

COMP 429/529: Project 3

Muhammed Esad SIMITCIOGLU

January 21, 2024

In Part 1, I parallelized the simulator using MPI by supporting one-dimensional processor geometry in this part under 1D geometry conditions (only -y is set). In Part 2, I introduced OpenMP to enhance the processing, achieving a hybrid programming approach. As a result, each process now utilizes local threads, resulting in improved speed and efficiency. In Part 3, I parallelized the simulator using MPI by supporting two-dimensional processor geometry in this part under 2D geometry conditions (-x and -y is set).

1 Part 1: MPI Parallelization with 1D Geometry

In the initial phase of this assignment, I created a parallelized simulator for 1D geometry. I implemented row-partitioning to ensure contiguous memory for ghost cells. The root process first collects data from the terminal, and then I broadcast all the parameters to the other processors.

```
1  if (rank == 0)
2  {
3      cmdLine(argc, argv, T, n, px, py, plot_freq, no_comm,
4          num_threads);
5  }
6  MPI_Bcast(&T, 1, MPI_DOUBLE, 0, MPLCOMM_WORLD);
7  MPI_Bcast(&n, 1, MPI_INT, 0, MPLCOMM_WORLD);
8  MPI_Bcast(&px, 1, MPI_INT, 0, MPLCOMM_WORLD);
9  MPI_Bcast(&py, 1, MPI_INT, 0, MPLCOMM_WORLD);
10 MPI_Bcast(&plot_freq, 1, MPI_INT, 0, MPLCOMM_WORLD);
11 MPI_Bcast(&no_comm, 1, MPI_INT, 0, MPLCOMM_WORLD);
12 MPI_Bcast(&num_threads, 1, MPI_INT, 0, MPLCOMM_WORLD);
```

I instantiated and initialized the input matrices E, E_prev, and R and then broadcasted them to the other processors.

```
1  if (rank == 0)
2  {
3
4      cout << "RP = " << rp << endl;
5      cout << "dte = " << dte << endl;
6      cout << "dtr = " << dtr << endl;
7      cout << "dt = " << dt << endl;
8      cout << "alpha = " << alpha << endl;
9      cout << "dx = " << dx << endl;
10
11
12     int i, j;
13     for (j = 1; j <= m; j++)
14         for (i = 1; i <= n; i++)
15             E_prev[j][i] = R[j][i] = 0;
16
17     for (j = 1; j <= m; j++)
18         for (i = n / 2 + 1; i <= n; i++)
19             E_prev[j][i] = 1.0;
20
21     for (j = m / 2 + 1; j <= m; j++)
22         for (i = 1; i <= n; i++)
23             R[j][i] = 1.0;
24
25     cout << "Grid Size      : " << n << endl;
26     cout << "Duration of Sim : " << T << endl;
27     cout << "Time step dt    : " << dt << endl;
28     cout << "Process geometry: " << px << " x " << py << endl;
29     if (no_comm)
30         cout << "Communication  : DISABLED" << endl;
31 }
```

```

32     cout << endl;
33 }
34
35 for (int i = 0; i < m + 2; i++)
36 {
37     MPI_Bcast(E[i], n + 2, MPI_DOUBLE, 0, MPLCOMM_WORLD);
38     MPI_Bcast(E_prev[i], n + 2, MPI_DOUBLE, 0, MPLCOMM_WORLD);
39     MPI_Bcast(R[i], n + 2, MPI_DOUBLE, 0, MPLCOMM_WORLD);
40 }

```

To address situations where the process geometry doesn't evenly divide the mesh size, I incorporate extra data into the last process and distribute it using Scatterv instead of Scatter. This ensures that the last process may receive more data compared to others when necessary, unlike Scatter, which evenly distributes data among all processors.

```

1  double **my_E;
2  double **my_E_prev;
3  double **my_R;
4
5  int y_size = 0;
6  if (n % py == 0)
7  {
8      y_size = n / py;
9  }
10 else
11 {
12     if (rank == (P - 1))
13     {
14         y_size = n / py + (n % py);
15     }
16     else
17     {
18         y_size = n / py;
19     }
20 }
21
22 int x_size = n / px;
23 int x_pos = rank % px;
24 int y_pos = rank / px;
25
26 my_E = alloc2D(y_size + 2, x_size + 2);
27 my_E_prev = alloc2D(y_size + 2, x_size + 2);
28 my_R = alloc2D(y_size + 2, x_size + 2);
29
30 int interval = (x_size + 2) * y_size;
31
32 int sendcounts[P];
33 int displs[P];
34 for (int i = 0; i < P; i++)
35 {
36     int y_size_i = n / py + (i == P - 1 ? n % py : 0);
37     int interval_i = (x_size + 2) * y_size_i;
38
39     sendcounts[i] = interval_i;
40     displs[i] = i * interval;
41 }

```

```

42 MPI_Scatterv(&E[1][0], sendcounts, displs, MPLDOUBLE, my_E[1],
43             interval, MPLDOUBLE, 0, MPLCOMM_WORLD);
44 MPI_Scatterv(&E_prev[1][0], sendcounts, displs, MPLDOUBLE,
45             my_E_prev[1], interval, MPLDOUBLE, 0, MPLCOMM_WORLD);
46 MPI_Scatterv(&R[1][0], sendcounts, displs, MPLDOUBLE, my_R[1],
47             interval, MPLDOUBLE, 0, MPLCOMM_WORLD);

```

I utilize immediate operations for sending and receiving from both the bottom and top, followed by a waitAll to ensure that all communication is completed before executing the simulation function.

```

1  if (rank != 0){
2      MPI_Irecv(my_E_prev[0], x_size + 2, MPLDOUBLE, upperRank
3      , 0, MPLCOMM_WORLD, &reqs[requestCount++]);
4      MPI_Isend(my_E_prev[1], x_size + 2, MPLDOUBLE, upperRank
5      , 0, MPLCOMM_WORLD, &reqs[requestCount++]);
6      }
7      // send bottom
8      if (rank != (P - 1)){
9          MPI_Isend(my_E_prev[y_size], x_size + 2, MPLDOUBLE,
10         bottomRank, 0, MPLCOMM_WORLD, &reqs[requestCount++]);
11         MPI_Irecv(my_E_prev[y_size + 1], x_size + 2, MPLDOUBLE,
12         bottomRank, 0, MPLCOMM_WORLD, &reqs[requestCount++]);
13     }
14     MPI_Waitall(requestCount, reqs, status);
15
16     t += dt;
17     niter++;
18
19     simulate(my_E, my_E_prev, my_R, alpha, x_size, y_size, kk,
20             dt, a, epsilon, M1, M2, b, x_pos, y_pos, px, py);

```

I employed identical parameters for both MPI_Gatherv and MPI_Scatterv, with the only variation being in the destination and source parameters.

```

1  MPI_Gatherv(my_E_prev[1], interval, MPLDOUBLE, &E_prev[1][0],
2  sendcounts, displs, MPLDOUBLE, 0, MPLCOMM_WORLD);
3  MPI_Gatherv(my_E[1], interval, MPLDOUBLE, &E[1][0], sendcounts,
4  displs, MPLDOUBLE, 0, MPLCOMM_WORLD);
5  MPI_Gatherv(my_R[1], interval, MPLDOUBLE, &R[1][0], sendcounts,
6  displs, MPLDOUBLE, 0, MPLCOMM_WORLD);

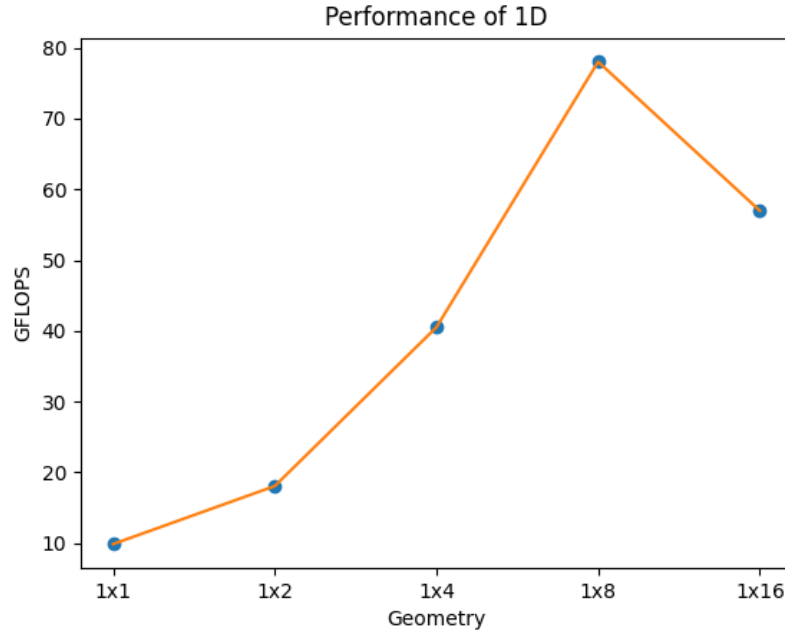
```

