

4. Algorithm Parallelization and Optimization on Distributed Memory Architectures

Team : parlab03

Ioannis Palaaios - el18184

I. INTRODUCTION

We were assigned to develop parallel and optimized versions of the K-Means algorithm and the Heat-Transfer in two dimensions (only the Jacobi Method) for a distributed memory architecture using MPI as the message-passing standard. For the K-means algorithm, we were assigned to benchmark it for a different number of MPI processes, while for the Jacobi method of the heat-transfer in two dimensions, we were assigned to benchmark it for the cases where we perform a convergence check and when we do not, again for different number of MPI processes and board sizes.

II. MPI PARALLELIZATION OF THE K-MEANS ALGORITHM

As we seen in the previous assignments, the compute intensive part of the K-means algorithm happens on the calculation of the new clusters in each iteration of the do while(). In order to efficiently parallelize this part of the algorithm, we split the objects array onto the MPI processes. This means that each process is only assigned to calculate the new distances, cluster memberships and the convergence variable for the specific objects it was assigned. After each iteration, each process has to communicate with *all* others in order to renew the *newClusters[]* and *newClusterSize[]* arrays and then each process calculates the average sum of each cluster with the new memberships. Then, each process communicates with *all* the other processes in order to get the total sum of the convergence delta variable and then divide it with the number of objects.

In order to assign different objects to each different process in the start of the program, the root MPI process (with the rank equal to 0) is used to calculate the new objects (stored in the *objects[]* array in the *file_io.c* file) and then send different parts of this array to different processes using MPI. For this exact reason, we will use the *MPI_Scatter()* function which achieves exactly this. For this function, we have to calculate the displacements of each data sent (the offset from the start - *displs[]* array) and the number of data sent to each process ((*sendcounts[]*) array). The process these arrays are calculated is self-evident and is shown in Listing 1, along with others in the *file_io.c* file, solely defining the *dataset_generation()* function.

```
1 double * dataset_generation(int numObjs, int numCoords, long *rank_numObjs)
2 {
3     double * objects = NULL, * rank_objects = NULL;
4     long i, j, k;
5
6     double val_range = 10;
7
8     int rank, size;
```

```
9     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
10    MPI_Comm_size(MPI_COMM_WORLD, &size);
11
12    if(numObjs%size == 0)
13        *rank_numObjs = numObjs/size;
14    else
15        *rank_numObjs = numObjs/size + 1;
16
17    int sendcounts[size], displs[size];
18    if (rank == 0) {
19        objects = (typeof(objects)) malloc(numObjs * numCoords * sizeof(*objects));
20        for(i = 0; i < size; i++) {
21            int dim = (*rank_numObjs)*numCoords;
22            displs[i] = dim*i;
23            if(dim*(i+1) >= numObjs*numCoords)
24                sendcounts[i] = (numObjs*numCoords) - dim*i;
25            else
26                sendcounts[i] = dim;
27        }
28    }
29
30    MPI_Bcast(sendcounts, size, MPI_INT, 0, MPI_COMM_WORLD);
31    MPI_Bcast(displs, size, MPI_INT, 0, MPI_COMM_WORLD);
32
33    rank_objects = (typeof(rank_objects)) malloc(sendcounts[rank] * sizeof(*
rank_objects));
34
35    if (rank == 0) {
36        for (i=0; i<numObjs; i++)
37        {
38            unsigned int seed = i;
39            for (j=0; j<numCoords; j++)
40            {
41                objects[i*numCoords + j] = (rand_r(&seed) / ((double) RAND_MAX)) *
val_range;
42                if (_debug && i == 0)
43                    printf("object [i=%ld][j=%ld]=%f\n",i,j,objects[i*numCoords + j])
44            }
45        }
46    }
47
48    int recvsize = (*rank_numObjs)*numCoords;
49    MPI_Scatterv(objects, sendcounts, displs, MPI_DOUBLE, rank_objects, recvsize,
MPI_DOUBLE, 0, MPI_COMM_WORLD);
50    if(rank == 0)
51        free(objects);
52
53    return rank_objects;
54 }
```

Listing 1: The source code of file_io.c. In this file the `dataset_generation()` method is defined which essentially generates the new objects and scatters them to all the different processes.

