# CS5300 Project Report
# Parallelizing Gauss-Siedel method

Aaryan, CO21BTECH11001
Aayush, CO21BTECH11002
Darpan, CO21BTECH11004

## 1 Introduction

In the context of a two-dimensional space (2D), the heat transfer problem involves analyzing the temporal evolution of surface temperatures at various points until a stable condition is achieved, where the temperatures become constant (referred to as equilibrium temperature). Initially set at specific temperatures, each point undergoes changes influenced by heat sources positioned along the boundaries of the plane, leading to the dispersion of heat across the surface.

### 1.1 Problem Modelling

The heat transfer problem, as mentioned earlier, involves the diffusion of heat across a surface until it reaches a stable state (equilibrium). In a two-dimensional plane with coordinates (x, y), the problem can be represented by the heat equation:

$$\frac{\partial u}{\partial t} = \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2}$$

This equation describes the temporal variation of temperature $u(x, y)$ at a specific point. To determine the temperature at each point when the steady state is reached (i.e., when heat transfer ceases, $\frac{\partial u}{\partial t} = 0$), the problem can be expressed as:

$$\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = 0$$

It's important to note that the steady-state equation doesn't imply uniform temperature across the entire plane. Temperature gradients can exist in the plate, but these gradients remain constant after reaching steady state as long as the heat sources are constant. To ascertain the temperature at any point in the plate, the presented equation needs to be applied to each point. A common approach is to use a system of linear equations, where each line represents the steady-state temperature at a specific point. Using the previous equation, the temperature at a given point when steady state is reached can be determined with the following expression:

$$\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \Rightarrow \frac{u(x + h/2, y) + u(x - h/2, y) - 2u(x, y)}{h^2} + \frac{u(x, y + k/2) + u(x, y - k/2) - 2u(x, y)}{k^2}$$

To calculate the temperature difference between two points, different steps $h$ and $k$ can be applied for the coordinates $x$ and $y$, introducing errors of magnitude $O(h^2)$ and $O(k^2)$ respectively.

## 1.2 Approach

The adoption of iterative methods for solving computational problems is driven by the necessity to address large-scale problems, particularly those involving systems of linear equations with a substantial number of variables. In such cases, employing a direct method, while yielding precise solutions, demands significant computational resources and time for execution. Consequently, iterative methods become more favorable for handling these scenarios. An illustrative example is the Gauss-Seidel iterative method, which continuously updates a matrix containing the coefficients of each equation.

The update for each element, denoted by $x_i^{k+1}$ in the $k$-th iteration, follows the expression:

$$x_i^{k+1} = \frac{b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{k+1} - \sum_{j=i+1}^{n} a_{ij} x_j^{k}}{a_{ii}}$$

Here, $k$ represents the current iteration. The method terminates when the differences between consecutive iterations fall below a predefined threshold or tolerance.

Given that the coefficients of numerous terms in the resulting linear equations are zero, the task simplifies to computing the mean of four neighboring points. These coefficients of the neighbouring points are denoted as $aE[i], aW[i], aN[i]$, and $aS[i]$.

Hence, our final equation becomes

$$x_{i,j}^{k+1} = \frac{b[i][j] + aE_{ij} x_{i+1,j}^{k} + aW_{ij} x_{i-1,j}^{k+1} + aN_{ij} x_{i,j+1}^{k} + aE_{ij} x_{i,j-1}^{k+1}}{aP_{i,j}}$$

This method enhances the Jacobi method by utilizing previously computed element values to update the current element, resulting in a faster convergence rate and eliminating the need to store two vectors. However, the reliance on the most recent values introduces data dependencies in each iteration, posing challenges for parallel development.

## 1.3 Discretization

Determining the temperature at every point across a continuous 2D plane is impractical due to the infinite number of points involved. Consequently, it becomes essential to transform the domain from continuous to discrete, a procedure known as discretization. This involves applying a mesh to the continuous domain and evaluating the function (under study) at each mesh point. In the context of the heat equation, the mesh comprises four boundaries referred to as heat sources, each with known temperature values, and a set of internal points with unknown temperatures. The discretization process introduces a truncation error because fewer points than those present in the initial problem are used to characterize the 2D plane. The magnitude of the truncation error is inversely proportional to the number of points in the mesh.

## 1.4 Exact Solution

We assume that homogeneous Dirichlet conditions are specified at the left, right and top boundaries, i.e. $T(x = 0) = T(x = Lx) = T(y = Ly) = 0$, and homogeneous Neumann condition is specified at the bottom, i.e. $(\partial T / \partial y)(y = 0) = 0$. The exact solution for the given problem and boundary conditions is expressed as

$$T(x,y) = x(1-x)cos(\pi y)$$

This analytical solution provides a precise description of the temperature distribution across the two-dimensional domain. In this formula, $x$ and $y$ represent the spatial coordinates within the domain.

