**Name:** *Ge Yuhao*
**NetID:** *Yuhaoge2@illinois.edu*
**Section:** *ZJ1/ZJ2*

# 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 | 0.174333ms | 0.630749ms | 1.216s | 0.86 |
| 1000 | 1.63886ms | 6.25089ms | 9.673s | 0.886 |
| 10000 | 16.0759ms | 62.9824ms | 1m34.672s | 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.

   *Optimize through using tiled shared memory convolution. Because the original method read image data and weight data from the global memory, the speed can be slow, and the total time may be limited by the memory bandwidth. By first storing the image data and weight data into the shared memory, the later access time will be smaller.*

   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?
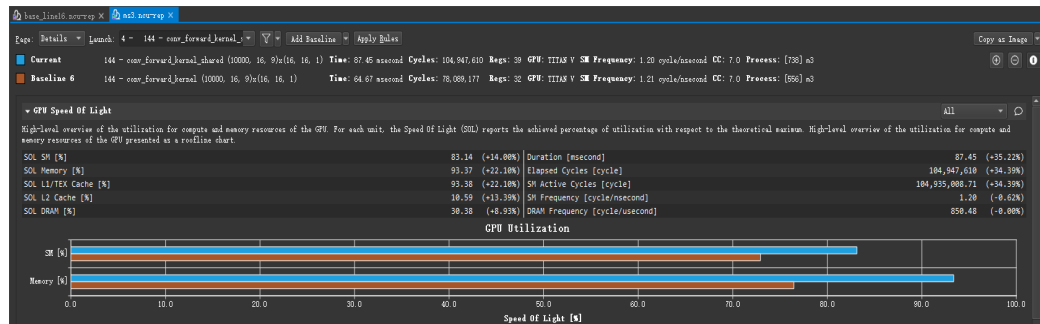
   *Divide each image into many tiles, and each block will calculate one tile. First copy the image date from the global memory into the shared memory, so each thread in a block can use the same shared memory to do calculation. Then the result will write back to the global memory. I think this optimization would increase performance of the forward convolution because the shared memory reading is much faster than the global memory. This is the first optimization I try.*

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 Layer Time | Total Execution Time | Accuracy |
|---|---|---|---|---|---|
| 100 | 0.2387ms | 0.6342ms | 11.3232ms | 1.324s | 0.86 |
| 1000 | 2.4932ms | 6.3892ms | 114.1242ms | 9.523s | 0.886 |
| 10000 | 24.007ms | 63.5762ms | 1120.903ms | 1m28.123s | 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).

*The performance is not improved, and the speed is even slower. From the Nsight-Compute, the utilization of SM and memory does increase while using shared memory, but the total time needed is increased. I think this is because the input images' size is not that large, and the shared memory is not fully used, or the original op time is not limited by the global memory bandwidth. Also, copying data from global memory to shared memory also need time.*



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

*I refer to the chapter16 of the textbook for this technique*

2. **Optimization 2: <Weight matrix (kernel values) in constant memory>**

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

*I choose to implement a new version with Weight matrix (kernel values) in constant memory. Because the kernel values are frequently used in the convolution, and if it is in constant memory the total speed may be faster.*

