Module 7 Lab

Text Histogram

GPU Teaching Kit - Accelerated Computing

OBJECTIVE

The purpose of this lab is to implement an efficient histogram algorithm for an input array of ASCII characters. There are 128 ASCII characters and each character will map into its own bin for a fixed total of 128 bins. The histogram bins will be unsigned 32-bit counters that do not saturate. Use the approach of creating a privitized histogram in shared memory for each thread block, then atomically modifying the global histogram.

INSTRUCTIONS

Edit the code in the code tab to perform the following:

- allocate device memory
- copy host memory to device
- initialize thread block and kernel grid dimensions
- invoke CUDA kernel
- copy results from device to host
- deallocate device memory

Instructions about where to place each part of the code is demarcated by the //@@ comment lines.

LOCAL SETUP INSTRUCTIONS

The most recent version of source code for this lab along with the build-scripts can be found on the <u>Bitbucket repository</u>. A description on how to use the <u>CMake</u> tool in along with how to build the labs for local development found in the <u>README</u> document in the root of the repository.

The executable generated as a result of compiling the lab can be run using the following command:

```
./Histogram_Template -e <expected.raw> \
   -i <input.txt> -o <output.raw> -t integral_vector
```

where <expected.raw> is the expected output, <input.txt> is the input dataset, and <output.raw> is an optional path to store the results. The datasets can be generated using the dataset generator built as part of the compilation process.

QUESTIONS

(1) Were there any difficulties you had with completing the optimization correctly.

ANSWER: The main optimization in the sample solution is for each threadblock to have a privitized histogram in shared memory. Atomic operations on shared memory are much faster than global memory.

(2) Which optimizations gave the most benifit.

ANSWER: Using shared memory requires thread syncrhonization.

(3) For the histogram kernel, how many global memory reads are being performed by your kernel? explain.

ANSWER: Using shared memory privitization provides significant performance improvement.

(4) For the histogram kernel, how many global memory writes are being performed by your kernel? explain.

ANSWER: One read per input element.

(5) For the histogram kernel, how many atomic operations are being performed by your kernel? explain.

ANSWER: Each threadblock does an atomic add into every global memory bin.

(6) Most text files will consist of only letters, numberm and whitespace characters. These 95 characters make up fall between ASCII numbers 32 - 126. What can we say about atomic access contention if more than 95 threads are simultaneously trying to atomically increment a private histogram?

ANSWER: One atomic operation per input element into shared memory, then NUM_BINS atomic operation per thread block to accumulate the results into the global bins.

CODE TEMPLATE

The following code is suggested as a starting point for students. The code handles the import and export as well as the checking of the solution. Students are expected to insert their code is the sections demarcated with //@@. Students expected the other code unchanged. The tutorial page describes the functionality of the wb* methods.

```
#include <wb.h>
   #define NUM BINS 4096
  #define CUDA_CHECK(ans)
    { gpuAssert((ans), __FILE__, __LINE__); }
   inline void gpuAssert(cudaError_t code, const char *file, int line,
                        bool abort = true) {
     if (code != cudaSuccess) {
       fprintf(stderr, "GPUassert: %s %s %d\n", cudaGetErrorString(code),
               file, line);
       if (abort)
12
         exit(code);
13
     }
   }
  int main(int argc, char *argv[]) {
     wbArg_t args;
     int inputLength;
     unsigned int *hostInput;
     unsigned int *hostBins;
     unsigned int *deviceInput;
     unsigned int *deviceBins;
     args = wbArg_read(argc, argv);
25
     wbTime_start(Generic, "Importing data and creating memory on host");
     hostInput = (unsigned int *)wbImport(wbArg_getInputFile(args, 0),
                                         &inputLength, "Integer");
     hostBins = (unsigned int *)malloc(NUM_BINS * sizeof(unsigned int));
     wbTime_stop(Generic, "Importing data and creating memory on host");
     wbLog(TRACE, "The input length is ", inputLength);
33
     wbLog(TRACE, "The number of bins is ", NUM_BINS);
34
     wbTime_start(GPU, "Allocating GPU memory.");
     //@@ Allocate GPU memory here
37
     CUDA_CHECK(cudaDeviceSynchronize());
     wbTime_stop(GPU, "Allocating GPU memory.");
     wbTime_start(GPU, "Copying input memory to the GPU.");
     //@@ Copy memory to the GPU here
42
     CUDA_CHECK(cudaDeviceSynchronize());
43
     wbTime_stop(GPU, "Copying input memory to the GPU.");
     // Launch kernel
     // -----
     wbLog(TRACE, "Launching kernel");
     wbTime_start(Compute, "Performing CUDA computation");
     //@@ Perform kernel computation here
     wbTime_stop(Compute, "Performing CUDA computation");
     wbTime_start(Copy, "Copying output memory to the CPU");
```

\

```
//@@ Copy the GPU memory back to the CPU here
    CUDA_CHECK(cudaDeviceSynchronize());
55
    wbTime_stop(Copy, "Copying output memory to the CPU");
    wbTime_start(GPU, "Freeing GPU Memory");
58
    //@@ Free the GPU memory here
    wbTime_stop(GPU, "Freeing GPU Memory");
    // Verify correctness
    // -----
    wbSolution(args, hostBins, NUM_BINS);
    free(hostBins);
    free(hostInput);
    return 0;
  }
69
```

