

Cloud Computing

Special Task 2 - Parallel Merge Sort with MPI – Summer Term 2018

Prachi Agrawal, Henry Cocos, David Merkl, Samuel Santos

{agrawal,cocos,merkl,santos}@stud.fra-uas.de

Computer Science
Faculty of Computer Science and Engineering
Frankfurt University of Applied Sciences

July 4th 2018

Contents

- 1 Theory
- 2 Implementation
- 3 Results
- 4 Conclusion
- 5 References

Merge Sort Pseudo Code I

Algorithm 1 Pseudocode for Merge Sort [1]

```

function mergesort(List m)
  if length of m  $\leq$  1 then
    return m
  end if
  left := emptylist
  right := emptylist
  for each x with index i in m do
    if i < (length of m) / 2 then
      add x to left
    else
      add x to right
    end if
  end for

  left := mergesort(left)
  right := mergesort(right)
  return merge(left, right)
end function

```

Merge Sort Pseudo Code II

Algorithm 2 Pseudocode for Merge Function [1]

```

function merge(List left, List right)
  result := emptylist
  while left is not empty and right is not empty do
    if first(left) <= first(right) then
      append first(left) to result
    else
      append first(right) to result
    end if
  end while

  while left is not empty do
    append first(left) to result
  end while

  while right is not empty do
    append right to result
  end while
return result
end function

```

Merge Sort Example

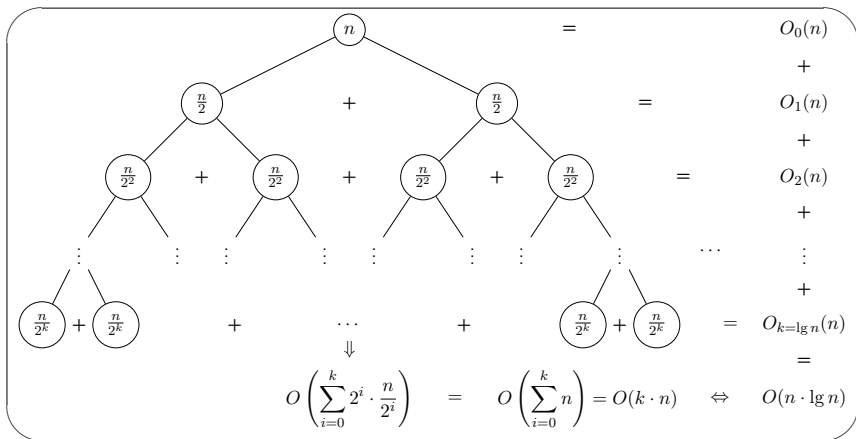


Figure: Recursion Tree Merge Sort [2]

Source Code MPI Merge Sort I

```
1  /***** Initialize MPI *****/
2  int world_rank;
3  int world_size;
4
5  MPI_Init(&argc, &argv);
6  MPI_Comm_rank(MPI_COMM_WORLD, &world_rank);
7  MPI_Comm_size(MPI_COMM_WORLD, &world_size);
8
9  /***** Divide the array in equal-sized chunks *****/
10 int size = n/world_size;
11
12 /***** Send each subarray to each process *****/
13 int *sub_array = malloc(size * sizeof(int));
14 MPI_Scatter(original_array, size, MPI_INT, sub_array, size, MPI_INT, 0,
15             MPI_COMM_WORLD);
16
17 // Start Timer Sorting Parallel
18 parallelTime1 = MPI_Wtime();
19
20 /***** Perform the mergesort on each process *****/
21 int *tmp_array = malloc(size * sizeof(int));
22 mergeSort(sub_array, tmp_array, 0, (size - 1));
```

Listing 1: Merge Sort with MPI [3]

Source Code MPI Merge Sort II

```
1  /***** Gather the sorted subarrays into one *****/
2  int *sorted = NULL;
3  if(world_rank == 0) {
4      sorted = malloc(n * sizeof(int));
5  }
6
7  // Start Timer Sorting Parallel
8  parallelTime2 = MPI_Wtime();
9
10 MPI_Gather(sub_array, size, MPI_INT, sorted, size, MPI_INT, 0,
11           MPI_COMM_WORLD);
12
13 // Sequential Last Merge Sort Start Timer
14 sequentialTime1 = MPI_Wtime();
15 /***** Make the final mergeSort call *****/
16 if(world_rank == 0) {
17     int *other_array = malloc(n * sizeof(int));
18     mergeSort(sorted, other_array, 0, (n - 1));
```

Listing 2: Merge Sort with MPI [3]

