



S62388

Achieving Higher Performance

From Data Center & Cloud Applications

Daniel Horowitz, NVIDIA, Senior Director of Engineering

Ankur Srivastava, Amazon, AWS Senior Solutions Architect



Agenda

Achieving Higher Performance From Data Center & Cloud Applications

Builds on GTC 2023 S51421 - Optimizing at Scale

Investigating Hidden Bottlenecks for Multi-Node Workloads

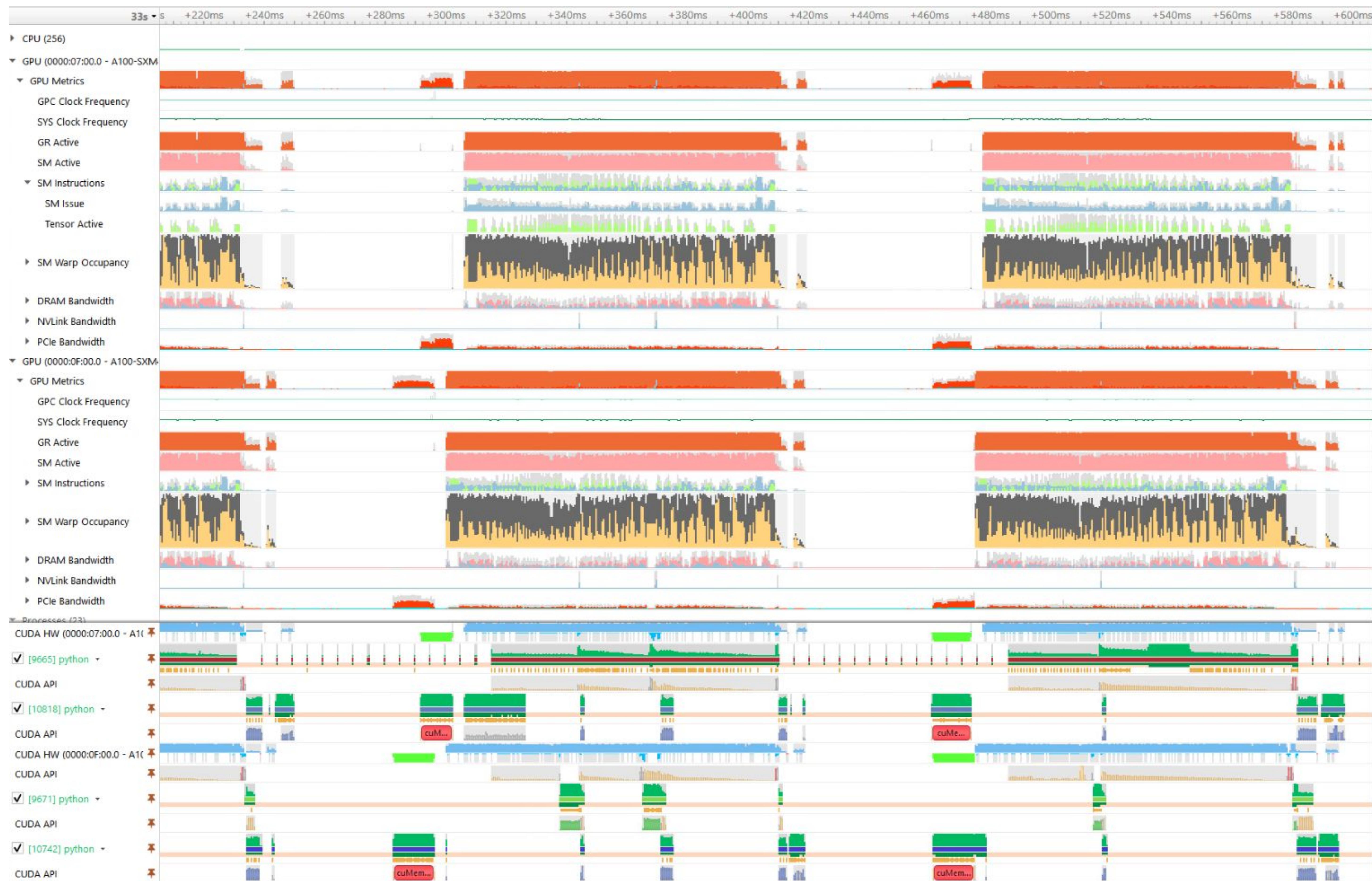
- Intro/Recap
- Data Centers
- Cloud & AWS
- Remote Servers

For those not familiar with...



Deep Learning Training Example

Zoomed Out - Graph Mode



CPU cores

GPU metrics samples

SMs active

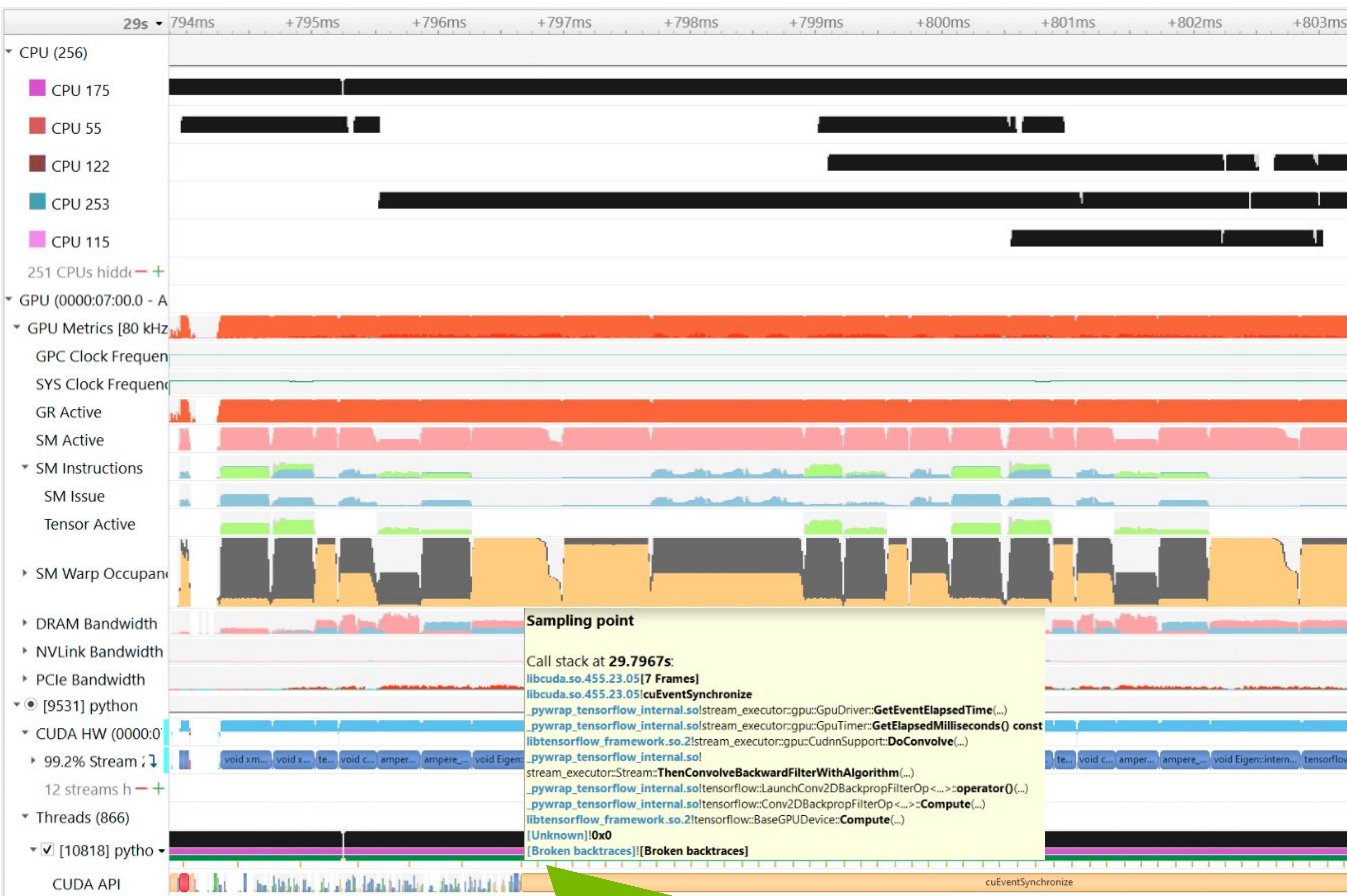
SM instructions
TensorCores

GPU bandwidths

GPU CUDA
kernel & memory
activities

CPU processes
& thread states

CUDA API trace

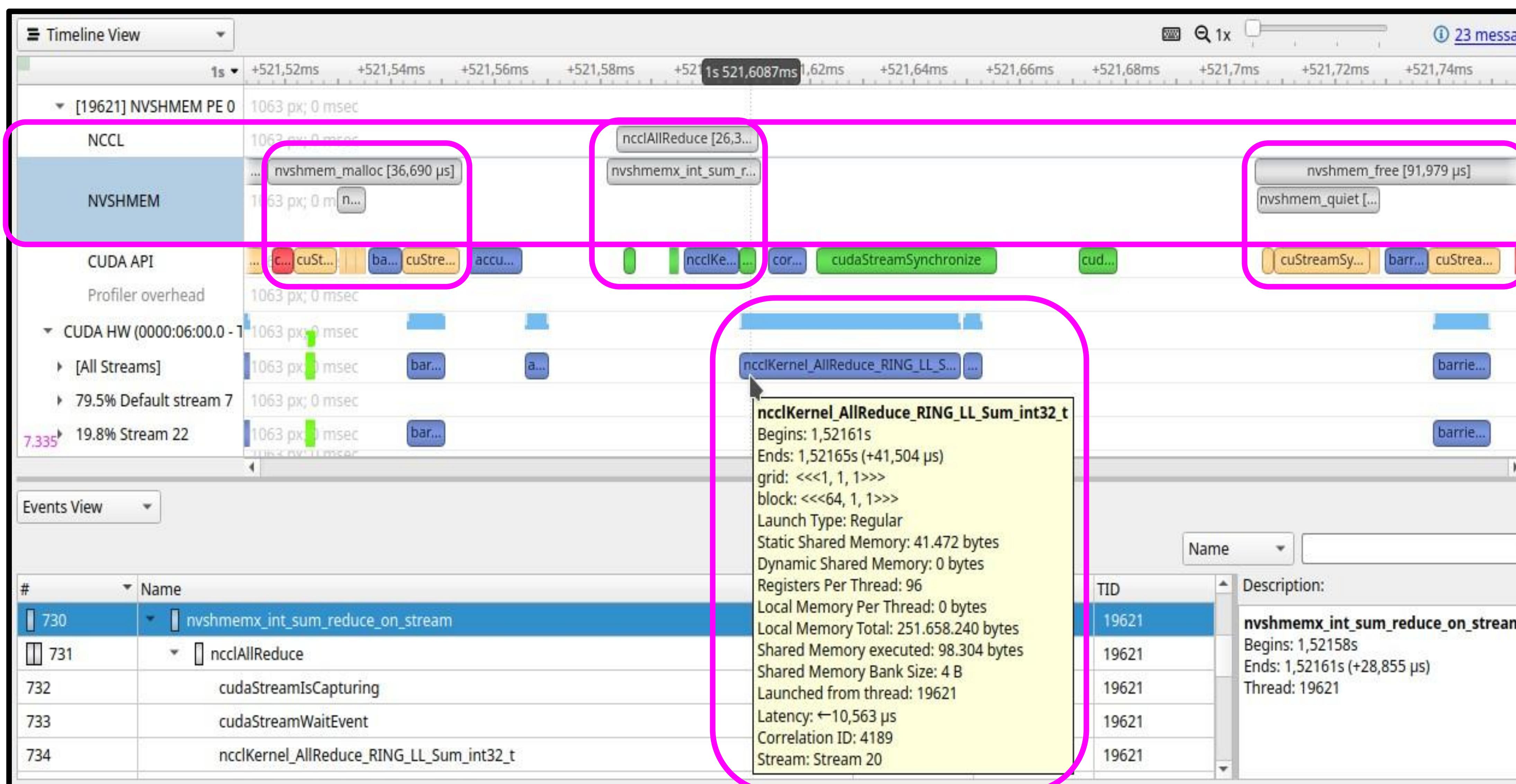
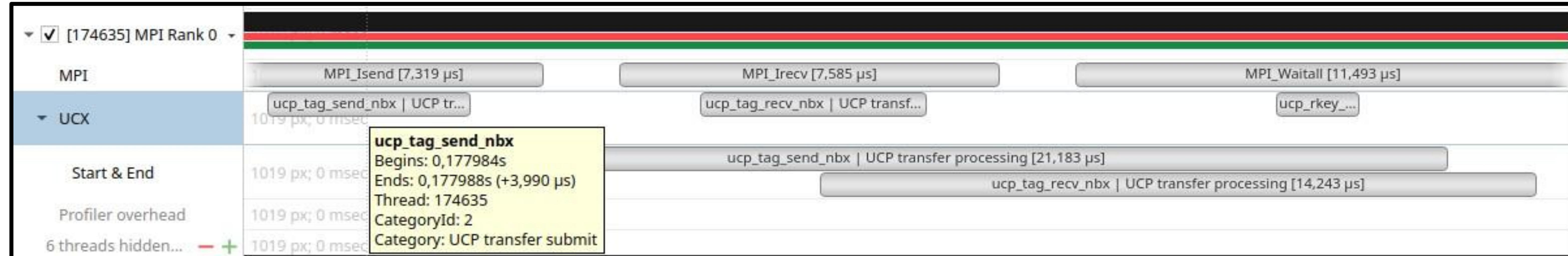


CPU call-stack samples

Networking

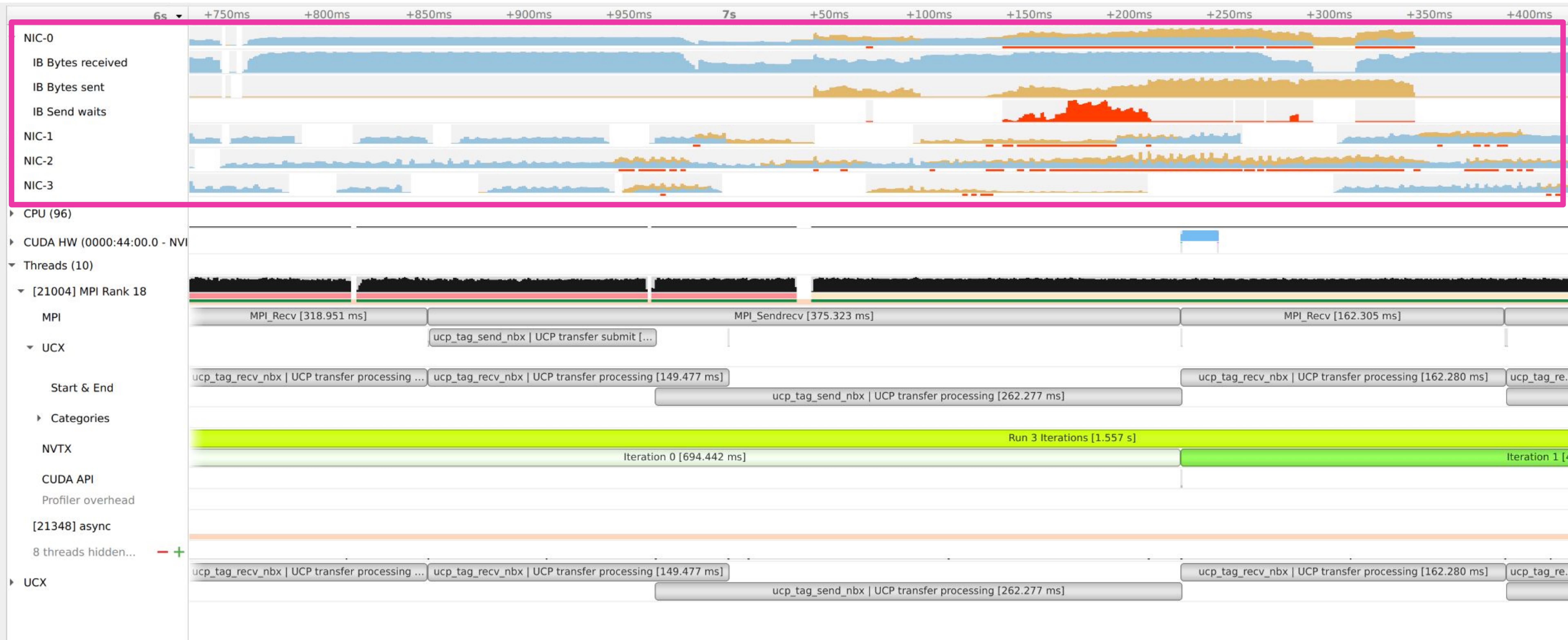
Communications Trace - LIBC, UCX, MPI, NCCL, & NVSHMEM

