# ❮ (/mp/5985?show=code)  List Scan Attempt

## Attempt Summary

### Submit Attempt for Grading

---

### Remember to answer the questions before clicking.

| | |
|---|---|
| Dataset Id: | 5 |
| Created: | less than a minute ago (2022-04-18 04:22:14 +0000 UTC) |
| Status: | Correct solution for this dataset. |

## Timer Output

| Kind | Location | Time (ms) | Message |
|---|---|---|---|
| **Generic** | main.cu::147 | 41.303532 | Importing data and creating memory on host |
| **GPU** | main.cu::154 | 1.119895 | Allocating GPU memory. |
| **GPU** | main.cu::161 | 0.028025 | Clearing output memory. |
| **GPU** | main.cu::165 | 0.048714 | Copying input memory to the GPU. |
| **Compute** | main.cu::173 | 0.072838 | Performing CUDA computation |
| **Copy** | main.cu::187 | 0.039268 | Copying output memory to the CPU |
| **GPU** | main.cu::191 | 0.167033 | Freeing GPU Memory |

# Logger Output

| Level | Location | Message |
|-------|----------|---------|
| **Trace** | main::152 | The number of input elements in the input is 9010 |

# Program Code

```
 1  /*
 2  Dataset Id: 5
 3  Created:  less than a minute ago
 4  Status: Correct solution for this dataset.
 5  Timer Output
 6  Kind   Location   Time (ms) Message
 7  Generic main.cu::127   54.583396  Importing data and creating memory on ho
 8  GPU  main.cu::134   1.679192   Allocating GPU memory.
 9  GPU  main.cu::141   0.043779   Clearing output memory.
10  GPU  main.cu::145   0.064188   Copying input memory to the GPU.
11  Compute main.cu::153   0.092417   Performing CUDA computation
12  Copy   main.cu::167   0.05866  Copying output memory to the CPU
13  GPU  main.cu::171   0.2507   Freeing GPU Memory
14  Logger Output
15  Level  Location   Message
16  Trace  main::132  The number of input elements in the input is 9010
17
18  */
19
20
21  // MP Scan
22  // Given a list (lst) of length n
23  // Output its prefix sum = {lst[0], lst[0] + lst[1], lst[0] + lst[1] + ..
24  // +
25  // lst[n-1]}
26
27  #include <wb.h>
28
29  #define BLOCK_SIZE 512 //@@ You can change this
30
31  #define wbCheck(stmt)
32      do {
```

```
33      cudaError_t err = stmt;
34      if (err != cudaSuccess) {
35        wbLog(ERROR, "Failed to run stmt ", #stmt);
36        wbLog(ERROR, "Got CUDA error ...  ", cudaGetErrorString(err));
37        return -1;
38      }
39    } while (0)
40
41  __global__ void scan(float *input, float *output, int len, int flag) {
42    //@@ Modify the body of this function to complete the functionality of
43    //@@ the scan on the device
44    //@@ You may need multiple kernel calls; write your kernels before this
45    //@@ function and call them from here
46    __shared__ float block_arrray_scan[2 * BLOCK_SIZE];
47
48    int loop_index = 0;
49    int stride = 1;
50    if (!flag){
51      loop_index = (2 * blockIdx.x * blockDim.x) + threadIdx.x;
52      stride = blockDim.x;
53    }
54    else{
55      loop_index = (threadIdx.x + 1) * (2 * blockDim.x) - 1;
56      stride = 2 * blockDim.x;
57    }
58
59    int storeIndex = (2 * blockIdx.x * blockDim.x) + threadIdx.x;
60
61    //data input
62    if (loop_index < len){
63      block_arrray_scan[threadIdx.x] = input[loop_index];
64    }
65    else{
66      block_arrray_scan[threadIdx.x] = 0;
67    }
68    if (loop_index + stride < len){
69      block_arrray_scan[threadIdx.x + blockDim.x] = input[loop_index + stri
70    }
71    else{
72      block_arrray_scan[threadIdx.x + blockDim.x] = 0;
73    }
74
75    //First Step: Reduction
76    for (int stride = 1; stride <= (2 * BLOCK_SIZE); stride *= 2) {
77      __syncthreads();
```

```
 78
 79        int loop_index = (threadIdx.x + 1) * 2 * stride - 1;
 80
 81      if ((loop_index < 2 * BLOCK_SIZE) && ((loop_index - stride) >= 0)){
 82          block_arrray_scan[loop_index] += block_arrray_scan[loop_index - s
 83      }
 84    }
 85
 86    //Use Distribution Tree method after Scanning
 87    for (int stride = 2 * BLOCK_SIZE / 4; stride > 0; stride /= 2) {
 88      __syncthreads();
 89
 90      int loop_index = (threadIdx.x + 1) * 2 * stride - 1;
 91      if ((loop_index + stride) < 2 * BLOCK_SIZE){
