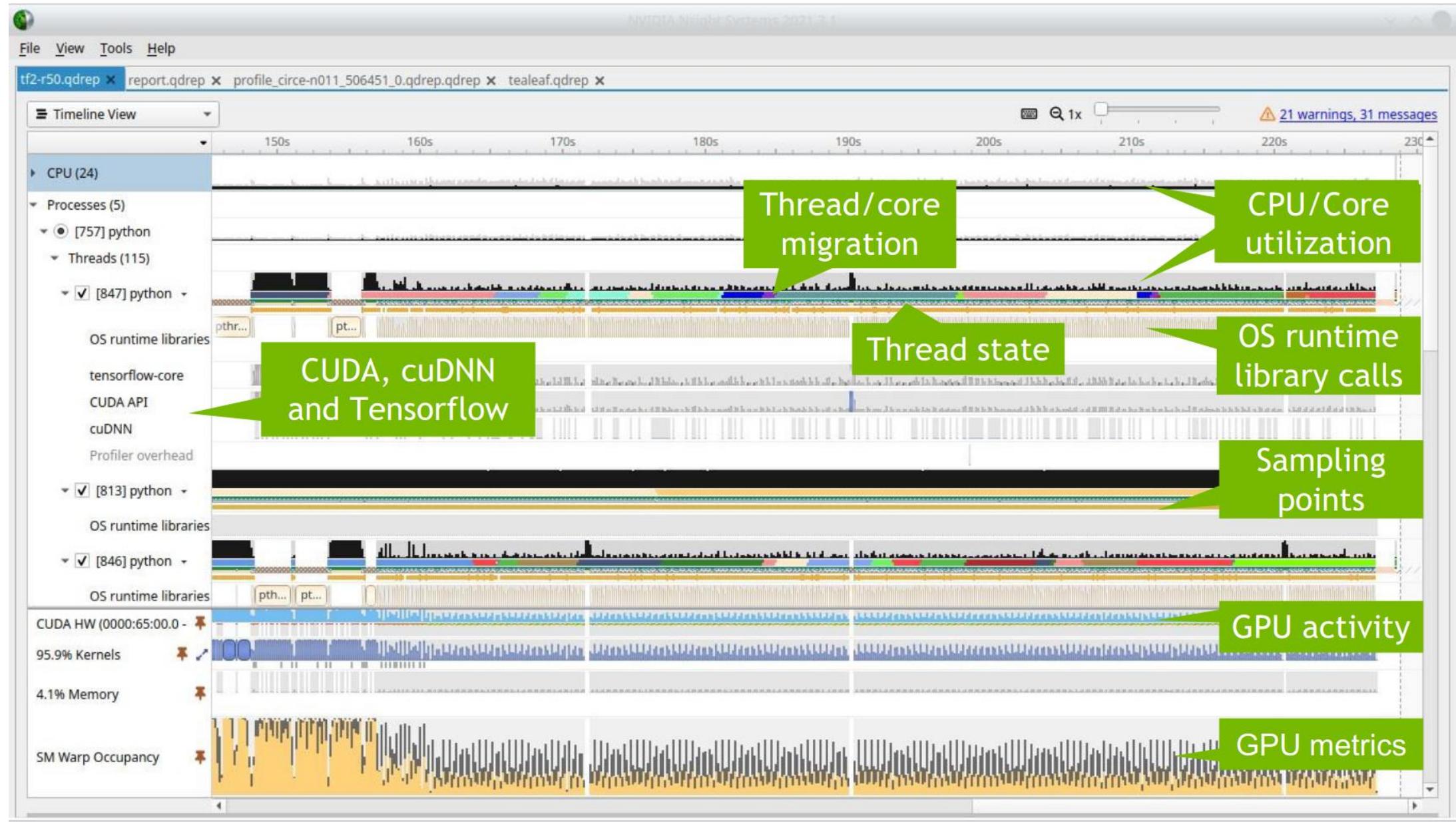
Mitglied der Helmholtz-Gemeinschaft

LEGACY TRANSITION



NSIGHT SYSTEM

- System-wide application tuning
- Locate optimization opportunities
 - Visualize millions of events on a timeline
 - See gaps of unused CPU and GPU time
- Balance workloads across multiple CPUs and GPUs
 - CPU utilization and thread state
 - GPU streams, kernels, memory transfers, etc.
- Multi-platform support
 - Linux, Windows and Mac OS X (host-only)
 - x86-64, Power9, ARM server, Tegra (Linux & QNX)



USAGE OF CLI

- Basic profiling session: `nsys profile ./app`
- Interactive session (scriptable): `nsys start|launch|stop|cancel`
- Useful flags:
 - API tracing: `-t, --trace={cuda,nvtx,mpi,openacc,openmp,none,...}`
 - Overwrite existing report: `-f, --force-overwrite=[true|false]`
 - Summary statistics (profile output on command line): `--stats=[true|false]`
 - Report file name: `-o, --output=report//pattern for hostname, PID, etc//`
 - Callstack sampling:
 - CUDA memory usage: `--cuda-memory-usage=[true|false]`

CLI PROFILING – MPI PROGRAMS

Single Node

```
nsys profile [nsys_args] mpirun [mpirun_args] your_executable
```

- ⇒ This creates one report file.

Multiple Nodes

```
mpirun [mpirun_args] nsys profile [nsys_args] your_executable
```

Set output report name with -o `report_name_%q{OMPI_COMM_WORLD_RANK}`
(for OpenMPI, `PMI_RANK` for MPICH and `SLURM_PROCID` for Slurm)

- ⇒ This creates one report file per MPI rank.

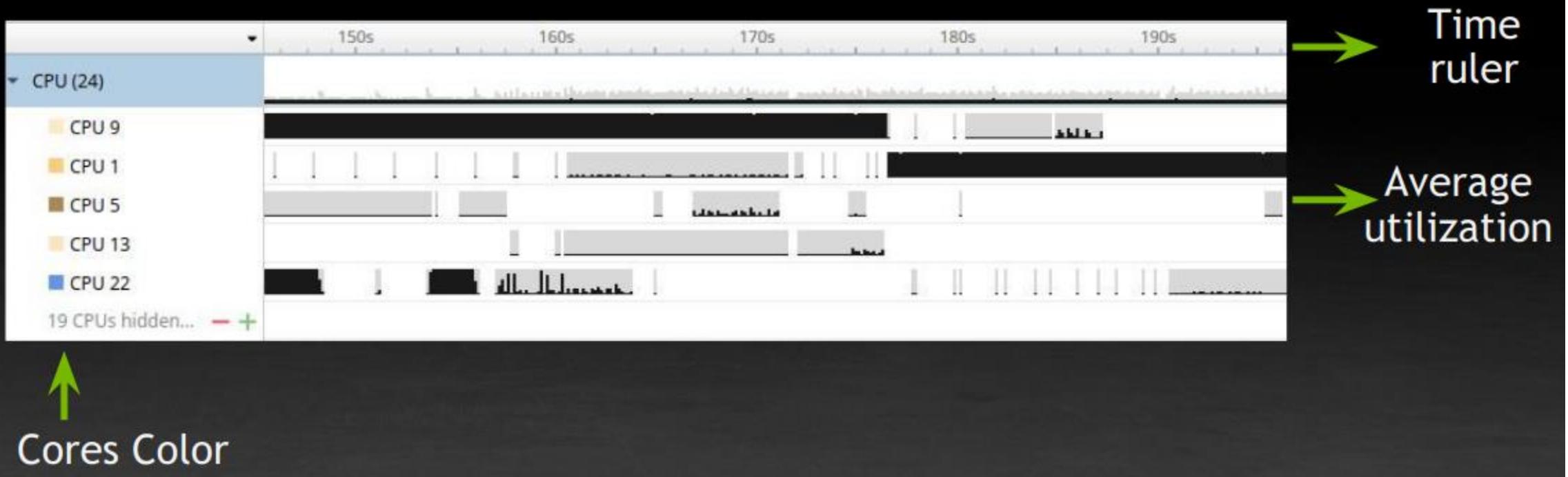
Profile only specific ranks. ⇒

```
#!/bin/bash
# OMPI_COMM_WORLD_LOCAL_RANK for node local rank
if [ $OMPI_COMM_WORLD_RANK -eq 0 ]; then
    nsys profile -t mpi "$@"
else
    "$@"
fi
```