Then the `dataset_generation()` function is called in the main function (`main.c`) and places the objects to each corresponding MPI process. Later in `main.c`, the `clusters[]` array is allocated and initialized in the root process and then broadcasted. Then, the `kmeans()` function is called, where the computations of the K-means algorithm happens and the clusters array is already calculated in each process so no other communication needs to happen for the clusters (the communication happens in the `kmeans()` function) but the root process has to gather the total

membership information from all the other processes, so we have to calculate the number of elements to receive from each other rank and the displacement of each rank's data (*recvcounts[]* and *displs[]* respectively). The important part of the source code of *main.c* is shown on Listing 2 below.

```
1 int main(int argc, char **argv) {
2     // variable declarations etc.
3     ...
4
5     MPI_Init(&argc, &argv);
6     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
7     MPI_Comm_size(MPI_COMM_WORLD, &size);
8
9     numObjs = (dataset_size * 1024 * 1024) / (numCoords * sizeof(double));
10
11     objects = dataset_generation(numObjs, numCoords, &rank_numObjs);
12
13     clusters = (double*) malloc(numClusters * numCoords * sizeof(double));
14
15     if (rank == 0) {
16         for (i=0; i<numClusters; i++)
17             for (j=0; j<numCoords; j++)
18                 clusters[i*numCoords + j] = objects[i*numCoords + j];
19
20         // check initial cluster centers for repetition
21         if (check_repeated_clusters(numClusters, numCoords, clusters) == 0) {
22             printf("Error: some initial clusters are repeated. Please select
distinct initial centers\n");
23             MPI_Finalize();
24             return 1;
25         }
26
27         /*
28         printf("Initial cluster centers:\n");
29         for (i=0; i<numClusters; i++) {
30             printf("(0) clusters[%ld] =", i);
31             for (j=0; j<numCoords; j++)
32                 printf(" %6.6f", clusters[i*numCoords + j]);
33             printf("\n");
34         } */
35
36     }
37
38     MPI_Bcast(clusters, numClusters*numCoords, MPI_DOUBLE, 0, MPI_COMM_WORLD);
39
40     membership = (int*) malloc(rank_numObjs * sizeof(int));
41     tot_membership = (int*) malloc(numObjs * sizeof(int));
42
43     double start_time = wtime();
44     kmeans(objects, numCoords, rank_numObjs, numClusters, threshold, loop_threshold,
membership, clusters);
45     double time = wtime() - start_time;
46
47     if (rank == 0) {
48         printf("Final cluster centers:\n");
49         for (i=0; i<2; i++) {
50             printf("clusters[%ld] = ", i);
51             for (j=0; j<numCoords; j++)
52                 printf("%6.6f ", clusters[i*numCoords + j]);
53             printf("\n");
54         }
55     }
```

```

54     }
55     printf("\n%d MPI Processes | Time Spent : %7.4f sec\n", size, time);
56 }
57
58 int recvcunts[size], displs[size];
59 if (rank == 0) {
60     for(i = 0; i < size; i++) {
61         long dim = rank_numObjs;
62         displs[i] = dim*i;
63         if(dim*(i+1) >= numObjs)
64             recvcunts[i] = numObjs - dim*i;
65         else
66             recvcunts[i] = dim;
67     }
68 }
69
70 MPI_Bcast(recvcunts, size, MPI_INT, 0, MPI_COMM_WORLD);
71 MPI_Bcast(displs, size, MPI_INT, 0, MPI_COMM_WORLD);
72
73 MPI_Gatherv(membership, rank_numObjs, MPI_INT, tot_membership, recvcunts,
74             displs, MPI_INT, 0, MPI_COMM_WORLD);
75
76 // free the dynamically allocated arrays
77 ...
78 return 0;
79 }

```

Listing 2: The source code of the main.c file.

