# CUDA Performance Measurement

## *Jiri Kraus (NVIDIA)*

Mitglied der Helmholtz-Gemeinschaft

# Why Performance Measurement Tools?

- You can only improve what you measure

    - Need to identify:

        - Hotspots: Which function takes most of the run time?

        - Bottlenecks: What limits the performance of the Hotspots?

- Manual timing is tedious and error prone

    - Possible for small application like jacobi and matrix multiplication

    - Impractical for larger/more complex application

- Access to hardware counters (PAPI, CUPTI)

# The command line profiler nvprof

- Simple launcher to get profiles of your application

- Profiles CUDA Kernels and API calls

```
> nvprof ./jacobi

======== NVPROF is profiling jacobi...

======== Command: jacobi

Jacobi (serial)

[…] snip

======== Profiling result:
 Time(%)      Time  Calls       Avg       Min       Max  Name
   72.14   352.65ms   1000   352.65us   350.48us   354.94us  Jacobi_86_gpu
   26.02   127.23ms   1000   127.23us    93.48us   128.34us  Jacobi_74_gpu
    0.84     4.09ms   1000     4.09us     4.04us     4.36us  Jacobi_96_gpu_red
    0.61     3.00ms   1009     2.97us     2.78us    56.16us  [CUDA memcpy HtoD]
    0.39     1.91ms   1002     1.91us     1.82us    52.41us  [CUDA memcpy DtoH]
```
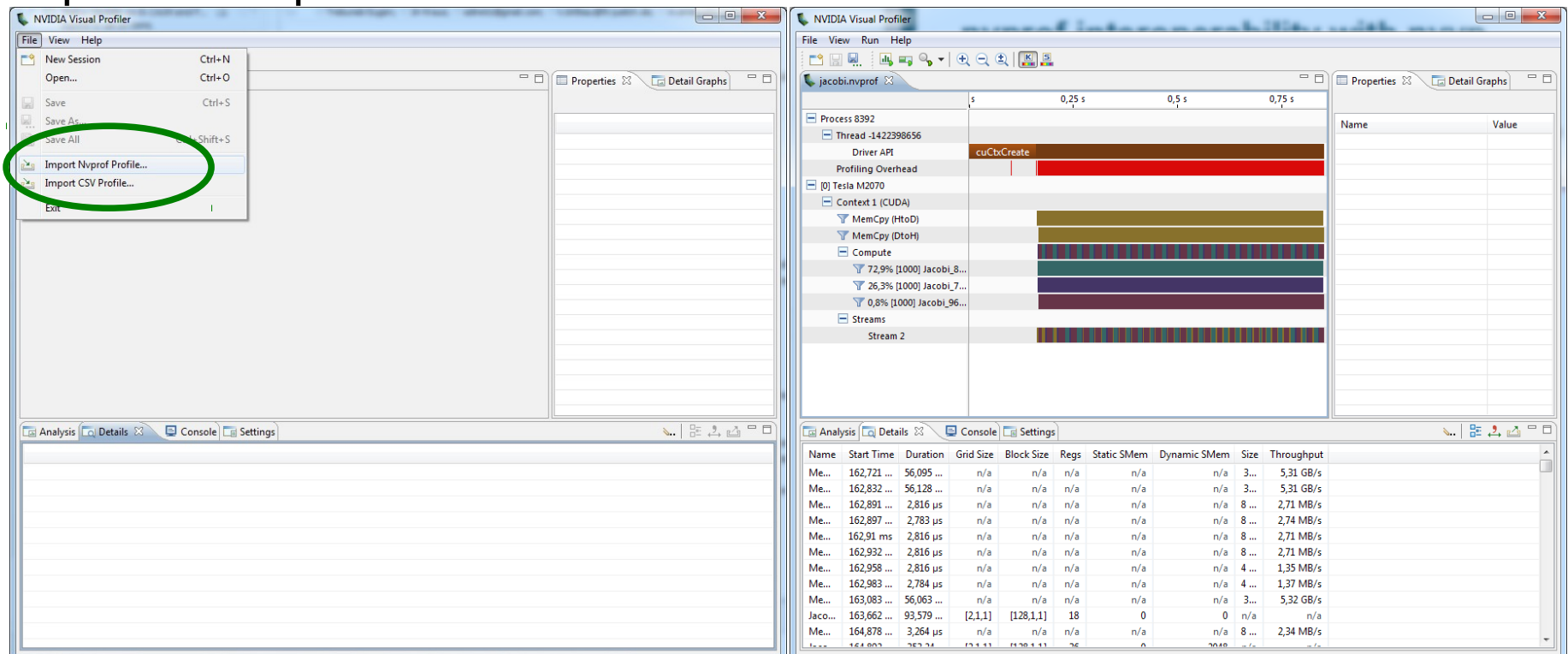
# nvprof interoperability with nvvp

- nvprof can write the application profile to nvvp compatible file:

  `nvprof -o jacobi.nvprof ./jacobi`

- Import in nvvp

# nvprof important command-line options

```
Options:

  -o,  --output-profile <filename>

                               Output the result file which can be imported later

                               or opened by the NVIDIA Visual Profiler.

  --events <event names>

                               Specify the events to be profiled on certain

                               device(s). Multiple event names separated by comma

                               can be specified. Which device(s) are profiled is

                               controlled by the '--devices' option. Otherwise

                               events will be collected on all devices.

                               For a list of available events, use

                               '--query-events'.

  --query-events               List all the events available on each device.

  -h,  --help                  Print this help information.
```
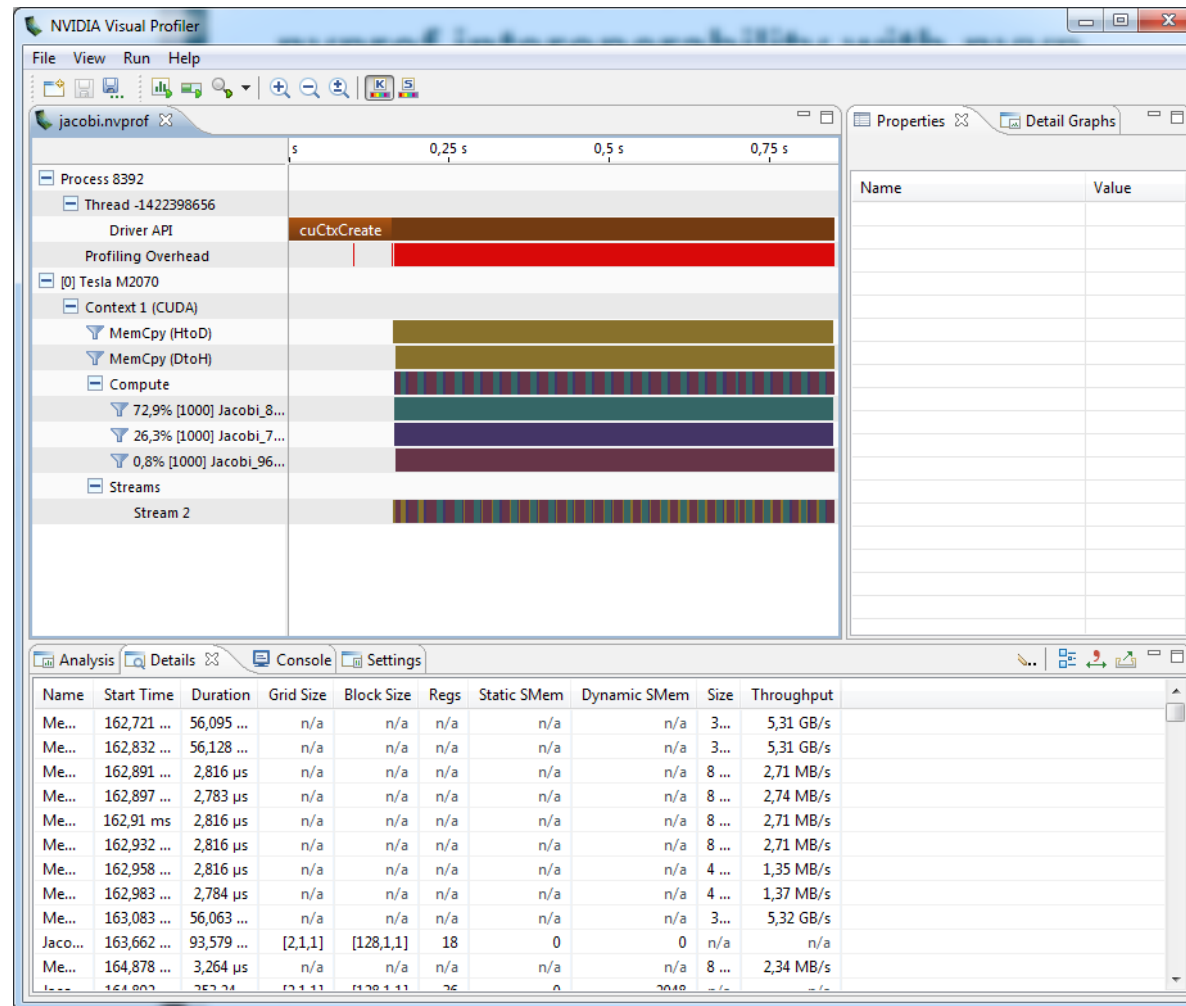
# nvvp introduction

# Task 1: Analyze Jacobi Timeline

- Start jacobi with nvprof and write profile to file

- Import profile into nvvp

- Compare the profiles with and without data region.

# Task 2: Analyze matrix multiplication example with nvvp

- Start new session in nvvp with the matrix multiplication example
- Run the "Uncoalesced Global Memory" experiment

# Task 3: Analyze matrix multiplication example with nvprof

- Start matrix multiplication with nvprof and collect gld_inst_32bit event

- Import profile into nvvp

- Read the value of gld_inst_32bit and compare it to the size of the input matrices and the number of executed floating point multiplications

  *Hint: M*N*K and M*N + N*K*

# Cheat Sheet

- Start nvprof

```
nvprof -o <output-profile> ./a.out
```

- Start nvvp

```
nvvp
```

- profiler users guide

```
http://docs.nvidia.com/cuda/profiler-users-guide/index.html
```