Università
della
Svizzera
italiana

**Institute of
Computing
CI**

**High-Performance Computing Lab**                    **Institute of Computing**

Student: Zitian Wang                    Discussed with: Ning Ding; Xiaorui Wang

# Solution for Project 4

## 1. Task: Ring maximum using MPI

```
1   int main(int argc, char *argv[]) {
2
3   // Initialize MPI, get size and rank
4   int size, rank;
5   MPI_Init(&argc, &argv);
6   MPI_Comm_size(MPI_COMM_WORLD, &size);
7   MPI_Comm_rank(MPI_COMM_WORLD, &rank);
8
9   // IMPLEMENT: Ring sum algorithm
10  int sum = rank; // initialize sum
11  int rank_send = rank;
12  int rank_recive;
13  int send, receive;
14
15  for (int i = 0; i < size - 1; i++){
16    send = (rank + 1) % size;
17    receive = (rank - 1 + size) % size;
18
```

```
19      MPI_Sendrecv(
20      &rank_send, 1, MPI_INT, send, 0,
21      &rank_recive, 1, MPI_INT, receive, 0,
22      MPI_COMM_WORLD, MPI_STATUS_IGNORE);
23
24    sum += rank_recive;
25    rank_send = rank_recive;
26    }
27    printf("Process %i: Sum = %i\n", rank, sum);
28    // Finalize MPI
29    MPI_Finalize();
30    return 0;
31 }
```

Listing 1: Ring sum function

To advoid deadlock issues, using `MPI_Sendrecv` instead of `MPI_Send` and `MPI_Recv`.

```
1      [wangzi@icsnode33 ring]$ mpirun ring_sum
2      Process 0: Sum = 6
3      Process 1: Sum = 6
4      Process 2: Sum = 6
5      Process 3: Sum = 6
```

Listing 2: Output of the Ring Program

## 2. Task: Ghost cells exchange between neighboring processes

### 2.1. Cartesian two-dimensional MPI communicator

```
1      // TODO: set the dimensions of the processor grid and periodic
           boundaries in both dimensions
2      dims[0] = 4;
3      dims[1] = 4;
4      periods[0] = 1;
5      periods[1] = 1;
6
7      // TODO: Create a Cartesian communicator (4*4) with periodic
           boundaries (we do not allow
8      // the reordering of ranks) and use it to find your neighboring
9      // ranks in all dimensions in a cyclic manner.
10     MPI_Cart_create(MPI_COMM_WORLD, 2, dims, periods, 0, &comm_cart);
```

Listing 3: Creating a Cartesian Communicator with Periodic Boundaries

### 2.2. derived data type

```
1      // TODO: find your top/bottom/left/right neighbor using the new
           communicator, see MPI_Cart_shift()
2      // rank_top, rank_bottom
3      // rank_left, rank_right
4      MPI_Cart_shift(comm_cart, 0, 1, &rank_top, &rank_bottom);
5      MPI_Cart_shift(comm_cart, 1, 1, &rank_left, &rank_right);
6
7      // part:1 start
8      // TODO: create derived datatype data_ghost, create a datatype for
           sending the column, see MPI_Type_vector() and MPI_Type_commit()
9      MPI_Type_vector(SUBDOMAIN, 1, DOMAINSIZE, MPI_DOUBLE, &data_ghost);
```

```
10      MPI_Type_commit (& data_ghost );
```