Now, we have to perform the calculations of the K-means algorithm. As we mentioned above, each process is responsible for calculating the distances from the current clusters etc. only for the objects it was assigned. The only reasons it must communicate with the other process is the renew of the *newClusters[]* and *newClusterSize[]* arrays and the *delta* variable. As we have to calculate the sum of the each of the corresponding element of each rank's local *rank_newClusters[]* and *rank_newClusterSize[]* arrays and then store the result on the *newClusters[]* and *newClusterSize[]* arrays respectively and then have them stored in each process (in order for each process to be able to calculate the same *clusters[]* array), we will use the *MPI_Allreduce()* function which does exactly that. The do while() of the *kmeans.c* file is shown on Listing 3 below.

```

1  do {
2      for (i=0; i<numClusters; i++) {
3          for (j=0; j<numCoords; j++) {
4              rank_newClusters[i*numCoords + j] = 0.0;
5              newClusters[i*numCoords + j] = 0.0;
6          }
7          rank_newClusterSize[i] = 0;
8          newClusterSize[i] = 0;
9      }
10     delta = 0.0;
11     rank_delta = 0.0;
12
13     for (i=0; i<numObjs; i++) {
14         index = find_nearest_cluster(numClusters, numCoords, &objects[i*numCoords],
15                                     clusters);
16
17         if (membership[i] != index)
18             rank_delta += 1.0;
19
20         membership[i] = index;

```

```
21     rank_newClusterSize[index]++;
22     for (j=0; j<numCoords; j++)
23         rank_newClusters[index*numCoords + j] += objects[i*numCoords + j];
24 }
25
26 MPI_Allreduce(rank_newClusters, newClusters, numClusters*numCoords, MPI_DOUBLE,
27 MPI_SUM, MPI_COMM_WORLD);
28 MPI_Allreduce(rank_newClusterSize, newClusterSize, numClusters, MPI_INT, MPI_SUM,
29 MPI_COMM_WORLD);
30
31 for (i=0; i<numClusters; i++) {
32     if (newClusterSize[i] > 0) {
33         for (j=0; j<numCoords; j++) {
34             clusters[i*numCoords + j] = newClusters[i*numCoords + j] /
35             newClusterSize[i];
36         }
37     }
38 }
39 MPI_Allreduce(&rank_delta, &delta, 1, MPI_DOUBLE, MPI_SUM, MPI_COMM_WORLD);
40 delta /= numObjs;
41 loop++;
42 } while (delta > threshold && loop < loop_threshold);
```

Listing 3: The source code of the do while() in the kmeans.c file which is responsible for the computations (and communications) of the K-means algorithm.

We benchmark our MPI implementation for the configuration {Size, Coords, Clusters, Loops} = {256, 16, 16, 10} and for different numbers of MPI processes, such as 1, 2, 4, 8, 16, 32 and 64. As we observe on the Figures 1 and 2 below, as we increase the number of MPI processes, we achieve significant performance benefits, almost in accordance with the number of Processes we create. The reason we achieve this performance scaling as we increase the number of processes is that the time spent for the communication between these processes is likely a lot less than the performance benefits of the parallelization of the compute intensive part of the algorithm. Another benefit in comparison with the shared memory architecture implementation (where we used OpenMP), is that each processes has a different memory address space and thus operates on its own data reducing memory conflicts, cache misses, cache contention problems, false sharing and other shared memory performance degrading phenomena. Another reason is that these MPI processes we create, are all on the same node which means that the communication between them, has lower latency and higher bandwidth than the case where these processes were located on another node. Of course we do not expect similar performance benefits (or scaling) when we create more processes than our node's virtual cores as context switching between them (by the Operating System) will likely introduce a large overhead. As we mentioned on the second assignment, the node we run these experiment has 32 physical cores and 64 virtual cores (by supporting hyper-threading). In the OpenMP implementation, we did not see sufficient scaling between increasing the number of threads from 32 to 64 as each of these threads (and specifically the threads occupying the same physical core) were using the same parts of the physical core and thus hyper-threading could not add performance benefits. This phenomenon does not happen with OpenMP as processes are likely having some benefits as resource isolation leading to less contention of shared resources and different memory access patterns compared to multi thread implementations.

