**Name:** Mingyuan Zhao
**NetID:** mz55
**Section:** OD1

# 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*").

```
*Running bash -c "time ./m2 1000"   \\ Output will appear after run is complete
.
Test batch size: 1000
Loading fashion-mnist data...Done
Loading model...Done
Conv-GPU==
Layer Time: 268.082 ms
Op Time: 5.46809 ms
Conv-GPU==
Layer Time: 229.296 ms
Op Time: 28.6412 ms

Test Accuracy: 0.886


real    0m11.504s
user    0m10.450s
sys     0m1.033s
```
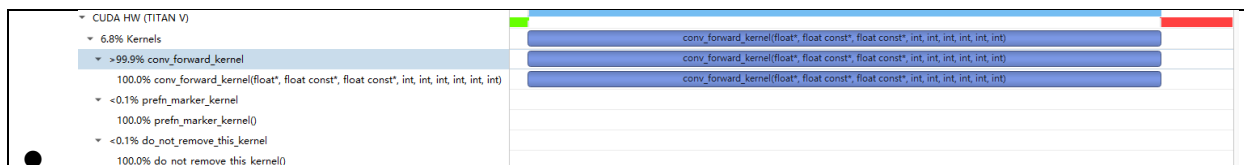
2. For the basic GPU implementation, list Op Times, whole program execution time, and accuracy for batch size of 100, 1k, and 10k images.

| Batch Size | Op Time 1 | Op Time 2 | Total Execution Time | Accuracy |
|---|---|---|---|---|
| 100 | 0.17625ms | 1.23577ms | 0m2.376s | 0.86 |
| 1000 | 5.46809ms | 28.6412ms | 0m11.504s | 0.886 |
| 10000 | 16.1087ms | 63.3525ms | 1m37.329s | 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).

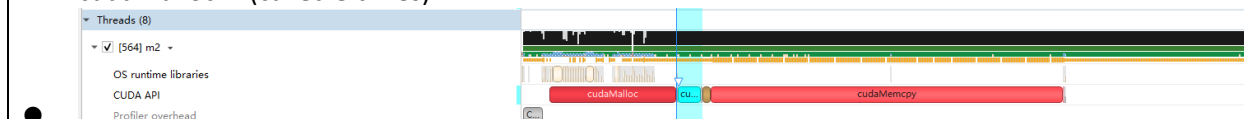This is the only kernel and being called twice:
- conv_forward_kernel

| | |
|---|---|
| ▼ CUDA HW (TITAN V) | |
| ▼ 6.8% Kernels | conv_forward_kernel(float*, float const*, float const*, int, int, int, int, int, int) |
| ▼ >99.9% conv_forward_kernel | conv_forward_kernel(float*, float const*, float const*, int, int, int, int, int, int) |
| 100.0% conv_forward_kernel(float*, float const*, float const*, int, int, int, int, int, int) | conv_forward_kernel(float*, float const*, float const*, int, int, int, int, int, int) |
| ▼ <0.1% prefn_marker_kernel | |
| 100.0% prefn_marker_kernel() | |
| ▼ <0.1% do_not_remove_this_kernel | |
| 100.0% 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).

There are 4 different CUDA API calls, the ones that consume 90% of the API time are:
- cudaMemcpy (called 8 times)
- cudaMalloc   (called 8 times)



5.  Explain the difference between kernels and CUDA API calls. Please give an example in your explanation for both.

A kernel is a C++ function. When called, "are executed N times in parallel by N different CUDA threads, as opposed to only once like regular C++ functions."
An API is a set of definitions, protocols and tools for building software. In CUDA, "it provides C and C++ functions that execute on the host to allocate and deallocate device memory, transfer data between host memory and device memory, manage systems with multiple devices, etc."
The CUDA API call's part is the section contains the time of CPU using, while the profiling result is the time that the GPUs taking. The total time of the CUDA API call is from the moment it is launched to the moment it is completed, so will overlap with executing kernels.

6.  Show a screenshot of the GPU SOL utilization

Kernel1:

Kernel2: