# List Reduction Grade

## Grade Summary (**History**) (/grade/history/5984)

| | |
|---|---|
| Created: | less than a minute ago (2022-04-11 03:16:06 +0000 UTC) |
| Total Score: | 100 out of 100 points |
| Coding Score: | 100 out of 100 points |
| Questions Score: | 0 |

## Program Code

```
1   // MP Reduction
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        wbLog(ERROR, "Got CUDA error ...  ", cudaGetErrorString(err));
15        return -1;
16      }
17    } while (0)
18
19  __global__ void reduction(float *input, float *output, int len) {
20    //@@ Load a segment of the input vector into shared memory
21    __shared__ float sum[BLOCK_SIZE * 2];
```

```
22
23    int thread_x  = threadIdx.x;
24
25    int load_in = 2*(BLOCK_SIZE * blockIdx.x) + thread_x;
26
27    sum[thread_x] = load_in < len ? input[load_in] : 0.0;
28    if (load_in + BLOCK_SIZE < len)
29        sum[BLOCK_SIZE + thread_x] = input[load_in + BLOCK_SIZE];
30    else
31        sum[BLOCK_SIZE + thread_x] = 0.0;
32
33    //@@ Traverse the reduction tree
34    for (int travel_tree = BLOCK_SIZE; travel_tree >= 1; travel_tree >>= 1)
35        __syncthreads();
36        if (thread_x < travel_tree)
37            sum[thread_x] += sum[thread_x + travel_tree];
38    }
39
40    //@@ Write the computed sum of the block to the output vector at the
41    //@@ correct index
42    if (thread_x == 0){
43      output[blockIdx.x] = sum[0];
44    }
45 }
46
47 int main(int argc, char **argv) {
48    int ii;
49    wbArg_t args;
50    float *hostInput;  // The input 1D list
51    float *hostOutput; // The output list
52    float *deviceInput;
53    float *deviceOutput;
54    int numInputElements;  // number of elements in the input list
55    int numOutputElements; // number of elements in the output list
56
57    args = wbArg_read(argc, argv);
58
59    wbTime_start(Generic, "Importing data and creating memory on host");
60    hostInput = (float *)wbImport(wbArg_getInputFile(args, 0), &numInputEle
61
62    numOutputElements = numInputElements / (BLOCK_SIZE << 1);
63    if (numInputElements % (BLOCK_SIZE << 1)){
64      numOutputElements++;
65    }
66
```

```
67    //2 variables to help with sizes
68    int Size_of_Input = numInputElements * sizeof(float);
69    int Size_of_Output = numOutputElements * sizeof(float);
70
71    hostOutput = (float *)malloc(numOutputElements * sizeof(float));
72
73    wbTime_stop(Generic, "Importing data and creating memory on host");
74
75    wbLog(TRACE, "The number of input elements in the input is ", numInputE
76    wbLog(TRACE, "The number of output elements in the input is ", numOutpu
77
78    wbTime_start(GPU, "Allocating GPU memory.");
79    //@@ Allocate GPU memory here
80    cudaMalloc(&deviceInput, Size_of_Input);
81    cudaMalloc(&deviceOutput, Size_of_Output);
82
83    wbTime_stop(GPU, "Allocating GPU memory.");
84
85    wbTime_start(GPU, "Copying input memory to the GPU.");
86    //@@ Copy memory to the GPU here
87    cudaMemcpy(deviceInput, hostInput, Size_of_Input, cudaMemcpyHostToDevic
88
89    wbTime_stop(GPU, "Copying input memory to the GPU.");
90    //@@ Initialize the grid and block dimensions here
91    dim3 dimGrid(numOutputElements, 1, 1);
92    dim3 dimBlock(BLOCK_SIZE, 1, 1);
93
94    wbTime_start(Compute, "Performing CUDA computation");
95    //@@ Launch the GPU Kernel here
96    reduction<<<dimGrid, dimBlock>>>(deviceInput, deviceOutput, numInputEle
97
98    cudaDeviceSynchronize();
99    wbTime_stop(Compute, "Performing CUDA computation");
100
101   wbTime_start(Copy, "Copying output memory to the CPU");
102   //@@ Copy the GPU memory back to the CPU here
103   cudaMemcpy(hostOutput, deviceOutput, Size_of_Output, cudaMemcpyDeviceTo
104
105   wbTime_stop(Copy, "Copying output memory to the CPU");
106
107   /******************************************************************
108    * Reduce output vector on the host
109    * NOTE: One could also perform the reduction of the output vector
110    * recursively and support any size input. For simplicity, we do not
111    * require that for this lab.
```

```
112        *****************************************************************/
113    for (ii = 1; ii < numOutputElements; ii++) {
114      hostOutput[0] += hostOutput[ii];
115    }
116
117    wbTime_start(GPU, "Freeing GPU Memory");
118    //@@ Free the GPU memory here
119    cudaFree(deviceInput);
120    cudaFree(deviceOutput);
121
122    wbTime_stop(GPU, "Freeing GPU Memory");
123
124    wbSolution(args, hostOutput, 1);
125
126    free(hostInput);
127    free(hostOutput);
128
129    return 0;
130 }
131
```

Designed and architected by Abdul Dakkak (https://www.dakkak.dev/).