



DEEP
LEARNING
INSTITUTE

Introduction into NVIDIA® Nsight™ Systems

Dr. Momme Allalen | LRZ | 12.07.2022



PRACE Training



LRZ as part of the Gauss Centre for Supercomputing (GCS) and IT4Innovations belong to the 14 **PRACE Training Centres** that started in 2012-2017-2020:

- Barcelona Supercomputing Center (Spain)
- CINECA Consorzio Interuniversitario (Italy)
- CSC – IT Center for Science Ltd (Finland)
- EPCC at the University of Edinburgh (UK)
- Gauss Centre for Supercomputing (Germany)
- Maison de la Simulation (France)
- GRNET – Greek Research and Technology Network (Greece)
- ICHEC – Irish Centre for High-End Computing (Ireland)
- IT4I – National Supercomputing Center VSB Technical University of Ostrava (Czech Republic)
- SURFsara (The Netherlands)
- TU Wien – VSC Research Center (Austria)
- University ANTWERPEN – VSC & CÉCI (Belgium)
- University of Ljubljana – HPC Center Slovenia (Slovenia)
- Swedish National Infrastructure for Computing (SNIC) (Sweden)



Mission: Serve as **European hubs and key drivers of advanced high-quality training** for researchers working in the computational sciences.

<http://www.training.prace-ri.eu/>



DEEP LEARNING INSTITUTE

DLI Mission: Help the world to solve the most challenging problems using AI and deep learning

We help developers, data scientists and engineers to get started in architecting, optimizing, and deploying neural networks to solve real-world problems in diverse industries such as autonomous vehicles, healthcare, robotics, media & entertainment and game development.

CUDA® PROFILING TOOLS

nvvp: NVIDIA visual profiler

nvprof: tool to understand and optimize the performance of your CUDA, OpenACC or OpenMP applications,

Application level opportunities

Overall application performance

Overlap CPU and GPU work, identify the bottlenecks (CPU or GPU)

Overall GPU utilization and efficiency

-Overlap compute and memory copies

-Utilize compute and copy engines effectively.



Nsight Systems
Nsight Compute

Kernel level opportunities

- Use memory bandwidth efficiently
- Use compute resources efficiently
- Hide instruction and memory latency

There are more features, example for Dependency Analysis

Command: **nvprof --dependency-analysis --cpu-thread-tracing on ./executable_cuda**

NSIGHT PRODUCT FAMILY

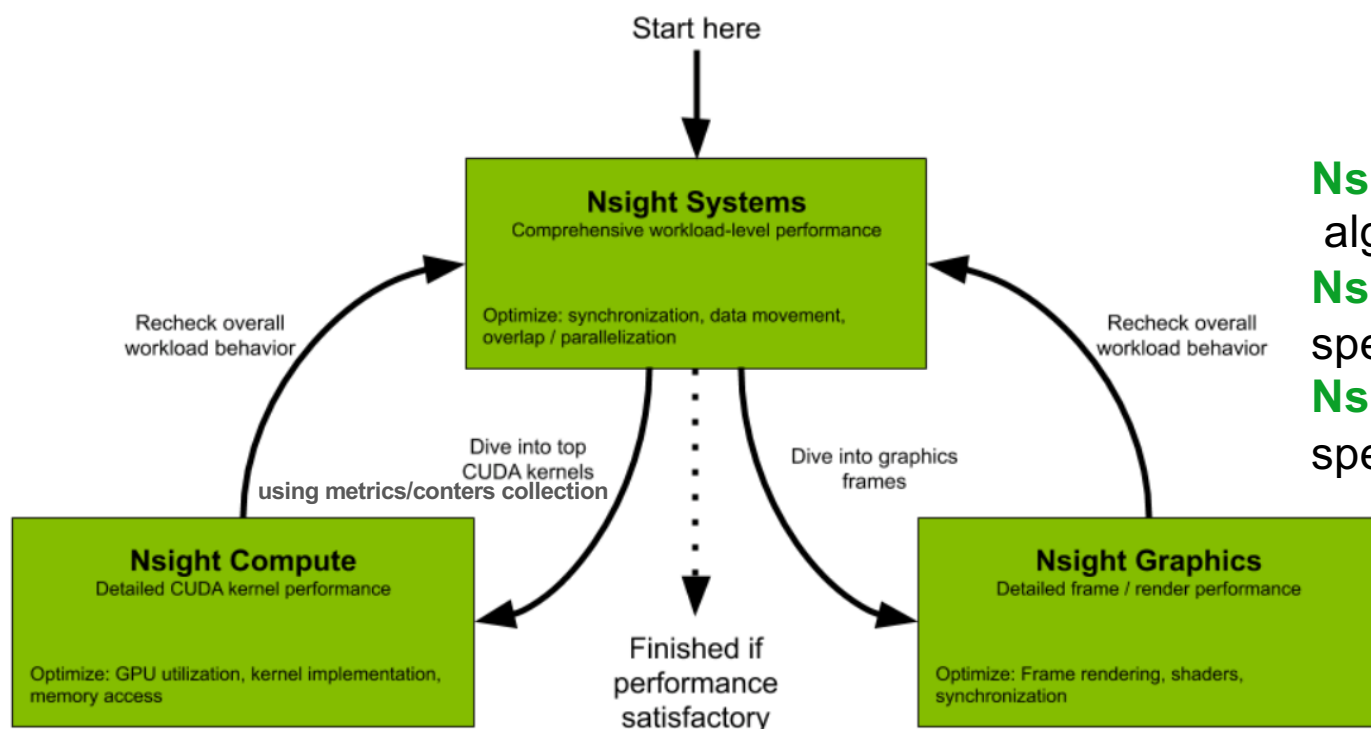


Figure 1. Flowchart describing working with new NVIDIA Nsight tools for performance optimization

Nsight Systems – Analyze application algorithm system wide
Nsight Compute – Debug/&optimize specific CUDA kernels
Nsight Graphics – Debug/&optimize specific graphics workloads

nvprof replaced with **nsys –profile....**

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

NVIDIA NSIGHT SYSTEMS



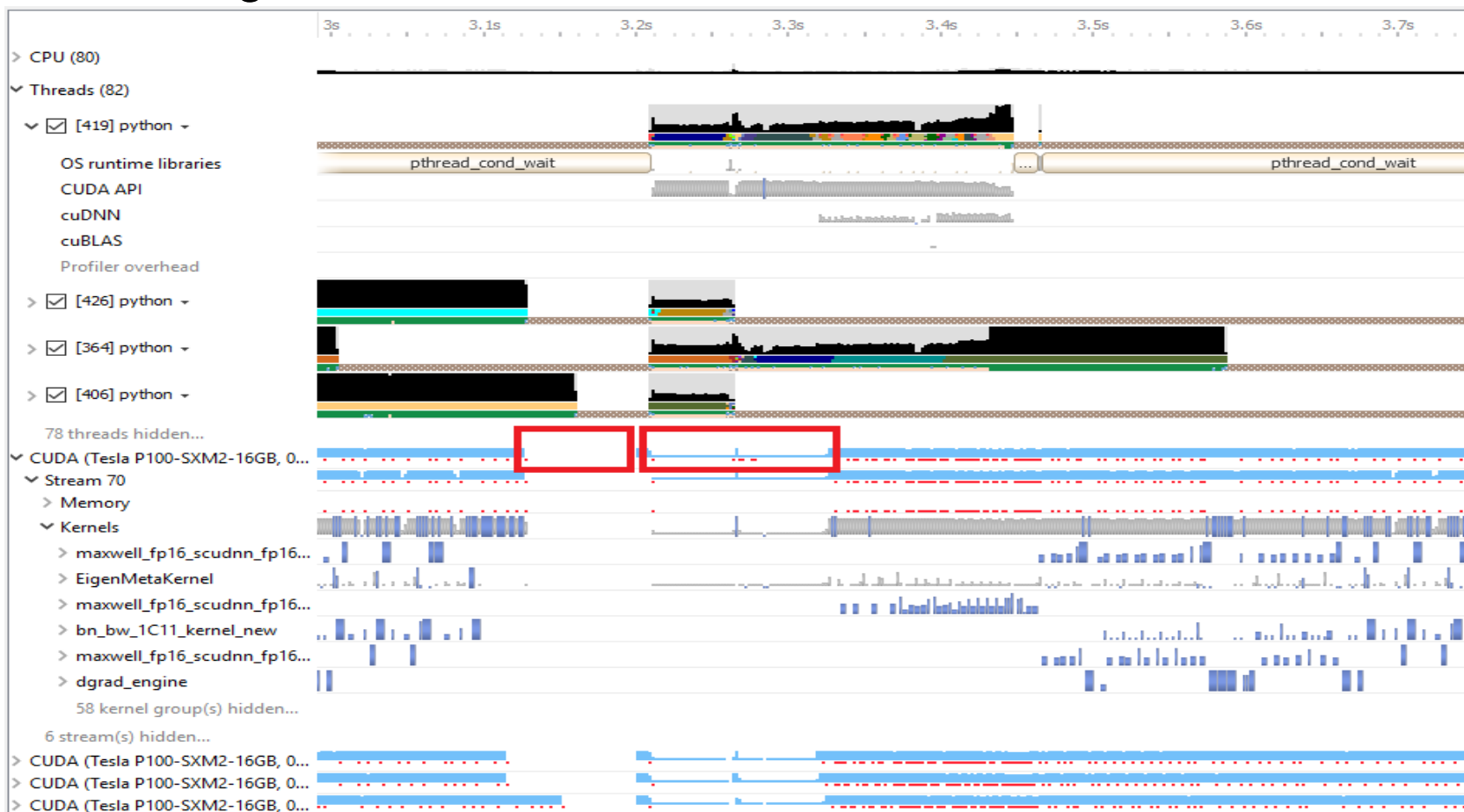
- Provides users with a more complete view of how their codes balance workload across multiple CPUs and GPUs
- Locate optimization opportunities, helps and allows to identify issues such as:
 - GPU starvation
 - Insufficient CPU parallelisation or pipelining
 - Unexpectedly expensive CPU or GPU algorithm
 - Unnecessary GPU synchronization
- The tool uses low overhead tracing and sampling techniques to collect process and thread activity and visualize millions of events on a very fast GUI timeline
- Correlates that data across CPU cores and GPU streams, allowing users to investigate bottlenecks.
- Multi-platform: Linux & Windows, x86-64, Tegra, Power, MacOSX (host only)

