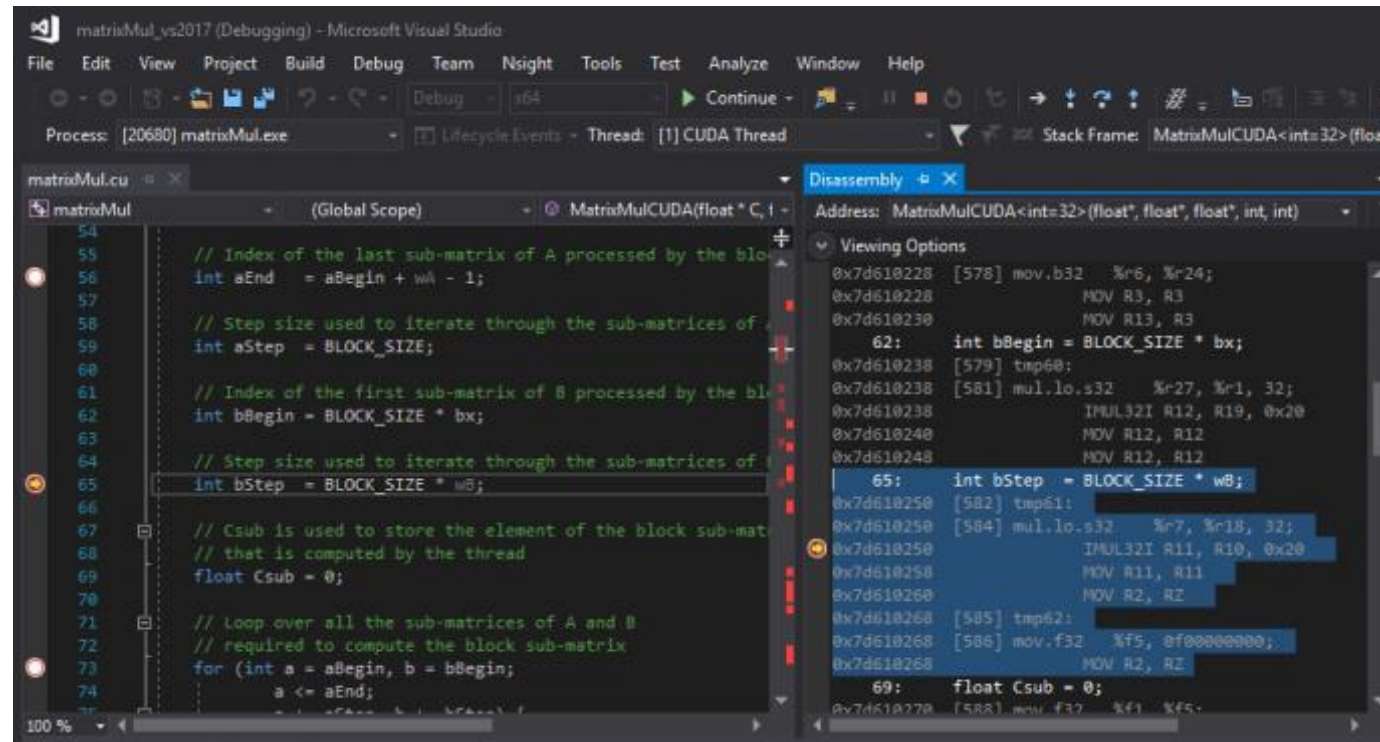**NVIDIA**

# NSIGHT DEVELOPER TOOLS

# DEVELOPER TOOLS

## Debuggers: cuda-gdb, Nsight Eclipse Edition, Nsight Visual Studio Edition, Nsight Visual Studio **Code** Edition
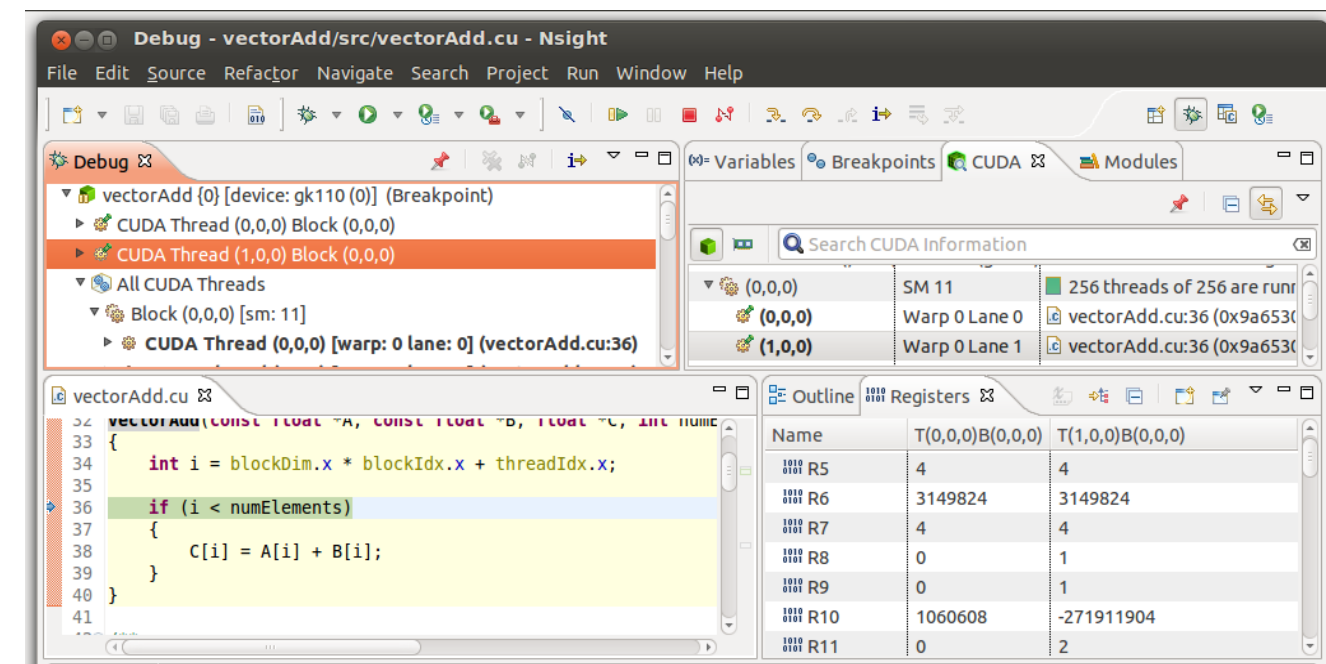


