

Programming Assignment 3 - Report

Ankit Srivastava & Cansu Tetik

February 15, 2017

1 Problem Description

Given a list of n integers, block distributed on p processors, we want to sort the integers using a parallel version of the sequential quicksort algorithm.

2 Algorithm Design Description

In this section, we give a brief description of the algorithms that we implemented for solving the given problem.

2.1 Sequential Algorithm

Since the problem description said that any algorithm can be used for sorting sequentially, we used C++ *algorithm* library's *std::sort* function for sorting when $p = 1$.

2.2 Parallel Algorithm

For sorting the given n elements in parallel on p processors, using quicksort, we partition the given array into two parts using a randomly chosen pivot: one consisting of elements less than or equal to the pivot and the other consisting of elements greater than the pivot. Then, based on the sizes of the two partitions, we split the p processors into two parts and block distribute the elements in both the partitions on the corresponding processors. We then sort the two partitions recursively and later redistribute the two sorted partitions on p processors.

We briefly explain our implementation of the given function, and a few other utility functions that we added, below.

- **parallel_sort**: This is the function that is called by the framework for sorting in parallel. We implemented the above defined algorithm primarily in this function. The function first checks if the communicator size is 1 and sorts sequentially if that is the case. If the communicator size is more than 1, it chooses a random pivot and partitions the array into two parts, using **partition** function described later, and splits the communicator based on the partition sizes. The partitioned data is then distributed on the corresponding processors, using **distribute_data** function. We then call **parallel_sort** again, for sorting the partitioned array recursively. Once the two partitions are sorted, we redistribute the sorted elements in the two partitions as per the original sizes, using **collect_data** function.
- **partition**: This function works like a sequential partition function for quicksort and is used by all the processors to partition their local data by the pivot, in place.
- **distribute_data**: This function moves data in a particular partition on all the processors to the processors assigned to sort elements in that partition. The function tries to minimize the number of required communications by only moving data whose movement is required i.e. data in this partition that resides on processors assigned to the other partition. It calls **MPI_Alltoallv** once for distributing the data.
- **collect_data**: This function collects sorted data in the processors assigned to the two partitions and block distributes it in the original configuration passed to **parallel_sort** that called this function. Again, it calls **MPI_Alltoallv** once for collecting the data.

Pseudocode for the parallel algorithm is shown in Algorithm 1.

