

Profiling GPU codes with Nsight

Sergey Mashchenko

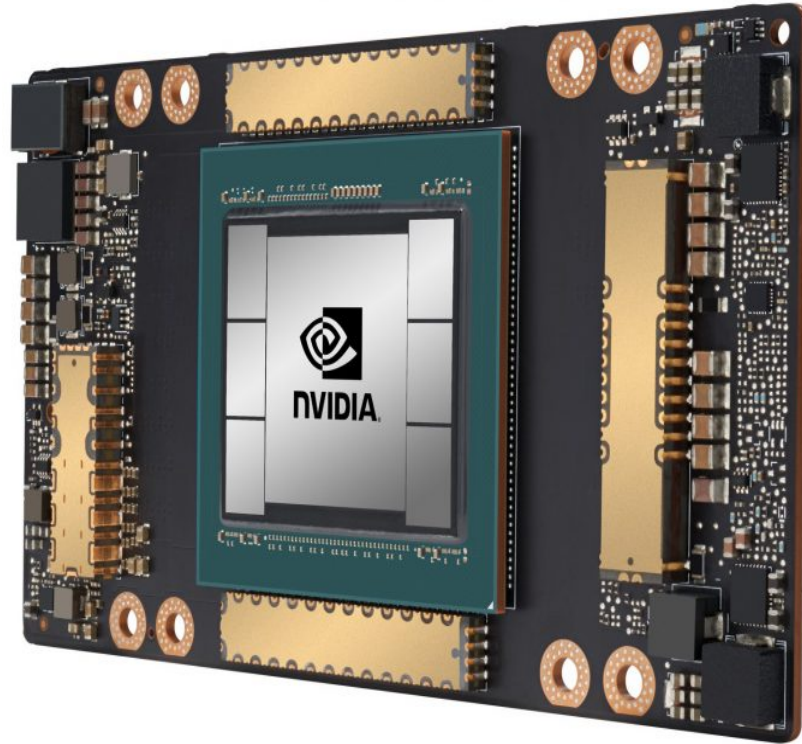
syam@sharcnet.ca

SHARCNET / Compute Ontario / Alliance

May 18, 2022

Overview

- Why?
- Where?
- How?
- Demo



Why to profile GPU codes?

- GPUs are significantly more expensive, and less available, than CPUs
- Not all research codes / algorithms are well suited for GPU acceleration.
- As a consequence, profiling is a critical step in a GPU code development, and should start from the very first kernel you write.

