Now, this is just a simulation of what the blocks will look like once they've assembled

- Profiling

# Developer Tools - Profilers



**NSIGHT**        **NVVP**        **NVPROF**

**NVIDIA Provided**

**TAU**        **VampirTrace**

**3rd Party**

https://developer.nvidia.com/performance-analysis-tools

# NVPROF EXERCISE

- Command Line Profiler
    - Compute time in each kernel
    - Compute memory transfer time
    - Collect metrics and events
    - Support complex process hierarchy's
    - Collect profiles for NVIDIA Visual Profiler
    - No need to recompile
- **$ nvprof ./myadd**
    - View available metrics

        **$ nvprof --query-metrics**

        - View global load/store efficiency

            $ nvprof --metrics gld_efficiency,gst_efficiency ./myadd

        - Store a timeline to load in NVVP

            **$ nvprof –o profile.timeline ./myadd**

        - Store analysis metrics to load in NVVP

            **$ nvprof –o profile.metrics --analysis-metrics ./myadd**
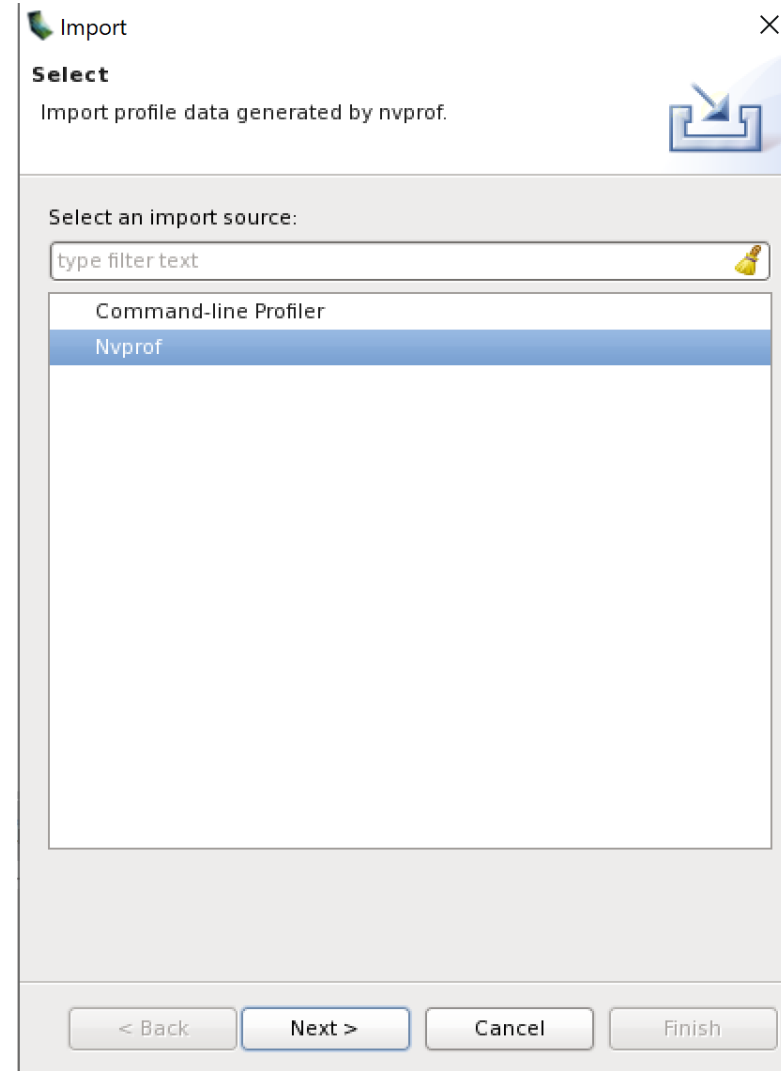
# NVPROF EXERCISE

- Either in interactive session

- or

- Use job submission script given in hw1 as a template and include nvprof as part of the "/usr/bin/time mpirun" command
  - /usr/bin/time mpirun -np 1 nvprof --log-file profile.txt ./myadd
  - /usr/bin/time mpirun -np 1 nvprof --metrics all -o nvvp_profile.nvvp ./myadd
  - /usr/bin/time mpirun -n 1 nvprof -o profile.timeline ./myadd
  - /usr/bin/time mpirun -n 1 nvprof –o profile.metrics --analysis-metrics ./myadd

