

---

## Homework #2

( Due: April 6 )

### Task 1. [ 250 Points ] Distributed Sample Sort with Multithreaded Merge Sort inside Each Compute Node.

This task asks you to implement a distributed sample sort algorithm that runs a multithreaded merge sort algorithm inside each compute node. You will be sorting double precision floating point numbers. We will assume that all numbers in the input array are distinct. For simplicity generate the input array as follows. First initialize the array by storing the number  $i$  (as a double precision float) at location  $i$  of the array, where  $1 \leq i \leq n$ . Then choose two random locations of the array and swap their contents, and repeat this swapping step  $n$  times. Use this permuted array as input to your sorting implementations. When timing the implementations do not include the time needed to generate the input array.

- (a) [ 30 Points ] Implement and optimize the shared-memory parallel merge sort algorithm (PAR-MERGE-SORT-SM) shown in Figure 1 which uses the standard serial algorithm (MERGE) for merging two sorted sequences.

Optimize for the base case size  $m_1$ , that is, empirically find a value of  $m_1$  that gives you the best or close to the best performance for large values of  $n$ , and use that value of  $m_1$  in all subsequent runs of the algorithm. For this task it is OK to check only powers of 2 for the potential best value of  $m_1$ . Use all cores during this base case optimize phase.

Let  $N$  be the largest power of 2 such that this algorithm can sort  $N$  numbers in less than 5 minutes when run on a single processing core. Find the value of  $N$ .

- (b) [ 30 Points ] Implement and optimize the shared-memory parallel merge sort algorithm (PAR-MERGE-SORT-PM) of Figure 2 which uses a parallel merge algorithm (PAR-MERGE).

First optimize for the base case size  $m_2$  of PAR-MERGE, and then use that value of  $m_2$  when you optimize for the base case size  $m_3$  of PAR-MERGE-SORT-PM.

- (c) [ 30 Points ] Plot the running times of PAR-MERGE-SORT-SM and PAR-MERGE-SORT-PM for 10 equispaced points between  $n = 1000$  and  $n = N$  when run on all cores of the compute node you are using. Compare and explain the results.

- (d) [ 30 Points ] Generate Cilkview scalability plots for PAR-MERGE-SORT-SM and PAR-MERGE-SORT-PM assuming  $n = N$ . Compare and explain the results.

- (e) [ 50 Points ] Implement the distributed-memory parallel sample sort algorithm (DISTRIBUTED-SAMPLE-SORT) shown in Figure 3. Make two implementations. In one implementation use the merge sort algorithm (PAR-MERGE-SORT-SM) you implemented in part (a) for the local sorts in steps 2 and 5, and call this DISTRIBUTED-SAMPLE-SORT-SM. In another implementation use the merge sort algorithm (PAR-MERGE-SORT-PM) you implemented in part (b) and

MERGE( $A, p, q, r$ )	{ Merge sorted segments $A[p : q]$ and $A[q + 1 : r]$ , and store in $A[p : r]$ . }
1. $n_1 \leftarrow q - p + 1, n_2 \leftarrow r - q$ 2. create array $L[1 : n_1 + 1]$ and $R[1 : n_2 + 1]$ 3. $L[1 : n_1] \leftarrow A[p : q], R[1 : n_2] \leftarrow A[q + 1 : r]$ 4. $L[n_1 + 1] \leftarrow \infty, R[n_2 + 1] \leftarrow \infty$ 5. $i \leftarrow 1, j \leftarrow 1$ 6. <b>for</b> $k \leftarrow p$ <b>to</b> $r$ <b>do</b> 7. <b>if</b> $L[i] \leq R[j]$ <b>then</b> $A[k] \leftarrow L[i], i \leftarrow i + 1$ 8. <b>else</b> $A[k] \leftarrow R[j], j \leftarrow j + 1$	
PAR-MERGE-SORT-SM( $A, p, r$ )	{ Sort the numbers in $A[p : r]$ . All numbers in $A$ are assumed to be distinct. }
1. $n \leftarrow r - p + 1$ 2. <b>if</b> $n \leq m_1$ <b>then</b> <span style="float: right;">{ <math>m_1</math> is the global base case size for PAR-MERGE-SORT-SM }</span> 3.     sort $A[p : r]$ using <i>insertion sort</i> 4. <b>else</b> 5. $q \leftarrow \left\lfloor \frac{p+r}{2} \right\rfloor$ 6. <b>spawn</b> PAR-MERGE-SORT-SM( $A, p, q$ ) 7.         PAR-MERGE-SORT-SM( $A, q + 1, r$ ) 8. <b>sync</b> 9.         MERGE( $A, p, q, r$ ) <span style="float: right;">{ merge <math>A[p : q]</math> and <math>A[q + 1 : r]</math>, and store in <math>A[p : r]</math> }</span>	

Figure 1: Parallel merge sort with serial merge.

call this DISTRIBUTED-SAMPLE-SORT-PM. In DISTRIBUTED-SAMPLE-SORT-SM, execute the two recursive calls to PAR-MERGE-SORT-SM in steps 5–7 of PAR-MERGE-SORT-SM serially (i.e., remove the **spawn** and **sync** keyword) instead of calling them in parallel. But do not change anything inside the PAR-MERGE-SORT-PM implementation used in DISTRIBUTED-SAMPLE-SORT-PM.

