

Name: Ruiqi Li
NetID: Ruiqili4
Section: AL1

ECE 408/CS483 Milestone 3 Report

0. List Op Times, whole program execution time, and accuracy for batch size of 100, 1k, and 10k images from your basic forward convolution kernel in milestone 2. This will act as your baseline this milestone.

Batch Size	Op Time 1	Op Time 2	Total Execution Time	Accuracy
100	0.338023 ms	1.21262 ms	0m1.141s	0.86
1000	3.08233 ms	11.9752 ms	0m9.566s	0.886
10000	30.435 ms	120.535 ms	1m38.909s	0.8714

1. Optimization 1: Tiled shared memory convolution

- a. Which optimization did you choose to implement and why did you choose that optimization technique.

I chose Tiled shared memory convolution. Because this is the most intuitive way to accelerate computation.

- b. How does the optimization work? Did you think the optimization would increase performance of the forward convolution? Why? Does the optimization synergize with any of your previous optimizations?

The optimization works very well. I think the optimization would increase performance of the forward convolution, because it replaces the global memory traffic with shared memory, therefore increases the execution speed. The execution speed of the kernel will be limited by the global memory bandwidth. But if we use shared memory tiling to temporarily store data in each block, the computation time can be significantly reduced. The optimization doesn't synergize with previous optimizations.

- c. List the Op Times, whole program execution time, and accuracy for batch size of 100, 1k, and 10k images using this optimization (including any previous optimizations also used).

Batch Size	Op Time 1	Op Time 2	Total Execution Time	Accuracy
100	0.226206 ms	1.01411 ms	0m1.601s	0.86
1000	1.99117 ms	9.94768 ms	0m9.905s	0.886
10000	19.2631 ms	98.5901 ms	1m34.265s	0.8714

- d. Was implementing this optimization successful in improving performance? Why or why not? Include profiling results from *nsys* and *Nsight-Compute* to justify your answer, directly comparing to your baseline (or the previous optimization this one is built off of).

Yes, it successfully improved performance. The execution time decreased with this optimization.

*Since I only modify the memory usage of each tile, namely, replace global memory access with shared memory, the memory copy and allocation times are roughly identical. Here is the comparison between *nsys* results of this optimization and that of the baseline.*

Time(%)	Total Time	Calls	Average	Minimum	Maximum	Name
81.2	1546321026	10	154632102.6	19306	598762736	cudaMemcpy
9.6	181871156	8	22733894.5	70365	176797902	cudaMalloc
7.9	151134746	8	18891843.2	1424	120623695	cudaDeviceSynchronize
1.1	21709731	6	3618288.5	23201	21577572	cudaLaunchKernel
0.1	2851037	8	356379.6	67387	1030270	cudaFree

Generating CUDA Kernel Statistics...

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

Time(%)	Total Time	Instances	Average	Minimum	Maximum	Name
100.0	151096614	2	75548307.0	30474316	120622298	conv_forward_kernel
0.0	2784	2	1392.0	1376	1408	do_not_remove_this_kernel
0.0	2656	2	1328.0	1312	1344	prefn_marker_kernel

CUDA Memory Operation Statistics (nanoseconds)

Time(%)	Total Time	Operations	Average	Minimum	Maximum	Name
66.9	1029784465	2	514892232.5	431950981	597833484	[CUDA memcpy DtoH]
33.1	510150122	8	63768765.2	1184	247256301	[CUDA memcpy HtoD]

CUDA Memory Operation Statistics (KiB)

Total	Operations	Average	Minimum	Maximum	Name
2261419.0	8	282677.0	0.004	1000000.0	[CUDA memcpy HtoD]
1722500.0	2	861250.0	722500.000	1000000.0	[CUDA memcpy DtoH]

Baseline ↑

Time(%)	Total Time	Calls	Average	Minimum	Maximum	Name
81.8	1615060821	10	161506082.1	22541	653661406	cudaMemcpy
10.9	215090157	8	26886269.6	95378	209962132	cudaMalloc
6.0	118675057	8	14834382.1	1139	99182254	cudaDeviceSynchronize
1.2	22807227	6	3801204.5	26244	22660964	cudaLaunchKernel
0.1	2794053	8	349256.6	67514	960841	cudaFree

Generating CUDA Kernel Statistics...

Generating CUDA Memory Operation Statistics...