Source Code MPI Merge Sort III

```
1 // Sequential Last Merge Sort Stop Timer
2 sequentialTime2 = MPI_Wtime();
3
4 /***** Clean up rest *****/
5 free(original_array);
6 free(sub_array);
7 free(tmp_array);
8
9 if(world_rank == 0)
10 {
11 // Time of Execution
12 printf("%i \t %.3f \t\t %f \t %f \t\t %f \n", numProc, (sequentialTime2 -
    sequentialMasterRead1), (sequentialMasterRead2 - sequentialMasterRead1),
    (parallelTime2 - parallelTime1), (sequentialTime2 - sequentialTime1) );
13 }
14 /***** Finalize MPI *****/
15 MPI_Barrier(MPI_COMM_WORLD);
16 MPI_Finalize();
17
18 }
```

Listing 3: Merge Sort with MPI [3]

Automated Tests - Shell Script

```

1 # Iterate over all Problem Sizes
2 for PROBLEM_SIZE in 1000 10000 100000 1000000 10000000 100000000 1000000000 10000000000 100000000000
   1000000000000 10000000000000
3 do
4
5 # Print Header for Result File
6 echo -e "#| Num of Proc | Total Time Elapsed | File Read Part on Master | Parallel Part | Sequential Last Merge
   Call |" >> "$RESULT_DIR"$RESULT_FILE_NAME_"$PROBLEM_SIZE.txt"
7
8 # Iterate over all Number of Processes
9 for NUM_PROC in 1 2 4 8 16 32 64 128 256 512
10 do
11
12 # Check if Number of Processes gets bigger than Number of Nodes
13 if [ "$NUM_PROC" -gt 128 ]
14 then
15     NUM_NODE=128
16 else
17     NUM_NODE=$NUM_PROC
18 fi
19
20 # Prompt Test Run on Console and Write to File
21 echo "Array Size: $PROBLEM_SIZE"
22 echo -e "Number of Processes: $NUM_PROC\nNumber of Nodes: $NUM_NODE\nProblem Size: $PROBLEM_SIZE\n"
23 echo -e "Number of Processes: $NUM_PROC\nNumber of Nodes: $NUM_NODE\nProblem Size: $PROBLEM_SIZE\n" >> "$LOG_DIR"
   "$LOGFILE_NAME_"$PROBLEM_SIZE.txt"
24
25 # Execute MPI with Parameters
26 mpiexec -np $NUM_PROC --hostfile "$HOSTFILE_DIR"cluster_$NUM_NODE.txt $MPI_EXEC_NAME $PROBLEM_SIZE $NUM_PROC >> "
   $RESULT_DIR"$RESULT_FILE_NAME_"$PROBLEM_SIZE.txt"
27
28
29 done
30 echo -e "Successful Test Run for $PROBLEM_SIZE !!!\n\n"
31 echo -e "Successful Test Run for $PROBLEM_SIZE !!!\n\n" >> "$LOG_DIR"$LOGFILE_NAME_"$PROBLEM_SIZE.txt"
32
33 done
34
35 echo -e "All Test runs completed successfully!!!"

```

Listing 4: Shell Script for automated Tests

Automated Test Runs

The Test runs were automated with the following Parameters:

- 1.000 entries
- 10.000 entries
- 100.000 entries
- 1.000.000 entries
- 10.000.000 entries
- 100.000.000 entries

The Test runs with 100.000.000 Entries produced errors!!!
Assumption is that Memory Size is not sufficient!!!

Memory

4 Bytes x 100.000.000 \approx 400 Mbytes

2 x 400 Mbytes \approx 800 Mbytes

Raspberry Pi 3 has only 1 Gbytes

Automated Tests - Result File

```
1 # | Num of Proc | Total Time Elapsed | File Read Part on
   Master | Parallel Part | Sequential Last Merge Call |
2 1   0.051      0.000231      0.000848      0.000724
3 2   0.074      0.000230      0.000391      0.000741
4 4   0.080      0.000230      0.000180      0.000761
5 8   0.094      0.000229      0.000097      0.000837
6 16  0.115      0.000229      0.000042      0.000801
7 32  0.220      0.000383      0.000024      0.000817
8 64  0.273      0.000232      0.000015      0.000842
9 128 0.511      0.000382      0.000011      0.000845
10 256 0.748      0.000384      0.000011      0.000841
11 512 1.947      0.000386      0.000008      0.001059
```

Listing 5: Example Results File for 1000 Entries

Automated Tests - Log File

```
1 Number of Processes: 1
2 Number of Nodes: 1
3 Problem Size: 1000
4
5 Number of Processes: 2
6 Number of Nodes: 2
7 Problem Size: 1000
8
9 Number of Processes: 4
10 Number of Nodes: 4
11 Problem Size: 1000
12
13 Number of Processes: 8
14 Number of Nodes: 8
15 Problem Size: 1000
16
17 Number of Processes: 16
18 Number of Nodes: 16
19 Problem Size: 1000
```

