

# Module 9 Lab

## Reduction

*GPU Teaching Kit – Accelerated Computing*

### OBJECTIVE

Implement a kernel that performs reduction of a 1D list. The reduction should give the sum of the list. You should implement the improved kernel discussed in week 4. Your kernel should be able to handle input lists of arbitrary length. However, for simplicity, you can assume that the input list will be at most  $2048 \times 65535$  elements so that it can be handled by only one kernel launch. The boundary condition can be handled by filling “identity value (0 for sum)” into the shared memory of the last block when the length is not a multiple of the thread block size. Further assume that the reduction sums of each section generated by individual blocks will be summed up by the CPU. Prerequisites

### PREREQUISITES

Before starting this lab, make sure that:

- You have completed week 4 lecture videos

### INSTRUCTION

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
- implement the improved reduction routine
- use shared memory to reduce the number of global accesses, handle the boundary conditions in when loading input list elements into the shared memory
- implement a CPU loop to perform final reduction based on the sums of sections generated by the thread blocks

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

## QUESTIONS

- (1) Name 3 applications of reduction.

ANSWER: **Determining min/max/average of a data series.**

- (2) How many floating operations are being performed in your reduction kernel? explain.

ANSWER: **Input length - 1.**

- (3) How many global memory reads are being performed by your kernel? explain.

ANSWER: **numInputs. Each input value is copied to shared memory.**

- (4) How many global memory writes are being performed by your kernel? explain.

ANSWER: **numInputs / (BLOCK\_SIZE \* 2). Only final outputs are written to global memory.**

- (5) What is the minimum, maximum, and average number of real operations that a thread will perform? Real operations are those that directly contribute to the final reduction value.

ANSWER: **Minimum is 1. Max is  $\log_2(\text{BLOCK\_SIZE})+1$ . Average is  $(\text{sum from } i = 0 \text{ to } \log_2(\text{BLOCK\_SIZE}) \text{ of } 2^i)/\text{BLOCK\_SIZE}$ .**

- (6) How many times does a single thread block synchronize to reduce its portion of the array to a single value?

ANSWER:  **$\log_2(\text{BLOCK\_SIZE})+1$**

- (7) Describe what optimizations were performed to your kernel to achieve a performance speedup.

ANSWER: **Active threads are adjacent so as high a fraction of each warp is active as possible.**

- (8) Describe what further optimizations can be implemented to your kernel to achieve a performance speedup.

ANSWER: **Use atomic operations once the reduction tree is small enough to accumulate a single final value.**

- (9) Suppose the input is greater than  $2048 \times 65535$ , what modifications are needed to your kernel?

ANSWER: **The work could be further coarsened, so a thread does more than two inputs in the beginning.**

- (10) Suppose you want to scan using a binary operator that's not commutative, can you use a parallel reduction for that?

ANSWER: **To parallelize a reduction the operator must be associative, not necessarily commutative.**

(11) Is it possible to get different results from running the serial version and parallel version of reduction? explain.

ANSWER: **The order of floating-point operations in the serial and parallel version is different, so it is possible.**

## 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 in the sections demarcated with `//@@`. Students expected the other code unchanged. The tutorial page describes the functionality of the `wb*` methods.

```

1
2 // Given a list (lst) of length n
3 // Output its sum = lst[0] + lst[1] + ... + lst[n-1];
4
5 #include <wb.h>
6
7 #define BLOCK_SIZE 512 //@@ You can change this
8
9 #define wbCheck(stmt)                                     \
10     do {                                                 \
11         cudaError_t err = stmt;                         \
12         if (err != cudaSuccess) {                       \
13             wbLog(ERROR, "Failed to run stmt ", #stmt); \
14             return -1;                                   \
15         }                                                 \
16     } while (0)
17
18 __global__ void total(float *input, float *output, int len) {
19     //@@ Load a segment of the input vector into shared memory
20
21     //@@ Traverse the reduction tree
22
23     //@@ Write the computed sum of the block to the output vector at the
24     //@@ correct index
25 }
26
27 int main(int argc, char **argv) {
28     int ii;
29     wbArg_t args;
30     float *hostInput; // The input 1D list
31     float *hostOutput; // The output list
32     float *deviceInput;
33     float *deviceOutput;
34     int numInputElements; // number of elements in the input list
35     int numOutputElements; // number of elements in the output list

```

```

36
37 args = wbArg_read(argc, argv);
38
39 wbTime_start(Generic, "Importing data and creating memory on host");
40 hostInput =
41     (float *)wbImport(wbArg_getInputFile(args, 0), &numInputElements);
42
43 numOutputElements = numInputElements / (BLOCK_SIZE << 1);
44 if (numInputElements % (BLOCK_SIZE << 1)) {
45     numOutputElements++;
46 }
47 hostOutput = (float *)malloc(numOutputElements * sizeof(float));
48
49 wbTime_stop(Generic, "Importing data and creating memory on host");
50
51 wbLog	TRACE, "The number of input elements in the input is ",
52         numInputElements);
53 wbLog	TRACE, "The number of output elements in the input is ",
54         numOutputElements);
55
56 wbTime_start(GPU, "Allocating GPU memory.");
57 ///@@ Allocate GPU memory here
58 wbTime_stop(GPU, "Allocating GPU memory.");
59
60 wbTime_start(GPU, "Copying input memory to the GPU.");
61 ///@@ Copy memory to the GPU here
62 wbTime_stop(GPU, "Copying input memory to the GPU.");
63 ///@@ Initialize the grid and block dimensions here
64 wbTime_start(Compute, "Performing CUDA computation");
65 ///@@ Launch the GPU Kernel here
66 cudaDeviceSynchronize();
67 wbTime_stop(Compute, "Performing CUDA computation");
68
69 wbTime_start(Copy, "Copying output memory to the CPU");
70 ///@@ Copy the GPU memory back to the CPU here
71 wbTime_stop(Copy, "Copying output memory to the CPU");
72
73 /*****
74  * Reduce output vector on the host
75  * NOTE: One could also perform the reduction of the output vector
76  * recursively and support any size input. For simplicity, we do not
77  * require that for this lab.
78 *****/
79
80 wbTime_start(GPU, "Freeing GPU Memory");
81 ///@@ Free the GPU memory here
82 wbTime_stop(GPU, "Freeing GPU Memory");
83
84 wbSolution(args, hostOutput, 1);
85
86 free(hostInput);
87 free(hostOutput);
88