- **At command line after loading the cuda module launch nvvp:**
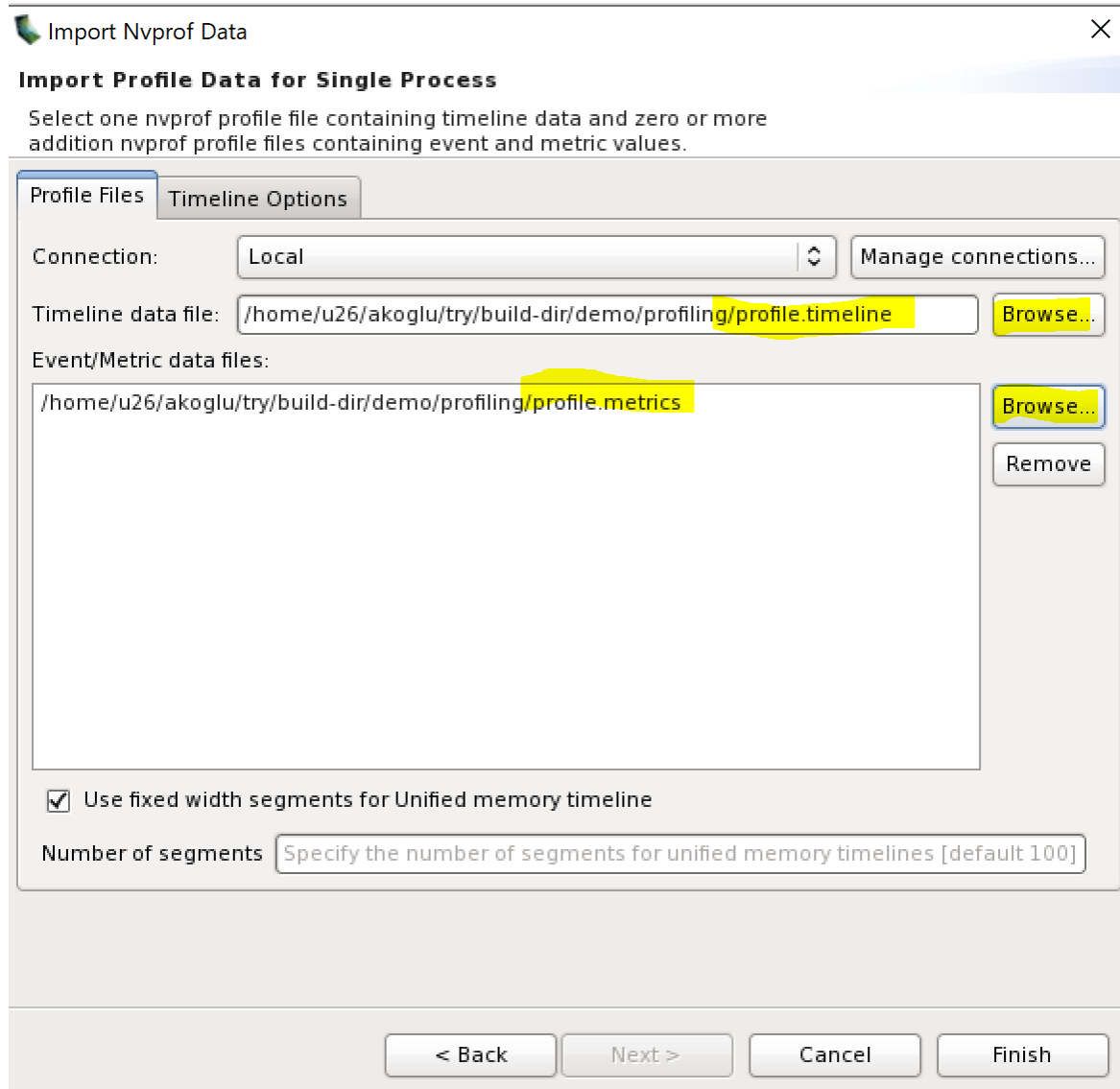
     **$ nvvp &**