CUDA Kernel Statistics (nanoseconds)

Time(%)	Total Time	Instances	Average	Minimum	Maximum	Name
100.0	118644791	2	59322395.5	19464881	99179910	conv_forward_kernel
0.0	2912	2	1456.0	1408	1504	do_not_remove_this_kernel
0.0	2688	2	1344.0	1344	1344	prefn_marker_kernel

CUDA Memory Operation Statistics (nanoseconds)

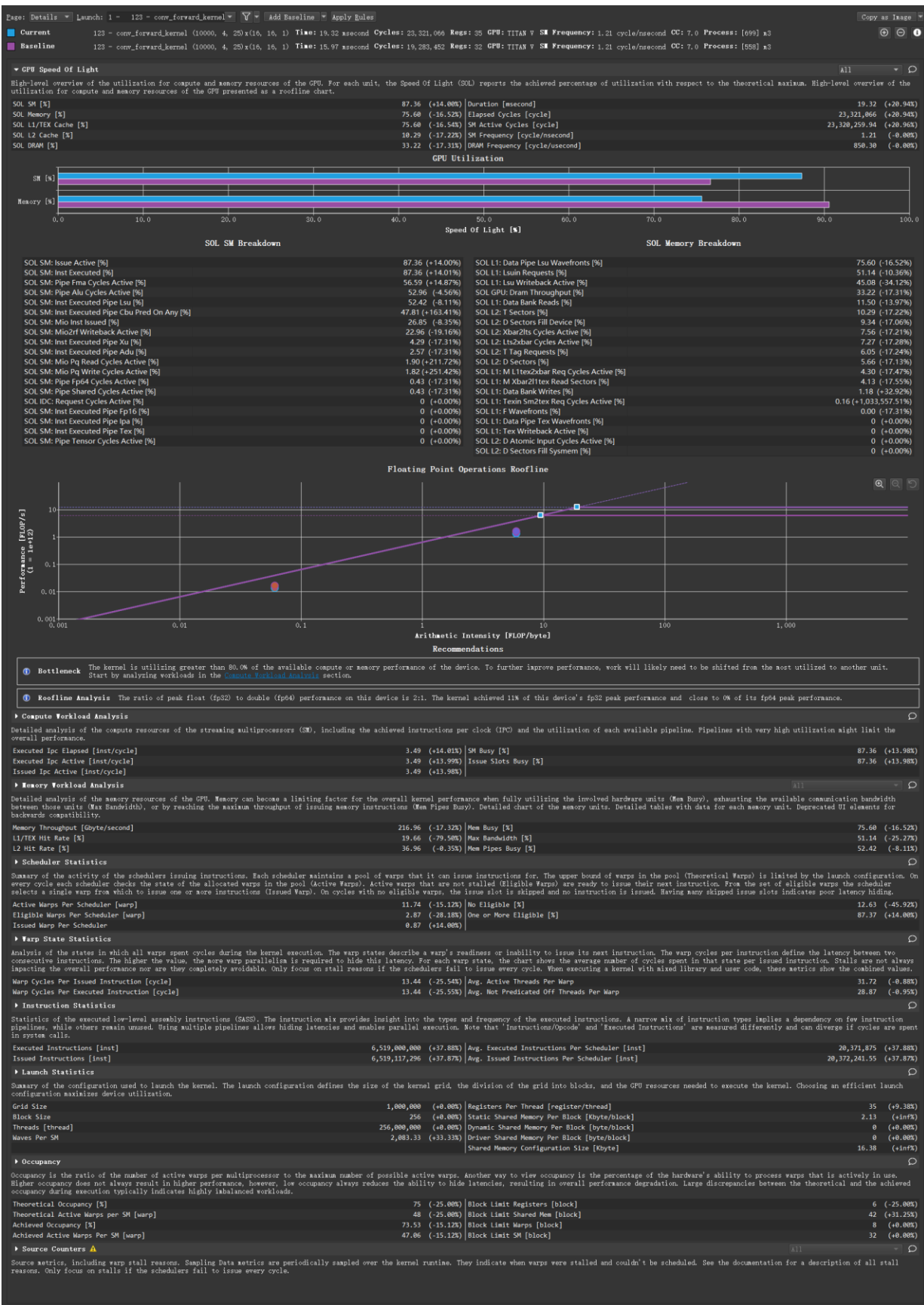
Time(%)	Total Time	Operations	Average	Minimum	Maximum	Name
68.2	1093096865	2	546548432.5	440407757	652689108	[CUDA memcpy DtoH]
31.8	509079359	8	63634919.9	1184	249950584	[CUDA memcpy HtoD]