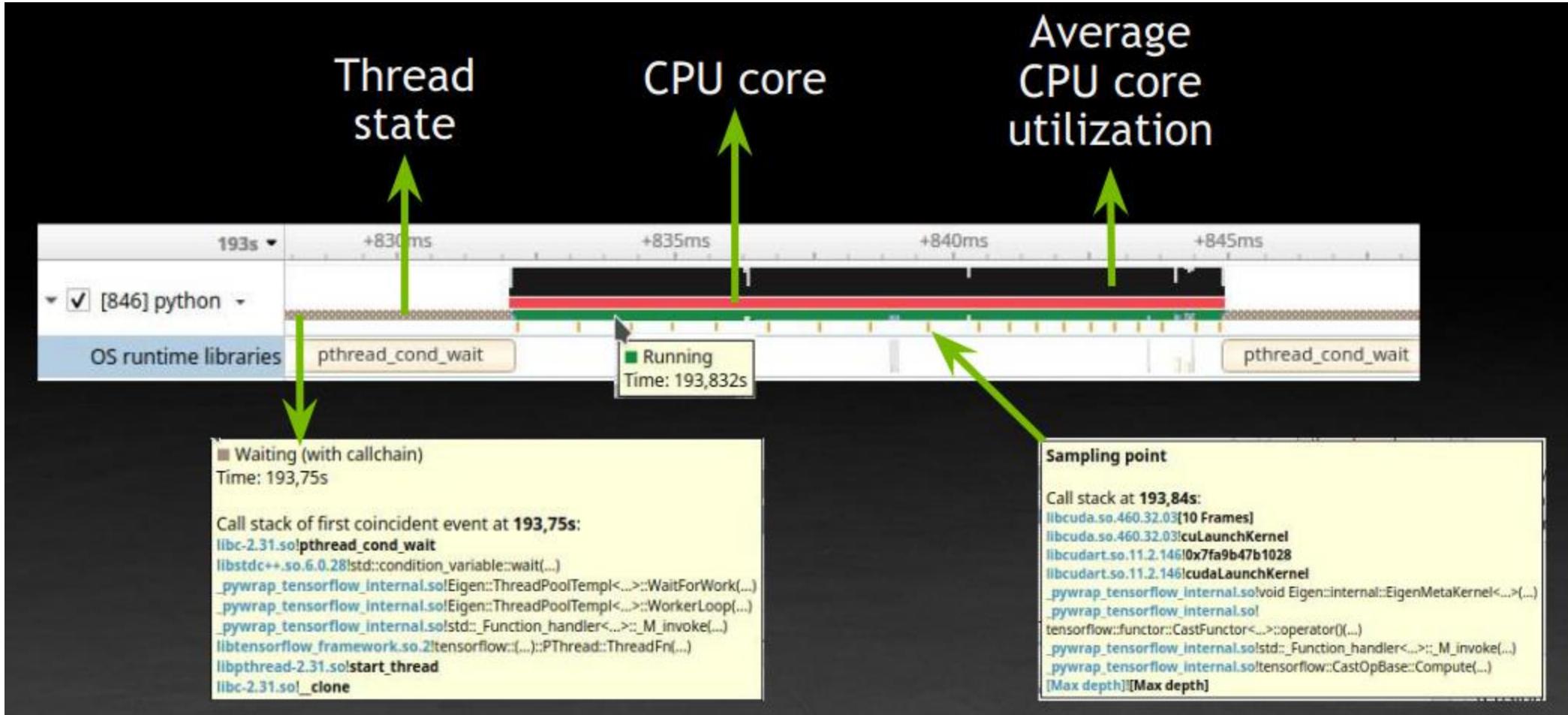
CPU CORES WORKLOAD

See CPU core utilization by application's threads

Locate idle time on CPU cores

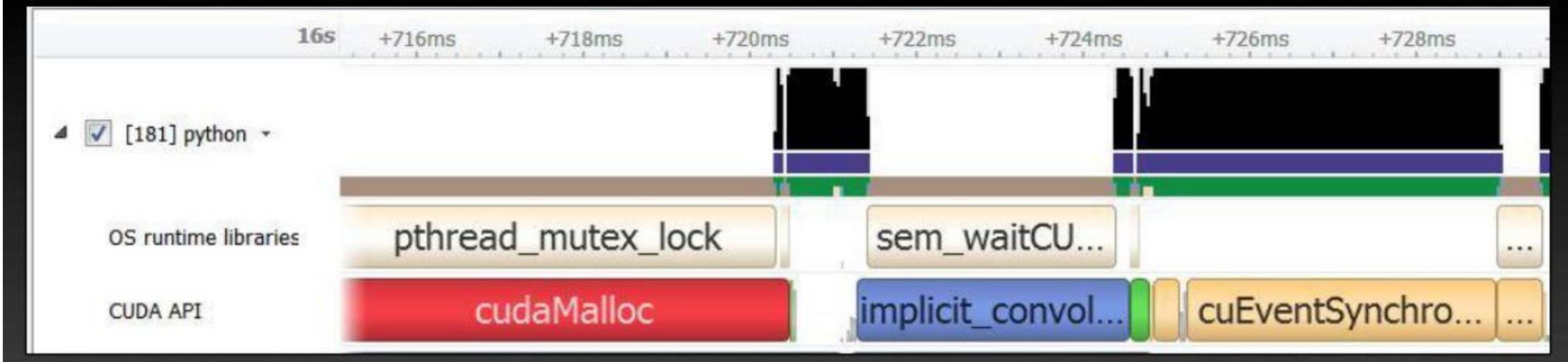


CPU THREAD ACTIVITY

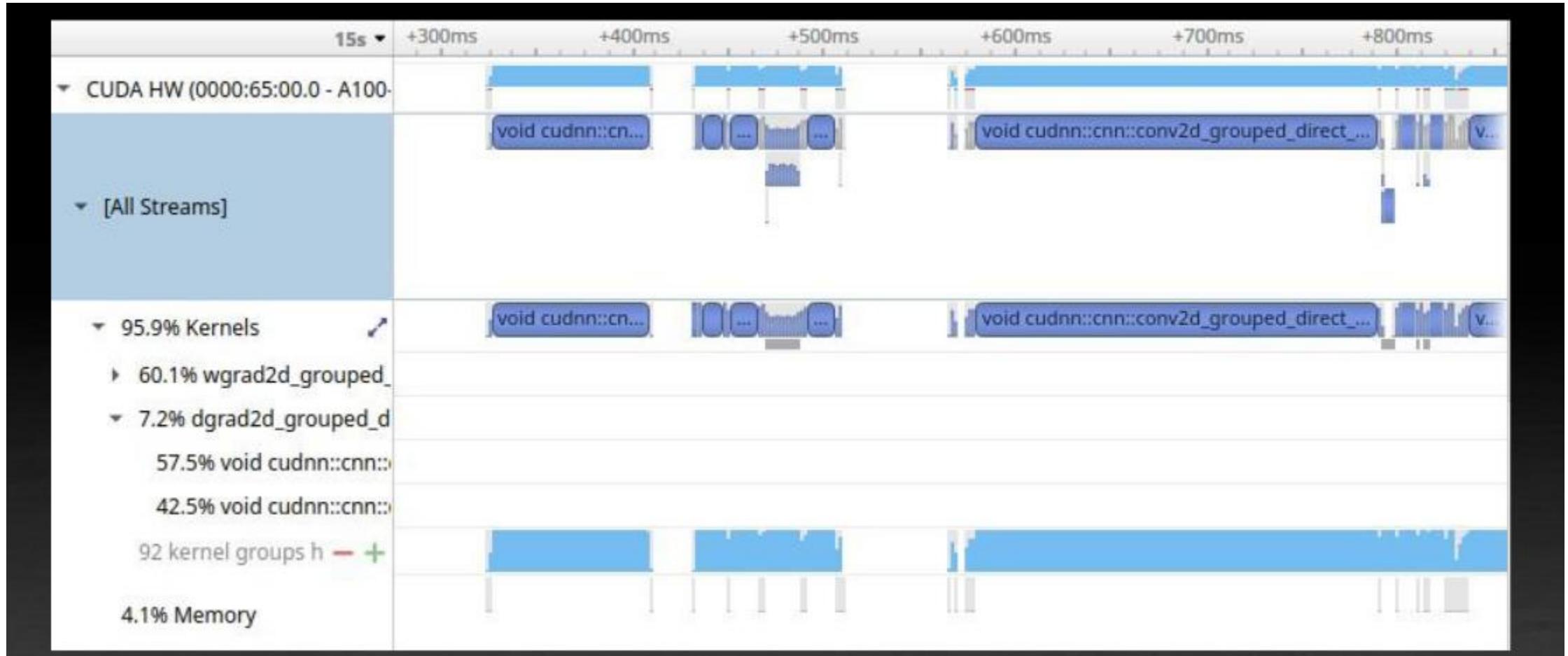


CUDA API

- Trace CUDA API Calls on OS thread
- See when kernels are dispatched
- See when memory operations are initiated
- Locate the corresponding CUDA workload on GPU