```

```

89     return 0;
90 }

```

## 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.

```

1  // Given a list (lst) of length n
2  // Output its sum = lst[0] + lst[1] + ... + lst[n-1];
3
4  #include <wb.h>
5
6  #define BLOCK_SIZE 512 ///@ You can change this
7
8  #define wbCheck(stmt) \
9      do { \
10         cudaError_t err = stmt; \
11         if (err != cudaSuccess) { \
12             wbLog(ERROR, "Failed to run stmt ", #stmt); \
13             return -1; \
14         } \
15     } while (0)
16
17  __global__ void total(float *input, float *output, int len) {
18     ///@ Load a segment of the input vector into shared memory
19     __shared__ float partialSum[2 * BLOCK_SIZE];
20     unsigned int t = threadIdx.x, start = 2 * blockIdx.x * BLOCK_SIZE;
21     if (start + t < len)
22         partialSum[t] = input[start + t];
23     else
24         partialSum[t] = 0;
25     if (start + BLOCK_SIZE + t < len)
26         partialSum[BLOCK_SIZE + t] = input[start + BLOCK_SIZE + t];
27     else
28         partialSum[BLOCK_SIZE + t] = 0;
29     ///@ Traverse the reduction tree
30     for (unsigned int stride = BLOCK_SIZE; stride >= 1; stride >>= 1) {
31         __syncthreads();
32         if (t < stride)
33             partialSum[t] += partialSum[t + stride];
34     }
35     ///@ Write the computed sum of the block to the output vector at the
36     ///@ correct index
37     if (t == 0)
38         output[blockIdx.x] = partialSum[0];
39 }
40
41 int main(int argc, char **argv) {
42     int ii;
43     wbArg_t args;

```

```

44 float *hostInput; // The input 1D list
45 float *hostOutput; // The output list
46 float *deviceInput;
47 float *deviceOutput;
48 int numInputElements; // number of elements in the input list
49 int numOutputElements; // number of elements in the output list
50
51 args = wbArg_read(argc, argv);
52
53 wbTime_start(Generic, "Importing data and creating memory on host");
54 hostInput =
55     (float *)wbImport(wbArg_getInputFile(args, 0), &numInputElements);
56
57 numOutputElements = numInputElements / (BLOCK_SIZE << 1);
58 if (numInputElements % (BLOCK_SIZE << 1)) {
59     numOutputElements++;
60 }
61 hostOutput = (float *)malloc(numOutputElements * sizeof(float));
62
63 wbTime_stop(Generic, "Importing data and creating memory on host");
64
65 wbLog	TRACE, "The number of input elements in the input is ",
66         numInputElements);
67 wbLog	TRACE, "The number of output elements in the input is ",
68         numOutputElements);
69
70 wbTime_start(GPU, "Allocating GPU memory.");
71 //@@ Allocate GPU memory here
72 cudaMalloc(&deviceInput, sizeof(float) * numInputElements);
73 cudaMalloc(&deviceOutput, sizeof(float) * numOutputElements);
74
75 wbTime_stop(GPU, "Allocating GPU memory.");
76
77 wbTime_start(GPU, "Copying input memory to the GPU.");
78 //@@ Copy memory to the GPU here
79 cudaMemcpy(deviceInput, hostInput, sizeof(float) * numInputElements,
80             cudaMemcpyHostToDevice);
81
82 wbTime_stop(GPU, "Copying input memory to the GPU.");
83 //@@ Initialize the grid and block dimensions here
84 dim3 dimGrid(numOutputElements, 1, 1);
85 dim3 dimBlock(BLOCK_SIZE, 1, 1);
86
87 wbTime_start(Compute, "Performing CUDA computation");
88 //@@ Launch the GPU Kernel here
89 total<<<dimGrid, dimBlock>>>(deviceInput, deviceOutput,
90                               numInputElements);
91
92 cudaDeviceSynchronize();
93 wbTime_stop(Compute, "Performing CUDA computation");
94
95 wbTime_start(Copy, "Copying output memory to the CPU");
96 //@@ Copy the GPU memory back to the CPU here

```

```

97     cudaMemcpy(hostOutput, deviceOutput, sizeof(float) * numOutputElements,
98               cudaMemcpyDeviceToHost);
99
100    wbTime_stop(Copy, "Copying output memory to the CPU");
101
102    /*****
103     * Reduce output vector on the host
104     * NOTE: One could also perform the reduction of the output vector
105     * recursively and support any size input. For simplicity, we do not
106     * require that for this lab.
107     *****/
108    for (ii = 1; ii < numOutputElements; ii++) {
109        hostOutput[0] += hostOutput[ii];
110    }
111
112    wbTime_start(GPU, "Freeing GPU Memory");
113    //@@ Free the GPU memory here
114    cudaFree(deviceInput);
115    cudaFree(deviceOutput);
116
117    wbTime_stop(GPU, "Freeing GPU Memory");
118
119    wbSolution(args, hostOutput, 1);
120
121    free(hostInput);
122    free(hostOutput);
123
124    return 0;
125 }

```