

Nsight Systems - Introduction Robert Dietrich - March 17, 2022

Typical Optimization Workflow



Iterate until desired performance is achieved

NVIDIA Performance Analysis Tools



System-Wide Application Tuning

Maximize your GPU Investment

Locate optimization opportunities

- Visualize millions of events on a timeline
- See gaps of unused CPU and GPU time
- Balance your workload across multiple CPUs and GPUs
- CPU utilization and thread state
- GPU streams, kernels, memory transfers, etc.
- NIC performance metrics & tracing of network libraries

Multi-platform support

- Linux, Windows and Mac OS X (host-only)
- x86-64, Power9, ARM server, Tegra (Linux & QNX)





File <u>View</u> <u>Tools</u> <u>H</u>elp



Command Line Interface (CLI) Statistics and Export to SQLite, JSON, etc.

Command Line Interface

The Nsight Systems CLI provides several different commands

- Basic profiling session nsys profile ./app
- Interactive sessions (scriptable) nsys start|launch|stop|cancel nsys session list nsys status|shutdown
- Statistics and export nsys stats | export (export to sqlite, hdf, text, json, info)



CLI Profiling - Some Useful Switches

API tracing

```
-t, --trace=cuda,nvtx,osrt,opengl
```

(cublas, cusparse, cudnn, mpi, oshmem, ucx, openacc, openmp, vulkan, 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#
```

```
(patterns for hostname, PID and environment variables)
```

CLI Profiling - Some Useful Switches

Callstack sampling

- -s, --sample=[**cpu**|none]
- --sampling-period=number of CPU Instructions Retired events
- -b, --backtrace=[**lbr**|fp|dwarf|none]
- --samples-per-backtrace={1..12}
 - (The number of CPU IP samples collected for every CPU IP sample backtrace collected.)

Set the paranoid level: "sudo sh -c 'echo 1 >/proc/sys/kernel/perf event paranoid'

CUDA memory usage

```
--cuda-memory-usage=[true |false]
```

```
(Use nsys profile --help for a list of available options.)
```

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)

CLI Profiling - Additional Output

WARNING: The command line includes a target application therefore the CPU context-switch scope has been set to process-tree. Collecting data...

... APP OUTPUT ...

Temporary data is written to **/tmp/nvidia/nsight_systems** by default. Set **TMPDIR** to specify another location.

Processing events...

Saving temporary "/tmp/nsys-report-2b96-9038-d0bd-2600.qdstrm" file to disk...

CLI Stats Output

- CUDA API, kernels and memory operations (by time and by size)
- OS runtime

CUDA

Time

00 1 0110110		NVTX P	ush-Pop Range	e Statisti	cs:					
NVTX		Time(%) Total Ti	ne (ns) I	nstances	Average	Minimum	Maximum	Range	
		80	,0 16.193.3 ,0 3.680.9	376.534 937.672	21.268 8.072	761.396,0 456.013,0	6.823 1.726	2.438.848 1.932.681	MPI:MPI_Waitall MPI:MPI_Allreduce	
		0	,0 88.5	587.708	21 268	22.146.927,0	21.714.183	22.579.671	MPI:MPI_Init	
		0	,0 24.2	282.234	21.268	1.141,0	872	16.931	MPI:MPI_Irecv	
		0	,0 8. .0 3.	793.048 314.158	4 28	2.198.262,0	2.181.292	2.215.232	MPI:MPI_Finalize MPI:MPI Barrier	
Kernel Statis <mark>ti</mark> cs:		Ũ	,0	266.630	32	8.332,0	832	104.980	MPI:MPI_Reduce	
(%)	Total Time (ns)	Instances	Average	Minimum	Maxim	um			Name	
7,0	7.344.350.347	8.568	857.183,0	851.54	6 876.	922 device_te	a_leaf_ppcg_	solve_calc_s	d_new(kernel_info_t	, d
6,0	7.229.217.310	8.568	843.746,0	839.32	3 899.	482 device_te	a_leaf_ppcg_	solve_update	_r(kernel_info_t, c	loub
0,0	2.050.961.573	2.010	1.020.378,0	1.005.84	9 1.036.	569 device_te	a_leaf_cg_so	lve_calc_ur(kernel_info_t, doub	le,
9,0	1.879.919.365	2.010	935.283,0	913.75	4 1.214.	808 device_te	a_leaf_cg_so	lve_calc_w(k	ernel_info_t, doubl	.e*,
5,0	1.008.596.179	1.980	509.392,0	501.46	9 521.	756 device_te	a_leaf_cg_so	lve_calc_p(k	ernel_info_t, doubl	.e, (
0,0	61,606.046	30	2.053.534,0	1.176.60	0 2.890.	541 device_te	a_leaf_calc_	rrn(kernel_i	nfo_t, double const	*.
0,0	0,0 38.094.843 12.336 3.088,0 1.8		1.88	7 13.	184 void redu	ction <double< td=""><td>, (REDUCTION</td><td>_TYPE)0>(int, doubl</td><td>e*)</td></double<>	, (REDUCTION	_TYPE)0>(int, doubl	e*)	

NVTX NVIDIA Tool Extension Library

NVTX NVIDIA Tool Extension Library

NVTX provides means to correlate the profile data with the application code.

