

Name: Joseph Thomas
 NetID: jpt7
 Section: AL2

ECE 408/CS483 Milestone 3 Report

- List Op Times, whole program execution time, and accuracy for batch size of 100, 1k, and 5k images from your basic forward convolution kernel in milestone 2. This will act as your baseline this milestone. Note: **Do not** use batch size of 10k when you profile in `--queue rai_amd64_exclusive`. We have limited resources, so any tasks longer than 3 minutes will be killed. Your baseline M2 implementation should comfortably finish in 3 minutes with a batch size of 5k (About 1m35 seconds, with nv-nsight).

Batch Size	Op Time 1	Op Time 2	Total Execution Time	Accuracy
100	0.198917ms	0.585584ms	0m2.964s	0.86
1000	2.77284ms	6.33284ms	0m12.430s	0.886
5000	14.748ms	32.8914ms	0m51.459s	0.871

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

Time(%)	Total Time	Calls	Average	Minimum	Maximum	Name
69.1	576325657	20	28816282.9	36210	285361046	cudaMemcpy
23.8	198959200	20	9947960.0	2913	194824399	cudaMalloc
4.9	40668641	10	4066864.1	4007	32019231	cudaDeviceSynchronize
1.9	15810312	10	1581031.2	19018	15562997	cudaLaunchKernel
0.3	2810686	20	140534.3	2102	580987	cudaFree

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

Time(%)	Total Time	Instances	Average	Minimum	Maximum	Name
100.0	40643429	6	6773904.8	7360	32017498	conv_forward_kernel
0.0	2816	2	1408.0	1376	1440	do_not_remove_this_kernel
0.0	2592	2	1296.0	1280	1312	prefn_marker_kernel

```
CUDA Memory Operation Statistics (nanoseconds)
```

Time(%)	Total Time	Operations	Average	Minimum	Maximum	Name
91.7	512753764	6	85458960.7	23328	284642563	[CUDA memcpy DtoH]
8.3	46333124	14	3309508.9	1152	24003756	[CUDA memcpy HtoD]

```
CUDA Memory Operation Statistics (KiB)
```

Total	Operations	Average	Minimum	Maximum	Name
862672.0	6	143778.7	148.535	500000.0	[CUDA memcpy DtoH]
276206.0	14	19729.0	0.004	144453.0	[CUDA memcpy HtoD]

- **Optimization 1: Weight matrix (kernel values) in constant memory (3 point)**

1.

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

I chose to implement weight matrix (kernel values) in constant memory because this technique allows a quicker access to constant memory comparing to accessing from global memory.

- 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?

This implementation works by making constant memory copy the host mask and make it readily available for blocks to access rather than blocks accessing global memory like it used to in M2 implementation. I think the optimization would work assuming that accessing constant memory is quicker than accessing global memory. No, this optimization doesn't synergizes with any other optimization.

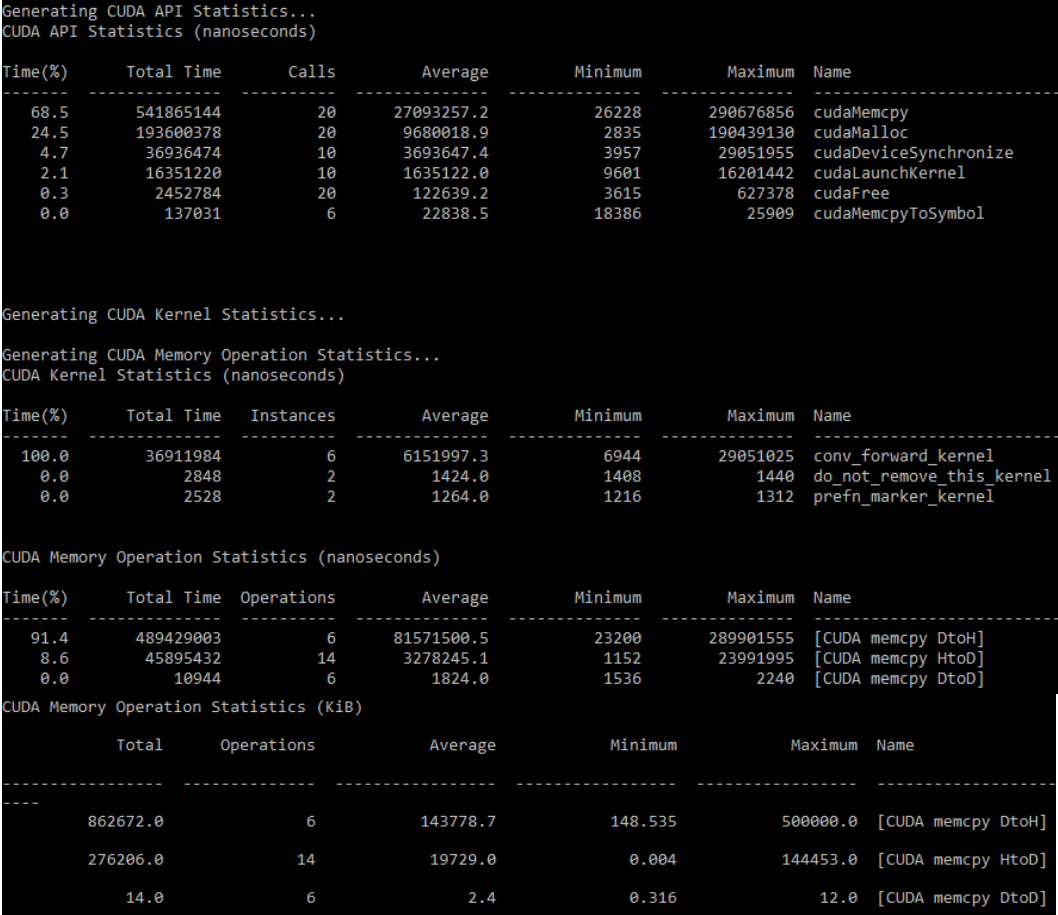
- c. List the Op Times, whole program execution time, and accuracy for batch size of 100, 1k, and 5k images using this optimization (including any previous optimizations also used).

