# Lab 6

George Onwubuya

November 15, 2018

## 1 Histogram

### 1.1 Output

```
Setting up the problem...0.002114 s
    Input size = 100000
    Number of bins = 4096
Allocating device variables...0.000280 s
Copying data from host to device...0.000272 s
Launching kernel...0.000452 s
Copying data from device to host...0.000039 s
Verifying results...TEST PASSED


Setting up the problem...0.003954 s
    Input size = 200000
    Number of bins = 4096
Allocating device variables...0.000304 s
Copying data from host to device...0.000484 s
Launching kernel...0.000466 s
Copying data from device to host...0.000040 s
Verifying results...TEST PASSED


Setting up the problem...0.008680 s
    Input size = 400000
    Number of bins = 4096
Allocating device variables...0.000328 s
Copying data from host to device...0.000805 s
Launching kernel...0.000458 s
Copying data from device to host...0.000039 s
Verifying results...TEST PASSED
```

```
Setting up the problem...0.017020 s
    Input size = 800000
    Number of bins = 4096
Allocating device variables...0.000258 s
Copying data from host to device...0.001406 s
Launching kernel...0.000510 s
Copying data from device to host...0.000043 s
Verifying results...TEST PASSED


Setting up the problem...0.022502 s
    Input size = 1000000
    Number of bins = 4096
Allocating device variables...0.000264 s
Copying data from host to device...0.001899 s
Launching kernel...0.000525 s
Copying data from device to host...0.000041 s
Verifying results...TEST PASSED


Setting up the problem...0.033799 s
    Input size = 1600000
    Number of bins = 4096
Allocating device variables...0.000311 s
Copying data from host to device...0.002794 s
Launching kernel...0.000543 s
Copying data from device to host...0.000040 s
Verifying results...TEST PASSED
```

## 1.2   Performance Analysis (4096 Bins)

| Execution Time (seconds) for Each Process | | | | | |
|---|---|---|---|---|---|
| Elements(m) | Setting Up | DeviceVar | HostToDevice | Kernel | DeviceToHost |
| 100000 | 0.002114 | 0.000280 | 0.000272 | 0.000452 | 0.000039 |
| 200000 | 0.003954 | 0.000304 | 0.000484 | 0.000466 | 0.000040 |
| 400000 | 0.008680 | 0.000328 | 0.000805 | 0.000458 | 0.000039 |
| 800000 | 0.017020 | 0.000258 | 0.001406 | 0.000510 | 0.000043 |
| 1000000 | 0.022502 | 0.000264 | 0.001899 | 0.000525 | 0.000041 |
| 1600000 | 0.033799 | 0.000311 | 0.002794 | 0.000543 | 0.000040 |

### 1.2.1   Comments

When setting up the shell script file, the number of elements were varied but I chose to use the same number of bins, 4096. The time taken to allocate device variables and copy from device to host are approximately the same. This is

because the same number of variables are allocated regardless of the size of the array of input elements. The times taken to set up the problem, copy from host to device and launching the kernel are all directly proportional to the size of the input array.

# 2 Histogram (Optimized)

## 2.1 Output

```
Setting up the problem...0.002145 s
    Input size = 100000
    Number of bins = 4096
Allocating device variables...0.000262 s
Copying data from host to device...0.000236 s
Launching kernel...0.000498 s
Copying data from device to host...0.000036 s
Verifying results...TEST PASSED


Setting up the problem...0.004105 s
    Input size = 200000
    Number of bins = 4096
Allocating device variables...0.000303 s
Copying data from host to device...0.000387 s
Launching kernel...0.000504 s
Copying data from device to host...0.000042 s
Verifying results...TEST PASSED


Setting up the problem...0.008779 s
    Input size = 400000
    Number of bins = 4096
Allocating device variables...0.000259 s
Copying data from host to device...0.000713 s
Launching kernel...0.000528 s
Copying data from device to host...0.000041 s
Verifying results...TEST PASSED


Setting up the problem...0.016875 s
    Input size = 800000
    Number of bins = 4096
Allocating device variables...0.000261 s
Copying data from host to device...0.001438 s
```

```
Launching kernel...0.000560 s
Copying data from device to host...0.000041 s
Verifying results...TEST PASSED


Setting up the problem...0.021037 s
    Input size = 1000000
    Number of bins = 4096
Allocating device variables...0.000290 s
Copying data from host to device...0.001656 s
Launching kernel...0.000589 s
Copying data from device to host...0.000040 s
Verifying results...TEST PASSED


Setting up the problem...0.033860 s
    Input size = 1600000
    Number of bins = 4096
Allocating device variables...0.000284 s
Copying data from host to device...0.002730 s
Launching kernel...0.000629 s
Copying data from device to host...0.000040 s
Verifying results...TEST PASSED
```

## 2.2   Performance Analysis

### 2.2.1   Array Size

| Execution Time (seconds) for Each Process | | | | | |
|---|---|---|---|---|---|
| Elements(m) | Setting Up | DeviceVar | HostToDevice | Kernel | DeviceToHost |
| 100000 | 0.002145 | 0.000262 | 0.000236 | 0.000498 | 0.000036 |
| 200000 | 0.004105 | 0.000303 | 0.000387 | 0.000504 | 0.000042 |
| 400000 | 0.008779 | 0.000259 | 0.000713 | 0.000528 | 0.000041 |
| 800000 | 0.016875 | 0.000261 | 0.001438 | 0.000560 | 0.000041 |
| 1000000 | 0.021037 | 0.000290 | 0.001656 | 0.000589 | 0.000040 |
| 1600000 | 0.033860 | 0.000284 | 0.002730 | 0.000629 | 0.000040 |

### 2.2.2   Comments

The execution times relate to the same way to the initial histogram kernel that
had no optimization.

## 2.3 Answers

### 2.3.1 Description of Optimization

To optimize the histogram kernel, the kernel code was written such that a location in shared memory would only be updated once. Three new variables were introduced to the kernel code. The first two variables current index and previous index were meant to keep track of the elements from the input array. The current index and previous were compared to check if the same element was still being evaluated because of possible repetitions in data. If they were found to be the same, a third variable accumulator would store the number of times this element occurred in succession in the array . However if the elements were not the same, the value in the accumulator would be written to a location in shared memory as indicated by the previous index, the previous index would then be set to the current index and the value in the accumulator would be set to one to indicate a shift to the next element. The accumulator and the previous index variables were initialized to zero outside the while loop and the current index variable was set to the current input element under evaluation by the histogram kernel code.
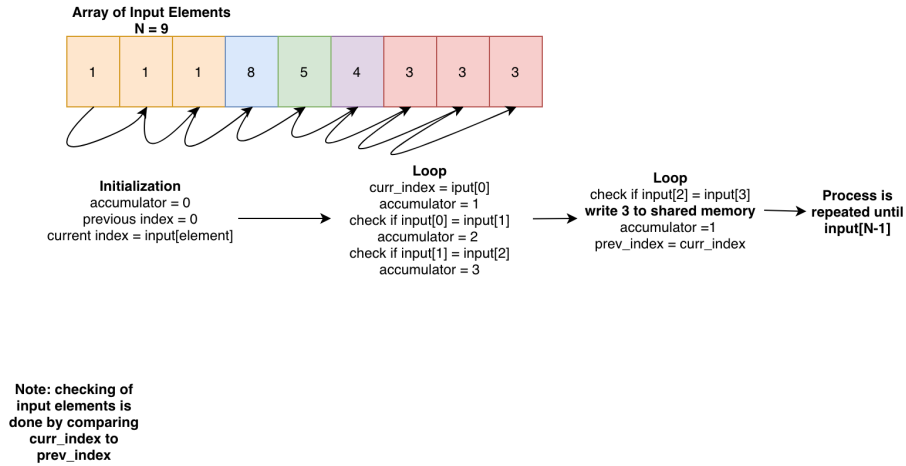


Figure 1: An illustration of what the kernel code does

5

### 2.3.2  Difficulties with Optimization

A couple of difficulties I experienced during this optimization were found in not analyzing the the second while loop structure of the code that is used to update the value in the accumulator to the shared memory location. I examined how the loop would behave at the beginning of the and during the middle part of the execution and neglected to check what happens when the loop ends. As a result my code was incomplete and was unable to simulate past the first bin. I also made the mistake of attempting to carry out the optimization without first checking if the code was properly written out.

