# Bubble Sort using Domain Decomposition with MPI

Claudio Scheer, Gabriell Araujo

## General Setup

We ran our *batch job* on two nodes (2x12 cores, 2x24 when considering hyper-threading) in the Cerrado cluster. All experiments were executed three times and then the average execution time and the standard deviation were calculated. Efficiency and speedup were based on the execution time reported by the sequential execution of the bubble sort algorithm.

For the implementation using MPI, we used the domain decomposition philosophy. In short, each process sorts $1/n$ of the vector using the bubble sort algorithm and shares its lowest values with the left neighbor. After that, the piece of the vector shared is interleaved with the vector held by the left neighbor. We tested the sharing of 10%, 30% and 50% of the vector. These steps are repeated until the vector distributed over the processes is sorted. In this report, we discuss three optimization that we perform in this workflow.

Considering the computational power available, we tested our implementation using the 24 physical cores with and without hyper-threading.
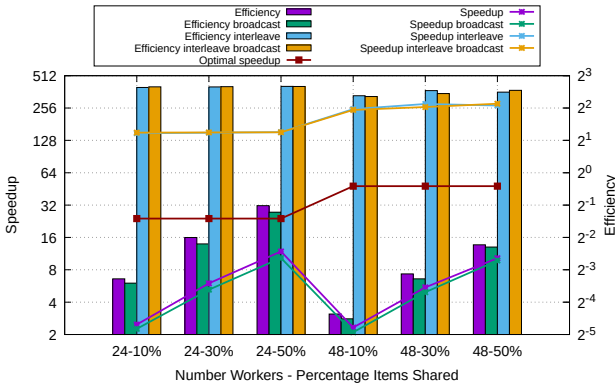


Figure 1: Speedup x Efficiency

## Bubble Sort without optimizations

The bubble sort problem addressed here consists of sorting one vector with 1000032 integers. This vector is divided among the processes. Therefore, each process have $n/processes$ integers. As we are sorting a vector in descending order, each process can generate its own vector, without the need to create the entire vector and share it between the processes.

The first phase of the bubble sort algorithm using domain decomposition structure is the application of the bubble sort algorithm over the part of the vector held by each process. After that, each processes communicates with all other processes to test whether they are sorted or not with their neighboring processes. Otherwise, the processes send their lowest numbers to the process on left and receive the highest values from the process on the right. This steps is repeated until the vector is sorted.

According to the results seen in Figure 1, the speedup of this implementation is deeply related to the percentage of numbers that are shared between the processes. It is important to note that we cannot shared more than 50% of the

vector between processes, since processes other than 0 and $processes - 1$ will send and receive numbers from the left and right.

## Broadcast optimization

The first optimization was to reduce the amount of broadcast messages sent by processes to check whether they are all sorted or not when compared with the other processes. Therefore, as soon as a broadcast message returns that a process is not sorted in relation to the left neighbor, we stop sending and receiving broadcast messages.

As shown in Figure 1, the results using this optimization, compared to results without optimization, were almost the same. Even when the percentage of numbers shared between the processes were higher, which results in faster convergence, the reduction in broadcast messages was not significant.

## Interleave instead of bubble sort

The second optimization was focused on reducing the use of the bubble sort algorithm as much as possible, since this algorithm has a time complexity of $O(n)$. Thus, as we can interleave the vector with the part of the vector shared between the processes, it is possible to execute the bubble sort algorithm only once in each process. To use interleaving instead of bubble sort, we first interleave the vector with the percentage of numbers received from right, ignoring the percentage of numbers received from the left. Finally, we interleave the ignored part with the entire vector. It is not possible to call just once the interleave method, as we can end up with three vector pieces in a single process.

This optimization lead to a speedup higher than the optimal speedup. Different from the implementations discussed previously, this optimization has a higher speedup when using hyper-threading. As the interleave algorithm has a linear time complexity, the percentage of items shared between the processes does not have a high impact on speedup or efficiency.

## Broadcast and interleave

We also tested using both optimization discussed above at the same time. The results showed that this can lead to almost the same results as using just the interleaving optimization.

## Discussion

Finally, the broadcast messages, even when widely required, have no significant impact on the domain decomposition performance with MPI. The main bottleneck in this process is the bubble sort algorithm.