## Profilers: Nsight Systems, Nsight Compute, CUPTI, NVIDIA Tools eXtension (NVTX)



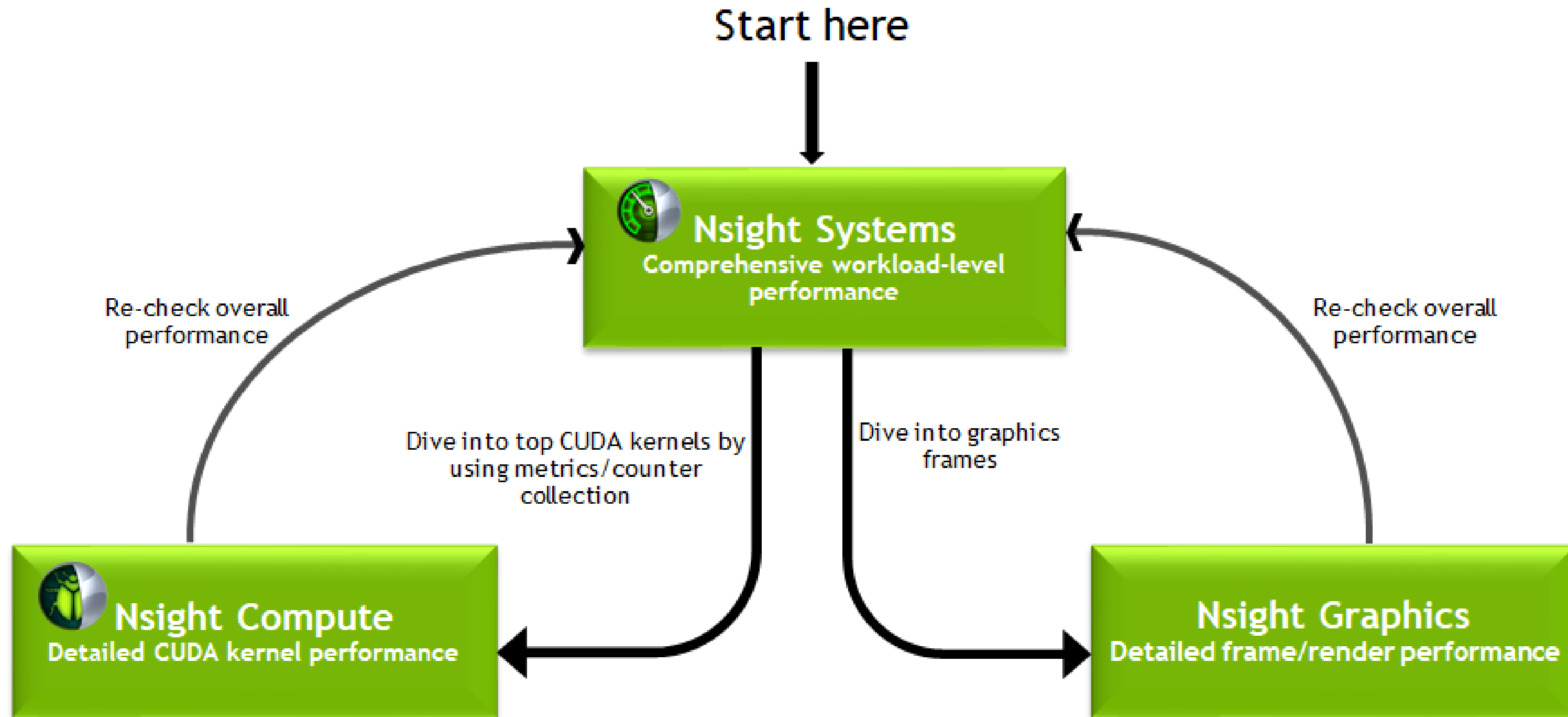## Correctness Checker: Compute Sanitizer

```
$ compute-sanitizer --leak-check full memcheck_demo
========= COMPUTE-SANITIZER
Mallocing memory
Running unaligned_kernel
Ran unaligned_kernel: no error
Sync: no error
Running out_of_bounds_kernel
Ran out_of_bounds_kernel: no error
Sync: no error
========= Invalid __global__ write of size 4 bytes
=========     at 0x60 in memcheck_demo.cu:6:unaligned_kernel(void)
=========     by thread (0,0,0) in block (0,0,0)
=========     Address 0x400100001 is misaligned
```

