

Milestone 1

Team Name: mallard

Names: Gauri Konanoor, Hamilton Silberg, Sankruth Kota

NetIDs: gmk2, hfs2, srkota2

Introduction

In milestone 1 of the final project, we run a batch forward pass on test data using the MXNet framework. We compare the performance of this operation on an AMD64 CPU as opposed to a Pascal GPU. We then checked which kernels consumed 90% of the computation time.

Kernel Calls

	Type	Time(%)	Time	Calls	Avg	Min	Max	Name
GPU activities:		38.53%	37.192ms	20	1.8596ms	1.0240us	35.087ms	[CUDA memcpy HtoD]
		21.60%	20.847ms	1	20.847ms	20.847ms	20.847ms	volta_scudnn_128x32_relu_interior_nn_v1
		19.79%	19.104ms	1	19.104ms	19.104ms	19.104ms	void cudnn::detail::implicit_convolve_sg
	, int, float, cudnn::detail::implicit_convolve_sgemm<float, float, int=1024, int=5, int=5, int=3, int=3, int=3, int							
		7.73%	7.4589ms	2	3.7294ms	24.543us	7.4343ms	void cudnn::detail::activation_fw_4d_ker
	::activation_fw_4d_kernel<float, float, int=128, int=1, int=4, cudnn::detail::tanh_func<float>>, cudnnTensorStruct*,							
		7.03%	6.7878ms	1	6.7878ms	6.7878ms	6.7878ms	volta_sgemm_128x128_tn
		4.60%	4.4392ms	1	4.4392ms	4.4392ms	4.4392ms	void cudnn::detail::pooling_fw_4d_kernel

The kernel calls adding up to greater than 90% include the combination of:

[CUDA memcpy HtoD]

volta_scudnn_128x32_relu_interior_nn_v1

cudnn::detail::implicit_convolve_sgemm

cudnn::detail::activation_fw_4d_kernel

volta_sgemm_128x128_tn

and,

[CUDA memcpy HtoD]

Volta_scudnn_128x32_relu_interior_nn_v1

cudnn::detail::implicit_convolve_sgemm

cudnn::detail::activation_fw_4d_kernel

cudn::detail::pooling_fw_4d_kernel.

API Calls

API calls:	39.88%	2.72075s	22	123.67ms	12.913us			
1.41769s	cudaStreamCreateWithFlags							
	34.21%	2.33356s	24	97.232ms	65.529us	2.32859s	cudaMemGetInfo	
	21.89%	1.49314s	19	78.586ms	288ns	389.22ms	cudaFree	

The API calls covering greater than 90 percent consumption by time is the combination of:

cudaStreamCreateWithFlags
cudaMemGetInfo
cudaFree

Kernel Calls v. API Calls

Kernels are C functions completed repeated N times in parallel by N different CUDA threads within the GPU device. API calls are specific actions called for the GPU to enact by the CPU, whether it be to allocate or transfer memory, or even to run an instance of a designated kernel with specific thread dimensions. Overall, the main difference is that the kernel deals with GPU tasks that are much more repeatable than those dealt with in API calls, with API calls being more generalized to GPU operation.

In general, it is beneficial to analyse the optimisations from parallelising computations using kernels, which is what we do in this milestone.

Test Results

First, we ran our forward pass test on a CPU. As we can see below, our runtime duration was **13.87 seconds**.

```
Successfully installed mxnet
$ Running /usr/bin/time python m1.1.py
Loading fashion-mnist data... done
Loading model... done
New Inference
EvalMetric: {'accuracy': 0.8177}
19.79user 4.22system 0:13.87elapsed 173%CPU (0avgtext+0avgdata 5954848maxresident)k
0inputs+2856outputs (0major+1578349minor)pagefaults 0swaps
```

Upon running the same computation on a GPU, we can see that our runtime duration is **4.68 seconds**.

```
Successfully installed mxnet
$ Running /usr/bin/time python m1.2.py
Loading fashion-mnist data... done
Loading model... done
New Inference
EvalMetric: {'accuracy': 0.8177}
4.11user 2.52system 0:04.68elapsed 141%CPU (0avgtext+0avgdata 2847636maxresident)k
8inputs+4568outputs (0major+707070minor)page
faults 0swaps
```

Parallelising our simple machine learning computation optimised our runtime by 66%.

Operation Time One:

Milestone 2

In this milestone, we implemented a sequential CPU version of a forward-propagation path convolutional layer.

```
Loading model... done
New Inference
Op Time: 24.896752
Op Time: 105.804736
Correctness: 0.8152 Model: ece408
142.33user 9.44system 2:15.91elapsed 111%CPU (0avgtext+0avgdata 5950984maxresident)k
```

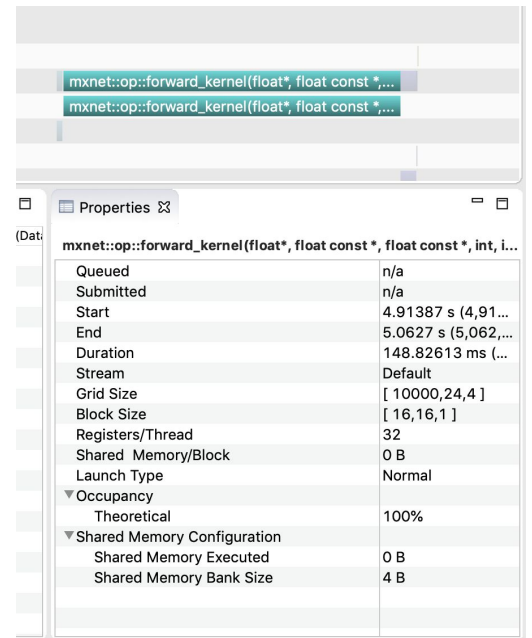
The operation times as shown above is **24.896 seconds** for the first time that layer is invoked, and approximately **105.8 seconds** for the second time the layer is invoked. The total runtime of our program is about **2 minutes and 15 seconds**.

Milestone 3

In this milestone, our team worked on developing a GPU convolution. Our convolution for this milestone had a correctness of **0.85 for 100 images**, a correctness of **0.827 for 1000 images**, and a correctness of **0.8171 for 10000 images** which aligned with the expected output.

After performing some analysis using nvvp and nvprof, we can better understand our kernel. We can see that the total time it took to execute the kernel was approximately **0.14883 seconds**. We can also see that we performed **no shared memory optimizations** since our shared memory usage was 0.

The visual profiler also provided key data to how our kernel executed. After taking a more in-depth look we can see that we launched our kernel with a **grid size of [10,000, 24, 4]** and a **block size of [16, 16, 1]**. With a **theoretical occupancy of 100 percent**, our team was fairly certain that our kernel utilized our hardware SM's well, and that we had the correct block and grid sizes.



Zooming out a little, the fact that we have no shared memory usage tells us that we have **poor bandwidth utilization**. The way our team determined this was by taking a look at the higher level graphical data shown above. As you can see, **cudaMemGetInfo** is the single cuda API call that takes the longest time to perform, and we're hoping to implement shared memory and some of the tiling strategies that we learned about in class.

Also, it seems as though our implementation was **not as parallel** as we had thought since we did receive an indication from nvprof informing us that we could have had stronger GPU usage to speed up our kernel, as shown in the image below.

 **Low Compute Utilization** [227.81303 ms / 5.11552 s = 4.5%]

The multiprocessors of one or more GPUs are mostly idle.

[More...](#)