CUDA Memory Operation Statistics (KiB)

Total	Operations	Average	Minimum	Maximum	Name
2261419.0	8	282677.0	0.004	1000000.0	[CUDA memcpy HtoD]
1722500.0	2	861250.0	722500.000	1000000.0	[CUDA memcpy DtoH]

Shared Memory Optimization ↑

Apparently, they are similar. However, if we take a look at results from Nsight-Compute:



It's obvious that global memory throughput has been reduced about 17.32%, and correspondingly SM throughput has been increased. This change can certainly improve GPU performance, since the global memory bandwidth is limited.

- e. What references did you use when implementing this technique?

Mainly Textbook chapter 16.

2. Optimization 2: Streams

- a. Which optimization did you choose to implement and why did you choose that optimization technique.

I chose Streams to overlap computation with data transfer. I choose this because as far as I know, this is the most potential way to reduce execution time significantly, since memory copy operations takes a lot of time.

- b. How does the optimization work? Did you think the optimization would increase performance of the forward convolution? Why? Does the optimization synergize with any of your previous optimizations?

I leveraged `cudaStream_t` classes to create `nStreams` cuda streams. This hyper parameter `nStreams` can be tuned in practice. Then I use `cudaMemcpyAsync()` to copy memory contents asynchronously and merge it with `nStreams` parts of kernel execution. All these merges can be done with cuda streams. Beyond that, I also overlap kernel execution with output memory copy asynchronously. All the code reference can be found on [How to Overlap Data Transfers in CUDA C/C++ | NVIDIA Developer Blog](#). The optimization doesn't synergize with previous optimizations.

- c. List the Op Times, whole program execution time, and accuracy for batch size of 100, 1k, and 10k images using this optimization (including any previous optimizations also used).

Batch Size	Op Time 1	Op Time 2	Total Execution Time	Accuracy
100	0.001077 ms	0.001201 ms	0m1.154s	0.86
1000	0.001233 ms	0.001186 ms	0m10.337s	0.886
10000	0.003007 ms	0.001301 ms	1m45.721s	0.8714

- d. Was implementing this optimization successful in improving performance? Why or why not? Include profiling results from `nsys` and `Nsight-Compute` to justify your answer, directly comparing to your baseline (or the previous optimization this one is built off of).

Yes, it improved performance a lot. Let's look at the results from nsys.

CUDA API Statistics (nanoseconds)

Time(%)	Total Time	Calls	Average	Minimum	Maximum	Name
85.1	1168393870	66	17702937.4	8895	43442016	cudaMemcpyAsync
13.6	186729720	8	23341215.0	77242	183299057	cudaMalloc
0.7	10286849	36	285745.8	4147	9926407	cudaLaunchKernel
0.3	3667023	2	1833511.5	14492	3652531	cudaMemcpy
0.3	3520463	8	440057.9	65264	1079808	cudaFree
0.0	169053	32	5282.9	1316	75246	cudaStreamCreate
0.0	158677	8	19834.6	1893	115068	cudaDeviceSynchronize
0.0	124683	32	3896.3	1799	20814	cudaStreamDestroy

Generating CUDA Kernel Statistics...

Generating CUDA Memory Operation Statistics...

CUDA Kernel Statistics (nanoseconds)

Time(%)	Total Time	Instances	Average	Minimum	Maximum	Name
100.0	78703226	32	2459475.8	1018081	3908840	conv_forward_kernel
0.0	2976	2	1488.0	1440	1536	do_not_remove_this_kernel
0.0	2784	2	1392.0	1376	1408	prefn_marker_kernel

CUDA Memory Operation Statistics (nanoseconds)

Time(%)	Total Time	Operations	Average	Minimum	Maximum	Name
93.3	1061203476	32	33162608.6	27101751	42764135	[CUDA memcpy DtoH]
6.7	76237228	36	2117700.8	1536	2530404	[CUDA memcpy HtoD]

CUDA Memory Operation Statistics (KiB)