APIs, Destinations, Sizes, Async Completion, & Related GPU Events



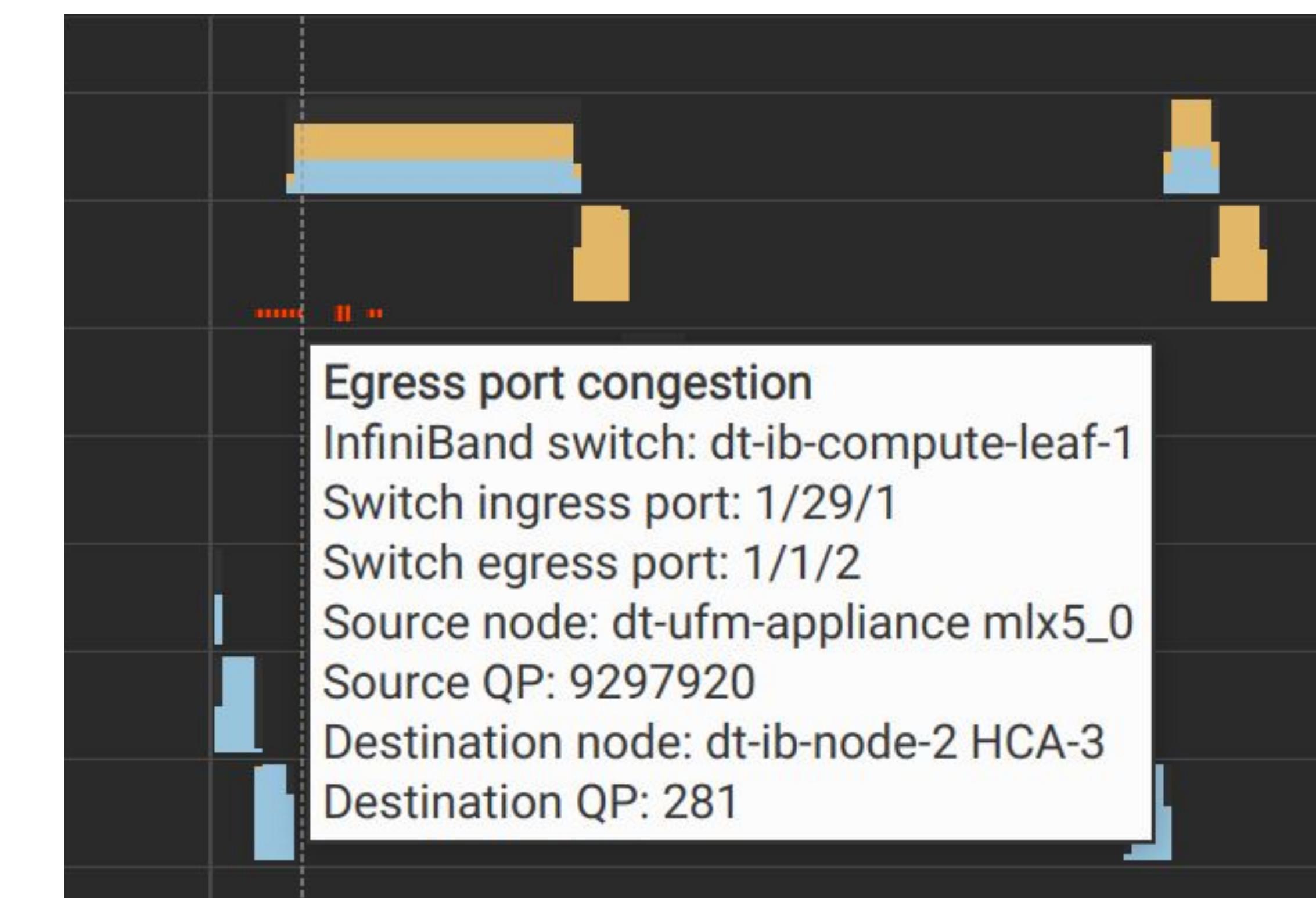
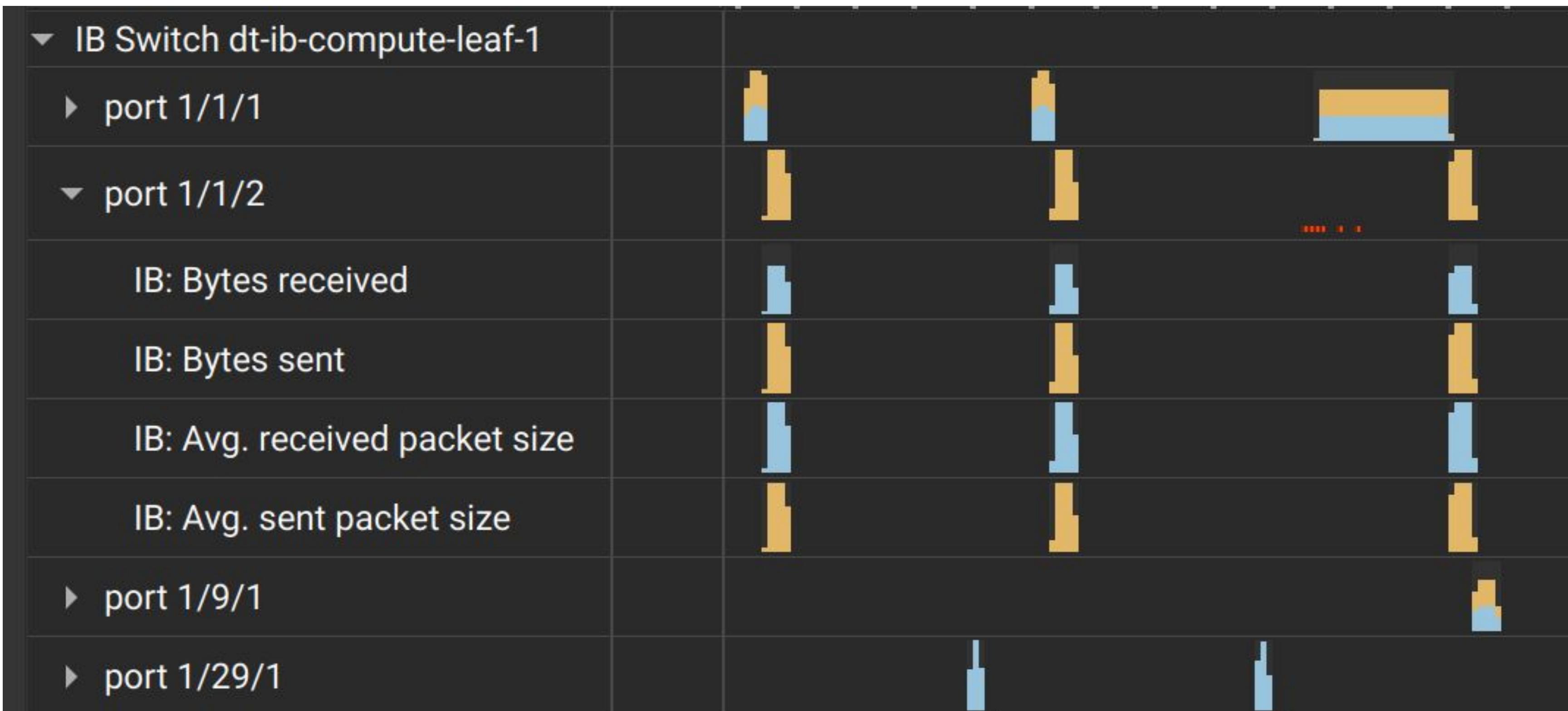
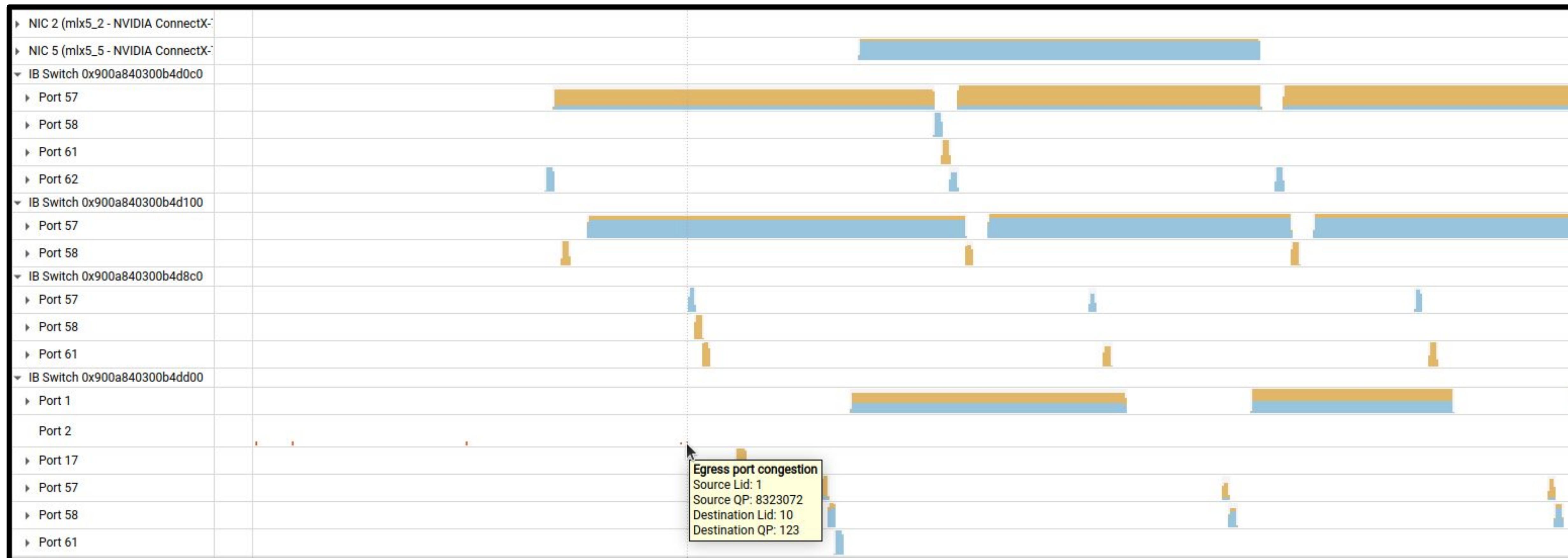
Description:	
MPI_Bcast	
Begins: 0,164935s	
Ends: 0,165045s (+109,210 µs)	
Thread: 1342434	
Bytes sent: 0	
Bytes received: 4	
Root: 0	
MPI_COMM_WORLD	

NVIDIA NIC/HCA Metrics Sampling (InfiniBand & RoCE)

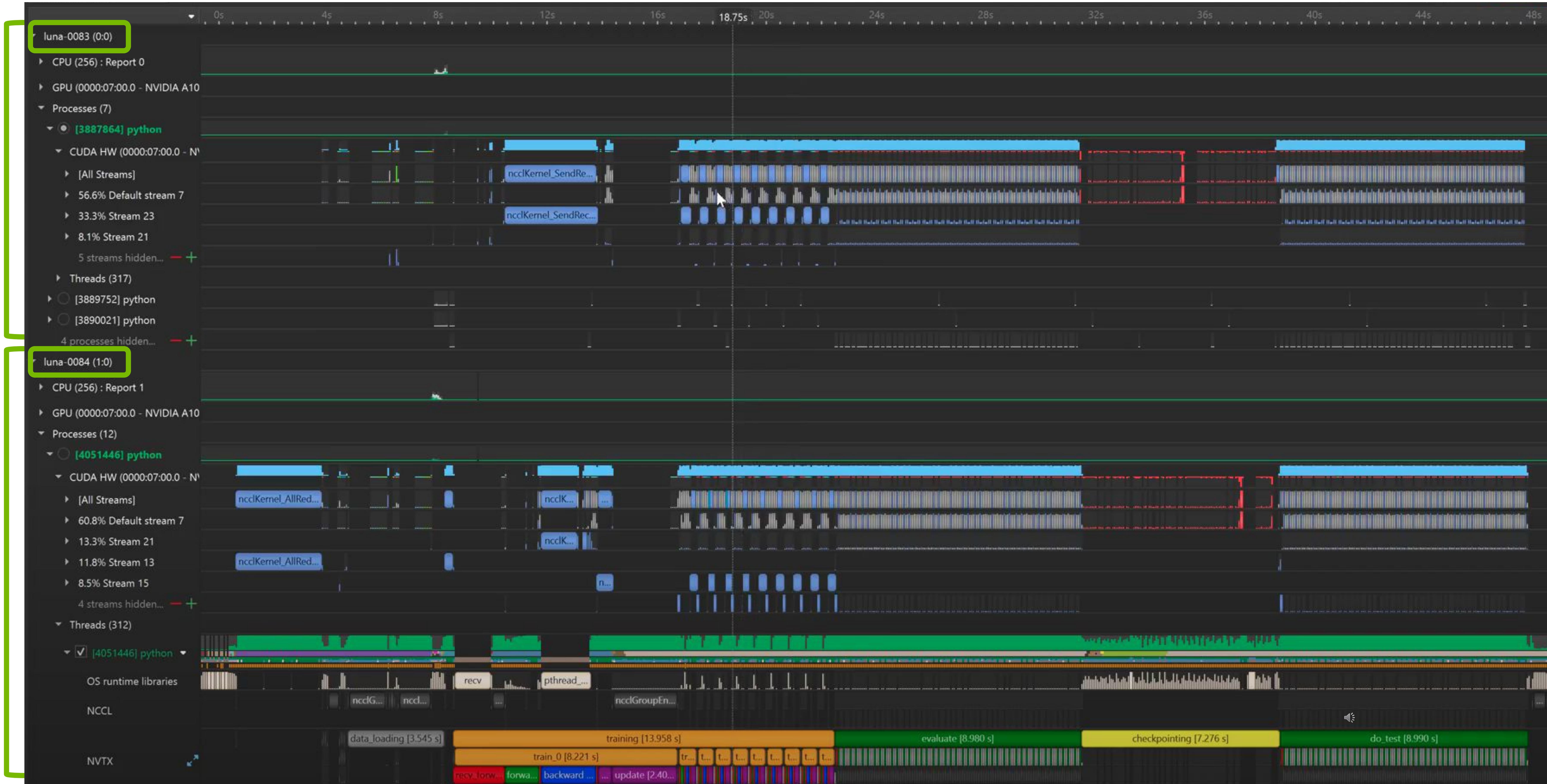


NVIDIA Quantum-2 Switch Sampling

Send, Receive, Congestion, & Now Per-Port Sampling



Combining Reports Dynamically - But Does That Scale?





Scaling Out
Multi-Report Analysis

Flexible Scale-Out Profiling

- Infrastructure diversity
 - On-premise or cloud service providers
 - Node hardware and networking fabrics
 - Single or multi-tenancy
 - Bare metal, virtual machines, containers
 - Job schedulers(Slurm), task scheduler(Dask), orchestrators(Kubernetes)
 - Micro-services vs homogenous tasks
- Often isolation with communication limitations
- No “single point of contact” required to setup & profile
- Pick and choose for your needs & evolving ecosystems
 - Many profiles taken independently

Multi-Node / Multi-Report Analysis

Reduce “getting started” time

Post-process instead of complicated realtime setups

Which might not be possible in your environment

Analyze thousands of reports

No manual measure & search

Avoid “guess which file”

Cluster-wide conclusion

Scale-out processing

Library of recipes → customizable

Results presentations (jupyter, graphs, ...)

Collaborate / Share



Analysis Recipes Library <= 2024.1

CLI Usage: nsys recipe [<args>] <recipe name> [<recipe args>]

cuda_api_sum -- CUDA API Summary

cuda_api_sync -- CUDA Synchronization APIs

cuda_gpu_kern_pace -- CUDA GPU Kernel Pacing

cuda_gpu_kern_sum -- CUDA GPU Kernel Summary

cuda_gpu_mem_size_sum -- CUDA GPU MemOps Summary (by Size)

cuda_gpu_mem_time_sum -- CUDA GPU MemOps Summary (by Time)

cuda_gpu_time_util_map -- CUDA GPU Time Utilization Heatmap

cuda_memcpy_async -- CUDA Async Memcpy with Pageable Memory

cuda_memcpy_sync -- CUDA Synchronous Memcpy

cuda_memset_sync -- CUDA Synchronous Memset

diff -- Statistics Diff

dx12_mem_ops -- DX12 Memory Operations

gpu_gaps -- GPU Gaps

gpu_metric_util_map -- GPU Metric Utilization Heatmap

gpu_time_util -- GPU Time Utilization

mpi_gpu_time_util_map -- MPI and GPU Time Utilization Heatmap

mpi_sum -- MPI Summary

nccl_gpu_proj_sum -- NCCL GPU Projection Summary

nccl_sum -- NCCL Summary

nvtx_gpu_proj_pace -- NVTX GPU Projection Pacing

nvtx_gpu_proj_sum -- NVTX GPU Projection Summary

nvtx_gpu_proj_trace -- NVTX GPU Projection Trace

nvtx_pace -- NVTX Pacing

nvtx_sum -- NVTX Range Summary

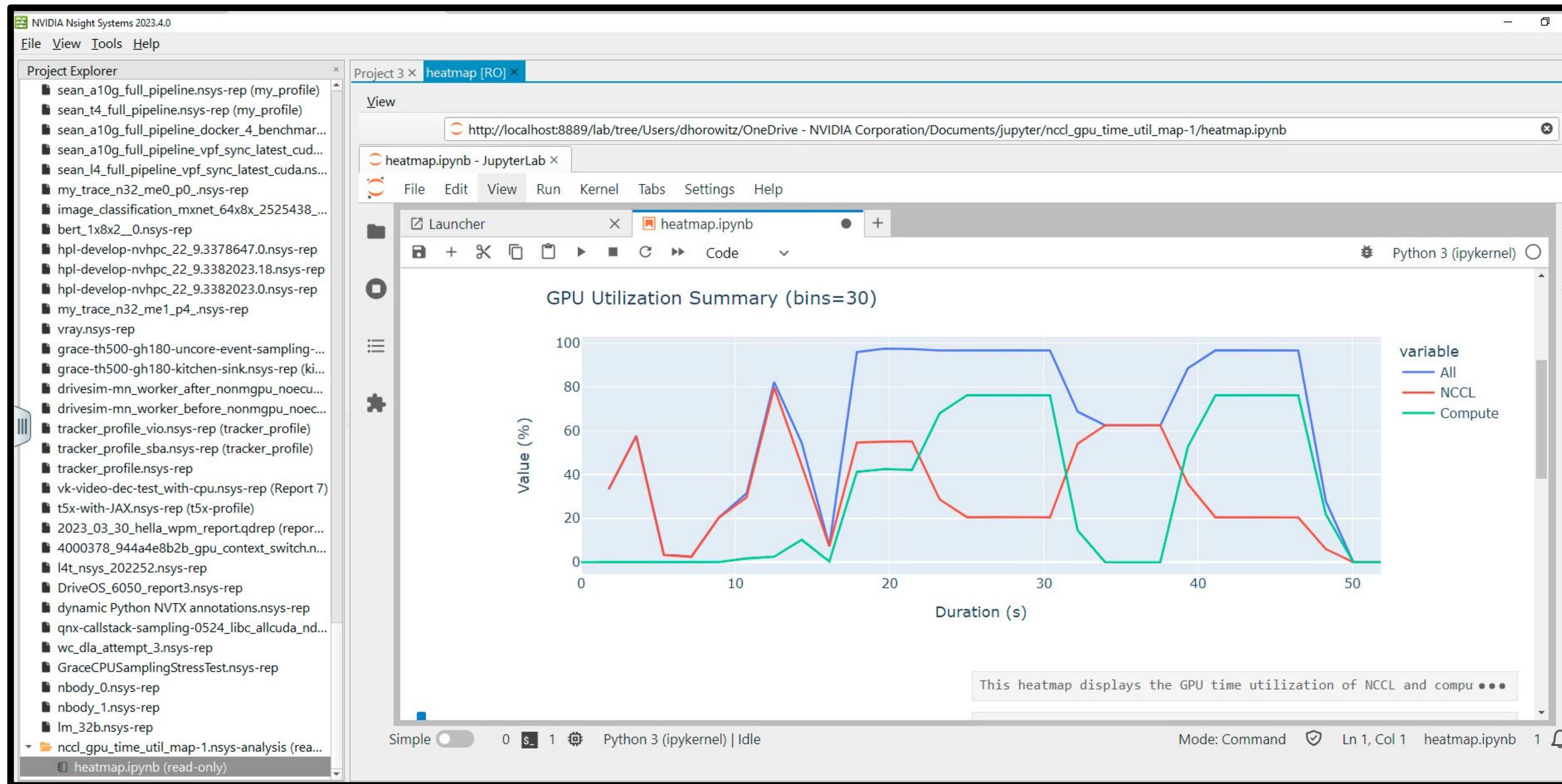
osrt_sum -- OS Runtime Summary

and growing...

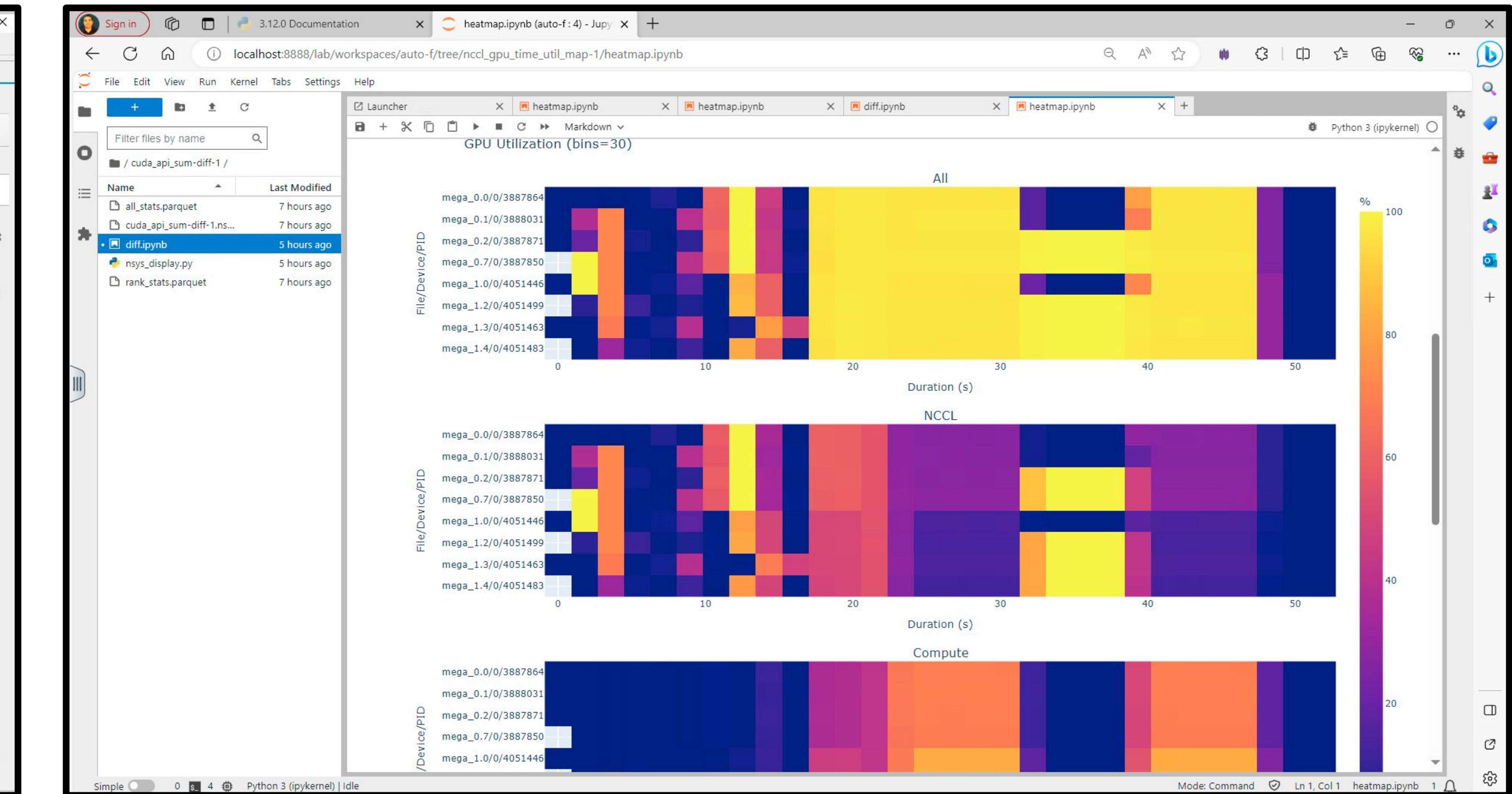
Jupyter Notebooks As Data Presentation

Loadable in...