Algorithm 1 Parallel

```
1: procedure LOCALPARTITION(begin, end, pivot)
2:   left  $\leftarrow$  0
3:   right  $\leftarrow$  (end - begin) - 1
4:   while left < right do
5:     if pivot  $\in$  [a[right], a[left]) then
6:       Swap a[left] and a[right]
7:     if a[left]  $\leq$  pivot then
8:       left  $\leftarrow$  left + 1
9:     if a[right] > pivot then
10:      right  $\leftarrow$  right - 1
11:   return left
12: procedure DISTRIBUTEDDATA(sbuf, current, rbuf, final, communicator)
13:   q  $\leftarrow$  Size of communicator
14:   r  $\leftarrow$  Rank in communicator
15:   extra  $\leftarrow$  Vector of size q
16:   need  $\leftarrow$  Vector of size q
17:   for p  $\in$  [0, q) do
18:     if final[p] < current[p] then
19:       extra[p]  $\leftarrow$  current[p] - final[p]
20:     else
21:       need[p]  $\leftarrow$  final[p] - current[p]
22:   cumulativeNeed  $\leftarrow$  Prefix sum of need
23:   cumulativeExtra  $\leftarrow$  Prefix sum of extra
24:   scounts  $\leftarrow$  Vector of size q
25:   rcounts  $\leftarrow$  Vector of size q
26:   extraMin  $\leftarrow$  cumulativeExtra[r] - extra[r]
27:   extraMax  $\leftarrow$  cumulativeExtra[r] - 1
28:   needMin  $\leftarrow$  cumulativeNeed[r] - need[r]
29:   needMax  $\leftarrow$  cumulativeNeed[r] - 1
30:   for p  $\in$  [0, q) do
31:     if r is p then
32:       continue
33:     if cumulativeExtra[p]  $\in$  (needMin, needMax + extra[p]) then
34:       procMin  $\leftarrow$  max(cumulativeExtra[p] - extra[p], needMin)
35:       procMax  $\leftarrow$  min(cumulativeExtra[p] - 1, needMax) + 1
36:       rcounts[p]  $\leftarrow$  procMax - procMin
37:     if cumulativeNeed[p]  $\in$  [extraMin, extraMax + need[p]) then
38:       procMin  $\leftarrow$  max(extraMin, cumulativeNeed[p] - need[p])
39:       procMax  $\leftarrow$  min(extraMax, cumulativeNeed[p] - 1) + 1
40:       scounts[p]  $\leftarrow$  procMax - procMin
```

```

40:   if  $extra[r] > 0$  then
41:        $scounts[r] \leftarrow final[r]$ 
42:        $rcounts[r] \leftarrow final[r]$ 
43:   else
44:        $scounts[r] \leftarrow current[r]$ 
45:        $rcounts[r] \leftarrow current[r]$ 
46:    $sdispls \leftarrow$  Offsets corresponding to  $scounts$ 
47:    $rdispls \leftarrow$  Offsets corresponding to  $rcounts$ 
48:    $MPI\_ALLTOALLV(sbuf, scounts, sdispls, rbuf, rcounts, rdispls)$ 
49: procedure COLLECTDATA( $sbuf, current, rbuf, final, communicator$ )
50:    $q \leftarrow$  Size of  $communicator$ 
51:    $r \leftarrow$  Rank in  $communicator$ 
52:    $cumulativeCur \leftarrow$  Prefix sum of  $current$ 
53:    $cumulativeFin \leftarrow$  Prefix sum of  $final$ 
54:    $scounts \leftarrow$  Vector of size  $q$ 
55:    $rcounts \leftarrow$  Vector of size  $q$ 
56:    $currentMin \leftarrow cumulativeCur[r] - current[r]$ 
57:    $currentMax \leftarrow cumulativeCur[r] - 1$ 
58:    $finalMin \leftarrow cumulativeFin[r] - final[r]$ 
59:    $finalMax \leftarrow cumulativeFin[r] - 1$ 
60:   for  $p \in [0, q)$  do
61:       if  $cumulativeCur[p] \in (finalMin, finalMax + current[p])$  then
62:            $procMin \leftarrow \max(cumulativeCur[p] - current[p], finalMin)$ 
63:            $procMax \leftarrow \min(cumulativeCur[p] - 1, finalMax) + 1$ 
64:            $rcounts[p] \leftarrow procMax - procMin$ 
65:       if  $cumulativeFin[p] \in [currentMin, currentMax + final[p])$  then
66:            $procMin \leftarrow \max(extraMin, cumulativeNeed[p] - need[p])$ 
67:            $procMax \leftarrow \min(currentMax, cumulativeFin[p] - 1) + 1$ 
68:            $scounts[p] \leftarrow procMax - procMin$ 
69:    $sdispls \leftarrow$  Offsets corresponding to  $scounts$ 
70:    $rdispls \leftarrow$  Offsets corresponding to  $rcounts$ 
71:    $MPI\_ALLTOALLV(sbuf, scounts, sdispls, rbuf, rcounts, rdispls)$ 
72: procedure PARALLELSORT( $begin, end, communicator$ )
73:    $q \leftarrow$  Size of  $communicator$ 
74:    $r \leftarrow$  Rank in  $communicator$ 
75:    $localSize \leftarrow end - begin$ 
76:    $m \leftarrow$  Sum of  $localSize$  across all the processors in  $communicator$ 
77:   if  $q$  is 1 then
78:        $std::sort(begin, end)$ 
79:    $minIndex \leftarrow$  Global minimum index in this processor

```

```

80:  Seed the random number generator with the level of recursion
81:   $pivot \leftarrow$  Random pivot element
82:   $boundary \leftarrow$  LOCALPARTITION( $begin, end, pivot$ )
83:   $partitionSizes \leftarrow \{boundary + 1, localSize - (boundary + 1)\}$ 
84:   $currentSizes \leftarrow$  Gather  $partitionSizes$  from all the processors
85:   $sumSizes \leftarrow$  Sum of both  $partitionSizes$  across all the processors
86:   $numProcessors \leftarrow$  Split  $q$  into two based on  $sumSizes$ 
87:   $color \leftarrow 0$  if  $r < numProcessors[0]$  else 1
88:   $newCommunicator \leftarrow$  MPI_COMM_SPLIT( $communicator, color, r$ )
89:   $finalSizes \leftarrow$  Block distribution of  $sumSizes$  on  $numProcessors$ 
90:   $finalSize \leftarrow$  Block distribute  $sumSizes$  on  $numProcessors$ 
91:   $finalBuf \leftarrow$  Temporary buffer of  $finalSize$ 
92:   $dataOffset \leftarrow 0$ 
93:  for  $i \in \{0, 1\}$  do
94:     $current \leftarrow currentSizes[color]$ 
95:     $final \leftarrow finalSizes[color]$ 
96:     $sbuf \leftarrow begin[dataOffset]$ 
97:     $rbuf \leftarrow finalBuf$  if  $i == color$  else NULL
98:    DISTRIBUTE_DATA( $sbuf, current, rbuf, final, newCommunicator$ )
99:     $dataOffset \leftarrow (boundary + 1)$ 
100:   $newBegin \leftarrow$  Initial pointer of  $finalBuf$ 
101:   $newEnd \leftarrow$  Final pointer of  $finalBuf$ 
102:  PARALLEL_SORT( $newBegin, newEnd, newCommunicator$ )
103:   $original \leftarrow$  Block decompose  $m$  on  $q$  processors
104:  COLLECT_DATA( $newBegin, current, begin, original, communicator$ )