Listing 4: Finding Neighbors and Creating a Derived Datatype

## 2.3. Exchange ghost cells

```
1    // TODO: ghost cell exchange with the neighbouring cells in all
            directions
2    // Part 1: to the top
3    MPI_Irecv (& data [1], DOMAINSIZE - 2, MPI_DOUBLE , rank_top , 0, comm_cart
            , &request );
4    MPI_Send (& data [ DOMAINSIZE + 1], DOMAINSIZE - 2, MPI_DOUBLE ,
            rank_bottom , 0, comm_cart );
5    MPI_Wait (& request , & status );
6
7    // Part 2: to the bottom
8    MPI_Irecv (& data [ DOMAINSIZE * ( DOMAINSIZE - 1) + 1], DOMAINSIZE - 2,
            MPI_DOUBLE , rank_bottom , 1, comm_cart , &request );
9    MPI_Send (& data [ DOMAINSIZE * ( DOMAINSIZE - 2) + 1], DOMAINSIZE - 2,
            MPI_DOUBLE , rank_top , 1, comm_cart );
10   MPI_Wait (& request , & status );
11
12   // Define a custom MPI datatype for a ghost column
13   MPI_Datatype data_ghost_column ;
14   MPI_Type_vector ( SUBDOMAIN , 1, DOMAINSIZE , MPI_DOUBLE , &
            data_ghost_column );
15   MPI_Type_commit (& data_ghost_column );
16
17   // Part 3: to the left
18   MPI_Irecv (& data [ DOMAINSIZE ], 1, data_ghost_column , rank_left , 2,
            comm_cart , &request );
19   MPI_Send (& data [ DOMAINSIZE + 1], 1, data_ghost_column , rank_right , 2,
            comm_cart );
20   MPI_Wait (& request , & status );
21
22   // Part 4: to the right
23   MPI_Irecv (& data [2 * DOMAINSIZE - 1], 1, data_ghost_column , rank_right ,
            3, comm_cart , &request );
24   MPI_Send (& data [2 * DOMAINSIZE - 2], 1, data_ghost_column , rank_left ,
            3, comm_cart );
25   MPI_Wait (& request , & status );
```

Listing 5: Ghost Cell Exchange with Neighboring Cells in All Directions

```
1    [ wangzi@icsnode33 ghost ]$ make
2    mpicc  ghost.c -o ghost
3    [ wangzi@icsnode33 ghost ]$ mpirun -np 16 ./ ghost
4    data of rank 9 after communication
5     9.0   5.0   5.0   5.0   5.0   5.0   5.0   9.0
6     8.0   9.0   9.0   9.0   9.0   9.0   9.0  10.0
7     8.0   9.0   9.0   9.0   9.0   9.0   9.0  10.0
8     8.0   9.0   9.0   9.0   9.0   9.0   9.0  10.0
9     8.0   9.0   9.0   9.0   9.0   9.0   9.0  10.0
10    8.0   9.0   9.0   9.0   9.0   9.0   9.0  10.0
```

Listing 6: Compilation and Execution of Ghost Communication

## 2.4. Bonus

```
// Bonus [10 Points]: Also exchange ghost values with the neighbors in
    ordinal directions (northeast, southeast, southwest and northwest)
    .
int coords[2];
int coords_nw[2], coords_ne[2], coords_sw[2], coords_se[2];
int rank_northwest, rank_northeast, rank_southwest, rank_southeast;
MPI_Cart_coords(comm_cart, rank, 2, coords);
coords_nw[0] = (coords[0] - 1 + dims[0]) % dims[0];
coords_nw[1] = (coords[1] - 1 + dims[1]) % dims[1];
MPI_Cart_rank(comm_cart, coords_nw, &rank_northwest);

coords_ne[0] = (coords[0] - 1 + dims[0]) % dims[0];
coords_ne[1] = (coords[1] + 1) % dims[1];
MPI_Cart_rank(comm_cart, coords_ne, &rank_northeast);

coords_sw[0] = (coords[0] + 1) % dims[0];
coords_sw[1] = (coords[1] - 1 + dims[1]) % dims[1];
MPI_Cart_rank(comm_cart, coords_sw, &rank_southwest);

coords_se[0] = (coords[0] + 1) % dims[0];
coords_se[1] = (coords[1] + 1) % dims[1];
MPI_Cart_rank(comm_cart, coords_se, &rank_southeast);
MPI_Sendrecv(&data[DOMAINSIZE + 1], 1, MPI_DOUBLE, rank_southeast, 4,
             &data[0], 1, MPI_DOUBLE, rank_northwest, 4,
             comm_cart, &status);

MPI_Sendrecv(&data[DOMAINSIZE + DOMAINSIZE - 2], 1, MPI_DOUBLE,
    rank_southwest, 5,
             &data[DOMAINSIZE - 1], 1, MPI_DOUBLE, rank_northeast, 5,
             comm_cart, &status);

MPI_Sendrecv(&data[(DOMAINSIZE - 2) * DOMAINSIZE + 1], 1, MPI_DOUBLE,
    rank_northeast, 6,
             &data[DOMAINSIZE * (DOMAINSIZE - 1)], 1, MPI_DOUBLE,
                 rank_southwest, 6,
             comm_cart, &status);

MPI_Sendrecv(&data[(DOMAINSIZE - 2) * DOMAINSIZE + DOMAINSIZE - 2], 1,
    MPI_DOUBLE, rank_northwest, 7,
             &data[DOMAINSIZE * DOMAINSIZE - 1], 1, MPI_DOUBLE,
                 rank_southeast, 7,
             comm_cart, &status);
```