Command Line Options nsys

Command	Description
profile	A fully formed profiling description requiring and accepting no further input. The command switch options used (see below table) determine when the collection starts, stops, what collectors are used (e.g. API trace, IP sampling, etc.), what processes are monitored, etc.
start	Start a collection in interactive mode. The start command can be executed before or after a launch command.
stop	Stop a collection that was started in interactive mode. When executed, all active collections stop, the CLI process terminates but the application continues running.
cancel	Cancels an existing collection started in interactive mode. All data already collected in the current collection is discarded.
launch	In interactive mode, launches an application in an environment that supports the requested options. The launch command can be executed before or after a start command.
shutdown	Disconnects the CLI process from the launched application and forces the CLI process to exit. If a collection is pending or active, it is cancelled
export	Generates an export file from an existing .nsys-rep file. For more information about the exported formats see the /documentation/nsys-exporter directory in your Nsight Systems installation directory.
stats	Post process existing Nsight Systems result, either in .nsys-rep or SQLite format, to generate statistical information.
analyze	Post process existing Nsight Systems result, either in .nsys-rep or SQLite format, to generate expert systems report.
status	Reports on the status of a CLI-based collection or the suitability of the profiling environment.
sessions	Gives information about all sessions running on the system.

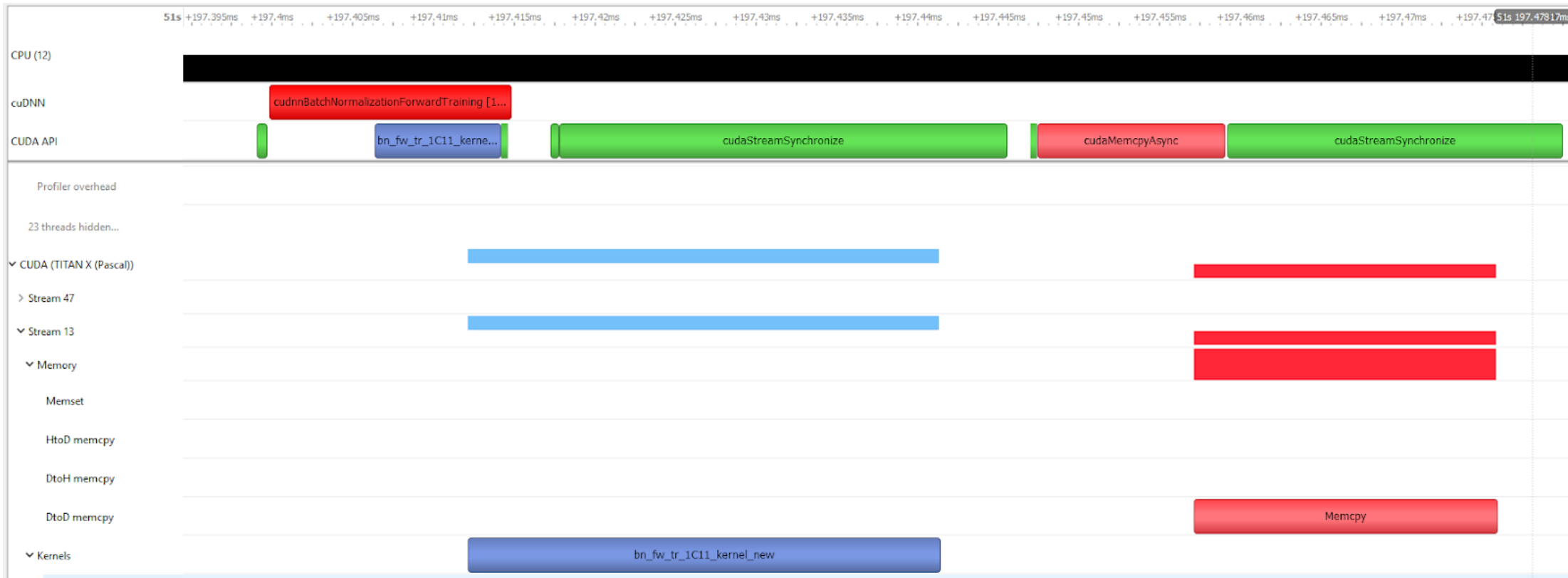
<https://docs.nvidia.com/nsight-systems/UserGuide/index.html>