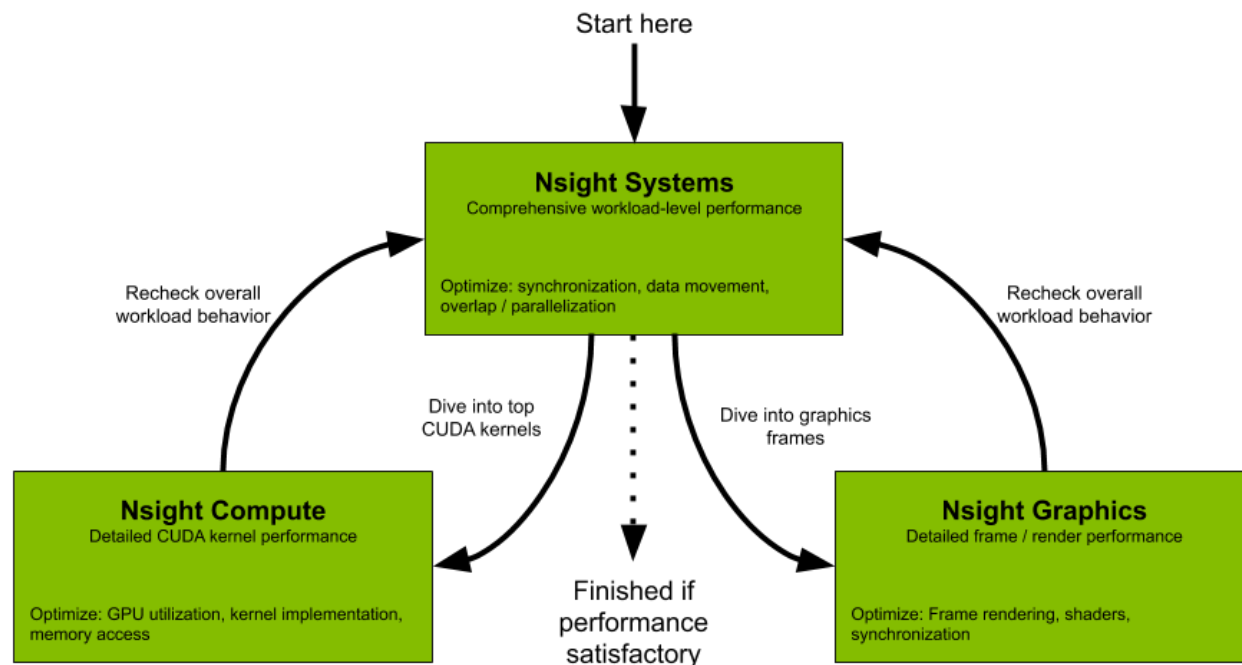
Tools

- NVIDIA has produced and retired quite a few GPU profilers.
- The following profilers are no longer maintained (though still available, work up to V100):
 - nvprof: command-line function-level profiler, for both GPU and CPU parts of the code
 - nvvp: graphical (GUI) profiler

Tools: Nsight

- NVIDIA also maintained for years their Nsight suite of products (IDE / debugger / profiler).
 - Nsight Eclipse edition (Linux, Mac)
 - Nsight Visual Studio Edition (Windows)
- Since 2018, Nsight profilers became also available as three stand alone packages: Nsight Compute, Nsight Systems, and Nsight Graphics.

Nsight packages



Where to run Nsight

- If you have a fairly capable recent GPU inside your laptop / PC, you can install and use Nsight suite (Eclipse for Linux/Mac, Visual Studio for Windows) on your own computer.
- But if the goal is to optimize your code for Alliance national systems, you definitely want to run Nsight remotely on our clusters.
 - Command line (CLI) tools can be submitted as jobs
 - Interactive GUI tools can be used on compute nodes allocated with `salloc`, with either X11 or VNC connections (more about that later).
- If your internet is too slow, you can run CLI Nsight tools on a cluster, then analyze the results on your computer using GUI Nsight.

How to use GUI tools on clusters

- X11 forwarding (MobaXterm for Windows, Xquartz for Mac):

```
$ ssh -Y graham.computecanada.ca  
$ salloc --x11 ...
```

- VNC on a compute node (requires two terminal windows)

```
[login_node]$ salloc ...  
[gra123]$ export XDG_RUNTIME_DIR=${SLURM_TMPDIR}  
[gra123]$ vncserver  
[your_PC]$ ssh graham.computecanada.ca -NL 5902:gra123:5901
```

- On your PC, launch TigerVNC viewer
- Enter the destination: localhost:5902



VNC helper script

```
#!/bin/bash
```

```
# Required for vncserver:
```

```
export XDG_RUNTIME_DIR=${SLURM_TMPDIR}
```

```
# Starting VNC server, recording the channel:
```

```
N=$(vncserver 2>&1 |grep "^New" |cut -d: -f3)
```

```
# Computing the remote port:
```

```
Rport=$((5900 + $N))
```

```
# Printing the command for the local computer:
```

```
echo "ssh $USER@${SLURM_CLUSTER_NAME}.computecanada.ca -NL 5902:${SLURMD_NODENAME}:${Rport}"
```

