

Name: Qian Liyang
NetID: 3190110719
Section: ZIUI

ECE 408/CS483 Milestone 2 Report

1. Show output of rai running Mini-DNN on the basic GPU convolution implementation for batch size of 1k images. This can either be a screen capture or a text copy of the running output. Please do not show the build output. (The running output should be everything including and after the line "*Loading fashion-mnist data...Done*").

```
Test batch size: 1000
Loading fashion-mnist data...Done
Loading model...Done
Conv-GPU==
Layer Time: 62.5606 ms
Op Time: 1.69811 ms
Conv-GPU==
Layer Time: 50.8151 ms
Op Time: 6.47096 ms

Test Accuracy: 0.886
```

2. For the basic GPU implementation, list Op Times, whole program execution time, and accuracy for batch size of 100, 1k, and 10k images.

Batch Size	Op Time 1	Op Time 2	Total Execution Time	Accuracy
100	0.182508 ms	0.6552 ms	0m1.143s	0.86
1000	1.69811ms	6.47096ms	0m9.645s	0.886
10000	16.7434 ms	64.4371 ms	1m37.553s	0.8714

3. List all the kernels that collectively consumed more than 90% of the kernel time and what percentage of the kernel time each kernel did consume (start with the kernel that consumed the most time, then list the next kernel, until you reach 90% or more).

Generating CUDA Memory Operation Statistics...
 CUDA Kernel Statistics (nanoseconds)

Time(%)	Total Time	Instances	Average	Minimum	Maximum	Name
100.0	81608506	2	40804253.0	16774710	64833796	conv_forward_kernel
0.0	2880	2	1440.0	1376	1504	prefn_marker_kernel
0.0	2560	2	1280.0	1280	1280	do_not_remove_this_kernel

4. List all the CUDA API calls that collectively consumed more than 90% of the API time and what percentage of the API time each call did consume (start with the API call that consumed the most time, then list the next call, until you reach 90% or more).

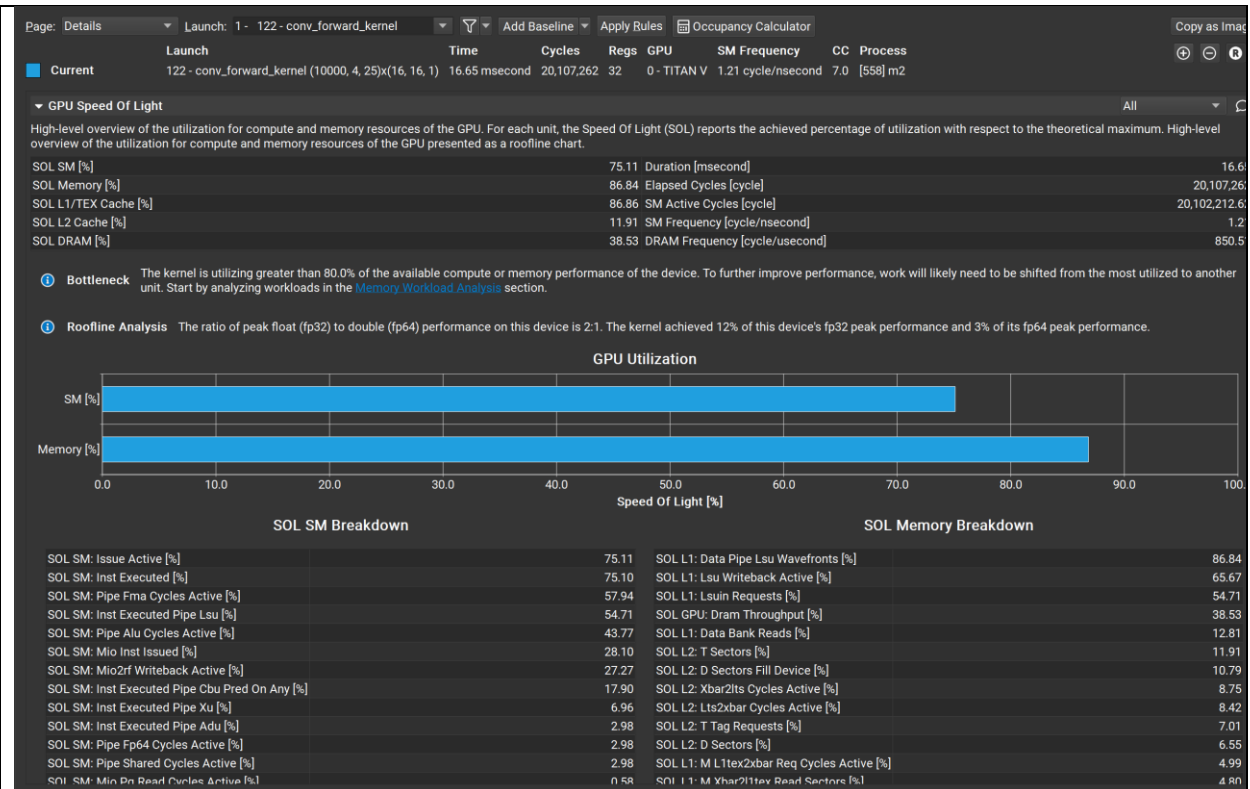
Generating CUDA API Statistics...
 CUDA API Statistics (nanoseconds)

Time(%)	Total Time	Calls	Average	Minimum	Maximum	Name
68.1	1183718378	8	147964797.3	43349	561717210	cudaMemcpy
26.0	451257818	8	56407227.2	201793	441706478	cudaMalloc
4.7	81646576	8	10205822.0	980	64837910	cudaDeviceSynchronize
1.0	17701844	6	2950307.3	18529	17551338	cudaLaunchKernel
0.2	2925108	8	365638.5	158512	757110	cudaFree

More than 90% is cudaMemcpy and cudaMalloc

5. Explain the difference between kernels and CUDA API calls. Please give an example in your explanation for both.
1. CUDA API calls for data transferring is larger than the kernel data transferring. For example, in above figures, the average for cudaMemcpy is 1.5e9 ns but for kernel, the total running time is just 4.1e8, which is only 27% of cudaMemcpy. That's because the gpu have larger bandwidth than the main memory for host on most machines.
 2. Cuda API is called by the CPU and kernel is called by GPU. The CUDA API interface in nsys shows the functions running on cpu, which are related to CUDA library, for example, the cudaMemcpy . The CUDA Kernel interface shows the functions running on GPU, for example, the conv_forward_kernel
6. Show a screenshot of the GPU SOL utilization

*Note: the screen shot is for 10,000 data size
 Layer 1:*



Layer 2:

