

ME759
High Performance Computing for Engineering Applications
Assignment 8
Parallel Prefix Scan

Date Assigned: October 28, 2013

Date Due: November 4, 2013 – 11:59 PM

The goals of this assignment are as follows:

A) Implement a tiled parallel solution of the exclusive scan operation discussed in class so that one can perform exclusive scan on more than 2048 elements. The number of elements however is going to be a power of two (see below).

B) Implement a work efficient algorithm to get either $O(n)$ or $O(n \log(n))$ performance.

For extra credit:

C) Implement your solution such that you can handle non-power-of-2 array sizes.

D) Implement your solution such that it minimizes shared memory bank conflicts.

To meet the goals of this assignment:

1) Edit the source files **scan_largearray.cu** to implement the necessary kernel functions. Use a tiled implementation so that computation can be performed on more than 2048 elements. Feel free to edit this file as you see fit. You might want to add, delete, or modify code; you might choose to ignore this altogether and start from scratch.

2) The modes of operation for the application are as follows:

a) No arguments: Randomly generate input data and compare against the host's result.

b) One argument: Randomly generate input data and write the result to a file, the name of which is specified by the first argument.

c) Two arguments: Read the first argument which indicates the size of the array (a power of 2), randomly generate input data and write the input data to the second argument (for generating random input data).

d) Three arguments: Read the first argument which indicate the size of the array, then input data from the file name specified by 2nd argument and write the output to a file name specified by the 3rd argument.

Note that if you wish to use the output of one run of the application as an input, you must delete the first line in the output file, which displays the accuracy of the values within the file. The value is not relevant for this application.

3) Please use the following steps for your performance analysis and report:

i) Near the top of the file **scan_largearray.cu**, set `#define DEFAULT_NUM_ELEMENTS` to 16777216 (which is 2^{24}). Set `#define MAX RAND` to 2.

ii) Include in your assignment report and also post on the ME759 Forum the performance results when running the executable without arguments on the CPU and GPU. Note that you are expected to report inclusive times in Release build mode.

iii) Using the NVIDIA profiler, generate and post on the forum and include in your report the **nvp** snapshot that shows how the execution time was spent by your program. Comment in your report on the amount of time spent in memory transactions and the bandwidth that your

program displays.

iv) Describe which of the four goals for this assignment (mentioned in the beginning) were achieved. How did you achieve them? Did you use any tricks to optimize for best performance? Please briefly discuss them if any.

The snapshot of the profiler

```
pkgupta3@euler01:~/prateek/HW08
[ pkgupta3@euler01 HW08 ]$ ./scan_large_20 16777216

**=====**
Processing 16777216 elements...
Host CPU Processing time: 66.026047 (ms)
CUDA Processing time: 4107.246094 (ms)
Speedup: 0.016076X
Test PASSED
device: 0.016076 host: 0.000000
[ pkgupta3@euler01 HW08 ]$ cat cuda_profile_0.log
# CUDA_PROFILE_LOG_VERSION 2.0
# CUDA_DEVICE 0 GeForce GTX 480
# CUDA_CONTEXT 1
# TIMESTAMPFACTOR 132c537ab5ae9929
method,gputime,cputime,occupancy
method=[ memcpyHtoD ] gputime=[ 11182.976 ] cputime=[ 17942.512 ]
method=[ _Z12prescanArrayPFS_i ] gputime=[ 8474.080 ] cputime=[ 9.006 ] occupancy=[ 0.667 ]
method=[ _Z13prescanArray2PFS_ ] gputime=[ 4065371.750 ] cputime=[ 3.789 ] occupancy=[ 0.667 ]
method=[ memcpyDtoH ] gputime=[ 14274.976 ] cputime=[ 4107337.000 ]
[ pkgupta3@euler01 HW08 ]$
```

Answer:

The file with named 'scan_large.cu' and scan_gold.cpp' are attached in the .zip folder.

The following goals are achieved:

A) Implementing a tiled parallel solution of the exclusive scan operation discussed in class so that one can perform exclusive scan on more than 2048 elements. The number of elements however is going to be a power of two.

This is achieved by using reduction method after implementing the algorithm for 2048 elements to further bring the results of rest 4096 blocks together and then performing the sum.

B) Implement a work efficient algorithm to get either $O(n)$ or $O(n \log(n))$ performance.

For extra credit:

The Hillis and Steele Algorithm is implemented with efficiency of $O(n \log(n))$ performance.

C) Implementing the solution such that you can handle non-power-of-2 array sizes.

D) Implement your solution such that it minimizes shared memory bank conflicts: **this is achieved by pinning the host memory.**