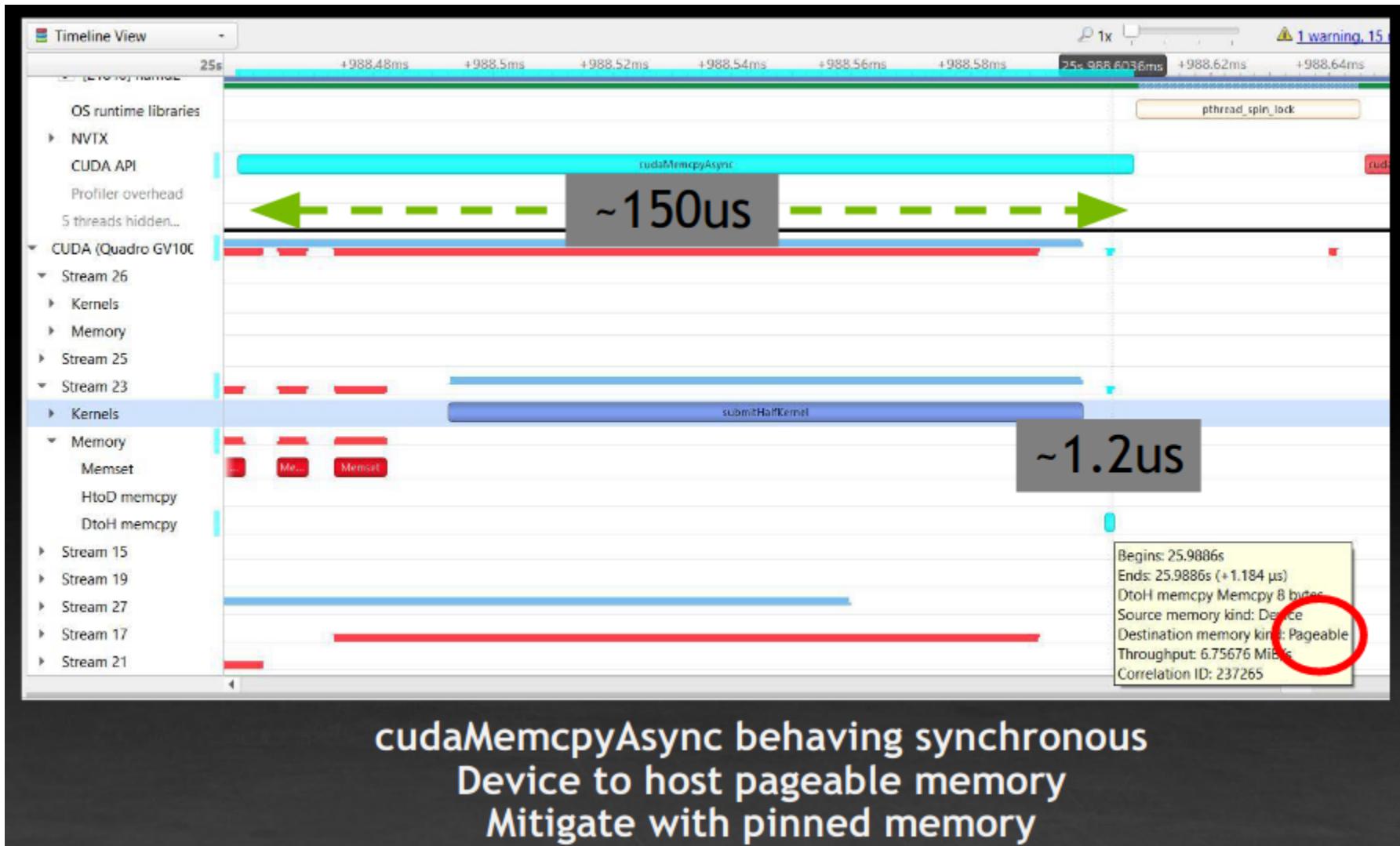
GPU UTILIZATION



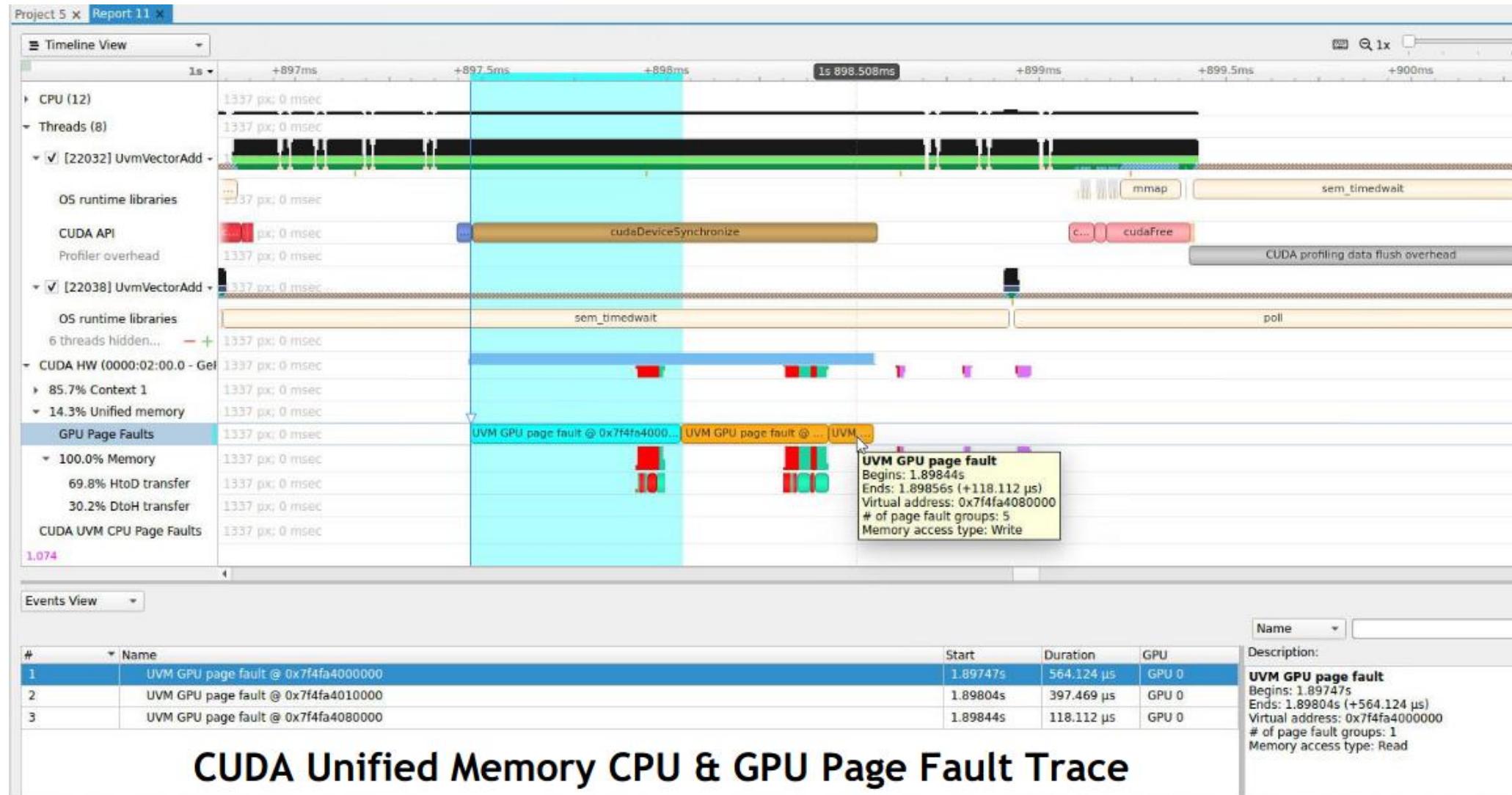
GPU UTILIZATION IN DETAIL



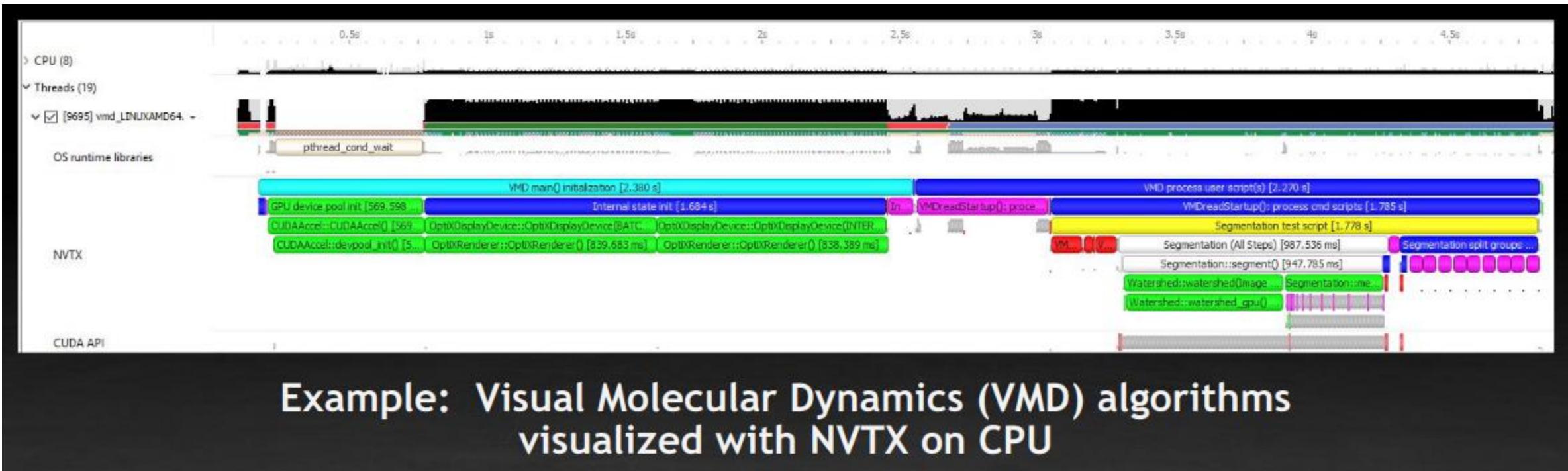
CUDA MEMORY TRANSFERS



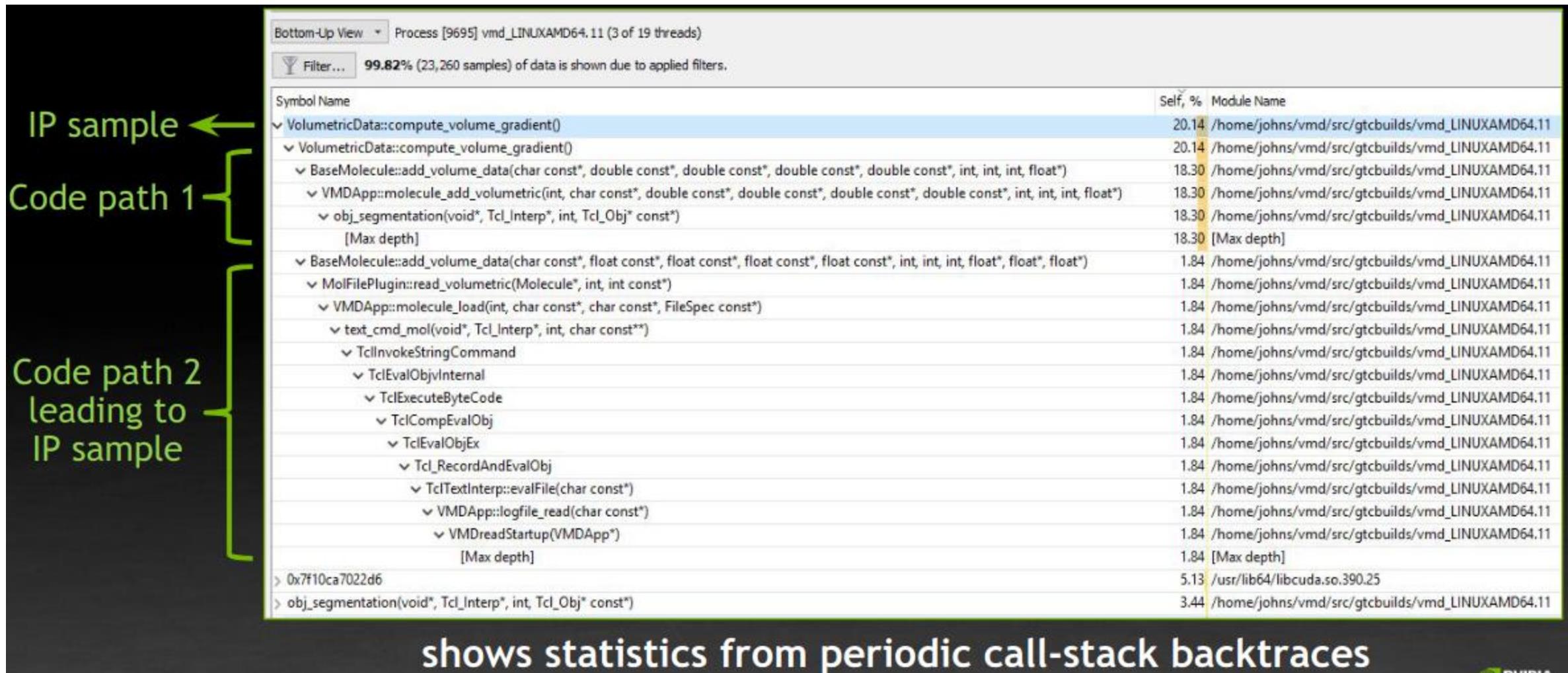
UNIFIED MEMORY ANALYSIS



USER ANNOTATIONS WITH NVTX



STATISTICAL SAMPLING

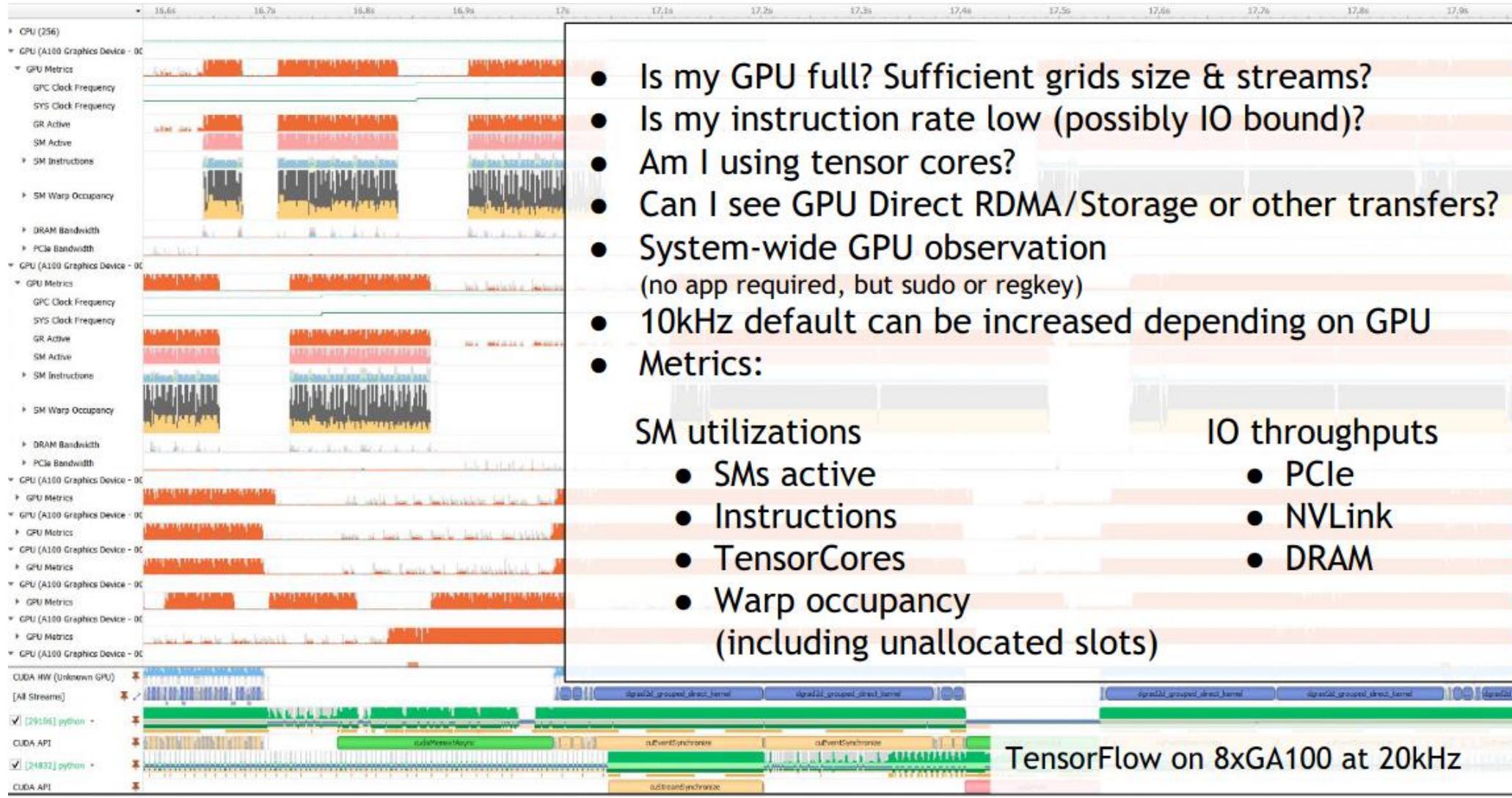


EVENTS VIEW

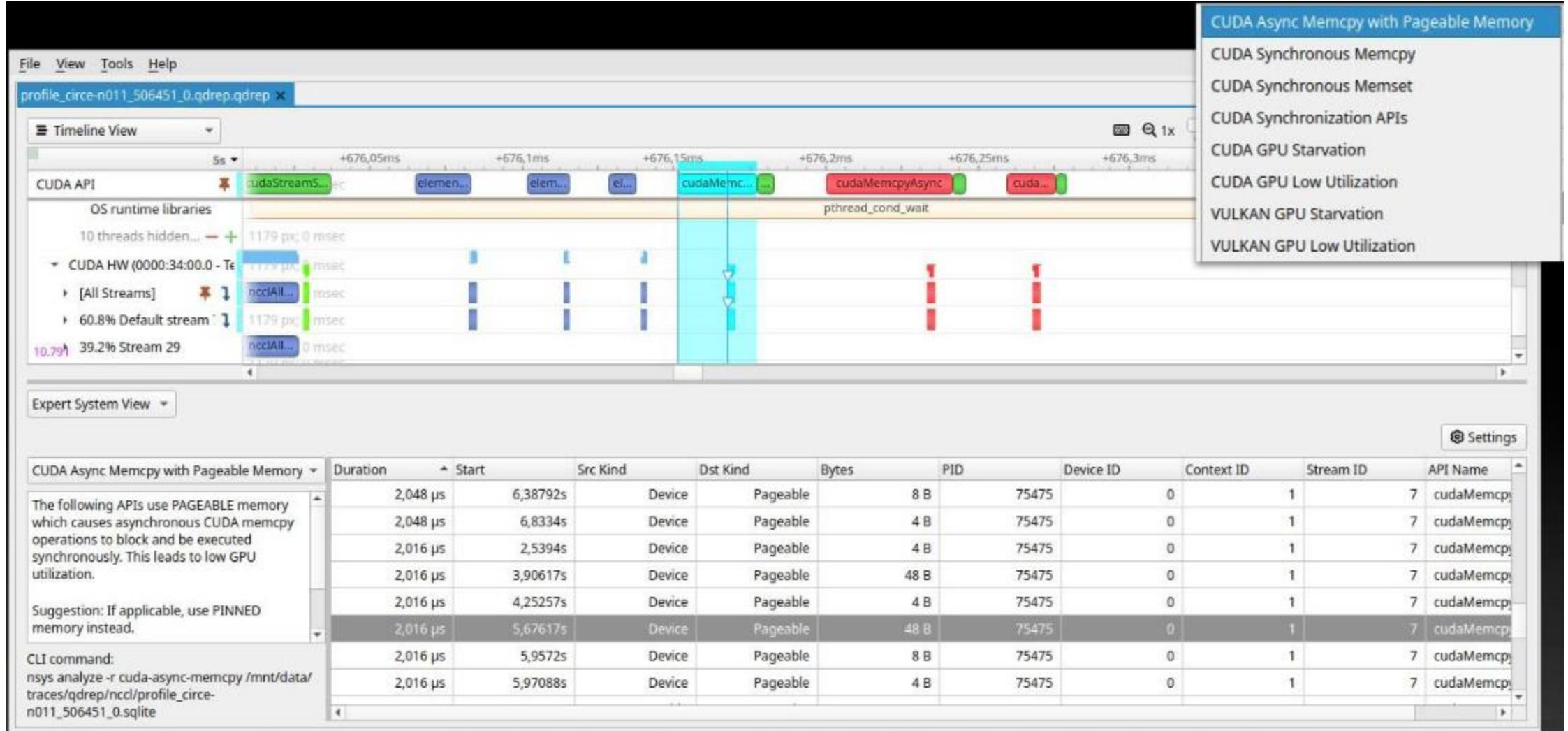
Events View ▾

#	Name	Start	Duration	TID	Description:
149153	device_tea_leaf_ppcg_solve_calc_sd_new	11,4088s	3,265 µs	390019	
149154	device_pack_bottom_buffer	11,4088s	3,247 µs	390019	
149155	cudaMemcpy	11,4088s	1,838 ms	390019	Call to cudaMemcpy Memory copies Begins: 11,4088s Ends: 11,4106s (+1,838 ms) Return value: 0 Correlation ID: 40146
149174	└ MPI_Isend	11,4106s	1,564 µs	390019	
149175	└ MPI_Irecv	11,4106s	1,420 µs	390019	
149176	└ MPI_Waitall	11,4106s	1,772 ms	390019	
149199	└ cudaMemcpy	11,4124s	77,788 µs	390019	
149201	device_unpack_bottom_buffer	11,4125s	5,232 µs	390019	
149202	device_tea_leaf_ppcg_solve_update_r	11,4125s	3,334 µs	390019	
149203	device_tea_leaf_ppcg_solve_calc_sd_new	11,4125s	3,268 µs	390019	
149204	device_pack_bottom_buffer	11,4125s	3,031 µs	390019	
149205	└ cudaMemcpy	11,4125s	1,853 ms	390019	
149224	└ MPI_Isend	11,4143s	2,081 µs	390019	
149225	└ MPI_Irecv	11,4143s	1,191 µs	390019	
149226	└ MPI_Waitall	11,4143s	1,786 ms	390019	
149227	libucp.so.0.0.0!lucp_worker_progress	11,4144s	-	390019	
149228	libucp.so.0.0.0!lucp_worker_progress	11,4144s	-	390019	
149229	libopen-pal.so.40.20.5!opal_progress	11,4145s	-	390019	
149230	libuct.so.0.0.0!uct_mm_iface_progress	11,4146s	-	390019	
149231	libonnen-nal.so.40.20.5!nal_nprogress	11,4147s	-	390019	