1.1 Experiment

Table 1: Study 1: Strong Scaling 1D

Configuration	Execution Time
Serial	8.10818
Geometry 1x1	9.86163
Geometry 1x2	18.0636
Geometry 1x4	40.5181
Geometry 1x8	78.0296
Geometry 1x16	56.9784

Results



(a) Study 1: Strong Scaling 1D Plot

Explanation of Speedup Curve The variance between the serial version and the 1x1 geometry primarily stems from MPI initialization and other overheads introduced by MPI gatherings. Although there is a slight difference between them, it is noticeable. Notably, the most optimal geometry is 1x8, achieving a GFLOPS rate of 78.0296. As evident, increasing the number of processes beyond 8 leads to a reduction in the GFLOP rate. Hence, the optimal choice does not necessarily involve the highest process count.

2 Part II: MPI + OpenMP Parallelism

In the second phase of this assignment, I incorporated hybrid programming by combining MPI and OpenMP. Each process now leverages multiple threads to accelerate the execution of its local tasks. The course materials indicate that the ODE solver exhibits 'No data dependency' and is inherently parallelizable, making it an ideal candidate for introducing parallelism.

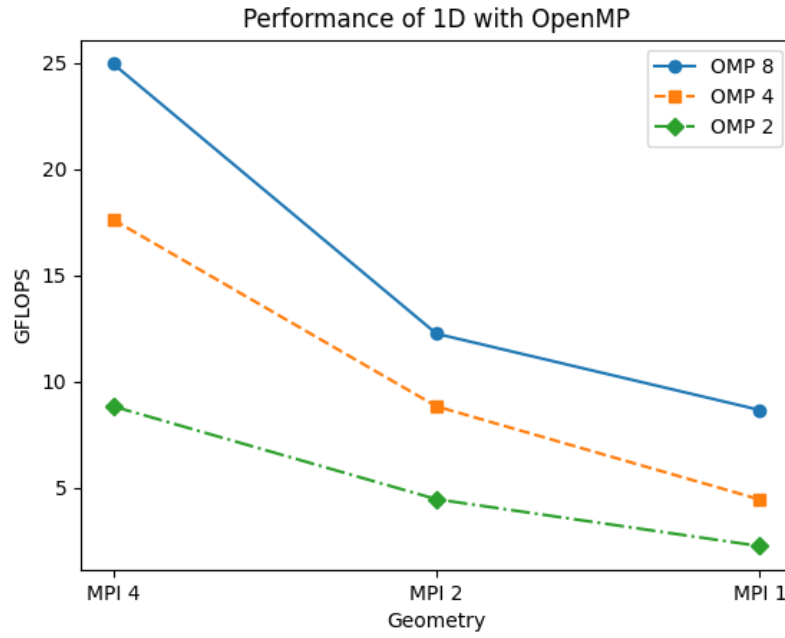
```
1  #pragma omp parallel
2  {
3      #pragma omp for collapse(2)
4      for (j=1; j<=m; j++){
5          for (i=1; i<=n; i++) {
6              E[j][i] = E_prev[j][i]+alpha*(E_prev[j][i+1]+E_prev[j][i-1]-4*
7                  E_prev[j][i]+E_prev[j+1][i]+E_prev[j-1][i]);
8          }
9      }
10
11     /*
12     * Solve the ODE, advancing excitation and recovery to the
13     * next timestep
14     */
15     #pragma omp for collapse(2)
16     for (j=1; j<=m; j++){
17         for (i=1; i<=n; i++)
18             E[j][i] = E[j][i] -dt*(kk* E[j][i]*(E[j][i] - a)*(E[j][i]-1)+ E[j][i] *R[j][i]);
19     }
20
21     #pragma omp for collapse(2)
22     for (j=1; j<=m; j++){
23         for (i=1; i<=n; i++)
24             R[j][i] = R[j][i] + dt*(epsilon+M1* R[j][i]/( E[j][i]+M2))*(-R[j][i]-kk* E[j][i]*(E[j][i]-b-1));
25     }
```

