

October 10-12, 2023

---



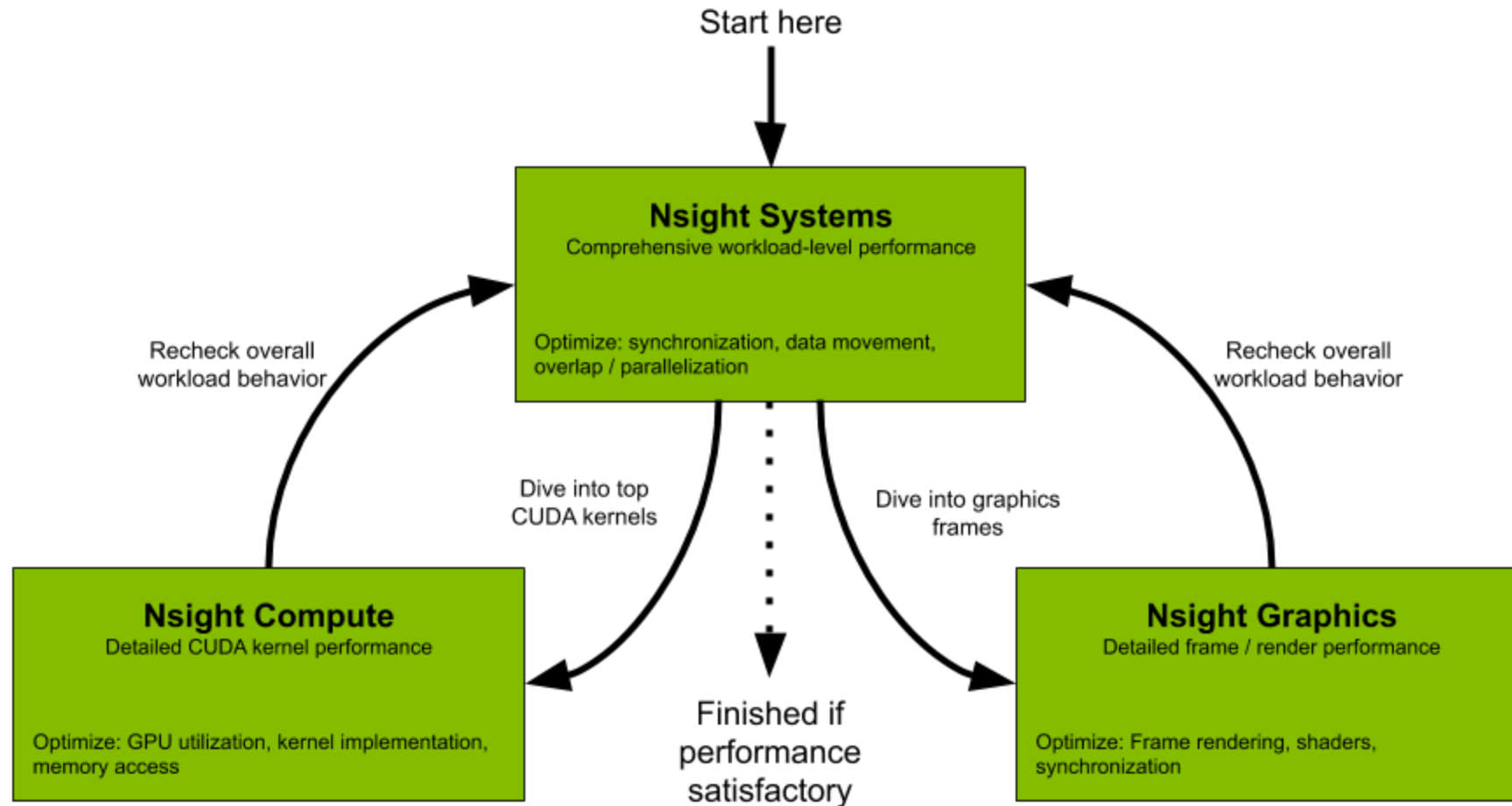
# ALCF Hands-on HPC Workshop

# Nsight Tools

Nsight Systems and Nsight Compute

Matt Stack, NVIDIA

# Overview

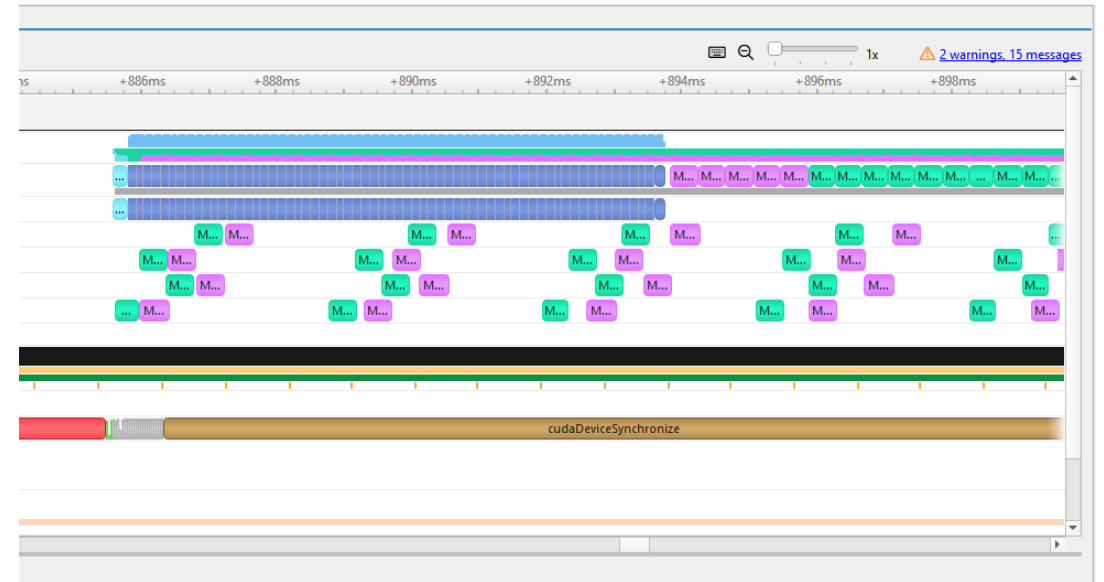




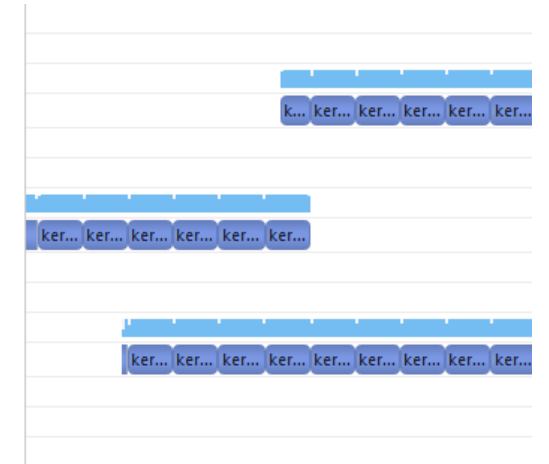
# NSIGHT SYSTEMS

## SYSTEM PROFILER

- System-wide view of application
  - Single timeline with GPU, CPU information
  - MPI, Memory patterns, multi-GPU
- Identify spots for performance tuning
  - Look for GPU idling, excessive synchronization, suboptimal memory traffic
- ALCF Nsight docs: <https://docs.alcf.anl.gov/theta-gpu/performance-tools/nvidia-nsight/>
- Nsight Systems User guide from Nvidia: <https://docs.nvidia.com/nsight-systems/UserGuide/index.html>
- Available on Linux (x86, Power9, ARM), Windows, and Mac (Host only)