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?
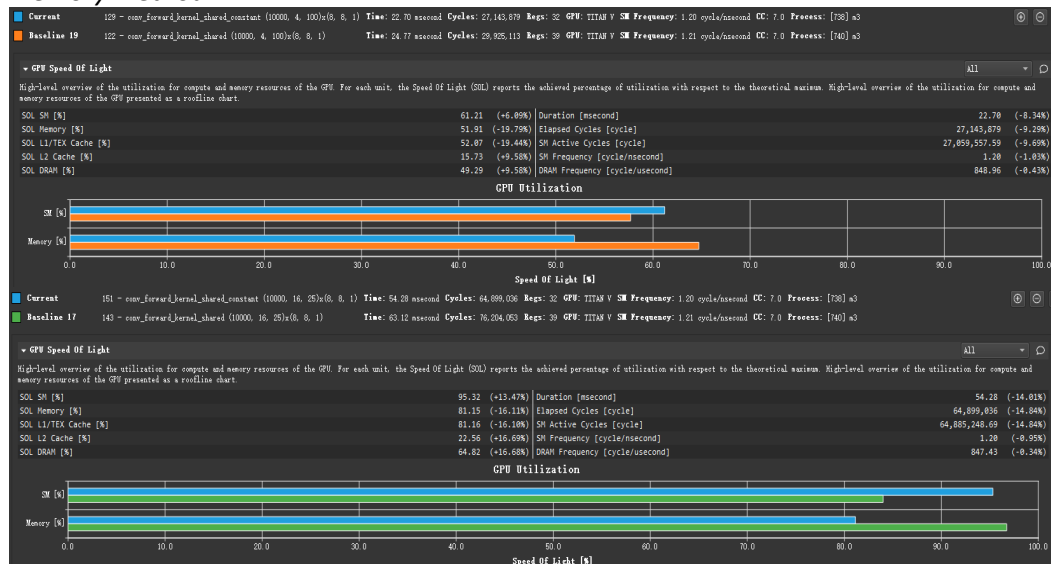
*In the code, I copy the weight matrix data into the constant memory and use this to do calculation in the kernel. I think this would increase performance, as we only need to copy the weight matric data once and can reach it quickly in later use. I think this will synergize with the former shared memory method.*

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 Layer Time | Total Execution Time | Accuracy |
|---|---|---|---|---|---|
| 100 | *0.2121ms* | *0.5421ms* | *10.782ms* | *1.258s* | *0.86* |
| 1000 | *2.2791ms* | *5.4101ms* | *113.217ms* | *9.781s* | *0.886* |
| 10000 | *22.6564ms* | *54.0805ms* | *1121.443ms* | *1m25.421s* | *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).

*This optimization successfully improved the performance. Because the optime become smaller than the former shared memory method. From the Nsight-Compute, the utilization of SM is increased and the use of memory is decreased, and this is corresponding to my theoretical analysis. Also, as the copying is conducted in the host, no more kernel time is needed for memory copying, so this is faster than the shared memory method.*



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

*I refer to the lecture talking about constant memory.*

3. **Optimization 3: <Shared memory matrix multiplication and input matrix unrolling>**

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

      *I use the shared memory matrix multiplication and input matrix unrolling method. Because I think the former shared memory method is not that good, maybe a brand-new algorithm can help.*

   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 method uses a brand-new algorithm, it first transfers the original image data into new arrangement. Then a matrix multiplication between the kernel and the newly rearranged data can give the result. I think this may increase the performance, because we can do many other optimizations basing on the matrix multiplication such as shared memory tiling. This method has nothing to do with the former methods.*
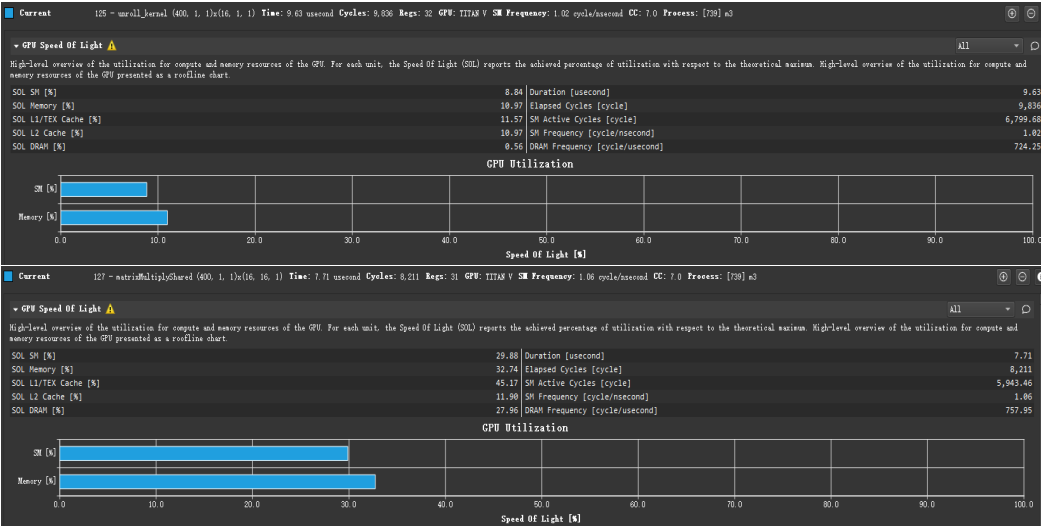
   c. List the Op Times, whole program execution time, and accuracy for batch size of 100, 1k, and 10k images using this optimization (including any previous optimizations also used).

      | Batch Size | Op Time 1 | Op Time 2 | Total Layer Time | Total Execution Time | Accuracy |
      |---|---|---|---|---|---|
      | 100 | *1.5268ms* | *2.291ms* | *13.291ms* | *1.316s* | *0.86* |
      | 1000 | *14.3975ms* | *23.101ms* | *136.712ms* | *9.671s* | *0.886* |
      | 10000 | *144.495ms* | *224.985ms* | *1353.618ms* | *1m31.671s* | *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).

      *The result is not that well. The op time increases a lot. I think the reason is that when we rearrange the data from the input image data into the matrix operands, the order of the data is messed up and the merit of memory burst may have gone.*

      *From the Nsight-Compute, we can see that the utilization of SM and Memory is extremely low for both two kernels. And there are many loops to do all the jobs which means the program is more serially rather than parallelly.*

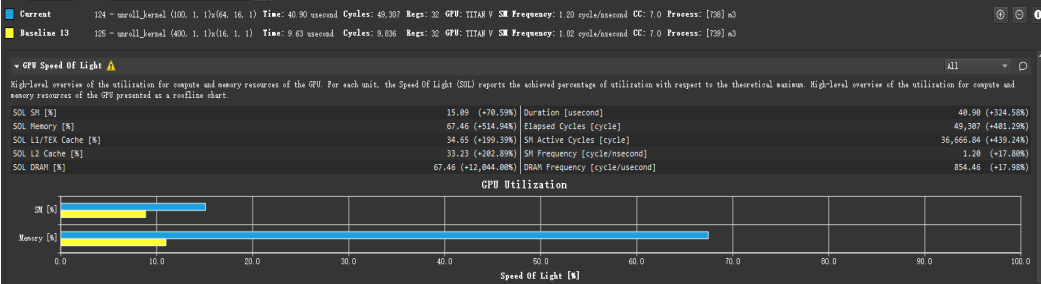e. What references did you use when implementing this technique?

*I refer to the lecture and the textbook talking about this technique.*
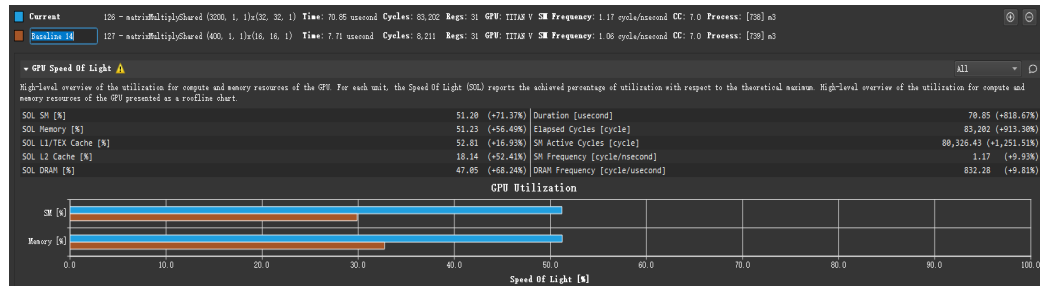
f. **Further modification (additional, not sure if this can count as points. Plz check it)**

*As the performance is extremely bad, I think I can do more modification on the data rearrangement. Originally, the method rearranges one image's data into a new matrix, but as our original data is considerably small, this rearrangement is not sufficient. So, I changed the algorithm a little bit, to make the rearrangement happen on several image together to produce one big matrix and do matrix multiplication. And do some index modification when storing the data back to the global memory (The code is in new-forward-mod-unroll.cu)*

| Batch Size | Op Time 1 | Op Time 2 | Total Layer Time | Total Execution Time | Accuracy |
|------------|-----------|-----------|------------------|----------------------|----------|
| 100 | *0.6667ms* | *0.5234ms* | *1.1002ms* | *1.429s* | *0.86* |
| 1000 | *6.5232ms* | *5.2612ms* | *12.3212ms* | *9.582s* | *0.886* |
| 10000 | *66.7534ms* | *58.247ms* | *1107.559ms* | *1m29.431s* | *0.8714* |