Nsight Systems

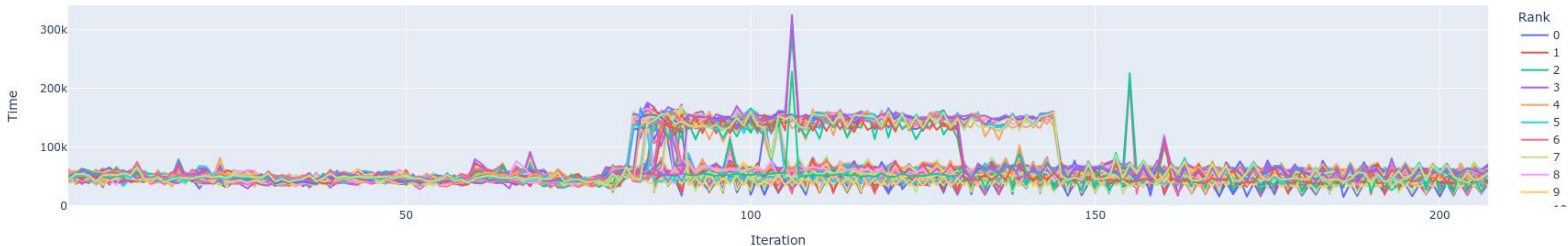


Jupyter Lab



Example From GTC 2023 - DL Training Barrier Time

Progress - Iterations defined by duration of ncclKernel_AllReduce_TREE_LL_Sum_int64_t



Consistency - Iterations defined by duration of ncclKernel_AllReduce_TREE_LL_Sum_int64_t



Heatmaps Overview

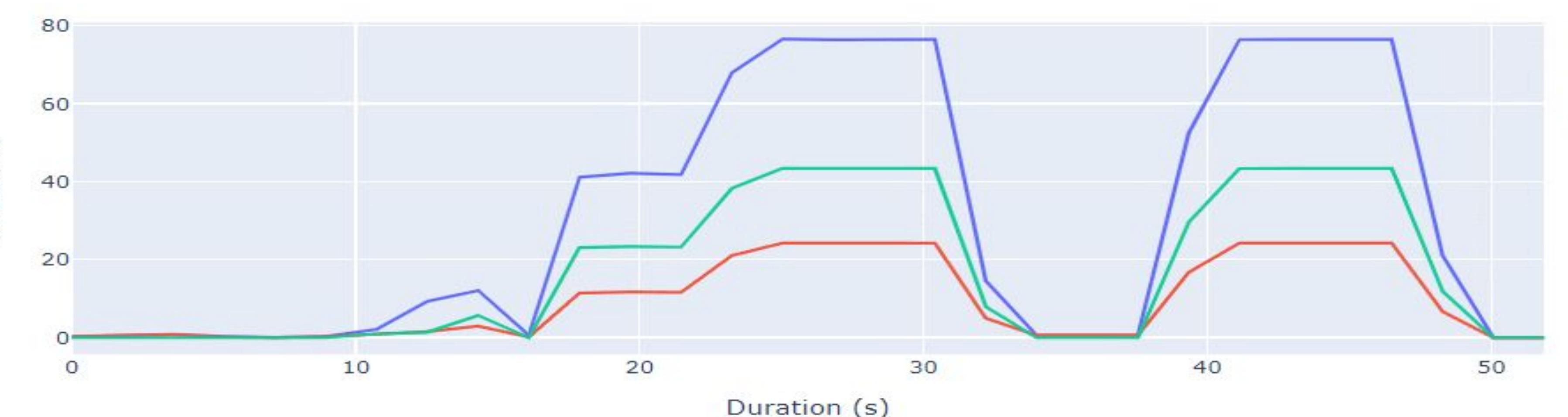
Ex: GPU metrics

- SM Active

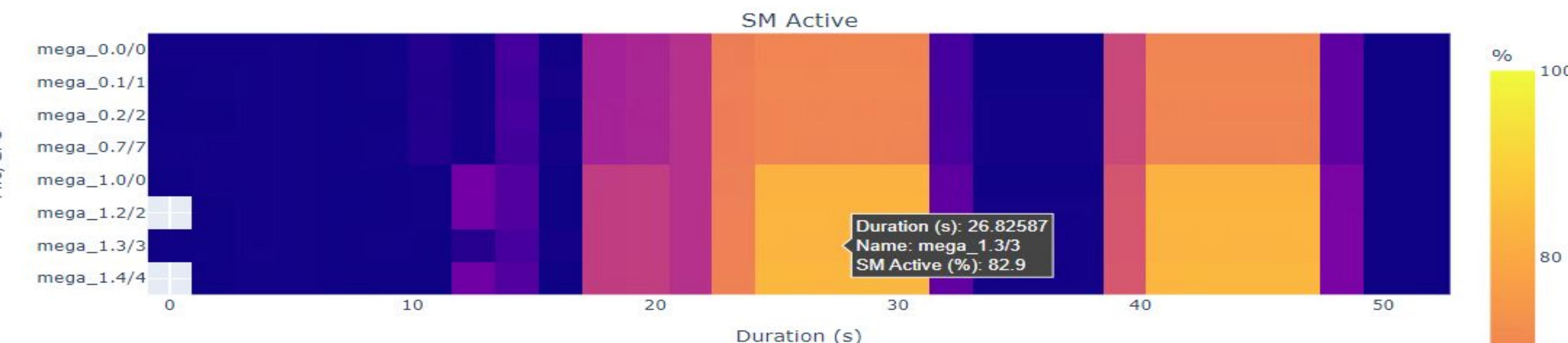
- SM Issue

- Tensor Active

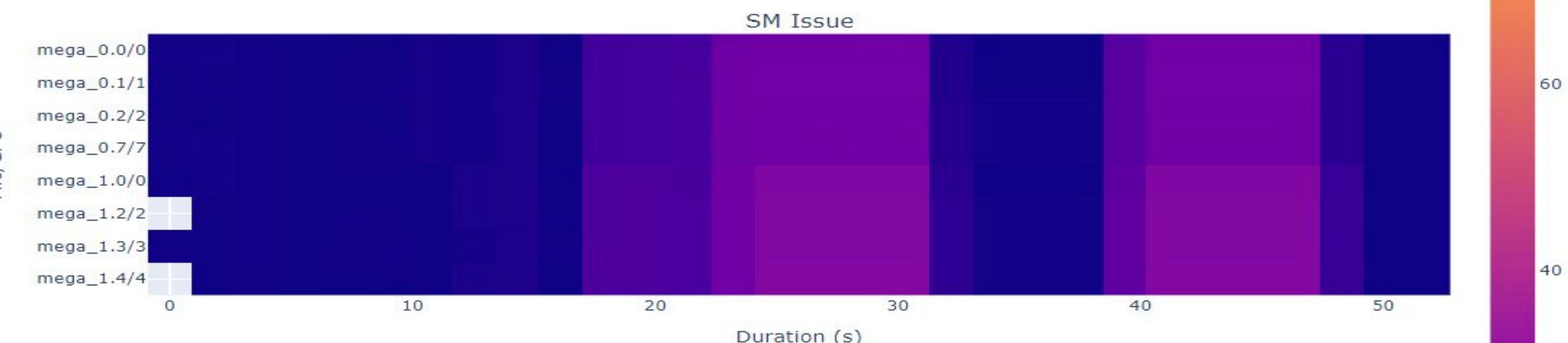
GPU Utilization Summary (bins=30)



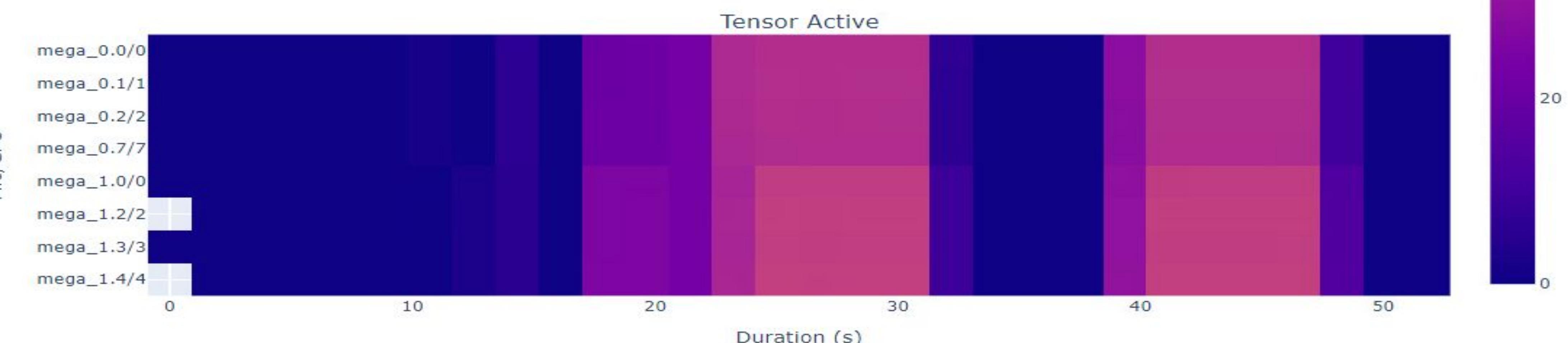
GPU Utilization (bins=30)



SM Active



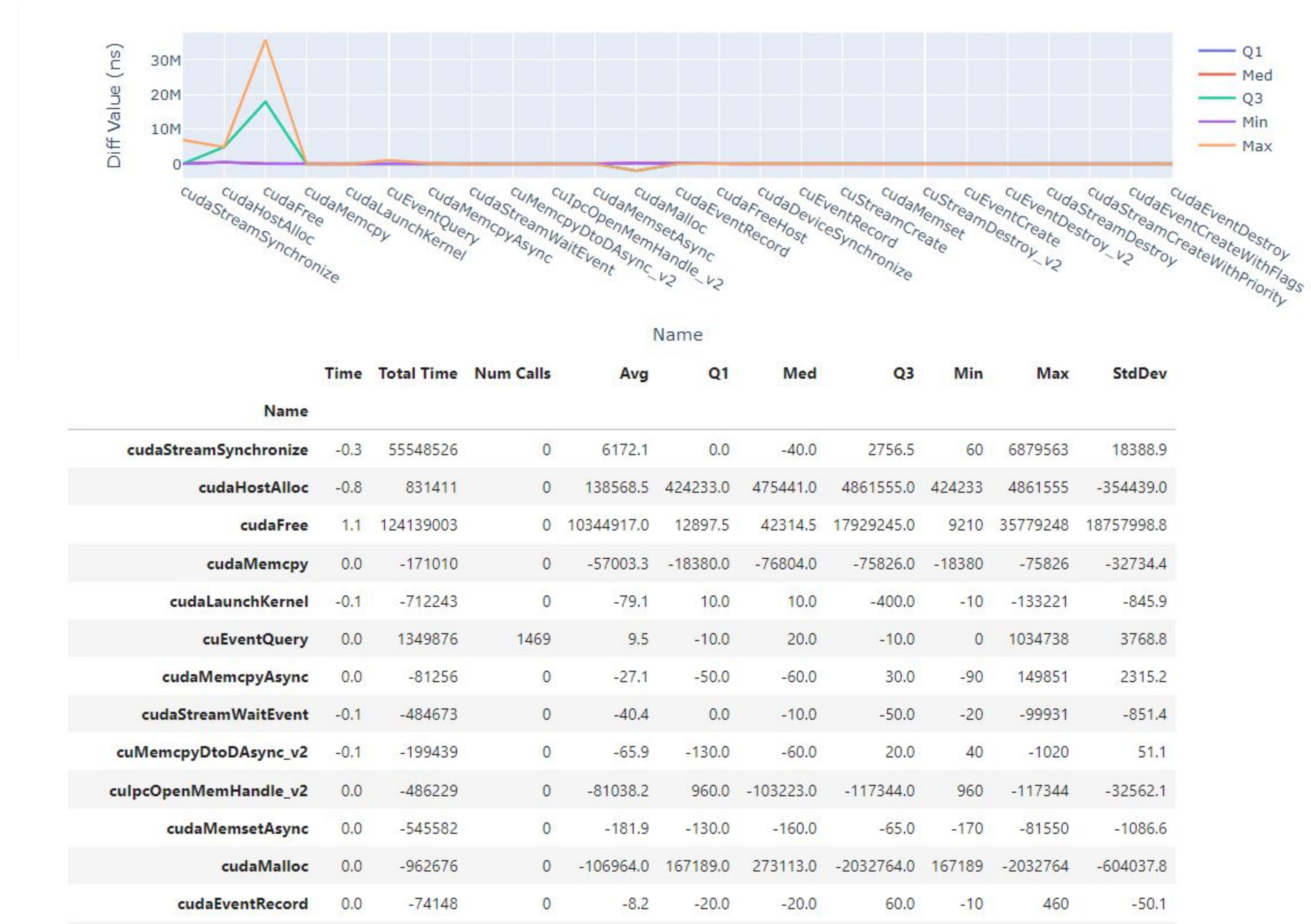
SM Issue



Tensor Active

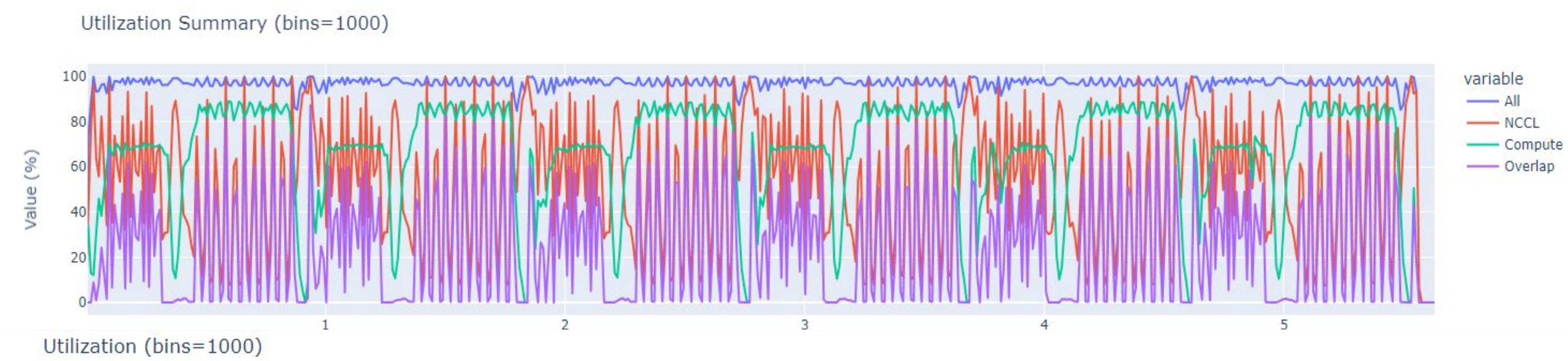
variable
— SM Active
— SM Issue
— Tensor Active

Recipe Report Differencing w/ Graphs

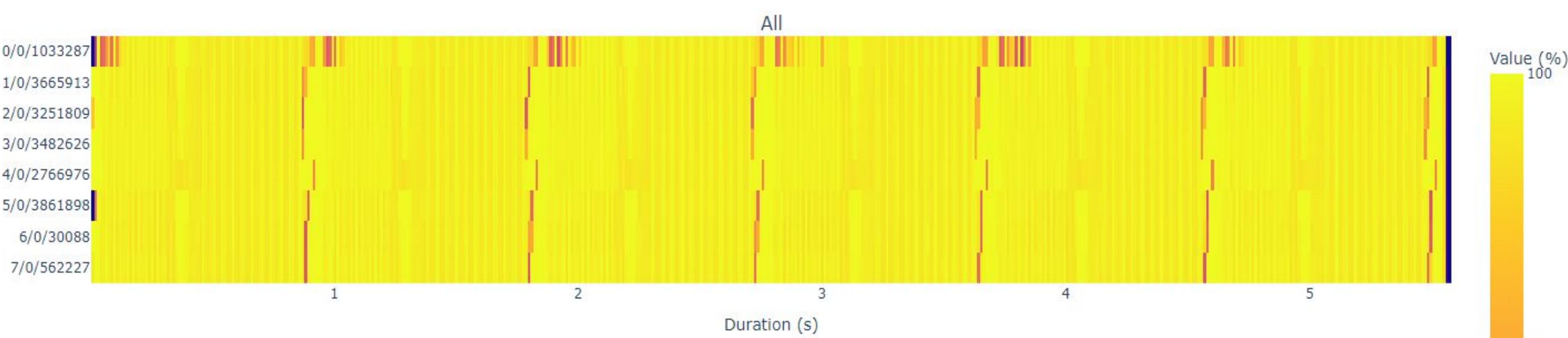


Compute-Comms Overlap

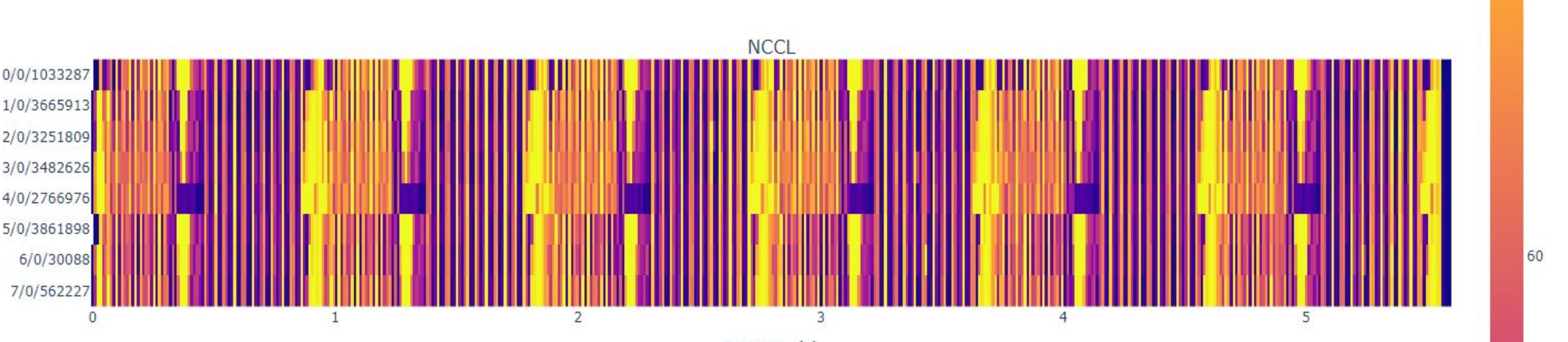
- GPU Utilization



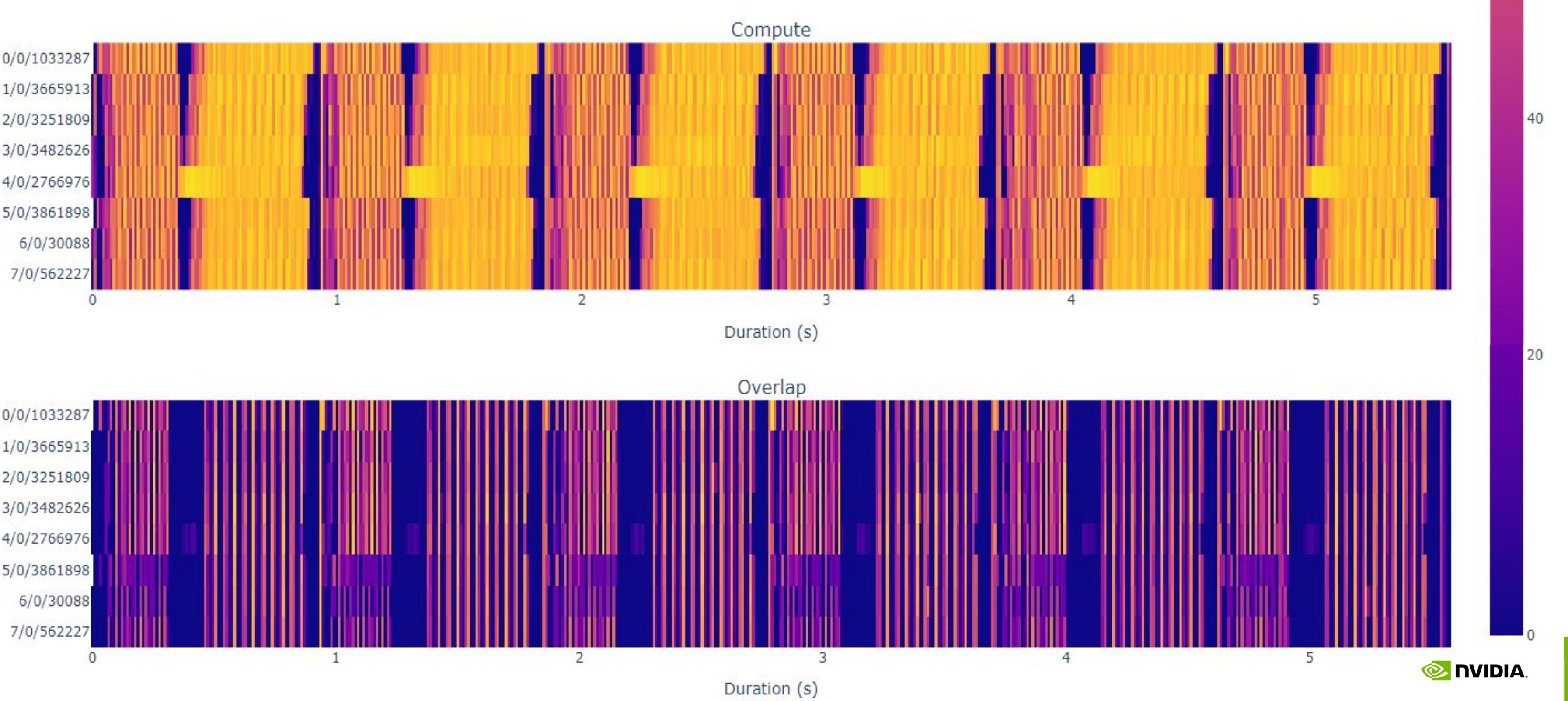
- Communications



- Compute



- Overlap



CUDA Kernel Overlap Statistics

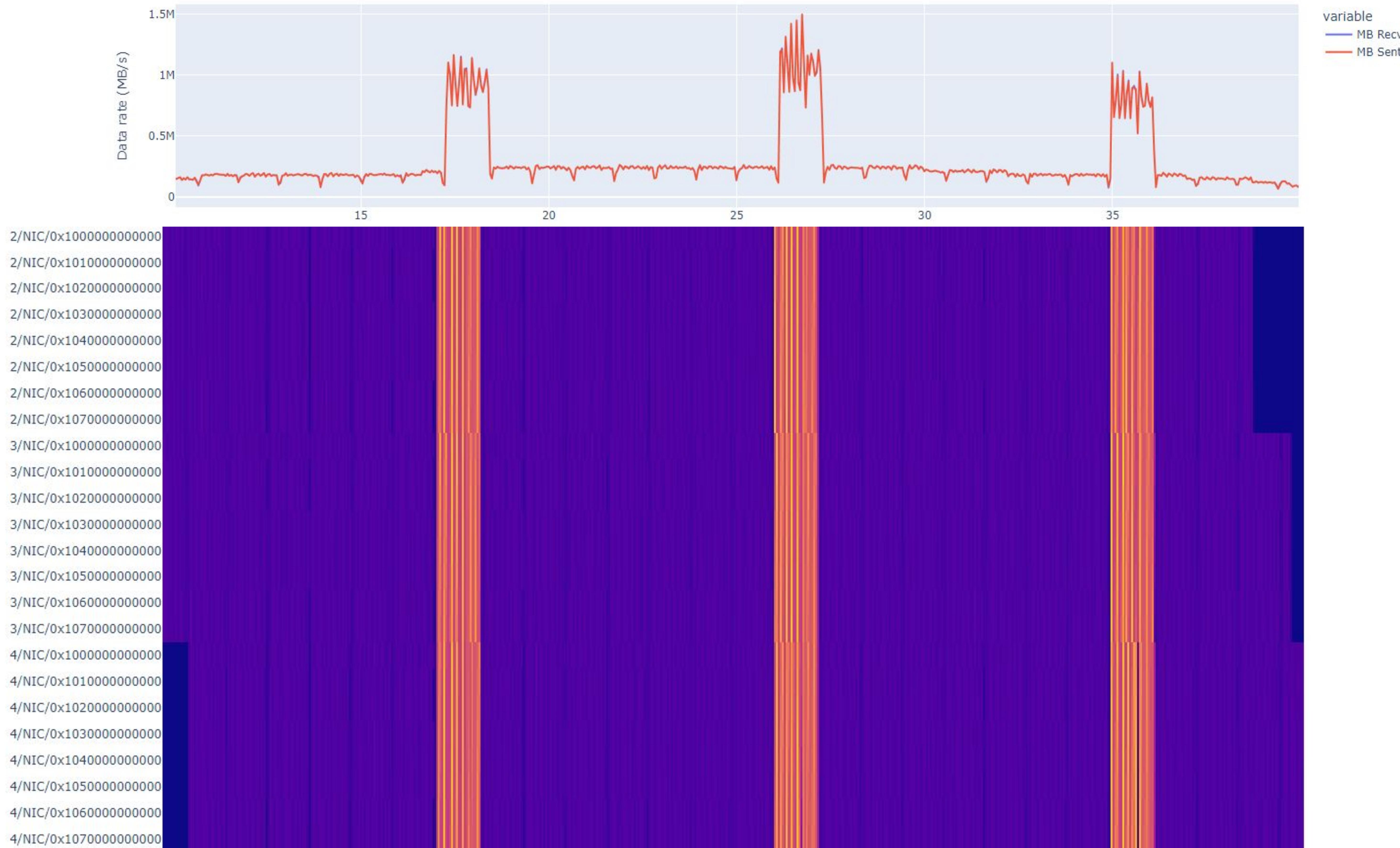
trace.ipynb +

Code Python 3 (ipykernel)

Name	Count	Communication Overlap	Compute Overlap
<code>CatArrayBatchedCopy</code>	48	0.2	0.0
<code>CatArrayBatchedCopy_aligned16_contig</code>	9222	0.3	0.0
<code>cast_transpose_dbias_dgelu_kernel</code>	4608	0.3	0.0
<code>cast_transpose_dbias_kernel</code>	4608	0.2	0.0
<code>cast_transpose_kernel</code>	6912	0.3	0.0
<code>cleanup</code>	48	0.0	0.0
<code>elementwise_kernel</code>	288	0.1	0.0
<code>elementwise_kernel_with_index</code>	9312	0.4	0.0
<code>embedding_backward_feature_kernel</code>	96	0.2	0.0
<code>flash_bwd_convert_dq_kernel</code>	4608	0.2	0.0
<code>flash_bwd_dot_do_o_kernel</code>	4608	0.3	0.0
<code>flash_bwd_dq_dk_dv_loop_seqk_parallel_kernel</code>	4608	0.2	0.0
<code>flash_fwd_kernel</code>	4608	0.4	0.0
<code>indexSelectLargeIndex</code>	96	0.4	0.0
<code>index_elementwise_kernel</code>	144	0.0	0.0
<code>kernel1</code>	9216	0.3	0.0
<code>kernel2</code>	9216	0.3	0.0
<code>kuserbuffers_pushrecv</code>	50688	0.3	0.9
<code>kuserbuffers_pushsend</code>	50688	0.3	0.9
<code>ln_bwd_finalize_kernel</code>	48	0.0	0.0
<code>ln_bwd_finalize_tuned_kernel</code>	9216	0.3	0.0
<code>ln_bwd_kernel</code>	48	0.0	0.0
<code>ln_bwd_tuned_kernel</code>	9216	0.3	0.0
<code>ln_fwd_kernel</code>	48	0.1	0.0

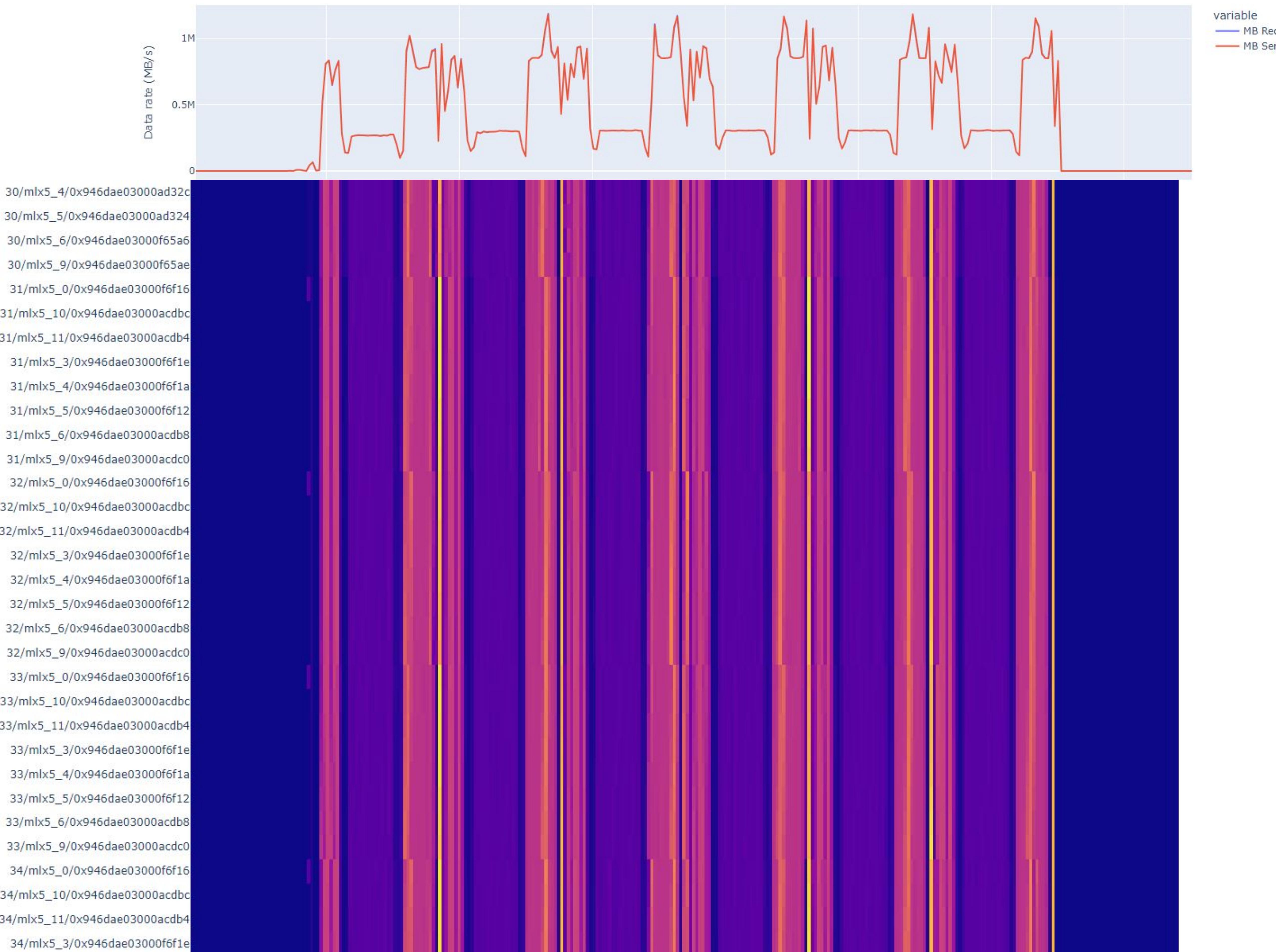
NVIDIA HCA/NIC Throughput Heatmaps

Example - Llamas training 64 HCAs - Many NICs off-screen



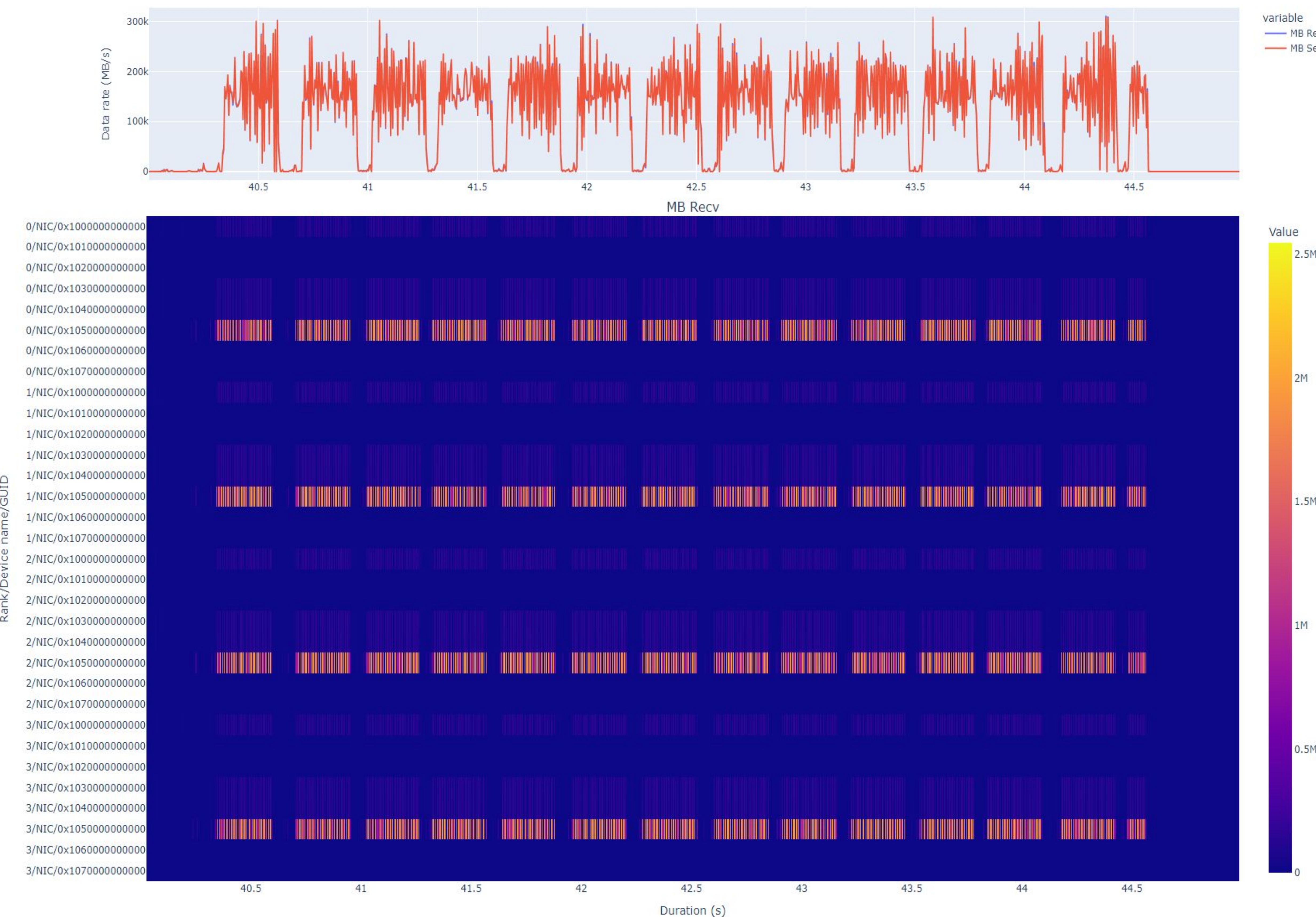
NVIDIA HCA/NIC Throughput Heatmaps

Example - LLM training 128 HCAs - Many NICs off-screen



NVIDIA HCA/NIC Throughput Heatmaps

Example - Gromacs 16 HCAs 4 nodes - Many NICs off-screen



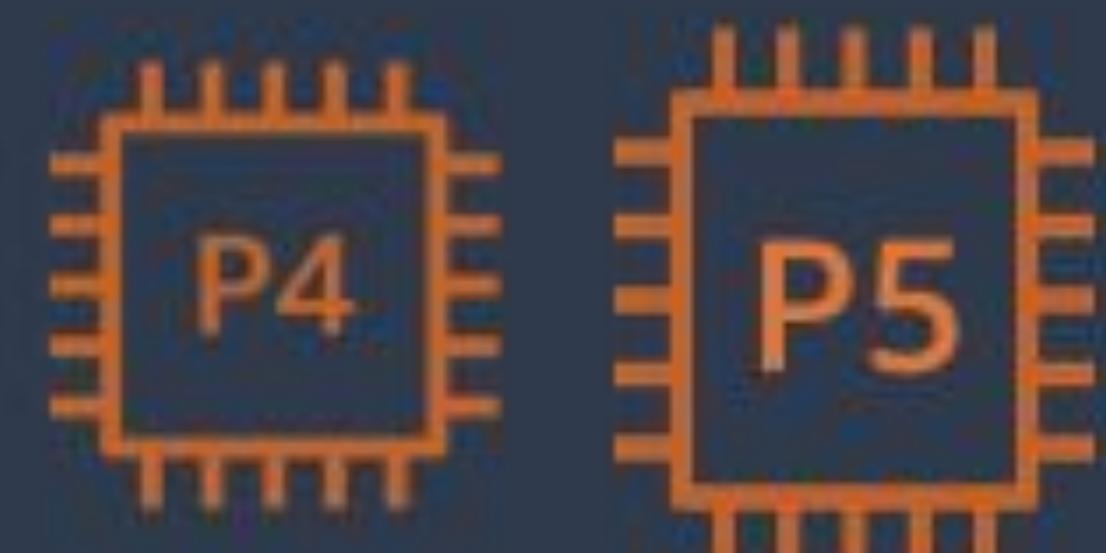
Per node, 1 GPU for long-range (PME) force calculations, which require all-to-all comms (1 NIC pinned per GPU for best perf). The remaining 3 GPUs per node are dedicated to short-range force calcs which have less intensive comm requirements. [Massively Improved Multi-node NVIDIA GPU Scalability with GROMACS | NVIDIA Technical Blog](#)

Cloud and Amazon AWS Collaborations

Ankur Srivastava, Amazon, AWS Senior Solutions Architect

Distributed training – HPC-ML cluster building blocks

Compute



Instances

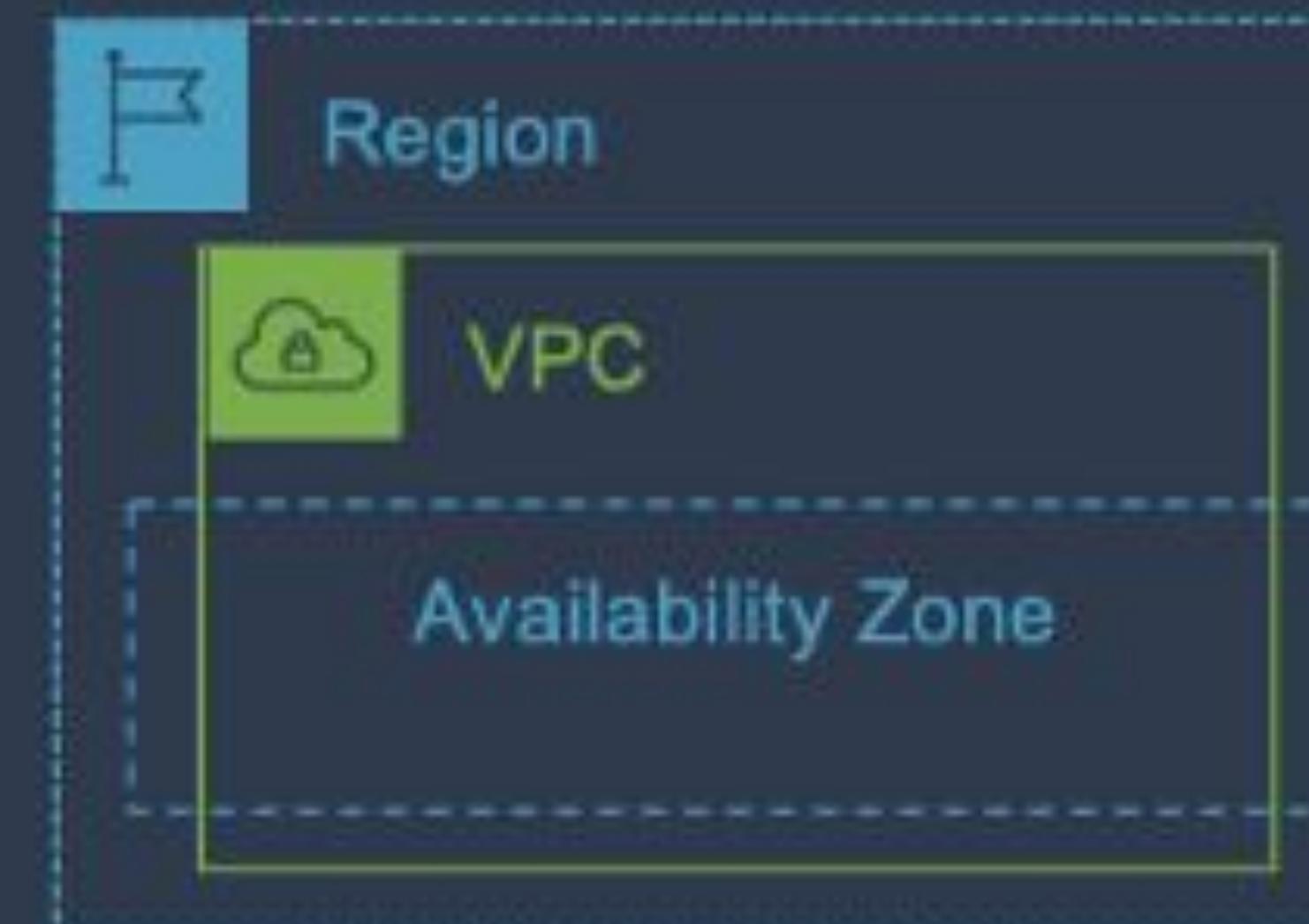


Containers

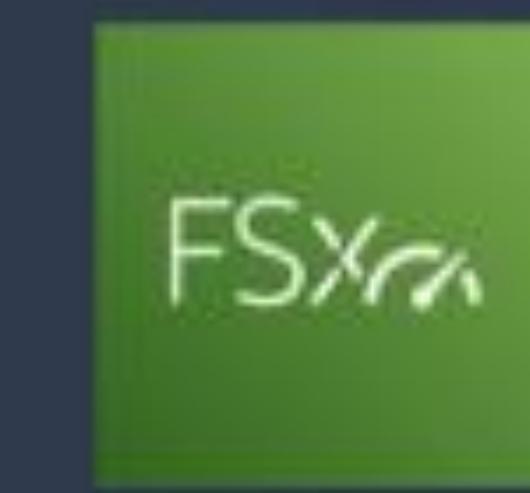
Network



Elastic Fabric Adapter



Storage



Amazon FSx
for Lustre

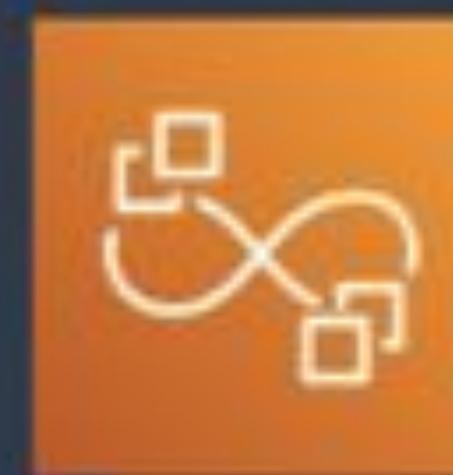


Amazon EC2
Instance Store



Amazon Elastic Block Store
(Amazon EBS)

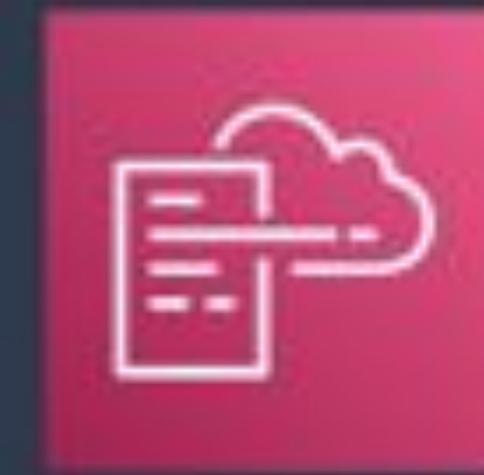
Architecture & Orchestration



AWS ParallelCluster



Amazon EKS



AWS CloudFormation

Distributed Training on AWS

Sagemaker
Jobs

Sagemaker
Hyperpod

AWS
ParallelCluster

Amazon EKS

Fully Managed ML Service

**Resilient and Persistent clusters dedicated
for distributed training**

Slurm based HPC clusters

Kubernetes based clusters

Best Practices for Large Scale Distributed Training

Step by step guides to create clusters:

- One-click VPC Deployments
- Mount FSx for Lustre Filesystems
- EFA enabled clusters

Recipes to customize AMIs

AWS optimized Dockerfiles

EFA cheatsheet

Distributed Training Examples:

- Slurm scripts/K8 manifests
- Working with Pyxis/Enroot
- Nemo (MegatronLM, Multimodal, Bionemo)
- MosaicML
- DDP, FSDP
- SMDP, SMMP
- Tensorflow/Jax

Validation (NCCL Tests etc)

Observability (Prometheus-Grafana etc)

Profiling (Nsight Product Family)



Running Nsight Systems with SLURM

- Nsight Systems comes pre-installed with DL AMIs
- NIC Sampler to collect EFA metrics will be available in the next Nsight release.
 - WIP to merge the latest Nsight Systems comes pre-installed with DL AMIs
- **Supports P5, P4, P4de, G5 and G4dn instances**
- EFA metrics are shared across GPUs on the same node
- Compatible with container plugins like [Pyxis](#)
- **srun -l "\${ARGS[@]}" /fsx/nsys-slurm-exec python3 <train.py>**

nsys-slurm-exec

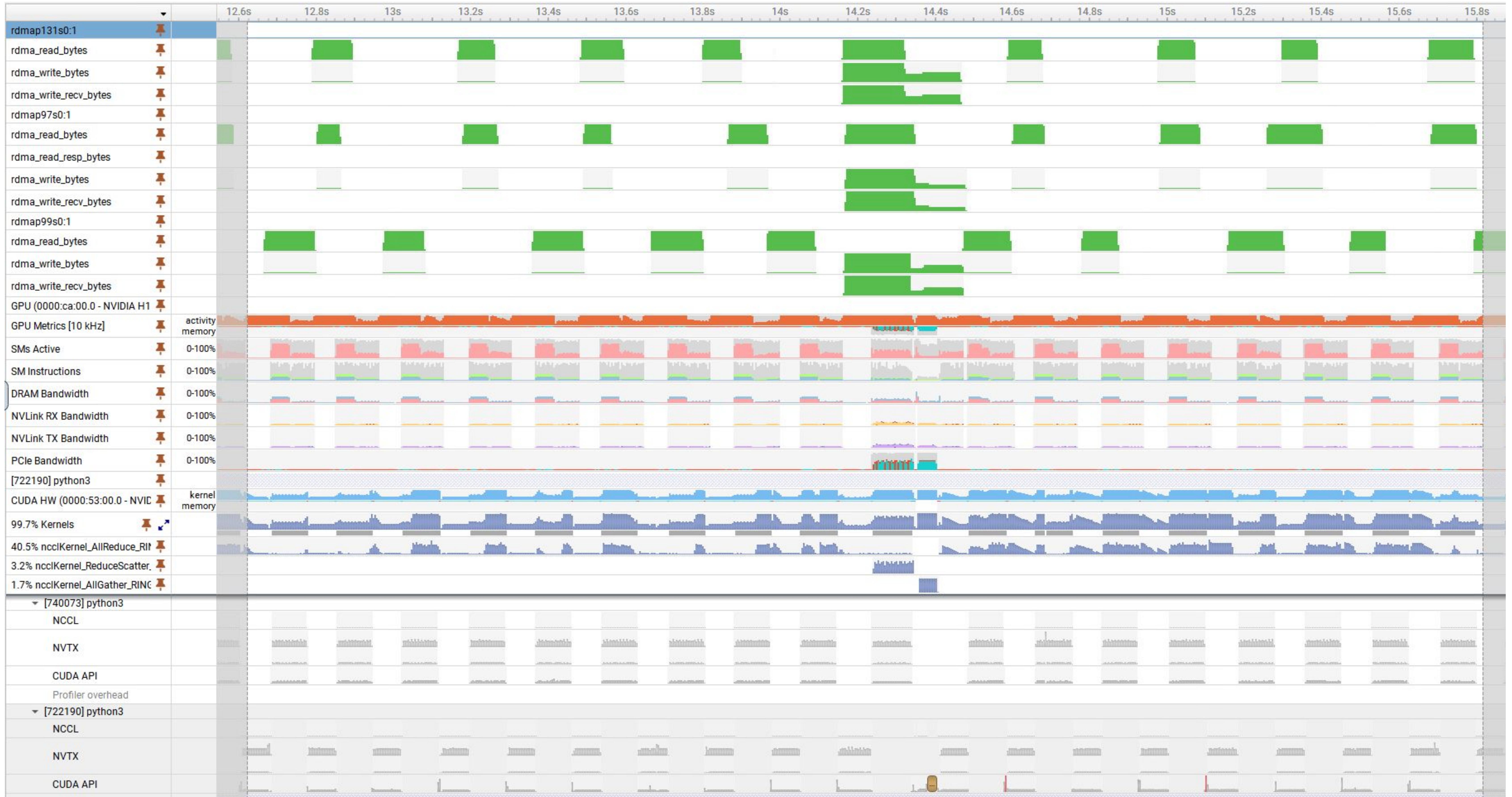
```
#!/bin/bash -x

NSYS_EXTRAS=""
if [ "$SLURM_LOCALID" == "0" ]; then
    NSYS_EXTRAS="--enable nic_sampler,-mode:counters,-struct:true,-efa:true"
fi

/fsx/nsight-efa/target-linux-x64/nsys profile $NSYS_EXTRAS \
    --gpu-metrics-device all --sample none --delay 30 \
    --force-overwrite true \
    --output report_nccl_efa_job%q{SLURM_JOB_ID}_rank%q{SLURM_PROCID}_on_%q{HOSTNAME}.nsys-report \
    "$@"
~
```

Amazon AWS Elastic Fabric Adapter(EFA) Sampling

Nemo training, in picture 3 of 32 EFA ports, 1 of 8 H100 GPUs ...



Amazon AWS Elastic Kubernetes Service (EKS)

Previously profiling in Kubernetes required

- Building Nsight Systems into the container
- Or mount a volume it into the container
- Changing the run or exec commands
- K8/Helm spec edits for profile vs production

Now profiling in Kubernetes requires

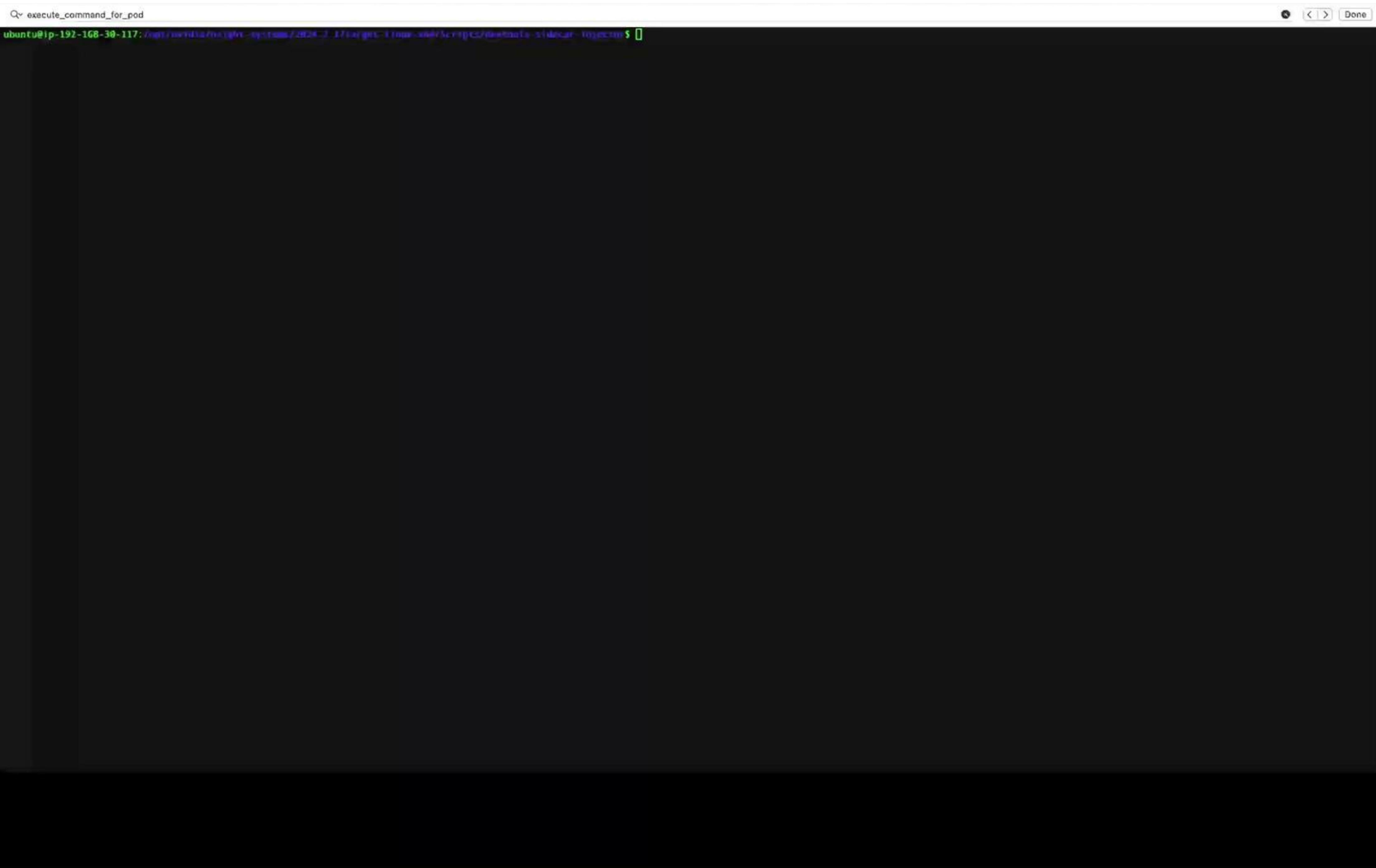
- Installing a helm chart & some custom values
- Labeling namespaces or pods of interest

NVIDIA's controller takes care of the rest!

Collaboration to support Amazon AWS EKS

Managed Kubernetes have host restrictions

Iterating to make it more powerful & user friendly

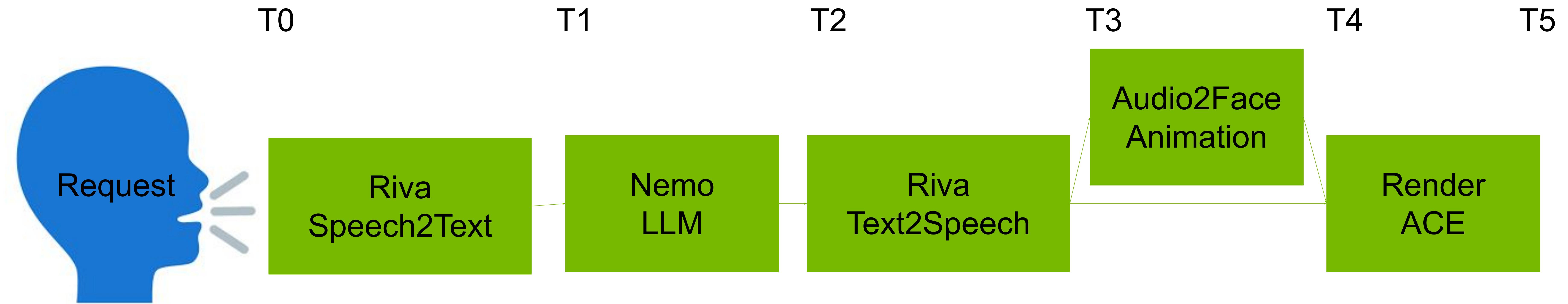


Cloud Microservices

Kubernetes & Helm

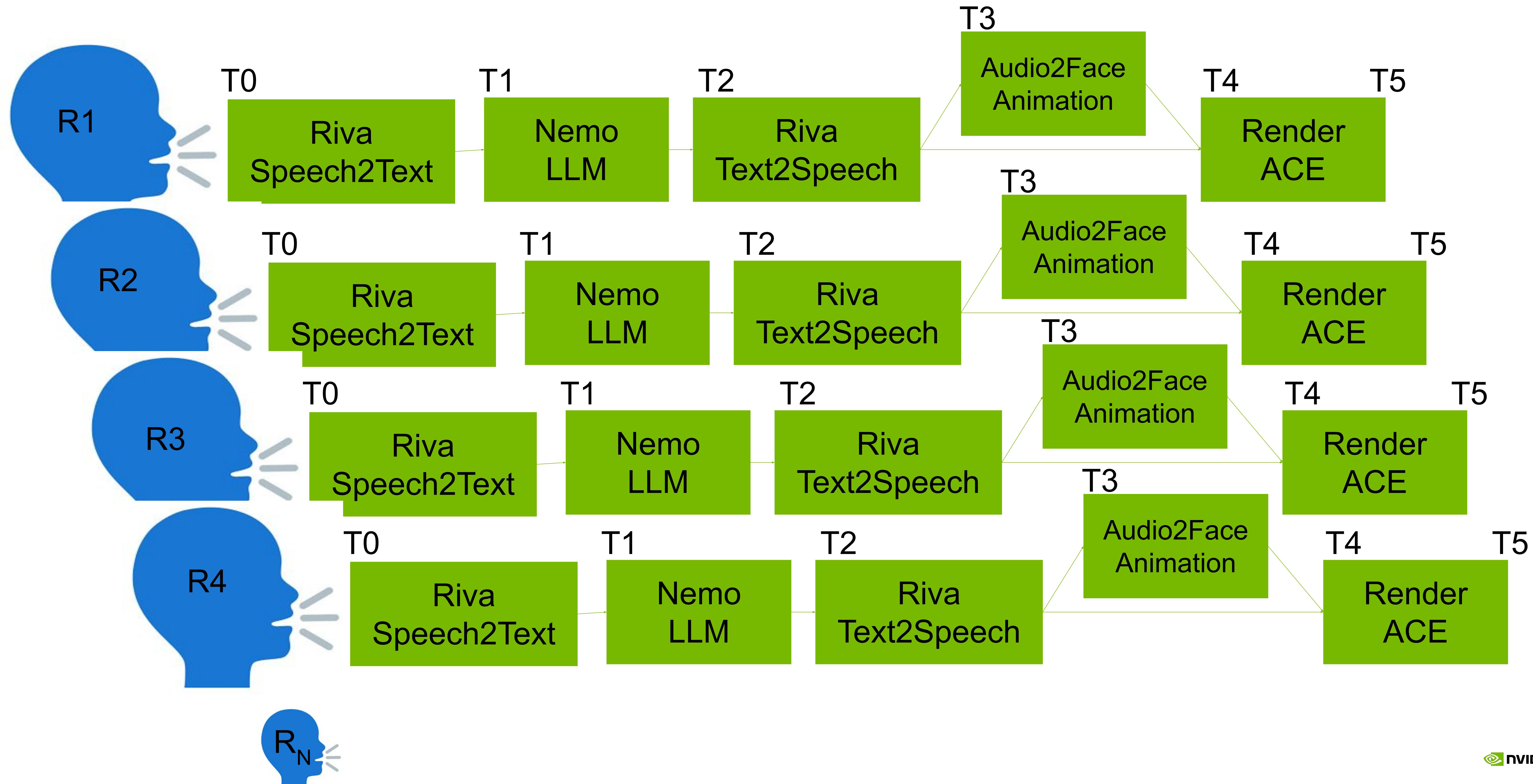
'My' cloud application with kubernetes & microservices!

Requests are a parallel pipelined across these stages/microservices



I want it to scale uniformly, consistent, & keep peek GPU performance

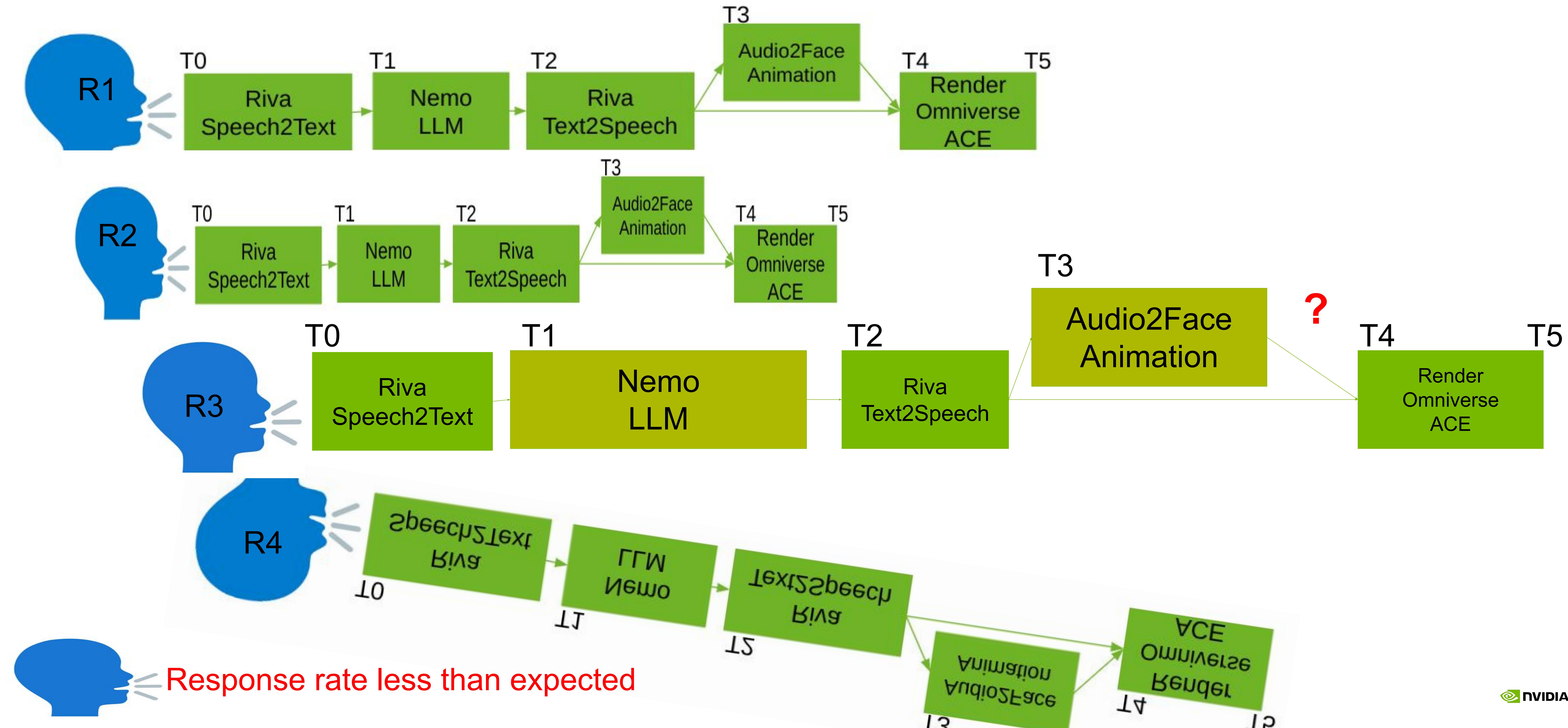
Ideally it behaves like this!!!



But what if I'm seeing this?!?! HELP

Find the needle in the haystack and investigate very deep!

Correlate to APIs, workloads, & source code algorithms, not just OpenTelemetry



NVIDIA DevTools Kubernetes Services

Deep investigations of container performance

Very high trace and sample rates vs OpenTelemetry

Injection into Pods/Container

Without modification of containers & k8/helm specs

Filtering

Kubernetes labels on namespaces or pods

Within-container process command-line regex

nsys_k8s for extra capabilities

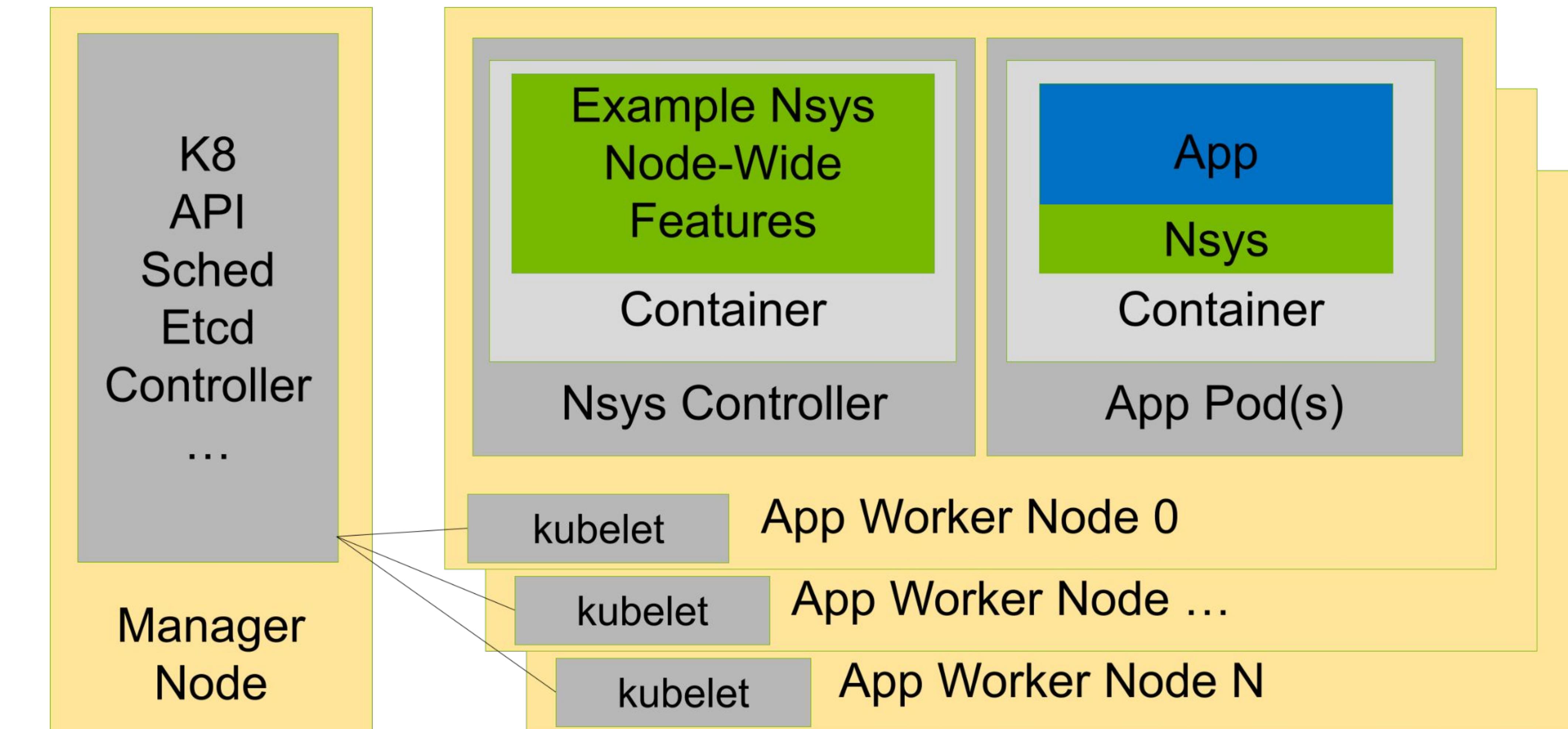
download without persistent volumes

Interactive start & stop

Available for CSP managed kubernetes services

Where to get it:

<https://catalog.ngc.nvidia.com/orgs/nvidia/teams/devtools/helm-charts/devtools-sidecar-injector>



Other Cloud and Remote Server Enhancements

Jupyter Lab

Jupyter Lab Extension

Cell profiling

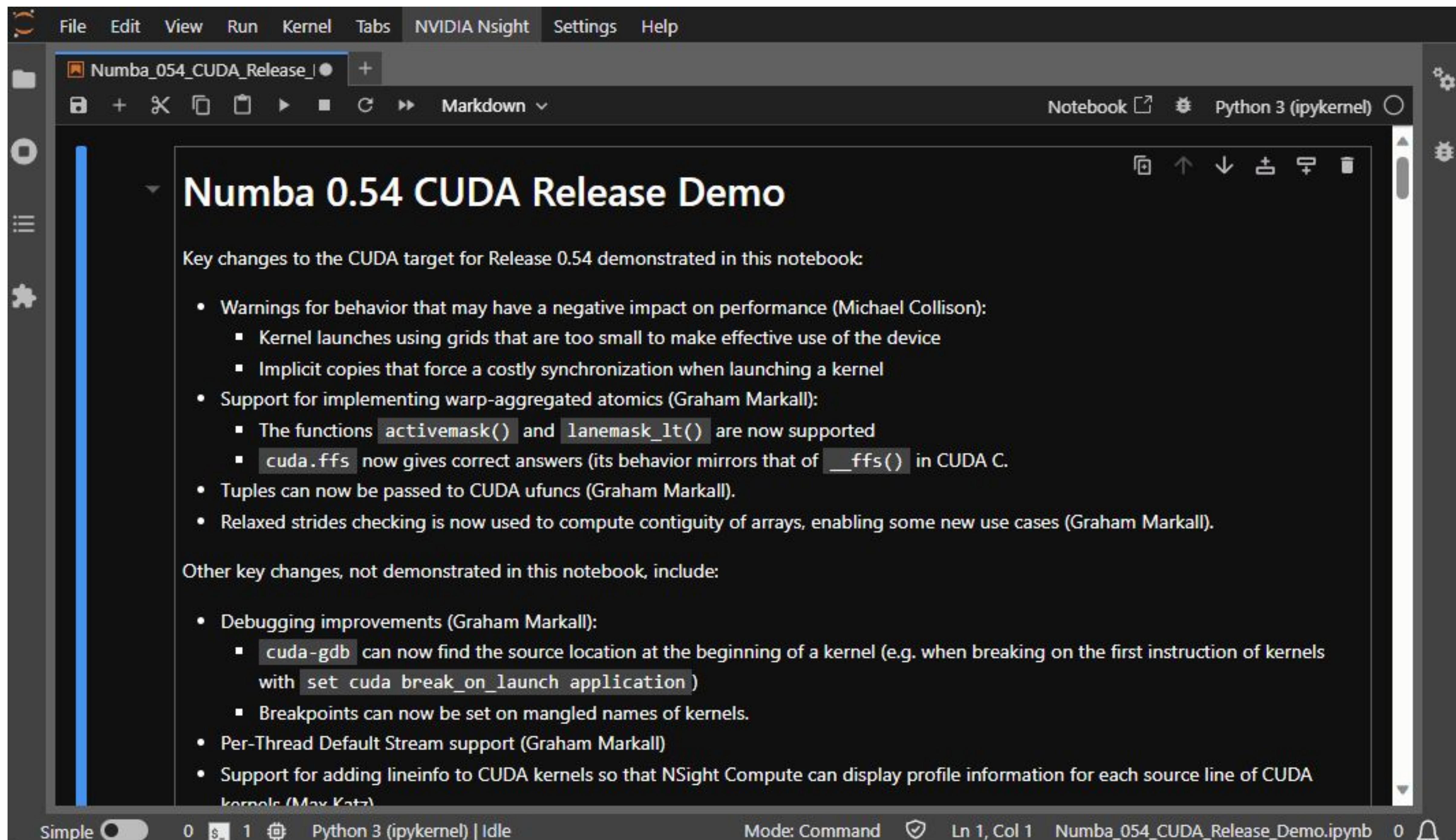
Launch remote streaming

Change options restarts kernel

Get it:

`pip install jupyterlab-nvidia-nsight`

<https://pypi.org/project/jupyterlab-nvidia-nsight/>



Remote UI Streaming

No need to download files

No need to match nsys versions

Access to more RAM vs laptop

Container with X, WM, & WebRTC server

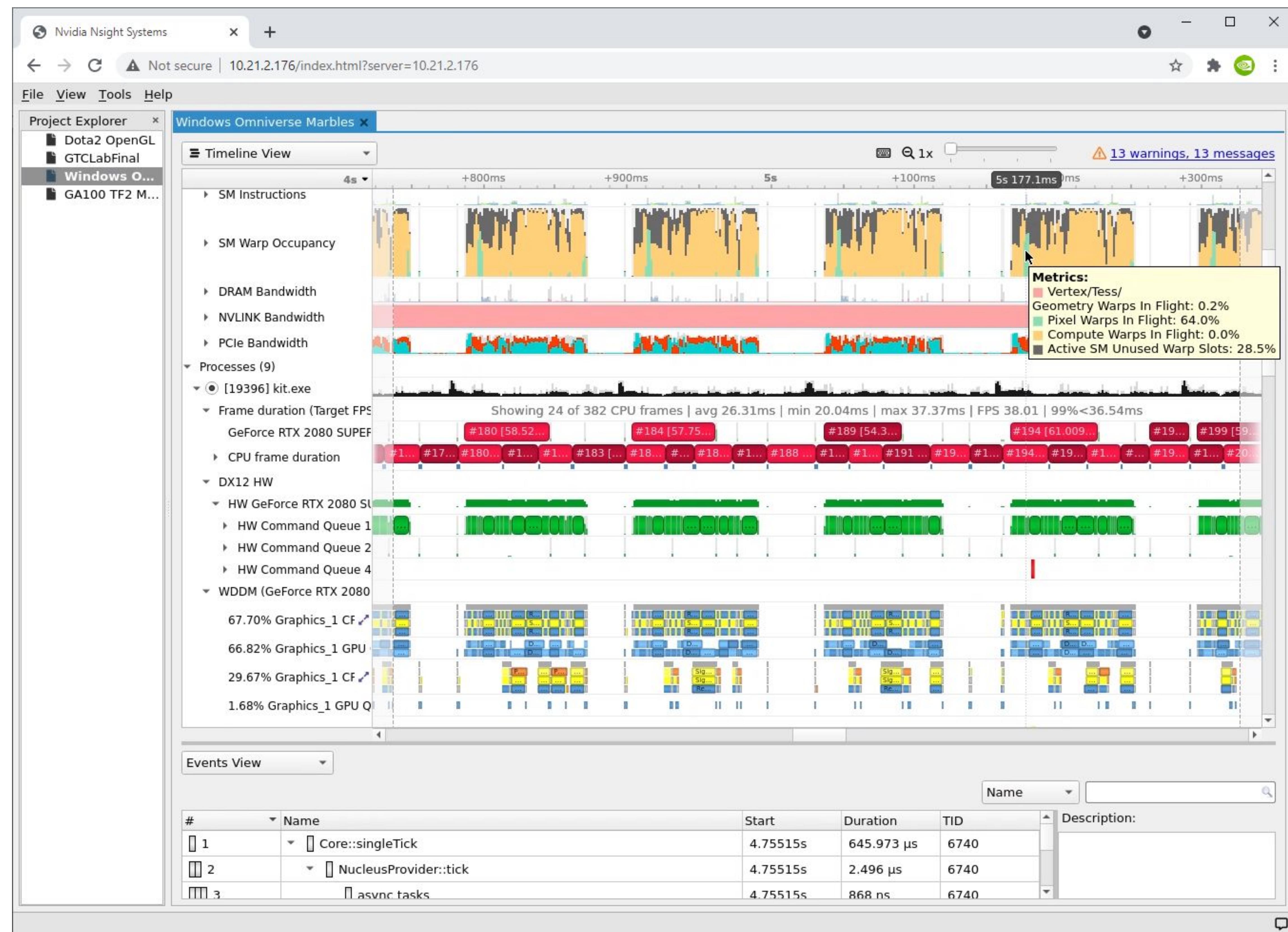
So your server doesn't have to!

Jupyter Lab extension can launch this!

Get it:

Dockerfile inside Nsight Systems Install dir

NVIDIA NCG Catalog soon



Python

Modern artificial intelligence and data analytics applications

Call-stack sampling

GIL trace

Code annotation library (NVTX)

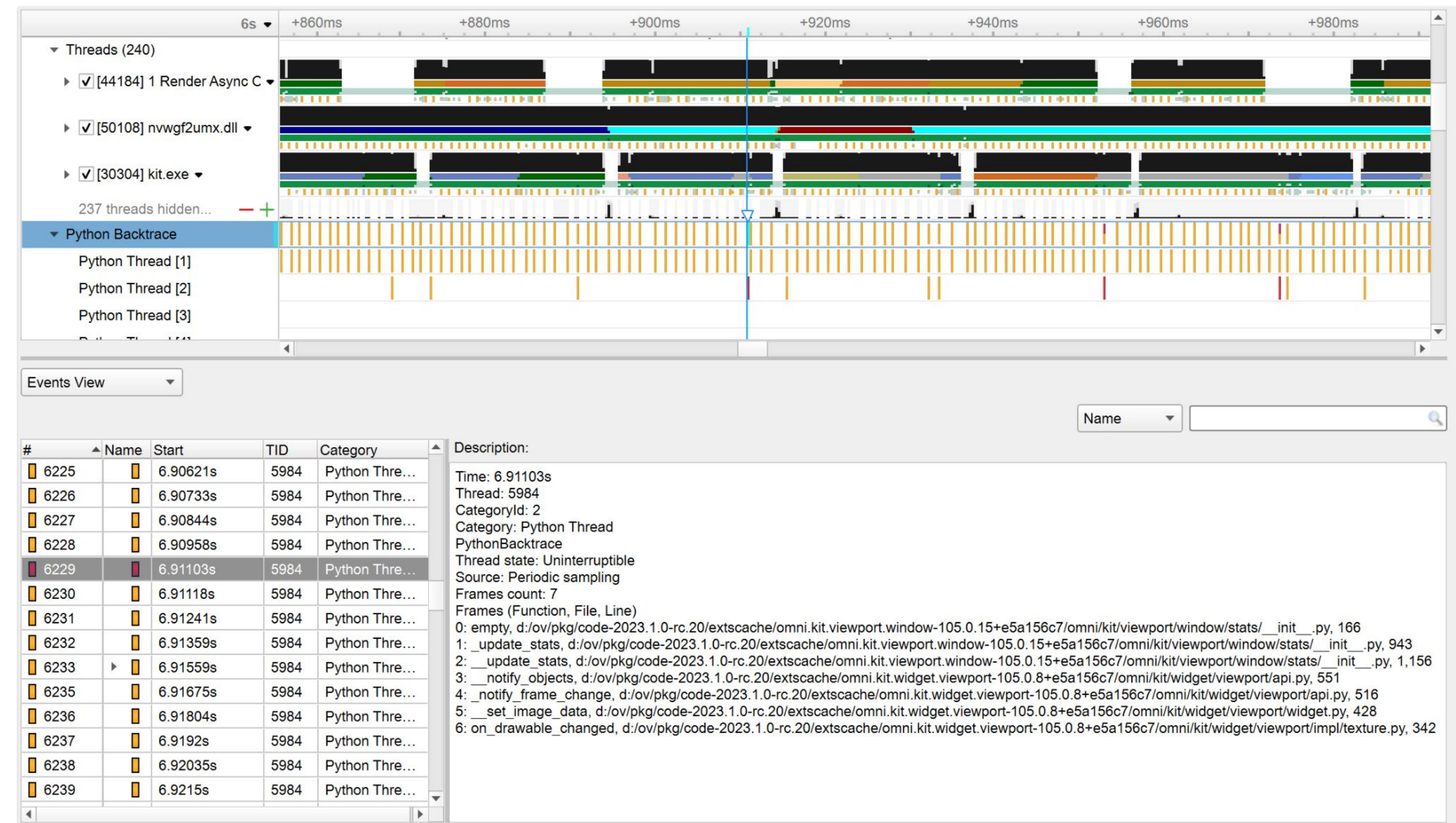
@ decorators

'With' statement

Download: 'pip install nvtx'

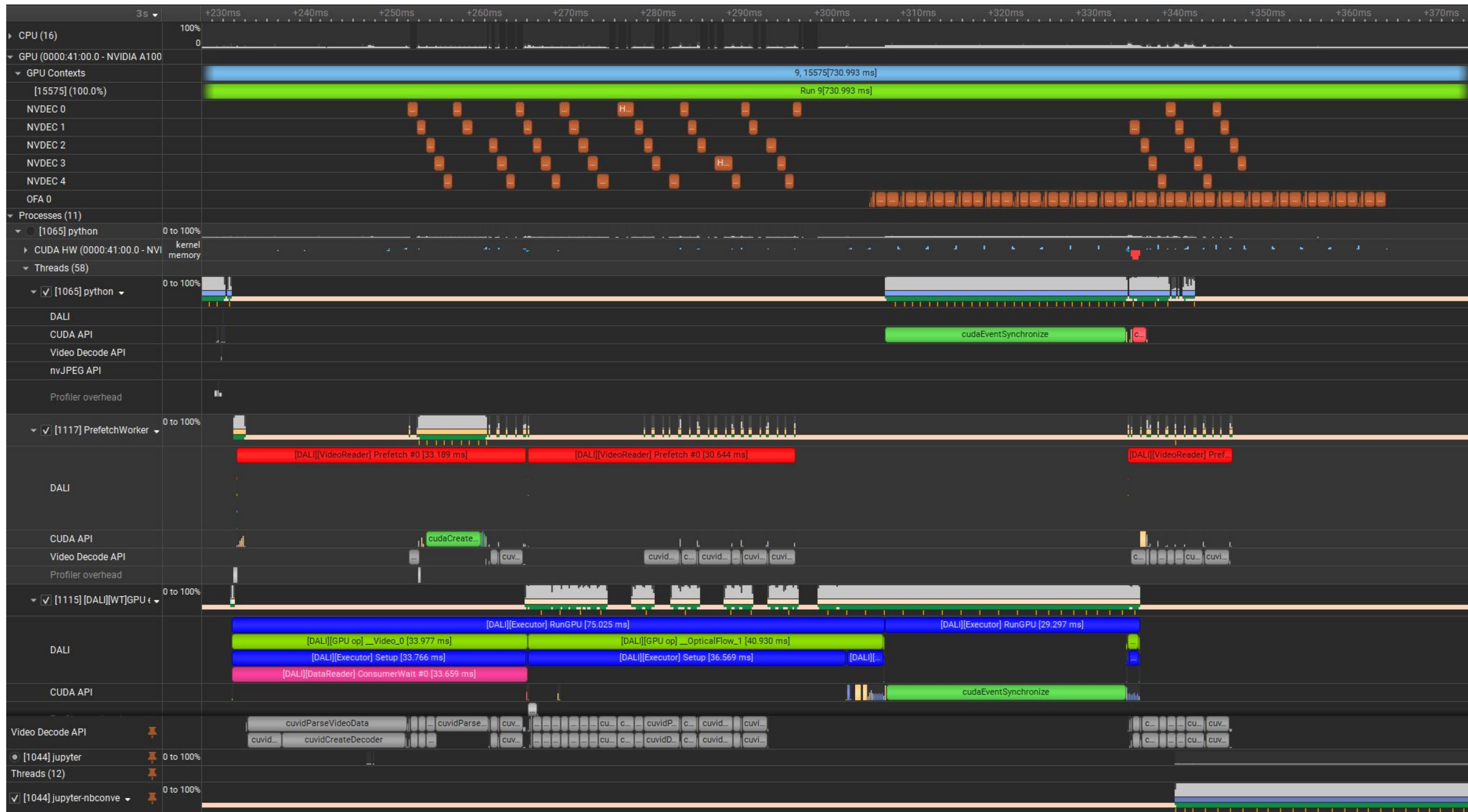
Dynamic function trace

Config file to target functions



ENC, DEC, OFA, & JPG Hardware Trace

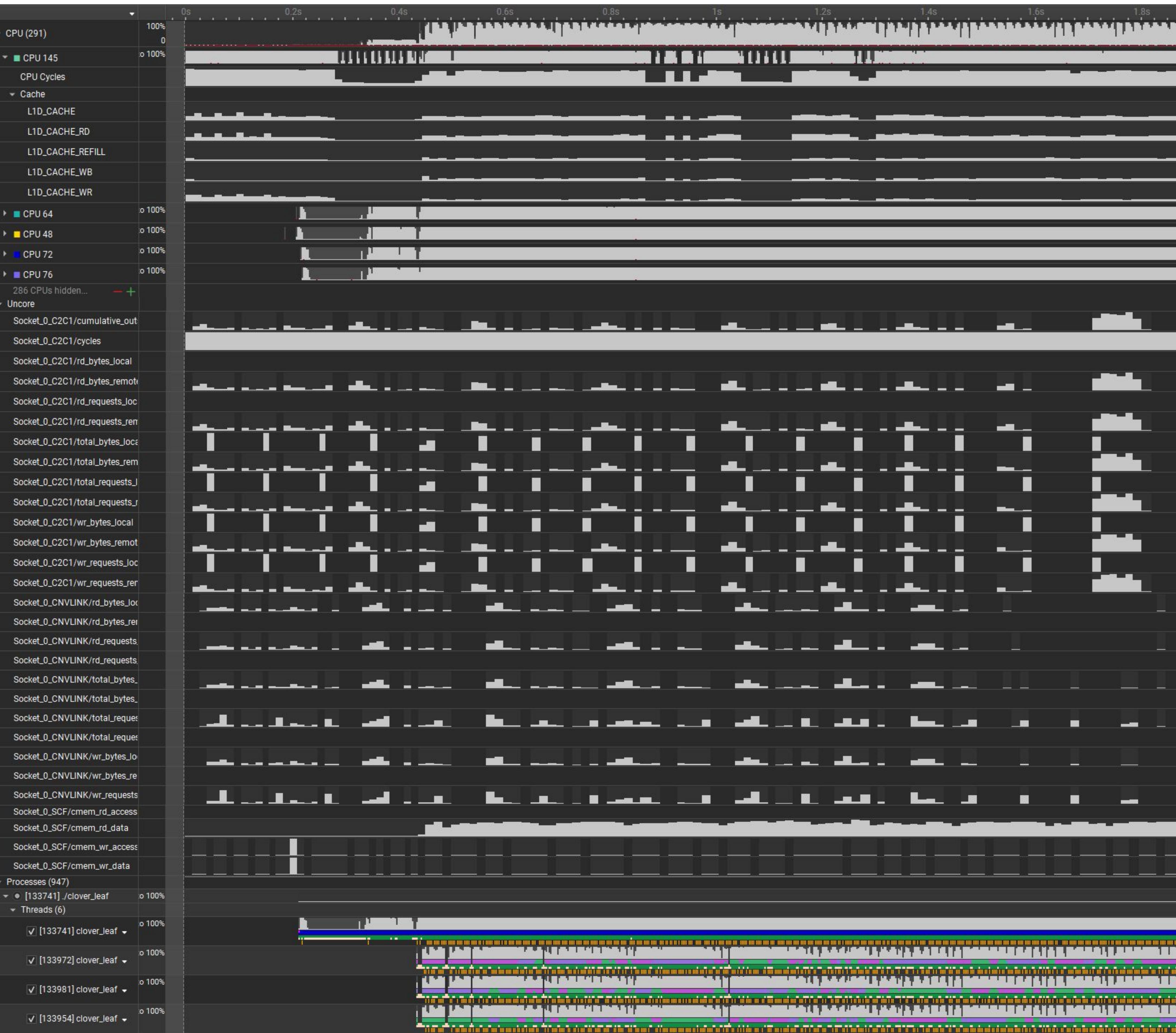
Example NVIDIA Data Loading Library (DALI)



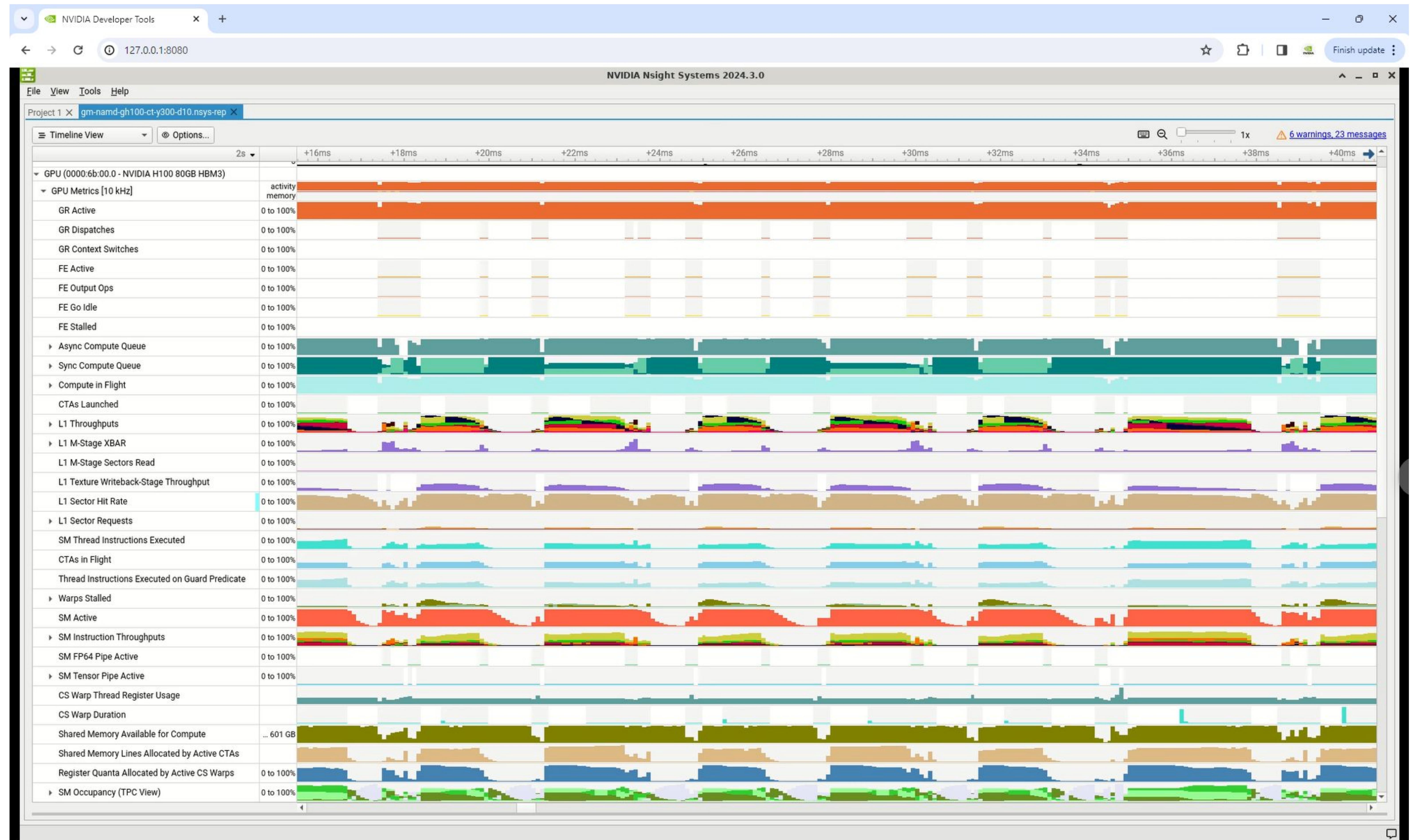
NVIDIA Grace CPU Performance Counters Sampling

Sampling Grace counters & metrics

See data from
Cores
Uncore
Coherency
Caches
Buses
Memory
NVLink-C2C
...