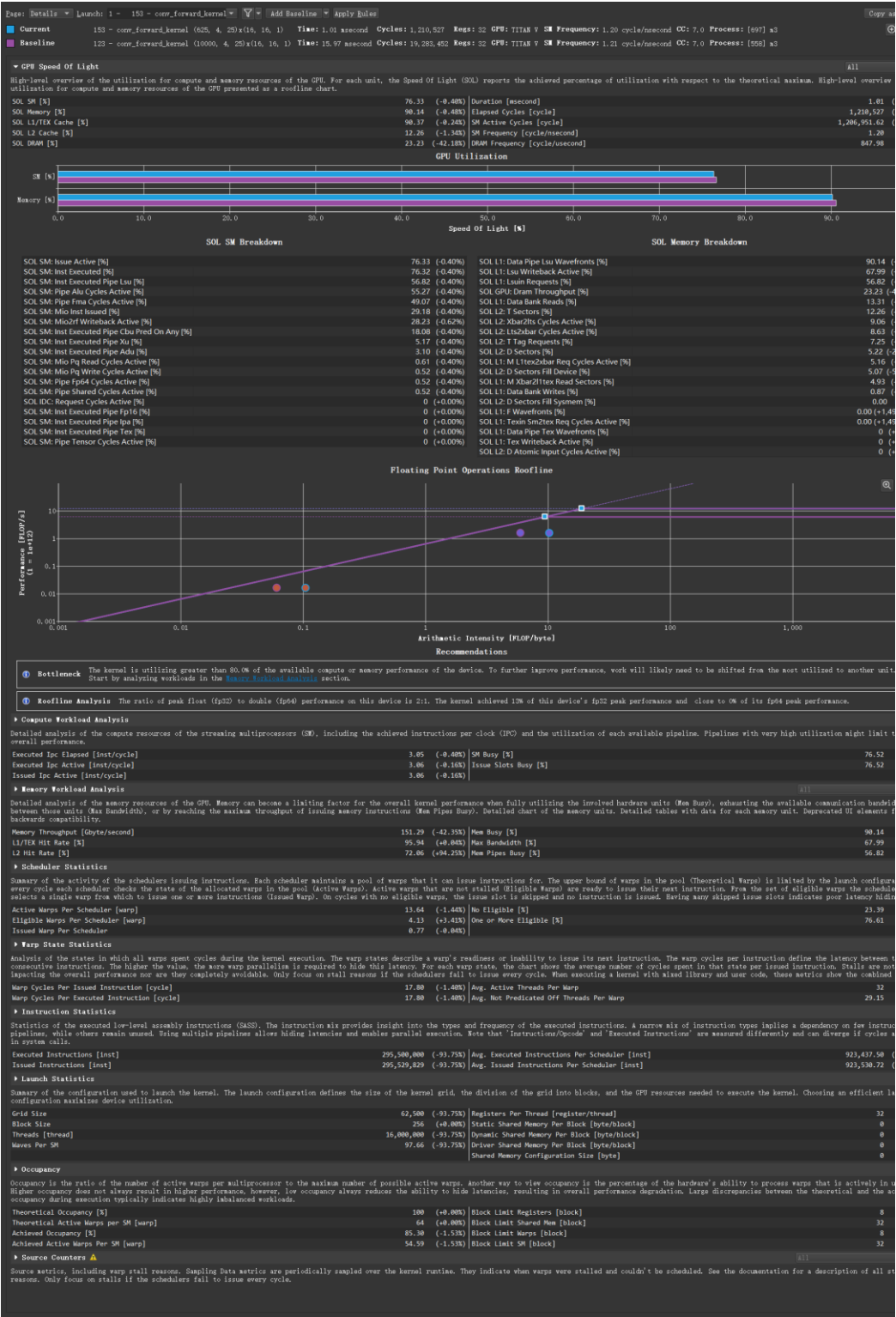
Total	Operations	Average	Minimum	Maximum	Name
1722500.0	32	53828.0	45156.250	62500.0	[CUDA memcpy DtoH]
538919.0	36	14970.0	0.004	18056.0	[CUDA memcpy HtoD]

Results that used Stream ↑

Compared with baseline, cudaMemcpyAsync() dominates the overall time consuming since we replace the normal cudaMemcpy() with it. Also, it's not hard to find out that the time used by CUDA memcpy HtoD is reduced dramatically while DtoH remains

roughly the same. It's largely because we overlap the output with input and kernel execution.

This is the result of Nsight-Compute of forward kernel call one of the Streams of the optimization:



Since the kernel part remains still, the SM memory condition remains the same. However, the global memory throughput has been reduced dramatically by about 42.35%. This is because the bulk memory copy operations are now distributed over time by cuda streams. From the very top of this picture, we can also find that one of cuda streams only process a small part of data set, not 10000 images anymore. The scheduling of data transportation and kernel execution indeed improves performance.