Batch Size	Op Time 1	Op Time 2	Total Execution Time	Accuracy
100	0.172613 ms	0.587303 ms	0m1.669s	0.86
1000	1.56391 ms	5.70719 ms	0m10.278s	0.886
5000	7.70291 ms	28.9042 ms	49.097s	0.871

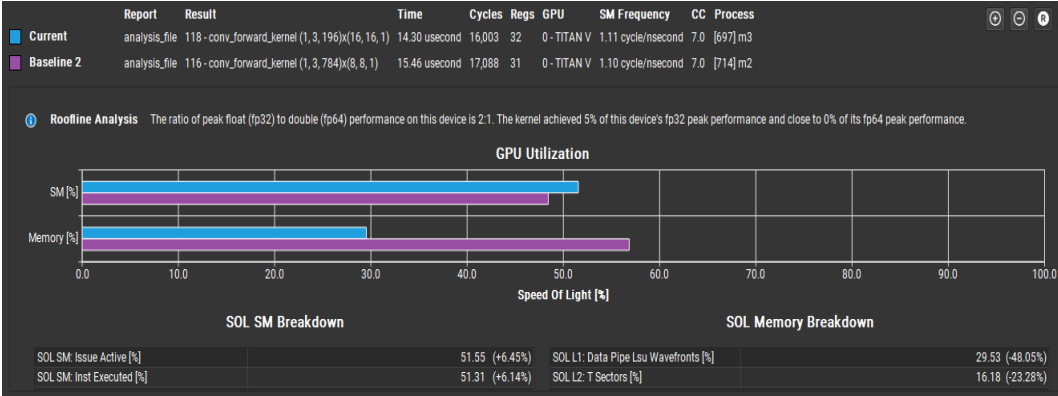
- d. Was implementing this optimization successful in improving performance? Why or why not? Include profiling results from *nsys* and *Nsight-Compute* to justify your answer, directly comparing to your baseline (or the previous optimization this one is built off of).

Yes, it was successful because it is quicker to access constant memory than global memory. While comparing the nsys of M2 and optimization, it is observable that cudaMemcpy and cudaMemcpyToSymbol combined time of the optimization is comparably below cudaMemcpy time of M2. It is also observable that this difference has also effect on kernel time, giving assurance that the optimization is successful.

Nsys:



Nsight-Compute:



- e. What references did you use when implementing this technique?
Lecture slides and Textbook

2. Optimization 2: Input channel reduction: tree

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

I chose Input channel reduction: tree as my optimization because input reduction tree increase parallelism.

- 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?

Optimization works by using tree reduction to increase parallelism by assigning each thread to a channel and then making use of reduction method to achieve higher parallelism. I think the optimization would increase performance of the forward convolution because it manages many channels at same time.