```

3 Plots and Observations

For running the experiments, we had only two parameters that we could vary: number of elements to be sorted (n) and the number of processors used for sorting (p). We obtained timings for different combinations of the two parameters and plotted and analyzed the run times so obtained. In the following sections, we discuss how we varied the parameters and the effect that the variation had on the performance.

We used *jinx4*, which has 4×6 cores, for running the experiments up to $p = 24$ and used 5 *sixcore* nodes for running the experiments for $p = \{26, 28, 30\}$.

3.1 Array Size (n)

First of all, we had to choose a suitable problem size for running our experiments. During some preliminary experimentation, we observed that the implementation started running out of memory while sorting 1 billion elements on higher number of processors. Since *jinx* cluster has low memory per node (maximum of 24 GB on one node) and 1 billion integers means 4 GB of memory, the memory could quickly run out if it requires a lot of recursions to solve the problem. Therefore, this was expected.

The implementation was able to sort arrays of size up to 900 million, up to $p = 24$ without any memory issues. However, for $p \geq 26$ the 900 million array size also starts failing because of insufficient memory. Therefore, we chose to vary the array size in the set $\{500000000, 600000000, 700000000, 800000000, 900000000\}$. The run times so obtained, for different problem sizes, are plotted in Figure 1 and the corresponding speedups, with respect to the run time of $p = 1$, are plotted in Figure 2.