```
#include <nvtx3/nvToolsExt.h>
...
nvtxMark("Point in time");
...
nvtxRangePush("Name of your code region");
// your code goes here
nvtxRangePop();
...
nvtxRangeId_t rid = nvtxRangeStart("Name of another code region");
// stop might be on another thread
nvtxRangeEnd(rid)
```

(Add nvtx to the tracing options of nsys profile.)

NVTX Domains and Resource Naming

```
nvtxDomainHandle_t yourDomain = nvtxDomainCreateA("My Domain");
nvtxEventAttributes_t evtAttr = {0};
evtAttr.version = NVTX_VERSION;
evtAttr.size = NVTX_EVENT_ATTRIB_STRUCT_SIZE;
evtAttr.messageType = NVTX_MESSAGE_TYPE_ASCII;
eventAttrib.message.ascii = "My Range";
nvtxDomainRangePushEx(yourDomain, &evtAttr);
// your code
nvtxDomainRangePop(yourDomain);
```

```
// #include <sys/syscall.h>
nvtxNameOsThreadA(syscall(SYS gettid), "My Thread");
```

NVTX Registered Strings

```
nvtxDomainRangePop(yourDomain);
```



The NVIDIA HPC SDK Fortran compiler provides NVTX bindings.

- libnvhpcwrapnvtx.[a]so] has to be linked
- Documentation can be found here: <u>https://docs.nvidia.com/hpc-sdk/compilers/fortran-cuda-interfaces/index.html#cfnvtx-runtime</u>

```
use nvtx
...
call
nvtxStartRange("YourRange")
! some Fortran code
call nvtxEndRange
```

```
For other compilers, write your own Fortran NVTX bindings or use existing ones, e.g. https://raw.githubusercontent.com/maxcuda/NVTX_example/master/nvtx.f90
```



NVTX Python Bindings

"Python developers can either use decorators <code>@nvtx.annotate()</code> or a context manager with <code>nvtx.annotate(..)</code> "

Get NVTX Python module:

python -m pip install nvtx

Also see this blog or the docs: https://nvtx.readthedocs.io/en/latest/index.html

PyTorch CUDA provides NVTX bindings

```
from torch.cuda import nvtx
nvtx.range_push("YourCode")
# your Python code
nvtx.range_pop()
```

Demo: Teafleaf

TeaLeaf Demo

			NVIDIA Nsight Systems	2021.3.1			
<u>V</u> iew <u>T</u> ools <u>H</u> elp							
r50.qdrep 🗙 report.qdrep 🗙 p	profile_circe-n011_506451	_0.qdrep.qdrep × laplace2	d_baseline.qdrep 🗙 laplace2	d_data.qdrep 🗙 tealeaf.qdre	p ×		
Timeline View 👻					📾 Q 1x 🗍	A warnings, 2	25 messa
 mreaos (11) 	+904ms +906ms	+908ms +91	0ms +912ms	+914ms +916ms	+918ms +920ms	+922ms +924m	5
👻 🗹 [52199] MPI Rank 1 👻							
MPI				MPI_Waitall [1111111111111111111		
NVTX	update_halo [3,701 ms]	update_halo [3,672 ms]	update_halo [3,557 ms]	update_halo [3,689 ms]	update_halo [3,695 ms]	update_halo [3,688 ms]	upd
CUDA API	cudaMemcpy	cudaMemcpy	cudaMemcpy	cudaMemcpy	cudaMemcpy	cudaMemcpy	cud.
Profiler overhead 10 threads hidden — +							
▼ CUDA HW (0000:65:00.0 - NV			1		T	1	
▶ 99.7% Kernels	devi)devic)	devi devic	devi devic	devi devic	devi]devic	devi devic)
▶ 0.3% Memory							
NVTX	updat	pdate_halo [1,]	update_halo [1,]	update_hal	lo [3,707 ms]	update_halo [1,	updat.
 [52198] ./tea_leaf Threads (11) 							
▼ 🔽 [52198] MPI Rank 0 →							
MPI	MPI_Waitall [MPI_Waitall [MPI_Waitall [MPI_Waitall [MPI_Waitall	
NVTX	update_halo [3,704 ms]	update_halo [3,673 ms]	update_halo [3,682 ms]	update_halo [3,565 ms]	update_halo [3,691 ms]	update_halo [3,691 ms]	upda.
	4						•

Sample from https://github.com/UK-MAC/TeaLeaf

Build TeaLeaf

Get the source code

git clone https://github.com/UK-MAC/TeaLeaf_CUDA.git; cd TeaLeaf_CUDA

Annotate the code