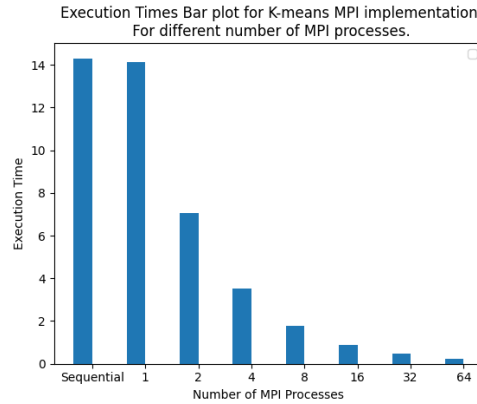


Figure 1: Bar Plot of the execution times of the K-Means Clustering Algorithm MPI implementation for the Configuration {Size, Coords, Clusters, Loops} = {256, 16, 16, 10} for different numbers of MPI processes.

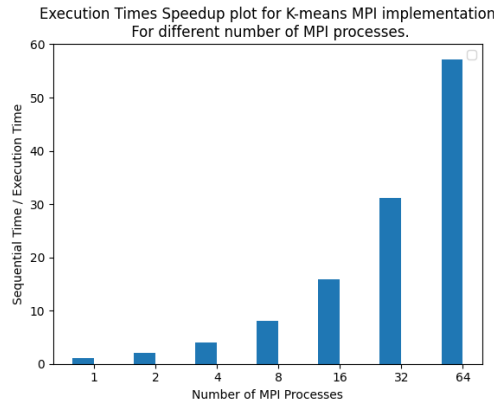


Figure 2: Speedup Plot of the execution times of the K-Means Clustering Algorithm MPI implementation for the Configuration {Size, Coords, Clusters, Loops} = {256, 16, 16, 10} for different numbers of MPI processes.

III. HEAT TRANSFER - JACOBI METHOD

For this version, we had to solve the 2D Heat equation using the Jacobi Method. We want to parallelize the compute intensive part of the heat equation using multiple MPI processes. In order to do that, we have to split the data to the multiple processes with each process taking a 'block' and then performing the calculations of the temperature of each block separately. But, when we want to calculate the temperature change of this specific block, we have to take the temperature values of its neighbouring blocks (north, south, west and east if they exist) which are stored in the 'neighbouring' processes and as a result processes must communicate in each step of the calculating for loop.

First of all, the root process calls the *init2d()* function in order to initialize the double array containing the temperatures in each cell (only the boundary values have a temperature different than 0). Then, we have to send to each process their corresponding blocks. For this splitting we

use the *MPI_Scatterv()* function which send different data to each process. For many reasons including code readability etc, we assign each process to an MPI Cartesian Communicator which gives each process a 2 dimensional rank, thus assigning it to a specific block. For similar reasons, specific data types are defined for the *MPI_Scatterv()* call. The source code of this specific part is shown on Listing 4 below.

```
1  if (rank==0) {
2      U=allocate2d(global_padded[0],global_padded[1]);
3      init2d(U,global[0],global[1]);
4  }
5
6  u_previous=allocate2d(local[0]+2,local[1]+2);
7  u_current=allocate2d(local[0]+2,local[1]+2);
8
9  MPI_Datatype global_block;
10 MPI_Type_vector(local[0],local[1],global_padded[1],MPI_DOUBLE,&dummy);
11 MPI_Type_create_resized(dummy,0,sizeof(double),&global_block);
12 MPI_Type_commit(&global_block);
13
14 MPI_Datatype local_block;
15 MPI_Type_vector(local[0],local[1],local[1]+2,MPI_DOUBLE,&dummy);
16 MPI_Type_create_resized(dummy,0,sizeof(double),&local_block);
17 MPI_Type_commit(&local_block);
18
19 int * scatteroffset, * scattercounts;
20 if (rank==0) {
21     scatteroffset=(int*)malloc(size*sizeof(int));
22     scattercounts=(int*)malloc(size*sizeof(int));
23     for (i=0;i<grid[0];i++)
24         for (j=0;j<grid[1];j++) {
25             scattercounts[i*grid[1]+j]=1;
26             scatteroffset[i*grid[1]+j]=(local[0]*local[1]*grid[1]*i+local[1]*j);
27         }
28 }
29 MPI_Scatterv(&(U[0][0]), scattercounts, scatteroffset, global_block, &(u_previous
    [1][1]), 1, local_block, 0, CART_COMM);
30
31 for(i = 0; i < local[0]+2; i++) {
32     for(j = 0; j < local[1]+2; j++) {
33         u_current[i][j] = u_previous[i][j];
34     }
35 }
```