Listing 6: Example Log File for 1000 Entries

Results

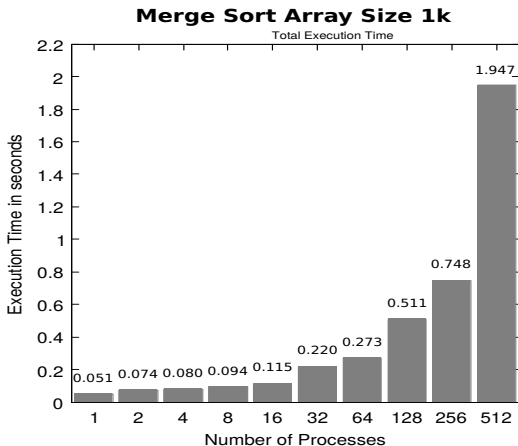


Figure: Plot for 1 Thousand Entries

Results

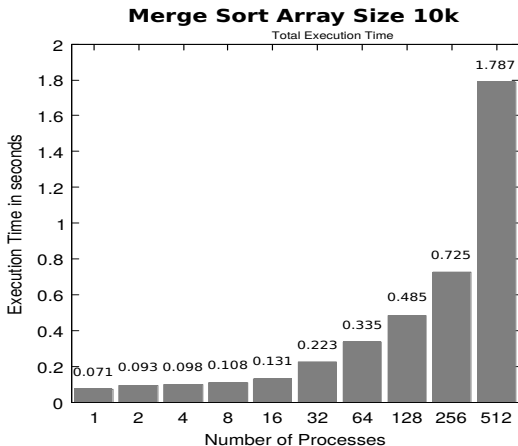


Figure: Plot for 10 Thousand Entries

Results

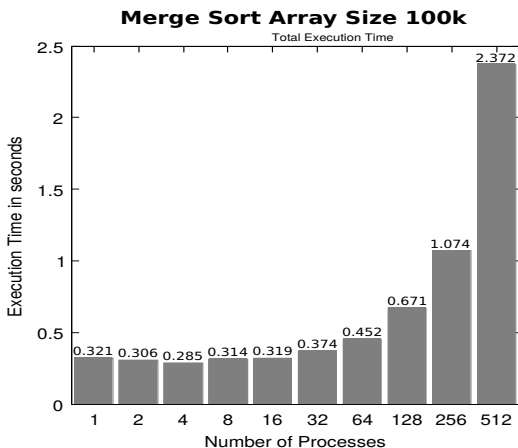


Figure: Plot for 100 Thousand Entries

Results

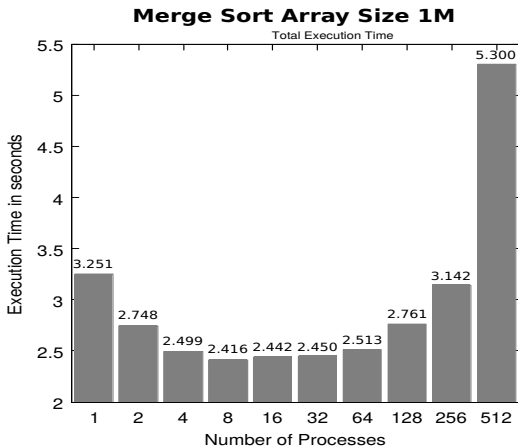


Figure: Plot for 1 Million Entries

Results

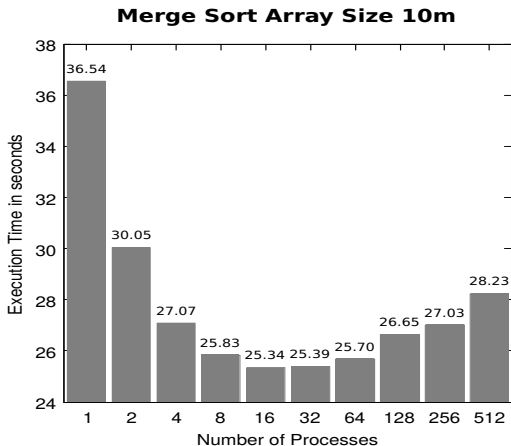


Figure: Plot for 10 Million Entries

Conclusion

Figure 6 demonstrates the following things:

- The Problem is scaling
- Increasing the Number of Processes on small problems leads to communication overhead
- Increasing the Problem Size leads to better performance

- [1] “Merge Sort Wikipedia,” [Accessed: July 4, 2018]. [Online]. Available: https://en.wikipedia.org/wiki/Merge_sort
- [2] “Merge Sort Recursion Tree,” [Accessed July 4, 2018]. [Online]. Available: <http://www.texample.net/tikz/examples/merge-sort-recursion-tree/>
- [3] “Merge Sort MPI Implementation,” [Accessed: July 4, 2018]. [Online]. Available: <https://github.com/racorretjer/Parallel-Merge-Sort-with-MPI>