2.1 Experiment

Table 2: Study 4 - Hybrid Programming

Geometry	GFLOPS
MPI4 + OMP8	24.9667
MPI2 + OMP8	12.2496
MPI1 + OMP8	8.65379
MPI4 + OMP4	17.6235
MPI2 + OMP4	8.82627
MPI1 + OMP4	4.43746
MPI4 + OMP2	8.82806
MPI2 + OMP2	4.44943
MPI1 + OMP2	2.24286

Results



(a) Study 4 - Hybrid Programming Plot

Observation

In my analysis, I investigated the performance of various configurations combining MPI and OpenMP (OMP) parallelization. The GFLOPS achieved for different combinations of MPI processes and OpenMP threads are presented in Table 2. Notably, increasing the number of MPI processes while maintaining 8 OpenMP threads resulted in a substantial GFLOPS improvement, signifying enhanced performance. Conversely, reducing the number of MPI processes while keeping 8 OpenMP threads led to diminished GFLOPS, underlining the inefficiency of fewer MPI processes. Furthermore, decreasing the number of OpenMP threads with the same MPI process count exhibited a noticeable GFLOPS decrease, emphasizing the significance of striking the right balance between these parameters for optimal performance in this specific workload. The best GFLOPS result was achieved with the MPI4 + OMP8 Geometry configuration. I couldn't do any more experiment because of the KUACC limitations. Maybe I would find a more optimal geometry configuration with different OMP and MPI counts.

3 Part III: MPI Parallelization with 2D Geometry

In this part, I add 2D geometry instead of 1D geometry part 1. In this part, the implementation must pack and unpack discontinuous memory locations to create messages. There are some different configurations in the code from the part 1. For example, dividing and collecting the data is different. Because we need to divide the data for 2D geometry. Also, for receiving and sending, I send and receive from 4 different sides (N,S,E,W). Let's begin by dividing the data

I've created an MPI Datatype in this code to facilitate data division and transfer. Initially, I defined a 2D matrix and specified a subarray within it. This subarray's dimensions and starting point were crucial for the subsequent Scatterv operation. The MPI Datatype was constructed to represent this subarray, taking into account its memory layout. After committing the datatype, I calculated the necessary parameters, such as the number of subarrays to send in the x and y dimensions (px and py). Lastly, I determined the send counts and displacements for Scatterv, which is essential for distributing data among MPI processes.


```

1  int sizes[2]      = {m, n};
2  int subsizes[2]   = {y, x};
3  int starts[2]     = {0,0};
4  MPI_Datatype type, box;
5  MPI_Type_create_subarray(2, sizes, subsizes, starts, MPI_ORDER_C,
6                           MPL_DOUBLE, &type);
7  MPI_Type_create_resized(type, 0, x*sizeof(double), &box);
8  MPI_Type_commit(&box);
9
10 int px = n / x;
11 int py = n / y;
12 int sendcounts[px*py];
13 int displs[px*py];
14
15 if (rank == 0)
16 {
17     for (int i=0; i<px*py; i++)
18         sendcounts[i] = 1;
19     int disp = 0;
20     for (int i=0; i<py; i++) {
21         for (int j=0; j<px; j++) {
22             disp = j + y * i * px;
23             displs[i*px+j] = disp;
24         }
25     }
26 }

```

After partitioning the data, I apply padding to accommodate ghost cells.

```

1  MPI_Scatterv(s_src[0], sendcounts, displs, box,
2              s_dst[0], x*y, MPL_DOUBLE, 0, MPL_COMM_WORLD);
3
4  addPadding(s_dst, dst, y, x);