 92          block_arrray_scan[loop_index + stride] += block_arrray_scan[loop_
 93      }
 94    }
 95
 96    __syncthreads();
 97    if (storeIndex < len){
 98      output[storeIndex] = block_arrray_scan[threadIdx.x];
 99    }
100
101    if (storeIndex + blockDim.x < len){
102      output[storeIndex + blockDim.x] = block_arrray_scan[threadIdx.x + blo
103    }
104  }
105
106  __global__ void add(float *input, float *output, float *array_sum, int le
107    __shared__ float move_loop;
108
109    int loop_index = threadIdx.x + (2 * blockIdx.x * blockDim.x);
110
111    if (threadIdx.x == 0){
112      if (blockIdx.x == 0){
113          move_loop = 0;
114      }
115      else{
116          move_loop = array_sum[blockIdx.x - 1];
117      }
118    }
119
120    __syncthreads();
121
122    if (loop_index < len){
```

```
123        output[loop_index] = input[loop_index] + move_loop;
124      }
125      if (loop_index + blockDim.x < len){
126        output[loop_index + blockDim.x] = input[loop_index + blockDim.x] + mo
127      }
128  }
129
130
131  int main(int argc, char **argv) {
132    wbArg_t args;
133    float *hostInput;  // The input 1D list
134    float *hostOutput; // The output list
135    float *deviceInput;
136    float *deviceOutput;
137    int numElements; // number of elements in the list
138
139    //Additional Variables
140    //store temporary results from scanning
141    //store block summations from scanning
142    float *device_temporary_value;
143    float *scanned_dev_temp_val;
144
145    args = wbArg_read(argc, argv);
146
147    wbTime_start(Generic, "Importing data and creating memory on host");
148    hostInput = (float *)wbImport(wbArg_getInputFile(args, 0), &numElements
149    hostOutput = (float *)malloc(numElements * sizeof(float));
150    wbTime_stop(Generic, "Importing data and creating memory on host");
151
152    wbLog(TRACE, "The number of input elements in the input is ", numElemen
153
154    wbTime_start(GPU, "Allocating GPU memory.");
155    wbCheck(cudaMalloc((void **)&deviceInput, numElements * sizeof(float)))
156    wbCheck(cudaMalloc((void **)&deviceOutput, numElements * sizeof(float))
157    wbCheck(cudaMalloc((void **)&device_temporary_value, numElements * size
158    wbCheck(cudaMalloc((void **)&scanned_dev_temp_val, 2 * BLOCK_SIZE * siz
159    wbTime_stop(GPU, "Allocating GPU memory.");
160
161    wbTime_start(GPU, "Clearing output memory.");
162    wbCheck(cudaMemset(deviceOutput, 0, numElements * sizeof(float)));
163    wbTime_stop(GPU, "Clearing output memory.");
164
165    wbTime_start(GPU, "Copying input memory to the GPU.");
166    wbCheck(cudaMemcpy(deviceInput, hostInput, numElements * sizeof(float),
167    wbTime_stop(GPU, "Copying input memory to the GPU.");
```

```
168
169      //@@ Initialize the grid and block dimensions here
170      dim3 dimGrid(ceil(numElements/(BLOCK_SIZE * 2.0)),    1, 1);
171      dim3 dimBlock(BLOCK_SIZE, 1, 1);
172
173      wbTime_start(Compute, "Performing CUDA computation");
174      //@@ Modify this to complete the functionality of the scan
175      //@@ on the deivce
176
177      //Here I store the temporary value in deviceOutput
178      scan<<<dimGrid, dimBlock>>>(deviceInput, device_temporary_value, numEle
179
180      dim3 postScanGrid(1, 1, 1);
181      scan<<<postScanGrid, dimBlock>>>(device_temporary_value, scanned_dev_te
182      add<<<dimGrid, dimBlock>>>(device_temporary_value, deviceOutput, scanne
183
184      cudaDeviceSynchronize();
185      wbTime_stop(Compute, "Performing CUDA computation");
186
187      wbTime_start(Copy, "Copying output memory to the CPU");
188      wbCheck(cudaMemcpy(hostOutput, deviceOutput, numElements * sizeof(float
189      wbTime_stop(Copy, "Copying output memory to the CPU");
190
191      wbTime_start(GPU, "Freeing GPU Memory");
192      cudaFree(deviceInput);
193      cudaFree(deviceOutput);
194      cudaFree(device_temporary_value);
195      cudaFree(scanned_dev_temp_val);
196      wbTime_stop(GPU, "Freeing GPU Memory");
197
198      wbSolution(args, hostOutput, numElements);
199
200      free(hostInput);
201      free(hostOutput);
202
203      return 0;
204  }
```