Nonetheless, I think the OP time measure in the output is not accurate. From the document, Op Time is time between the last cudaMemcpy call before the first kernel call and the first cudaMemcpy after the last kernel call. However, it measures this time by simply measure the execution time of new-forward.cu->conv_forward_gpu() function. So, it assumes that we all implement this function properly. Unfortunately, to implement CUDA Stream boosted optimization, we have to write the memory copy code and the kernel calls together to merge them. This makes new-forward.cu->conv_forward_gpu() useless and I simply returns once the program enters it. This will make the OP time measure meaningless. Yet, comparing Layer times still proves that Streams method improved performance.

- e. What references did you use when implementing this technique?

[CUDA stream · CUDA Little Book \(gitbooks.io\)](https://gitbooks.io/CUDA-stream-CUDA-Little-Book)
[How to Overlap Data Transfers in CUDA C/C++ | NVIDIA Developer Blog](#)
[GPU Pro Tip: CUDA 7 Streams Simplify Concurrency | NVIDIA Developer Blog](#)

3. Optimization 3: Fixed point (FP16) arithmetic

- a. Which optimization did you choose to implement and why did you choose that optimization technique.

I chose 16-bit floating point (FP16) data arithmetic. 16-bit “half-precision” floating point types are useful in applications that can process larger datasets or gain performance by choosing to store and operate on lower-precision data.

- b. How does the optimization work? Did you think the optimization would increase performance of the forward convolution? Why? Does the optimization synergize with any of your previous optimizations?

The optimization works well. I only leveraged FP16 data arithmetic but not storage. So based on the baseline method, I only cast the input X values and kernel values into half-precision representation, then do the computation, then cast them back to 32-bit floating point values and stored them back. The optimization doesn't synergize with previous optimizations.

- c. List the Op Times, whole program execution time, and accuracy for batch size of 100, 1k, and 10k images using this optimization (including any previous optimizations also used).

Batch Size	Op Time 1	Op Time 2	Total Execution Time	Accuracy
100	0.239754 ms	0.850697 ms	0m1.169s	0.86
1000	2.10516 ms	8.22387 ms	0m10.344s	0.887
10000	20.7124 ms	81.9974 ms	1m37.546s	0.8716

- d. Was implementing this optimization successful in improving performance? Why or why not? Include profiling results from *nsys* and *Nsight-Compute* to justify your answer, directly comparing to your baseline (or the previous optimization this one is built off of).

The nsys results of this optimization is similar to the baseline.

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

Time(%)	Total Time	Calls	Average	Minimum	Maximum	Name
78.5	1575841657	10	157584165.7	15022	603009293	cudaMemcpy
9.8	196049067	8	24506133.4	74468	193653203	cudaMalloc
5.8	115602581	8	14450322.6	69378	66851279	cudaFree
5.2	103575808	8	12946976.0	1552	81840786	cudaDeviceSynchronize
0.8	15617410	6	2602901.7	23104	15456344	cudaLaunchKernel

Generating CUDA Kernel Statistics...

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

Time(%)	Total Time	Instances	Average	Minimum	Maximum	Name
100.0	102474157	2	51237078.5	20643336	81830821	conv_forward_kernel
0.0	2784	2	1392.0	1376	1408	do_not_remove_this_kernel
0.0	2687	2	1343.5	1311	1376	prefn_marker_kernel