GPU Starvation Investigations



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

Unnecessary GPU Synchronization Calls



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

NVIDIA NSIGHT SYSTEMS

- Support: **MPI**, **OpenACC**, **OpenMP**
- Complex data mining capabilities, enables to go beyond basic statistics.
- Support multiple simultaneous sessions.
- **MPI** trace feature enables to analyse when the threads are busy or blocked in long-running functions of the **MPI** standard, available on **OpenMPI**, **MPICH** and **NVShmem**.
- **OpenACC** trace enables to see where code has been offload and parallelized onto the GPU, which helps you to analyse the activities executing on the CPUs and GPUs in parallel.
- Tracing **OpenMP** code is available for compilers supporting **OpenMP 5** and **OMPT interface**. This capability enables tracing of the parallel regions of code that are distributed either across multiple threads or to the GPU.
- Provides support for **CUDA** graphs. To understand the execution of the source of **CUDA** kernels and execution of **CUDA** graphs, kernels can be correlated back through the graph launch, instantiation, and all the way back to the node creation, to identify the origin of the kernel execution on the GPU.

NVIDIA NSIGHT COMPUTE (ncu)



Interactive CUDA Kernel profiler

Targeted metric sections for various performance aspects (Debug/Profile)
API debugging via a user interface command line tool.

Very high freq GPU perf counter, customizable data collection and presentation
(tables, charts ...)

Python-based rules for guided analysis (or postprocessing)

Provides a customizable and data-driven user interface and metric collection and
can be extended with analysis scripts for post-processing results.

<https://developer.nvidia.com/nsight-compute>

NVIDIA NSIGHT COMPUTE

Important Features



- Result comparison across one or multiple reports within the tool
- Graphical profile report
- Interactive kernel profiler and API debugger: debugging CPU and GPU simultaneously and capable of handling thousands of simultaneous threads.
- Fast data collection
- GUI and command line interface
- Fully customizable reports and analysis rules.