*It's a surprise that the op time is much smaller than the original unrolling method. Although the time is still larger than the base line, this modification is a success. This can also be proved be the Nsight-Compute. For both kernel, there is a boost in SM and Memory utilization.*

```
Current        128 - matrixMultiplyShared (3200, 1, 1)x(32, 32, 1)  Time: 70.85 usecond  Cycles: 83,202  Regs: 31  GPU: TITAN V  SM Frequency: 1.17 cycle/nsecond  CC: 7.0  Process: [738] m3
Baseline 14    127 - matrixMultiplyShared (400, 1, 1)x(16, 16, 1)  Time: 7.71 usecond  Cycles: 8,211  Regs: 31  GPU: TITAN V  SM Frequency: 1.06 cycle/nsecond  CC: 7.0  Process: [739] m3

▼ GPU Speed Of Light ⚠                                                                                                                                                          All        ▼  ○
High-level overview of the utilization for compute and memory resources of the GPU. For each unit, the Speed Of Light (SOL) reports the achieved percentage of utilization with respect to the theoretical maximum. High-level overview of the utilization for compute and
memory resources of the GPU presented as a roofline chart.
SOL SM [%]                                       51.20  (+71.37%)  Duration [usecond]                                 70.85  (+818.67%)
SOL Memory [%]                                   51.23  (+56.49%)  Elapsed Cycles [cycle]                             83,202  (+913.30%)
SOL L1/TEX Cache [%]                             52.81  (+16.93%)  SM Active Cycles [cycle]                        80,326.43  (+1,251.51%)
SOL L2 Cache [%]                                 18.14  (+52.41%)  SM Frequency [cycle/nsecond]                         1.17  (+9.93%)
SOL DRAM [%]                                     47.05  (+68.24%)  DRAM Frequency [cycle/usecond]                     832.28  (+9.81%)
```

4.  **Optimization 4: *<Multiple kernel implementations for different layer sizes>***

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

        *I use multiple kernel implementations for different layer sizes to optimize. When I am testing the former methods and trying with different TILE_WIDTH, I found that a big tile width is better for the first layer and a small tile width is better for the second layer. So, I think I can use different kernel configuration for different layer.*

    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 add an "if" logic in my host code. If the size of the image is larger than 64, I will use 16 as the TILE_WIDTH and block dimension, or I will use 8. I think this will increase performance because if the input image is large, a larger tile may have better use of the threads, if the input image is small, a smaller tile may cause less control divergence. I think this method can add to any other methods because it only changes the size of kernel. I tested this method in many former implementations and the speed is always faster. When profiling, I use this method together with the shared memory method and constant weight method.*
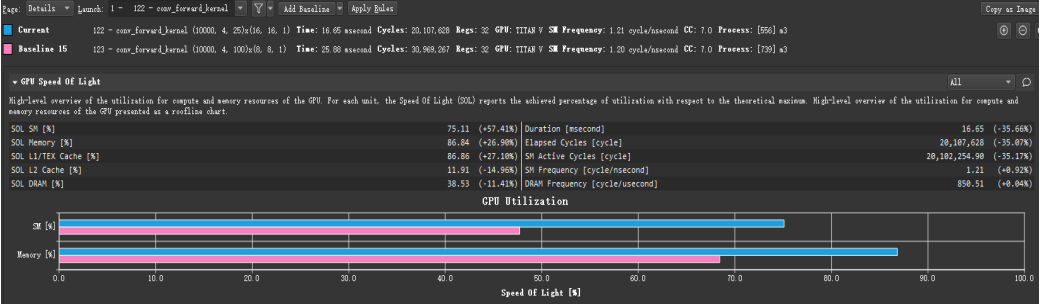
    c.  List the Op Times, whole program execution time, and accuracy for batch size of 100, 1k, and 10k images using this optimization (including any previous optimizations also used).

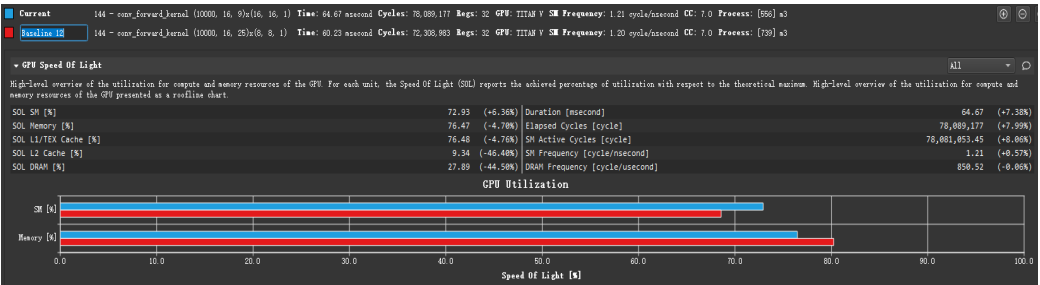| Batch Size | Op Time 1 | Op Time 2 | Total Layer Time | Total Execution Time | Accuracy |
|---|---|---|---|---|---|
| 100 | 0.1454ms | 0.5541ms | 11.2353ms | 1.361s | 0.86 |
| 1000 | 1.5394ms | 5.6181ms | 113.131ms | 9.212s | 0.886 |
| 10000 | 15.1784ms | 56.6811ms | 1132.293ms | 1m29.119s | 0.8714 |

    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).

*This optimization is successful as the total op time is smaller. Below is the profiling data I get when testing the different configurations for different layers.*
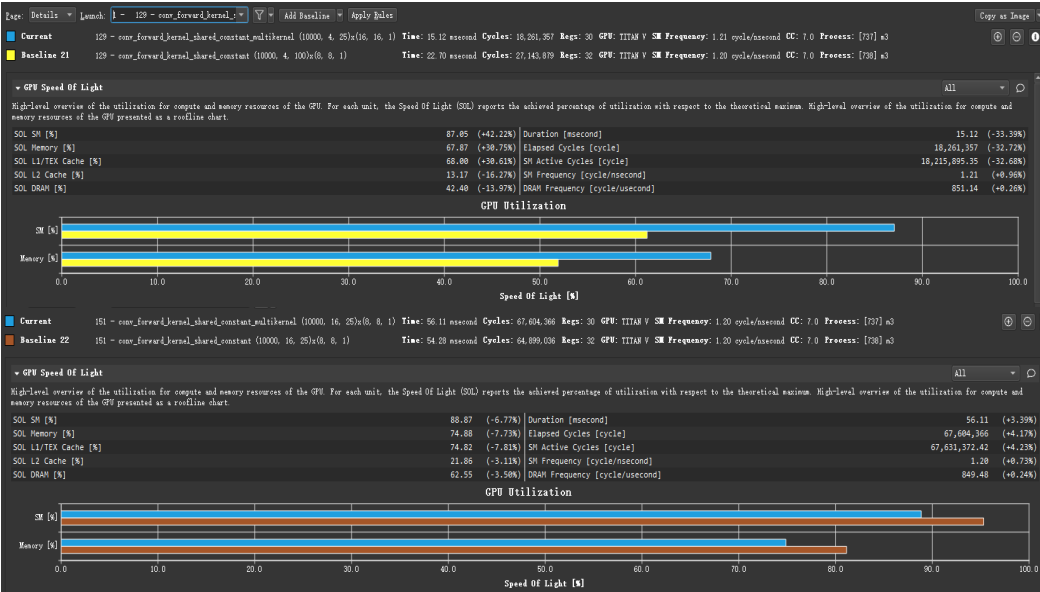
*For the first layer, when using the block size of 16\*16\*1, the utilization of SM is 57% larger and the memory use is 27% larger comparing with the size of 8\*8\*1, and the time for op time1 reduced from 25.8 ms to 16.65 ms.*



*For the second layer, the utilization of SM and memory is almost the same for bock size of 16\*16\*1 and 8\*8\*1, but the op time 2 is 4ms smaller when using the size of 8\*8\*1.*