### 2.3.3  Change in Execution Time

| Kernel Times in microseconds | | |
| --- | --- | --- |
| Elements(m) | Non-Optimized | Optimized |
| 100000 | 11.776 | 76.414 |
| 200000 | 18.367 | 82.398 |
| 400000 | 32.614 | 99.646 |
| 800000 | 63.039 | 129.12 |
| 1000000 | 77.356 | 144.09 |
| 1600000 | 120.38 | 190.17 |

```
==16836== NVPROF is profiling process 16836, command: /home/onwubuyag/mpplabs/lab6-histogram/histogram 1600000
==16836== Profiling application: /home/onwubuyag/mpplabs/lab6-histogram/histogram 1600000
==16836== Profiling result:
            Type  Time(%)      Time     Calls       Avg       Min       Max  Name
 GPU activities:   94.55%   2.3343ms         1   2.3343ms   2.3343ms   2.3343ms  [CUDA memcpy HtoD]
                    4.88%   120.38us         1   120.38us   120.38us   120.38us  histogram_kernel(unsigned int*, unsigned int*, unsigned int, unsigned int)
                    0.28%   6.9760us         2   3.4880us   2.2720us   4.7040us  [CUDA memset]
                    0.17%   4.3200us         1   4.3200us   4.3200us   4.3200us  [CUDA memcpy DtoH]
                    0.11%   2.7520us         1   2.7520us   2.7520us   2.7520us  convert_kernel(unsigned int*, unsigned char*, unsigned int)
      API calls:   97.41%   141.24ms         4   35.310ms   102.12us   140.26ms  cudaFree
                    1.90%   2.7523ms         2   1.3762ms   33.967us   2.7183ms  cudaMemcpy
                    0.27%   392.60us         3   130.87us   118.50us   155.02us  cudaMalloc
                    0.11%   157.60us         2   78.799us   15.451us   142.15us  cudaLaunch
                    0.11%   154.40us        94   1.6420us     451ns   38.678us  cuDeviceGetAttribute
                    0.10%   149.68us         1   149.68us   149.68us   149.68us  cuDeviceTotalMem
                    0.05%   70.629us         4   17.657us   4.6120us   47.116us  cudaDeviceSynchronize
                    0.02%   28.445us         2   14.222us   11.024us   17.421us  cudaMemset
                    0.02%   23.885us         1   23.885us   23.885us   23.885us  cuDeviceGetName
                    0.01%   9.9800us         7   1.4250us     439ns   6.4630us  cudaSetupArgument
                    0.00%   6.1310us         3   2.0430us     585ns   4.3260us  cuDeviceGetCount
                    0.00%   2.6600us         2   1.3300us     887ns   1.7730us  cudaConfigureCall
                    0.00%   1.7650us         2      882ns     608ns   1.1570us  cuDeviceGet
```

Figure 2: Nvidia profiler result for non-optimized histogram kernel for m = 1600000

```
==17003== NVPROF is profiling process 17003, command: /home/onwubuyag/mpplabs/lab6-histogram/histogram 1600000
==17003== Profiling application: /home/onwubuyag/mpplabs/lab6-histogram/histogram 1600000
==17003== Profiling result:
            Type  Time(%)      Time     Calls       Avg       Min       Max  Name
 GPU activities:   91.79%  2.2817ms         1  2.2817ms  2.2817ms  2.2817ms  [CUDA memcpy HtoD]
                    7.65%  190.17us         1  190.17us  190.17us  190.17us  histogram_kernel_optimized(unsigned int*, unsigned int*, unsigned int, unsigned i
nt)
                    0.28%  6.8800us         2  3.4400us  2.2400us  4.6400us  [CUDA memset]
                    0.18%  4.4800us         1  4.4800us  4.4800us  4.4800us  [CUDA memcpy DtoH]
                    0.11%  2.6560us         1  2.6560us  2.6560us  2.6560us  convert_kernel(unsigned int*, unsigned char*, unsigned int)
     API calls:   97.40%  142.05ms         4  35.513ms  104.83us  140.98ms  cudaFree
                    1.84%  2.6820ms         2  1.3410us  33.133us  2.6489ms  cudaMemcpy
                    0.27%  387.24us         3  129.08us  103.39us  151.38us  cudaMalloc
                    0.14%  205.15us         1  205.15us  205.15us  205.15us  cuDeviceTotalMem
                    0.14%  202.01us        94  2.1490us     442ns  56.735us  cuDeviceGetAttribute
                    0.11%  159.70us         2  79.850us  16.012us  143.69us  cudaLaunch
                    0.05%  76.425us         4  19.106us  4.9700us  51.634us  cudaDeviceSynchronize
                    0.02%  29.115us         2  14.557us  10.949us  18.166us  cudaMemset
                    0.02%  28.463us         1  28.463us  28.463us  28.463us  cuDeviceGetName
                    0.01%  10.258us         7  1.4650us     452ns  6.7510us  cudaSetupArgument
                    0.00%  5.9090us         3  1.9690us     527ns  4.3270us  cuDeviceGetCount
                    0.00%  2.7080us         2  1.3540us     907ns  1.8010us  cudaConfigureCall
                    0.00%  1.6430us         2     821ns     545ns  1.0980us  cuDeviceGet
```

Figure 3: Nvidia profiler result for optimized histogram kernel for m = 1600000

### 2.3.4 Explanation for Optimization

The optimization increased the amount of time spent executing the kernel. I suspect that might be the result of the input array being generated with less repetition than originally anticipated. The contention rate was low and hence the addition of more variables would only serve to increase the execution time.