```

When transferring data from neighboring processes, the process is similar to Part 1, with the additional consideration of left and right neighbors. Initially, I copy the cells from the left and right neighbors into an local array. After completing each transfer using a WaitAll Call, I update my_E_prev by copying the left and right cells.

```

1  if (y_pos != 0) {
2      MPI_Irecv(my_E_prev[0], x_size+2, MPL_DOUBLE, upperRank,
3               BOTTOMTAG, MPL_COMM_WORLD, &reqs[requestCount++]);
4      MPI_Isend(my_E_prev[1], x_size+2, MPL_DOUBLE, upperRank,
5               UPPERTAG, MPL_COMM_WORLD, &reqs[requestCount++]);
6  }
7
8  if (y_pos != ((P - 1) / px) ) {
9      MPI_Irecv(my_E_prev[y_size+1], x_size+2, MPL_DOUBLE,
10               bottomRank, UPPERTAG, MPL_COMM_WORLD, &reqs[requestCount++]);
11      MPI_Isend(my_E_prev[y_size], x_size+2, MPL_DOUBLE,
12               bottomRank, BOTTOMTAG, MPL_COMM_WORLD, &reqs[requestCount++]);
13  }
14
15  if (x_pos != 0 ) {
16      int i;
17      for (i = 0; i < y_size; i++) {
18          leftSend[i] = my_E_prev[i+1][1];
19      }
20  }

```

```

15     }
16     MPI_Irecv(&left[0], y_size, MPLDOUBLE, rightRank, LEFTTAG,
17     MPLCOMM_WORLD, &reqs[requestCount++]);
18     MPI_Isend(&leftSend[0], y_size, MPLDOUBLE, rightRank,
19     RIGHTTAG, MPLCOMM_WORLD, &reqs[requestCount++]);
20
21     }
22     if(x_pos != ((P-1) % px)) {
23         int i;
24         for(i = 0; i < y_size; i++) {
25             rightSend[i] = my_E_prev[i+1][x_size];
26         }
27         MPI_Irecv(&right[0], y_size, MPLDOUBLE, leftRank, RIGHTTAG,
28         , MPLCOMM_WORLD, &reqs[requestCount++]);
29         MPI_Isend(&rightSend[0], y_size, MPLDOUBLE, leftRank,
30         LEFTTAG, MPLCOMM_WORLD, &reqs[requestCount++]);
31
32     }
33     //
34     MPI_Waitall(requestCount, reqs, status);
35
36     int i;
37     if(x_pos != 0) {
38         for(i = 0; i < y_size; i++) {
39             my_E_prev[i+1][0] = left[i];
40         }
41     }
42     if(x_pos != ((P-1) % px)) {
43         for(i = 0; i < y_size; i++) {
44             my_E_prev[i+1][x_size+1] = right[i];
45         }
46     }

```

Gathering data from other processes to the root process shares a common approach with data division. The primary distinction lies in modifying the destination and receive parameters. We utilize gatherV to assemble the data.

3.1 Experiment

3.1.1 Study 2 - Strong Scaling in 2D

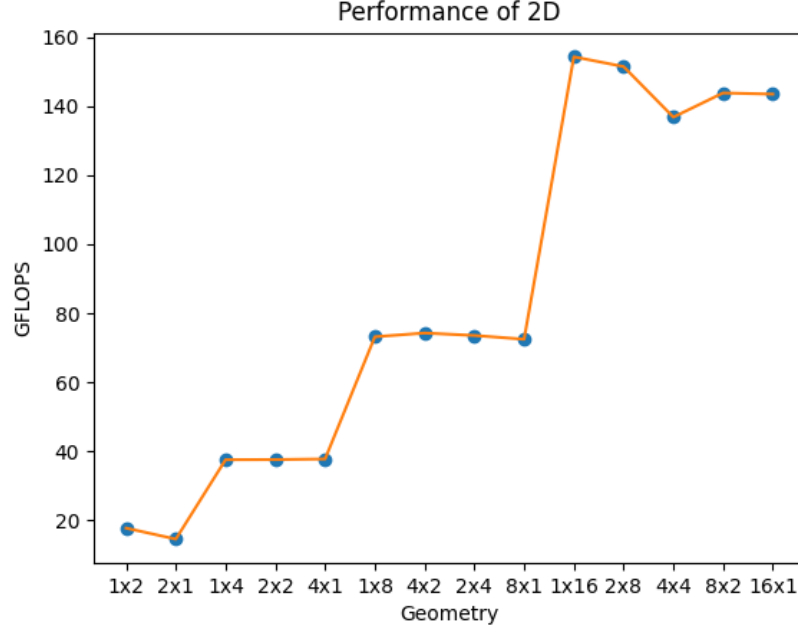
Geometry	GFLOPS
1x2	17.6851
2x1	14.4999
1x4	37.5146
2x2	37.5534
4x1	37.728
1x8	73.1636
4x2	74.2463
2x4	73.5394
8x1	72.4491
1x16	154.342
2x8	151.507
4x4	136.844
8x2	143.86
16x1	143.515