- ▼ Processes (30)
  - ▼ [2433679] ./test
    - ▼ CUDA HW (0004:04:00.0 - Tesla V100-SXM2-16GB)
      - ▶ 100.0% Kernels
      - ▶ Threads (4)
    - ▼ [2433680] ./test
      - ▼ CUDA HW (0004:04:00.0 - Tesla V100-SXM2-16GB)
        - ▶ 100.0% Kernels
        - ▶ Threads (4)
      - ▼ [2433681] ./test
        - ▼ CUDA HW (0004:04:00.0 - Tesla V100-SXM2-16GB)
          - ▶ 100.0% Kernels
          - ▶ Threads (4)
        - ▼ [2433682] ./test
          - ▼ CUDA HW (0004:04:00.0 - Tesla V100-SXM2-16GB)



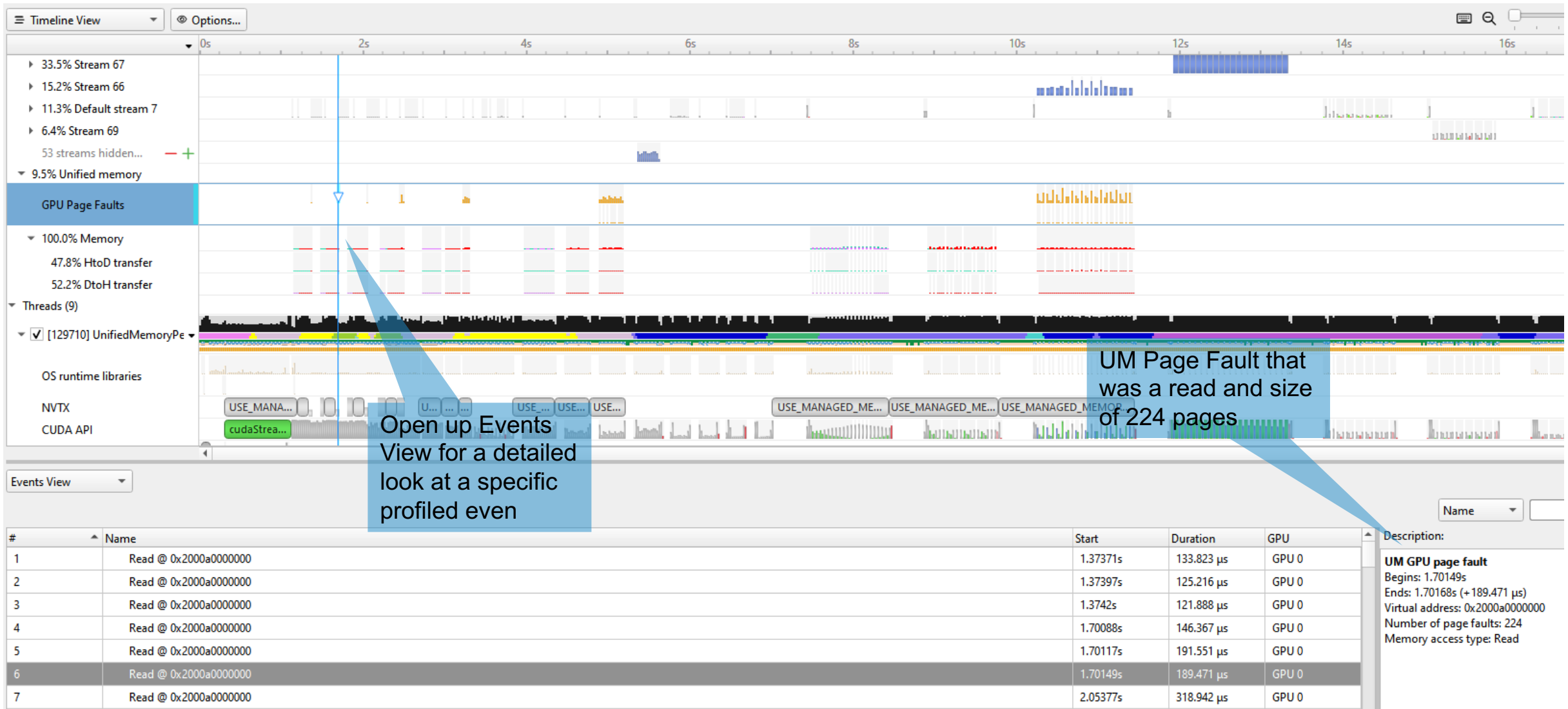
Time ->



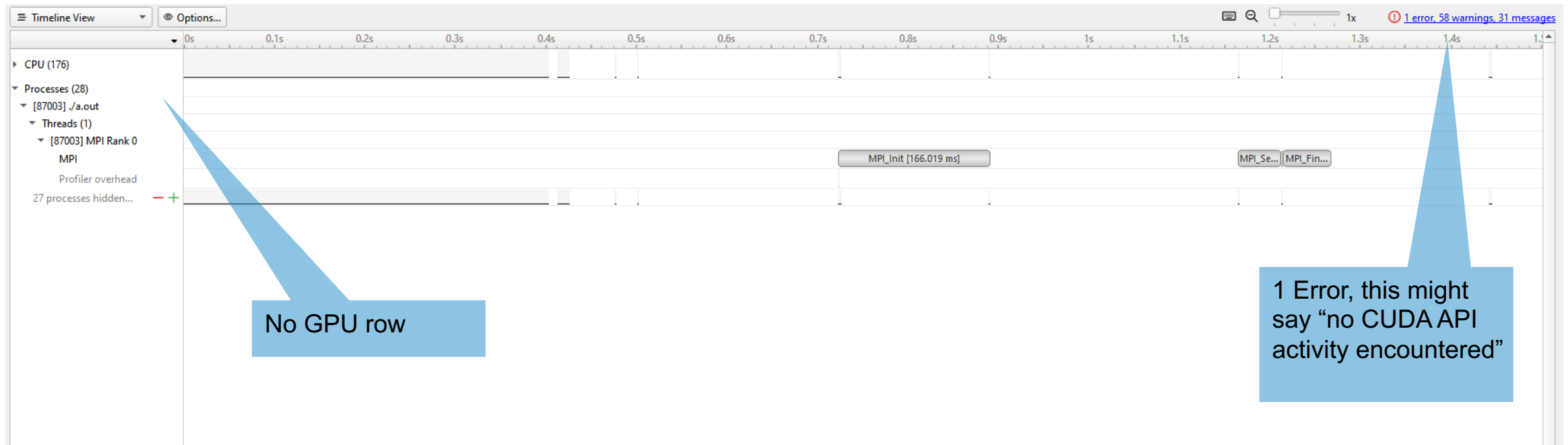
Device row, V100

Track memory transfers, and UM page faulting

NXTV ranges, user defined (Or included in your libraries)



- Trial and error, your first profile might look like this.. What went wrong?



# How to run Nsight Systems on Polaris

- ALCF Beginners Guide: [https://github.com/argonne-lcf/ALCFBeginnersGuide/blob/master/polaris/02\\_profiling.md](https://github.com/argonne-lcf/ALCFBeginnersGuide/blob/master/polaris/02_profiling.md)

- Generally, Nsight Systems command looks like

```
nsys profile -o <profile_name> --stats true ./<app exe>
```

(+ any additional flags, here are two useful ones:

- 1) If using OpenACC or OpenMP, consider `--trace openacc,openmp,cuda,nvtx,osrt`. By default nsys traces API calls from CUDA, NVTX, opengl, and osrt),
- 2) If using UM, consider adding `--cuda-um-gpu-page-faults true` and `--cude-um-cpu-page-faults true`

```
stack6@polaris-login-04:~> ml load nvhpc/23.3
The following have been reloaded with a version change:
 1) nvhpc/21.9 => nvhpc/23.3
stack6@polaris-login-04:~> which nsys
/opt/nvidia/hpc_sdk/Linux_x86_64/23.3/compilers/bin/nsys
stack6@polaris-login-04:~> nsys --version
NVIDIA Nsight Systems version 2023.1.1.127-32365746v0
```

```
stack6@polaris-login-04:~> ml load cudatoolkit-standalone/12.0.0
stack6@polaris-login-04:~> which nsys
/soft/compilers/cudatoolkit/cuda-12.0.0/bin/nsys
stack6@polaris-login-04:~> nsys --version
NVIDIA Nsight Systems version 2022.4.2.18-32044700v0
```

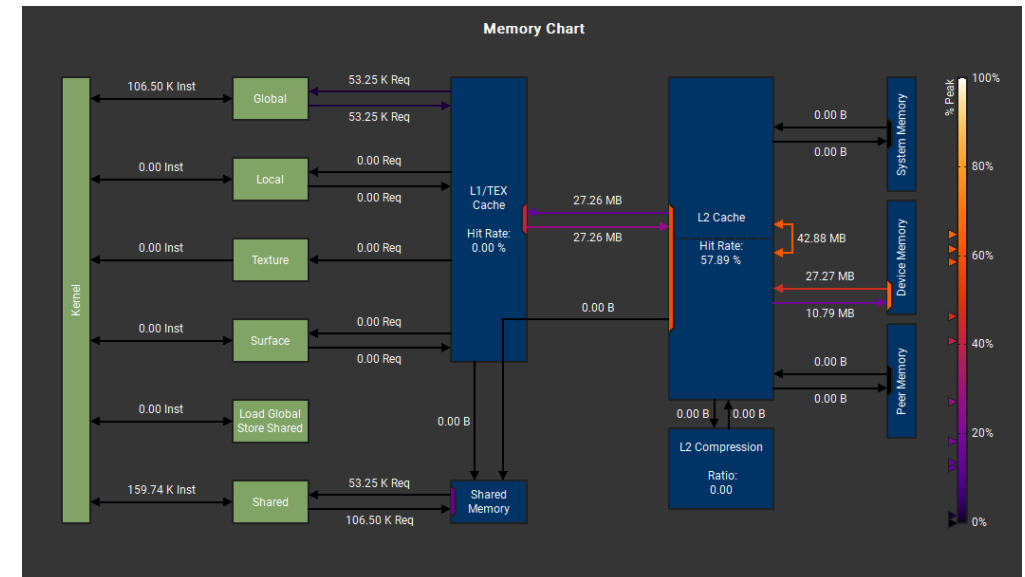
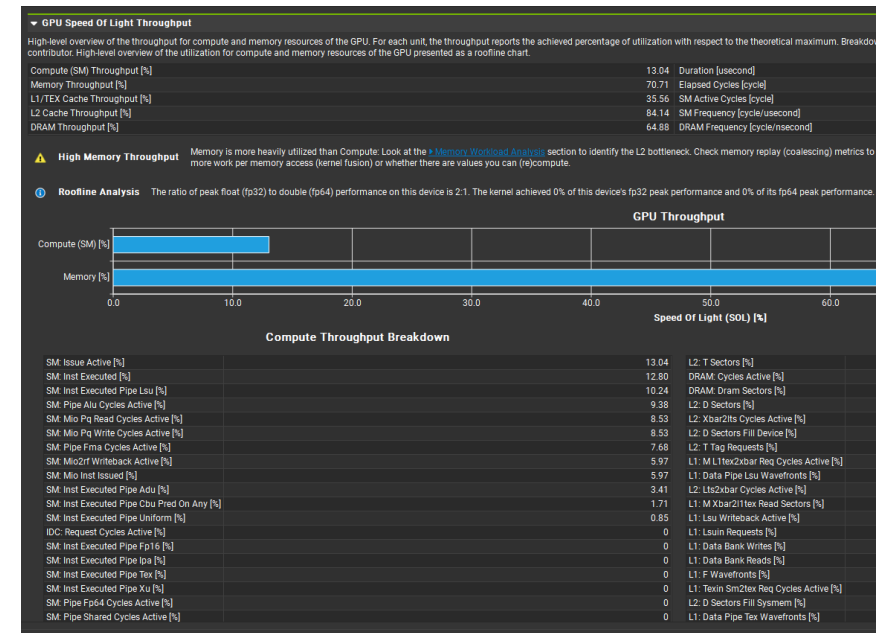




# NSIGHT COMPUTE

## KERNEL PROFILING TOOL

- “Zoomed in” profiling on a CUDA Kernel
  - Typically after Nsight Systems
- Compare multiple kernels with delta performance difference
  - “Did my code optimization make a change in the right direction?”
- View your kernel code with debugging information
  - Opt-in source code profiling
- Nsight Compute Docs: <https://docs.nvidia.com/nsight-compute/NsightCompute/index.html>
- Available on Linux (x86, Power9, ARM), Windows, and Mac (Host only)



Page: Details Result: 0 - 512 - MatrixMulCUDA Add Baseline Apply Rules Occupancy Calculator Copy as Image

**Current** 512 - MatrixMulCUDA (20, 10, 1)x(32, 32, 1) 956.80 usecond 1,046,054 32 0 - NVIDIA A100-SXM4-40GB 1.09 cycle/nsecond 8.0 [3578520] matrixMul

The report contains imported source files.

**GPU Speed Of Light Throughput**

High-level overview of the throughput for compute and memory resources of the GPU. For each contributor, the achieved percentage of utilization with respect to the theoretical maximum. Breakdowns show the throughput for each individual sub-metric of Compute and Memory to clearly identify the highest contributor. High-level overview of the utilization for compute and memory resources of the GPU presented as a roofline chart.

Metric	Value	Unit
Compute (SM) Throughput [%]	83.39	Duration [usecond]
Memory Throughput [%]	3.93	Elapsed Cycles [cycle]
L1/TEX Cache Throughput [%]	4.24	SM Active Cycles [cycle]
L2 Cache Throughput [%]	0.33	SM Frequency [cycle/nsecond]
DRAM Throughput [%]	0.08	DRAM Frequency [cycle/nsecond]

**High Throughput** The kernel is utilizing greater than 80.0% of the available compute or memory performance of the device. To further improve performance, work will likely need to be shifted from the most utilized to another unit. Start by analyzing workloads in the [Compute Workload Analysis](#) section.

**Roofline Analysis** The ratio of peak float (fp32) to double (fp64) performance on this device is 2:1. The kernel achieved close to 1% of this device's fp32 peak performance and 0% of its fp64 peak performance. See the [Kernel Profiling Guide](#) for more details on roofline analysis.

### GPU Throughput

**Compute Throughput Breakdown**

Metric	Value
SM: Pipe Alu Cycles Active [%]	83.39
SM: Issue Active [%]	47.53
SM: Inst Executed [%]	47.52
SM: Inst Executed Pipe Adu [%]	15.79
SM: Mio Pq Write Cycles Active [%]	7.89
SM: Mio Pq Read Cycles Active [%]	7.89
SM: Mio Inst Issued [%]	5.89
SM: Inst Executed Pipe Lsu [%]	3.89
SM: Mio2rf Writeback Active [%]	1.90
SM: Pipe Fma Cycles Active [%]	1.90
SM: Inst Executed Pipe Cbu Pred On Any [%]	0.12
IDC: Request Cycles Active [%]	0.06
SM: Inst Executed Pipe Tex [%]	0.01
SM: Inst Executed Pipe Fp16 [%]	0
SM: Inst Executed Pipe Ipa [%]	0
SM: Inst Executed Pipe Uniform [%]	0
SM: Inst Executed Pipe Xu [%]	0
SM: Pipe Fp64 Cycles Active [%]	0

**Memory Throughput Breakdown**

Metric	Value
L1: Data Pipe Lsu Wavefronts [%]	3.93
L1: Lsuin Requests [%]	3.89
L1: Lsu Writeback Active [%]	3.76
L1: Data Bank Reads [%]	1.00
L2: Lts2xbar Cycles Active [%]	0.33
L2: T Sectors [%]	0.26
L2: T Tag Requests [%]	0.24
L1: M Xbar2l1tex Read Sectors [%]	0.23
L2: Xbar2lts Cycles Active [%]	0.22
L2: D Sectors [%]	0.21
L1: M L1tex2xbar Req Cycles Active [%]	0.12
L1: Data Bank Writes [%]	0.12
DRAM: Cycles Active [%]	0.08
DRAM: Dram Sectors [%]	0.05
L2: D Sectors Fill Device [%]	0.05
L1: Texin Sm2tex Req Cycles Active [%]	0.01
L1: F Wavefronts [%]	0.00
L1: Data Pipe Tex Wavefronts [%]	0