NVIDIA profilers on our clusters

- Loaded when a cuda module is loaded, e.g.
`$ module load cuda/11.4`
- Old tools (up to Volta: P100, V100):
 - nvprof (CLI)
 - nvvp (GUI)
- Nsight tools (work on P100*, V100, T4, A100):
 - ncu, ncu-ui: Nsight Compute (CLI / GUI; V100 and up)
 - nsys, nsys-ui: Nsight Systems (CLI / GUI; P100 and up)

Compiling code

- Compile the code as usual (after loading the cuda module)
 - The only extra compiler switch required is `-lineinfo` (do not use `-G` – that one is for debugging)

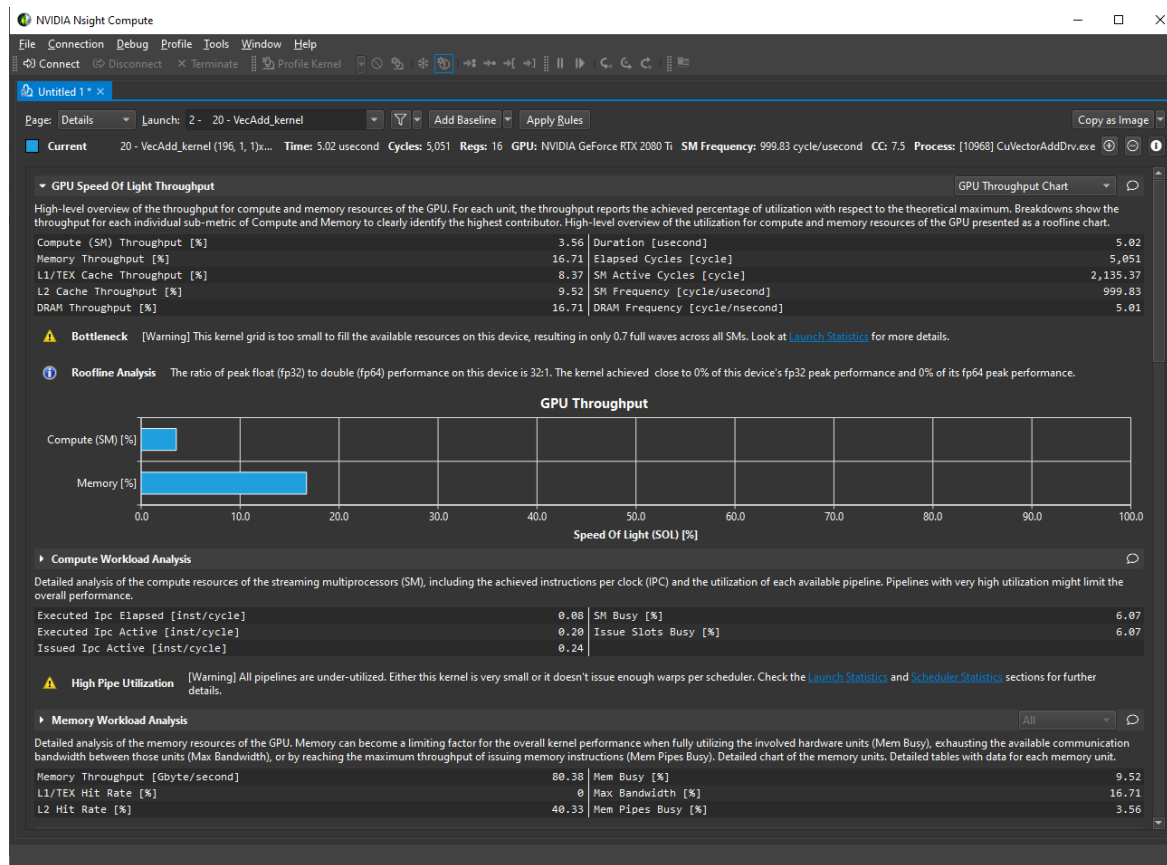
```
$ module load cuda
```

```
$ nvcc -O2 -arch=sm_70 -lineinfo code.cu
```

