**Name:** *Ding-Hsiang Huang*
**NetID:** *dhhuang3*
**Section:** *AL2*

# 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*").



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 | *0.176167ms* | *0.637506 ms* | *0m1.245s* | *0.86* |
| 1000 | *1.65375ms* | *6.32831ms* | *9.752s* | *0.886* |
| 10000 | *16.1505 ms* | *63.687 ms* | *1m35.159s* | *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).

| |
|---|
| *conv_forward_kernel 100.0%* |
| 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). |
| cudaMemcpy 77.4%<br>cudaMalloc 17.8% |
| 5. Explain the difference between kernels and CUDA API calls. Please give an example in your explanation for both. |
| Kernels are code that executed in device(kernel) like" conv_forward_kernel", which contained the calculation that executed in parallel.<br><br>CUDA API calls are those code that calling between host and device. CudaMemcpy, cudaMalloc, cudaDeviceSynchronize… all of them are calls made by the code into the CUDA driver. |
| 6. Show a screenshot of the GPU SOL utilization |