**Note: we have source files**

**First place to look**

Page: Source Result: 0 - 512 - MatrixMulCUDA Add Baseline Apply Rules Occupancy Calculator Copy as Image

Result Time Cycles Regs GPU SM Frequency CC Process  
 Current 512 - MatrixMulCUDA (20, 10, 1)x(32, 32, 1) 956.80 usecond 1,046,054 32 0 - NVIDIA A100-SXM4-40GB 1.09 cycle/nsecond 8.0 [3578520] matrixMul

View: Source and SASS

Source: matrixMul.cu Find... Navigation: Instructions Executed

Navigation: Live Registers

Register Dependencies - new(-ish) feature!

Stall reasons, instructions executed, hot spots

# Source	Live Registers	irp Stall (All Cycles)	Sampling (Not-Issued Cycles)	# Address	Source	Live Registers	Register Dependencies	Predicate Dependencies	Unit
105				1	00007f78 7327ba00	MOV R1, c[0x0][0x28]	1		
106	27	1,828	1,048	2	00007f78 7327ba10	MOV R14, c[0x0][0x118]	2		
107	27	1,771	1,031	3	00007f78 7327ba20	MOV R15, c[0x0][0x11c]	3		
108				4	00007f78 7327ba30	MOV R2, 0x160	4		
109				5	00007f78 7327ba40	LDC.64 R2, c[0x0][R2]	5		
110	21	45	25	6	00007f78 7327ba50	MOV R10, R2	6		
111				7	00007f78 7327ba60	MOV R11, R3	6		
112				8	00007f78 7327ba70	MOV R10, R10	5		
113				9	00007f78 7327ba80	MOV R11, R11	5		
114				10	00007f78 7327ba90	MOV R2, 0x168	6		
115				11	00007f78 7327baa0	LDC.64 R2, c[0x0][R2]	7		
116				12	00007f78 7327bab0	MOV R4, R2	8		
117	23	14,279	6,609	13	00007f78 7327bac0	MOV R9, R3	8		
118	28	81,699	39,582	14	00007f78 7327bad0	MOV R4, R4	7		
119				15	00007f78 7327bae0	MOV R9, R9	7		
120				16	00007f78 7327baf0	MOV R2, 0x170	8		
121				17	00007f78 7327bb00	LDC.64 R2, c[0x0][R2]	9		
122				18	00007f78 7327bb10	MOV R5, R2	10		
123				19	00007f78 7327bb20	MOV R6, R3	10		
124	21	249	112	20	00007f78 7327bb30	MOV R5, R5	9		
125				21	00007f78 7327bb40	MOV R6, R6	9		
126				22	00007f78 7327bb50	MOV R0, 0x178	10		
127				23	00007f78 7327bb60	LDC R0, c[0x0][R0]	10		
128				24	00007f78 7327bb70	MOV R7, R0	11		
129	12	14	6	25	00007f78 7327bb80	MOV R0, 0x17c	11		
130	10	32	24	26	00007f78 7327bb90	LDC R0, c[0x0][R0]	11		
131				27	00007f78 7327bba0	MOV R8, R0	12		
				28	00007f78 7327bbb0	MOV R0, R10	12		
				29	00007f78 7327bbc0	MOV R2, R11	12		
				30	00007f78 7327bbd0	MOV R3, R4	12		
				31	00007f78 7327bbe0	MOV R4, R9	12		
				32	00007f78 7327bbf0	MOV R5, R5	11		
				33	00007f78 7327bc00	MOV R6, R6	11		
				34	00007f78 7327bc10	MOV R7, R7	11		

Inline Functions

#	Source	Inline Function Address	Live Registers	Instructions Executed	irp Stall (All Cycles)	Sampling (Not-Issued Cycles)