Listing 4: The source code of most of the matrix initializations and the initial scatter.

Then, when the blocks are sent to their corresponding processes, each process has to find the coordinates of it's neighbours while also taking into account the possibility that the process is in the boundary of the global block and does not have one or two neighbours. Then, each block has to find the ranks of the neighbouring processes (from their coordinates) in order to communicate with them. The source code of this specific part is shown on Listing 5 below.

```
1  int row_idx = rank_grid[0], col_idx = rank_grid[1];
2  int *num_of_rows = (int*) malloc(sizeof(int));
3  int *num_of_cols = (int*) malloc(sizeof(int));
4  *num_of_rows = local[0];
5  *num_of_cols = local[1];
6  int *row_size = (int*) malloc(sizeof(int));
7  int *col_size = (int*) malloc(sizeof(int));
8  *row_size = local[0]+2;
9  *col_size = local[1]+2;
```

```
10
11 // send element arrays
12 double *north_arr = (double*) malloc(*num_of_cols * sizeof(double));
13 double *south_arr = (double*) malloc(*num_of_cols * sizeof(double));
14 double *east_arr = (double*) malloc(*num_of_rows * sizeof(double));
15 double *west_arr = (double*) malloc(*num_of_rows * sizeof(double));
16
17 // receive element arrays
18 double *north_arr_recv = (double*) malloc((*num_of_cols) * sizeof(double));
19 double *south_arr_recv = (double*) malloc((*num_of_cols) * sizeof(double));
20 double *east_arr_recv = (double*) malloc((*num_of_rows) * sizeof(double));
21 double *west_arr_recv = (double*) malloc((*num_of_rows) * sizeof(double));
22
23 int north, south, east, west;
24 north = row_idx - 1;
25 south = row_idx + 1;
26 east = col_idx + 1;
27 west = col_idx - 1;
28
29 int i_min, i_max, j_min, j_max;
30
31 /*Fill your code here*/
32
33 // if is west boundary
34 if(west < 0) {
35     j_min = 2;
36 } else {
37     j_min = 1;
38 }
39
40 // if is east boundary
41 if(east >= grid[1]) {
42     j_max = *num_of_cols-1;
43 } else {
44     j_max = *num_of_cols;
45 }
46
47 // if is north boundary
48 if(north < 0) {
49     i_min = 2;
50 } else {
51     i_min = 1;
52 }
53
54 // if is south boundary
55 if(south >= grid[0]) {
56     i_max = *num_of_rows-1;
57 } else {
58     i_max = *num_of_rows;
59 }
60
61 int *cart_coords_north = (int*) malloc(2 * sizeof(int));
62 int *cart_coords_south = (int*) malloc(2 * sizeof(int));
63 int *cart_coords_west = (int*) malloc(2 * sizeof(int));
64 int *cart_coords_east = (int*) malloc(2 * sizeof(int));
65
66 cart_coords_north[0] = row_idx-1;
67 cart_coords_north[1] = col_idx;
68 cart_coords_south[0] = row_idx+1;
69 cart_coords_south[1] = col_idx;
70 cart_coords_east[0] = row_idx;
```



```
71 cart_coords_east[1] = col_idx+1;
72 cart_coords_west[0] = row_idx;
73 cart_coords_west[1] = col_idx-1;
74
75 //*****//
76
77 MPI_Status *north_stats = (MPI_Status*) malloc(2 * sizeof(MPI_Status));
78 MPI_Status *south_stats = (MPI_Status*) malloc(2 * sizeof(MPI_Status));
79 MPI_Request *north_reqs = (MPI_Request*) malloc(2 * sizeof(MPI_Request));
80 MPI_Request *south_reqs = (MPI_Request*) malloc(2 * sizeof(MPI_Request));
81 MPI_Status *east_stats = (MPI_Status*) malloc(2 * sizeof(MPI_Status));
82 MPI_Status *west_stats = (MPI_Status*) malloc(2 * sizeof(MPI_Status));
83 MPI_Request *east_reqs = (MPI_Request*) malloc(2 * sizeof(MPI_Request));
84 MPI_Request *west_reqs = (MPI_Request*) malloc(2 * sizeof(MPI_Request));
85
86 int *tags = (int*) malloc(4 * sizeof(int));
87 for(i = 0; i < 4; i++)
88     tags[i] = 0;
89 int *c = (int*) malloc(1 * sizeof(int));
90
91 int *cart_rank_north = (int*) malloc(1 * sizeof(int));
92 int *cart_rank_south = (int*) malloc(1 * sizeof(int));
93 int *cart_rank_east = (int*) malloc(1 * sizeof(int));
94 int *cart_rank_west = (int*) malloc(1 * sizeof(int));
95
96 *cart_rank_north = -1;
97 *cart_rank_south = -1;
98 *cart_rank_east = -1;
99 *cart_rank_west = -1;
100
101 int coord_var = grid[0];
102 if(cart_coords_north[0] >= 0) {
103     *cart_rank_north = coord_var*cart_coords_north[0] + cart_coords_north[1];
104 }
105 if(cart_coords_south[0] < grid[0]) {
106     *cart_rank_south = coord_var*cart_coords_south[0] + cart_coords_south[1];
107 }
108 if(cart_coords_west[1] >= 0) {
109     *cart_rank_west = coord_var*cart_coords_west[0] + cart_coords_west[1];
110 }
111 if(cart_coords_east[1] < grid[1]) {
112     *cart_rank_east = coord_var*cart_coords_east[0] + cart_coords_east[1];
113 }
```