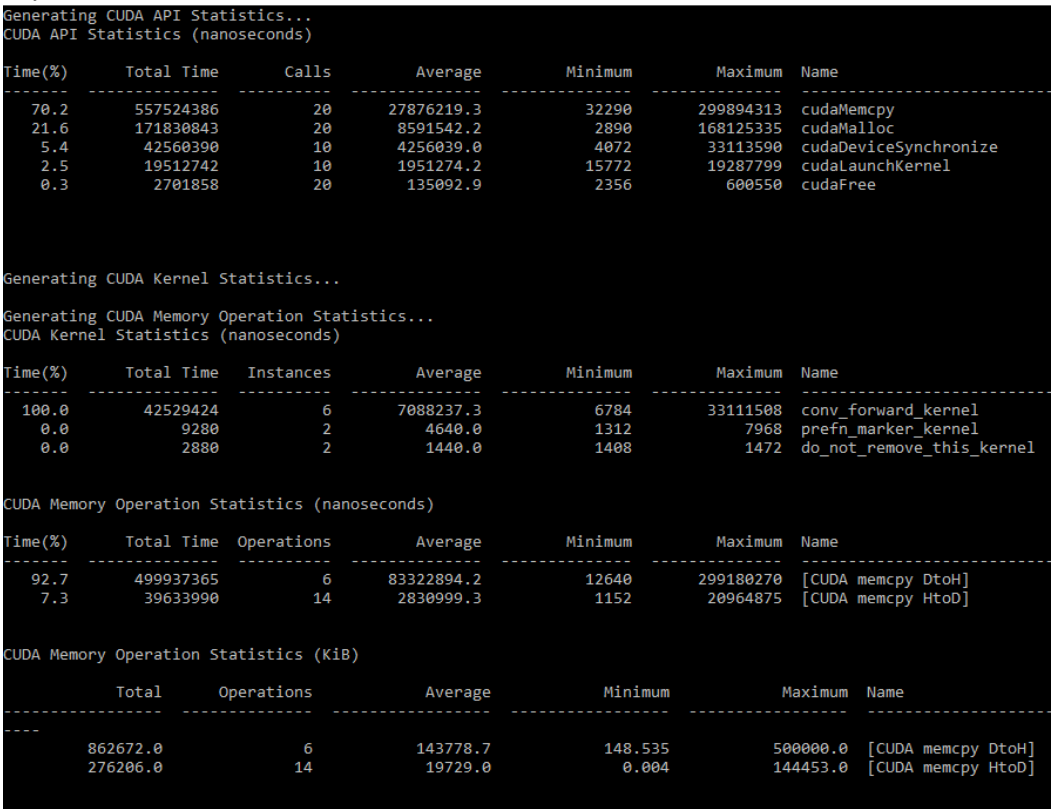
- c. List the Op Times, whole program execution time, and accuracy for batch size of 100, 1k, and 5k images using this optimization (including any previous optimizations also used).

Batch Size	Op Time 1	Op Time 2	Total Execution Time	Accuracy
100	0.205646 ms	0.678121 ms	0m1.581s	0.86
1000	1.89966 ms	6.63198 ms	0m10.352s	0.886
5000	9.42865 ms	33.073 ms	0m51.558s	0.871

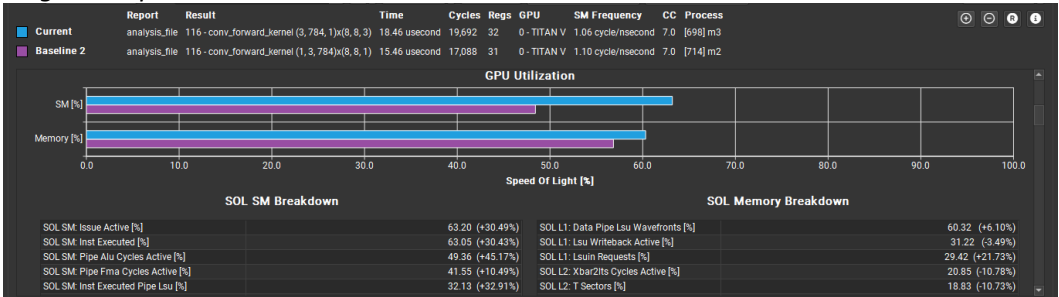
- d. Was implementing this optimization successful in improving performance? Why or why not? Include profiling results from *nsys* and *Nsight-Compute* to justify your answer, directly comparing to your baseline (or the previous optimization this one is built off of).

No, it was not successful and it is comparably worse than M2 baseline. I think the reason could be because of the use of synctreads, which is observable when you compare the difference in kernel time shown in nsys.