## IDE integrations: Nsight Eclipse Edition
Nsight Visual Studio Edition
Nsight Visual Studio **Code** Edition

# NSIGHT PROFILING TOOLS WORKFLOW
## ITERATIVE OPTIMIZATION FOR COMPUTE AND GRAPHICS

Start here

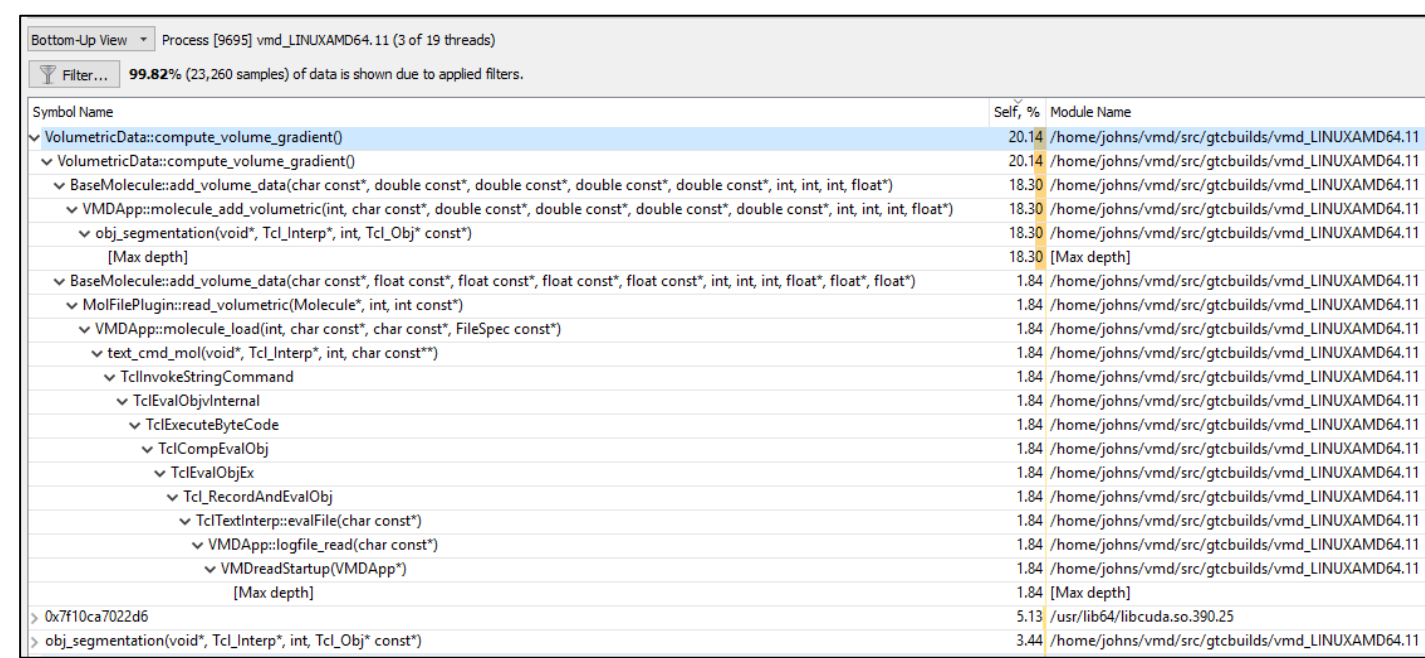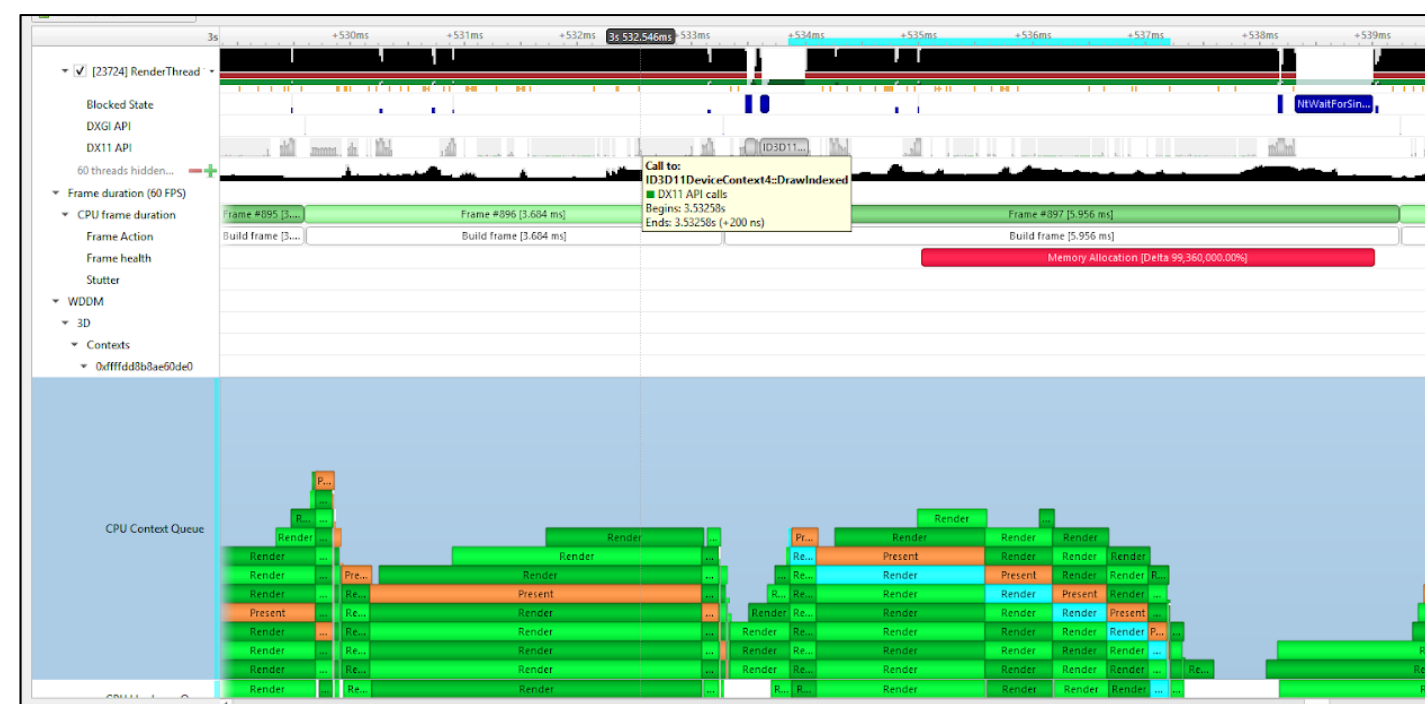**Nsight Systems**
Comprehensive workload-level performance