Listing 5: The source code of the assignment of coordinates and ranks to the neighbouring processes etc.

After, we have the neighbouring information, we move onto the compute intensive part which is the for loop (of default size equal to 256 iterations). We use non blocking Sends and Receives (*MPI_Isend()* and *MPI_Irecv()* respectively) and then we use the *MPI_Waitall()* method in order to make sure each process has received its corresponding boundary parts. Then, we perform the convergence check in each process and if all processes have converged then the for loop finishes. We use *MPI_Allreduce()* and the MPI operation *MPI_MIN*, reducing the converged value with the minimum operation in order to make sure that if one process block has not converged (*converged* variable equal to 0), all the other processes/blocks have not converged either. The source code of the for loop is shown on Listing 7 below.

```
1 gettimeofday(&tts, NULL);
2 #ifdef TEST_CONV
```

```
3  for (t=0;t<T && !global_converged;t++) {
4  #endif
5  #ifndef TEST_CONV
6  #undef T
7  #define T 256
8  for (t=0;t<T;t++) {
9  #endif
10 // initialize boundary arrays
11 for(i = 1; i < (*num_of_cols)+1; i++) {
12     north_arr[i-1] = u_previous[1][i];
13     south_arr[i-1] = u_previous[*num_of_rows][i];
14 }
15
16 for(i = 1; i < (*num_of_rows)+1; i++) {
17     west_arr[i-1] = u_previous[i][1];
18     east_arr[i-1] = u_previous[i][*num_of_cols];
19 }
20
21 // sending phase
22 if(*cart_rank_north > -1) {
23     MPI_Isend(north_arr, *num_of_cols, MPI_DOUBLE, *cart_rank_north, tags
24 [0], CART_COMM, &(north_reqs[0]));
25 }
26 if(*cart_rank_south > -1) {
27     MPI_Isend(south_arr, *num_of_cols, MPI_DOUBLE, *cart_rank_south, tags
28 [1], CART_COMM, &(south_reqs[0]));
29 }
30 if(*cart_rank_west > -1) {
31     MPI_Isend(west_arr, *num_of_rows, MPI_DOUBLE, *cart_rank_west, tags[2],
32 CART_COMM, &(west_reqs[0]));
33 }
34 if(*cart_rank_east > -1) {
35     MPI_Isend(east_arr, *num_of_rows, MPI_DOUBLE, *cart_rank_east, tags[3],
36 CART_COMM, &(east_reqs[0]));
37 }
38
39 // receiving phase
40 if(*cart_rank_north > -1) {
41     MPI_Irecv(north_arr_recv, *num_of_cols, MPI_DOUBLE, *cart_rank_north,
42 tags[0], CART_COMM, &north_reqs[1]);
43 }
44 if(*cart_rank_south > -1) {
