

ECE 408
Spring 2019

Final Project Lab Report

Team: sphavwv
Xinyi Guo (xinyig2, City Scholars)
Kevin So (kevinso2, On Campus)
Jiamin Zhu (jiaminz2, City Scholars)

MILESTONE 1

1. Include a list of all kernels that collectively consume more than 90% of the program time.

- [CUDA memcpy HtoD]
- void cudnn::detail::implicit_convolve_sgemm...
- volta_cgemm_64x32_tn
- void op_generic_tensor_kernel
- Volta_sgemm_128x128_tn
- void fft2d_c2r_32x32

2. Include a list of all CUDA API calls that collectively consume more than 90% of the program time.

- cudaStreamCreateWithFlags
- cudaMemGetInfo
- cudaFree

3. Include an explanation of the difference between kernels and API calls.

One of the highest percentages from API calls came from cudaStreamCreateWithFlags which is a host device function. API calls seem to deal with things like syncthreads and organizing the GPU's execution while the GPU deals with the actual computing tasks. Regarding the time percentages for the GPU, we see that the memory copy from host to device is taking up the most due to global memory read and write cycles and the fact that it is called 20 times.

4. Show output of rai running MXNet on the GPU & CPU

```
[Running dist-install] -user -e /mxnet/python  
DEPRECATION: Python 2.7 will reach the end of its life on January 1st, 2020. Please upgrade your Python as Python 2.7 won't be maintained after that date. A future port for Python 2.7.  
Obtaining file:///mxnet/python  
Requirement already satisfied: graphviz<0.9.0,>=0.8.1 in /root/.local/lib/python2.7/site-packages (from mxnet==1.3.1) (0.8.4)  
Requirement already satisfied: numpy<=1.15.0,>=1.8.2 in /root/.local/lib/python2.7/site-packages (from mxnet==1.3.1) (1.14.0)  
Requirement already satisfied: requests<2.19.0,>=2.18.4 in /root/.local/lib/python2.7/site-packages (from mxnet==1.3.1) (2.18.4)  
Requirement already satisfied: idna<2.7,>=2.5 in /root/.local/lib/python2.7/site-packages (from requests<2.19.0,>=2.18.4->mxnet==1.3.1) (2.6)  
Requirement already satisfied: urllib3<1.23,>=1.21.1 in /root/.local/lib/python2.7/site-packages (from requests<2.19.0,>=2.18.4->mxnet==1.3.1) (1.22)  
Requirement already satisfied: certifi>=2017.4.17 in /root/.local/lib/python2.7/site-packages (from requests<2.19.0,>=2.18.4->mxnet==1.3.1) (2018.11.29)  
Requirement already satisfied: chardet<3.0.0,>=3.0.2 in /root/.local/lib/python2.7/site-packages (from requests<2.19.0,>=2.18.4->mxnet==1.3.1) (3.0.4)  
Installing collected packages: mxnet  
  Running setup.py develop for mxnet  
Successfully installed mxnet  
[Running /usr/bin/time python m1.1.py  
Loading fashion-mnist data... done  
Loading model... done  
New Inference  
EvalMetric: {'accuracy': 0.8236}  
8.75User 3.41System 0:04.95elapsed 245%CPU (0avgtext+  
0avgdata 2468868maxresident)k  
0inputs+2824outputs (0major+69669minor)pagefaults 0swaps  
[Running /usr/bin/time python m1.2.py  
Loading fashion-mnist data... done  
Loading model... done  
New Inference  
EvalMetric: {'accuracy': 0.8236}  
4.48User 3.30System 0:04.20elapsed 185%CPU (0avgtext+0avgdata 2861040maxresident)k  
0inputs+  
1728outputs (0major+664979minor)pagefaults 0swaps
```

CPU: List program run time

8.75cpu seconds (user mode) 3.41 cpu seconds system (kernel mode) 4.95s elapsed real time

GPU: List program run time

4.48 cpu seconds (user mode) 3.30cpu seconds (kernel mode) 4.20s elapsed real time

MILESTONE 2

```
Successfully installed mxnet
$ Running /usr/bin/time python m2.1.py
Loading fashion-mnist data... done
Loading model... done
New Inference
Op Time: 2.505609
Op Time: 7.466405
Correctness: 0.8397 Model: ece408
15.41user 4.29system 0:11.54elapsed 170%CPU (0avgtext+0avgdata 1620160maxresident)k
0inputs+2824outputs (0major+617337minor)pagefaults 0swaps
$ Running python m2.1.py
Loading fashion-mnist data... done
Loading model... done
New Inference
Op Time: 2.481413
Op Time: 7.484842
Correctness: 0.8397 Model: ece408
$ Running python m2.1.py 1000
Loading fashion-mnist data... done
Loading model... done
New Inference
Op Time: 0.245897
Op Time: 0.752752
Correctness: 0.852 Model: ece408
$ Running python m2.1.py 100
Loading fashion-mnist data... done
Loading model... done
New Inference
Op Time: 0.031788
Op Time: 0.076578
Correctness: 0.84 Model: ece408
$ The build folder has been uploaded to http://s3.amazonaws.com/files.rai-project.com/userda
```

Runtimes: User 15.41, System 4.29, Elapsed 11.5

MILESTONE 3

```


# Running ./mnist_time python m3.1.py
Loading Fashion-mnist data... done
Loading model... done
New Inference
Op Time: 0.018376
Op Time: 0.070265
Correctness: 0.8397 Model: ece408
4.57user 3.17system 0:04.47elapsed 172%CPU (0avgtext+0avgdata 2828240maxresident)
0i
inputs+4576outputs (0major+659360minor)pagefaults 0swaps
# Running python m3.1.py 100
Loading Fashion-mnist data... done
Loading model... done
New Inference
Op Time: 0.000209
Op Time: 0.000696
Correctness: 0.84 Model: ece408
# Running python m3.1.py 1000
Loading Fashion-mnist data... done
Loading model... done
New Inference
Op Time: 0.001868
Op Time: 0.006977
Correctness: 0.852 Model: ece408
# Running nvprof python m3.1.py
Loading Fashion-mnist data... done
==546== NVPDF is profiling process 546, command: python m3.1.py


```



```


# running nvprof python m3.1.py
Loading fashion-mnist data... done
==368== NVPDF is profiling process 368, command: python m3.1.py
Loading model... done
New Inference
Op Time: 0.018457
Op Time: 0.071863
Correctness: 0.8397 Model: ece408
==368== Profiling application: python m3.1.py
==368== Profiling result:
      Type    Time    Calls      Avg      Min      Max   Name
GPU activities:  78.58%  90.248ms      2  45.124ms  18.414ms  71.834ms  mxnet::op::forward kernel<float*, float const *, float const *, int, int, int, int, int, int>
               14.53%  16.693ms      20  834.63us  1.1200us  16.257ms  [CUDA memcpy HtoD]
               2.18%  2.5048ms      2  1.2524ms  21.408us  2.4834ms  volta sgemm_32x128 tn
               2.11%  2.4244ms      2  1.2122ms  735.48us  1.6889ms  void mshadow::expr::Plan<mshadow::gpu, int=4, float>, float, mshadow::expr::ScalarExp<float>, mshadow::Tensor<mshadow::gpu, int=4, float>, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=4, int)
               1.41%  1.621ms      2  810.54us  21.888us  1.5992ms  void op_generic_tensor_kernel<int=2, float, float, float, float, int=256, cudnnGenericOp_t>, cudnnNanPropagation_t=0, cudnnDimOrder_t=0, int=1>(cudnnTensorStruct, float*, cudnnTensorStruct, float const *, cudnnTensorStruct, float const *, cudnnTensorStruct, float const *, float, float, float, float, dim3<ray, reducedDivisorArray>)
               0.92%  1.0526ms      1  1.0526ms  1.0526ms  1.0526ms  void cudnn::detail::pooling_fw_4d_kernel<float, float, cudnn::detail::maxPooling_func<float, cudnnNanPropagation_t=0>, int=0, bool=0>, cudnnTensorStruct, cudnnPoolingStruct, float, cudnnPoolingStruct, int, cudnn::reduced divisor, float>
               0.14%  157.76us      1  157.76us  157.76us  157.76us  void mshadow::expr::Plan<mshadow::gpu, int=2, float>, float, mshadow::expr::Plan<mshadow::gpu, int=2, float>, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2, int)
               0.07%  75.295us      1  75.295us  75.295us  75.295us  void mshadow::expr::ScalarExp<float>, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2)
               0.07%  28.192us      13  2.1680us  1.1840us  6.5928us  void mshadow::expr::Plan<mshadow::gpu, int=2, float>, float, mshadow::expr::Plan<mshadow::gpu, int=2, float>, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2)
               0.02%  24.032us      2  12.016us  2.5280us  21.304us  void mshadow::expr::Plan<mshadow::gpu, int=2, float>, float, mshadow::expr::Plan<mshadow::gpu, int=2, float>, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2)
               0.01%  11.455us      10  1.1450us  991ns  1.6960us  [CUDA memset]
               0.01%  5.7600us      1  5.7600us  5.7600us  5.7600us  [CUDA memcpy DtoH]
               0.00%  4.8630us      1  4.8630us  4.8630us  4.8630us  void mshadow::expr::ReduceWithAxis<xp>(mshadow::red::maximum, mshadow::Tensor<mshadow::gpu, int=3, float>, float, int=3, bool=1, int=2>, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2)
API calls:  42.94%  3.11387s      22  141.54ms  12.916us  1.62078s  cudaStreamCreateWithFlags
           31.88%  2.30554s      22  104.80ms  90.630us  2.30082s  cudaMemGetInfo
           21.39%  1.55128s      18  86.182ms  834ns  419.05ms  cudaFree
1.33%  96.657ms      912  105.98us  310ns  19.753ms  cudaFuncSetAttribute
1.28%  92.703ms      6  15.451ms  5.5350us  71.841ms  cudaDeviceSynchronize
0.47%  34.216ms      9  3.8018ms  24.097us  16.410ms  cudaMemcpy2DAsync
0.23%  16.579ms      66  251.19us  5.6160us  7.7308ms  cudaMalloc
0.19%  13.798ms      216  63.879us  891ns  12.162ms  cudaEventCreateWithFlags


```

MILESTONE 4

Optimization 1: Loop Unrolling with restrict variable and variable tuning

This is based on the concept that the fastest memory transfers and processing happen at the register level. By unrolling the loop into local registers, we surpass the speeds of using shared memory and increase thread activity. When memory accesses are independent, instead of doing 1 read at a time per thread, we can do many more reads per thread depending on the input size. This speeds things up especially when there is local data re-use. On the downside, memory in our registers is relatively much lower than any other storage in the GPU and CPU and that is our limitation.

Restricted is a pointer aliasing optimization. It is similar to const in a sense that the programmer makes an agreement with the program to never point to the same memory location from 2 different arrays. By doing so, we relieve the compiler from checking it's pointer accesses which boosts the performance speed of the kernel.

Optimization with all restricted variables and inner/middle for loop unrolls:

Optimization with all restricted variables and inner for loop unroll:

Optimization with outer loop unroll:

```

f Running /usr/bin/time python m3.1.py
Loading fashion-mnist data... done
Loading model... done
New Inference
Op Time: 0.006358
Op Time: 0.015087
Correctness: 0.8397 Model: ece408
4.43usec 3.32system 0:04.27elapsed 181%CPU (0avgtext+0avgdata 2830628maxresident)
0inputs+4584outputs (0major+659469minor)pagefaults 0swaps
3 Running nvpprof python m3.1.py
Loading fashion-mnist data... done
==369== NVPROF is profiling process 369, command: python m3.1.py
Loading model... done
New Inference
Op Time: 0.006474
Op Time: 0.015145
Correctness: 0.8397 Model: ece408
==369== Profiling application: python m3.1.py
==369== Profiling result:
      Type    time(%)   Time    Calls     Avg     Min     Max   Name
GPU activities: 46.70% 21.545ms   2  10.772ms  6.4312ms 15.114ms mxnet::op::forward_kernel<float>, float const *, float const *, int, int, int, int, int, int
36.73% 16.949ms   20  847.0ms  1.0560us 16.438ms [CUDA memcpy DtoD]
5.15% 2.3740ms   2  1.1870ms  725.24us 1.6487ms void mshadow::cuda::MapPlanLargeKernel<mshadow::svs::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=4, float>, float<int>, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=4, int)
      5.01% 2.3128ms   2  1.1564ms  21.280us 2.2910ms volta_sgemm_32x128_tn
3.49% 1.6123ms   2  806.15us  22.016us 1.5903ms void op_generic_tensor_kernel<int=2, float, float, float, float, int=256, cudnnGenericOp_t=7, cudnnNanPropagation_t=0>
cudnnDimOrder_t=0, int=>(cudnnTensorStruct, float*, cudnnTensorStruct, float const *, cudnnTensorStruct, float const *, float, float, float, float, dim3Array, reducedDivisorArray)
2.26% 1.0405ms   1  1.0405ms  1.0405ms 1.0405ms void cudnn::detail::pooling_fw_4d_kernel<float, float, float, cudnn::detail::maxpooling_func<float, cudnnNanPropagation_t=0>, int=0, bool=>, cndnTensorStruct*, cudnnPoolingStruct, float, cudnnPoolingStruct, int, cudnn::reduced divisor, float)
      0.33% 154.24us   1  154.24us  154.24us void mshadow::cuda::MapPlanLargeKernel<mshadow::svs::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2, int)
      0.16% 75.327us   1  75.327us  75.327us void mshadow::cuda::SoftmaxKernel<int=8, float, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>>(mshadow::gpu, int=2, unsigned int)
      0.06% 27.968us   13  2.1510us  2.1610us 6.5280us void mshadow::cuda::MapPlankernel<mshadow::svs::saveto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2)
      0.05% 23.935us   2  11.967us  2.5608us 21.3719us void mshadow::cuda::MapPlankernel<mshadow::svs::plusto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=1, float>, float, int=2, int=1>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2)
      0.02% 11.392us   10  1.1390us  960ns 1.9200us [CUDA memset]
      0.02% 7.8720us   1  7.8720us  7.8720us 7.8720us [CUDA memcpy DtoH]
      0.01% 5.0240us   1  5.0240us  5.0240us 5.0240us void mshadow::cuda::MapPlanKernel<mshadow::svs::saveto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2)

```

Optimization with only x and k arrays restricted and inner for loop unroll:

```
Successfully installed mxnet
Running /usr/bin/time python m3.1.py
Loading fashion-mnist data... done
Loading model... done
New Inference
Op Time: 0.006581
Op Time: 0.016771
Correctness: 0.8397 Model: ece408
4.47user 3.12system 0:04:29elapsed 176%CPU (0avgtext+0avgdata 2845384maxresident)k
0inputs+4576outputs
s (0major+663584minor)pagefaults 0swaps
m (0major+0minor)memcopy 0memmove
Loading fashion-mnist data... done
==370== NVPROF is profiling process 370, command: python m3.1.py
Loading model... done
New Inference
Op Time: 0.006699
Op Time: 0.016928
Correctness: 0.8397 Model: ece408
==370== Profiling application: python m3.1.py
==370== Profiling result:
      Type  Time(%)   Time    Calls    Avg     Min     Max   Name
GPU activities:  48.59% 23.552ms      2  11.776ms  6.6552ms 16.897ms mxnet::op::forward_kernel<float>, float const *, float const *, int, int, int, int, int)
          35.64% 17.276ms      20  863.81us 1.1208us 16.884ms [CUDA memcpy HtoD]
          4.88% 2.3665ms      2  1.1833ms 723.68us 1.6428ms void mshadow::expr::MapPlanLargeKernel<mshadow::sv::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=4, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=4, int)
          4.78% 2.3166ms      2  1.1583ms 20.672us 2.2959ms volta_sgemm_32x128_tn
          3.33% 1.6132ms      2  806.59us 22.239us 1.5909ms void op_generic_tensor_kernel<int=2, float, float, float, float, int=256, cudnnGenericOp_t=7, cudnnNanPropagation_t=0, cudnnDimOrder_t=0, int=1>(cudnnTensorStruct, float*, cudnnTensorStruct, float const *, cudnnTensorStruct, float const *, float, float, float, float, dimArray, reducedDivisorArray)
          2.15% 0.4033ms      1  1.0433ms 1.0433ms void cudnn::detail::pooling_fw_4d kernel<float, float, cudnn::detail::maxpooling_func<float, cudnnNanPropagation_t=0>, int=0, bool=0>, cudnnTensorStruct*, cudnnPoolingStruct, float, cudnnPoolingStruct, int, cudnn::reduced divisor, float)
          0.32% 152.86us      1  152.86us 152.86us 152.86us void mshadow::expr::MapPlanLargeKernel<mshadow::sv::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2)
          0.16% 75.551us      1  75.551us 75.551us 75.551us void mshadow::expr::SoftmaxKernel<int=8, float, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>>(mshadow::gpu, mshadow::Shape<int=2>, float, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>>(mshadow::gpu, int=2, float), float>>(mshadow::gpu, int=2, float, unsigned int, mshadow::Shape<int=2>, int=2)
          0.06% 27.840us      13  2.1410us 1.1840us 6.4640us void mshadow::expr::MapPlanKernel<mshadow::sv::saveto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>>(mshadow::gpu, mshadow::expr::Plan<mshadow::ScalarExp<float>, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2)
          0.05% 23.872us      2  11.936us 2.5600us 21.312us void mshadow::expr::MapPlanKernel<mshadow::sv::plusto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>>(mshadow::gpu, mshadow::expr::Broadcast1DExp<mshadow::Tensor<mshadow::gpu, int=1, float>, float, int=2, int=1, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2)
          0.02% 11.264us      10  1.1260us 992ns 1.6000us [CUDA memset]
          0.01% 5.8560us      1  5.8560us 5.8560us 5.8560us [CUDA memcpy DtoH]
```

Optimization with only y restriction:

```
Running /usr/bin/time python m3.1.py
Loading fashion-mnist data... done
Loading model... done
New Inference
B is: 10000 M is: 6 H is: 48 C is: 1 K is: 5 Op Time: 0.006352
B is: 10000 M is: 22 W is: 22 C is: 6 K is: 5 Op Time: 0.015088
Correctness: 0.8397 Model: ece408
4.36user 3.41system 0:04:29elapsed 181%CPU (0avgtext+0avgdata 2828144maxresident)
0inputs+4584outputs (6
major+659026minor)pagefaults 0swaps
Running nvprof python m3.1.py
Loading fashion-mnist data... done
==368== NVPROF is profiling process 368, command: python m3.1.py
Loading model... done
New Inference
B is: 10000 M is: 6 H is: 48 W is: 48 C is: 1 K is: 5 Op Time: 0.006493
B is: 10000 M is: 16 H is: 22 W is: 22 C is: 6 K is: 5 Op Time: 0.015150
Correctness: 0.8397 Model: ece408
==368== Profiling application: python m3.1.py
==368== Profiling result:
      Type  Time(%)   Time    Calls    Avg     Min     Max   Name
GPU activities:  47.23% 21.553ms      2  10.777ms  6.4386ms 15.115ms mxnet::op::forward_kernel<float>, float const *, float const *, int, int, int, int, int)
          36.05% 16.453ms      20  823.63us 1.0560us 16.070ms [CUDA memcpy HtoD]
          5.19% 2.3693ms      2  1.1847ms 721.82us 1.6479ms void mshadow::expr::MapPlanLargeKernel<mshadow::sv::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=4, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=4, int)
          5.05% 2.3043ms      2  1.5222ms 21.440us 2.2822ms volta_sgemm_32x128_tn
          3.53% 1.6122ms      2  806.12us 21.664us 1.5906ms void op_generic_tensor_kernel<int=2, float, float, float, float, int=256, cudnnGenericOp_t=7, cudnnNanPr
cudnnDimOrder_t=0, int=1>(cudnnTensorStruct, float*, cudnnTensorStruct, float const *, cudnnTensorStruct, float const *, float, float, float, float, dimArray, reducedDivisorArray)
          2.29% 1.0446ms      1  1.0446ms 1.0446ms void cudnn::detail::pooling_fw_4d kernel<float, float, cudnn::detail::maxpooling_func<float, cudnnNanPropagation_t=0>, int=0, bool=0>, cudnnTensorStruct*, cudnnPoolingStruct, float, cudnnPoolingStruct, int, cudnn::reduced divisor, float)
          0.34% 153.41us      1  153.41us 153.41us 153.41us void mshadow::expr::MapPlanLargeKernel<mshadow::sv::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2)
          0.17% 75.359us      1  75.359us 75.359us 75.359us void mshadow::expr::SoftmaxKernel<int=8, float, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>>(mshadow::gpu, mshadow::Shape<int=2>, float, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>>(mshadow::gpu, int=2, float), float>>(mshadow::gpu, int=2, unsigned int)
          0.06% 28.032us      13  2.1560us 1.1840us 6.4960us void mshadow::expr::MapPlanKernel<mshadow::sv::saveto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>>(mshadow::gpu, mshadow::expr::Plan<mshadow::ScalarExp<float>, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2)
          0.05% 24.064us      2  12.032us 2.5600us 21.504us void mshadow::expr::MapPlanKernel<mshadow::sv::plusto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>>(mshadow::gpu, mshadow::expr::Broadcast1DExp<mshadow::Tensor<mshadow::gpu, int=1, float>, float, int=2, int=1, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2)
          0.03% 11.776us      10  1.1770us 960ns 1.9840us [CUDA memset]
          0.01% 5.7600us      1  5.7600us 5.7600us 5.7600us [CUDA memcpy DtoH]
          0.01% 4.7040us      1  4.7040us 4.7040us 4.7040us void mshadow::expr::MapPlanKernel<mshadow::sv::saveto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2)
```

Size 8 tile width with all restricted array variables but no unrolling:

```

Successfully installed mxnet
E: Running nvprof python m4.1.py
Loading fashion-mnist data... done
==279== NVPROF is profiling process 279, command: python m4.1.py
Loading model... done
New Inference
B is: 10000 M is: 6 H is: 48 W is: 48 C is: 1 K is: 5 Op Time: 0.011965
B is: 10000 M is: 16 H is: 22 W is: 22 C is: 6 K is: 5 Op Time: 0.028503
Correctness: 0.8397 Model: ece408
==279== Profiling application: python m4.1.py
==279== Profiling result:
      Type  Time(%)     Time    Calls      Avg      Min      Max  Name
GPU activities:  61.87%  40.373ms        2  20.187ms  11.910ms  28.463ms  mxnet::op::forward_kernel(float*, float const *, float const *, int, int, int, int)
                  26.00%  16.965ms       20  848.26us  1.0880us  16.417ms  [CUDA memcpy HtoD]
                  3.85%  2.5107ms        2  1.2554ms  20.928us  2.4898ms  volta_sgemm_32x128_tn
                  3.71%  2.4208ms        2  1.2104ms  734.27us  1.6863ms  void mshadow::cuda::MapPlanLargeKernel<mshadow::sv::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=4, float>, float, int=1>, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=4, int)
                  2.49%  1.6232ms        2  811.63us  21.888us  1.6614ms  void op_generic_tensor_kernel<int=2, float, float, float, int=256, cudnnGenericOp_t=7, cudnnNanPropagation_t=0, cudnnDimOrder_t=0, int=1>(cudnnTensorStruct, float*, cudnnTensorStruct, float const *, cudnnTensorStruct, float const *, float, float, float, dimArray, reducedDivisor)
                  1.62%  1.0546ms        1  1.0546ms  1.0546ms  1.0546ms  void cudnn::detail::pooling_fw_4d_kernel<float, float, cudnn::detail::maxpooling_func<float, cudnnNanPropagation_t=0>, int=0, bool=0>(cudnnTensorStruct, float const *, cudnn::detail::pooling_fw_4d_kernel<float, float, cudnn::detail::maxpooling_func<float, cudnnNanPropagation_t=0>, int=0, bool=0>, cudnnTensorStruct*, cudnnPoolingStruct, float, cudnnPoolingStruct, int, cudnn::reduced_divisor, float)
                  0.24%  157.76us        1  157.76us  157.76us  157.76us  void mshadow::cuda::MapPlanLargeKernel<mshadow::sv::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>, mshadow::expr::Plan<mshadow::expr::ScalarExp<float>, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2, int)
                  0.12%  75.072us        1  75.072us  75.072us  75.072us  void mshadow::cuda::SoftmaxKernel<int=8, float, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>>(mshadow::gpu, int=2, unsigned int)

```

Final Optimization: 2 loop unrolls, size 8 tile width, and all restricted array variables

```

E: Running nvprof python m4.1.py
Loading fashion-mnist data... done
==278== NVPROF is profiling process 278, command: python m4.1.py
Loading model... done
New Inference
B is: 10000 M is: 6 H is: 48 C is: 1 K is: 5 Op Time: 0.009558
B is: 10000 M is: 16 H is: 22 W is: 22 C is: 6 K is: 5 Op Time: 0.020876
Correctness: 0.8397 Model: ece408
==278== Profiling application: python m4.1.py
==278== Profiling result:
      Type  Time(%)     Time    Calls      Avg      Min      Max  Name
GPU activities:  55.33%  30.366ms        2  15.183ms  9.5152ms  20.851ms  mxnet::op::forward_kernel(float*, float const *, float const *, int, int, int, int)
                  30.74%  16.868ms       20  843.41us  1.0880us  16.462ms  [CUDA memcpy HtoD]
                  4.33%  2.3775ms        2  1.1888ms  728.09us  1.6494ms  void mshadow::cuda::MapPlanLargeKernel<mshadow::sv::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::gpu, int=4, float>, float>, mshadow::expr::Plan<mshadow::expr::BinaryMapExp<mshadow::expr::mul, mshadow::expr::ScalarExp<float>, mshadow::Tensor<mshadow::gpu, int=4, float>, float, int=1>, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=4, int)
                  4.26%  2.3065ms        2  1.1533ms  21.344us  2.2852ms  volta_sgemm_32x128_tn
                  2.94%  1.6127ms        2  806.36us  21.280us  1.5914ms  void op_generic_tensor_kernel<int=2, float, float, float, int=256, cudnnGenericOp_t=7, cudnnNanPropagation_t=0, int=1>(cudnnTensorStruct, float*, cudnnTensorStruct, float const *, float, float, float, float, dimArray, reducedDivisor)
                  1.91%  1.0457ms        1  1.0457ms  1.0457ms  1.0457ms  void cudnn::detail::pooling_fw_4d_kernel<float, float, cudnn::detail::maxpooling_func<float, cudnnTensorStruct, float const *, cudnn::detail::pooling_fw_4d_kernel<float, float, cudnn::detail::maxpooling_func<float, cudnnNanPropagation_t=0, int=nnTensorStruct*, cudnnPoolingStruct, float, cudnnPoolingStruct, int, cudnn::reduced_divisor, float>
                  0.28%  151.68us        1  151.68us  151.68us  151.68us  void mshadow::cuda::MapPlanLargeKernel<mshadow::sv::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::gpu, int=2, float>, float>, mshadow::expr::Plan<mshadow::expr::ScalarExp<float>, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2, int)
                  0.14%  75.199us        1  75.199us  75.199us  75.199us  void mshadow::cuda::SoftmaxKernel<int=8, float, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>>(mshadow::gpu, int=2, unsigned int)
                  0.05%  28.032us       13  2.1560us  1.2160us  6.4320us  void mshadow::cuda::MapPlanKernel<mshadow::sv::saveto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>, mshadow::expr::Plan<mshadow::expr::ScalarExp<float>, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2)
                  0.04%  23.872us       2  11.936us  2.5600us  21.312us  void mshadow::cuda::MapPlanKernel<mshadow::sv::saveto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>, mshadow::expr::Plan<mshadow::expr::ReduceWithAxisExp<mshadow::expr::reduce::maximum, mshadow::Tensor<mshadow::gpu, int=3, float>, float, int=3, bool=1, int=2>, ::gpu, unsigned int, mshadow::Shape<int=2>, int=2>
                  API calls:  42.72%  2.9988bs        22  135.95ms  14.317us  1.53094s  cudaStreamCreateWithFlags
                  33.84%  2.36879s       22  107.67ms  90.835us  2.36399s  cudaMemGetInfo
                  21.10%  1.47682s       18  82.046ms  826ns  392.91ms  cudaFree
                  0.88%  61.723ms        301ns  18.472ms  cudaFuncSetAttribute

```

Analysis:

Original op-time: .018, .072

Post optimization op-time: .009, .021

Unrolling seems to allow for more gpu compute time and decrease the op time. This is because we reuse the variable y[] so many times and unrolling allows us to be faster than the

usual iterative indexing process. Using restrict on the variables also improved performance. Furthermore, changing the tile width to 8 really boosted the block utilizations.

Optimization 2: Weight Matrix (kernel values) in Constant Memory

The outputs of the neural network are calculated by convoluting the input matrices with the a weight matrix, which remains constant throughout the entirety of the execution. In milestone 3, the weight matrix pointer is passed into the forward_kernel along with the input and output matrix pointers. Whenever the forward_kernel needs values from the weight matrix, it must access it through global memory, which is one of the slowest device memories in CUDA. This optimization places the weight matrix into constant memory instead, speeding up the access times for weight kernels.

```
Successfully installed m4.1
Running /usr/bin/time python m4.1.py 100
Loading fashion-mnist data... done
Loading model... done
New Inference
Op Time: 0.000223
Op Time: 0.000698
Correctness: 0.84 Model: ece408
4.24user 3.14system 0:04.15elapsed 177%CPU (0avgtext+0avgdata 2745916maxresident)k
0inputs+4624outputs (0major+615844minor)pagefaults 0swaps
Running /usr/bin/time python m4.1.py 1000
Loading fashion-mnist data... done
Loading model... done
New Inference
Op Time: 0.001876
Op Time: 0.006885
Correctness: 0.852 Model: ece408
4.05user 3.42system 0:03.97elapsed 188%CPU (0avgtext+0avgdata 2767488maxresident)k
0inputs+0outputs (0major+621402minor)pagefaults
0swaps
Running /usr/bin/time python m4.1.py
Loading fashion-mnist data... done
Loading model... done
New Inference
Op Time: 0.018420
Op Time: 0.071111
Correctness: 0.8397 Model: ece408
4.18user 3.24system 0:04.11elapsed 180%CPU (0avgtext+0avgdata 2840940maxresident)k
0inputs+0outputs
(0major+662092minor)pagefaults 0swaps
```

```

==548== NVPROF is profiling process 548, command: python m4.1.py
Loading model... done
New Inference
Op Time: 0.018448
Op Time: 0.071179
Correctness: 1.8397 Mdccl: ece408
==548== Profiling application: python m4.1.py
==548== Profiling result:
      Type  Time(%)    Time   Calls    Avg     Min     Max   Name
GPU activities:  78.74% 89.485ms   2 44.743ms 18.340ms 71.139ms mxnet::op::forward_kernel(float*, float const *, int, int, int, int, int, int)
               14.29% 16.240ms   20 812.02us 1.0560us 15.866ms [CUDA memcpy HtoD]
               2.20% 2.5023ms   2 1.2511ms 21.446us 2.4880ms volta_sgemm_32x128_tn
               2.13% 2.4222ms   2 1.2111ms 733.02us 1.6891ms void mshadow::op::mshadow::BinaryOpExp:mshadow::op::mul, mshadow::expr::ScalarExp<float>, mshadow::Tensor<mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=4>, int>
float>, float>, mshadow::expr::Plan:mshadow::expr::BinaryOpExp:mshadow::op::mul, mshadow::expr::ScalarExp<float>, mshadow::Tensor<mshadow::gpu, int=4, float>, float>, float, int=1>, float>)(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=4)
               1.43% 1.6237ms   2 811.83us 22.080ms void op_generic_tensor_kernel<int=2, float, float, float, float, int=256, cudnnGenericOp_t>, cudnnNanPropagation_t=0, cudmDimOrder_t=0, int=1>(cudnnTensorStruct, float*, cudnnTensorStruct, float const *, cudnnTensorStruct, float const *, float, float, float, dimArray, reducedDivisorArray)
               0.93% 1.0533ms   1 1.0533ms void cudnn::detail::pooling_fw_4d_kernel<float>, float, cudnn::detail::maxPooling_func<float, cudnnNanPropagation_t=0>, int=0, bool=0>(cudnnTensorStruct, float const *y, cudnn::detail::pooling_fw_4d_kernel<float>, float, cudnn::detail::maxPooling_func<float, cudnnNanPropagation_t=0>, int=0, bool=0>, cudnnTensorStruct, cudnnPoolingStruct, int=1, cudnn::detail::reduced_divisorArray)
               0.14% 160.54us   1 160.54us 160.54us void mshadow::op::mshadow::BinaryOpExp:mshadow::op::mul, mshadow::expr::Plan:mshadow::expr::ScalarExp<float>, float>)(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2, int>)
float>, float>, mshadow::expr::Plan:mshadow::expr::ScalarExp<float>, float>)(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2, int>)
               0.02% 27.936us   12 2.1480us 1.1840us 6.4960us void mshadow::op::mshadow::BinaryOpExp:mshadow::op::mul, mshadow::expr::Plan:mshadow::expr::ScalarExp<float>, float>)(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2)
               0.02% 23.584us   2 11.792us 2.5000us 21.024us void mshadow::op::mshadow::BinaryOpExp:mshadow::op::mul, mshadow::expr::Plan:mshadow::expr::ScalarExp<float>, float>)(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2)
               0.01% 11.744us   10 1.1740us 992s 2.1440us [CUDA memset]
               0.01% 6.3360us   2 3.1680us 3.1840us 3.2320us [CUDA memcpy DtoD]
               0.01% 5.9840us   1 5.9840us 5.9840us 5.9840us [CUDA memcpy DtoH]
               0.00% 4.8320us   1 4.8320us 4.8320us void mshadow::op::mshadow::BinaryOpExp:mshadow::op::mul, mshadow::expr::Plan:mshadow::expr::ScalarExp<float>, float>)(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2)
mshadow::expr::Plan:mshadow::expr::ReduceWithAxisExp:mshadow::red::maximum, mshadow::Tensor<mshadow::gpu, int=2, float>, float>, float, int=3, bool=1, int=2>, float>)(mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2)
API calls: 41.64% 2.93283s   22 133.31ms 14.058us 1.54611s cudaStreamCreateWithFlags
               32.72% 2.30446s   22 104.75ms 89.036us 2.29968s cudaMemGetInfo
               21.23%
1.49526s   18 83.070ms 845ns 401.39ms cudaFree
               1.49% 104.70ms   912 114.80us 296ns 50.792ms cudaFuncSetAttribute
               1.31% 91.913ms   6 15.322ms 4.6500us 71.143ms cudaDeviceSynchronize
               0.71% 50.185ms   216 232.34us 847ns 22.377ms cudaEventCreateWithFlags
               0.47% 33.294ms   9 3.6994ms 20.474us 16.051ms cudaMemCopy2DAsync
               0.14% 11.617ms   67 1.397ms 1.397us 1.397us cudaEventDestroy

```

Analysis (For the default dataset containing 10000 images):

Milestone 3 Optimes: 0.1847, 0.07168

Post Optimization Optimes: 0.1842, 0.071179

Optimization 3: Shared Memory Convolution

For forward convolution, we noticed that each block corresponds with a tile within one input image convolved with one mask, which means that for every thread in one block, it needs to perform the same global memory reads for k^2 elements in that mask, and it would be much faster if we divide the loading mask into shared memory to each thread and make the process parallel, which would give a shorter time for memory reads for each block since shared memory reads are faster than global memory reads.

For each block, since it performs the convolution for a tile of one input image, we use the tiling optimization methods for preloading the tile of input elements into the shared memory with all threads in parallel and make the access to the input element faster by reading from shared memory rather than from global memory.

Thus this optimization reduces the memory read times and speed up the kernel by utilizing shared memory properties.

-- analysis:

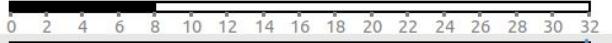
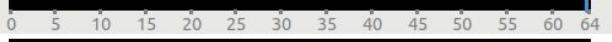
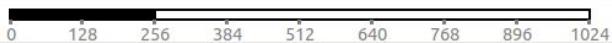
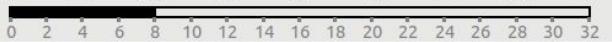
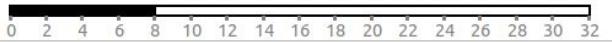
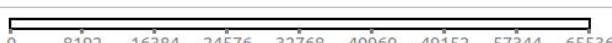
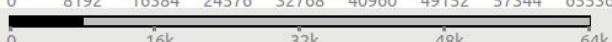
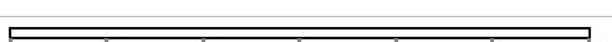
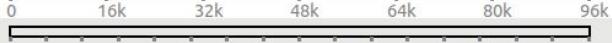
Milestone 3 without optimization : Op Time: 0.018446 - Op Time: 0.071832

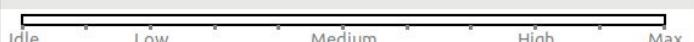
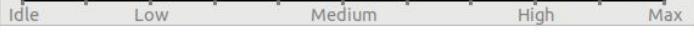
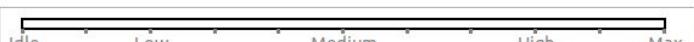
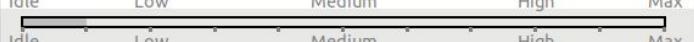
After optimization : Op Time: 0.006718 - Op Time: 0.034986

```
* Running pip2 install --user -e /mxnet/python
DEPRECATION: Python 2.7 will reach the end of its life on January 1st, 2020. Please upgrade your Python as Python 2.7 won't be maintained after that date. A future version of pip will drop support for Python 2.7.
Obtaining file:///mxnet/python
Requirement already satisfied: graphviz<0.9.0,>=0.8.1 in /root/.local/lib/python2.7/site-packages (from mxnet==1.3.1) (0.8.4)
Requirement already satisfied: numpy<1.16.0,>=1.1.8.2 in /root/.local/lib/python2.7/site-packages (from mxnet==1.3.1) (1.14.0)
Requirement already satisfied: requests<2.19.0,>=2.18.4 in /root/.local/lib/python2.7/site-packages (from mxnet==1.3.1) (2.18.4)
Requirement already satisfied: idna<2.7,>=2.5 in /root/.local/lib/python2.7/site-packages (from requests<2.19.0,>=2.18.4->mxnet==1.3.1) (2.6)
Requirement already satisfied: urllib3<1.23,>=1.21.1 in /root/.local/lib/python2.7/site-packages (from requests<2.19.0,>=2.18.4->mxnet==1.3.1) (1.22)
Requirement already satisfied: certifi>=2017.4.17 in /root/.local/lib/python2.7/site-packages (from requests<2.19.0,>=2.18.4->mxnet==1.3.1) (2018.11.29)
Requirement already satisfied: charsetet<3.1.0,>=3.0.2 in /root/.local/lib/python2.7/site-packages (from requests<2.19.0,>=2.18.4->mxnet==1.3.1) (3.0.4)
Installing collected packages: mxnet
  Running setup.py develop for mxnet
Successfully installed mxnet
* Running nvprof --o timeline_nvprof python m3.1.py
Loading fashion-mnist data... done
==280== NVPROF is profiling process 280, command: python m3.1.py
Loading model... done
New Inference
Op Time: 0.006787
Op Time: 0.034936
Correctness: 0.8397 Model: ece@08
==280== Generated result file: /build/timeline_nvprof
* Running nvprof --kernels "::forward1" --analysis-metrics -o forward1_analysis.nvprof python m3.1.py
Loading fashion-mnist data... done
==377== NVPROF is profiling process 377, command: python m3.1.py
Loading model... done
New Inference
==377== Some kernel(s) will be replayed on device 0 in order to collect all events/metrics.
Replaying kernel "mxnet::op::forward_kernel(float*, float const *, float const *, int, int, int, int, int, int)" (done)
Op Time: 3569581 events
Op Count: 0.00011
Correctness: 0.8397 Model: ece@08
==377== Generated result file: /build/forward1_analysis.nvprof
* Running nvprof --kernels "::forward2" --analysis-metrics -o forward2_analysis.nvprof python m3.1.py
Loading fashion-mnist data... done
==475== NVPROF is profiling process 475, command: python m3.1.py
Loading model... done
New Inference
Op Time: 0.006676
==475== Some kernel(s) will be replayed on device 0 in order to collect all events/metrics.
Replaying kernel "mxnet::op::forward_kernel(float*, float const *, float const *, int, int, int, int, int, int)" (done)
Op Time: 34.94178 events
Correctness: 0.8397 Model: ece@08
==475== Generated result file: /build/forward2_analysis.nvprof
* The build folder has been uploaded to http://s3.amazonaws.com/files.rai-project.com/userdata/build-5cb8edb27c204656045bc4c1.tar.gz. The data will be present for only a short duration of time.
* Server has ended your request.
wirelessrp-10-192-203-94:uiuc xinyiguosu
```

NVVP Analysis:

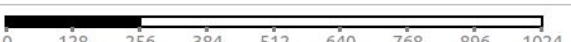
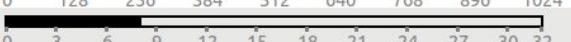
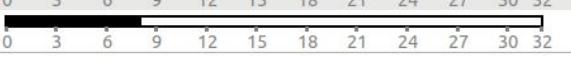
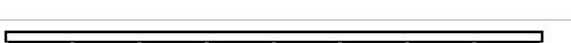
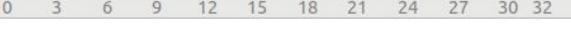
Milestone 3 (Unoptimized)

i Occupancy Is Not Limiting Kernel Performance				
The kernel's block size, register usage, and shared memory usage allow it to fully utilize all warps on the GPU.				
Variable	Achieved	Theoretical	Device Limit	Grid Size: [10000,6,9] (540000 blocks) Block Size: [16,16,1] (256 threads)
Occupancy Per SM				
Active Blocks	8	32		
Active Warps	63.46	64	64	
Active Threads		2048	2048	
Occupancy	99.2%	100%	100%	
Warp				
Threads/Block	256	1024		
Warps/Block	8	32		
Block Limit	8	32		
Registers				
Registers/Thread	28	65536		
Registers/Block	8192	65536		
Block Limit	8	32		
Shared Memory				
Shared Memory/Block	0	98304		
Block Limit	0	32		

i Memory Bandwidth And Utilization				
The following table shows the memory bandwidth used by this kernel for the various types of memory on the device. The table also shows the utilization of each memory type relative to the maximum throughput supported by the memory.				
	Transactions	Bandwidth	Utilization	
Shared Memory				
Shared Loads	0	0 B/s		
Shared Stores	0	0 B/s		
Shared Total	0	0 B/s		
L2 Cache				
Reads	27256512	52.557 GB/s		
Writes	480480498	926.477 GB/s		
Total	507737010	979.034 GB/s		
Unified Cache				
Local Loads	0	0 B/s		
Local Stores	0	0 B/s		
Global Loads	499476236	963.105 GB/s		
Global Stores	480480000	926.476 GB/s		
Texture Reads	128550696	991.501 GB/s		
Unified Total	1108506932	2,881.082 GB/s		
Device Memory				
Reads	33603390	64.795 GB/s		
Writes	21577031	41.605 GB/s		
Total	55180421	106.401 GB/s		
System Memory [PCIe configuration: Gen3 x8, 8 Gbit/s]				
Reads	0	0 B/s		
Writes	5	9.641 kB/s		

Constant Memory

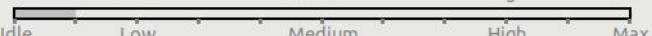
Results

i Occupancy Is Not Limiting Kernel Performance				
The kernel's block size, register usage, and shared memory usage allow it to fully utilize all warps on the GPU.				
Variable	Achieved	Theoretical	Device Limit	Grid Size: [10000,6,9] (540000 blocks) Block Size: [16,16,1] (256 threads)
Occupancy Per SM				
Active Blocks		8	32	
Active Warps	63.46	64	64	
Active Threads		2048	2048	
Occupancy	99.2%	100%	100%	
Warp				
Threads/Block		256	1024	
Warps/Block		8	32	
Block Limit		8	32	
Registers				
Registers/Thread		28	65536	
Registers/Block		8192	65536	
Block Limit		8	32	
Shared Memory				
Shared Memory/Block		0	98304	
Block Limit		0	32	

i Memory Bandwidth And Utilization

The following table shows the memory bandwidth used by this kernel for the various types of memory on the device. The table also shows the utilization of each memory type relative to the maximum throughput supported by the memory.

[More...](#)

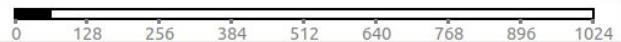
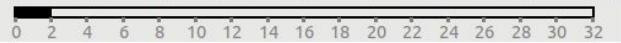
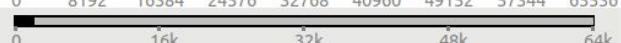
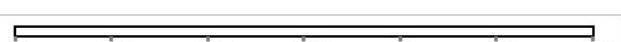
	Transactions	Bandwidth	Utilization
Shared Memory			
Shared Loads	0	0 B/s	
Shared Stores	0	0 B/s	
Shared Total	0	0 B/s	
L2 Cache			
Reads	27276447	51.967 GB/s	
Writes	480480102	915.404 GB/s	
Total	507756549	967.371 GB/s	
Unified Cache			
Local Loads	0	0 B/s	
Local Stores	0	0 B/s	
Global Loads	499475103	951.593 GB/s	
Global Stores	480480000	915.404 GB/s	
Texture Reads	128563020	979.746 GB/s	
Unified Total	1108518123	2,846.744 GB/s	
Device Memory			
Reads	33602900	64.02 GB/s	
Writes	21601630	41.155 GB/s	
Total	55204530	105.175 GB/s	
System Memory [PCIe configuration: Gen3 x8, 8 Gbit/s]			
Reads	0	0 B/s	
Writes	5	9.525 kB/s	

Unrolling

Occupancy Is Not Limiting Kernel Performance

The kernel's block size, register usage, and shared memory usage allow it to fully utilize all warps on the GPU.

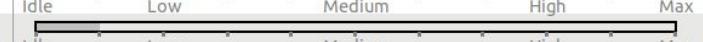
[More...](#)

Variable	Achieved	Theoretical	Device Limit	Grid Size: [10000,6,36] (2160000 blocks) Block Size: [8,8,1] (64 threads)
Occupancy Per SM				
Active Blocks		32	32	
Active Warps	42.15	64	64	
Active Threads		2048	2048	
Occupancy	65.9%	100%	100%	
Warp Metrics				
Threads/Block		64	1024	
Warps/Block		2	32	
Block Limit		32	32	
Registers				
Registers/Thread		32	65536	
Registers/Block		2048	65536	
Block Limit		32	32	
Shared Memory				
Shared Memory/Block		0	98304	
Block Limit		0	32	

Memory Bandwidth And Utilization

The following table shows the memory bandwidth used by this kernel for the various types of memory on the device. The table also shows the utilization of each memory type relative to the maximum throughput supported by the memory.

[More...](#)

	Transactions	Bandwidth	Utilization
Shared Memory			
Shared Loads	0	0 B/s	
Shared Stores	0	0 B/s	
Shared Total	0	0 B/s	
L2 Cache			
Reads	46994088	136.784 GB/s	
Writes	22440016	65.316 GB/s	
Total	69434104	202.1 GB/s	
Unified Cache			
Local Loads	0	0 B/s	
Local Stores	0	0 B/s	
Global Loads	757622591	2,205.192 GB/s	
Global Stores	22440000	65.316 GB/s	
Texture Reads	341576365	3,976.869 GB/s	
Unified Total	1121638956	6,247.376 GB/s	
Device Memory			
Reads	66173944	192.611 GB/s	
Writes	22486769	65.452 GB/s	
Total	88660713	258.062 GB/s	
System Memory [PCIe configuration: Gen3 x16, 8 Gbit/s]			
Reads	0	0 B/s	
Writes	5	14.553 kB/s	

Shared Convolution

i Memory Bandwidth And Utilization

The following table shows the memory bandwidth used by this kernel for the various types of memory on the device. The table also shows the utilization of each memory type relative to the maximum throughput supported by the memory.

[More..](#)

	Transactions	Bandwidth	Utilization				
Shared Memory							
Shared Loads	330191009	6,995.915 GB/s					
Shared Stores	15161650	321.237 GB/s					
Shared Total	345352659	7,317.152 GB/s	Idle	Low	Medium	High	Max
L2 Cache							
Reads	26882967	142.396 GB/s					
Writes	18480016	97.886 GB/s					
Total	45362983	240.282 GB/s	Idle	Low	Medium	High	Max
Unified Cache							
Local Loads	0	0 B/s					
Local Stores	0	0 B/s					
Global Loads	26891524	142.441 GB/s					
Global Stores	18480000	97.886 GB/s					
Texture Reads	235145723	4,982.145 GB/s					
Unified Total	280517247	5,222.472 GB/s	Idle	Low	Medium	High	Max
Device Memory							
Reads	33594060	177.943 GB/s					
Writes	18517670	98.086 GB/s					
Total	52111730	276.029 GB/s	Idle	Low	Medium	High	Max
System Memory [PCIe configuration: Gen3 x8, 8 Gbit/s]							
Reads	0	0 B/s	Idle	Low	Medium	High	Max
Writes	5	26.484 kB/s	Idle	Low	Medium	High	Max
Occupancy Per SM							
Active Blocks	6	32	0	2	4	6	8
Active Warps	46.65	64	0	5	10	15	20
Active Threads	1536	2048	0	256	512	768	1024
Occupancy	72.9%	75%	0%	15%	30%	45%	60%
			75%	90%	100%		
Warp							
Threads/Block	256	1024	0	128	256	384	512
Warps/Block	8	32	0	2	4	6	8
Block Limit	8	32	0	2	4	6	8
Registers							
Registers/Thread	39	65536	0	8192	16384	24576	32768
Registers/Block	10240	65536	0	16k	32k	48k	64k
Block Limit	6	32	0	2	4	6	8
Shared Memory							
Shared Memory/Block	1700	98304	0	16k	32k	48k	64k
Block Limit	54	32	0	2	4	6	8

Based on the NVVP analysis shown above, we can see that constant memory has similar statistics as the milestone 3 kernel, and offers very little performance boosts. Meanwhile,

unrolling provides a lot more, as can be seen by the amount of active threads and warps, which is maxed out. Finally, shared memory convolution is also an effective enhancement -- since shared memory one of the fastest memories, using it for computation greatly enhances performance.

MILESTONE 5:

Optimization 4: Variable Sweeping

Sometimes, the best way to find the best optimization is by trying every combination of tuning based off of many possible areas of improvement or realizing the effect of certain variables on the optime. We started off by doing a variable sweep and tile difference for shared convolutions because some of the optimes were the lowest we'd had. When we tried introducing new kernels sometimes, we might get a slower run time. However, if we adjusted the tile widths and switched out some unroll numbers, we could sometimes halve the time it took. Additionally, we restricted the array pointers because read-only memory is faster than normal global memory.

	Tile width 4	Tile width 8	Tile width 16	Tilewidth 24	Tilewidth 32
Shared memory conv	Op Time: 0.021967 Op Time: 0.042524	Op Time: 0.008990 Op Time: 0.020861	Op Time: 0.005638 Op Time: 0.028150	Op Time: 0.007272 Op Time: 0.020014	Op Time: 0.015271 Op Time: 0.036846
Kernel Fusion	Op Time: 0.030086 Op Time: 0.050335	Op Time: 0.007826 Op Time: 0.015902	Op Time: 0.010030 Op Time: 0.015930	Op Time: 0.015342 Op Time: 0.014376	Op Time: 0.024249 Op Time: 0.019810

nvprof Analysis: Sweeping helps us find the optimal tile widths for resource utilization

Optimization 5: Unroll Matrix Multiply Fusion (w/Constant Kernel)

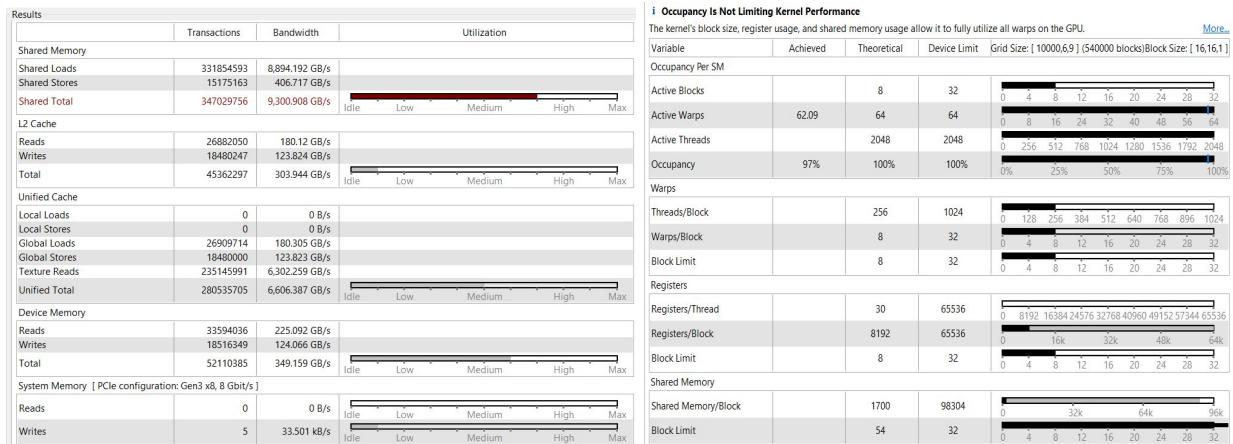
We expected Optimization 4 to help much more than it actually did which lead us to fusing the kernels. Due to the fact that calling multiple kernels for one pass involves many transfers of memory, the original unroll/Matrix multiply method will not be able to perform at it's optimized time. By fusing the two kernels together, we eliminate the transfers of memory to 2 separate kernels. Additionally, by adding the filter to Constant memory, we hope to optimize the memory transfers when we load the filter values.

```

Loading fashion-mnist data... done
==373== NVPROF is profiling process 373, command: python m4.1.py
Loading model... done
New Inference
Op Time: 0.017461
Op Time: 0.011492
Correctness: 8397 Model: ece408
==373== Profiling application: python m4.1.py
==373== Profiling result:
      Type  Time(%)    Calls    Avg    Min    Max   Name
GPU activities: 53.56%  28.879ms     2  14.440ms  11.465ms  17.415ms mxnet::op::forward_kernel<float*, float const *, float const *, int, int, int, int, int, int>
31.73%  17.105ms    20  85.27us  1.0560us  16.571ms [CUDA memcpy HtoD]
4.65%  3.5057ms     2  1.2528ms  22.015us  2.4837ms volta_sgemm_32x128_tn
4.51%  2.4320ms     2  1.2160ms  735.61us  1.6964ms void mshadow::cuda::MapPlanLargeKernel<mshadow::sv::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=4, float>, float, int=1>, float>(mshadow::gpu, unsigned int, mshadow::Shape<int>*, int=4, int)
3.03%  1.6347ms     2  81.37us  22.048us  1.6127ms void op_genericOp<tensor_kernelt<int>, float, float, float, float, int=256, cudnnGenericOp_t=7, cudnnNanPropagation_t=0, cudnnDimOrder_t=0, int=1>(cudnnTensorStruct, float, cudnnTensorStruct, float const *, cudnnTensorStruct, float const *, float, float, float, float, dimArray, reducedDivisorArray)
1.96%  1.0546ms     1  1.0546ms  1.0546ms  1.0546ms void cudnn::detail::pooling_fw_4d_kernel<float, float, cudnn::detail::maxpooling_func<float, cudnnNanPropagation_t=0>, int=0, bool=0>(cudnnTensorStruct, float const *, cudnn::detail::pooling_fw_4d_kernel<float, float, cudnn::detail::maxpooling_func<float, cudnnNanPropagation_t=0>, int=0, bool=0>, cudnnTensorStruct, cudnnPoolingStruct, float, int, cudnn::detail::reduced divisor, float)
0.28%  153.38us     1  153.38us  153.38us  153.38us void mshadow::cuda::MapPlanLargeKernel<mshadow::sv::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>(mshadow::gpu, int=2, float, float, mshadow::expr::Plan<mshadow::expr::ScalarExp<float>, float>(mshadow::gpu, unsigned int, mshadow::Shape<int>*, int=2, int)
0.14%  75.360us     1  75.360us  75.360us  75.360us void mshadow::cuda::SoftMaxKernel<int=8, float, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>(mshadow::gpu, int=2, float, float, mshadow::expr::Plan<mshadow::expr::ScalarExp<float>, float>(mshadow::gpu, unsigned int, mshadow::Shape<int>*, int=2))
0.05%  27.872us     13  2.1440us  1.1520us  6.5280us void mshadow::cuda::MapPlanKernel<mshadow::sv::saveto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>(mshadow::gpu, int=2, unsigned int)
0.04%  24.096us     2  12.048us  2.4960us  21.600us void mshadow::cuda::MapPlanKernel<mshadow::sv::plusto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>(mshadow::gpu, int=2, float, float, mshadow::expr::Plan<mshadow::expr::Broadcast1DExp<mshadow::Tensor<mshadow::gpu, int=1, float>, float, int=2, int=1>, float>(mshadow::gpu, unsigned int, mshadow::Shape<int>*, int=2))
0.02%  12.000us     10  1.2000us  992ns  1.98400us [CUDA memset]
0.01%  7.7760us     1  7.7760us  7.7760us  7.7760us [CUDA memcpy DtoH]
0.01%  5.1840us     1  5.1840us  5.1840us  5.1840us void mshadow::cuda::MapPlanKernel<mshadow::sv::saveto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>(mshadow::gpu, unsigned int, mshadow::Shape<int>*, int=2)
API calls: 42.12% 3.13130s    22  142.33ms  14.260us  1.61165s cudaStreamCreateWithFlags
32.06%  3.28323s    22  108.33ms  90.308us  2.37832s cudaMemAllocInfo
23.38%  1.73172s    18  96.207ms  824ns  448.14ms cudaFree
1.01%  74.746ms    912  81.958us  309ns  45.990us cudaFuncSetAttribute
0.46%  34.396ms     9  3.8218ms  28.493us  16.606ms cudaMemcpy2DAsync
0.42%  31.339ms     6  5.2231ms  4.5980us  17.419ms cudaDeviceSynchronize
0.18%  13.572ms    216  62.831us  875ns  12.383ms cudaEventCreateWithFlags
0.13%  9.3435ms     66  141.57us  6.0200us  1.9787ms cudaMalloc

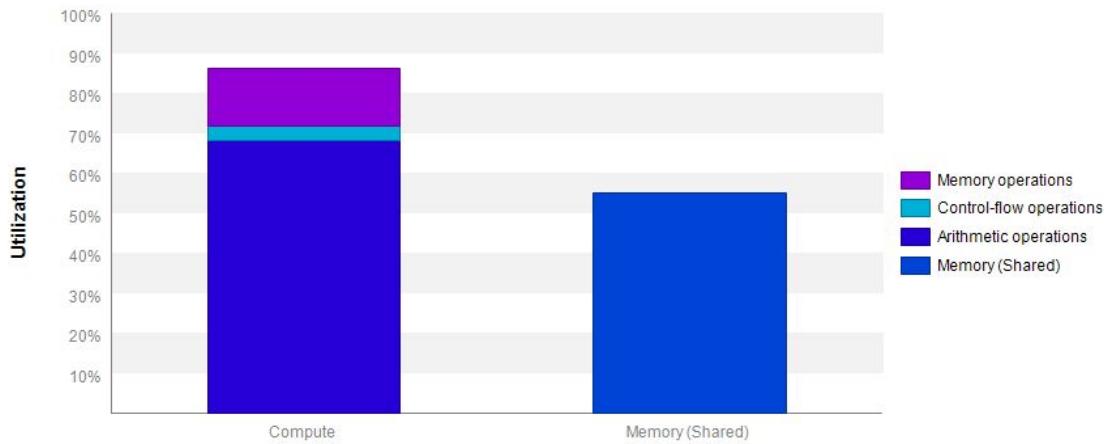
```

NVVP:



i Kernel Performance Is Bound By Compute

For device "TITAN V" the kernel's memory utilization is significantly lower than its compute utilization. These utilization levels indicate that the performance of the kernel is most likely being limited by computation on the SMs.



NVVP Analysis: The kernel seems to be limited by the amount of math done per call but that is because we are mapping each thread to one output. The math behind that is quite large. Maybe if we mapped every thread to an input, the runtime would be faster.

Optimization 6: Multiple Kernel implementations for different layer sizes

As we were doing our optimizations, we learned that different layer sizes have faster times with certain kernels. So far we had 3 different promising kernels so we decided to tune the tile widths and combinations of the 3 kernels that we saw some good op times in. During this process, inputs of each layer size was paired with a different kernel to sweep every combination of 2 kernels there is. On top of this, we adjusted the tile widths for the optimal performance of the combinations.

We know that the fused matrix multiply w/constant memory did the best with the second pass because it can handle the larger amount of layers with smaller inputs to take advantage of the ability to coalesce memory. On the other hand, the shared memory convolution worked best with the first pass as the input sizes were larger and there were less channels. We wanted to know if there was an optimal combination of kernels for each pass. To test this, we ran the following combinations:

Kernel type->first pass (6) -second pass (16) V	Fused Matrix Multiply w/constant memory	Restrict and unroll Convolution	Shared memory convolution
Fused Matrix Multiply w/constant memory:	.0204, .0156 (TILEWIDTHS 16)	.0057, .0159 (TILE WIDTHS 16,24)	.0051, .0159 (TILE WIDTHS(24,16)

Restrict and unroll convolution	.01, .023 (TILEWIDTHS 16, 8)	.0066, .021(TILEWIDTHS 16,8)	.0058, .017 (TILEWIDTHS 16,24)
Shared memory convolution	.01, .18 (TILEWIDTHS 24,16)	.0056, .016 (TILEWIDTHS 16,24)	0.006718, 0.034986 (TILEWIDTH 16)

Restrict and unroll convolution (1st) Fused Matrix Multiply w/constant mem (2nd)

Shared mem convolution (1st), Fused Matrix Multiply w/constant mem (2nd)

```

Running nvprof python m4.1.py
Loading fashion-mnist data... done
==373== NVPROF is profiling application 373, command: python m4.1.py
Loading model... done
New Inference
Op Time: 0.005703
Op Time: 0.015948
Correctness: 0.8397 Model: ece408
==373== Profiling application: python m4.1.py
==373== Profiling result:
      Type  Time(%)   Time    Calls   Avg     Min     Max   Name
GPU activities:  36.05% 16.600ms    20 829.99us 1.0880us 16.051ms [CUDA memcpy HtoD]
34.52% 15.896ms      1 15.896ms 15.896ms mxnet::op::forward kernel2<float*, float*, int, int, int, int, int, int>
12.24% 5.6374ms      1 5.6374ms 5.6374ms 5.6374ms mxnet::op::forward kernel<float*, float const *, float const *, int, int, int, int, int>
5.44% 2.5066ms       2 1.2533ms 21.504ms 2.4851ms volta sgemm 32x128 tn
5.25% 2.4151ms       2 1.2976ms 730.81us 1.6843ms void mshadow::cuda::MapPlanLargeKernel<mshadow::sv::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, unsigned int, mshadow::Shape<int>, int>, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int>, int=4, int)
3.52% 1.6217ms       2 810.83us 22.112us 1.5995ms void op::generic_tensor<kernelInt2>, float, float, float, float, int=256, cudnnGenericOp_t=t, cudnnNanPropagation_t=0, cudnnDimOrder_t=0, int=1>(cudnnTensorStruct, float const *, cudnnTensorStruct, float const *, float, float, float, dimArray, reducedDivisorArray)
2.88% 1.0514ms       1 1.0514ms 1.0514ms void cudnn::detail::pooling_fw_4d_kernel<float, float, float, cudnn::detail::maxPooling_func<float, cudnnNanPropagation_t=>, int=0, bool=>(cudnnTensorStruct, cudnnPoolingStruct, float, cudnnPoolingStruct, int, cuunn::reduced_divisor, float)
0.34% 158.72us       1 158.72us 158.72us 158.72us void mshadow::cuda::MapPlanLargeKernel<mshadow::sv::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, unsigned int, mshadow::Shape<int>, int>, float>>(mshadow::gpu, unsigned int, mshadow::Shape<int>, int=2, int)
0.16% 75.167us       1 75.167us 75.167us 75.167us void mshadow::cuda::SoftmaxKernel<int8, float, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>, float>(mshadow::gpu, int=2, unsigned int)
0.06% 27.936us       13 2.1480us 1.1840us 6.5280us void mshadow::cuda::MapPlanKernel<mshadow::sv::saveto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>, float>(mshadow::gpu, unsigned int, mshadow::Shape<int>, int=2)
0.05% 24.096us       2 12.048us 2.4960us 21.600us void mshadow::cuda::MapPlanKernel<mshadow::sv::plusto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>, float>(mshadow::gpu, unsigned int, mshadow::Shape<int>, int=2)
0.03% 12.255us       10 1.2250us 992ns 1.9840us [CUDA memset]
0.02% 7.8720us       1 7.8720us 7.8720us 7.8720us [CUDA memcpy DtoH]
0.01% 5.0560us       2 2.5280us 2.4640us 2.5920us [CUDA memcpy DtoB]
0.01% 4.8000us       1 4.8000us 4.8000us 4.8000us void mshadow::cuda::MapPlanKernel<mshadow::sv::saveto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, unsigned int, mshadow::Shape<int>, int>, float>, float>(mshadow::gpu, unsigned int, mshadow::Shape<int>, int=2)
API calls: 43.563s 3.08525s      22 140.24ms
14.131us 1.58313s  cudaStreamCreateWithFlags
31.89% 2.25827s      22 102.65ms 95.054us 2.25328s cudaMemGetInfo
22.19% 1.57184s      18 87.324ms 880ns 427.44ms cudaFree
0.57% 40.468ms       216 187.35us 894ns 38.828ms cudaEventCreateWithFlags

```

Shared mem convolution (1st) with Restrict tuning (2nd)

```

Op Time: 0.005568
Op Time: 0.016178
Correctness: 0.8397 Model: ece408
Running numpy python m4.1.py
Loading fashion-mnist data... done
==373== NVPROF is profiling process 373, command: python m4.1.py
Loading model... done
New Inference
Op Time: 0.005681
Op Time: 0.016222
Correctness: 0.8397 Model: ece408
==373== Profiling application: python m4.1.py
==373== Profiling result:
      Type    Time(%)  Time   Calls    Avg     Min     Max   Name
GPU activities: 47.29% 21.841ms           2  10.920ms  5.6453ms 16.195ms mxnet::op::forward_kernel(float*, float const *, float const *, int, int, int, int, int, int)
35.56% 16.422ms          20  821.09us  1.0560us 16.001ms [CUDA memcpy HtoD]
5.43% 2.5071ms          2  1.2536ms  21.376us 2.4858ms volta_sgmem_32x128_tn
5.25% 2.4243ms          1  1.0560ms  21.376us 1.6884ms void msshadow::cuda::MapPlanLargeKernel<mshadow::sv::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::gpu, int=4, float>, float, int=1>, float>(mshadow::gpu, unsigned int, mshadow::Shape<int>*, int=4, int)
3.52% 1.6248ms          2  812.39us  22.400us 1.6024ms void op::generic_tensor_kernel<int=2, float, float, float, float, int=256, cudnnGeneticop_t=7, cudnnNanPropagation_t=0, cudnnDimOrder_t=0, int=1>(cudnnTensorStruct, float*, cudnnTensorStruct, float const *, cudnnTensorStruct, float const *, float, float, float, dimArray, reducedDivisorArray)
2.79% 1.0568ms          1  1.0560ms  1.0560ms 1.0560ms void cudnn::detail::pooling_fw_4d_kernel<float, float, float, float>(cudnn::detail::maxpooling_func<float, cudnnNanPropagation_t=0, int=0>, int=0, bool=0)(cudnnTensorStruct, float const *, cudnn::detail::pooling_fw_4d_kernel<float, float, float, float>, cudnn::detail::maxpooling_func<float, cudnnNanPropagation_t=0>, cudnnTensorStruct, cudnnPoolingStruct, float, cudnn::detail::maxpooling_func<float, cudnnNanPropagation_t=0>, int=0, bool=0)
0.34% 157.44us          1  157.44us  157.44us 157.44us void msshadow::cuda::MapPlanLargeKernel<mshadow::sv::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>, float>(mshadow::gpu, unsigned int, mshadow::Shape<int>*, int=2, int)
0.16% 75.775us          1  75.775us  75.775us 75.775us void msshadow::cuda::SoftmaxKernel<int=8, float, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>, float, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>(mshadow::gpu, int=2, unsigned int)
0.06% 27.930us          13  2.1480us  2.1480us 6.5280us void msshadow::cuda::MapPlanKernel<mshadow::sv::saveto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>(mshadow::gpu, unsigned int, mshadow::Shape<int>*, int=2, int)
0.05% 23.903us          2  11.951us  2.5280us 21.375us void msshadow::cuda::MapPlanKernel<mshadow::sv::plusto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>, mshadow::expr::BroadcastExpXpmshadow::Tensor<mshadow::gpu, int=1, float>, float, int=2, int>(float>(mshadow::gpu, unsigned int, mshadow::Shape<int>*, int=2, int))
0.03% 11.648us          10  1.1640us  992ns 1.6000ms [CUDA memset]
0.01% 5.7600us          1  5.7600us  5.7600us 5.7600us [CUDA memcpy DtoH]
0.01% 4.7350us          1  4.7350us  4.7350us 4.7350us void msshadow::cuda::MapPlanKernel<mshadow::sv::saveto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>(mshadow::gpu, unsigned int, mshadow::Shape<int>*, int=2, int)
API calls: 43.18% 3.16502ms          22  143.87ms  13.857us 1.66533s cudaStreamCreateWithFlags
31.74% 2.32677s          22  105.76ms  99.164ms 2.32179s cudaMemGetInfo
21.65% 1.58674s          18  88.152ms  827ns 427.97ms cudaFree
1.93% 141.84ms          912  155.520s  303ns 103.78ms cudaFuncSetAttribute
0.46% 33.700ms          9  3.7445ms  22.003us 16.150ms cudaLemcyP2ASync

```

Restrict Unroll tuning Convolution(1st) Shared mem conv (2nd)

Fused Matrix Multiply w/constant mem (1st) Restrict and unroll convolution (2nd)

```

Loading model... done
New Inference
Op Time: 0.010055
Op Time: 0.017861
Correctness: 0.8397 Model: ece408
# Running nvprof python m4_1.py
Loading Fashion-mnist data... done
==373== NVPROF is profiling process 373, command: python m4_1.py
Loading model... done
New Inference
Op Time: 0.010296
Op Time: 0.017905
Correctness: 0.8397 Model: ece408
==373== Profiling application: python m4_1.py
==373== Profiling result:
      Type  Time(%)    Time   Calls    Avg     Min     Max  Name
GPU activities:  34.00%  17.861ms           1  17.861ms  17.861ms  mxnet::op::forward kernel(float*, float const *, float const *, int, int, int, int, int, int)
               31.20%  16.392ms          20  819.61us  1.1200us  16.000ms  [CUDA memcpy HtoD]
               19.51%  10.223ms          1  10.223ms  10.223ms  mxnet::op::forward kernel2(float*, float*, int, int, int, int, int, int)
               4.79%  2.5079ms          2  1.2540ms  21.408us  2.4865ms  volta_sgmemm_32x128_tn
               4.62%  2.4220ms          2  1.2110ms  741.27us  1.6800ms  void mshadow::cuda::MapPlanLargeKernel<mshadow::sv::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::gpu, int=4, float>, float, int=1>, float>>(<mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=4>, int)
               3.10%  1.6265ms          2  813.23us  22.592us  1.6839ms  void op_generic_tensor_kernel<int=2, float, float, float, float, int=256, cudnnGenericOp_t=7, cudnnPropagationToken_t=0>, int=0, bool=0>(<mshadow::gpu, float>, cudnnTensorStruct, float const *, cudnnTensorStruct, float const *, float, float, float, float, dim3Array, reducedDivisorArray)
               2.00%  1.0476ms          1  1.0476ms  1.0476ms  void cudnn::detail::pooling_fw_4d_kernel<float, float, cudnn::detail::maxpooling_func<float, cudnnPropagationToken_t=0>, cudnnTensorStruct*, cudnnPoolingStruct, float, cudnnPoolingStruct, int, cudnn::detail::maxpooling_func<float, cudnnPropagationToken_t=0>, int=0, bool=0>
               0.31%  160.77us          1  160.77us  160.77us  160.77us  void mshadow::cuda::MapPlanLargeKernel<mshadow::sv::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::Ten
sor<mshadow::gpu, int=2>, float>, float>, mshadow::expr::Plan<mshadow::sv::saveto, float>, float>>(<mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2>, int)
               0.14%  75.231us          1  75.231us  75.231us  75.231us  void mshadow::cuda::SoftmaxKernel<int=8, float>, float, mshadow::expr::Plan<mshadow::gpu, int=2, float>, float>, mshadow::expr::Plan<mshadow::gpu, int=2, float>, float>>(<mshadow::gpu, int=2, float>, float)
               0.05%  27.903us          13  2.1460us  1.1840us  6.4960us  void mshadow::cuda::MapPlanKernel<mshadow::sv::saveto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>, mshadow::expr::Plan<mshadow::gpu, int=2, float>, float>>(<mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2>)
               0.05%  23.872us          2  11.936us  2.5600us  21.312us  void mshadow::cuda::MapPlanKernel<mshadow::sv::plusto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>, mshadow::expr::ReduceWithAxisExp<mshadow::red::maximum, mshadow::Tensor<mshadow::gpu, int=3, float>, float, int=3, bool=1, int=2>, float>>(<mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2>)
               0.02%  11.744us          10  1.1740us  992ns  2.1440us  [CUDA memset]
               0.01%  5.8880us          1  5.8880us  5.8880us  5.8880us  [CUDA memcpy DtoH]
               0.01%  2.5800us          2  2.6400us  2.3680us  2.9120us  [CUDA memcpy DtoB]
               0.01%  4.9600us          1  4.9600us  4.9600us  4.9600us  void mshadow::cuda::ReduceWithAxisExp<mshadow::red::maximum, mshadow::Tensor<mshadow::gpu, int=3, float>, float, int=3, bool=1, int=2>, float>>(<mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2>)
API calls: 42.09%  3.02711s          22  137.60ms

```

Fused Matrix Multiply(1st) Convolution with restrict and unrolling tuning(2nd)

```

# Running python m4_1.py 10000
Loading Fashion-mnist data... done
Loading model... done
New Inference
Op Time: 0.010054
Op Time: 0.023503
Correctness: 0.8397 Model: ece408
# Running nvprof python m4_1.py
Loading Fashion-mnist data... done
==373== NVPROF is profiling process 373, command: python m4_1.py
Loading model... done
New Inference
Op Time: 0.010312
Op Time: 0.023672
Correctness: 0.8397 Model: ece408
==373== Profiling application: python m4_1.py
==373== Profiling result:
      Type  Time(%)    Time   Calls    Avg     Min     Max  Name
GPU activities:  40.00%  23.629ms           1  23.629ms  23.629ms  mxnet::op::forward_kernel(float*, float const *, float const *, int, int, int, int, int, int)
               29.29%  17.298ms          20  864.92us  1.0880us  16.913ms  [CUDA memcpy HtoD]
               17.31%  10.225ms          1  10.225ms  10.225ms  mxnet::op::forward kernel2(float*, float*, int, int, int, int, int, int)
               4.23%  2.4964ms          2  1.2480ms  31.568us  3.4740ms  volta_sgmemm_32x128_tn
               4.11%  2.4267ms          2  1.2133ms  734.04us  1.6927ms  void mshadow::cuda::MapPlanLargeKernel<mshadow::sv::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=4, float>, float, int=1>, float>>(<mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=4>, int)
               2.75%  1.6239ms          2  811.93us  22.144us  1.6617ms  void op_generic_tensor_kernel<int=2, float, float, float, float, int=256, cudnnGenericOp_t=7, cudnnPropagationToken_t=0>, int=0, bool=0>(<mshadow::gpu, float>, cudnnTensorStruct, float const *, cudnnTensorStruct, float const *, float, float, float, float, dim3Array, reducedDivisorArray)
               0, cudnnDimOrder_t=0, int=1>(<mshadow::gpu, float>, cudnnTensorStruct, float const *, cudnnTensorStruct, float const *, float, float, float, float, dim3Array, reducedDivisorArray)
               0, cudnnTensorStruct*, cudnnPoolingStruct, float, cudnnPoolingStruct, int, cudnn::detail::maxpooling_func<float, cudnnPropagationToken_t=0>, int=0, bool=0>
               0.27%  158.72us          1  158.72us  158.72us  158.72us  void mshadow::cuda::MapPlanLargeKernel<mshadow::sv::saveto, int=8, int=1024, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>, float>>(<mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2>)
               0.13%  75.456us          1  75.456us  75.456us  75.456us  void mshadow::cuda::SoftmaxKernel<int=8, float>, float, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>>(<mshadow::gpu, int=2, float>, float)
               0.05%  28.192us          13  2.1600us  1.1840us  6.5920us  void mshadow::cuda::MapPlanKernel<mshadow::sv::saveto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=2, float>, float>>(<mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2>)
               0.04%  24.064us          2  12.032us  2.5600us  21.5040us  void mshadow::cuda::MapPlanKernel<mshadow::sv::plusto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=1, float>, float, int=2, int=1>, float>>(<mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2>)
               0.02%  11.328us          10  1.1320us  992ns  1.6000us  [CUDA memset]
               0.01%  5.9200us          1  5.9200us  5.9200us  5.9200us  [CUDA memcpy DtoH]
               0.01%  5.2800us          2  2.6400us  2.4960us  2.7840us  [CUDA memcpy DtoB]
               0.01%  4.6400us          1  4.6400us  4.6400us  4.6400us  void mshadow::cuda::MapPlanKernel<mshadow::sv::saveto, int=8, mshadow::expr::Plan<mshadow::Tensor<mshadow::gpu, int=1, float>, float>>(<mshadow::gpu, unsigned int, mshadow::Shape<int=2>, int=2>)

```

NVVP

Kernel fusion pass

Results		Transactions	Bandwidth	Utilization
Shared Memory				
Shared Loads	331854593	8,894.192 GB/s		
Shared Stores	15175163	406.717 GB/s		
Shared Total	347029756	9,300.908 GB/s	Idle Low Medium High Max	
L2 Cache				
Reads	26882050	180.12 GB/s		
Writes	18480247	123.824 GB/s		
Total	45362297	305.944 GB/s	Idle Low Medium High Max	
Unified Cache				
Local Loads	0	0 B/s		
Local Stores	0	0 B/s		
Global Loads	26909714	180.305 GB/s		
Global Stores	18480000	123.823 GB/s		
Texture Reads	235145991	6,302.259 GB/s		
Unified Total	280535705	6,606.387 GB/s	Idle Low Medium High Max	
Device Memory				
Reads	33594036	225.092 GB/s		
Writes	18516349	124.066 GB/s		
Total	52110385	349.159 GB/s	Idle Low Medium High Max	
System Memory [PCIe configuration: Gen3 x8, 8 Gbit/s]				
Reads	0	0 B/s	Idle Low Medium High Max	
Writes	5	33.501 kB/s	Idle Low Medium High Max	

i Occupancy Is Not Limiting Kernel Performance

The kernel's block size, register usage, and shared memory usage allow it to fully utilize all warps on the GPU.

[More...](#)

Variable	Achieved	Theoretical	Device Limit	Grid Size: [10000,6,9] (540000 blocks) Block Size: [16,16,1]
Occupancy Per SM				
Active Blocks		8	32	
Active Warps	62.09	64	64	
Active Threads	2048	2048	2048	
Occupancy	97%	100%	100%	0% 25% 50% 75% 100%
Warp Statistics				
Threads/Block		256	1024	
Warp/Block		8	32	
Block Limit		8	32	
Registers				
Registers/Thread		30	65536	
Registers/Block		8192	65536	
Block Limit		8	32	
Shared Memory				
Shared Memory/Block		1700	98304	
Block Limit		54	32	

i Memory Bandwidth And Utilization

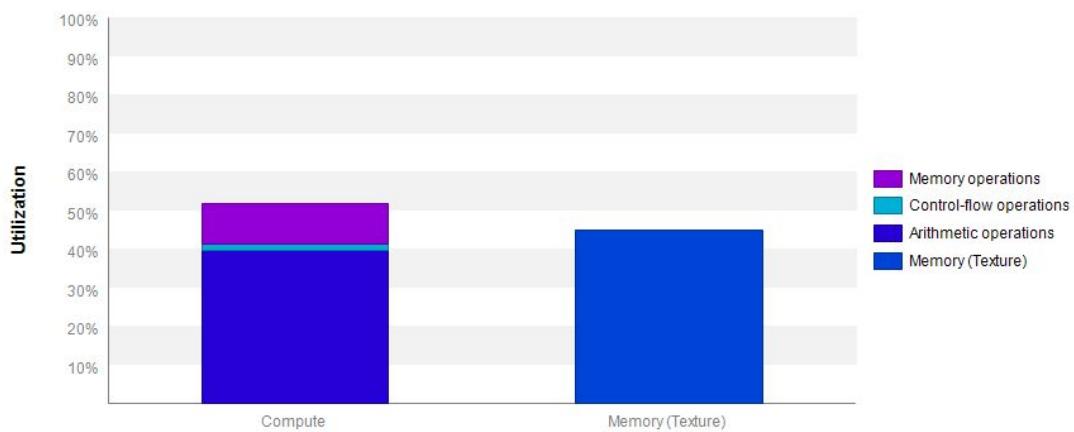
The following table shows the memory bandwidth used by this kernel for the various types of memory on the device. The table also shows the utilization of each memory type relative to the maximum throughput supported by the memory.

[More...](#)

Results		Transactions	Bandwidth	Utilization
Shared Memory				
Shared Loads	330191009	6,995.915 GB/s		
Shared Stores	15161650	321.237 GB/s		
Shared Total	345352659	7,317.152 GB/s	Idle Low Medium High Max	
L2 Cache				
Reads	26882967	142.396 GB/s		
Writes	184800016	97.886 GB/s		
Total	45362983	240.282 GB/s	Idle Low Medium High Max	
Unified Cache				
Local Loads	0	0 B/s		
Local Stores	0	0 B/s		
Global Loads	26891524	142.441 GB/s		
Global Stores	18480000	97.886 GB/s		
Texture Reads	235145723	4,982.145 GB/s		
Unified Total	280517247	5,222.472 GB/s	Idle Low Medium High Max	
Device Memory				
Reads	33594060	177.943 GB/s		
Writes	18517670	98.086 GB/s		
Total	52111730	276.029 GB/s	Idle Low Medium High Max	
System Memory [PCIe configuration: Gen3 x8, 8 Gbit/s]				
Reads	0	0 B/s	Idle Low Medium High Max	
Writes	5	26.484 kB/s	Idle Low Medium High Max	

i Kernel Performance Is Bound By Instruction And Memory Latency

This kernel exhibits low compute throughput and memory bandwidth utilization relative to the peak performance of "TITAN V". These utilization levels indicate that the performance of the kernel is most likely limited by the latency of arithmetic or memory operations. Achieved compute throughput and/or memory bandwidth below 60% of peak typically indicates latency issues.



Analysis: This combination actually uses less shared memory than the fusion kernel alone but runs faster perhaps because of the smaller amount of arithmetic operations. There seems to be a memory latency issue which might be a result of the redundant loads from X and the many loads we have in general even if they are coalesced/avoid bank conflicts in most cases. The NVVPs are very similar to that of the single kernels. Both the block limits have been reached and the use of shared memory is pretty high for both of the kernels. When you compare this to the NVVPs for the respective kernel calls of the original however, we see that the pass and kernel combinations had better overall usage of resources and the fastest runtime.

Additional insights:

Something we've learned about CUDA is that you cannot dynamically allocate shared memory unless you declare it as extern. When we were fusing the matrix multiply kernel, we tried to change the tile width based on an if statement before the kernel call. After further research, we found that when the program runs, cuda only supports a single dynamically declared memory allocation for each block. If we need more, we would have needed to use pointers within that allocation

When we originally tried Unrolled matrix multiply: Convolution operation essentially performs dot products between the filters and local regions of the input. By unrolling the forward pass into a large matrix, we can run through all the calculations in bulk, leaving most of the work in the inner for loops of the channels. Additionally, the memory is all stored together giving potential for memory coalescing within that kernel. The downside to matrix multiply is that there are a lot of redundant loads from X. This implies that if there are many channels and layers, these could outweigh the fact that all the memory for this kernel is bunched together.

We expected the unroll to really help with the memory usage because we had the opportunity to coalesce. However, we realized the the transfers of memory were actually incredibly heavy from this. On top of this, the redundant loads from X also hit our runtime hard.

Work division (All optimizations): We made sure everyone did 2 optimizations each and split up the work evenly across the first 3 milestones. Only Kevin was able to install cuda toolkit so he got everyone the nvvp results. Overall, the work was split up very evenly.

Works cited (All optimizations):

H, Albert. *Devtalk.nvidia.com*. N.p., 23 Dec. 2016. Web. 25 Apr. 2019.

Lee, Daren, et al. "CUDA Optimization Strategies for Compute- and Memory-Bound Neuroimaging Algorithms." *Computer Methods and Programs in Biomedicine*, U.S. National Library of Medicine, June 2012, www.ncbi.nlm.nih.gov/pmc/articles/PMC3262956/. Web. 30 April 2019.

Shafkat, Irhum. "Intuitively Understanding Convolutions for Deep Learning." *Towards Data Science*, Towards Data Science, 1 June 2018, towardsdatascience.com/intuitively-understanding-convolutions-for-deep-learning-1f6f42faee1.

Woolley, Cliff. "GPU Optimization Fundamentals." *Bluewaters Ncsa*, 2013, bluewaters.ncsa.illinois.edu. Web. 21 April 2019.

"An Effective Method for Better Power Efficiency on Multithreaded GPU." *Kernel Fusion*. IEEE Computer Society, 2010. Web. 17 April 2019.

"Loop Unrolling." *GeeksforGeeks*. N.p., 19 Feb. 2018. Web. 3 May 2019.

http://homepages.math.uic.edu/~jan/mcs572/memory_coalescing.pdf

<http://cseweb.ucsd.edu/classes/wi12/cse260-a/Lectures/Lec08.pdf>