Listing 7: Exchanging Ghost Values with Ordinal Neighbors

```
[wangzi@icsnode35 ghost]$ mpirun -np 16 ./ghost
data of rank 9 after communication
 4.0   5.0   5.0   5.0   5.0   5.0   5.0   6.0
 8.0   9.0   9.0   9.0   9.0   9.0   9.0  10.0
 8.0   9.0   9.0   9.0   9.0   9.0   9.0  10.0
 8.0   9.0   9.0   9.0   9.0   9.0   9.0  10.0
 8.0   9.0   9.0   9.0   9.0   9.0   9.0  10.0
 8.0   9.0   9.0   9.0   9.0   9.0   9.0  10.0
 8.0   9.0   9.0   9.0   9.0   9.0   9.0  10.0
12.0  13.0  13.0  13.0  13.0  13.0  13.0  14.0
```

Listing 8: Output of Ghost Cell Communication

## 3. Task: Parallelizing the Mandelbrot set using MPI

### 3.1. update functions createPartition, updatePartition and createDomain

```
1   Partition createPartition(int mpi_rank, int mpi_size) {
2       Partition p;
3
4       // TODO: determine size of the grid of MPI processes (p.nx, p.ny),
            see MPI_Dims_create()
5       int dims[2];
6       MPI_Dims_create(mpi_size, 2, dims);
7       p.ny = dims[1];
8       p.nx = dims[0];
9
10      // TODO: Create cartesian communicator (p.comm), we do not allow
            the reordering of ranks here, see MPI_Cart_create()
11      int periods[2] = {0, 0};
12      MPI_Comm comm_cart;
13      MPI_Cart_create(MPI_COMM_WORLD, 2, dims, periods, 0, &comm_cart);
14      p.comm = comm_cart;
15
16      // TODO: Determine the coordinates in the Cartesian grid (p.x, p.y
            ), see MPI_Cart_coords()
17      int coords[2];
18      MPI_Cart_coords(p.comm, mpi_rank, 2, coords);
19      p.x = coords[0];
20      p.y = coords[1];
21
22      return p;
23  }
```

Listing 9: Cartesian Partition Creation

```
1   Partition updatePartition(Partition p_old, int mpi_rank) {
2       Partition p;
3
4       // copy grid dimension and the communicator
5       p.ny = p_old.ny;
6       p.nx = p_old.nx;
7       p.comm = p_old.comm;
8
9       // TODO: update the coordinates in the cartesian grid (p.x, p.y)
            for given mpi_rank, see MPI_Cart_coords()
10      int coords[2];
11      MPI_Cart_coords(p.comm, mpi_rank, 2, coords);
12      p.x = coords[0];
13      p.y = coords[1];
14
15      return p;
16  }
```