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
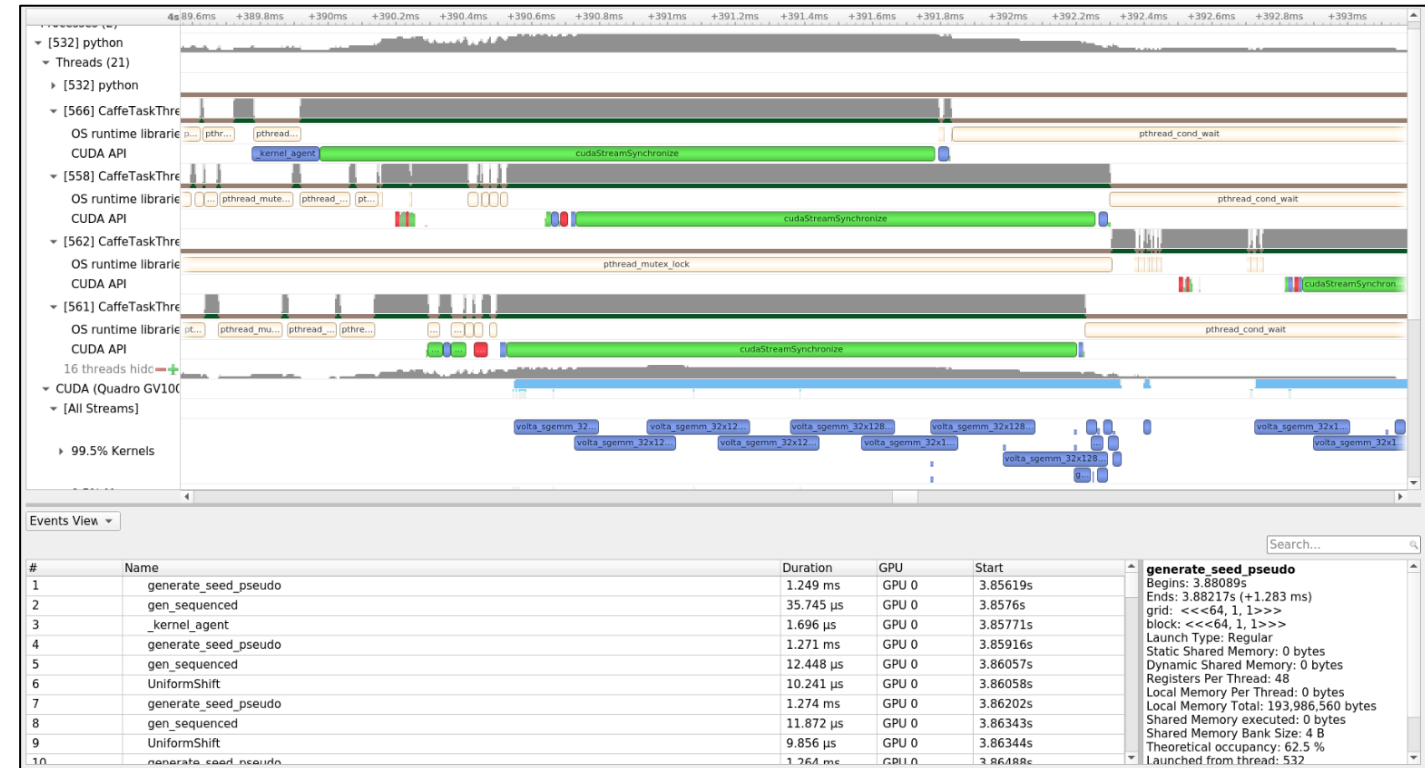
3