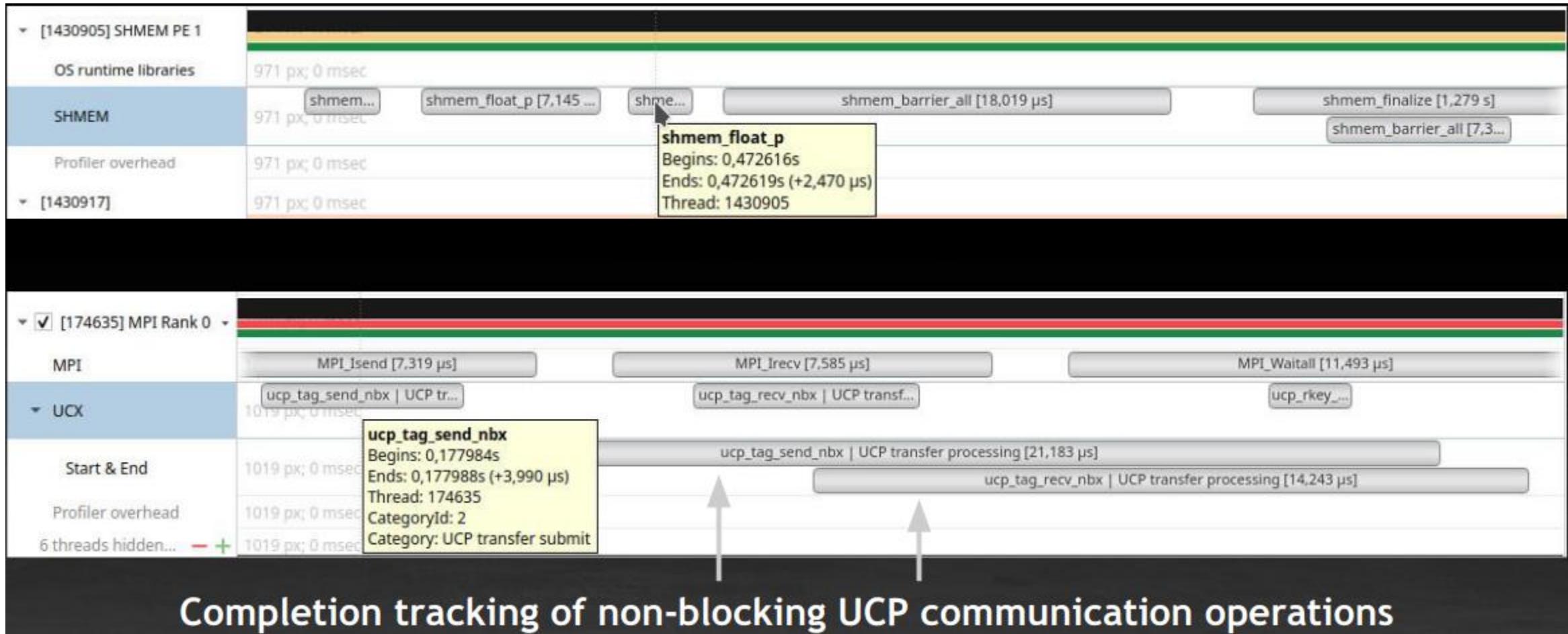
GPU METRIC SAMPLING



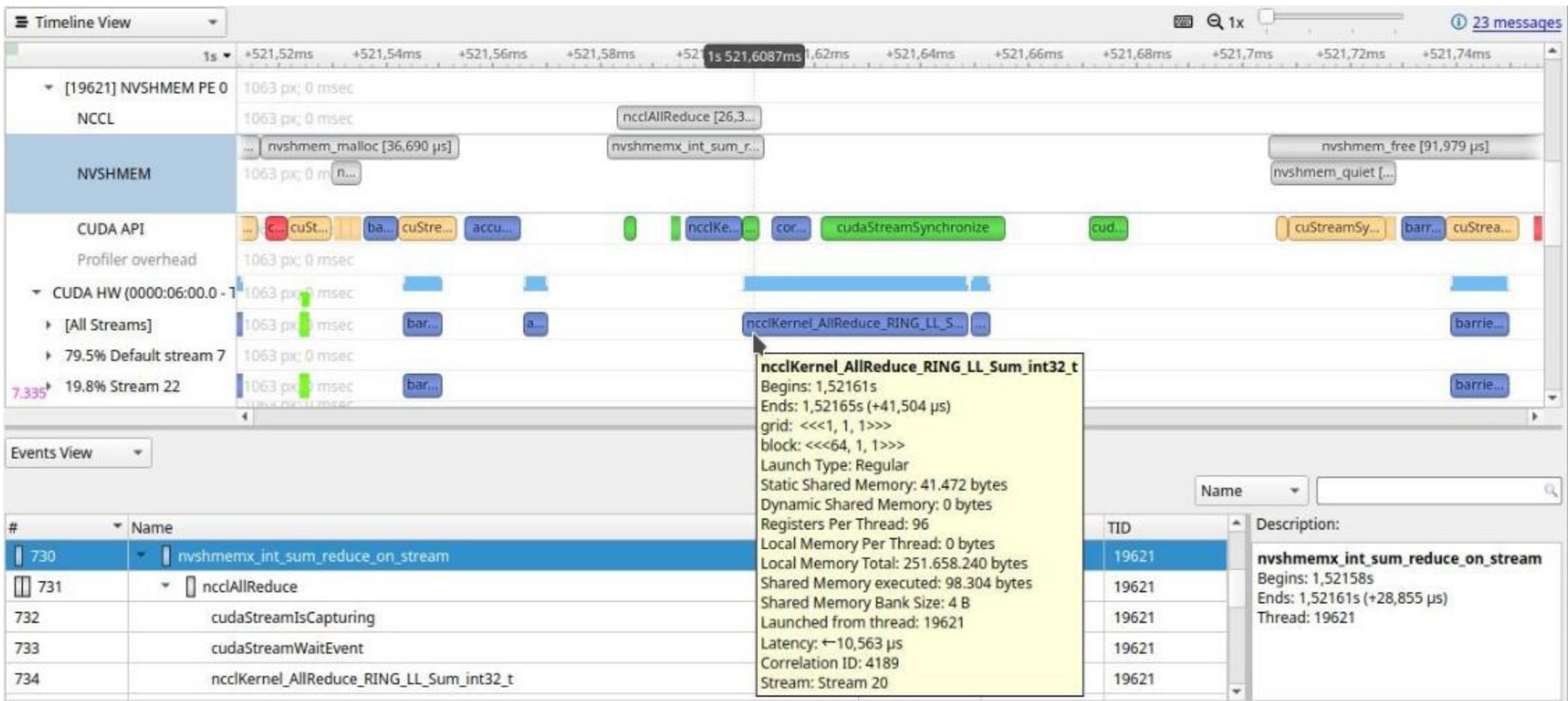
EXPERT SYSTEM



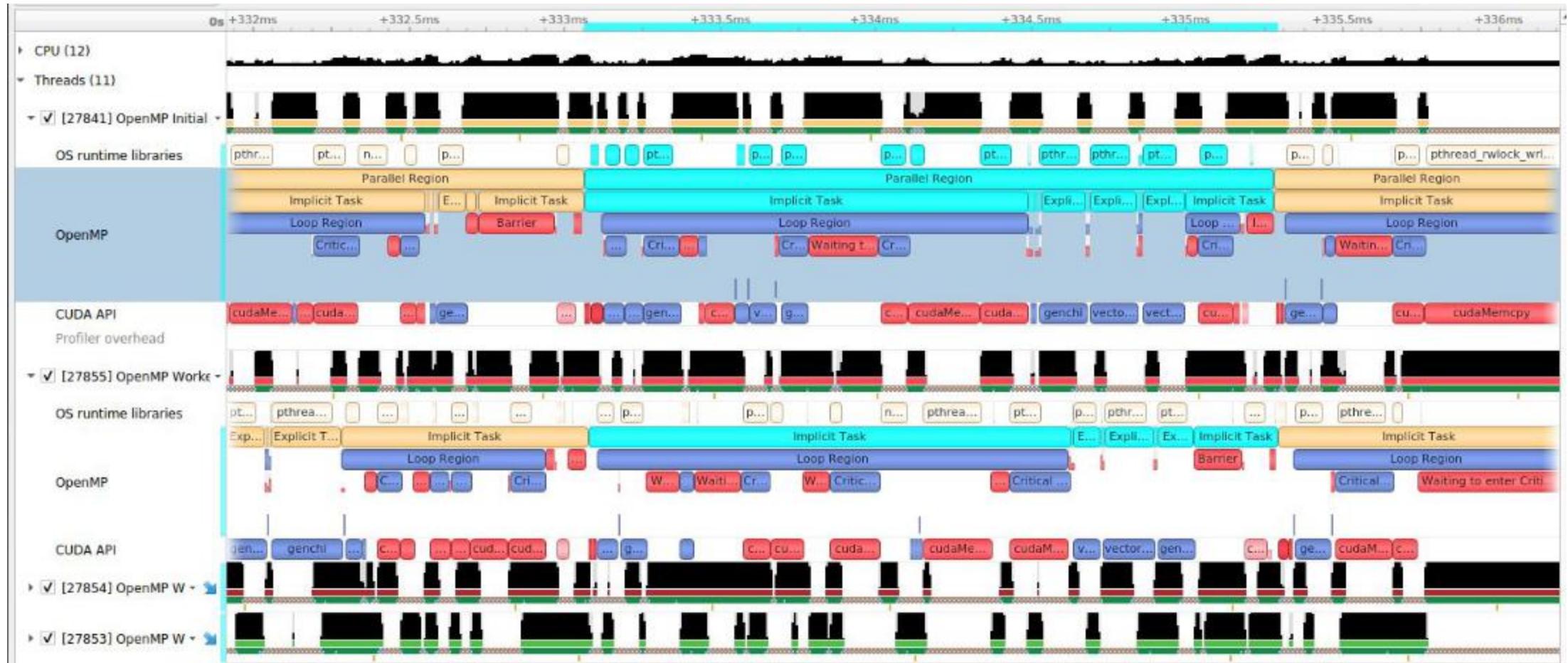
SHMEM, MPI AND UCX



NVSHMEM AND NCCL

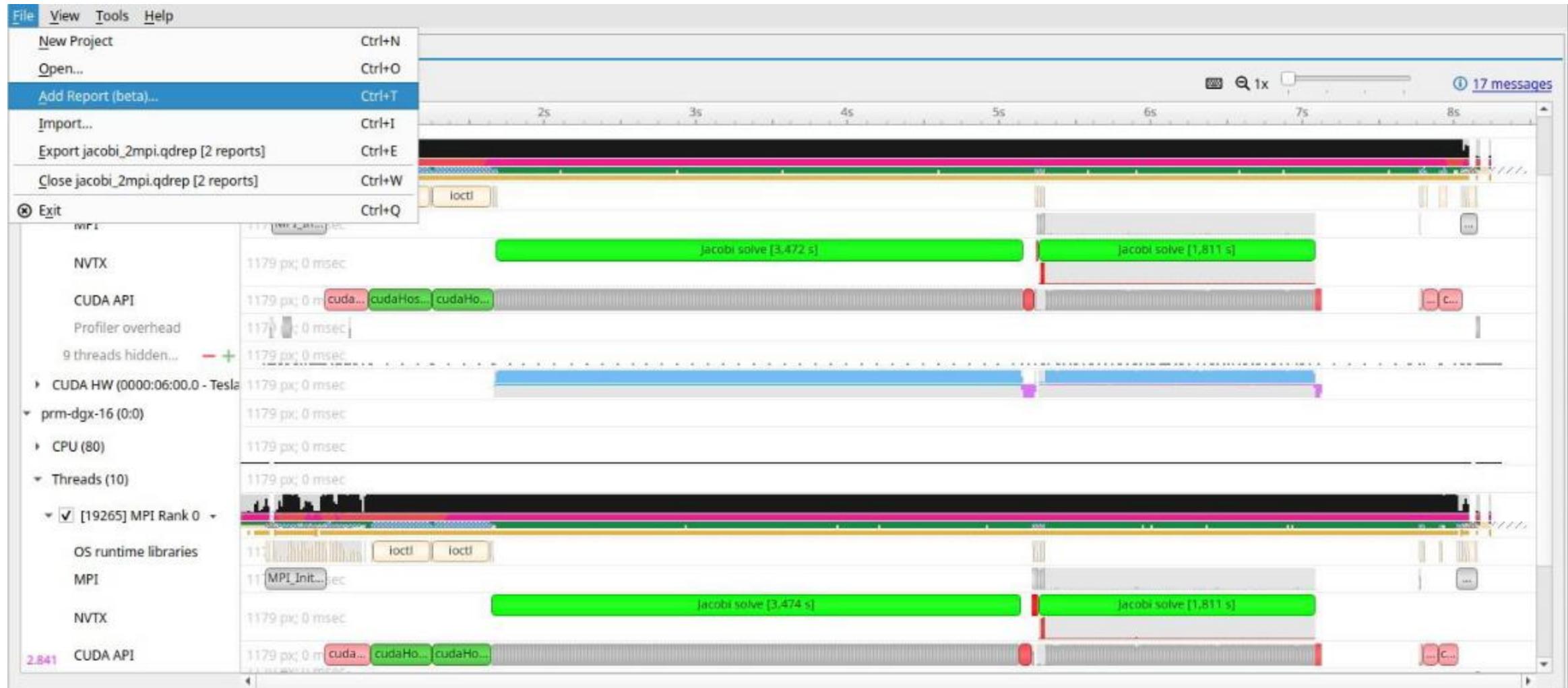


OPENMP



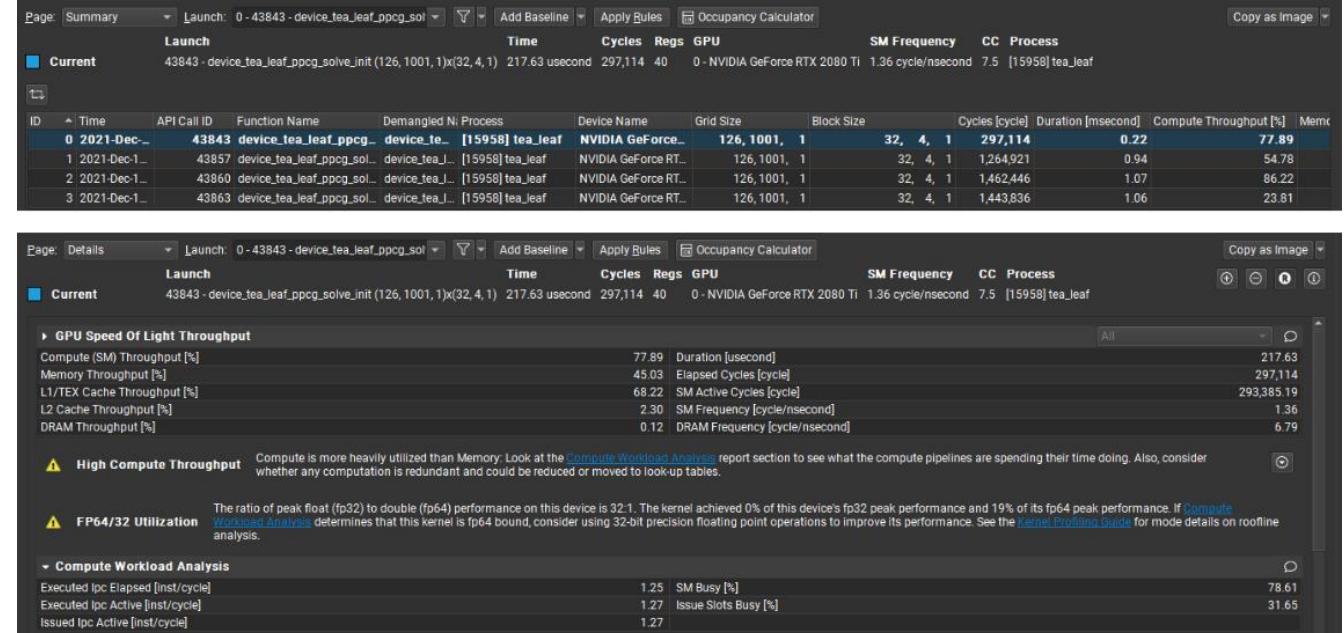
OMPT-capable OpenMP runtime required

LOADING MULTIPLE REPORTS