<https://developer.nvidia.com/nsight-compute>

Nsight Compute Feature Spotlight in CUDA Toolkit 11 and A100

- **Roofline Analysis**

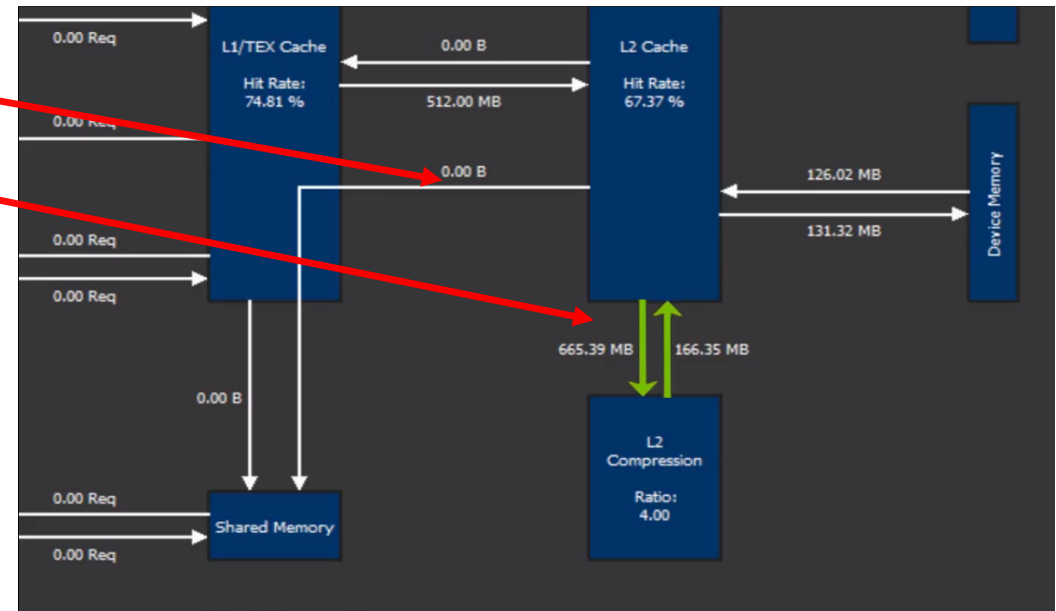
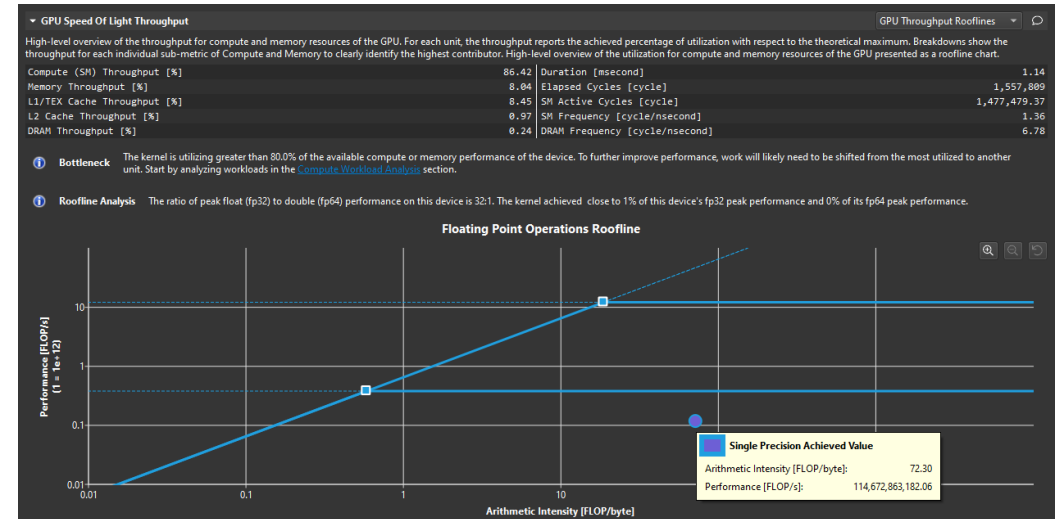
Arithmetic intensity = Compute/Memory

FLOPS = Floating Points Ops/Second

- **Asynchronous copy**

- **Sparse Data Compression**

Shows the amount of data compressed through this feature and the compression ratio, helps on kernels with bandwidth or cache issues.



NVIDIA® Tools Extension SDK (NVTX)

- C-based Application Programming Interface (API) for annotating events, code ranges, and resources in your applications
- Codes which integrate NVTX can use NVIDIA Nsight, Tegra System Profiler, and Visual Profiler to capture and visualize these events and ranges.

```

[allalen1@jwlogin22 v2]$ ncu -h | grep nvtx
--nvtx                Enable NVTX support.
--nvtx-include arg    Adds include statement to the NVTX filter, which allows selecting kernels to
--nvtx-exclude arg    Adds exclude statement to the NVTX filter, which allows selecting kernels to
--print-nvtx-rename arg (=none) Select how NVTX should be used for renaming:
                        per-nvtx
Usage of --nvtx-include and --nvtx-exclude:
ncu --nvtx --nvtx-include "Domain A@Range A"
ncu --nvtx --nvtx-exclude "Range A]"
ncu --nvtx --nvtx-include "Range A" --nvtx-exclude "Range B"
  
```

<https://docs.nvidia.com/nsight-visual-studio-edition/nvtx/index.html>

NVIDIA® Tools Extension SDK (NVTX)

```
#include <nvToolsExt.h>
#include <sys/syscall.h>
#include <unistd.h>
```

```
static void wait(int seconds) {
    nvtxRangePush(__FUNCTION__);
    nvtxMark("Waiting...");
    sleep(seconds);
    nvtxRangePop();
}
```

```
int main(void) {
    nvtxNameOsThread(syscall(SYS_gettid), "Main Thread");
    nvtxRangePush(__FUNCTION__);
    wait(1);
    nvtxRangePop();
}
```