*So using block size of 16\*16\*1 in the first layer and 8\*8\*1 in the second will improve the performance. So I try to add this method to the former shared memory and constant weight optimization, and the result is good. The total op time is 71.3ms which is very close to the 70ms requirement.*



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

*I refer to the data of my testing to raise up this method.*

5. **Optimization 5: *<Input channel reduction: atomics>***

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

    *I use input channel reduction method, because the dimension of the kernel is not fully used in the original code, and I think I can use z-dimension threads to make the computation in different channel parallelized. And I also need to use atomic addition to get rid of confliction.*

    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?
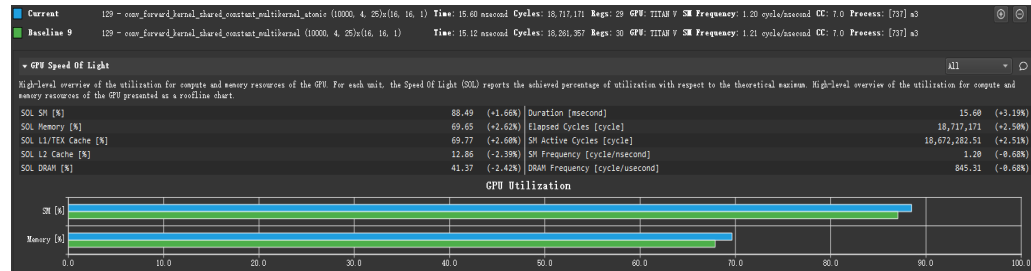
    *I change the configuration of the kernel and add new z dimension threads to make the computation in different channel parallelized. When adding the result back to the global memory, I use atomic addition to get rid of confliction. I think this would increase performance because it makes the calculation more paralleled. This optimization can build on any former method with a channel loop. As the performance for the optimization "shared memory, constant weight, Multiple kernel for different layer" is good, I add this new optimization onto the former optimization.*

    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).
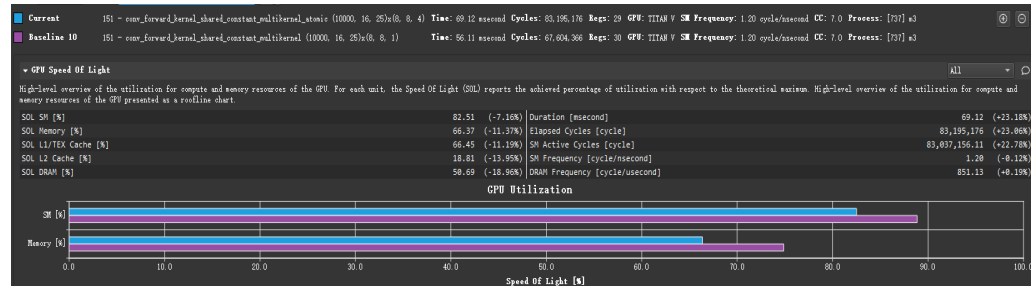
| Batch Size | Op Time 1 | Op Time 2 | Total Execution Time | Accuracy |
|---|---|---|---|---|
| 100 | *0.281ms* | *0.661ms* | *1.419s* | *0.86* |
| 1000 | *2.792ms* | *6.819ms* | *9.172s* | *0.886* |
| 10000 | *28.218ms* | *69.0906ms* | *1m31.291s* | *0.8714* |

    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).

    *As for the first layer, the input channel is 1, so there is no difference when using input channel reduction*

*As for the second layer, the input channel is 4, so the dimension of z will increase to 4. But from the Nsignt-Compute, we can see that the performance becomes worse, and more time is needed.*



*I think this is because although the input channel is computed parallelly, the threads also need to wait in line when doing the atomic addition. Also, the total number of threads in one block is limited by 1024, if more are allocated in the z dimension, less will be allocated in the x-y plate.*

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

*I refer to the lecture talking about reduction and histogram.*

6.  **Optimization 6: <Using Streams to overlap computation with data transfer>**

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

*I use streams to overlap computation with data transfer. I use this method because I think continue to optimize op time is kind of hard and there are some optimization spaces for the data transferring time in host code.*

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 use many streams to make the memory copy and calculation asynchronized, and this will make the computation and data transferring of different streams occur at the same time. I think this optimization can add to any former methods. When profiling, I use this method above the constant weight method.*

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 | Layer Time 1 | Layer Time 2 | Total Execution Time | Accuracy |
|---|---|---|---|---|
| 100 | 5.719ms | 4.499ms | 1.398s | 0.86 |
| 1000 | 58.129ms | 43.142ms | 9.019s | 0.886 |
| 10000 | 588.248ms | 429.997ms | 1m28.111s | 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).

