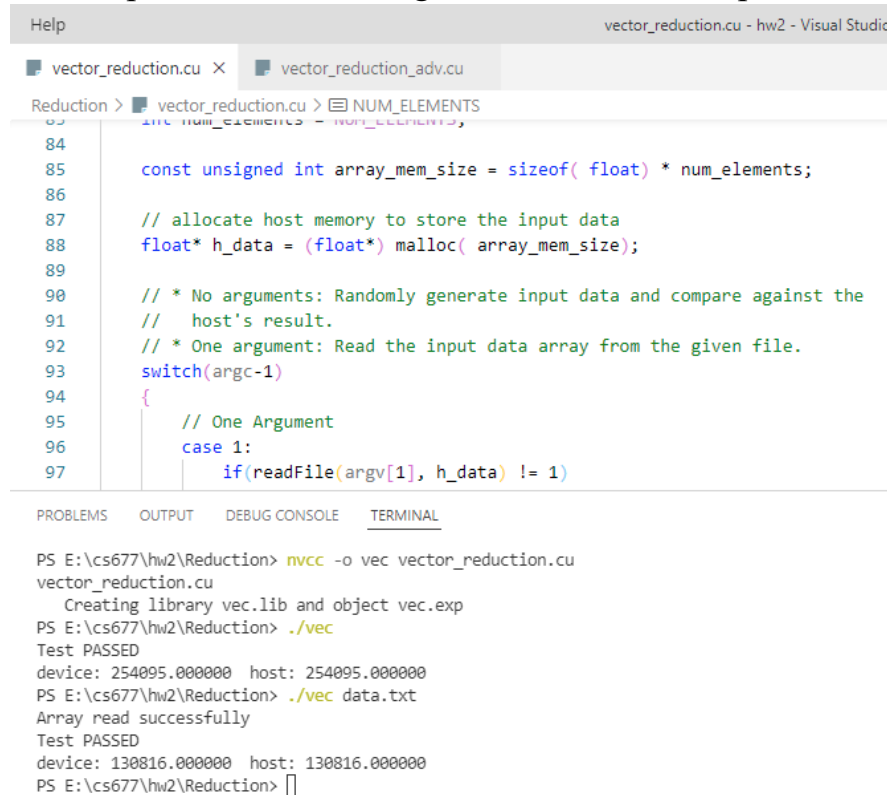


1. Problem1

- 2 times synchronize
- Maximum number of “real” operations: 8; Minimum number of “real” operations: 0; Average number of “real” operations: 4.



```
Help vector_reduction.cu - hw2 - Visual Studio
vector_reduction.cu X vector_reduction_adv.cu
Reduction > vector_reduction.cu > NUM_ELEMENTS
83 int num_elements = NUM_ELEMENTS;
84
85 const unsigned int array_mem_size = sizeof( float ) * num_elements;
86
87 // allocate host memory to store the input data
88 float* h_data = (float*) malloc( array_mem_size );
89
90 // * No arguments: Randomly generate input data and compare against the
91 //   host's result.
92 // * One argument: Read the input data array from the given file.
93 switch( argc-1 )
94 {
95     // One Argument
96     case 1:
97         if( readFile( argv[1], h_data ) != 1 )
98             return 1;
99 }
100
101 // Run the kernel
102 cudaError_t err = cudaRun( h_data, num_elements );
103 if( err != cudaSuccess )
104     return 1;
105
106 // Print the result
107 printf( "Result: %f\n", h_data[0] );
108 return 0;
109 }
110
111 // Main function
112 int main( int argc, char* argv[] )
113 {
114     return Reduction( argc, argv );
115 }
```

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL

```
PS E:\cs677\hw2\Reduction> nvcc -o vec vector_reduction.cu
vector_reduction.cu
Creating library vec.lib and object vec.exp
PS E:\cs677\hw2\Reduction> ./vec
Test PASSED
device: 254095.000000 host: 254095.000000
PS E:\cs677\hw2\Reduction> ./vec data.txt
Array read successfully
Test PASSED
device: 130816.000000 host: 130816.000000
PS E:\cs677\hw2\Reduction> []
```

2. Problem2

- I add another kernel in ‘vector_reduction_kernel.cu’ named reductionLarge() which take case of vector size greater than 512.
- In reductionLarge() function,

```
int temp = arr_size / 512;
int remainder = arr_size % 512;
int division = remainder == 0 ? temp : temp+1;
```

it can find out how many times this vector’s size is 512. Then, I use a for-loop to load data to threads and calculate the sub-value. After the calculation, storing them into a shared memory.

