

Homework 2

Due Friday, April 22nd via GradeScope

Problem 1: Implement a parallel function that sums separately the odd and even values of a vector.

Idea: Need to implement `parallelSum` using `omp parallel for` with reductions for both even and odd accumulators.

```
1  std::vector<uint> parallelSum(const std::vector<uint> &v)
2  {
3      omp_set_num_threads(4);
4      std::vector<uint> sums(2);
5      uint sum0 = 0; uint sum1 = 0;
6
7      #pragma omp parallel for reduction(+:sum0) reduction(+:sum1)
8      for(uint i=0; i<v.size(); i++) {
9          if (v[i] % 2 == 0) {
10             sum0 += v[i];
11         }
12         else {
13             sum1 += v[i];
14         }
15     }
16     sums[0] = sum0; sums[1] = sum1;
17     return sums;
18 }
```

Console logs from `main_q1.cpp`.

```
$ make main_q1
g++ -std=c++11 -g -Wall -O3 -fopenmp main_q1.cpp -o main_q1

$ ./main_q1
Parallel
Sum Even: 757361650
Sum Odd: 742539102
Time: 0.00433168
Serial
Sum Even: 757361650
Sum Odd: 742539102
Time: 0.106256
main_q1.cpp:60:main      TEST PASSED.
```

Problem 2: Implement Radix Sort algorithm in parallel.

- **Question 1: computeBlockHistograms()** Idea: using `openMP` to parallelize computation across "blocks" when creating local histograms. Code must pass `Test1()`.

```
$ make main_q2
g++ -std=c++11 -g -Wall -O3 -fopenmp main_q2.cpp -o main_q2

$ ./main_q2
tests_q2.h:22:Test1      TEST PASSED.
```

- **Question 2: reduceLocalHistoToGlobal()** Idea: accumulate values based on the remainder of the index divided by `bucketSize`. Code must pass `Test2()`.

```
$ make main_q2
g++ -std=c++11 -g -Wall -O3 -fopenmp main_q2.cpp -o main_q2

$ ./main_q2
tests_q2.h:38:Test2      TEST PASSED.
```

- **Question 3: scanGlobalHisto()** Idea: implement cumulative sum using `std::partial_sum` standard library function. Note, needed to adjust `Output Iterator` and `Input Iterator` to ensure we begin at zero and do not inadvertently overflow.

```
$ make main_q2
g++ -std=c++11 -g -Wall -O3 -fopenmp main_q2.cpp -o main_q2

$ ./main_q2
tests_q2.h:50:Test3      TEST PASSED.
```

- **Question 4: computeBlockExScanFromGlobalHisto()** Idea: populate first using `globalHistoScan` and then increment using `blockHistograms` for subsequent blocks. This has the effect of splitting the global histogram into blocks need to update our sorting algorithm (next step).

```
$ make main_q2
g++ -std=c++11 -g -Wall -O3 -fopenmp main_q2.cpp -o main_q2

$ ./main_q2
tests_q2.h:67:Test4      TEST PASSED.
```

- **Question 5: populateOutputFromBlockExScan()** Idea: use pre-computed `blockExScan` to help target where entries of our unsorted input vector should map to in `sorted`. We can parallelize this operation by block using `openMP`. Note, we still need to re-compute which "bucket" each of our unsorted entries are from at each step since this information is not stored and passed to the function.

Also, this step only succeeds in sorting our input up to the `numBits`'th bit (in this case 8 bits of sorting per pass). Subsequent "passes" are needed to complete our radix sort algorithm since many input entries are greater than 256 (limit of 8 bits).

```
$ make main_q2
g++ -std=c++11 -g -Wall -O3 -fopenmp main_q2.cpp -o main_q2

$ ./main_q2
tests_q2.h:84:Test5      TEST PASSED.
```