Comparing the domain decomposition, with its better optimization, with the divide and conquer approach shows that divide and conquer has a much better speedup and efficiency when using hyper-threading.

# Bubble Sort Source Code

Listing 1: Dataset generator

```
1   #include <stdio.h>
2
3   int *get_vector(int vector_size, int *v) {
4       for (int i = 0; i < vector_size; i++) {
5           v[i] = vector_size - i;
6       }
7       return v;
8   }
9
10  int *get_vector_offset(int vector_size, int subvector_size, int *v,
11                         int offset) {
12      int init = subvector_size * offset;
13      int end = init + subvector_size;
14      for (int i = init; i < end; i++) {
15          v[i - init] = vector_size - i;
16      }
17      return v;
18  }
```

Listing 2: Bubble Sort Sequential

```
1   #include "dataset-generator.h"
2   #include <stdio.h>
3   #include <stdlib.h>
4   #include <sys/time.h>
5
6   float time_difference_msec(struct timeval t0, struct timeval t1) {
7       return (t1.tv_sec - t0.tv_sec) * 1000.0f +
8              (t1.tv_usec - t0.tv_usec) / 1000.0f;
9   }
10
11  int is_vector_sorted(int *vector, int vector_size) {
12      int last_element = -1;
13      for (int i = 0; i < vector_size; i++) {
14          if (last_element > vector[i]) {
15              return 0;
16          }
17          last_element = vector[i];
18      }
19      return 1;
20  }
21
22  int *bubble_sort(int vector_size, int *vector_unsorted) {
23      int swapped = 1;
24
25      int i = 0;
26      while ((i < (vector_size - 1)) & swapped) {
27          swapped = 0;
28          for (int j = 0; j < vector_size - i - 1; j++)
29              if (vector_unsorted[j] > vector_unsorted[j + 1]) {
30                  int temp = vector_unsorted[j];
31                  vector_unsorted[j] = vector_unsorted[j + 1];
32                  vector_unsorted[j + 1] = temp;
33                  swapped = 1;
34              }
35          i++;
36      }
37
38      return vector_unsorted;
39  }
40
41  int main(int argc, char **argv) {
42      int vector_size = atoi(argv[1]);
43
44      int vector_unsorted[vector_size];
45      get_vector(vector_size, vector_unsorted);
46
47      struct timeval t0;
```

```
48      struct timeval t1;
49
50      gettimeofday(&t0, 0);
51      int *vector_sorted = bubble_sort(vector_size, vector_unsorted);
52      gettimeofday(&t1, 0);
53
54      float total_time = time_difference_msec(t0, t1);
55
56      printf("Vector sorted: %d\n", is_vector_sorted(vector_sorted, vector_size));
57      printf("Vector size: %d\n", vector_size);
58      printf("Time sort (ms): %f\n", total_time);
59
60      return 0;
61  }
```