```
use nvtx
...
call nvtxStartRange("Update Halo")
! some fortran code
call nvtxEndRange
```

Setup environment (MPI, CUDA, compiler, etc.)

Build TeaLeaf # add AMPERE architecture and -std=c++14 flag for nvcc make -e COMPILER=PGI

Profile TeaLeaf

```
# Run TeaLeaf without profiling
srun/mpirun [args] ./tea leaf
```

Profile with NVTX

```
srun/mpirun [args] \
nsys profile -t cuda,mpi,ucx,osrt,nvtx \
-y 1 -d 8 --kill=none \
-o tea-nvtx.%q{SLURM_PROCID} --stats=true \
./tea_leaf
```

```
# Profile with CPU and GPU sampling
srun/mpirun [args] \
nsys profile -t cuda,mpi,ucx,osrt,nvtx \
-y 1 -d 8 --kill=none -o tea-sampling%q{SLURM_PROCID} \
--backtrace=dwarf --sampling-period=3000000 \
--gpu-metrics-set=ga100 --gpu-metrics-device=0,1 \
--gpu-metrics-frequency=15000 \
./tea_leaf
```

Analyze the Profile

Investigate CLI stats output (--stats)

Open the report file in the Nsight Systems GUI: nsys-ui tea-nvtx.0.nsys-rep

≡ 🖃 🖸 🔌 t System: ile View Tools Help	s 2022.1.1
ea-full.0 [2 reports] ×	
■ Timeline View ・	📾 Q 1x 🗋 📥 4 warnings, 17 messages
	0s 0.5s 1s 1.5s 2s 2.5s 3s 3.5s 4s 4.5s 5s 5.5s 6s 6.5s 7s 7.5s 8s 8.5s 9s
 jwb0057.juwels (0:0) 	
▶ CPU (96)	
 Threads (5) 	
👻 🔽 [25027] MPI Rank 0 🔹	
OS runtime libraries	
MPI	
► UCX	
4 threads hidden	
L LICY	
, ocx	- 1 of a full initial as which in the initial as a second second of the second s
CUDA HW (0000:03:00.0 - 1	
 jwb0057.juwels (1:0) 	
 CPU (96): Report 1 	
 Threads (5) 	
▼ ▼ [25026] MPI Rank 1	
OS runtime libraries	
MPI	
+ UCX	
	CONTRACTOR AND A DESCRIPTION OF A DESCRIPT
4 threads hidden – –	
> UCX	
, dex	
CUDA HW (0000:44:00.0 - 1	
 GPU (0000:03:00.0 - NVIDI) 	
GPU Metrics [15 kHz]	
SYS Clock Frequency	
GR Active	
SM Active	
SM Instructions	

Eile View Tools Help

tea-full.0 [2 reports] 🗙						
≡ Timeline View 👻					🔤 Q 1x 🗌	🕂 🔥 4 warnings, 17 messages
2s •	+135ms +140ms +145ms +150ms +155	ims 2s 160.70ms +165	5ms +170ms +175ms +	180ms +185ms +190ms +1	95ms +200ms +205ms	+210ms +215ms
 jwb0057.juwels (0:0) 						
 CPU (96) 						
 Threads (5) 			a a a		1 N	
👻 🔽 [25027] MPI Rank 0 👻						
OS runtime libraries						
MPI	MPI_Waitall [23.327 ms]		MPI_Waitall [17.795 ms]		MPI_Waitall [17.800 ms]	MPI
→ UCX						
NVTX 000	update_halo [23.885 ms]		update_halo [18.347 ms]		update_halo [18.772 ms]	UUUUUupdate
CUDA API		MUMMA				
4 threads hidden +	urn tag cand phy 1100 transfer processing [32,330 ms]	- ·	s cand phy LUCB transfer processing 117 797 ms		too cood phy LUCD transfer processing [17]	from t
► UCX	ucp_tag_recv_nbx UCP transfer processing [23.327 ms]	ucp_tag	recv_nbx UCP transfer processing [17.797 ms		tag_recv_nbx UCP transfer processing [17.7	96 ms] ucp_t
+ CUDA HW (0000:03:00.0 - N				ARAINARAINARA DISTANG		
✓ jwb0057.juwels (1:0)						
+ CPU (96): Report 1		CDU utilizati	1			
 Threads (5) 	_	Average: 1.	.0%			
- ✓ [25026] MPI Rank 1 -		1111E. 2.1303	13			
OS runtime libraries		THURIDIN		Tomununfununun		THURLING
MPI						
+ UCX						
NVTX 000	update_halo [23.928 ms]		update_halo [18.340 ms]		update_halo [18.769 ms]	update
CUDA API	cudaMemcpy		cudaMemcpy		cudaMemcpy	cudaMe
4 threads hidden +		Tummer				THEFT
+ UCX				- 1 × · · · · · · · · · · · · · · · · · ·		· · · · · · · · · · · · · · ·
► CUDA HW (0000:44:00.0 - N		LARBINGTON.				
▼ GPU (0000:03:00.0 - NVIDIA						
 GPU Metrics [15 kHz] 		Addeniated.		And the second s		ALL
GPC Clock Frequency						
SYS Clock Frequency						