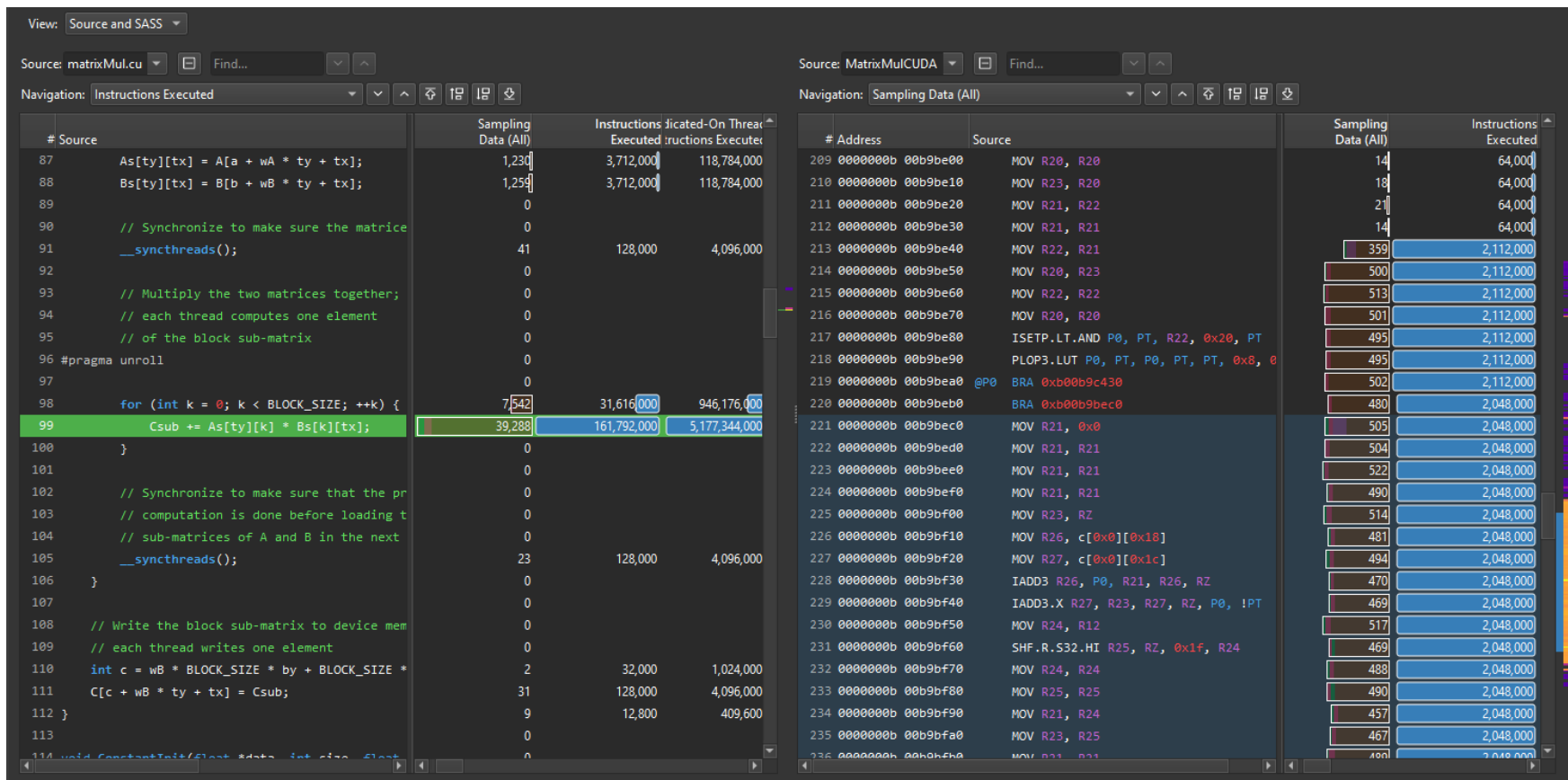
Nsight Compute

- Command: ncu, ncu-ui
\$ ncu -o output_file code
- Typically the first step in profiling a GPU code
- Allows one to maximize the performance of individual kernels
- Full documentation:
<https://docs.nvidia.com/nsight-compute/NsightCompute>