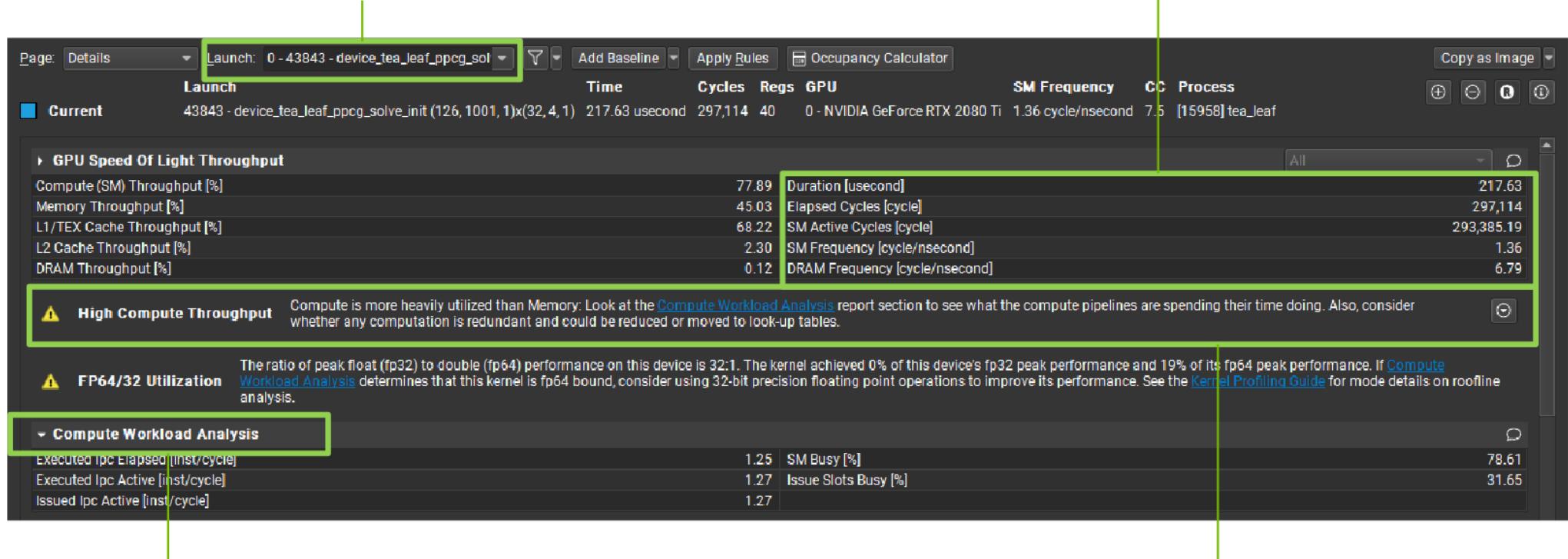
NSIGHT COMPUTE

- Interactive CUDA kernel profiler
- Targeted metric sections for various performance aspects
- Customizable data collection and presentation (tables, charts, ...)
- GUI and CLI
- Python-based API for guided analysis and post-processing
- Support for remote profiling across machines and platforms