H100 Compute Metrics-Set (Coming Soon)





Wrapping up

- Meeting developers where they are developing
- More visibility into critical path & utilization
- Paint the pictures to:
 - Swiftly guide your investigation
 - Ultimately get the most from your hardware

DEVELOPER TOOLS ACROSS GTC

Sessions

- [S62388](#): Achieving Higher Performance From Your Data Center and Cloud Application
- [S62256](#): Demystify CUDA debugging and performance with powerful developer tools
- [SE62128](#): Exploring AI-Assisted Developer Tools for Accelerated Computing
- [S62398](#): Advances in Ray Tracing Developer Tools

Labs

- [DLIT61667](#): Profilers, Python, and Performance: Nsight Tools for Optimizing Modern CUDA Workloads

Connect with the Experts

- [CWE61532](#): What's in Your CUDA Toolbox? CUDA Profiling, Optimization, and Debugging Tools
- [CWE61581](#): Using Nsight Graphics Tools to Transform Your Application to a Next-Gen Powerhouse
- [CWE61231](#): Connect With the Experts: GPU Compute Performance Analysis and Optimizations

Demos

Come and visit the Developer Tools pod during show floor hours!

The latest version of Nsight Systems is available at:

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

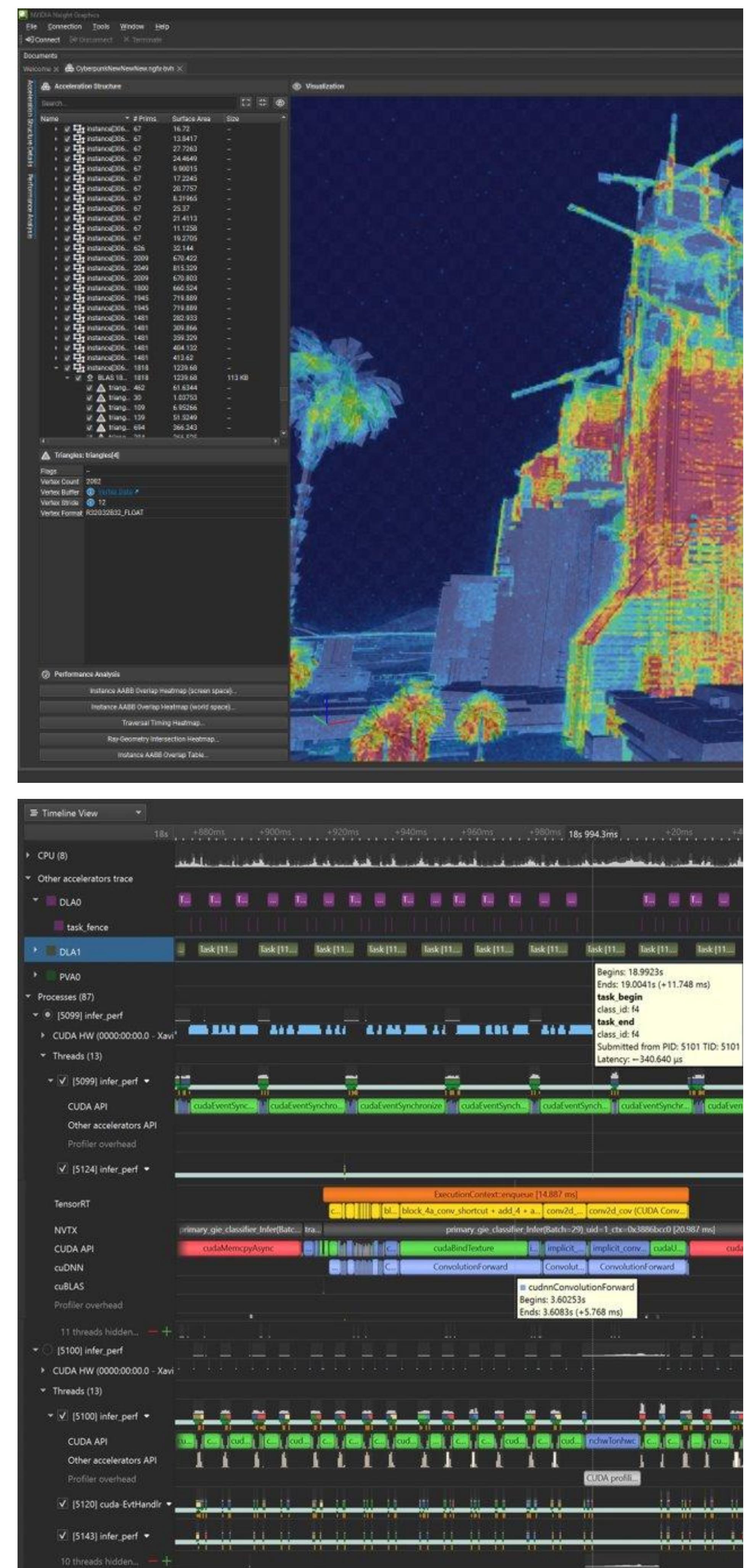
More developer tools at:

<https://developer.nvidia.com/tools-overview>

Support is available at:

<https://forums.developer.nvidia.com/c/developer-tools>

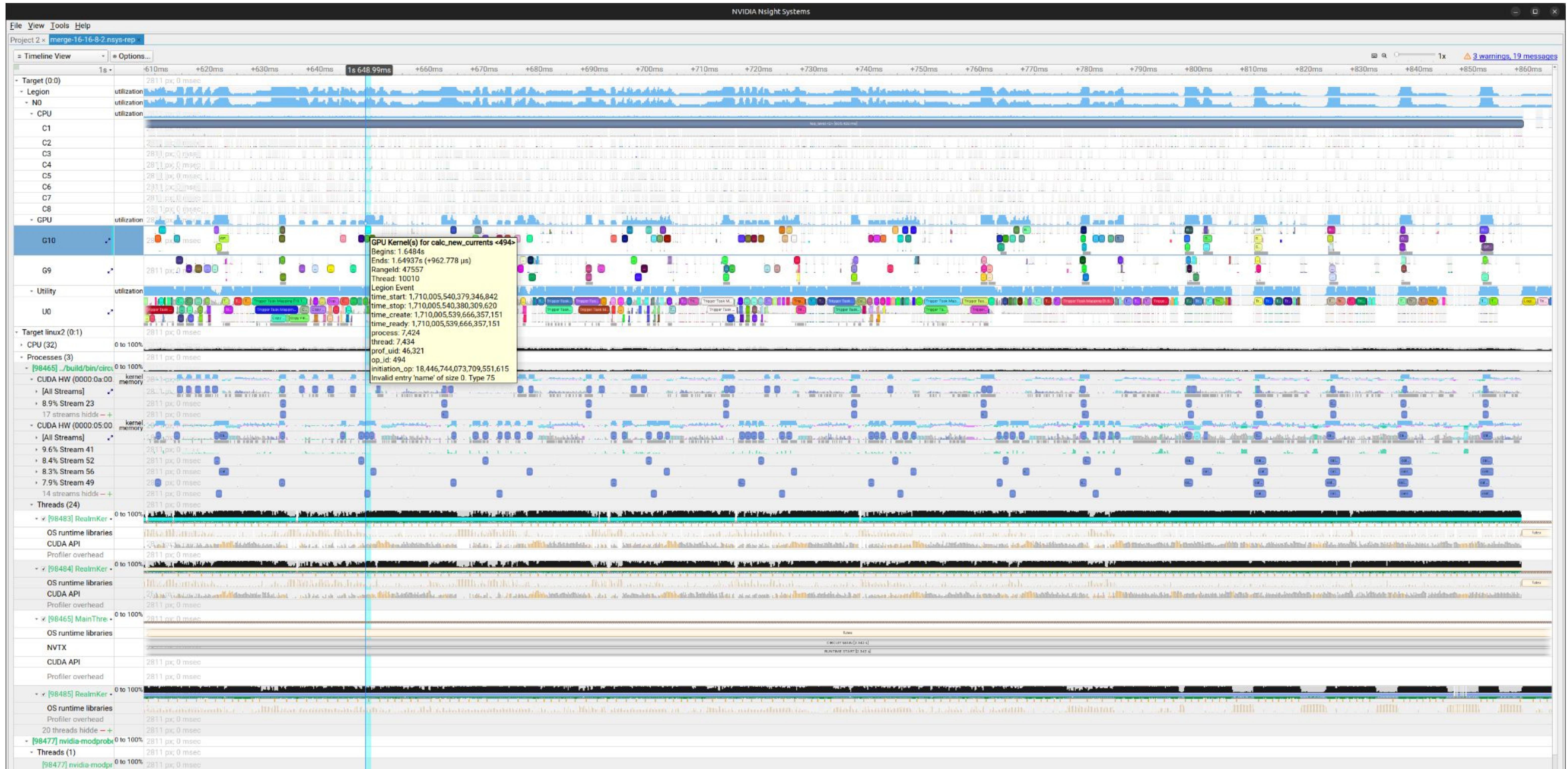
Interested in working on Developer Tools? We are hiring! Scan the QR code 





Legate

Unified timeline of Legate high-level logical task decomposition alongside low-level hardware, OS, and software details



Logical view:
Tasks, regions,
scheduling,
dependencies

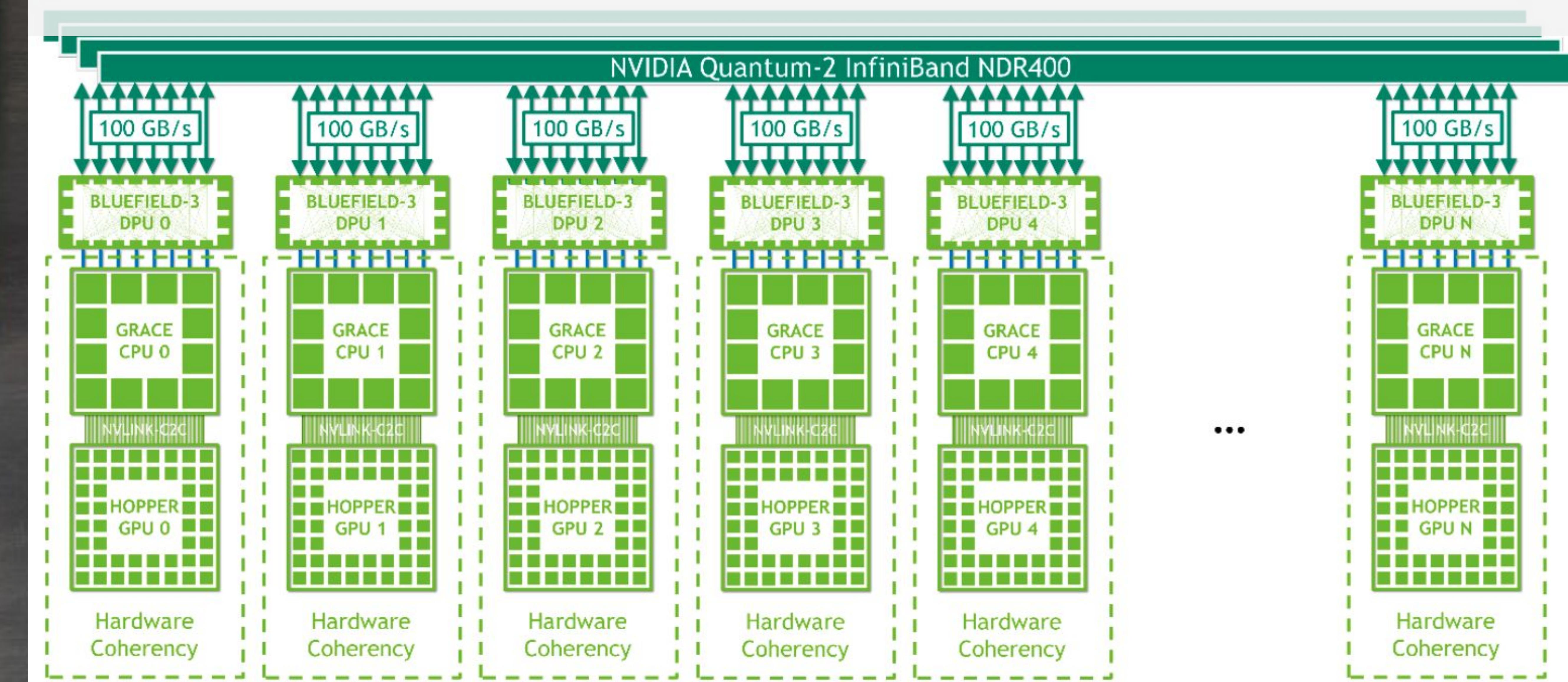
Physical view:
CPU, GPU, NIC,
OS, processes,
threads, libraries

...





Scale to 1000s of GPUs



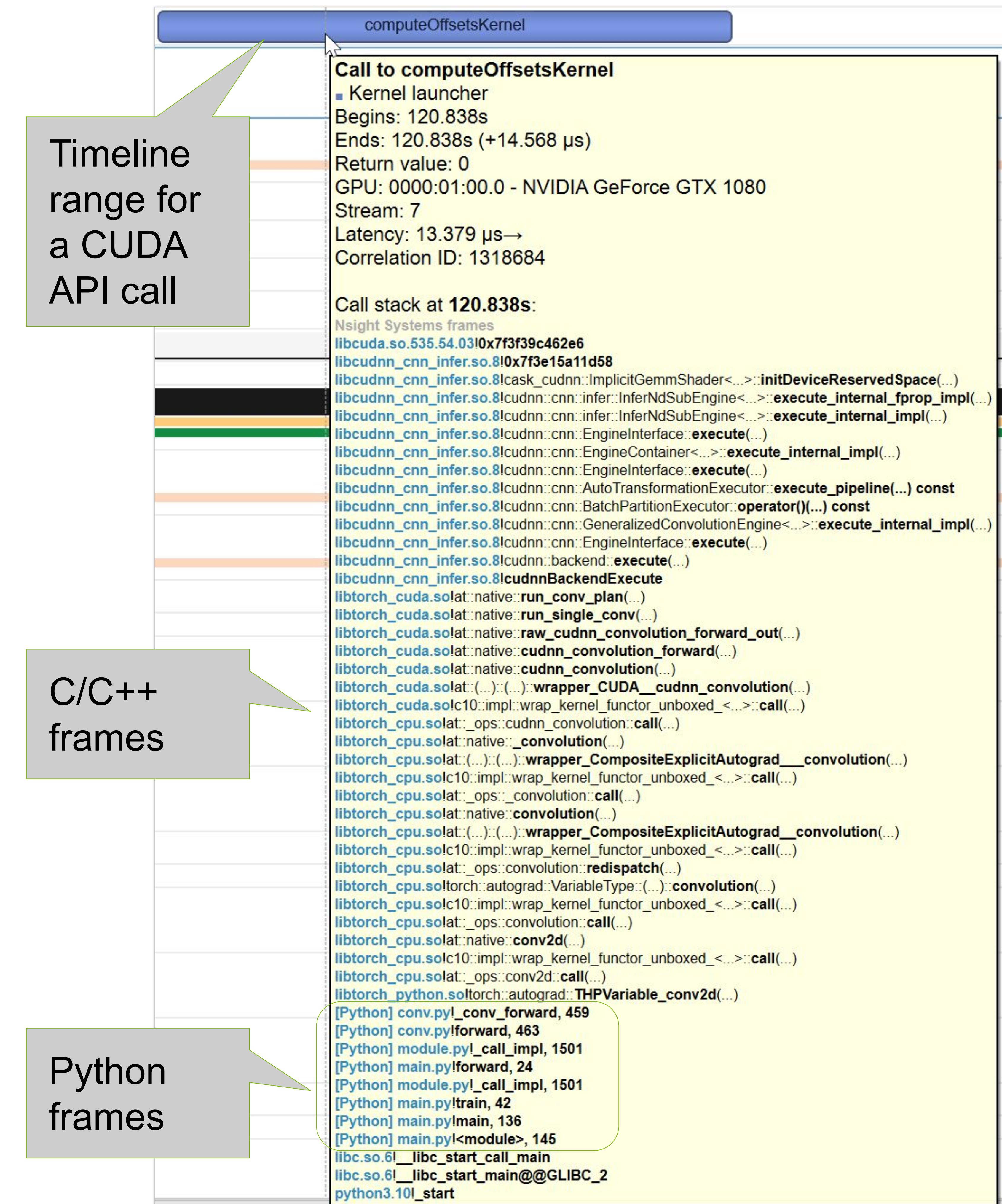
More Resources

- GTC Talks
 - 2023 - [Optimizing at Scale: Investigating Hidden Bottlenecks for Multi-Node Workloads](#)
 - 2022 - [Optimizing Communication with Nsight Systems Network Profiling](#)
 - 2021 - [Tuning GPU Network and Memory Usage in Apache Spark](#)
 - 2020 - [Rebalancing the Load: Profile-Guided Optimization of the NAMD Molecular Dynamics Program for Modern GPUs using Nsight Systems](#)
 - 2019 - [Using Nsight Tools to Optimize the NAMD Molecular Dynamics Simulation Program](#)
 - 2018 - [Optimizing HPC Simulation and Visualization Codes Using NVIDIA Nsight Systems](#)
 - 2018 - [Boost DNN Training Performance using NVIDIA Tools](#)
- DLI Labs
 - [Fundamentals of Accelerated Computing with CUDA C/C++](#)
 - [Optimizing CUDA Machine Learning Codes With Nsight Profiling Tools](#)
- More presentations, videos, blog posts, ...
 - https://docs.nvidia.com/nsight-systems/UserGuide/index.html#other_resources