# NSIGHT SYSTEMS
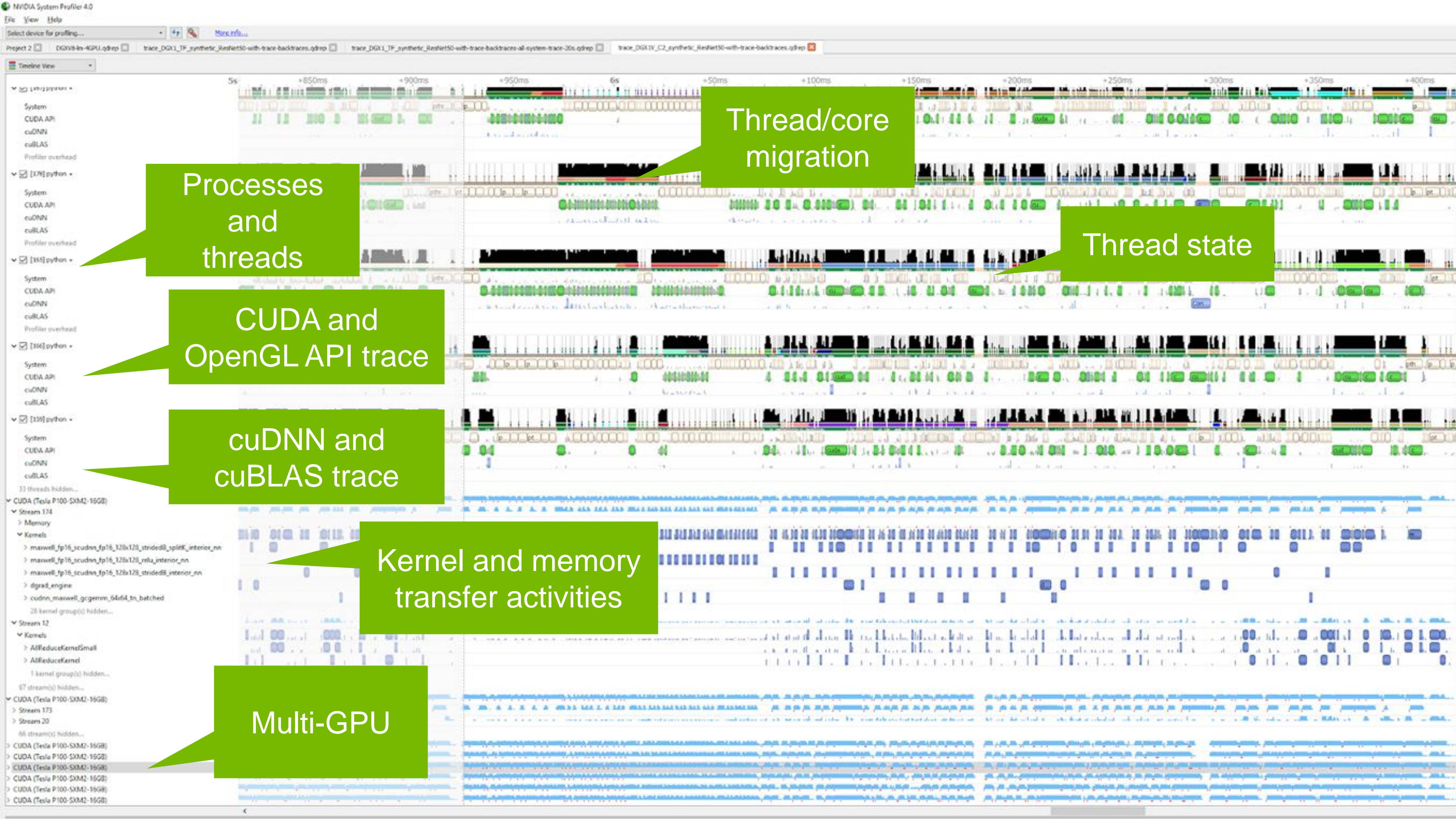## SYSTEM PROFILER

Key Features:

- 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

- Command Line, Standalone, IDE Integration

OS: Linux (x86, Power, Arm SBSA, Tegra), Windows, MacOSX (host)

GPUs: Pascal+

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

# ZOOM/FILTER TO EXACT AREAS OF INTEREST

# COLLECT A PROFILE WITH NSIGHT SYSTEMS

```
$ nsys profile -o report --stats=true ./myapp.exe
```

Generated file: `report.qdrep` (or `report.nsys-rep`); open for viewing in the Nsight Systems UI

When using MPI, I recommend using nsys after mpirun/srun/jsrun/etc.:

```
$ mpirun -n 4 nsys profile ./myapp.exe
```
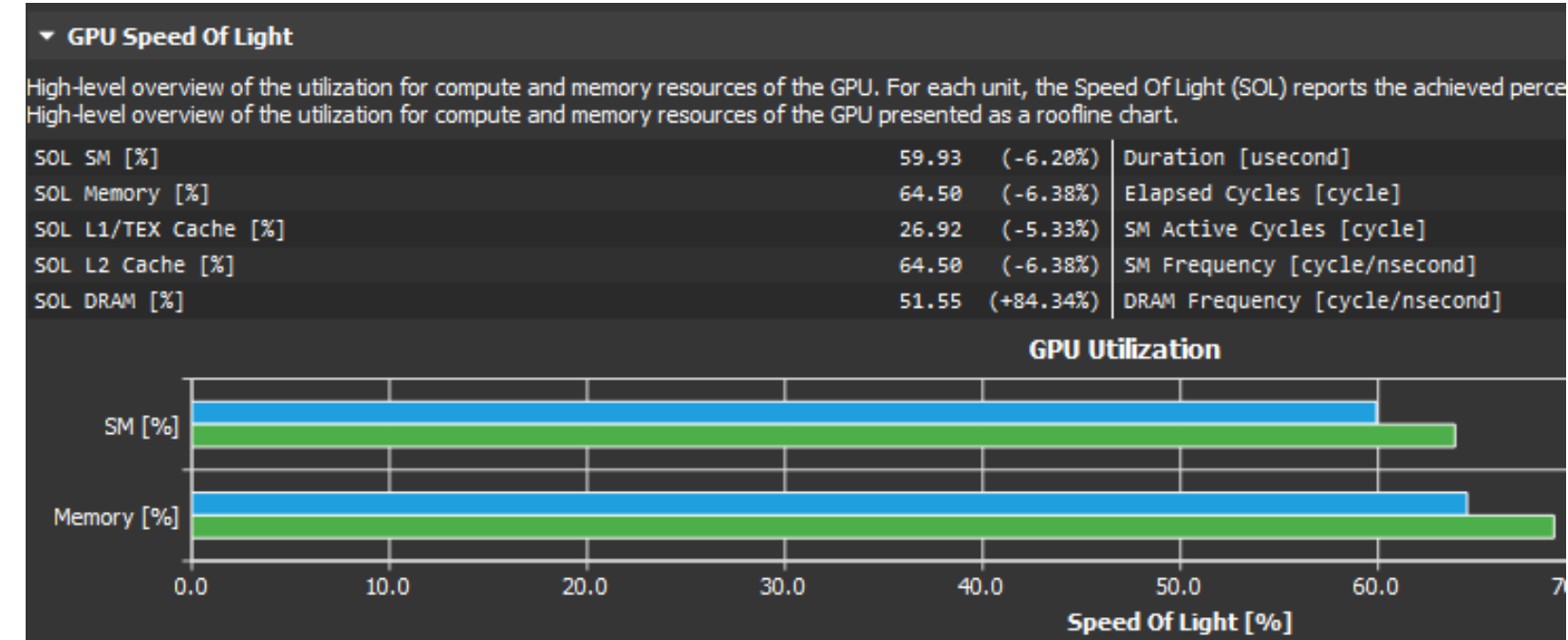
# NSIGHT COMPUTE
## KERNEL PROFILING TOOL

Key Features:

- Interactive CUDA API debugging and kernel profiling

- Built-in rules expertise

- Fully customizable data collection and display

- Command Line, Standalone, IDE Integration, Remote Targets

OS: Linux (x86, Power, Tegra, Arm SBSA), Windows, MacOSX (host only)

GPUs: Volta, Turing, Ampere GPUs

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

Targeted metric sections

Customizable data collection and presentation

Built-in expertise for Guided Analysis and optimization

Visual memory analysis chart

Metrics for peak performance ratios

All

Detailed analysis of the memory resources of the GPU. Memory can become a limiting factor for the overall kernel performance when fully utilizing the involved hardware units (Mem Busy), exhausting the available communication bandwidth between those units (Max Bandwidth), or by reaching the maximum throughput of issuing memory instructions (Mem Pipes Busy). Detailed chart of the memory units. Detailed tables with data for each memory unit.

| | | | |
|---|---|---|---|
| Memory Throughput [Gbyte/second] | 310.08 | Mem Busy [%] | 42.60 |
| L1 Hit Rate [%] | 46.75 | Max Bandwidth [%] | 44.73 |
| L2 Hit Rate [%] | 94.03 | Mem Pipes Busy [%] | 42.23 |