Listing 3: Bubble Sort MPI

```
1  // TO DO
2  // - Count iterations number.
3
4  #include "dataset-generator.h"
5  #include <mpi.h>
6  #include <stdio.h>
7  #include <stdlib.h>
8  #include <sys/time.h>
9  #define DEBUG 0
10 #define OPTIMIZE_BROADCAST 1
11 #define BUBBLE_SORT_ONLY_ONCE 1
12
13 void interleave_vector(int *vector, int a_init_index, int a_end_index,
14                        int b_init_index, int b_end_index, int *result) {
15     if (a_end_index - a_init_index > b_end_index - b_init_index) {
16         int temp = a_init_index;
17         a_init_index = b_init_index;
18         b_init_index = temp;
19         temp = a_end_index;
20         a_end_index = b_end_index;
21         b_end_index = temp;
22     }
23
24     int i = 0, a = a_init_index, b = b_init_index;
25     while (a <= a_end_index  b <= b_end_index) {
26         if (a <= a_end_index && (vector[a] < vector[b]  b > b_end_index)) {
27             result[i] = vector[a];
28             a++;
29         } else {
30             result[i] = vector[b];
31             b++;
32         }
33         i++;
34     }
35 }
36
37 void bubble_sort(int vector_size, int *vector_unsorted) {
38     int swapped = 1;
39
40     int i = 0;
41     while ((i < (vector_size - 1)) & swapped) {
42         swapped = 0;
43         for (int j = 0; j < vector_size - i - 1; j++)
44             if (vector_unsorted[j] > vector_unsorted[j + 1]) {
45                 int temp = vector_unsorted[j];
46                 vector_unsorted[j] = vector_unsorted[j + 1];
47                 vector_unsorted[j + 1] = temp;
48                 swapped = 1;
49             }
50         i++;
51     }
52 }
53
54 int is_vector_sorted(int *vector, int vector_size) {
55     int last_element = -1;
```

```c
56      for (int i = 0; i < vector_size; i++) {
57          if (last_element > vector[i]) {
58              return 0;
59          }
60          last_element = vector[i];
61      }
62      return 1;
63  }
64
65  int main(int argc, char **argv) {
66      int vector_size = atoi(argv[1]);
67      float percentage_items_exchange = atof(argv[2]);
68
69      MPI_Status status;
70      int my_rank;
71      int num_processes;
72
73      MPI_Init(&argc, &argv);
74      MPI_Comm_rank(MPI_COMM_WORLD, &my_rank);
75      MPI_Comm_size(MPI_COMM_WORLD, &num_processes);
76
77      if (vector_size % num_processes) {
78          printf("Error: vector size must be a multiple of the number of "
79                  "processes.\n");
80          MPI_Abort(MPI_COMM_WORLD, -1);
81      }
82
83      double t0;
84      double t1;
85      t0 = MPI_Wtime();
86
87      int subvector_size = vector_size / num_processes;
88      int number_items_shared = subvector_size * percentage_items_exchange;
89      int subvector[subvector_size + number_items_shared];
90      get_vector_offset(vector_size, subvector_size, subvector, my_rank);
91
92      // Avoid allocating a new vector in each iteration.
93      int process_status[num_processes];
94      int interleaved_shared_vector_pre_defined[number_items_shared * 2];
95  #if BUBBLE_SORT_ONLY_ONCE == 1
96      int interleaved_subvector_pre_defined[subvector_size];
97      int bubble_sort_executed = 0;
98  #endif
99
100 #if DEBUG == 1
101     int loop_index = 0;
102 #endif
103
104     int done = 0;
105     while (!done) {
106         // --------------------FIRST PHASE--------------------
107 #if BUBBLE_SORT_ONLY_ONCE == 1
108         if (!bubble_sort_executed) {
109             bubble_sort(subvector_size, subvector);
110             bubble_sort_executed = 1;
111         } else {
112             interleave_vector(subvector, number_items_shared,
113                             subvector_size - number_items_shared - 1,
114                             subvector_size - number_items_shared,
115                             subvector_size - 1,
116                             interleaved_subvector_pre_defined);
117         for (int i = 0; i < subvector_size - number_items_shared; i++) {
118             subvector[i + number_items_shared] =
119                 interleaved_subvector_pre_defined[i];
120         }
121
122         interleave_vector(subvector, 0, number_items_shared - 1,
123                         number_items_shared, subvector_size - 1,
124                         interleaved_subvector_pre_defined);
125         for (int i = 0; i < subvector_size; i++) {
126             subvector[i] = interleaved_subvector_pre_defined[i];
```

```c
127                }
128            }
129    #else
130            bubble_sort(subvector_size, subvector);
131    #endif
132            // -------------------FIRST PHASE-------------------
133
134            // -------------------SECOND PHASE-------------------
135            if (my_rank < num_processes - 1) {