Listing 10: Updating Cartesian Partition Coordinates

```
1   Domain createDomain(Partition p) {
2       Domain d;
3
4       // TODO: compute size of the local domain
5       d.nx = IMAGE_WIDTH / p.nx;
6       d.ny = IMAGE_HEIGHT / p.ny;
7
```

```
8          if (p.x == p.nx - 1) {
9              d.nx += IMAGE_WIDTH % p.nx;
10         }
11         if (p.y == p.ny - 1) {
12             d.ny += IMAGE_HEIGHT % p.ny;
13         }
14
15         // TODO: compute index of the first pixel in the local domain
16         d.startx = p.x * (IMAGE_WIDTH / p.nx);
17         d.starty = p.y * (IMAGE_HEIGHT / p.ny);
18
19         // TODO: compute index of the last pixel in the local domain
20         d.endx = d.startx + d.nx - 1;
21         d.endy = d.starty + d.ny - 1;
22
23         return d;
24     }
```

Listing 11: Creating a Local Domain

## 3.2. Send the data to the master

```
1      // TODO: send local partition c to the master process
2      MPI_Send(c, d.nx * d.ny, MPI_INT, 0, 0, MPI_COMM_WORLD);
3      // TODO: receive partition of the process proc into array c (overwrite
           its data)
4      MPI_Recv(c, d1.nx * d1.ny, MPI_INT, proc, 0, MPI_COMM_WORLD,
         MPI_STATUS_IGNORE);
```
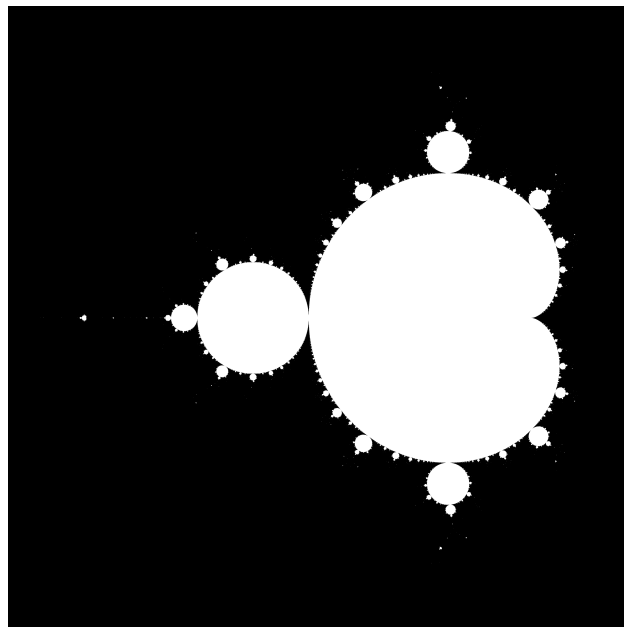
Listing 12: Sending and Receiving Local Partition Data