45     MPI_Irecv(south_arr_recv, *num_of_cols, MPI_DOUBLE, *cart_rank_south,
46 tags[1], CART_COMM, &south_reqs[1]);
47 }
48 if(*cart_rank_west > -1) {
49     MPI_Irecv(west_arr_recv, *num_of_rows, MPI_DOUBLE, *cart_rank_west, tags
50 [2], CART_COMM, &west_reqs[1]);
51 }
52 if(*cart_rank_east > -1) {
53     MPI_Irecv(east_arr_recv, *num_of_rows, MPI_DOUBLE, *cart_rank_east, tags
54 [3], CART_COMM, &east_reqs[1]);
55 }
56
57 // process synchronization
58 if(*cart_rank_north > -1) {
59     MPI_Waitall(2, north_reqs, north_stats);
60 }
61 if(*cart_rank_south > -1) {
62     MPI_Waitall(2, south_reqs, south_stats);
63 }
64 }
```

```
56     if(*cart_rank_west > -1) {
57         MPI_Waitall(2, west_reqs, west_stats);
58     }
59     if(*cart_rank_east > -1) {
60         MPI_Waitall(2, east_reqs, east_stats);
61     }
62
63     /*Compute and Communicate*/
64
65     /*Add appropriate timers for computation*/
66
67     gettimeofday(&tcs, NULL);
68     for(i = i_min; i < i_max; i++) {
69         for(j = j_min; j < j_max; j++) {
70             u_current[i][j] = (u_previous[i-1][j] + u_previous[i][j-1]+
71             u_previous[i+1][j]+u_previous[i][j+1])/4;
72         }
73     }
74     gettimeofday(&tcf, NULL);
75     tcomp += (tcf.tv_sec-tcs.tv_sec)+(tcf.tv_usec-tcs.tv_usec)*0.000001;
76
77     #ifdef TEST_CONV
78     if (t%C==0) {
79         gettimeofday(&tconvs, NULL);
80         converged=1;
81         for(i = i_min; i < i_max; i++) {
82             for(j = j_min; j < j_max; j++) {
83                 if(abs(u_current[i][j]-u_previous[i][j])>0.01) {
84                     converged=0;
85                     i = i_max;
86                     j = j_max;
87                     break;
88                 }
89             }
90         }
91         gettimeofday(&tconvf, NULL);
92         tconv += (tconvf.tv_sec-tconvs.tv_sec)+(tconvf.tv_usec-tconvs.tv_usec)
93         *0.000001;
94     }
95     #endif
96     for(i = i_min; i < i_max; i++) {
97         for(j = j_min; j < j_max; j++) {
98             u_previous[i][j] = u_current[i][j];
99         }
100     }
101     //*****//
102     MPI_Allreduce(&converged, &global_converged, 1, MPI_DOUBLE, MPI_MIN,
103     CART_COMM);
104 }
```

Listing 6: The source code of the compute intensive for loop.