Helpvector_reduction_kernel.cu - hw2 - V

vector_reduction.cuvector_reduction_adv.cuvector_reduction_kernel.cu X

Reduction > vector_reduction_kernel.cu > reductionLarge(float *, int)

```
47     __shared__ float partialSum[1024];
48     unsigned int t = threadIdx.x;
49     unsigned int start = 2*blockIdx.x*blockDim.x;
50     partialSum[t] = g_data[t+start];
51     partialSum[blockDim.x+t] = g_data[start+blockDim.x+t];
52
53     for(unsigned int stride = blockDim.x/2; stride >= 1; stride >>= 1)
```

PROBLEMSOUTPUTDEBUG CONSOLETERMINAL

```
PS E:\cs677\hw2\Reduction> nvcc -o vec2 vector_reduction_adv.cu
    Creating library vec2.lib and object vec2.exp
PS E:\cs677\hw2\Reduction> ./vec2
length of vector: 512
Test PASSED
PS E:\cs677\hw2\Reduction> ./vec2 0
length of vector: 512
Test PASSED
device: 254095.000000 host: 254095.000000
PS E:\cs677\hw2\Reduction> ./vec2 1
length of vector: 512
Test PASSED
device: 254095.000000 host: 254095.000000
PS E:\cs677\hw2\Reduction> ./vec2 0 data.txt
Array read successfully
length of vector: 512
Test PASSED
device: 130816.000000 host: 130816.000000
PS E:\cs677\hw2\Reduction> ./vec2 1 data.txt
Array read successfully
length of vector: 512
Test PASSED
device: 130816.000000 host: 130816.000000
PS E:\cs677\hw2\Reduction> █
```

c.

```
File Edit Selection View Go Run Terminal Help
vector_reduction_adv.cu vector_reduction_kernel.cu x matrixmul.cu
Reduction > vector_reduction_kernel.cu > reductionLarge(float *, int)
72 int temp = arr_size / 512;
73 int remainder = arr_size % 512;
74 int division = remainder == 0 ? temp : temp+1;
75
76 __shared__ float partialSum[1024];
77 //__shared__ float* temp;
78 unsigned int t = threadIdx.x;
79 unsigned int start = 2*blockIdx.x*blockDim.x;
80
81 for(int i = 0; i < division; ++i)
82 {
83     if(start + t < 512){
84         partialSum[t] += g_data[t+start + i * 512];
85     }
86     if(start+t+blockDim.x < 512){
87         partialSum[t+blockDim.x] += g_data[start+t+blockDim.x + i * 512];
88     }
89 }
90
91 for(unsigned int stride = blockDim.x/2; stride >= 1; stride >>= 1)
92 {
93     __syncthreads();
94     if(t < stride )
95     {
96         partialSum[t] += partialSum[t+stride];
97     }
98 }
99 __syncthreads();

PROBLEMS OUTPUT DEBUG CONSOLE TERMINAL
PS E:\cs677\hw2\Reduction> nvcc -o vec2 vector_reduction_adv.cu
vector_reduction_adv.cu
Creating library vec2.lib and object vec2.exp
PS E:\cs677\hw2\Reduction> ./vec2 1
(length of vector, kernel) --> (6000, 1)
Test PASSED
device: 3018652.000000 host: 3018652.000000
PS E:\cs677\hw2\Reduction> nvcc -o vec2 vector_reduction_adv.cu
vector_reduction_adv.cu
Creating library vec2.lib and object vec2.exp
PS E:\cs677\hw2\Reduction> ./vec2 1
(length of vector, kernel) --> (10000, 1)
Test PASSED
device: 5032845.000000 host: 5032845.000000
PS E:\cs677\hw2\Reduction> █
```

d.
3. Problem3

```

PROBLEMS  OUTPUT  DEBUG CONSOLE  TERMINAL

GPU computation complete
CPU computation complete
Test PASSED
PS E:\cs677\hw2\MatMul> nvcc -o mat matrixmul.cu
matrixmul.cu
    Creating library mat.lib and object mat.exp
PS E:\cs677\hw2\MatMul> ./mat
GPU computation complete
CPU computation complete
Test PASSED
PS E:\cs677\hw2\MatMul> █

```

a. █

4. Problem 4

- a. $1024 * 256 = 262,144$
- b. Two global memory loads and one global memory store for each thread.
- c. $256 * 2 + (128 + 64 + 32 + 16 + 8 + 4 + 2 + 1) * 3 = 1277$
- d. None
- e. $128 / 2 / 2 / 2 / 2 / 2 / 2 / 2 / 2 = 0$; So, there is 8 iterations of the for loop in total and 5 iterations will have branch divergence.
- f. There are 256 times to write back to the global memory. By reduction, we can store all value to thread 0 or the last thread then let that thread access the global memory, which can eliminate $1024 * 255$ accesses.