GR Active		United at		and the state of t		ALCO AND AND A
SM Active		Anternation		an a		all and a state of the
SM Instructions						100 100 100 100 100 100 100 100 100 100
	*	* ARX *** * * *		Enterne and an excense		<u></u>

$\equiv \bigsqcup :: \bigotimes t \text{ Systems 2022.1.1}$ <u>File View Tools Help</u>

tea-full.0 [2 reports] 🗙

🖳 Diagnostics Summary 👻 Report: 1 💌

Messages

	-			
	Source	Process ID	Time	Description
0	Injection	25026	-00:00.766	OS runtime libraries injection initialized successfully.
0	Injection	25026	-00:00.766	NVTX injection initialized successfully.
0	Injection	25026	-00:00.766	Buffers holding CUDA trace data will be flushed on CudaProfilerStop() call.
0	Injection	25026	-00:00.766	CUDA injection initialized successfully.
0	Injection	25026	-00:00.766	Common injection library initialized successfully.
Â	Daemon		-00:00.428	Unable to collect CPU kernel IP samples and backtraces. perf event paranoid level is 2. Change the paranoid level to 1 to enable kernel sample collection. Try sudo sh -c 'echo 1 >/proc/sys/kernel/perf_event_paranoid' to change the paranoid level to 1.
0	Daemon		-00:00.412	Dwarf backtraces collected.
0	Daemon		-00:00.412	Hardware event 'instructions', with sampling period 3000000, used to trigger sample collection.
0	Daemon		00:00.640	4 CPU IP samples collected for every CPU IP backtrace collected.
0	Analysis		00:01.110	Profiling has started.
0	Analysis		00:01.412	Number of GPU Metrics events collected: 199,298.
	Analysis	25026	00:01.874	Not all NVTX events might have been collected.
0	Analysis	25026	00:01.874	Number of NVTX events collected: 40,871.
	Analysis	25026	00:01.874	Not all CUDA events might have been collected.
0	Analysis	25026	00:01.874	Number of CUDA events collected: 55,019.
	Analysis	25026	00:01.874	Not all OS runtime libraries events might have been collected.
0	Analysis	25026	00:01.874	Number of OS runtime libraries events collected: 204.
0	Analysis		00:09.134	Profiling has stopped.
0	Injection	25026	00:09.167	Number of CUPTI events produced: 55,331, CUPTI buffers: 50.
0	Daemon		00:09.793	Number of IP samples collected: 7,894.
0	Daemon	25026	08:27:01.988	Process was launched by the profiler, see /tmp/nvidia/nsight_systems/guadd_session_1125009/streams/pid_25026_stdout.log and stderr.log for program output

≡ 🖃 🖸 💘 t Systems 2022.1.1

File View Tools Help



Launch Application via GUI

•		NVIDIA Nsight Systems 2021.3.1		~ ^ (
<u>File View Tools H</u> elp				
Project Explorer SimpleCUBLAS	Localhost connection	✓ Target is ready More info		
mpirestsuite	SSH connections (3)			
	Solution 000,00 Solution 000,00 <td< td=""><td>0 events tions Retired' events counted before a CPU instruction p ected. The smaller the sample period, the higher the san icantly increase the size of result file(s).</td><td>ointer (IP) sample is npling rate. Lower</td><td>Start profiling manually Start profiling after 10,0 \$ seconds Start profiling after 100 \$ frames Limit profiling to 10,0 \$ seconds</td></td<>	0 events tions Retired' events counted before a CPU instruction p ected. The smaller the sample period, the higher the san icantly increase the size of result file(s).	ointer (IP) sample is npling rate. Lower	Start profiling manually Start profiling after 10,0 \$ seconds Start profiling after 100 \$ frames Limit profiling to 10,0 \$ seconds