Figure 1: Mandelbrot set

## 3.3. Performance observe

The highest computation times are visible for the smaller number of processes (MPI size = 1 and 2), which significantly decrease as the number of processes increases. This decrease flattens somewhat
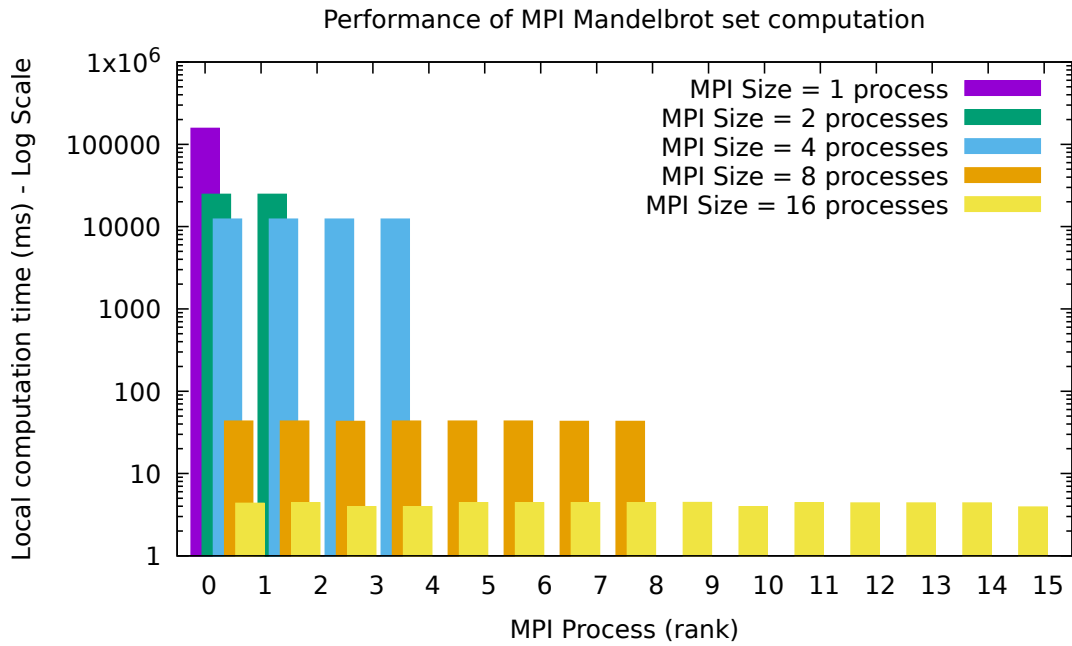
Figure 2: Performance Report

for higher numbers of processes (MPI size = 8 and 16), indicating a diminishing return in performance gains. The graph shows expected scalability up to a point, where adding more processes yields lesser improvements in performance, likely due to overhead associated with managing more processes and increased communication between them.

Blocking `MPI_Send` and `MPI_Recv` calls can lead to waiting if one process is delayed. Replace them with `MPI_Isend` and `MPI_Irecv`. Besides, each process is single-threaded, so using vectorized instructions (SIMD) and multithreading (OpenMP) can also help for enhance performance.

## 4. Task: Parallel matrix-vector multiplication and the power method

```
// To do: Broadcast the random initial guess vector to all MPI
    processes.
// Hint: MPI_Bcast.
MPI_Bcast(y, n, MPI_DOUBLE, 0, MPI_COMM_WORLD);

// Prepare for MPI_Allgatherv
int* recvcounts = (int*)malloc(size * sizeof(int));
int* displs = (int*)malloc(size * sizeof(int));

for (int i = 0; i < size; ++i) {
  if (i < remainder) {
    recvcounts[i] = nrows_base + 1;
    displs[i] = i * recvcounts[i];
  } else {
    recvcounts[i] = nrows_base;
    displs[i] = i * nrows_base + remainder;
  }
}
```

Listing 13: Broadcasting and Preparing for MPI_Allgatherv

```
1    int nrows_base = n / size;       // Base number of rows per process
2    int remainder = n % size;        // Extra rows to distribute
3    int nrows_local;
4    int row_beg_local;
5    int row_end_local;
6
7    if (rank < remainder) {
8      nrows_local = nrows_base + 1;
9      row_beg_local = rank * nrows_local;
10   } else {
11     nrows_local = nrows_base;
12     row_beg_local = rank * nrows_base + remainder;
13   }
14   row_end_local = row_beg_local + nrows_local - 1;
15
16   printf("[Proc␣%3d]␣Doing␣rows␣%d␣to␣%d\n", rank, row_beg_local,
          row_end_local);
```

Listing 14: Determining Local Row Ranges for MPI Processes
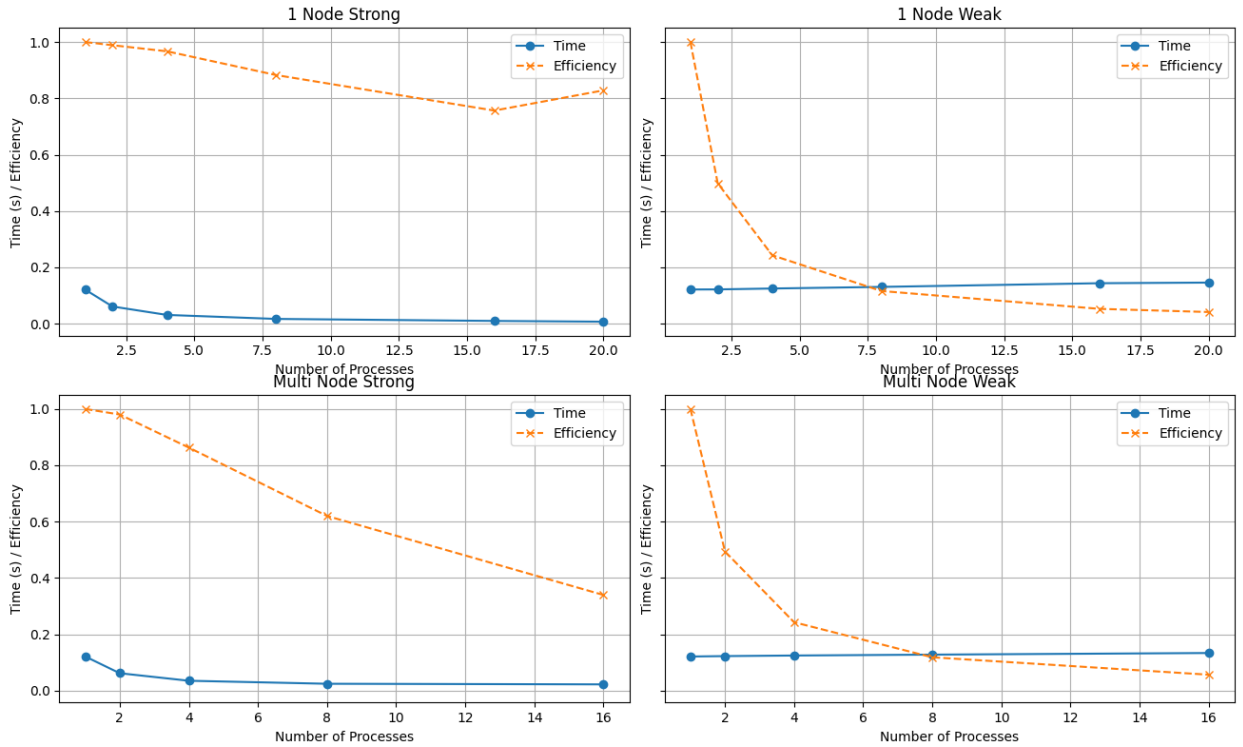


Figure 3: Analysis Report

In our MPI scalability tests, we focused on the time to solution and efficiency across different configurations and scaling strategies. Notably, in strong scaling on a single node, the time to solution improved significantly as we increased the number of processors from 1 to 20. For instance, it decreased from 0.121156 seconds with one processor to just 0.007312 seconds with 20 processors, showcasing a dramatic improvement in performance.

However, this scaling efficiency tends to decrease as we add more processors. While ideal efficiency would be close to 1, the overhead from communication and synchronization among processors dampens this figure. This is especially pronounced when expanding to multiple nodes, where the time to solution slightly worsens due to inter-node communication. For example, with 16 processors spread across different nodes, the efficiency is lower compared to them being on a single node.

For weak scaling, where the problem size grows with the square root of the number of processors, we noticed that the time to solution remains relatively stable, indicating a well-balanced load among

processors. However, as we extend to multiple nodes, such as moving from 1 to 16 processors, slight increases in execution times from 0.121217 to 0.133656 seconds suggest emerging bottlenecks in network communication and data management.

These results underline the importance of considering both algorithmic and system architecture optimizations. Enhancing communication patterns and refining problem partitioning strategies could further improve our application's scalability and efficiency.