- **Question 6: Serial vs parallel benchmarking** Use ICME GPU cluster to compare time estimates for different numbers of threads and blocks. In order to get this coded to run I needed to handle exception case where the number of blocks did not cleanly divide the number of elements in keys. In these cases, the derived block size implied we would not sort the end of the keys list. My solution was to pad the keys list with zeros until it evenly divided the next block number and remove this number of zero elements after sorting.

```
jelc@icme-gpu:~/cme213-para/hw2$ make
g++ -std=c++11 -g -Wall -O3 -fopenmp main_q1.cpp -o main_q1
g++ -std=c++11 -g -Wall -O3 -fopenmp main_q2.cpp -o main_q2
g++ -std=c++11 -g -Wall -O3 -fopenmp main_q2.cpp -D QUESTION6 -o main_q2_part6
```

```
jelc@icme-gpu:~/cme213-para/hw2$ sbatch script.sh
Submitted batch job 41096
```

```
jelc@icme-gpu:~/cme213-para/hw2$ cat job_41096.out
Date = Wed Apr 20 06:31:24 UTC 2022
Hostname = icmet01
Working directory = /home/jelc/cme213-para/hw2
```

```
Number of nodes allocated = 1
Number of cores/task allocated = 16
tests_q2.h:22:Test1 TEST PASSED.
tests_q2.h:38:Test2 TEST PASSED.
tests_q2.h:50:Test3 TEST PASSED.
tests_q2.h:67:Test4 TEST PASSED.
tests_q2.h:84:Test5 TEST PASSED.
```

```

Serial Radix Sort: PASS
Parallel Radix Sort: PASS
stl: 0.310788
serial radix: 0.0461267
parallel radix: 0.031369
Threads Blocks / Timing

```

	1	2	4	8	12	16	24	32	40	48
1	0.058	0.052	0.046	0.048	0.059	0.056	0.059	0.068	0.063	0.0
2	0.045	0.033	0.038	0.035	0.036	0.043	0.042	0.041	0.057	0.0
4	0.050	0.034	0.022	0.022	0.028	0.025	0.030	0.030	0.036	0.0
8	0.057	0.031	0.021	0.015	0.019	0.018	0.022	0.022	0.027	0.0
12	0.069	0.035	0.026	0.020	0.020	0.024	0.024	0.025	0.028	0.0
16	0.070	0.033	0.027	0.020	0.019	0.019	0.026	0.024	0.030	0.0
24	0.065	0.043	0.028	0.027	0.023	0.022	0.035	0.031	0.034	0.0
32	0.066	0.041	0.034	0.031	0.023	0.023	0.029	0.029	0.033	0.0
40	0.064	0.042	0.038	0.029	0.024	0.028	0.030	0.029	0.031	0.0
48	0.065	0.036	0.028	0.029	0.025	0.026	0.031	0.029	0.032	0.0

```
Benchmark runs: PASS
```

Optimal time of 0.015 seconds achieved at 8 threads / 8 blocks on the icme gpu cluster. See benchmark results above.

In general, our performance improves as we increase number of threads up available to compute our sorting program in parallel. However, there is a trade-off between extra parallelization and overhead required to generate and manage threads. For the scenarios we ran, the best performance was typically around 8-16 threads. Also since the icme gpu cluster allocates a maximum of 16 cpu cores per task, we should expect to hit a hardware limit at 32 threads (2x threads per core).

Similarly, for number of blocks the optimal is driven by a tradeoff between increased ability to parallelize sorting operations and overhead associated with combining sorted blocks. For the scenarios we ran, the best performance was typically around 8-16 blocks.

Submission information logs.

```

jelc@cardinal1:~$ /afs/ir.stanford.edu/class/cme213/script/submit.py hw2 private/cme213-
Submission for assignment 'hw2' as user 'jelc'
Attempt 1/10
Time stamp: 2022-04-13 20:09
List of files being copied:
  private/cme213-jelc53/hw2/sum.h  [768 bytes]
  private/cme213-jelc53/hw2/parallel_radix_sort.h  [7625 bytes]

```

Your files were copied successfully.

Directory where files were copied: /afs/ir.stanford.edu/class/cme213/submissions/hw2/jelc

List of files in this directory:

```
sum.h [768 bytes]
parallel_radix_sort.h [7625 bytes]
metadata [137 bytes]
```

This completes the submission process. Thank you!

```
jelc@cardinal1:~$ ls /afs/ir.stanford.edu/class/cme213/submissions/hw2/jelc/1
metadata parallel_radix_sort.h sum.h
```