Kernel metrics



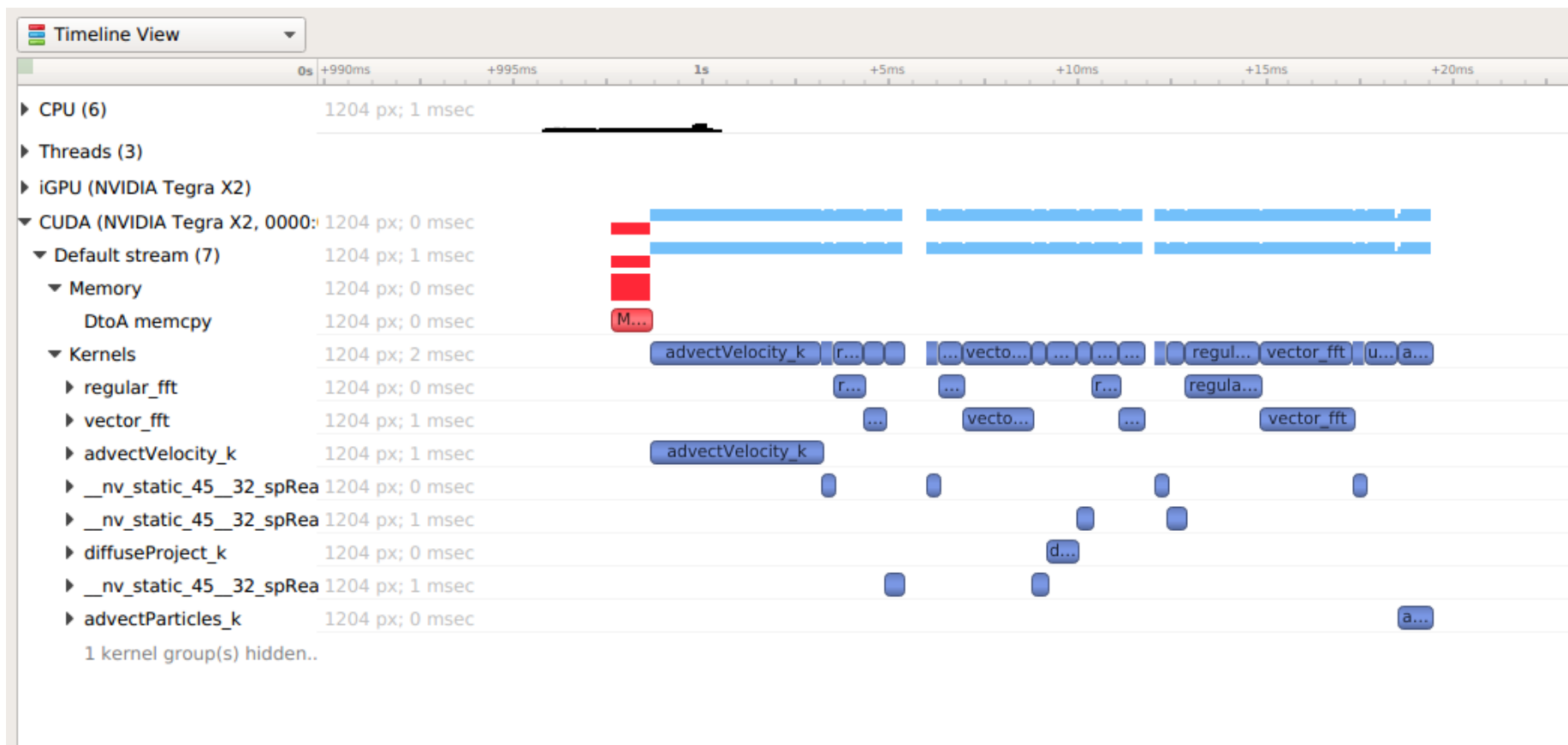
Line-by-line info



Nsight Systems

- Command: `nsys`, `nsys-ui`
`$ nsys profile -o output_file code`
- Typically the last step in profiling a GPU code
- Fine-tuning the interactions between kernels, memcopies, CPU code etc, both synchronous and asynchronous
- Full documentation:
<https://docs.nvidia.com/nsight-systems/UserGuide>