PROFILER REPORT

Selected result



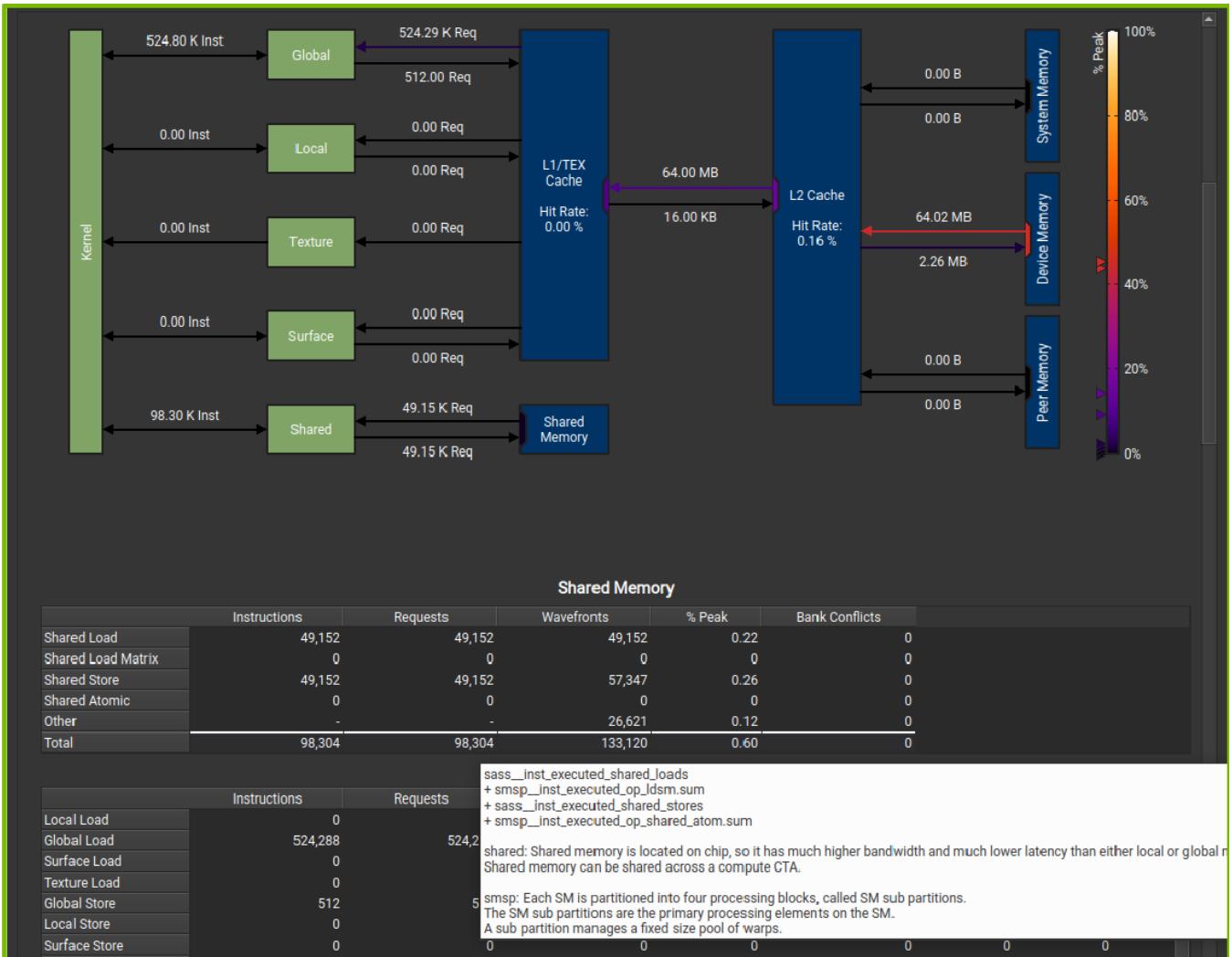
Expandable
Sections

Metric values

Expert Analysis
(Rules)

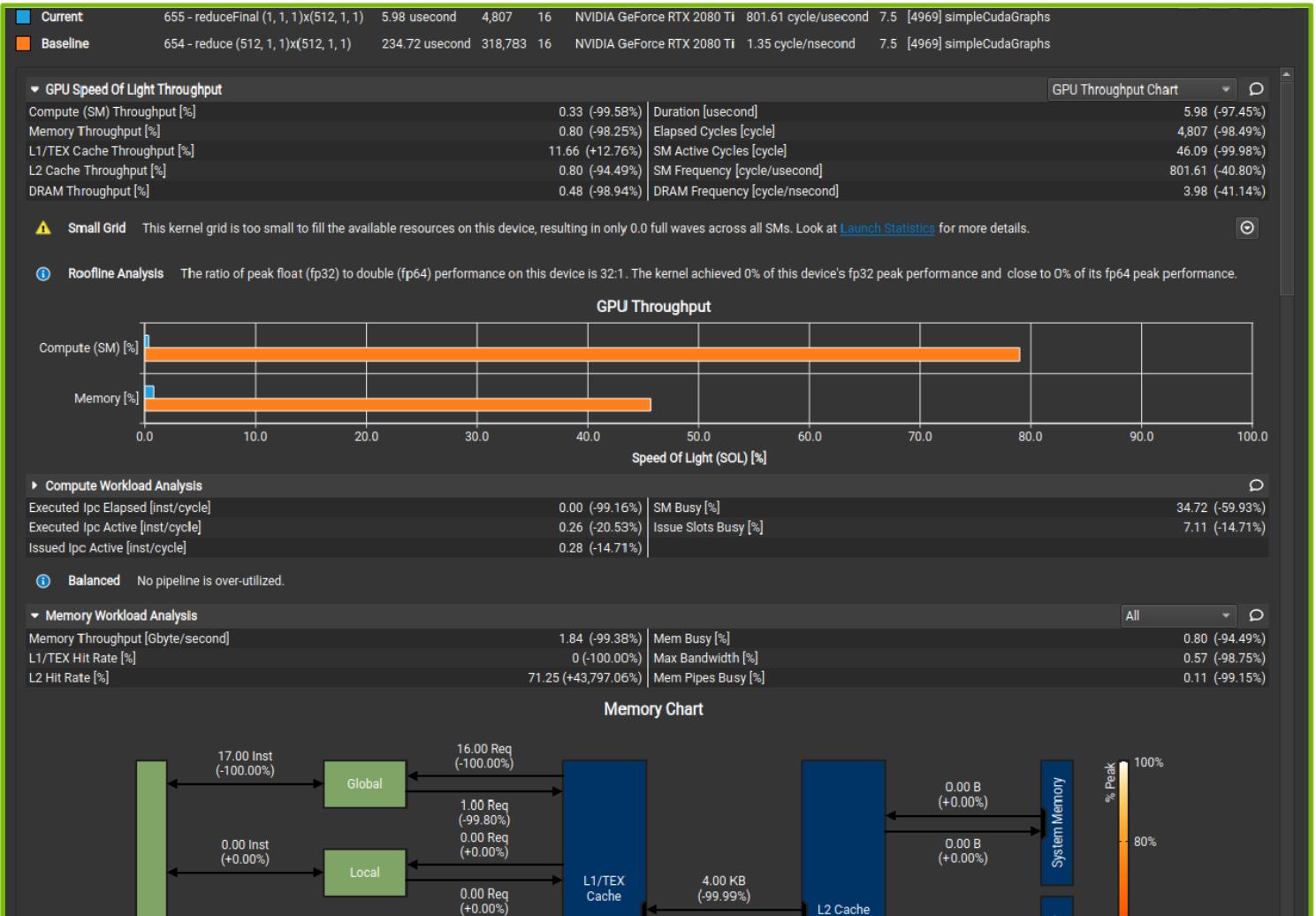
DATA TRANSFER ANALYSIS

- Detailed memory workload analysis chart and tables
- Shows transferred data or throughputs
- Tooltips provide metric names, calculation formulas and detailed background info