Nsys:



Nsight Compute:



- e. What references did you use when implementing this technique?
Lecture slides and Textbook

3. **Optimization 3: Multiple kernel implementations for different layer sizes**
(Delete this section blank if you did not implement this many optimizations.)

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

I chose Multiple kernel implementation as my optimization because I think that kernel incorporating different layer sizes will improve performance.

- 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?

It has two separate kernels with the same algorithm but optimized to incorporate different layer sizes. I think that this optimization would increase performance of the forward convolution because this implementation modifies its parameters to incorporate unique inputs. This optimization synergizes with Weight matrix (kernel values) in constant memory.

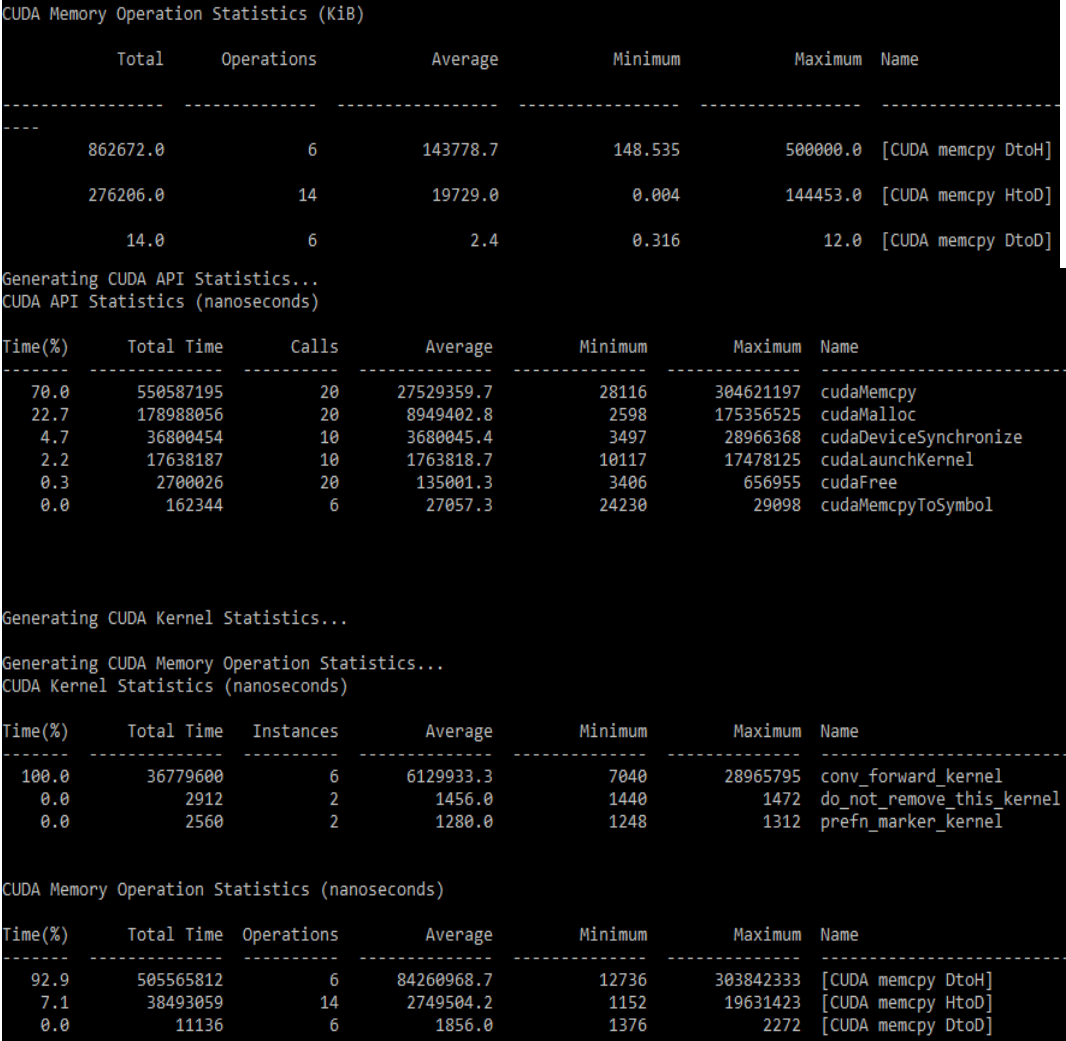
- c. List the Op Times, whole program execution time, and accuracy for batch size of 100, 1k, and 5k images using this optimization (including any previous optimizations also used).