Don't ignore these Expert Analysis messages

#### ► 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 [%]	100	Block Limit Registers [block]	10
Theoretical Active Warps per SM [warp]	64	Block Limit Shared Mem [block]	21
Achieved Occupancy [%]	83.79	Block Limit Warps [block]	8
Achieved Active Warps Per SM [warp]	53.63	Block Limit SM [block]	32

#### ⚠ Occupancy Limiters

This kernel's theoretical occupancy is not impacted by any block limit. The difference between calculated theoretical (100.0%) and measured achieved occupancy (83.8%) can be the result of warp scheduling overheads or workload imbalances during the kernel execution. Load imbalances can occur between warps within a block as well as across blocks of the same kernel. See the [CUDA Best Practices Guide](#) for more details on optimizing occupancy.

#### ► Scheduler Statistics

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

Active Warps Per Scheduler [warp]	12.15	No Eligible [%]	48.69
Eligible Warps Per Scheduler [warp]	2.15	One or More Eligible [%]	51.31
Issued Warp Per Scheduler	0.51		

#### ⚠ Issue Slot Utilization

Every scheduler is capable of issuing one instruction per cycle, but for this kernel each scheduler only issues an instruction every 1.9 cycles. This might leave hardware resources underutilized and may lead to less optimal performance. Out of the maximum of 16 warps per scheduler, this kernel allocates an average of 12.15 active warps per scheduler, but only an average of 2.15 warps were eligible per cycle. Eligible warps are the subset of active warps that are ready to issue their next instruction. Every cycle with no eligible warp results in no instruction being issued and the issue slot remains unused. To increase the number of eligible warps, reduce the time the active warps are stalled by inspecting the top stall reasons on the [Warp State Statistics](#) and [Source Counters](#) sections.

# How to run Nsight Compute on Polaris

- ALCF Beginners Guide: [https://github.com/argonne-lcf/ALCFBeginnersGuide/blob/master/polaris/02\\_profiling.md](https://github.com/argonne-lcf/ALCFBeginnersGuide/blob/master/polaris/02_profiling.md)

- Generally, Nsight Compute command looks like

```
ncu -o <profile_name> -k <kernel_name> -c 1 ./<app exe>
```

(+ the addition of any additional flags, here are two useful ones:

- 1) `--import-source true` (IF `nvcc` had the flag `-lineinfo`)
- 2) `--set detailed` or `--set full` to get Memory Analysis section and Roofline, default is `--set basic`

# Compute Sanitizer

- Correctness checking tool
- Drop-in replacement for depreciated cuda-memcheck

Memcheck

Out-of-bounds  
and misaligned  
memory access

Racecheck

Race conditions  
for *shared*  
memory

Initcheck

Uninitialized  
device global  
memory detector

Synccheck

Thread  
synchronization  
hazard detector

Docs: <https://docs.nvidia.com/compute-sanitizer/ComputeSanitizer/index.html>

Compute-sanitizer --tool <which tool> ./app <app options>

# Resources

- Nsight Compute memory analysis deep dive: <https://www.nvidia.com/en-us/on-demand/session/gtcspring21-s32089/>
- Nsight Systems user guide: <https://docs.nvidia.com/nsight-systems/UserGuide/index.html>
- Nsight Compute CLI user guide: <https://docs.nvidia.com/nsight-compute/NsightComputeCli/index.html>
- Nsight Compute GUI user guide: <https://docs.nvidia.com/nsight-compute/NsightCompute/index.html>
- Analysis Driven Optimization (ADO) with Nsight Compute Dev Blog: <https://developer.nvidia.com/blog/analysis-driven-optimization-preparing-for-analysis-with-nvidia-nsight-compute-part-1/>
- “What The Profiler Is Telling You” GTC talk: <https://www.nvidia.com/en-us/on-demand/session/gtcsj20-s22141/>
- Roofline analysis- ALCF 2021: <https://www.youtube.com/watch?v=fsC3QeZHM1U>