Python CUDA API Backtrace

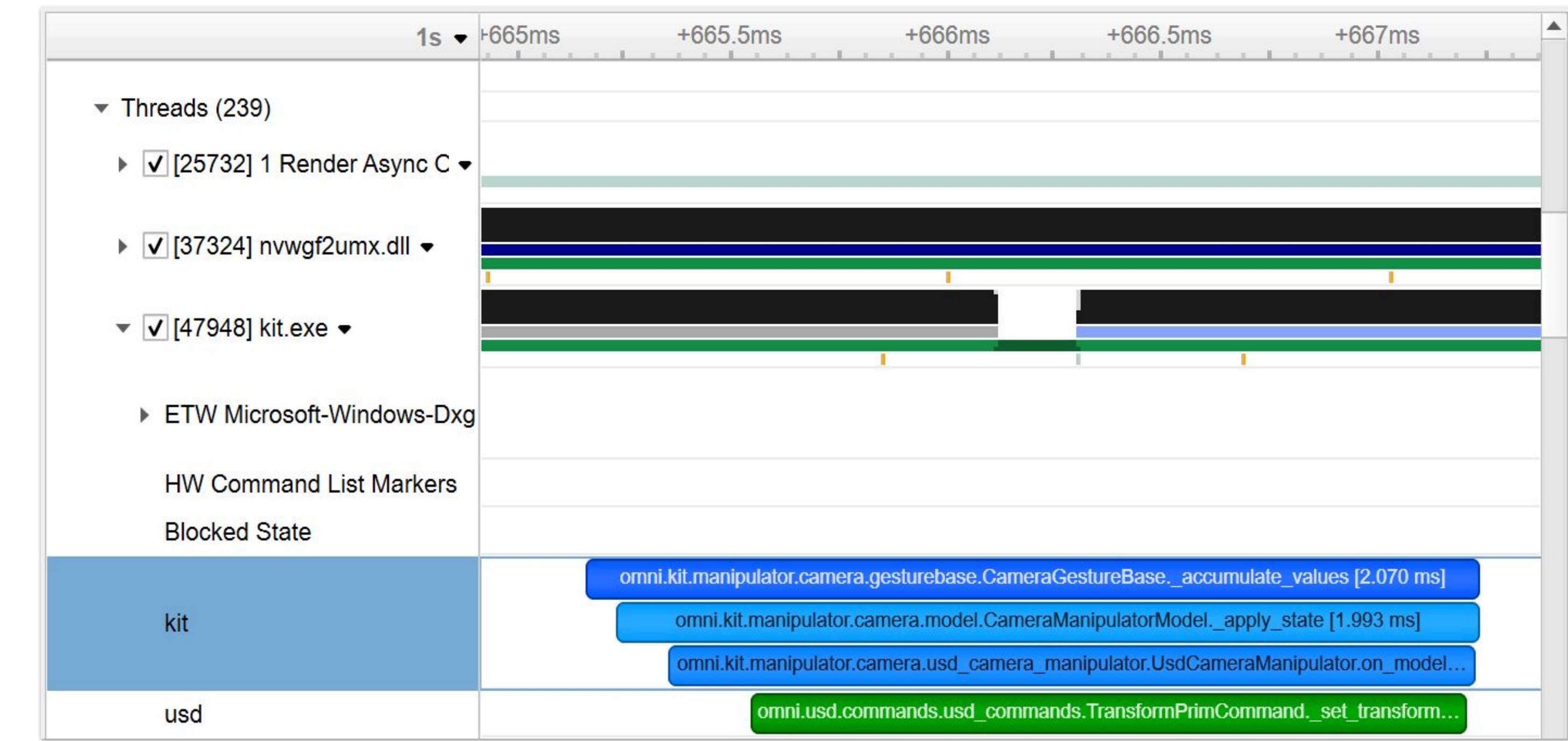
Combined backtrace of C/C++ and Python frames

- Extends the CUDA API Trace feature
- Time threshold
- Notes:
 - Requires CUDA trace
 - Target application must be launched by Nsys



Python Dynamic Annotations

- JSON config file
 - which functions
 - NVTX placement info
- No change to source code!
- Low overhead.
- Notes:
 - Target application must be launched by Nsys



```
[  
 {  
   "domain": "usd",  
   "color": "0x008000",  
   "module": "omni.usd.commands.usd_commands",  
   "functions": ["TransformPrimCommand._set_transform_matrix"]  
 },  
 {  
   "domain": "kit",  
   "module": "omni.kit.manipulator.camera.usd_camera_manipulator",  
   "color": "0x0080FF",  
   "functions": ["UsdCameraManipulator.on_model_updated"]  
 },  
 {  
   "domain": "kit",  
   "module": "omni.kit.manipulator.camera.gesturebase",  
   "color": "0x0055FF",  
   "functions": ["CameraGestureBase._accumulate_values"]  
 },  
 {  
   "domain": "kit",  
   "module": "omni.kit.manipulator.camera.model",  
   "color": "0x009cff",  
   "functions": ["CameraManipulatorModel._apply_state"]  
 }]  
 ]
```

Traditional Job & Resource Scheduler

SLURM

Profile on every task/rank and keep each file unique in network storage

Add special features to one task/rank per node

```
NSYS_RANK_OPT="--trace cuda,nvtx ... --output /netstorage/%q{SLURM_PROCID}.nsys-rep"
NSYS_NODE_OPT=$(if [ "$SLURM_LOCALID" == "0" ]; then echo "--nic-sample=true --gpu-metrics-device=all"; fi)

NSYS_SESSION=job%q{SLURM_JOB_ID}_q{SLURM_LOCALID}

srun nsys profile --session-new $NSYS_SESSION $NSYS_RANK_OPT $NSYS_NODE_OPT <your-app> <args>
```

The yellow above is for optional interactive control, paired with features like --start-later, --delay, --duration, etc

```
srun --jobid <jobId> nsys start --session $NSYS_SESSION ...
srun --jobid <jobId> nsys stop --session $NSYS_SESSION ...
```

If pyxis plugin, slightly simpler session name

```
NSYS_SESSION=job%q{SLURM_JOB_ID}
```

RAPIDS & Dask

DaskMPI (or custom SLURM or MPI scripts)

Start on every workers (optionally delayed until manually started)

Example with mpirun, easily translatable to SLURM srun/sbatch

```
mpirun|srun nsys profile <nargs> dask-mpi --scheduler-file <schedOnNet>.json
```

If interactive start/stop

Again, add `--start-later` & `--session-new` with SLURM_LOCALID or OMPI_COMM_WORLD_LOCAL_RANK

DaskSSH

Need a helper script to provide a local ID if more than one worker per OS at receiving IP address+port (ie not containerized)

Interactive Start/Stop

Options for coordinated start|stop

- a) Dask API: client.run()
- b) Tool like ‘parallel-ssh’ & address all local IDs

OpenMPI-Compatible Scheduler

Alternative Environment Variables

Profile on every rank and keep each file unique

```
srun nsys profile --output /netstorage/<user>/<jobname>/$OMPI_COMM_WORLD_RANK_$HOSTNAME.nsys-rep <app> <args>
```

SLURMD_NODEID, SLURMD_NODENAME, SLURM_JOBID, SLURM_JOB_NAME

- Add special features to rank 0

```
if [ "$OMPI_COMM_WORLD_LOCAL_RANK" == "0" ]
```

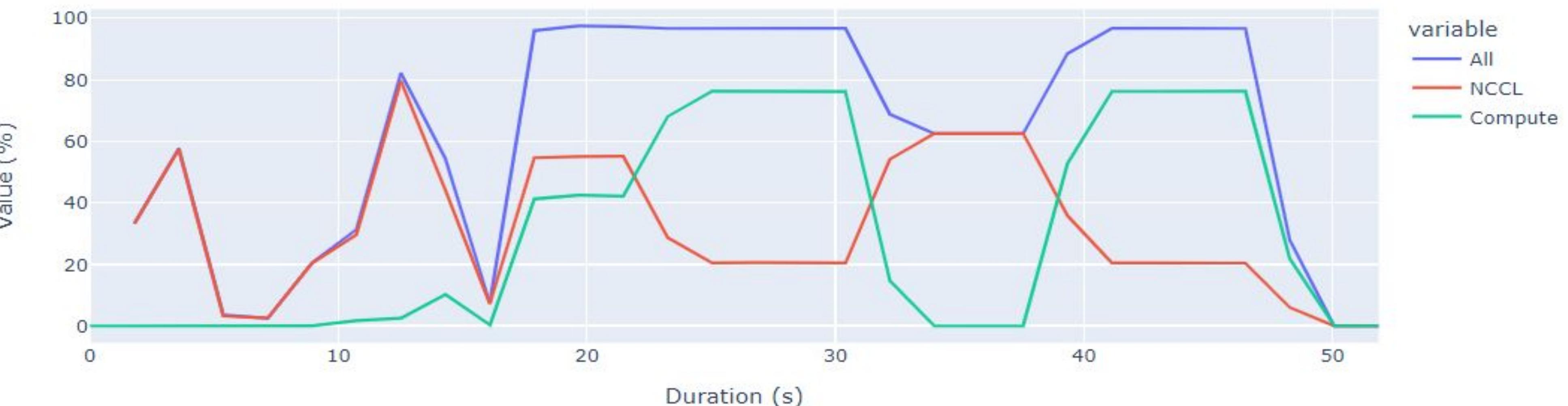
Compute-Only Zones

- GPU Utilization

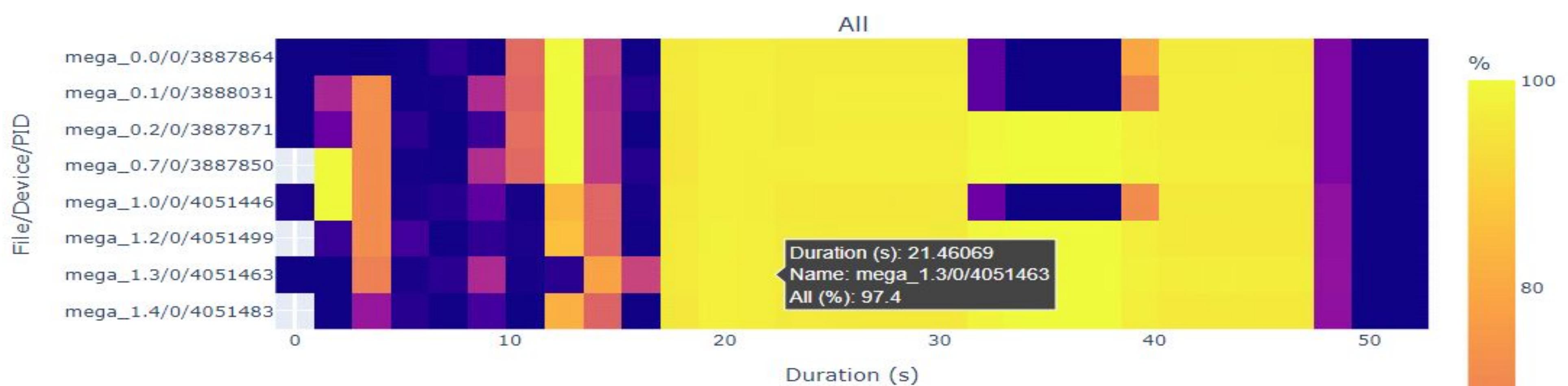
- Communications

- Compute-Only

GPU Utilization Summary (bins=30)

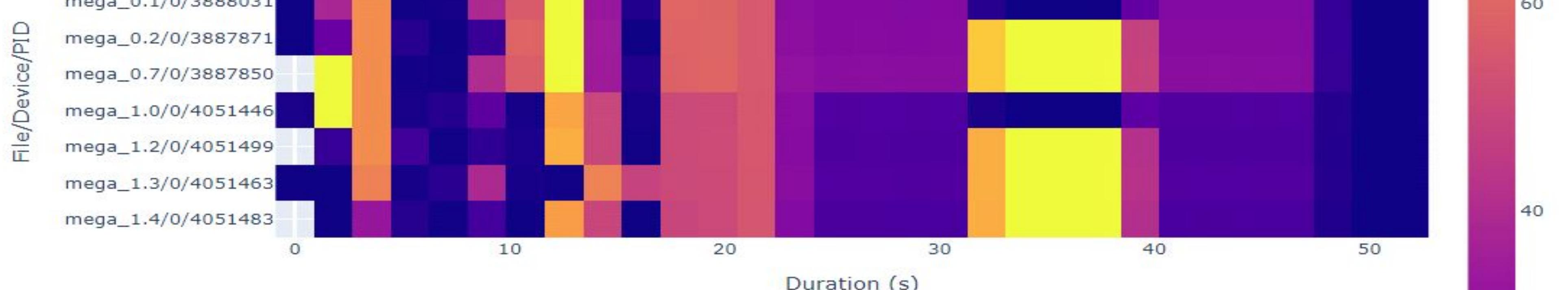


GPU Utilization (bins=30)

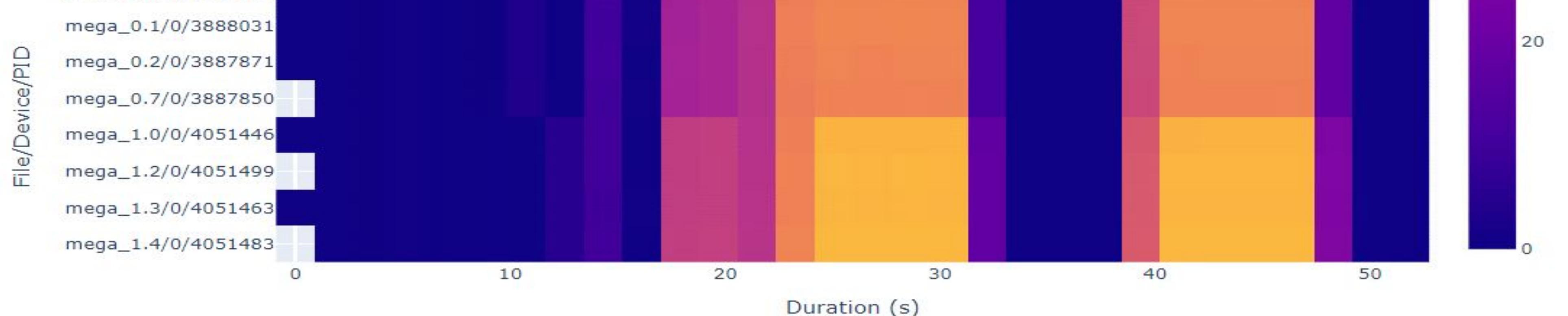


All

Duration (s): 21.46069
Name: mega_1.3/0/4051463
All (%): 97.4



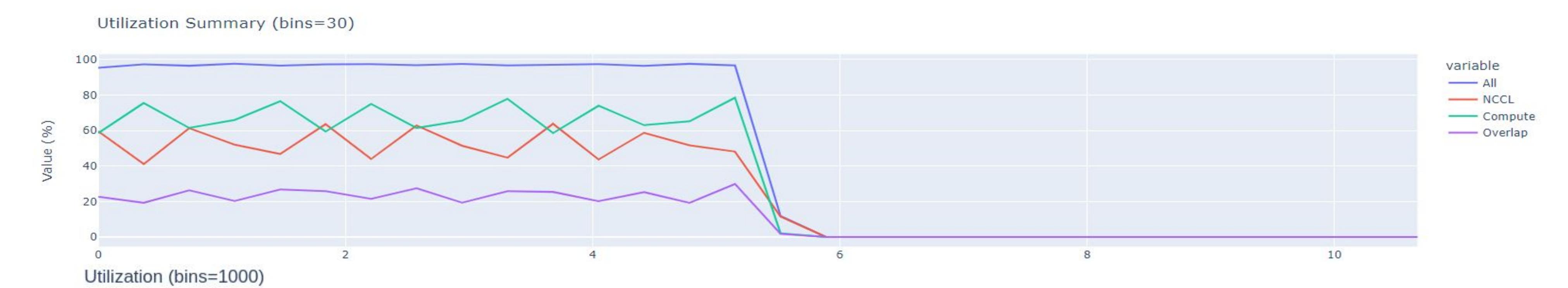
NCCL



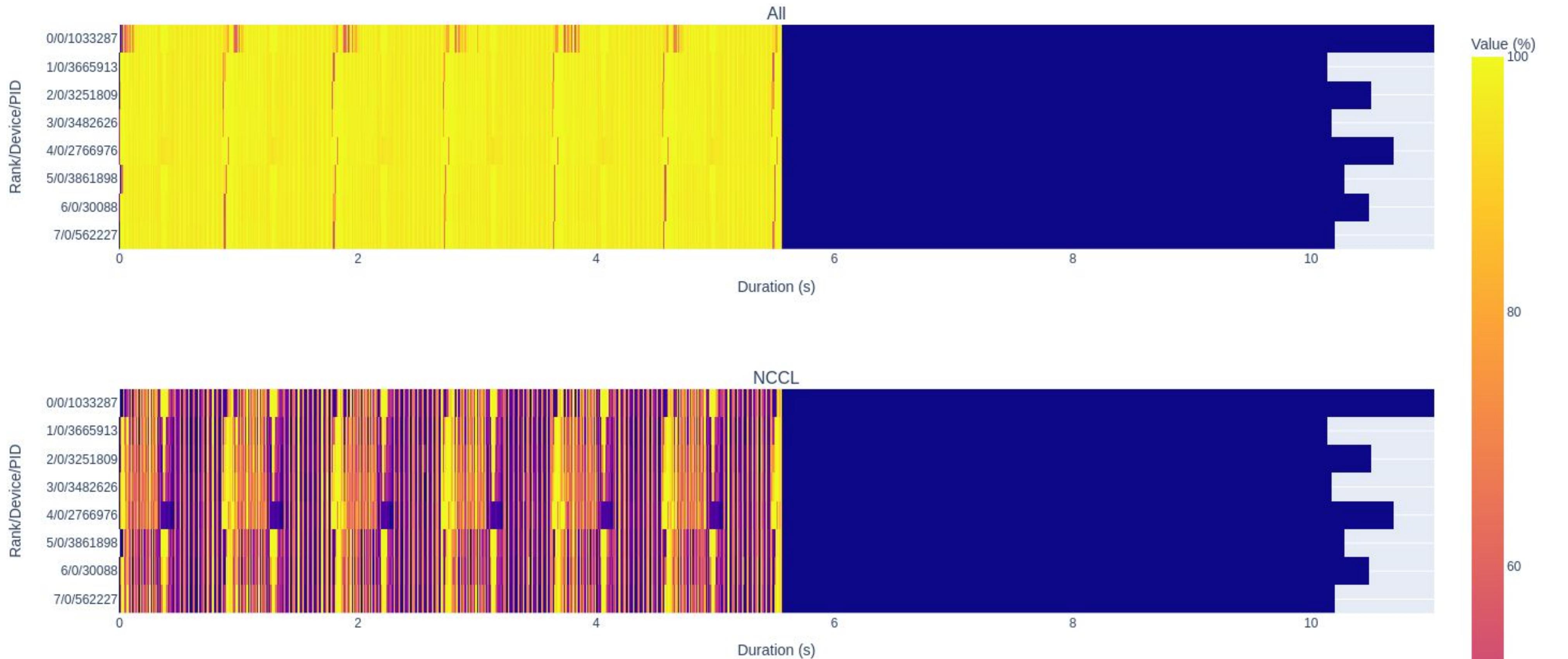
Compute

Compute-Comms Overlap

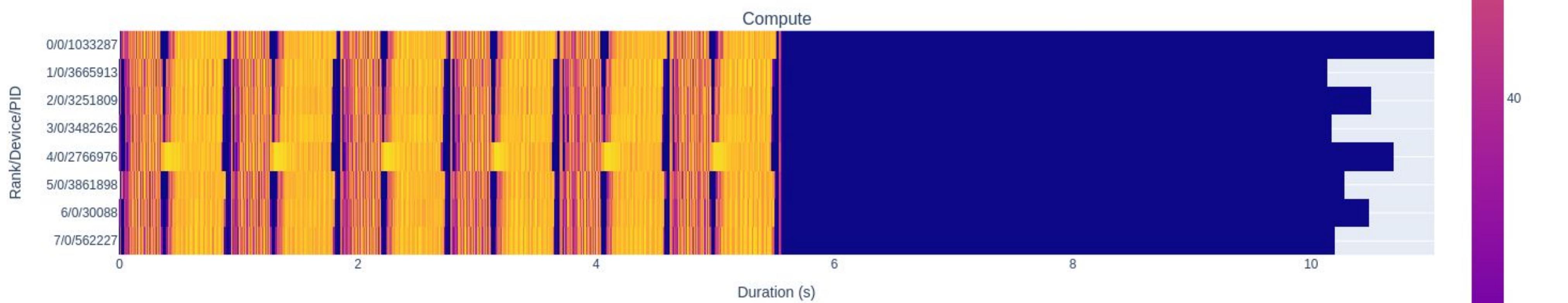
- GPU Utilization



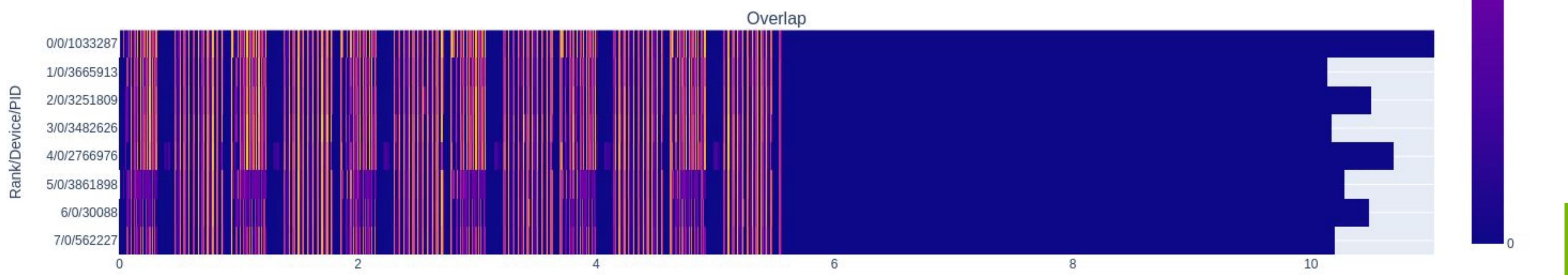
- Communications



- Compute



- Overlap



Scaling Out Multi-Report Analysis



Best Practices for Large Scale Distributed Training

Step by step guides to create clusters:

- One-click VPC Deployments
- Mount FSx for Lustre Filesystems
- EFA enabled clusters

Recipes to customize AMIs

AWS optimized Dockerfiles

EFA cheatsheet

Distributed Training Examples:

- Slurm scripts/K8 manifests
- Working with Pyxis/Enroot
- Nemo (MegatronLM, Multimodal, Bionemo)
- MosaicML
- DDP, FSDP
- SMDP, SMMP
- Tensorflow/Jax

Validation (NCCL Tests etc)

Observability (Prometheus-Grafana etc)

Profiling (Nsight Product Family)



Best Practices for Large Scale Distributed Training

Step by step guides to create clusters:

- One-click VPC Deployments
- Mount FSx for Lustre Filesystems
- EFA enabled clusters

Recipes to customize AMIs

AWS optimized Dockerfiles

EFA cheatsheet

Distributed Training Examples:

- Slurm scripts/K8 manifests
- Working with Pyxis/Enroot
- Nemo (MegatronLM, Multimodal, Bionemo)
- MosaicML
- DDP, FSDP
- SMDP, SMMP
- Tensorflow/Jax

Validation (NCCL Tests etc)

Observability (Prometheus-Grafana etc)

Profiling (Nsight Product Family)