Batch Size	Op Time 1	Op Time 2	Total Execution Time	Accuracy
100	0.176561 ms	0.586961 ms	0m4.156s	0.86
1000	1.5644 ms	5.70981 ms	0m11.248s	0.886
5000	7.7109 ms	28.9908 ms	0m52.801s	0.871

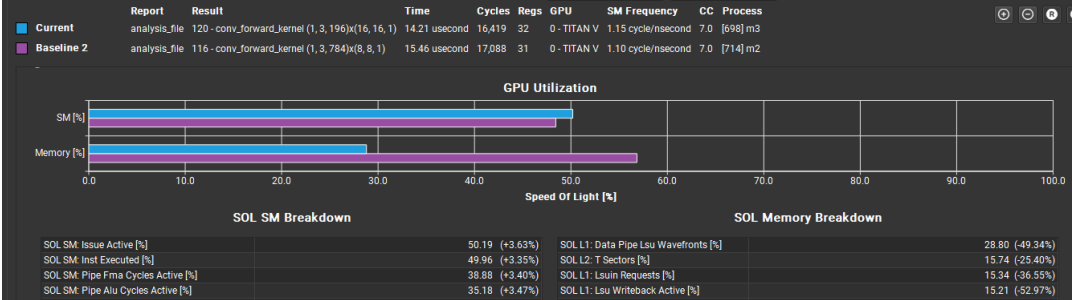
- d. Was implementing this optimization successful in improving performance? Why or why not? Include profiling results from *nsys* and *Nsight-Compute* to justify your answer, directly comparing to your baseline (or the previous optimization this one is built off of).

*Yes, implementing this optimization was successful in improving performance because we successfully incorporated parameters to match with specific unique inputs. It is observable in *nsys* that the kernel time has a significant improvement when compared to M2 baseline and Weight matrix (kernel values) in constant memory optimization which was synergized with this optimization.*

Nsys:



Nsight Compute:



- e. What references did you use when implementing this technique?
Lecture slides and Textbook

4. Optimization 4: Tuning with restrict and loop unrolling (3 points)

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

I chose to implement Tuning with restrict and loop unrolling as my optimization because restrict and unrolling improves performance by implementing a single iteration rather than multiple iteration.

- 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?

In the implementation every data pointer is changed to restricted pointers. A loop is then introduced to iterate through elements of the matrices to allow an advantage to small K values. I think that this optimization would increase performance of the forward convolution because restrict and unrolling improves performance by implementing a single iteration rather than multiple iteration. This optimization don't synergize with any of the other previous optimizations.

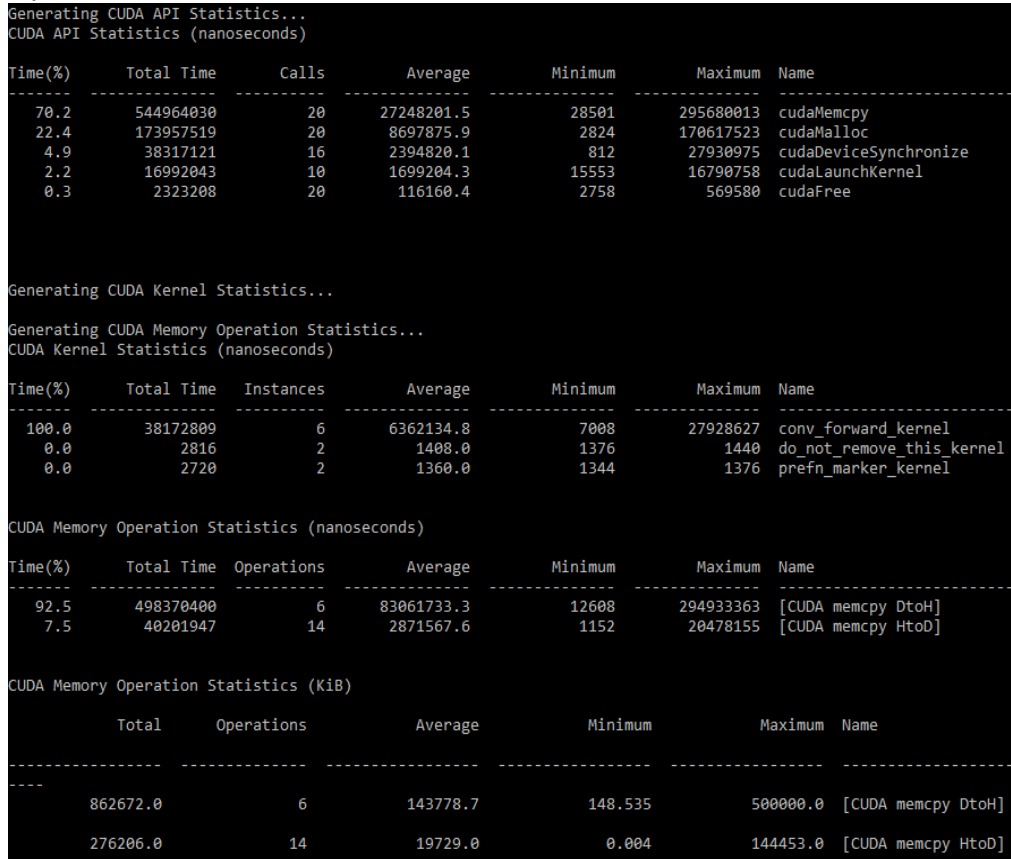
- c. List the Op Times, whole program execution time, and accuracy for batch size of 100, 1k, and 5k images using this optimization (including any previous optimizations also used).

Batch Size	Op Time 1	Op Time 2	Total Execution Time	Accuracy
100	0.215549 ms	0.572844 ms	0m1.777s	0.86
1000	2.02565 ms	5.59096 ms	0m10.812s	0.886
5000	10.036 ms	27.8985 ms	0m51.108s	0.871

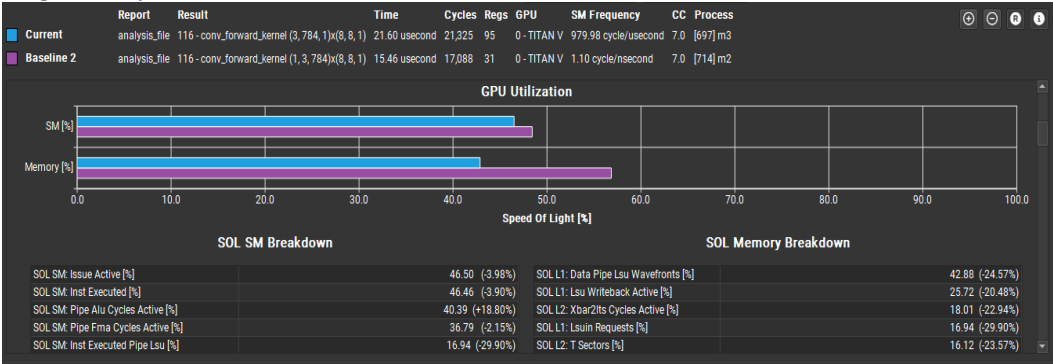
- d. Was implementing this optimization successful in improving performance? Why or why not? Include profiling results from *nsys* and *Nsight-Compute* to justify your answer, directly comparing to your baseline (or the previous optimization this one is built off of).

This optimization was successful in improving performance because restrict and unrolling has increased performance by implementing a single iteration rather than multiple iteration. It is observable in nsys that both cudaMemcpy and cudaMalloc times are significantly lower while comparing to baseline M2. Which is also reflected in the difference in kernel execution time. It is also observable in the Nsight compute that there is a major difference in Memory(%) indicating it is running efficiently than M2 baseline.

Nsys:



Nsight Compute:



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

Lecture and Textbook

5. **Optimization 5:** Tiled shared memory convolution (2 points)

(Delete this section if you did not implement this many optimizations.)

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

I chose Tiled shared memory convolution as my optimization technique because using shared memory will reduce the time consumed for accessing input features.

- 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?

We use shared memory to improve memory reuse rather than accessing global memory which tends to be slow. I think that this optimization would increase performance of the forward convolution, because it is quicker to access shared memory than global memory. This optimization doesn't synergize with any other previous optimizations.

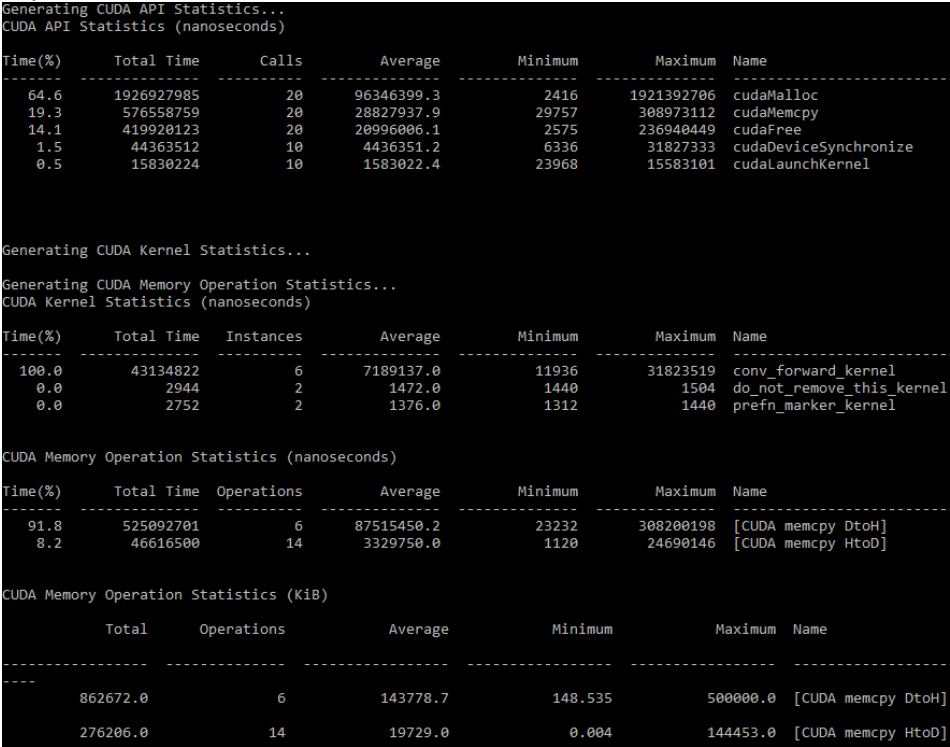
- c. List the Op Times, whole program execution time, and accuracy for batch size of 100, 1k, and 5k images using this optimization (including any previous optimizations also used).

Batch Size	Op Time 1	Op Time 2	Total Execution Time	Accuracy
100	0.751419 ms	1.19764 ms	0m2.830s	0.86
1000	2.54745 ms	6.93911 ms	0m11.514s	0.886
5000	11.7706 ms	32.3134 ms	0m52.301s	0.871

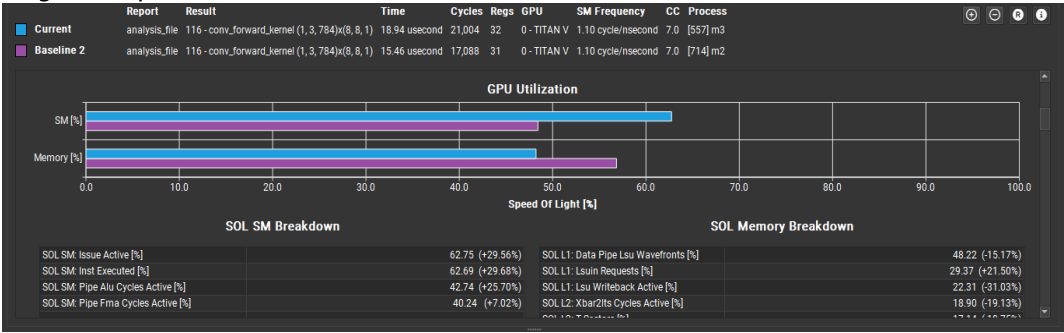
- d. Was implementing this optimization successful in improving performance? Why or why not? Include profiling results from *nsys* and *Nsight-Compute* to justify your answer, directly comparing to your baseline (or the previous optimization this one is built off of).

No, the optimization was not successful in improving performance because there was not much reuse of the shared memory, which would have led to the misuse of resources in shared memory. It is also observable with increased kernel time while comparing it with M2 baseline. While it is evident also in Nsight Compute with the difference in SM(%) while comparing with M2 baseline.

Nsys:



Nsight Compute:



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

6. **Optimization 6:** Sweeping various parameters to find best (0.5 point)

(Delete this section if you did not implement this many optimizations.)

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

I chose Sweeping various parameters as my optimization because it needs very little changes to the code and can result to improvement in performance.

- 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?

For this optimization I tested Tile Sizes and fixed on Tile Size of 24. I think that the optimization would increase performance of the forward convolution because trying different Tile Size can result to a better performance. This optimization doesn't synergize with any of the other optimizations.