Finally, in order to gather the local matrices that are stored in each process back to the global array, we use *MPI_Gatherv()* with similar offsets and counts as the original scatter performed in the beginning of the program.

```
1 // identical to the scatter above
2 int * gatheroffset, * gathercounts;
3 if (rank==0) {
4
5     gatheroffset=(int*)malloc(size*sizeof(int));
```

```

6  gathercounts=(int*)malloc(size*sizeof(int));
7
8  for (i=0;i<grid[0];i++)
9      for (j=0;j<grid[1];j++) {
10         gathercounts[i*grid[1]+j]=1;
11         gatheroffset[i*grid[1]+j]=(local[0]*local[1]*grid[1]*i+local[1]*j);
12     }
13     gatheroffset[(grid[0]-1)*grid[1]+grid[1]-1] -= 1;
14 }
15
16 MPI_Bcast(gatheroffset, size, MPI_INT, 0, CART_COMM);
17 MPI_Bcast(gathercounts, size, MPI_INT, 0, CART_COMM);
18
19 MPI_Gatherv(&(u_current[1][1]), 1, local_block, &(U[0][0]), gathercounts,
    gatheroffset, global_block, 0, CART_COMM);

```

Listing 7: The source code of most of the matrix initializations and the initial scatter.

We were assigned to benchmark the Jacobi method for 1024x1024 board size and 64 MPI processes. We also leave the convergence check variable called C at it's default value of 100 although it might be beneficial to reduce it. The execution time of this benchmark is 0.206033, the computation time is 0.025395, the communication time is 0.180638 and time spent on the converge check is 0.000319. It converges after the first 100 iterations and this means that each process gets into the convergence check exactly two times. We observe that the time spent for communication is almost 7x bigger than the time spent for computation. This is reasonable as most processes have to communicate with 4 other processes, in order to receive and send data (degrading performance due to the waiting) , while the computational part is performing a double nested for loop with 16384 iterations.

For the benchmarks with no convergence, we were assigned to benchmark this version for a constant number of 256 iterations for board sizes of 2048x2048, 4096x4096, 6144x6144 and 1, 2, 4, 8, 16, 32 and 64 MPI processes. We observe that our program increases it's performance when we increase the number of MPI Processes. This happens due to that while the workload that each process gets is smaller than when we had less processes, the communication overhead also gets proportionately less as the size of the data that get transferred are also less. The execution time bar plot is shown on Figure 2.

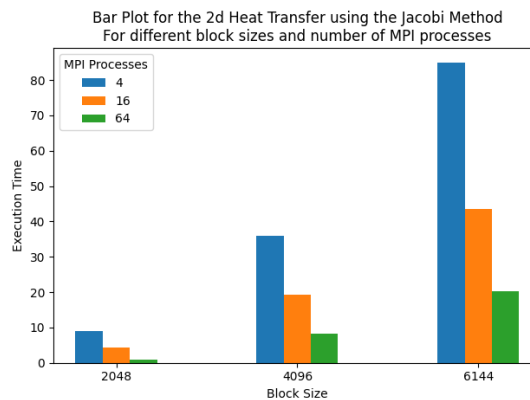


Figure 3: Bar Plot of the execution times of the 2d Heat Transfer with the Jacobi Method using 4, 16 and 64 MPI processes for the block sizes : 2048x2048, 4096x4096, 6144x6144.

We also observe that, as we increase the number of processes the fraction of communication

time divided by the total time gets bigger. Namely, for 4 processes the fraction is around 50%, for 16 processes around 55% and for 64 processes it is around 60%. This happens due to that applying the Jacobi method to smaller blocks decreases the computation time, while the communication time does not decrease as drastically.

IV. CONCLUSION

In this assignment, we observed how different algorithm can be optimized and parallelized in a distributed memory environment, using MPI. First, we parallelized the K-means algorithm with MPI exploring how to efficiently distribute the data to multiple processes and communicate in each iteration in order to reduce an array and make it available to all the processes. We observed significant performance benefits in comparison to the shared memory architecture implementation for reason we also observed and explained. Lastly, we implemented and optimized a parallel version of the 2d Heat Equation using the Jacobi method. After a number of weird program behavior and *a lot of segmentation faults*, we figured out how to efficiently distribute and gather data in a fault tolerant manner while also understanding how inter process communication can be used in a geometric fashion (2d grid) in order to solve a partial differential equation in a parallel and distributed memory fashion. Lastly, we gathered our results and plotted them, visualizing our results and comparing them.