CUDA trace

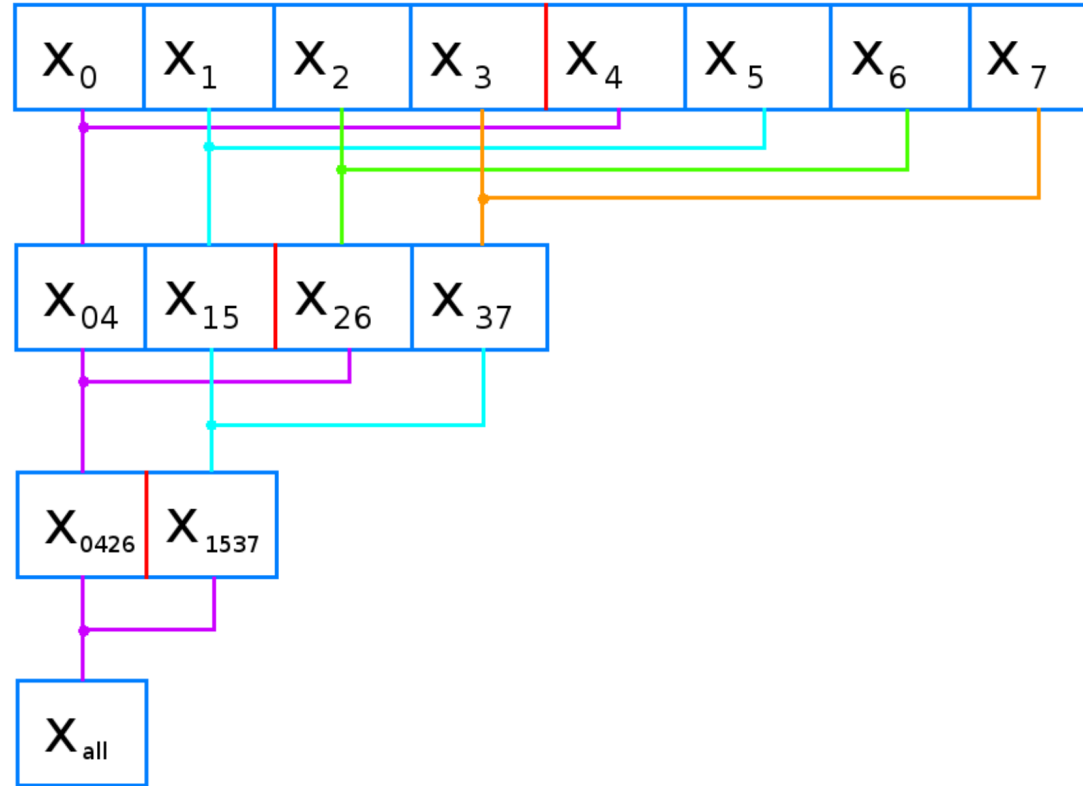


Live demo

Reduction code

- The code adds up elements of a long vector on GPU.
- The first version uses a very slow method: serialized summation using `atomicAdd()` function.
- The second version uses the proper way: parallel summation using binary reduction.

Binary reduction



Staged copy/compute code

- The first version of the code first copies a long vector to GPU, then carries out independent per-element computations.
- The second version hides some of the costs of copying data to GPU by running copying and computing in parallel, using two GPU streams.

Staged copy/compute algorithm

One stream scenario:



You need two streams for this:



Thank you!