**Memory Chart**

Kernel

Global — 24.58 K Inst — 18.43 K Req / 6.14 K Req
Local — 43.01 K Inst — 12.29 K Req / 36.86 K Req
Texture — 0.00 Inst — 0.00 Req
Surface — 0.00 Inst — 0.00 Req / 0.00 Req
Shared — 110.59 K Inst — 65.29 K Req / 49.15 K Req

Unified Cache 46.75 %
384.38 KB / 3.56 MB
L2 Cache 94.03 %

System Memory — 0.00 B / 0.00 B
Device Memory — 478.34 KB / 4.20 MB

Shared Memory

**Shared Memory**

| | Instructions | Requests | % Peak | Bank Conflicts |
|---|---|---|---|---|
| Shared Load | 61,440 | 65,289 | 6.59 | 3,698 |
| Shared Store | 49,152 | 49,152 | 4.96 | 0 |
| Shared Atomic | 0 | - | - | - |
| Total | 110,592 | 114,441 | 11.55 | 3,698 |

**First-Level (Unified) Cache**

| | Instructions | SM->TEX Requests | % Peak | Hit Rate | TEX->L2 Requests | % Peak | L2->TEX Returns | % Peak |
|---|---|---|---|---|---|---|---|---|
| Global Load Cached | 18,432 | 18,432 | 1.86 | 66.65 | - | - | | |
| Global Load Uncached | - | - | - | - | - | - | 12,300 | 1.24 |
| Local Load Cached | 12,288 | 12,288 | 1.24 | 100 | | | | |

NVIDIA.

# KERNEL PROFILES WITH NSIGHT COMPUTE

```
$ ncu -k mykernel -o report ./myapp.exe
```

Generated file: `report.ncu-rep;` open for viewing in the Nsight Compute UI

(Without the -k option, Nsight Compute will profile everything and take a long time!)