	Mode: Specify process launch options below Command line with arguments: mpirun -np 1 ./mpi_test_suite Working directory: /home/rdietrich/testing/mpi/mpi-test-suite		Edit arguments	Limit profiling to 600 🗘 frames Hotkey Start/Stop F12 • (not available in console apps)
	Environment variables Trace fork before exec	Eile <u>V</u> iew <u>T</u> ools <u>H</u> elp		٨
	 Collect CPU context switch trace Collect OS runtime libraries trace 	<u>N</u> ew Project Open <u>A</u> dd Report (beta)	Ctrl+N MpiTes Ctrl+O Ctrl+T Selec	ect target for profiling
	Collect CUDA trace Collect OpenMP trace	Import Export	Ctrl+I Last Ctrl+E Selec	used target: marvel (localhost). <u>Select</u> ct a target to see available options.
	Collect GPU context switch trace	Close MpiTestSuite	Ctrl+W Ctrl+Q	

Fosturo Soloction		► ✓ Collect CPU context switch trace				
l'eature selection	۲	Collect OS runtime libraries trace				
 ✓ Sample target process 	•	✓ Collect CUDA trace				
Sampling Period: 1,000,000 events The sampling period is the number of 'CPU Instructions Retired' events counted bef collected. If configured, call stacks may also be collected. The smaller the sample periods will increase overhead and significantly increase the size of result Collect call stacks of executing threads Backtracing algorithm Current settings: use DWARF debug information Symbol locations No directories with symbol files. When stripped libraries (e.g. *.so files) are used on the target, specify here director symbols resolved. For best backtraces, specify the following compiler flags: on x86_64: -g Note that stripped binaries typically do not contain the debug information. Conside	•	 Flush data periodically 10,00 \$ seconds Skip some API calls Collect GPU memory usage Collect UM CPU page faults Collect UM GPU page faults Collect cuDNN trace Collect cuBLAS trace Collect OpenACC trace Collect CUDA backtraces 				
Target application Mode: Specify process launch options below Command line with arguments: mpirun -n 2 /home/rdietrich/testing/mpi/hello_mpi/hello_mpi Working directory:	> > >	Collect GPU metrics Collect NVENC trace Collect NVTX trace Collect OpenGL trace				
/home/rdietrich/testing/mpi/hello_mpi	•	Collect Vulkan trace				
	•	Communication profiling options (MPI, SHMEM, UCX)				

Nsight Systems - Summary

Today

- Mainly used for single server profiling runs
- Focused on GPU workloads, API trace, thread samples
- Cluster profiling
 - Load multiple reports into one timeline
 - MPI (including communication parameters), OpenSHMEM, UCX tracing
 - NIC metrics sampling (congestions, read/write bandwidth)

Roadmap

- Improve multi-node support, e.g. correlation across nodes, switch statistics, etc.
- Determine latency, congestion, hot devices (NICs, switches, network storage)

THANK YOU!

Download	https://developer.nvidia.com/nsight-systems NOTE: website version is newer than CUDA Toolkit version
Docs	https://docs.nvidia.com/nsight-systems/index.html
Forums	https://devtalk.nvidia.com
Email	nsight-systems@nvidia.com
Blogs	 <u>https://developer.nvidia.com/blog/nvidia-nsight-systems-containers-cloud</u> <u>https://developer.nvidia.com/blog/nsight-systems-exposes-gpu-optimization</u> <u>https://developer.nvidia.com/blog/understanding-the-visualization-of-overhead-and-latency-in-nsight-systems</u> <u>https://developer.nvidia.com/blog/nvidia-tools-extension-api-nvtx-annotation-tool-for-profiling-code-in-python-and-c-c/</u>