# Import nvprof into NVVP

- Launch nvvp
- Click
  - File/ Import/ Nvprof/ **Next**/ Single process/ Next / Browse
  (next to Timeline data file)
  - Select profile.timeline
- Add Metrics to timeline
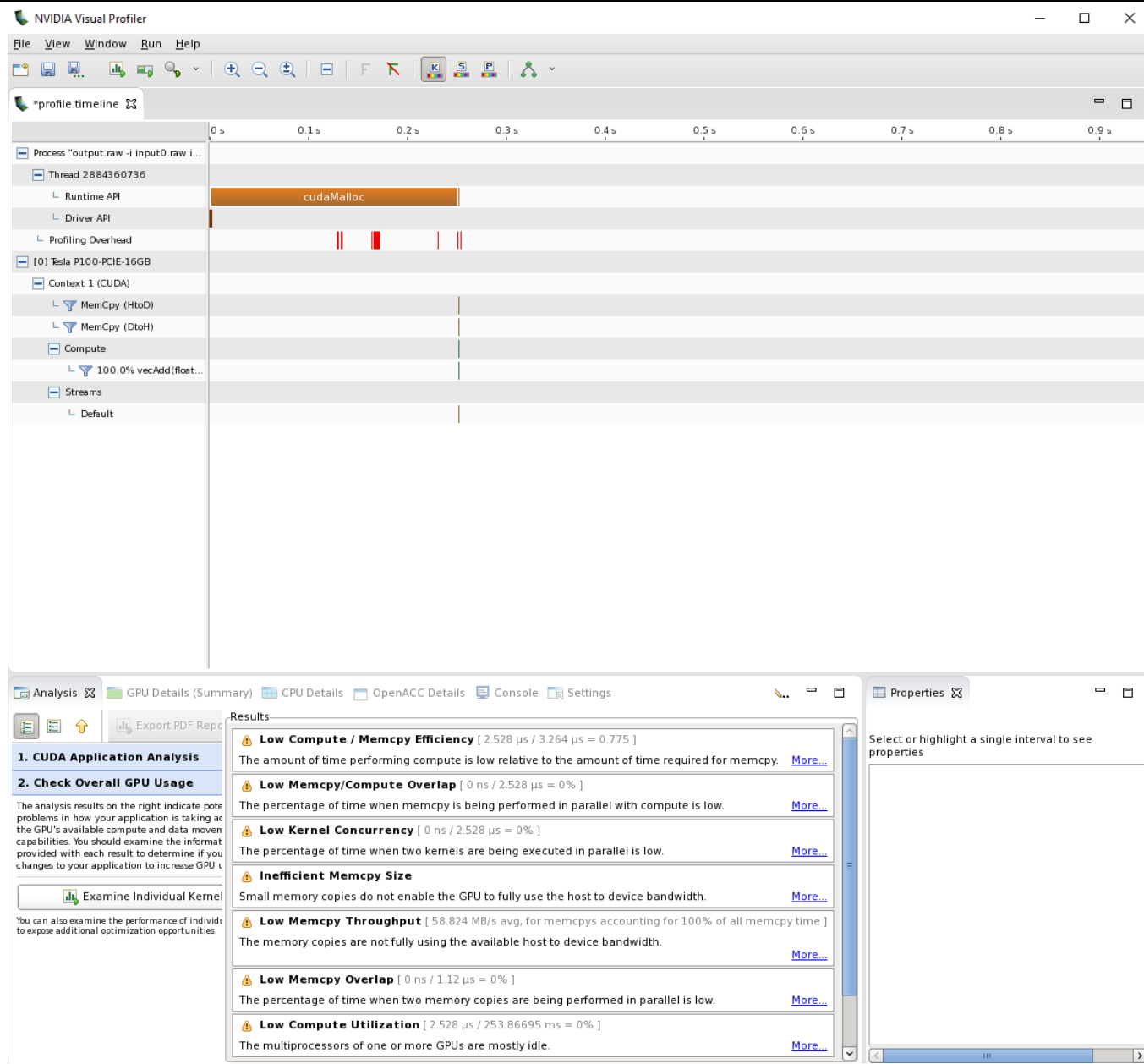  - Click on 2nd Browse (for Event/Metric data files)
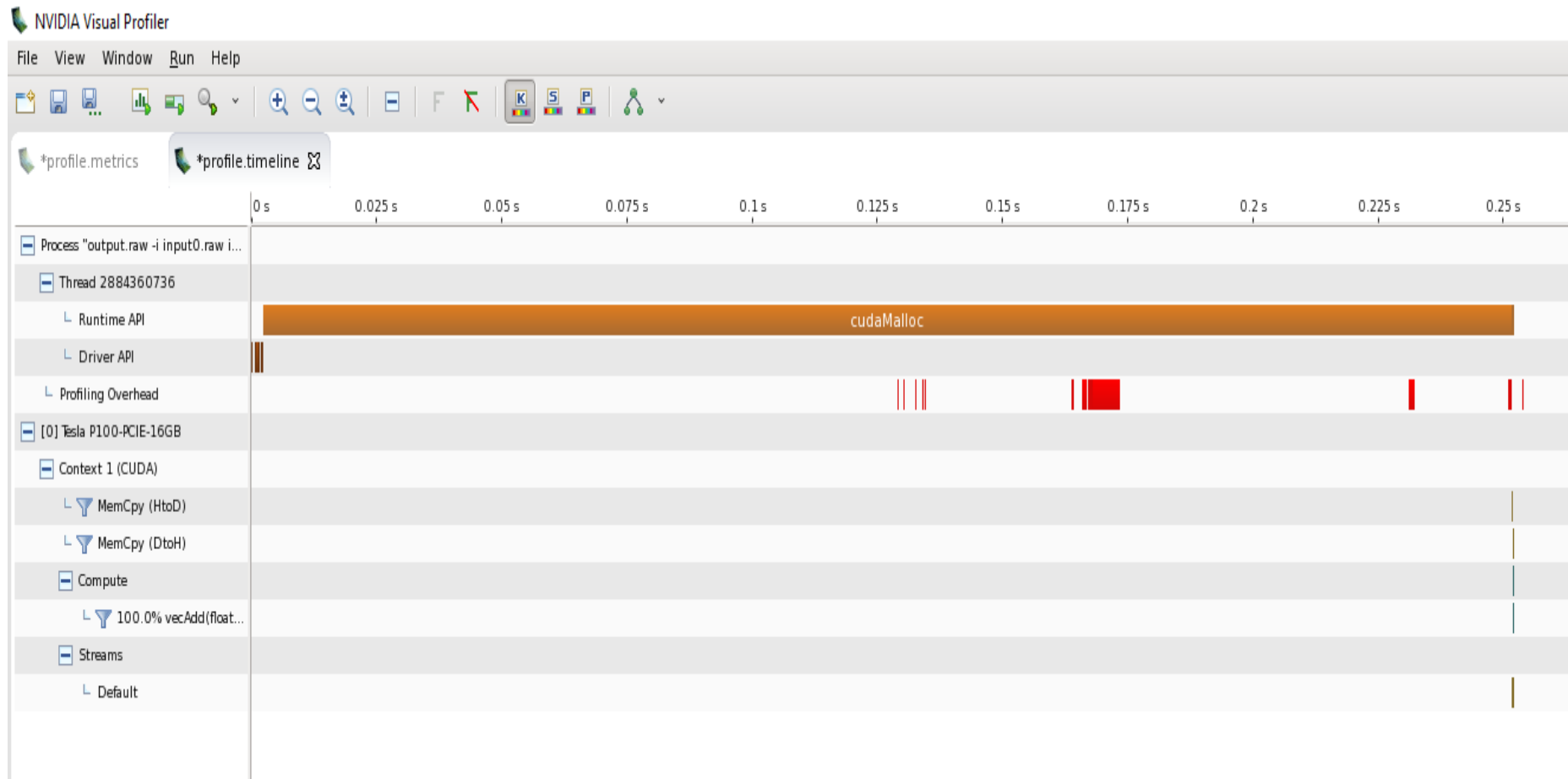  - Select profile.metrics