Table 3: Table of Geometry and GFLOPS

Results

Observations

Notably, configurations with more processors, such as 1x16, 8x2, 2x8, and 16x1, exhibit superior computational throughput, indicating efficient parallel processing. Exceptionally high GFLOPS values are observed in the 1x16 and 2x8 geometries, suggesting their potential optimization for parallel workloads. Conversely, the 4x4 geometry yields slightly lower GFLOPS values, indicating potentially less efficient processing for this specific workload. It is crucial to consider that GFLOPS is influenced by various factors, including workload specifics and parallelization efficiency, emphasizing the importance of selecting an appropriate geometry to achieve optimal performance.



(a) Study 2 - Strong Scaling in 2D Plot

3.1.2 Study 3 - Measure Communication Overhead in 2D

In order to compute the communication overhead, I have used the 'no_comm' variable. Whenever I use communication call with MPI, I disabled it with this variable like in the below:

```

1 if (no_comm == 0)
2 {
3     divideData(E, my_E, x_size, y_size, n, m, rank);
4     divideData(E_prev, my_E_prev, x_size, y_size, n, m, rank);
5     divideData(R, my_R, x_size, y_size, n, m, rank);
6 }

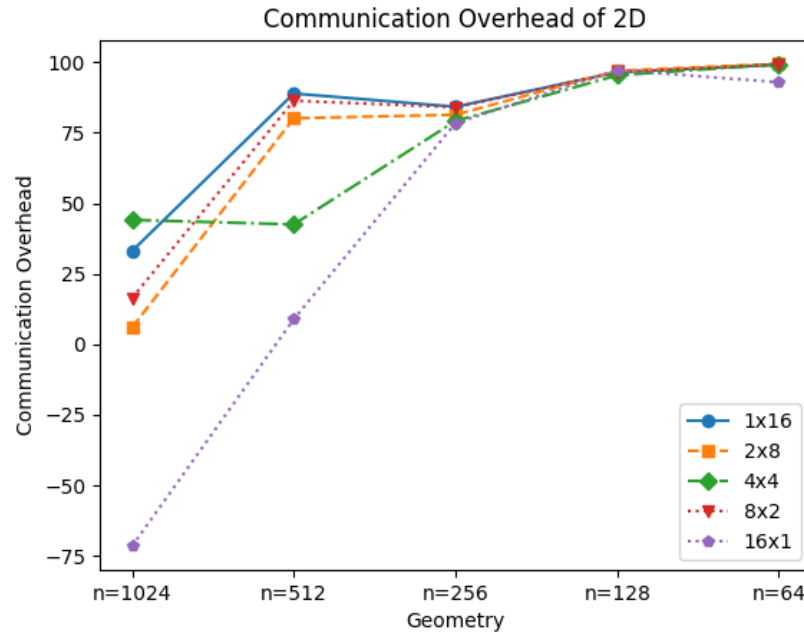
```

Results

Observations

As the problem size decreases, the communication overhead tends to increase. This phenomenon occurs because the calculation-to-communication ratio decreases for each core as the problem size diminishes. With less computational work per unit of communication, the overhead becomes more prominent, leading to increased inefficiencies in data exchange.

It's important to consider that when conducting experiments to calculate the communication within the algorithm, the communication tends to decrease as we reduce the problem count. This reduction is primarily because less data needs to be transferred in this configuration.



(a) Study 2 - Measure Communication Overhead in 2D Plot

An anomaly occurred with the cluster in the 16x1 configuration for $n=1024$. The specific cause of this anomaly remains unclear, and despite attempts to re-experiment with that data point, unfortunately, the issue could not be resolved.

4 Formulas Used

a. *Communication Overhead*

$$CommunicationOverhead = (1 - \frac{ExecutionTimewithNoCommunication}{ExecutionTimewithCommunication} * 100)$$