136                MPI_Send(&subvector[subvector_size - 1], 1, MPI_INT, my_rank + 1, 0,
137                        MPI_COMM_WORLD);
138            }
139            if (my_rank > 0) {
140                int largest_number_left;
141                MPI_Recv(&largest_number_left, 1, MPI_INT, my_rank - 1, 0,
142                        MPI_COMM_WORLD, &status);
143                process_status[my_rank] = subvector[0] > largest_number_left;
144            } else {
145                process_status[my_rank] = 1;
146            }
147
148            for (int i = 0; i < num_processes; i++) {
149                MPI_Bcast(&process_status[i], 1, MPI_INT, i, MPI_COMM_WORLD);
150    #if OPTIMIZE_BROADCAST == 1
151                if (!process_status[i]) {
152                    break;
153                }
154    #endif
155            }
156
157            int vector_sorted = 1;
158            for (int i = 0; i < num_processes; i++) {
159                vector_sorted &= process_status[i];
160                if (!vector_sorted) {
161                    break;
162                }
163            }
164            if (vector_sorted) {
165                done = 1;
166                break;
167            }
168            // -------------------SECOND PHASE-------------------
169
170            // -------------------THIRD PHASE-------------------
171            if (my_rank > 0) {
172                MPI_Send(&subvector[0], number_items_shared, MPI_INT, my_rank - 1,
173                        0, MPI_COMM_WORLD);
174            }
175            if (my_rank < num_processes - 1) {
176                MPI_Recv(&subvector[subvector_size], number_items_shared, MPI_INT,
177                        my_rank + 1, 0, MPI_COMM_WORLD, &status);
178
179                interleave_vector(subvector, subvector_size - number_items_shared,
180                            subvector_size - 1, subvector_size,
181                            subvector_size + number_items_shared - 1,
182                            interleaved_shared_vector_pre_defined);
183
184                MPI_Send(
185                    &interleaved_shared_vector_pre_defined[number_items_shared],
186                    number_items_shared, MPI_INT, my_rank + 1, 0, MPI_COMM_WORLD);
187
188                for (int i = 0; i < number_items_shared; i++) {
189                    subvector[subvector_size - number_items_shared + i] =
190                        interleaved_shared_vector_pre_defined[i];
191                }
192            }
193            if (my_rank > 0) {
194                MPI_Recv(&subvector[0], number_items_shared, MPI_INT, my_rank - 1,
195                        0, MPI_COMM_WORLD, &status);
196            }
197            // -------------------THIRD PHASE-------------------
```

```c
#if DEBUG == 1
        if (loop_index == 0) {
            for (int i = 0; i < subvector_size; i++) {
                printf("%d ", subvector[i]);
            }
            printf("----- %d", my_rank);
            printf("\n");
            fflush(stdout);
        }
        loop_index++;
#endif
    }

    t1 = MPI_Wtime();
    double total_time = (t1 - t0) * 1000;

#if DEBUG == 1
    if (my_rank > 0) {
        MPI_Send(&subvector[0], subvector_size, MPI_INT, 0, 0, MPI_COMM_WORLD);
    } else {
        int final_vector[vector_size];
        for (int i = 0; i < subvector_size; i++) {
            final_vector[i] = subvector[i];
        }

        for (int i = 1; i < num_processes; i++) {
            int subvector_received[subvector_size];
            MPI_Recv(&subvector_received[0], subvector_size, MPI_INT, i, 0,
                    MPI_COMM_WORLD, &status);

            for (int j = 0; j < subvector_size; j++) {
                final_vector[subvector_size * i + j] = subvector_received[j];
            }
        }

        printf("Final vector:\n");
        for (int i = 0; i < vector_size; i++) {
            printf("%d ", final_vector[i]);
        }
        printf("\n");

        printf("DEBUG: %d\n", DEBUG);
        printf("OPTIMIZE_BROADCAST: %d\n", OPTIMIZE_BROADCAST);
        printf("BUBBLE_SORT_ONLY_ONCE: %d\n", BUBBLE_SORT_ONLY_ONCE);
        printf("Vector sorted: %d\n",
                is_vector_sorted(subvector, subvector_size));
        printf("Vector size: %d\n", vector_size);
        printf("Time sort (ms): %f\n", total_time);
    }
#else
    if (my_rank == 0) {
        printf("DEBUG: %d\n", DEBUG);
        printf("OPTIMIZE_BROADCAST: %d\n", OPTIMIZE_BROADCAST);
        printf("BUBBLE_SORT_ONLY_ONCE: %d\n", BUBBLE_SORT_ONLY_ONCE);
        printf("Vector sorted: %d\n",
                is_vector_sorted(subvector, subvector_size));
        printf("Vector size: %d\n", vector_size);
        printf("Time sort (ms): %f\n", total_time);
    }
#endif

    MPI_Finalize();

    return 0;
}
```