# Load your profile files

- – Select **profile.timeline**
- – Click on 2nd Browse (for Event/Metric data files)
- – Select **profile.metrics**

- Click Finish

- Explore Timeline
  - – Control + mouse drag in timeline to zoom in
  - – Control + mouse drag in measure bar (on top) to measure time

### Import Nvprof Data ✕

**Import Profile Data for Single Process**

Select one nvprof profile file containing timeline data and zero or more addition nvprof profile files containing event and metric values.

| Profile Files | Timeline Options |

Connection: Local    Manage connections...

Timeline data file: /home/u26/akoglu/try/build-dir/demo/profiling/profile.timeline   Browse...

Event/Metric data files:

/home/u26/akoglu/try/build-dir/demo/profiling/profile.metrics

Browse...

Remove

☑ Use fixed width segments for Unified memory timeline

Number of segments   Specify the number of segments for unified memory timelines [default 100]

< Back   Next >   Cancel   Finish

# NVVP

# Timing Analysis Window

# Guided Analysis Window

# Guided Analysis – GPU Usage

# Unguided Analysis



Instructions:
1. Click on a kernel
2. On Analysis tab click on the unguided analysis


2. Click Analyze All
   Explore metrics and properties

# Unguided Analysis

# NSIGHT

- **CUDA enabled IDE**
  - Source code editor: syntax highlighting
  - Build Manger
  - Visual Debugger
  - Visual Profiler
- **Linux/Macintosh**
  - Editor = Eclipse
  - Debugger = cuda-gdb with a visual wrapper
  - Profiler = NVVP
- **Windows**
  - Integrates directly into Visual Studio
  - Profiler is NSIGHT VSE

# Nsight Eclipse Edition

- **After loading the cuda module run:**

  **$ nsight &**

https://docs.nvidia.com/cuda/nsight-eclipse-edition-getting-started-guide/index.html

# NSIGHT: Importing project

Instructions:
1. Run nsight
     Select default workspace
2. Click File / New / Makefile Project With Existing CodeTest
3. Enter Project Name and select the project directory
4. Click Finish
5. Right Click On Project / Properties / Run Settings / New / C++ Application
6. Browse for executable
7. In Project Explorer double click on .cu and explore source
8. Click on the build icon
9. Click on the run icon
10. Click on the profile icon

# Before profiling

- Certain kinds of errors cause CUDA programs to complete, but crash under profiling

- Check your program with cuda-memcheck if code behaves incorrectly under profiling

```
cuda-memcheck ./my-program ...
```
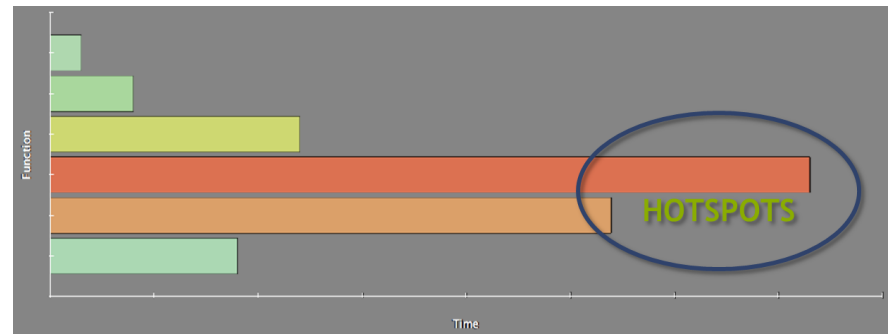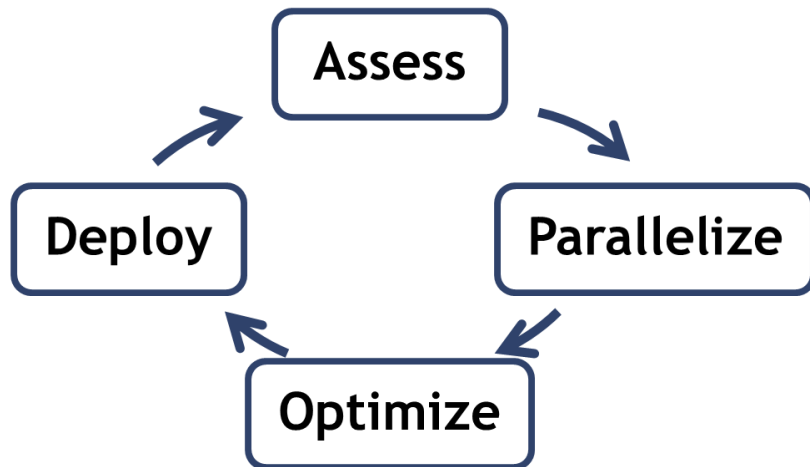
# Before profiling

- **Compile device code with optimizations**
  - Optimizations dramatically change performance
  - Remove "-G" device-debug flag from nvcc

- **Compile device code with line information**
  - Add "-lineinfo" flag to nvcc

```
nvcc –G main.cu          nvcc –lineinfo main.cu
```

# Profiling, Debugging Tools Summary

- Performance programmers rely on debugging and profiling tools

  – You should setup your environment and choose the tools most convenient for you.

  – For assignments you may not need them much but for your project you will need these tools

    - quantify how well you are utilizing compute resources, memory bandwidth

# Next

- **Threads, Thread blocks, workload management**