- (f) [ 40 Points ] Plot the running times of DISTRIBUTED-SAMPLE-SORT-SM and DISTRIBUTED-SAMPLE-SORT-PM for 10 equispaced points between  $n = 1000$  and  $n = N$  when run on 5 compute nodes. For DISTRIBUTED-SAMPLE-SORT-PM use all cores inside each compute node. Compare and explain the results. Also plot the speedups w.r.t. PAR-MERGE-SORT-SM run on a single core of a single compute node.
- (g) [ 40 Points ] Repeat part (f), but without including the time needed for the initial distribution of input (step 1 of DISTRIBUTED-SAMPLE-SORT) and the final collection of output (step 6) in the running time.

PAR-MERGE( $T, p_1, r_1, p_2, r_2, A, p_3$ )	{ Merge disjoint sorted segments $T[p_1 : r_1]$ and $T[p_2 : r_2]$ , and store result in $A[p_3 : r_3]$ , where $r_3 = p_3 + n_1 + n_2 - 1$ , $n_1 = r_1 - p_1 + 1$ , and $n_2 = r_2 - p_2 + 1$ . All numbers in $T$ are assumed to be distinct. }
1. $n_1 \leftarrow r_1 - p_1 + 1, n_2 \leftarrow r_2 - p_2 + 1, r_3 = p_3 + n_1 + n_2 - 1$ 2. <b>if</b> $n_1 + n_2 \leq m_2$ <b>then</b> <span style="float: right;">{ <math>m_2</math> is the global base case size for PAR-MERGE }</span> 3.     use the standard serial merge algorithm to merge $T[p_1 : r_1]$ and $T[p_2 : r_2]$ , and store result in $A[p_3 : r_3]$ <b>else</b> 4. <b>if</b> $n_1 < n_2$ <b>then</b> $p_1 \leftrightarrow p_2, r_1 \leftrightarrow r_2, n_1 \leftrightarrow n_2$ 5. $q_1 \leftarrow \left\lfloor \frac{p_1 + r_1}{2} \right\rfloor$ 6. $q_2 \leftarrow \text{BINARY-SEARCH}( T[q_1], T, p_2, r_2 )$ <span style="float: right;">{ each number in <math>T[p_2 : q_2 - 1]</math> is smaller than <math>T[q_1]</math> and each in <math>T[q_2 : r_2]</math> is larger than <math>T[q_1]</math> }</span> 7. $q_3 \leftarrow p_3 + (q_1 - p_1) + (q_2 - p_2)$ 8. $A[q_3] \leftarrow T[q_1]$ 9. <b>spawn</b> PAR-MERGE( $T, p_1, q_1 - 1, p_2, q_2 - 1, A, p_3$ ) 10.     PAR-MERGE( $T, q_1 + 1, r_1, q_2 + 1, r_2 - 1, A, q_3$ ) 11. <b>sync</b>	
PAR-MERGE-SORT-PM( $A, p, r$ )	{ Sort the numbers in $A[p : r]$ . All numbers in $A$ are assumed to be distinct. }
1. $n \leftarrow r - p + 1$ 2. <b>if</b> $n \leq m_3$ <b>then</b> <span style="float: right;">{ <math>m_3</math> is the global base case size for PAR-MERGE-SORT-PM }</span> 3.     sort $A[p : r]$ using <i>insertion sort</i> <b>else</b> 4. $q \leftarrow \left\lfloor \frac{p + r}{2} \right\rfloor$ 5. <b>spawn</b> PAR-MERGE-SORT-PM( $A, p, q$ ) 6.     PAR-MERGE-SORT-PM( $A, q + 1, r$ ) 7. <b>sync</b> 8. $T[p : r] \leftarrow A[p : r]$ <span style="float: right;">{ <math>T</math> is a global array of the same size as <math>A</math> }</span> 9.     PAR-MERGE( $T, p, q, q + 1, r, A, p$ ) <span style="float: right;">{ merge <math>T[p : q]</math> and <math>T[q + 1 : r]</math>, and store in <math>A[p : r]</math> }</span>	

Figure 2: Parallel merge sort with parallel merge.