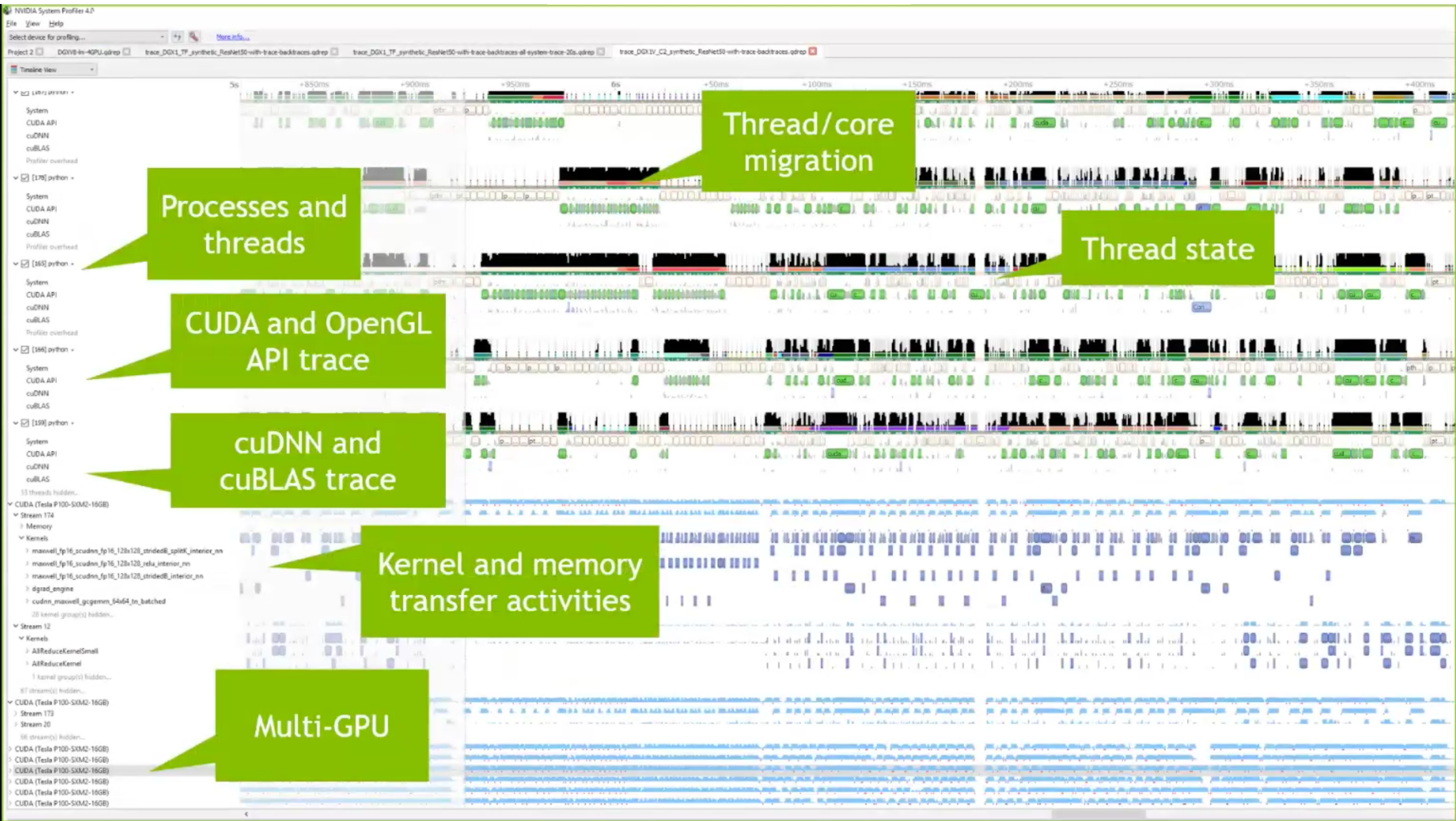
nsys profile -t nvtx --stats=true ...