*I run the test many times and find that this method will increase the speed and synergize with any of previous optimizations. Also, I find that the performance is the best when the stream number is 5. From the Nsight-system, we can see that by using the method of overlap, different streams can do the computation and copying at the same time. From the layer time data we can find that the total layer time is reduced by using this method.*



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

*I use this technique basing on the knowledge from the class.*

7. **Optimization 7: *<Combine different methods together (Pin memory method is added)>***

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

*I use constant memory for weight, Multiple kernel implementations for different layer sizes, Streams to overlap computation with data transfer, sweeping various parameters to find best values, and pin memory to implement the fastest version. I choose this combination as my final result basing on tons of running test, the log is in my code folder.*

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?

*First, the weight data is stored into the constant memory, so each time I need the weight data, it can be fetched from the fast constant memory. Then, the kernel size of each layer is specially chosen (16 for big picture, 8 for small picture). Next, use five streams to do overlapping, and loop 10 times, so each kernel will calculate 200 images at one time (this configuration is determined by many tests, logs are also in the ./log folder). Lastly, pin the device memory to avoid further memory copying.*
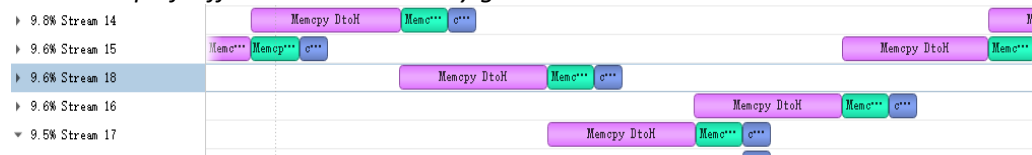
c. List the Op Times, whole program execution time, and accuracy for batch size of 100, 1k, and 10k images using this optimization (including any previous optimizations also used).

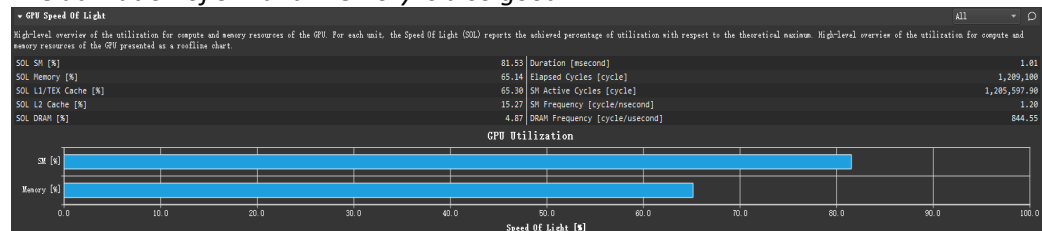| Batch Size | Layer time 1 | Layer Time 2 | Total Execution Time | Accuracy |
|---|---|---|---|---|
| 100 | 2.518ms | 1.841ms | 1.419s | 0.86 |
| 1000 | 26.234ms | 18.192ms | 9.121s | 0.886 |
| 10000 | 260.575ms | 187.767ms | 1m26.341s | 0.8714 |

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).

*The optimization is extremely successful in improving performance. The total layer time is in top three of the class. Basing on all the former optimization trails, each method I use in this compound optimization can somehow optimize the speed. No wonder this combination can boost the final speed. By looking at the nsys, I find that originally, the time is mostly consumed on the memory copying (97% for first layer, 89% for second layer). By using the pin memory method, this time is reduced visibly (87% for first layer, 62% for second layer).*
*The overlap of different streams is very good.*



*The utilization of SM and Memory is also good.*



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

*Experience from all other optimizations above.*