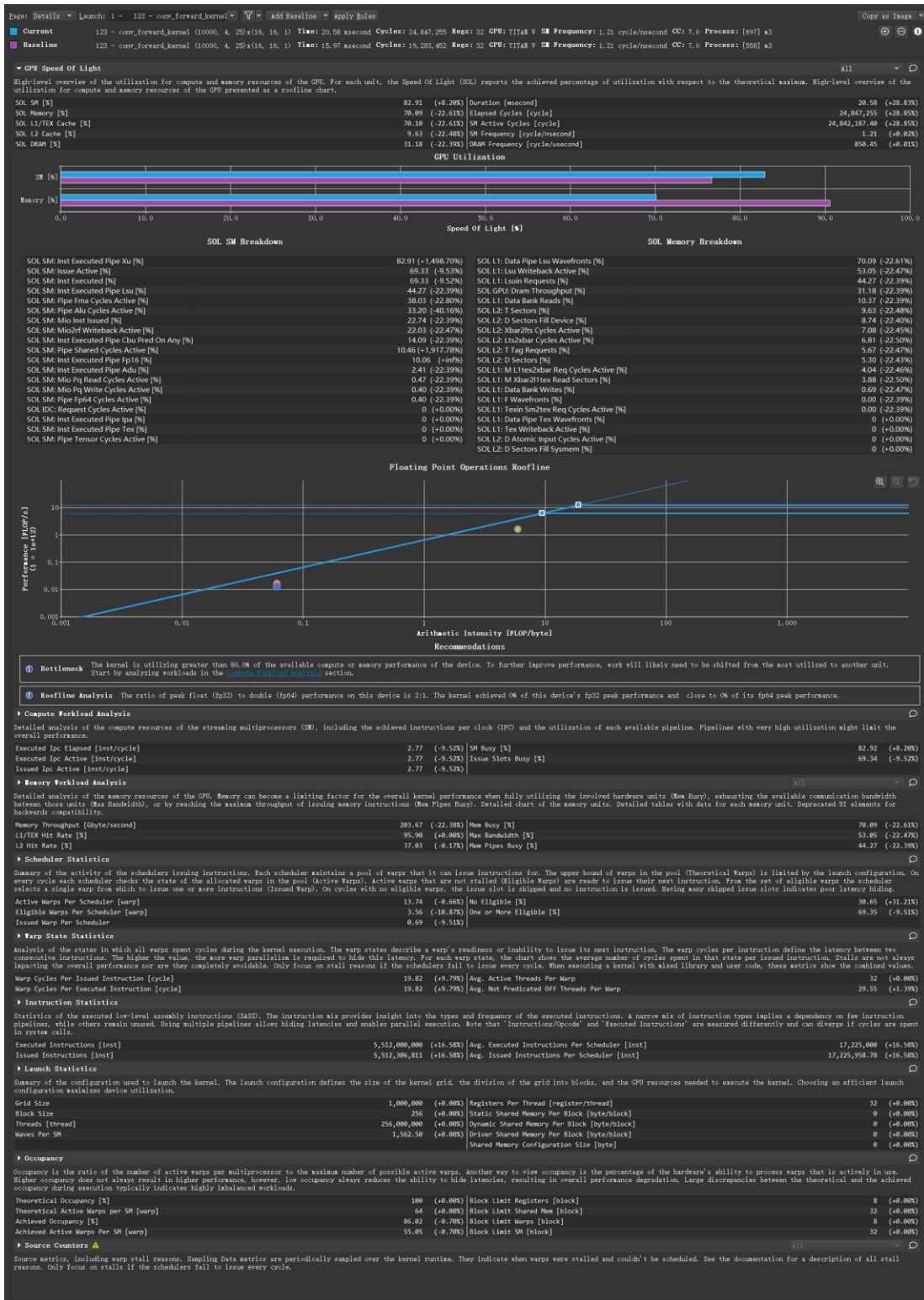
CUDA Memory Operation Statistics (nanoseconds)

Time(%)	Total Time	Operations	Average	Minimum	Maximum	Name
66.7	1047597210	2	523798605.0	445479694	602117516	[CUDA memcpy DtoH]
33.3	522153882	8	65269235.2	1184	250909320	[CUDA memcpy HtoD]

CUDA Memory Operation Statistics (KiB)

Total	Operations	Average	Minimum	Maximum	Name
2261419.0	8	282677.0	0.004	1000000.0	[CUDA memcpy HtoD]
1722500.0	2	861250.0	722500.000	1000000.0	[CUDA memcpy DtoH]

This is the results comparison of Nsight-Compute ↓.



The biggest difference is that the global memory throughput is decreased about 22.38%. In fact, I'm not 100% sure why this would happen, maybe casting from 16-bit floating point to 32-bit can save some space or time when assign the results back to output pointer. Also, in the middle it says "The kernel achieved 0% of this device's fp32

peak performance and close to 0% of its fp64 peak performance”, which means we were indeed 100% utilizing FP16 arithmetic. Beyond that, there is a measure called “Compute Workload Analysis” showing that the computation overhead is actually reduced by a little.

In conclusion, this optimization surely improves computation performance.

- e. What references did you use when implementing this technique?

[New Features in CUDA 7.5 | NVIDIA Developer Blog](#)

[Support FP16 in CUDA · Issue #699 · apache/tvm \(github.com\)](#)

[c++ - How to allocate cuda half-precision arrays correctly? - Stack Overflow](#)