Or for Julia code:

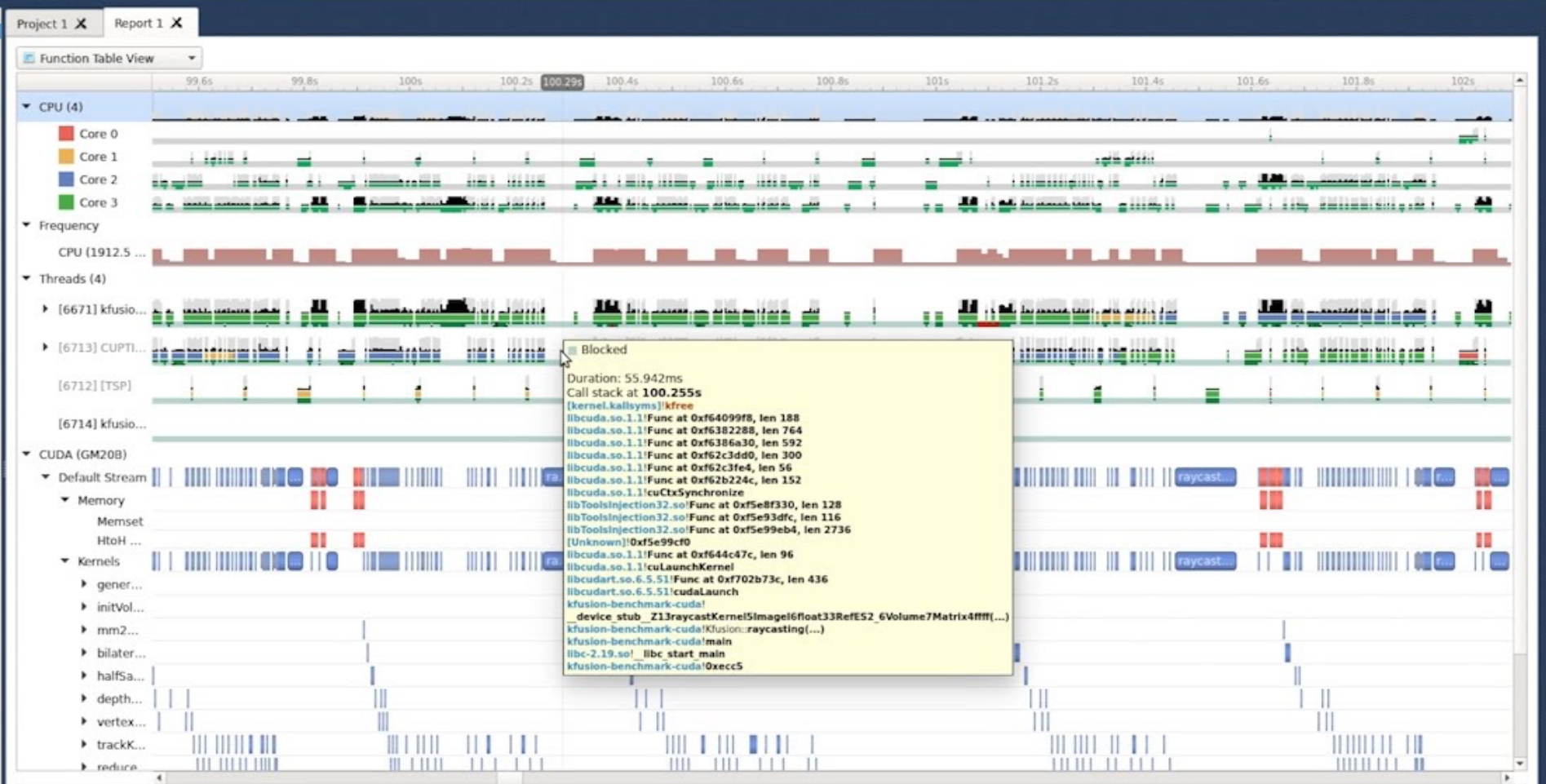
**nsys profile -t nvtx,cuda -o output_file.qdrep
julia --project=../.. script.jl**