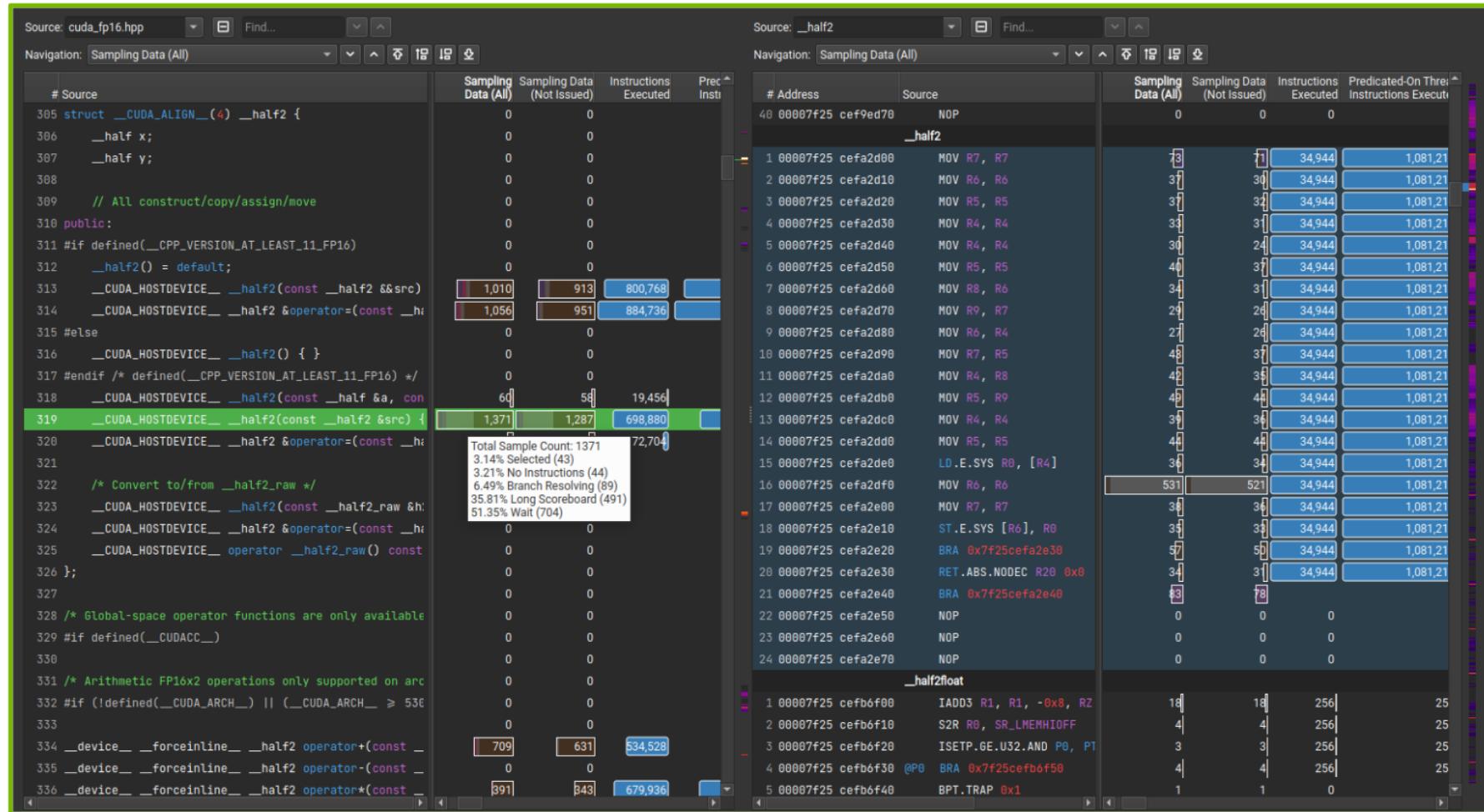
BASELINE COMPARISON

- Comparison of results directly within the tool with "Baselines"
- Supported across kernels, reports, and GPU architectures



SOURCE CORRELATION

- Source/PTX/SASS analysis and correlation
- Source metrics per instruction and aggregated (e.g. PC sampling data)
- Metric heatmap



OCCUPANCY CALCULATOR

- In-depth analysis of GPU occupancy
- Quickly determine if computational resources are wasted

The screenshot shows the Occupancy Calculator interface. At the top, there are configuration dropdowns for Compute Capability (7.5), Shared Memory Size Config (65536 bytes), CUDA version (11.0), and Global Load Cache Mode (L1+L2 (ca)). Below these are three tabs: Tables (selected), Graphs, and GPU Data.

Occupancy Data:

Property	Value
Active Threads per Multiprocessor	1024
Active Warps per Multiprocessor	32
Active Thread Blocks per Multiprocessor	8
Occupancy of each Multiprocessor	100 %

Allocated Resources:

Resources	Per Block	Limit Per SM	Allocatable Blocks Per SM
Warp (Threads Per Block / Threads Per Warp)	4	32	8
Registers (Warp limit per SM due to per-warp reg count)	4	48	12
Shared Memory (Bytes)	3072	65536	21

Occupancy Limiters:

Limited By	Blocks per SM	Warps Per Block	Warps Per SM
Max Warps or Max Blocks per Multiprocessor	8	4	32
Registers per Multiprocessor	12		
Shared Memory per Multiprocessor	21		

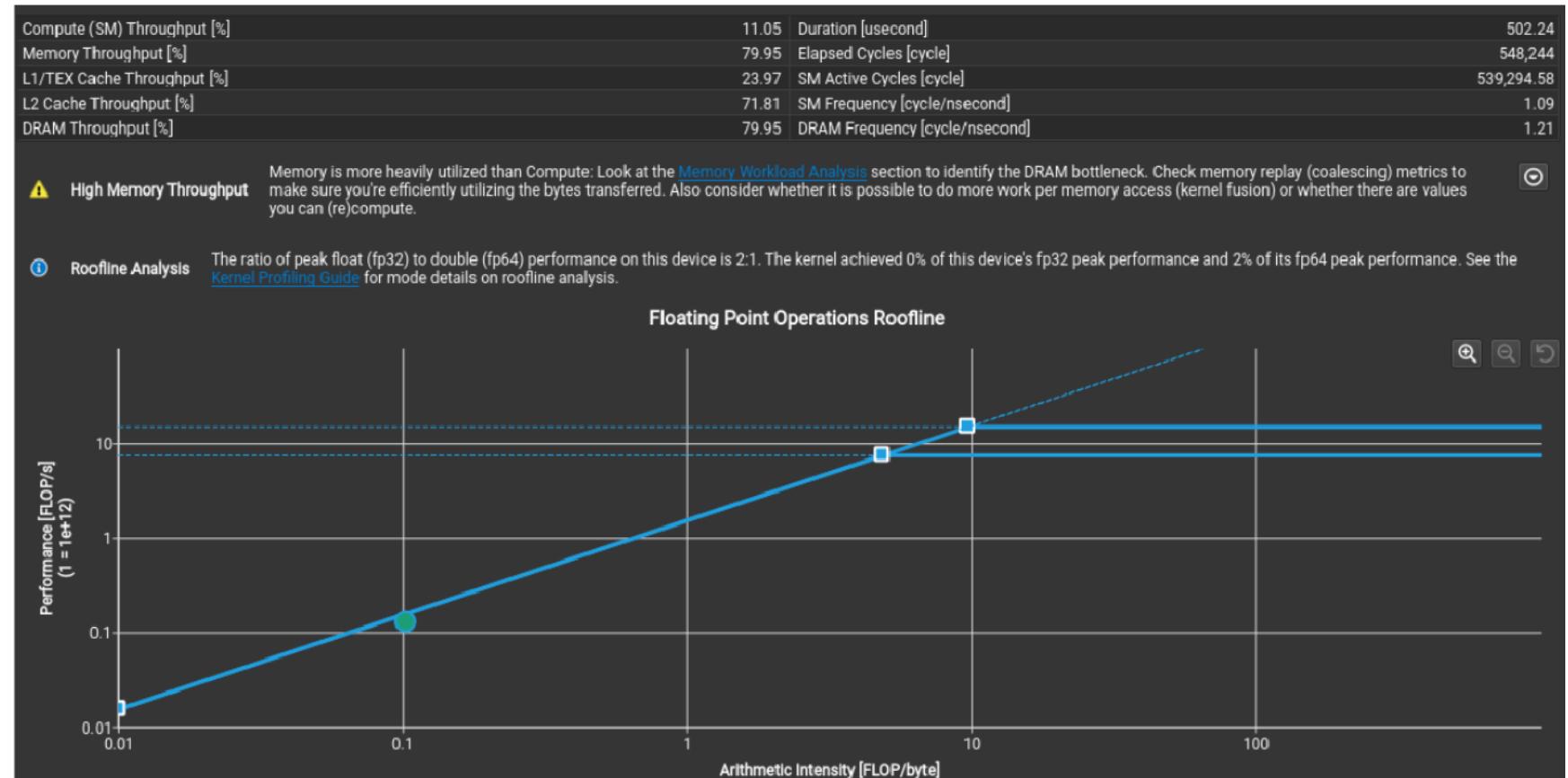
Physical Limit of GPU (7.5):

Property	Limit
Threads per Warp	32
Max Warps per Multiprocessor	32
Max Thread Blocks per Multiprocessor	16
Max Threads per Multiprocessor	1024
Maximum Thread Block Size	1024
Registers per Multiprocessor	65536
Max Registers per Thread Block	65536
Max Registers per Thread	255
Shared Memory per Multiprocessor (bytes)	65536
Max Shared Memory per Block	65536
Register Allocation Unit Size	256
Register Allocation Granularity	warp
Shared Memory Allocation Unit Size	256
Warp Allocation Granularity	4
Shared Memory Per Block (bytes) (CUDA runtime use)	0

A tooltip message "The occupancy is limited by block size" appears over the "Blocks per SM" column of the Occupancy Limiters table.

ROOFLINE ANALYSIS

- Determine whether the application is memory bound or compute bound
- Guided analysis points to detailed analysis of the most severe problem



QUESTIONS?



- <http://www.scalasca.org>
- scalasca@fz-juelich.de



- <http://www.score-p.org>
- support@score-p.org