## 2.4 Kernel

```
1  /*****************************************************************************
2   *cr
3   *cr            (C) Copyright 2010 The Board of Trustees of the
4   *cr                        University of Illinois
5   *cr                         All Rights Reserved
6   *cr
7
↪   *****************************************************************************/
8
9  // Define your kernels in this file you may use more than one
↪   kernel if you
10 // need to
11 __global__ void histogram_kernel(unsigned int* input, unsigned
↪   int* bins,
12     unsigned int num_elements, unsigned int num_bins){
13
14         extern __shared__ unsigned int bins_s[];
15
16         //Shared Memory
17         int thid = threadIdx.x;
18         while(thid < num_bins){
19
20                 bins_s[thid] = 0u;
21                 thid += blockDim.x;
```

7

```
22              }
23              __syncthreads();
24
25
26              //Histogram calculation
27              unsigned int element = blockIdx.x * blockDim.x +
   ↪   threadIdx.x;
28
29              while(element < num_elements){
30
31                      atomicAdd(&(bins_s[input[element]]), 1);
32                      element += blockDim.x * gridDim.x;
33              }
34              __syncthreads();
35
36              //Global Memory
37              thid = threadIdx.x;
38              while(thid < num_bins){
39
40                      atomicAdd(&(bins[thid]), bins_s[thid]);
41                      thid += blockDim.x;
42              }
43      }
44
45
46      __global__ void histogram_kernel_optimized(unsigned int* input,
   ↪   unsigned int* bins,
47          unsigned int num_elements, unsigned int num_bins) {
48
49                      // INSERT CODE HERE
50              extern __shared__ unsigned int bins_s[];
51
52              //Shared memory
53              int thid = threadIdx.x;
54              while ( thid < num_bins){
55
56                      bins_s[thid] = 0u;
57                      thid += blockDim.x;
58              }
59              __syncthreads();
60
61              //Histogram calculation
62              unsigned int element = blockIdx.x * blockDim.x +
   ↪   threadIdx.x;
63              unsigned int accumulator = 0;
64              unsigned int prev_index = 0;
```

```
65
66        while(element < num_elements){
67
68                unsigned int curr_index = input[element];
69
70                if(curr_index != prev_index){
71
72                        atomicAdd(&(bins_s[prev_index]),
                          ↪   accumulator);
73                        accumulator = 1;
74                        prev_index = curr_index;
75
76                }
77
78                else{
79                        accumulator++;
80                }
81                element += blockDim.x * gridDim.x;
82        }
83        if(accumulator > 0){
84                atomicAdd(&(bins_s[prev_index]), accumulator);
85        }
86        __syncthreads();
87
88        //Global memory
89        thid = threadIdx.x;
90        while(thid < num_bins){
91
92                atomicAdd(&(bins[thid]), bins_s[thid]);
93                thid += blockDim.x;
94        }
95
96  }
97
98  __global__ void convert_kernel(unsigned int *bins32, uint8_t
    ↪   *bins8,
99      unsigned int num_bins) {
100
101     // INSERT CODE HERE
102         int thid = blockIdx.x * blockDim.x + threadIdx.x;
103
104         while (thid < num_bins){
105
106                //Use local  register value (avoids copying from
                   ↪   global twice)
107                unsigned int reg_bin = bins32[thid];
```

9

```
108
109                    if(reg_bin > 255){
110                            bins8[thid] = 255u;
111                    }
112
113                    else{
114                            bins8[thid] = (uint8_t) reg_bin;
115                    }
116                    thid += blockDim.x * gridDim.x;
117            }
118
119    }
120
121    /******************************************************************************
122    Setup and invoke your kernel(s) in this function. You may also
       ↪   allocate more
123    GPU memory if you need to
124    ******************************************************************************/
125    void histogram(unsigned int* input, uint8_t* bins, unsigned int
       ↪   num_elements,
126            unsigned int num_bins) {
127
128        // Create 32 bit bins
129        unsigned int *bins32;
130        cudaMalloc((void**)&bins32, num_bins * sizeof(unsigned int));
131        cudaMemset(bins32, 0, num_bins * sizeof(unsigned int));
132
133        // Launch histogram kernel using 32-bit bins
134        dim3 dim_grid, dim_block;
135        dim_block.x = 512; dim_block.y = dim_block.z = 1;
136        dim_grid.x = 30; dim_grid.y = dim_grid.z = 1;
137
138        //Comment out the kernel not used
139        //histogram_kernel<<<dim_grid, dim_block,
       ↪   num_bins*sizeof(unsigned int)>>>
140            // (input, bins32, num_elements, num_bins);
141      histogram_kernel_optimized<<<dim_grid, dim_block,
       ↪   num_bins*sizeof(unsigned int)>>>
142            (input, bins32, num_elements, num_bins);
143
144        // Convert 32-bit bins into 8-bit bins
145        dim_block.x = 512;
146        dim_grid.x = (num_bins - 1)/dim_block.x + 1;
147        convert_kernel<<<dim_grid, dim_block>>>(bins32, bins,
       ↪   num_bins);
148
```

```
149        // Free allocated device memory
150        cudaFree(bins32);
151
152    }
```