DISTRIBUTED-SAMPLE-SORT( $n, p, q$ )	<i>{Sort <math>n</math> distinct numbers using <math>p</math> processors. One node contains all <math>n</math> input numbers initially, and also the numbers in the final sorted order.}</i>
<ol style="list-style-type: none"> <li>1. <b>Initial Distribution:</b> The master node scatters the <math>n</math> keys among <math>p</math> processing nodes as evenly as possible.</li> <li>2. <b>Pivot Selection:</b> Each node sorts its local keys, and selects <math>q - 1</math> evenly spaced keys from its sorted sequence. The master node gathers these local pivots from all nodes, locally sorts those <math>p(q - 1)</math> keys, selects <math>p - 1</math> evenly spaced global pivots from them, and broadcasts them to all nodes.</li> <li>3. <b>Local Bucketing:</b> Each node inserts the global pivots into its local sorted sequence using binary search, and thus divides the keys among <math>p</math> buckets.</li> <li>4. <b>Distribute Local Buckets:</b> For <math>1 \leq i \leq p</math>, each node sends bucket <math>i</math> to node <math>i</math>.</li> <li>5. <b>Local Sort:</b> Each node locally sorts the elements it received in step 4.</li> <li>6. <b>Final Collection:</b> The master node collects all sorted keys from all nodes, and for <math>1 \leq i &lt; p</math>, places all keys from node <math>i</math> ahead of all keys from node <math>i + 1</math>.</li> </ol>	

Figure 3: Distributed sample sort.

## APPENDIX 1: Calling Cilk++ Functions from MPI Code

ncr.cilk	ncr-mpi.cpp
<pre> #include &lt;cilk.h&gt;  int nCr( int n, int r ) {     if ( r &gt; n ) return 0;     if ( ( r == 0 )    ( r == n ) ) return 1;      int x, y;      x = cilk_spawn nCr( n - 1, r - 1 );     y = nCr( n - 1, r );      cilk_sync;      return ( x + y ); }  extern "C++" int nCr_CPP( int n, int r ) {     return cilk::run( nCr, n, r ); } </pre>	<pre> #include &lt;mpi.h&gt;  extern "C++" int nCr_CPP( int n, int r );  int main( int argc, char *argv[ ] ) {     MPI_Init( &amp;argc, &amp;argv );      int rank;     MPI_Comm_rank( MPI_COMM_WORLD, &amp;rank );      printf( "C( %d, %d ) = %d\n", 30, 15 + rank,             nCr_CPP( 30, 15 + rank ) );      MPI_Finalize( );      return 0; } </pre>

In `ncr.cilk` we have a Cilk++ function called `nCr` which we would like to call from within the MPI code `ncr-mpi.cpp`. Since we do not have a `cilk_main` function in `ncr-mpi.cpp`, we do not have a Cilk++ context, and so `nCr` cannot be called directly from within `ncr-mpi.cpp`. Instead we create a function (named `nCr_CPP`) callable from C++ which starts a Cilk++ environment through

`cilk::run` and calls `nCr`.

You can compile and link the files as follows on Lonestar/Stampede. The first command creates a shared library named `libncr.so` from `ncr.cilk`, and the second one compiles `ncr-mpi.cpp` and links it with `libncr.so`.

```
cilk++ -m64 -fPIC -shared -o libncr.so ncr.cilk
mpicxx ncr-mpi.cpp -L. -L$CILKHOME/lib64 -Wl,-rpath=.
        -lncr -lcilk_main -lcilkrts -lcilkutil
```

The resulting MPI program (`a.out`) can be run as follows (from your job script).

```
ibrun tacc_affinity a.out
```

If you want to run your MPI program on  $t$  compute nodes on Lonestar, and launch  $k \in \{1, 2, 3, 4, 6, 12\}$  parallel processes on each node, then include the following line in your job script with  $m = 12t$ .

```
#$ -pe kway m
```

If  $k$  parallel processes are launched on each node, then Cilk++ functions called from each process will be able to launch at most  $12/k$  concurrent threads. Recall that when multiple processes are launched on the same node then the total memory is divided among the processes and no process is able to access the memory allocated to other processes, but all threads running under a process share the memory allocated to that process.

## APPENDIX 2: What to Turn in

One compressed archive file (e.g., zip, tar.gz) containing the following items.

- Source code, makefiles and job scripts.
- A PDF document containing all answers and plots.

## APPENDIX 3: Things to Remember

- **Please never run anything that takes more than a minute or uses multiple cores on TACC login nodes.** TACC policy strictly prohibits such usage. They reserve the right to suspend your account if you do so. All runs must be submitted as jobs to compute nodes (even when you use Cilkview or PAPI).
- Please store all data in your work folder (`$WORK`), and not in your home folder (`$HOME`).
- When measuring running times please exclude the time needed for reading the input and writing the output. Measure only the time needed by the algorithm.