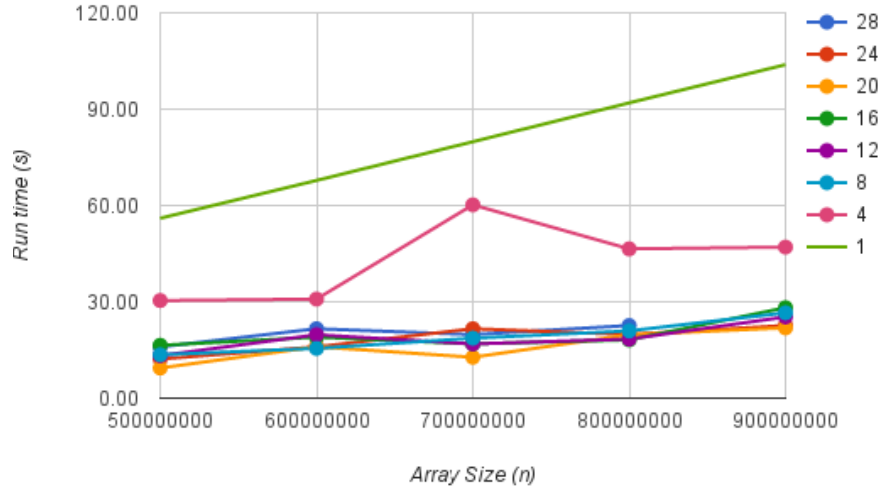


Figure 1: Run time v/s Array Size

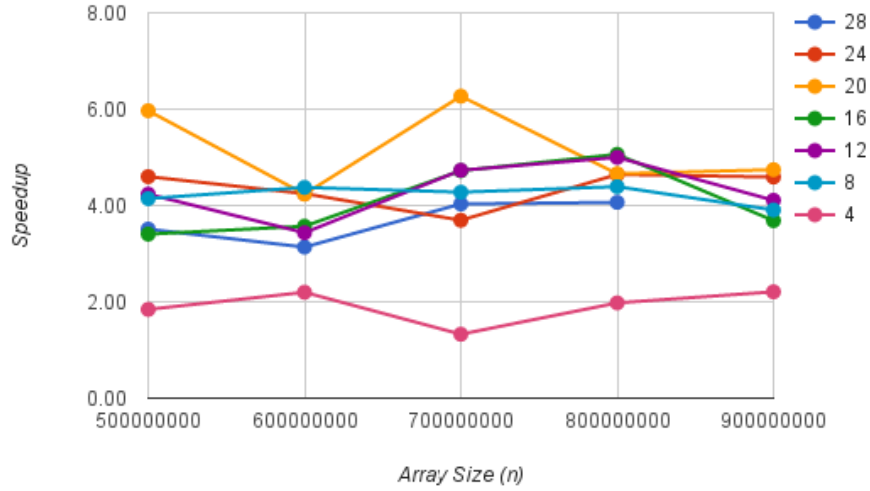


Figure 2: Speedup v/s Array Size

It can be seen from the run time plot that the run time generally increases on increasing the problem size and is maximum for $p = 1$, followed by $p = 4$ case. However, the distinction between run times on higher number of processors is less clear. Therefore, we refer to the speedup plot where we can see that the speedup remains approximately the same on increasing the problem size. This was expected since the speedup depends on the number of processors used. We

did observe that the maximum speedup was obtained for $n = 700000000$, using $p = 20$.

3.2 Number of Processors (p)

We varied number of processors in the range $[1, 30]$, in jumps of 2, and obtained the run times for all the problem sizes discussed in Section 3.1. The run time and speedup plots so obtained are plotted in Figure 3 and Figure 4. As discussed in Section 3.1, we didn't have any data for $p \geq 26$ and $n = 900000000$ and therefore the corresponding data points are missing from the plots.

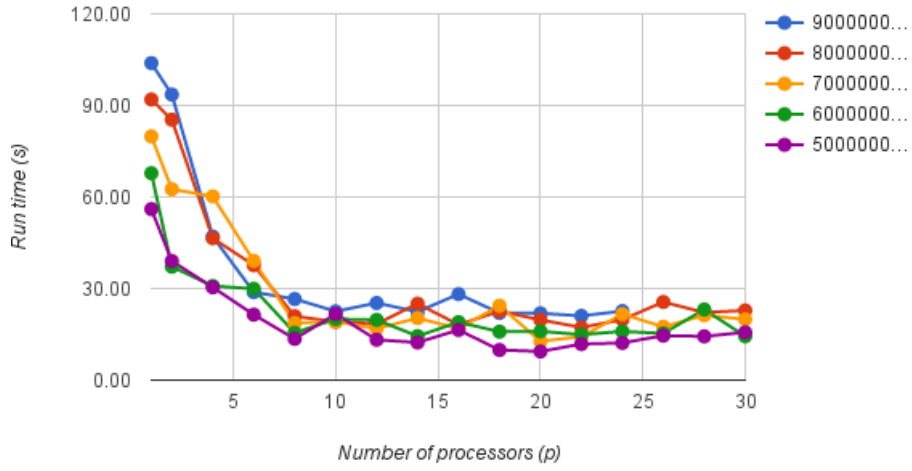


Figure 3: Run time v/s Number of Processors

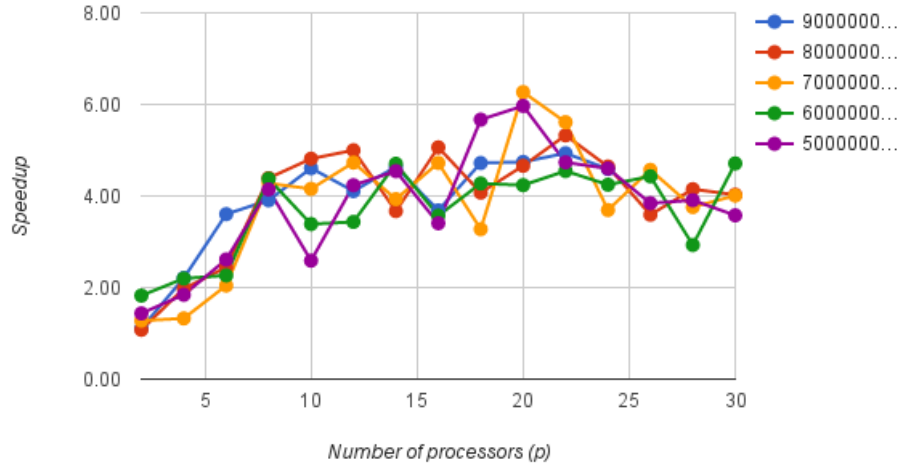


Figure 4: Speedup v/s Number of Processors

From the above plots, we can see that on increasing the number of processors the runtime decreases up to $p = 20$ and then starts increasing. Thus the observed speedup is highest, for most problem sizes, for $p = 20$. This wasn't expected and could be because of the choice of pivot or because of the low memory on the nodes.