

Solution for Project 5

HPC Lab — Submission Instructions

(Please, notice that following instructions are mandatory:
submissions that don't comply with, won't be considered)

- Assignments must be submitted to iCorsi (i.e. in electronic format).
- Provide both executable package and sources (e.g. C/C++ files, Matlab). If you are using libraries, please add them in the file. Sources must be organized in directories called:
Project_number_lastname_firstname
and the file must be called:
project_number_lastname_firstname.zip
project_number_lastname_firstname.pdf
- The TAs will grade your project by reviewing your project write-up, and looking at the implementation you attempted, and benchmarking your code's performance.
- You are allowed to discuss all questions with anyone you like; however: (i) your submission must list anyone you discussed problems with and (ii) you must write up your submission independently.

This project will introduce you a parallel space solution of a nonlinear PDE using MPI.

1. Task 1 - Initialize and finalize MPI

```

1  [wangzi@icsnode33 mini_app]$ mpirun ./main 128 100 0.005
2  =====
3                                     Welcome to mini-stencil!
4  version      :: C++ MPI
5  processes    :: 4
6  mesh         :: 128 * 128 dx = 0.00787402
7  time         :: 100 time steps from 0 .. 0.005
8  iteration    :: CG 300, Newton 50, tolerance 1e-06
9  =====
10 -----
11 1513 conjugate gradient iterations, at rate of 6599.61 iters/
    second
12 300 newton iterations
13 -----
14 Goodbye!
```

Listing 1: Execution of mini-stencil Simulation

Limited condition `rank == 0` to make sure only process 0 report in terminal.

2. Task 2 - Create a Cartesian topology

```
1 // TODO: determine the number of sub-domains in the x and y dimensions
2 // using MPI_Dims_create
3 int dims[2] = {1, 1};
4 MPI_Dims_create(mpi_size, 2, dims);
5 ndomy = dims[0];
6 ndomx = dims[1];
7
8 // TODO: create a 2D non-periodic Cartesian topology using
9 // MPI_Cart_create
10 int periods[2] = {0, 0};
11 MPI_Comm comm_cart;
12 MPI_Cart_create(MPI_COMM_WORLD, 2, dims, periods, 0, &comm_cart);
13
14 // TODO: retrieve coordinates of the rank in the topology using
15 // MPI_Cart_coords
16 int coords[2] = {0, 0};
17 MPI_Cart_coords(comm_cart, mpi_rank, 2, coords);
18 domy = coords[0] + 1;
19 domx = coords[1] + 1;
20
21 // TODO: set neighbours for all directions using MPI_Cart_shift
22 MPI_Cart_shift(comm_cart, 0, 1, &neighbour_south, &neighbour_north);
23 MPI_Cart_shift(comm_cart, 1, 1, &neighbour_west, &neighbour_east);
```

Listing 2: Setting Up a 2D Cartesian Topology with MPI

I used a **2D Cartesian domain decomposition** to divide the grid across MPI processes. `MPI_Dims_create` determined a balanced decomposition, and `MPI_Cart_create` established a non-periodic Cartesian topology for identifying neighboring subdomains.

Each process handles a specific subdomain, ensuring consistent workload distribution and efficient computation.

2.1. Method selected reason

1. **Load Balancing:** Dividing the domain into equal-sized subdomains ensures each process has similar work, preventing load imbalance.
2. **Communication Efficiency:** The 2D Cartesian topology allows efficient communication with only immediate neighbors, minimizing overhead.
3. **Scalability:** This strategy scales well with increasing grid size or number of processes, maintaining balance and efficiency.

2.2. Implications on Performance

1. **Load Balancing:** `MPI_Dims_create` ensures even domain division, minimizing idle time and maximizing resource use.
2. **Communication Overhead:** Communication is limited to neighboring processes, reducing data transfer and overall cost.
3. **Computational Efficiency:** Each process independently handles a subdomain, allowing effective parallelization and overlapping of communication and computation.

3. Task 3 - Change linear algebra functions

```
1  double hpc_dot(Field const& x, Field const& y) {
2      double local_result = 0;
3      int N = y.length();
4
5      // Compute local result
6      for (int i = 0; i < N; i++) {
7          local_result += x[i] * y[i];
8      }
9
10     double global_result = 0;
11     // Perform an MPI_Allreduce to sum up results across all processes
12     MPI_Allreduce(&local_result, &global_result, 1, MPI_DOUBLE,
13                  MPI_SUM, MPI_COMM_WORLD);
14
15     return global_result;
16 }
```

Listing 3: Implementation of hpc_dot for Dot Product with MPI

```
1  double hpc_norm2(Field const& x) {
2      double local_result = 0;
3      int N = x.length();
4
5      // Compute local result
6      for (int i = 0; i < N; i++) {
7          local_result += x[i] * x[i];
8      }
9
10     double global_result = 0;
11     // Perform an MPI_Allreduce to sum up the squared values across
12     // all processes
13     MPI_Allreduce(&local_result, &global_result, 1, MPI_DOUBLE,
14                  MPI_SUM, MPI_COMM_WORLD);
15
16     // Return the square root of the global result
17     return sqrt(global_result);
18 }
```

Listing 4: Implementation of hpc_norm2 for 2-Norm Calculation with MPI

```
1  void hpc_cg(Field& deltay, Field const& y_old, Field const& y_new,
2              Field const& f, const int maxiters, const double tol,
3              bool& success) {
4      // This is the dimension of the linear system to solve
5      int nx = data::domain.nx;
6      int ny = data::domain.ny;
7
8      if (!cg_initialized) cg_init(nx, ny);
9
10     // Epsilon value for matrix-vector approximation
11     double eps = 1.e-4;
12     double eps_inv = 1. / eps;
13
14     // Initialize to zero
15     hpc_fill(fv, 0.0);
16
17     // Jacobian times deltay matrix-vector multiplication
18     // approximation
19 }
```

```

18 // J*deltay = 1/epsilon * ( f(y_new + epsilon*deltay) - f(y_new) )
19
20 // v = y_new + epsilon*deltay
21 hpc_lcomb(v, 1.0, y_new, eps, deltay);
22
23 // fv = f(v)
24 diffusion(y_old, v, fv);
25
26 // r = b - A*x = f(y_new) - J*deltay
27 // where A*x = (f(v) - f) / eps
28 hpc_add_scaled_diff(r, f, -eps_inv, fv, f);
29
30 // p = r
31 hpc_copy(p, r);
32
33 // r_old_inner = <r, r>
34 double r_old_inner = hpc_dot(r, r), r_new_inner = r_old_inner;
35 MPI_Allreduce(MPI_IN_PLACE, &r_old_inner, 1, MPI_DOUBLE, MPI_SUM,
36               MPI_COMM_WORLD);
37
38 // Check for convergence
39 success = false;
40 if (sqrt(r_old_inner) < tol) {
41     success = true;
42     return;
43 }
44
45 int iter;
46 for (iter = 0; iter < maxiters; iter++) {
47     // Ap = A*p
48     hpc_lcomb(v, 1.0, y_new, eps, p);
49     diffusion(y_old, v, fv);
50     hpc_scaled_diff(Ap, eps_inv, fv, f);
51
52     // alpha = r_old_inner / p' * Ap
53     double pAp = hpc_dot(p, Ap);
54     MPI_Allreduce(MPI_IN_PLACE, &pAp, 1, MPI_DOUBLE, MPI_SUM,
55                   MPI_COMM_WORLD);
56     double alpha = r_old_inner / pAp;
57
58     // deltay += alpha * p
59     hpc_axpy(deltay, alpha, p);
60
61     // r -= alpha * Ap
62     hpc_axpy(r, -alpha, Ap);
63
64     // Find new norm
65     r_new_inner = hpc_dot(r, r);
66     MPI_Allreduce(MPI_IN_PLACE, &r_new_inner, 1, MPI_DOUBLE,
67                   MPI_SUM, MPI_COMM_WORLD);
68
69     // Test for convergence
70     if (sqrt(r_new_inner) < tol) {
71         success = true;
72         break;
73     }
74
75     // p = r + (r_new_inner / r_old_inner) * p
76     double beta = r_new_inner / r_old_inner;
77     hpc_lcomb(p, 1.0, r, beta, p);

```

```

75         r_old_inner = r_new_inner;
76     }
77     stats::iters_cg += iter + 1;
78
79     if (!success) std::cerr << "ERROR: CG failed to converge" << std::
80         endl;
81 }

```

Listing 5: Implementation of `hpc_cg` for Conjugate Gradient Method

3.1. Functions Modified and Not Modified

I modified several functions to enable parallel computation. The `hpc_dot` function, responsible for inner product calculation, was updated to compute the local inner product within each process and then use `MPI_Allreduce` to aggregate results globally. Similarly, the `hpc_norm2` function for 2-norm calculation was modified to calculate the local sum of squares, aggregate using `MPI_Allreduce`, and compute the square root of the global result. The `hpc_cg` function is a Conjugate Gradient solver for solving large, sparse, symmetric positive-definite systems of linear equations from PDE discretization. For Jacobian matrix-vector multiplication, I implemented halo cell exchanges using non-blocking communication (`MPI_Isend` and `MPI_Irecv`) to ensure efficient handling of subdomain boundaries. Additionally, I updated the convergence check to use `MPI_Allreduce` for global residual norm computation.

I did not modify the vector operations (`hpc_axpy`, `hpc_add_scaled_diff`, `hpc_scaled_diff`, `hpc_scale`, `hpc_lcomb` and `hpc_copy`). These functions only involve local calculations within each process and do not require cross-process communication, so the original implementation was sufficient.

3.2. Modification Method and Analysis

For the inner product calculation (`hpc_dot`), I ensured that each process computes the inner product of its local vector segment and then aggregates the results using `MPI_Allreduce` to obtain the global value. This approach efficiently parallelizes the computation while ensuring consistency.

For the 2-norm calculation (`hpc_norm2`), I applied a similar strategy. Each process calculates the sum of squares of its local vector elements, and `MPI_Allreduce` aggregates these sums. The square root of the global result is then computed to finalize the norm calculation.

For the `hpc_cg`, replaced standard dot product with an MPI-based version using `MPI_Allreduce` for global aggregation. Added halo cell exchanges using non-blocking communication (`MPI_Isend` and `MPI_Irecv`) to handle subdomain boundaries efficiently. Convergence Check (`hpc_norm2`), Used `MPI_Allreduce` to compute the global residual norm after each iteration.

The modifications were designed to balance load and minimize communication overhead. Each process performs its share of the computation, ensuring that the workload is evenly distributed. Communication overhead is introduced by `MPI_Allreduce` and halo cell exchanges, but these operations are essential for maintaining accuracy and consistency in parallel computations. Using non-blocking communication for boundary exchanges further optimizes performance by overlapping communication with computation.

4. Task 4 - Exchange ghost cells [45 Points]

4.1. Ghost cell exchange function

```

1     MPI_Request requests[8];
2     int req_count = 0;
3
4     // Send to and receive from the northern neighbor

```

```

5 MPI_Isend(&s_new(0, jend-1), nx, MPI_DOUBLE, domain.neighbour_north,
6 0, MPI_COMM_WORLD, &requests[req_count++]);
7 MPI_Irecv(&buffN[0], nx, MPI_DOUBLE, domain.neighbour_north, 1,
8 MPI_COMM_WORLD, &requests[req_count++]);
9
10 // Send to and receive from the southern neighbor
11 MPI_Isend(&s_new(0, 0), nx, MPI_DOUBLE, domain.neighbour_south, 2,
12 MPI_COMM_WORLD, &requests[req_count++]);
13 MPI_Irecv(&buffS[0], nx, MPI_DOUBLE, domain.neighbour_south, 3,
14 MPI_COMM_WORLD, &requests[req_count++]);
15
16 // Send to and receive from the eastern neighbor
17 MPI_Isend(&s_new(iend-1, 0), ny, MPI_DOUBLE, domain.neighbour_east, 4,
18 MPI_COMM_WORLD, &requests[req_count++]);
19 MPI_Irecv(&buffE[0], ny, MPI_DOUBLE, domain.neighbour_east, 5,
20 MPI_COMM_WORLD, &requests[req_count++]);
21
22 // Send to and receive from the western neighbor
23 MPI_Isend(&s_new(0, 0), ny, MPI_DOUBLE, domain.neighbour_west, 6,
24 MPI_COMM_WORLD, &requests[req_count++]);
25 MPI_Irecv(&buffW[0], ny, MPI_DOUBLE, domain.neighbour_west, 7,
26 MPI_COMM_WORLD, &requests[req_count++]);
27
28 // Wait for all communication requests to complete
29 MPI_Waitall(req_count, requests, MPI_STATUSES_IGNORE);

```

Listing 6: Non-Blocking Communication for Neighbor Exchange

To exchange ghost cells between neighboring processes, I used `MPI_Isend` and `MPI_Irecv`, which allow overlapping communication with computation, improving efficiency. Each subdomain sends its boundary values to neighboring subdomains and receives boundary values from them. Communication is required for the north, south, east, and west boundaries.

For each direction:

- Initiate an `MPI_Isend` to send the boundary values.
- Initiate an `MPI_Irecv` to receive the ghost cell values.

After initiating all communications, I used `MPI_Waitall` to ensure all communication had completed before proceeding with further calculations.

4.2. Implement parallel I/O

```

1 [wangzi@icsnode30 mpi_io]$ mpirun -np 4 ./mpi_io 100 100
2 [Process 3 ( 1, 1)] global dims: 100 x 100; local dims: 50 x
3 50 (50:99 50:99);
4 [Process 1 ( 0, 1)] global dims: 100 x 100; local dims: 50 x
5 50 (0:49, 50:99);
6 [Process 0 ( 0, 0)] global dims: 100 x 100; local dims: 50 x
7 50 (0:49, 0:49);
8 [Process 2 ( 1, 0)] global dims: 100 x 100; local dims: 50 x
9 50 (50:99, 0:49);

```

Listing 7: Output of MPI Program for Distributed IO

The demonstration code showcases the use of **MPI-IO** to write a distributed 2D grid into a single binary output file efficiently. The grid is decomposed across MPI processes, with each process handling a subdomain. The grid values are initialized to the MPI rank, and the data is written collectively to a binary file. Below are the main components:

Domain Decomposition: The global grid is divided into subdomains using a 2D Cartesian decomposition. `MPI_Dims_create` determines the decomposition dimensions, while `MPI_Cart_create` creates a Cartesian topology to simplify communication and indexing. The `decompId` function ensures balanced workload distribution across processes.

Memory Allocation and Initialization: Each process allocates memory for its subdomain and initializes values to its MPI rank, enabling easy identification of ownership in the output.

MPI-IO for Data Output:

- **Subarray Types:** `MPI_Type_create_subarray` defines the layout of each process's subdomain within the global grid.
- **File View:** `MPI_File_set_view` sets the global file view for each process, ensuring correct placement of data.
- **Collective Writing:** `MPI_File_write_all` performs a collective I/O operation, ensuring synchronization and optimized performance.

BOV File Creation: After data writing, the root process generates a **BOV (Brick of Values)** header file, enabling visualization tools like `VisIt` to interpret the binary data correctly.

Performance Considerations: The domain decomposition ensures load balancing, while the use of MPI-IO and Cartesian topology supports scalability. Collective I/O minimizes file system contention, and derived datatypes ensure efficient data placement.

4.3. Strong scaling

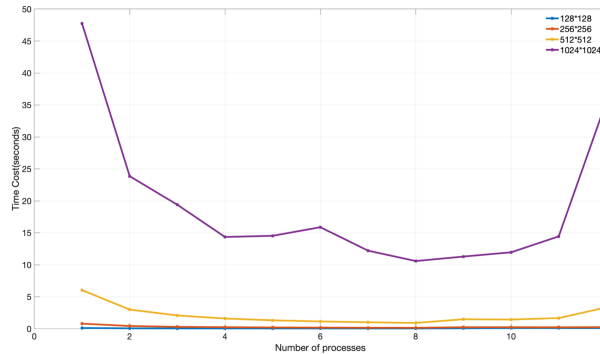


Figure 1: Strong scaling

Smaller resolutions (128×128 , 256×256): When utilizing the parallel version of the PDE solver at these resolutions, performance does not improve and may even degrade. The reason for this condition might be that numerical procedures at this resolutions are computationally inexpensive. The parallel approach may even reduce performance owing to the expense of generating and managing threads.

Higher resolutions (e.g., 512×512 and 1024×1024) improve speed by employing several threads. The improvement is most noticeable at 1024×1024 resolution. Larger grid sizes are computationally more costly than smaller grids. In certain cases, the parallel version can greatly improve performance. The non-linear time cost can be ascribed to the trade-off between the expense of managing processes and the benefits derived from employing numerous processes.

4.4. Weak scaling

The result is that as the number of processes increases, the time cost will rise proportionally. The reason for this is that, while increasing the number of processes, the communication overhead associated with the greater problem size is computationally expensive, resulting in an increase in

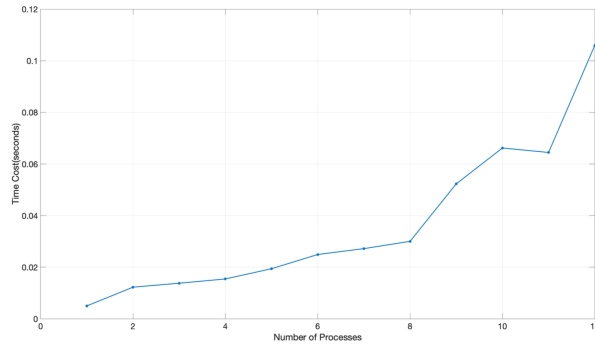


Figure 2: Weak scaling

time consumed. Furthermore, increasing the number of nonlinear iterations can degrade the PDE solver's performance.

5. Task 5 - Testing

5.1. Sum of ranks: MPI collectives

```

1  sum_of_ranks = comm.allreduce(rank, op=MPI.SUM)
2  if rank == 0:
3      print(f"Sum of all ranks (pickle-based): {sum_of_ranks}")
4  comm.Barrier()
5
6  rank_array = np.array(rank, dtype='i')
7  sum_array = np.zeros(1, dtype='i')
8  comm.Allreduce([rank_array, MPI.INT], [sum_array, MPI.INT], op=MPI.SUM
    )

```

Listing 8: MPI-based Sum of Ranks Calculation in Python

5.2. Ghost cell exchange between neighboring processes

```

1  dims = MPI.Compute_dims(size, [0, 0])
2  rows, cols = dims
3  periodic = [True, True]
4  reorder = False
5  cart_comm = comm.Create_cart(dims, periods=periodic, reorder=reorder)
6  coords = cart_comm.Get_coords(rank)
7  north, south = cart_comm.Shift(0, 1)
8  east, west = cart_comm.Shift(1, 1)

```

Listing 9: Creating a 2D Cartesian Topology in MPI

5.3. A self-scheduling example: Parallel Mandelbrot

```

1  def manager(comm, tasks):
2      completed = []
3      num_workers = comm.Get_size() - 1
4      for worker_id in range(1, num_workers + 1):
5          comm.send(tasks[worker_id - 1], dest=worker_id, tag=TAG_TASK)
6      for task_index in range(num_workers, len(tasks)):
7          finished = comm.recv(source=MPI.ANY_SOURCE, tag=MPI.ANY_TAG)
8          completed.append(finished)

```



```

9         comm.send(tasks[task_index], dest=finished.Get_source(), tag=
           TAG_TASK)
10     for _ in range(1, num_workers + 1):
11         final_task = comm.recv(source=MPI.ANY_SOURCE, tag=MPI.ANY_TAG)
12         completed.append(final_task)
13         comm.send(None, dest=final_task.Get_source(), tag=TAG_DONE)
14     return completed
15
16 def worker(comm):
17     task_flag = -1
18     tasks_processed = 0
19     while task_flag != TAG_DONE:
20         task = comm.recv(source=MANAGER, tag=MPI.ANY_TAG)
21         task_flag = task.Get_tag()
22         if task_flag != TAG_DONE:
23             task.process()
24             comm.send(task, dest=MANAGER, tag=TAG_TASK_DONE)
25             tasks_processed += 1
26     return tasks_processed

```

Listing 10: Manager-Worker Model in MPI