<https://docs.nvidia.com/nsight-visual-studio-edition/2020.1/nvtx/index.html>

<https://docs.nvidia.com/nsight-visual-studio-edition/nvtx/index.html>

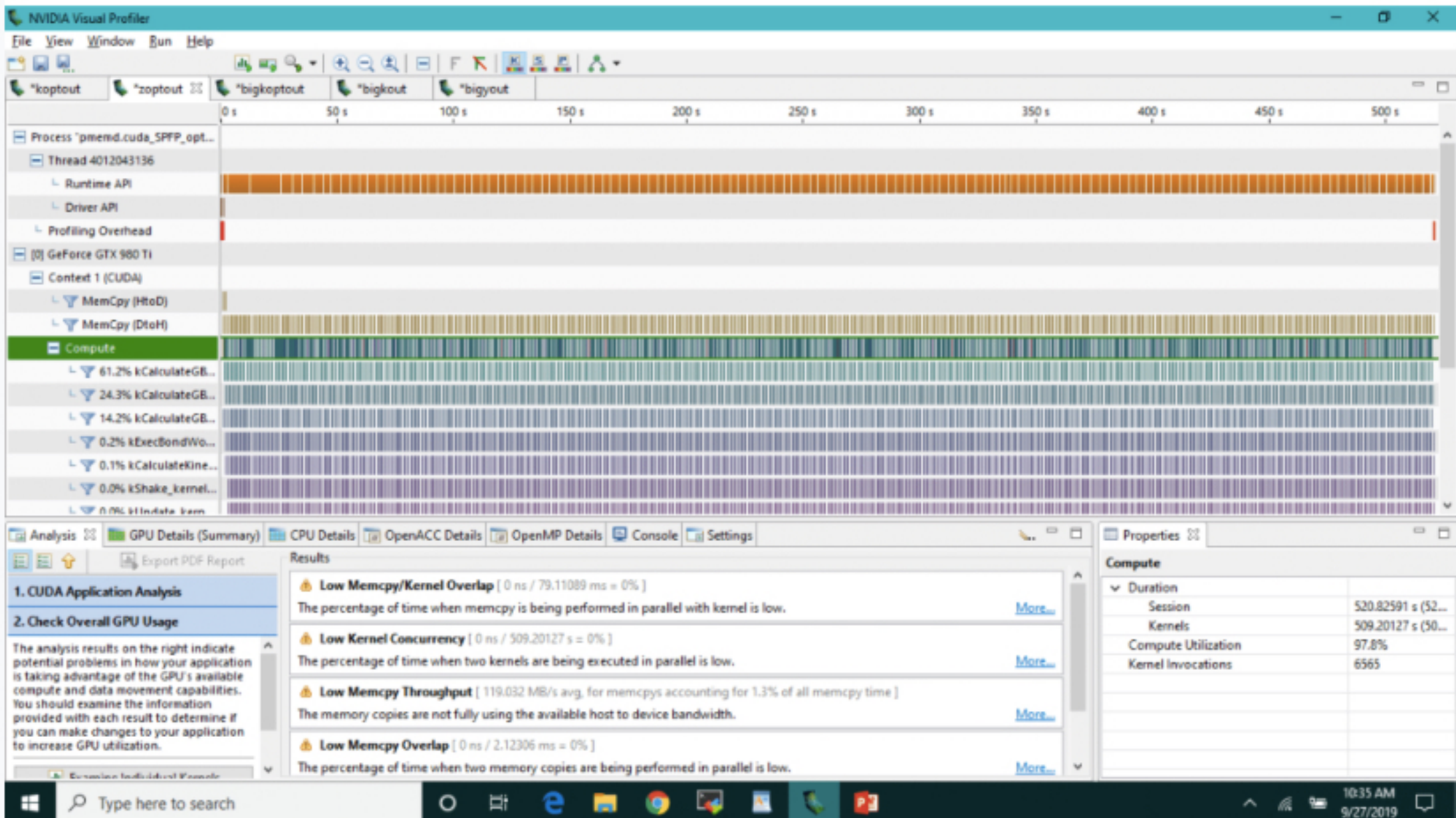


Project 1
 Report 1



Top-Down View Filter... 51.53% of data is shown due to applied filters. Search...

Symbol Name	Self, %	App, %	Total, %	Module Name
0xeccc5	-	62.22	62.22	/home/ubuntu/sdcard/slambench-demo/build/kfusion/kfusion-benchmark-cuda
0xf6f4320d	-	-	30.50	/lib/arm-linux-gnueabi/libc-2.19.so
[Broken backtraces]	0.02	-	7.27	[Broken backtraces]



Demo using NVTX and Nsight

```

1 using BenchmarkTools
2 using CUDA
3
4 using QXContexts
5
6 function main(args)
7     file_path = @__DIR__
8     dsl_file = joinpath(dirname(dirname(file_path)), "examples/ghz/ghz_5.qx")
9     input_file = joinpath(dirname(dirname(file_path)), "examples/ghz/ghz_5.jld2")
10
11     cg, _ = parse_dsl_files(dsl_file, input_file)
12
13     # get time on gpu
14     ctx_gpu = QXContext{CuArray{ComplexF32}}(cg)
15     set_open_bonds!(ctx_gpu)
16     # run to ensure all is precompiled
17     t = NVTX.@range "Warm up" begin @elapsed ctx_gpu() end
18     @info "GPU warmup ran in $t"
19     CUDA.@profile NVTX.@range "Run iteration" begin
20         ctx_gpu()
21     end
22     nothing
23 end
24
25 main(ARGS)

```

THANK YOU

Instructor: Dr. Momme Allalen
www.nvidia.com/dli