CODE SOLUTION

The following is a possible implementation of the lab. This solution is intended for use only by the teaching staff and should not be distributed to students.

```
#include <wb.h>
   #define NUM_BINS 128
   #define CUDA_CHECK(ans)
     { gpuAssert((ans), __FILE__, __LINE__); }
   inline void gpuAssert(cudaError_t code, const char *file, int line,
                          bool abort = true) {
     if (code != cudaSuccess) {
       fprintf(stderr, "GPUassert: %s %s %d\n", cudaGetErrorString(code),
                file, line);
11
       if (abort)
12
         exit(code);
13
     }
   }
15
   __global__ void histogram_kernel(const char *input, unsigned int *bins,
17
                                     unsigned int num_elements,
                                     unsigned int num_bins) {
20
     unsigned int tid = blockIdx.x * blockDim.x + threadIdx.x;
21
     // Privatized bins
     extern __shared__ unsigned int bins_s[];
24
     for (unsigned int binIdx = threadIdx.x; binIdx < num_bins;</pre>
25
          binIdx += blockDim.x) {
       bins_s[binIdx] = 0;
28
     __syncthreads();
```

```
30
     // Histogram
31
     for (unsigned int i = tid; i < num_elements;</pre>
          i += blockDim.x * gridDim.x) {
       atomicAdd(&(bins_s[(unsigned int)input[i]]), 1);
34
35
     }
     __syncthreads();
     // Commit to global memory
     for (unsigned int binIdx = threadIdx.x; binIdx < num_bins;</pre>
          binIdx += blockDim.x) {
       atomicAdd(&(bins[binIdx]), bins_s[binIdx]);
42
   }
43
   void histogram(const char *input, unsigned int *bins,
                   unsigned int num_elements, unsigned int num_bins) {
46
47
     // zero out bins
     CUDA_CHECK(cudaMemset(bins, 0, num_bins * sizeof(unsigned int)));
     // Launch histogram kernel on the bins
51
       dim3 blockDim(256), gridDim(30);
       histogram_kernel<<<gridDim, blockDim,
                           num_bins * sizeof(unsigned int)>>>(
54
           input, bins, num_elements, num_bins);
55
       CUDA_CHECK(cudaGetLastError());
       CUDA_CHECK(cudaDeviceSynchronize());
     }
58
   }
59
   int main(int argc, char *argv[]) {
     wbArg_t args;
62
     int inputLength;
63
     char *hostInput;
     unsigned int *hostBins;
     char *deviceInput;
     unsigned int *deviceBins;
67
     args = wbArg_read(argc, argv);
     wbTime_start(Generic, "Importing data and creating memory on host");
71
     hostInput =
          (char *)wbImport(wbArg_getInputFile(args, 0), &inputLength, "Text");
     hostBins = (unsigned int *)malloc(NUM_BINS * sizeof(unsigned int));
74
     wbTime_stop(Generic, "Importing data and creating memory on host");
     wbLog(TRACE, "The input length is ", inputLength);
77
     wbLog(TRACE, "The number of bins is ", NUM_BINS);
     wbTime_start(GPU, "Allocating GPU memory.");
     //@@ Allocate GPU memory here
81
     CUDA_CHECK(cudaMalloc((void **)&deviceInput, inputLength));
82
```

```
CUDA_CHECK(
83
         cudaMalloc((void **)&deviceBins, NUM_BINS * sizeof(unsigned int)));
84
     CUDA_CHECK(cudaDeviceSynchronize());
85
     wbTime_stop(GPU, "Allocating GPU memory.");
     wbTime_start(GPU, "Copying input memory to the GPU.");
     //@@ Copy memory to the GPU here
     CUDA_CHECK(cudaMemcpy(deviceInput, hostInput, inputLength,
                          cudaMemcpyHostToDevice));
     CUDA_CHECK(cudaDeviceSynchronize());
     wbTime_stop(GPU, "Copying input memory to the GPU.");
     // Launch kernel
     // -----
     wbLog(TRACE, "Launching kernel");
     wbTime_start(Compute, "Performing CUDA computation");
     // @@ Insert code here
     histogram(deviceInput, deviceBins, inputLength, NUM_BINS);
     wbTime_stop(Compute, "Performing CUDA computation");
     wbTime_start(Copy, "Copying output memory to the CPU");
103
     //@@ Copy the GPU memory back to the CPU here
104
     CUDA_CHECK(cudaMemcpy(hostBins, deviceBins,
                          NUM_BINS * sizeof(unsigned int),
106
                          cudaMemcpyDeviceToHost));
107
     CUDA_CHECK(cudaDeviceSynchronize());
     wbTime_stop(Copy, "Copying output memory to the CPU");
     wbTime_start(GPU, "Freeing GPU Memory");
111
     //@@ Free the GPU memory here
112
     CUDA_CHECK(cudaFree(deviceInput));
     CUDA_CHECK(cudaFree(deviceBins));
     wbTime_stop(GPU, "Freeing GPU Memory");
115
     // Verify correctness
     // -----
     wbSolution(args, hostBins, NUM_BINS);
119
120
     free(hostBins);
     free(hostInput);
122
     return 0;
123
   }
124
```