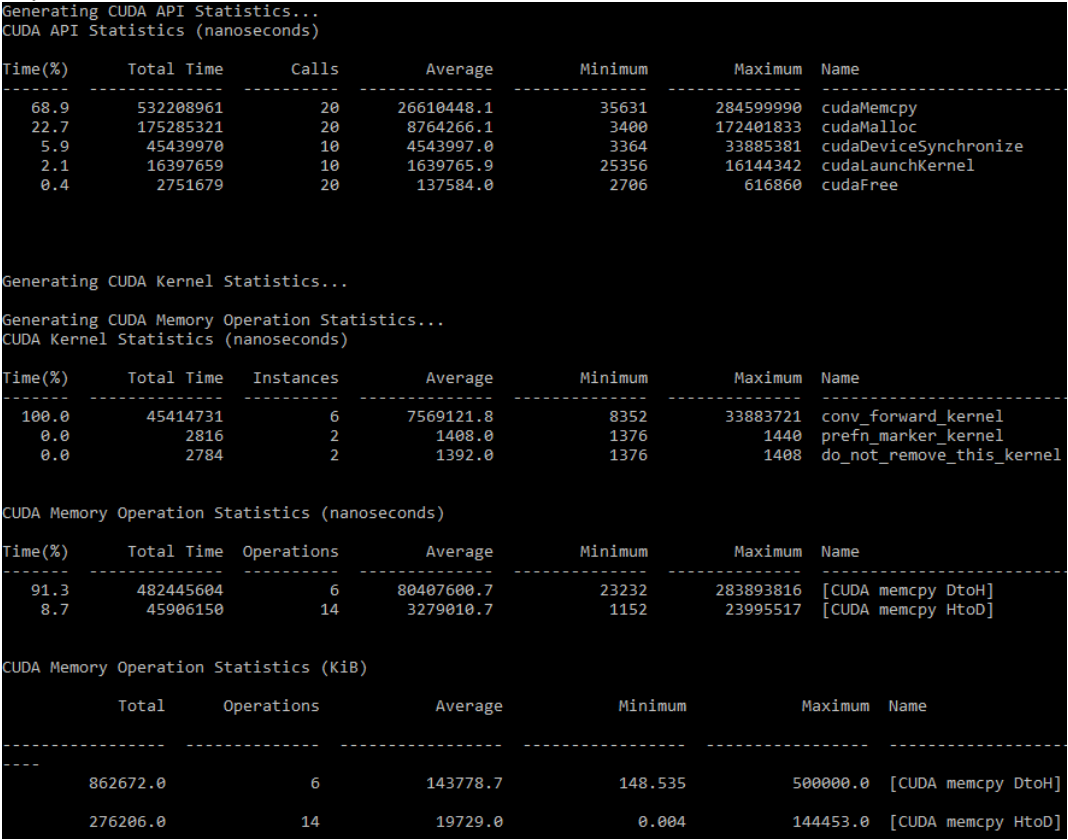
- c. List the Op Times, whole program execution time, and accuracy for batch size of 100, 1k, and 5k images using this optimization (including any previous optimizations also used).

Batch Size	Op Time 1	Op Time 2	Total Execution Time	Accuracy
100	0.243166 ms	0.687001 ms	0m1.658s	0.86
1000	2.29288 ms	6.77237 ms	0m10.609s	0.886
5000	11.337 ms	34.2141 ms	0m51.766s	0.871

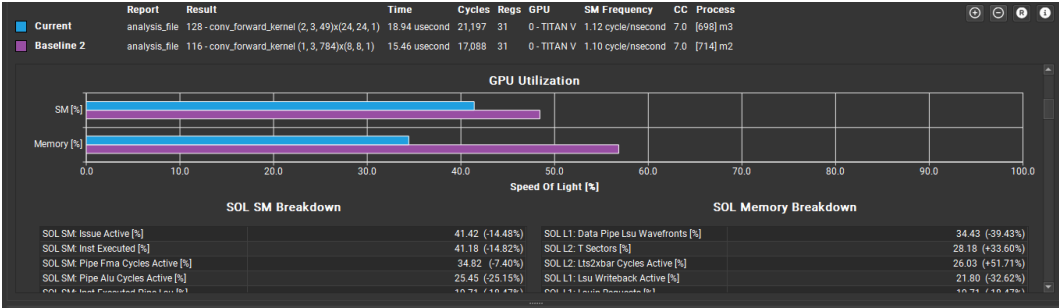
- d. Was implementing this optimization successful in improving performance? Why or why not? Include profiling results from *nsys* and *Nsight-Compute* to justify your answer, directly comparing to your baseline (or the previous optimization this one is built off of).

Implementing this optimization is successful in improving performance because as seen with Op time which has decreased significantly from M2 baseline. It is also observable in nsys where the kernel time has improved significantly show that the performance has improved. It is also represented in Nsight Compute where both SM(%) and Memory(%) are lower compared to M2 baseline.

Nsys:



Nsight Compute:



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