# PROFILING RESOURCES

Nsight Systems, Nsight Compute product pages

NVIDIA Developer Blog: NVTX

NVIDIA Developer Blog: Transitioning from nvprof to nsys

NVIDIA Developer Blog: Using Nsight Compute to Inspect Your Kernels

OLCF Training, March 2020: Nsight Compute, Nsight Systems

# CUDA GDB

## COMMAND LINE AND IDE BACKEND DEBUGGER

- Unified CPU and CUDA Debugging

- CUDA-C/PTX/SASS support

- Built on GDB and uses many of the same CLI commands

```
(cuda-gdb)  info cuda threads breakpoint all
   BlockIdx ThreadIdx          Virtual PC Dev SM Wp Ln       Filename   Line
Kernel 0
    (1,0,0)    (0,0,0) 0x0000000000948e58   0 11  0  0 infoCommands.cu     12
    (1,0,0)    (1,0,0) 0x0000000000948e58   0 11  0  1 infoCommands.cu     12
    (1,0,0)    (2,0,0) 0x0000000000948e58   0 11  0  2 infoCommands.cu     12
    (1,0,0)    (3,0,0) 0x0000000000948e58   0 11  0  3 infoCommands.cu     12
    (1,0,0)    (4,0,0) 0x0000000000948e58   0 11  0  4 infoCommands.cu     12
    (1,0,0)    (5,0,0) 0x0000000000948e58   0 11  0  5 infoCommands.cu     12


(cuda-gdb)  info cuda threads breakpoint 2 lane 1
   BlockIdx ThreadIdx          Virtual PC Dev SM Wp Ln       Filename   Line
Kernel 0
    (1,0,0)    (1,0,0) 0x0000000000948e58   0 11  0  1 infoCommands.cu     12
```
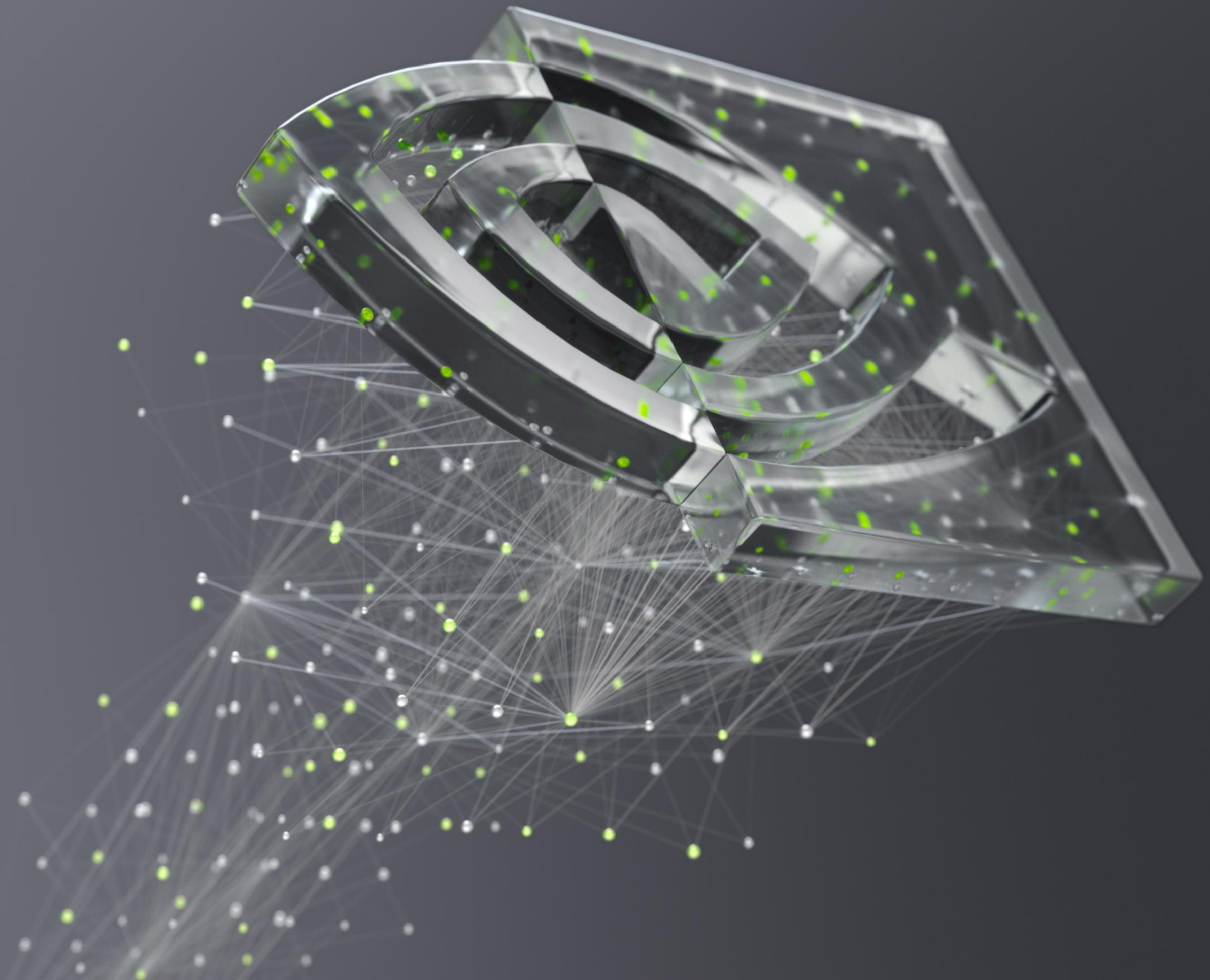
# COMPUTE SANITIZER
## AUTOMATICALLY SCAN FOR BUGS AND MEMORY ISSUES

- Compute Sanitizer checks correctness issues via sub-tools:

- *Memcheck* – The memory access error and leak detection tool.

- *Racecheck* – The shared memory data access hazard detection tool.

- *Initcheck* – The uninitialized device global memory access detection tool.

- *Synccheck* – The thread synchronization hazard detection tool.

```
~/W/m/c/build $ cmake ../ && cmake --build .
-- Configuring done
-- Generating done
-- Build files have been written to: /home/rmaynard/Work/misc/cuda_sanitizer_ctest/build
[2/2] Linking CUDA executable demo
~/W/m/c/build $ ctest -D MemoryCheck
    Site: RMAYNARD-DT
    Build name: Linux-unknown
Create new tag: 20210325-1346 - Experimental
Configure project
    Each . represents 1024 bytes of output
    . Size of output: 0K
Build project
    Each symbol represents 1024 bytes of output.
    '!' represents an error and '*' a warning.
    . Size of output: 0K
    0 Compiler errors
    0 Compiler warnings
Performing coverage
 Cannot find any coverage files. Ignoring Coverage request.
Memory check project /home/rmaynard/Work/misc/cuda_sanitizer_ctest/build
    Start 1: verify
1/1 MemCheck #1: verify ........................    Passed    6.77 sec

100% tests passed, 0 tests failed out of 1

Total Test time (real) =   6.77 sec
-- Processing memory checking output:
1/1 MemCheck: #1: verify ........................    Defects: 4
MemCheck log files can be found here: (<#> corresponds to test number)
/home/rmaynard/Work/misc/cuda_sanitizer_ctest/build/Testing/Temporary/MemoryChecker.<#>.log
Memory checking results:
Invalid __global__ read - 1
cudaErrorLaunchFailure - 3
Submit files
    SubmitURL: http://my.cdash.org/submit.php?project=CMakeTutorial
    Uploaded: /home/rmaynard/Work/misc/cuda_sanitizer_ctest/build/Testing/20210325-1346/Config
    Uploaded: /home/rmaynard/Work/misc/cuda_sanitizer_ctest/build/Testing/20210325-1346/Build.
    Uploaded: /home/rmaynard/Work/misc/cuda_sanitizer_ctest/build/Testing/20210325-1346/Dynami
    Uploaded: /home/rmaynard/Work/misc/cuda_sanitizer_ctest/build/Testing/20210325-1346/Done.x
    Submission successful
~/W/m/c/build $ 
```