# 2 Sequential Algorithm

A pseudocode for sequential solver is given in Listing 1

```
       Listing 1: Pseudocode for sequential implementation of Gauss-Siedel method
1   // T -> (nx+2, ny+2) matrix
2   // aW, aS, aP, aN, aE, b -> each of size (nx, ny)
3
4   aW, aS, aP, aN, aE, b = get_coefficients(nx, ny, xst, xen, yst, yen);
5
6   void sequentialGS() {
7       for iter in [1, max_iter] {
8           err = 0.0;
9           for i in [0, nx) {
10              for j in [0, ny) {
11                  ip = i + 1, jp = j+1;
12
13                  T_prev = T[ip][jp];
14
15                  T[ip][jp] = ( b[i][j] + aE[i][j]*Tpnew[ip+1][jp] + aW[i][j]*
                        Tpnew[ip-1][jp] +
16                              aN[i][j]*Tpnew[ip][jp+1] + aS[i][j]*Tpnew[ip][jp-1] )
                                / aP[i][j];
17
18                  err += (T[ip][j] - T_prev)^2;
19              }
20          }
21
22          if (err < threshold) break;
23      }
24  }
```

# 3   Parallel Algorithms using Barrier

The general structure for implementing multithreaded algorithms using barriers is given in Listing 2.

```
       Listing 2: Pseudocode for parallel implementation of Gauss-Siedel method
1   err = MAX_VALUE; // global variable
2   thread_err[num_threads]; // global array
3
4   void updateVals() {
5       temp_err = 0.0; // Thread local variable
6       for iter in [1, max_iter] {
7           // Update T and calculate temp_err
8
9           thread_err[ThreadId.get()] = temp_err;
10
11          Barrier_1.wait();
12          // Now calculate error thread will run
13
14          Barrier_2.wait();
15
16          if (err < threshold) break;
17      }
18  }
19
20  void calcError() {
21      for iter in [1, max_iter] {
22          // Wait for updating threads to complete the iteration
23          Barrier_1.wait();
24
25          // Calculate error and update the variable err
26
```

```
27            Barrier_2.wait();
28
29            if (err < threshold) break;
30        }
31    }
32
33    void parallelGS(n) {
34        // Create n updateVals threads
35
36        // Create 1 calcError thread
37    }
```

## 3.1 Wave Method

A grid of size $(nx, ny)$ has $nx + ny - 1$ number of diagonals. Each diagonal is allocated to one thread. Precisely, if the number of threads are $n$, then a thread with id $i$ is allocated the diagonals $\{i, i+n, i+2n, ...\}$.
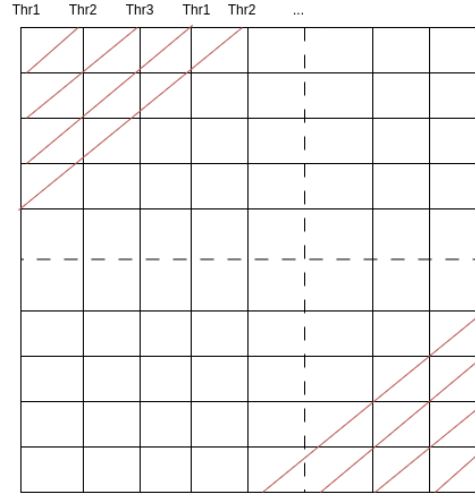


Figure 1: Representation of the distribution of diagonals to 3 threads

A pseudocode for wave method is given in Listing 3.

```
                 Listing 3: Pseudocode for updating values using Wave Method
1   thrId = ThreadId.get();
2
3   for (iter = 0; iter < max_iter; iter++) {
4       diagonal_number = thrId;
5
6       while (diagonal_number <= num_diagonals) {
7           // Update every point in this diagonal
8
9           diagonal_number += num_threads;
10      }
11  }
```

## 3.2 Red Black Method

We split the grid points into red and black color scheme (shown in Figure 2).
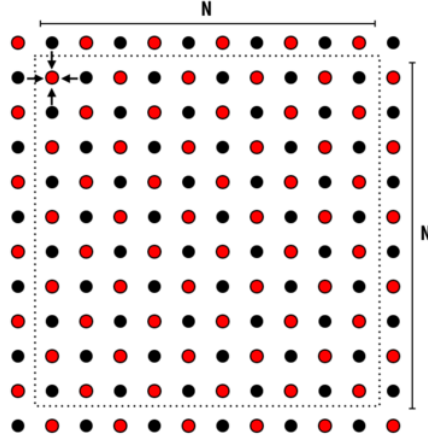
4

Figure 2: Splitting of grid points into red and black points

Notice that the calculation for a red point is dependent only on black points and vice versa. Each iteration of this algorithm can be divided into these steps:

- Distribute all the red points among the threads. Since the calculation of red points are dependent on black points, which are not calculated yet, the values of the previous iteration will be used for them.

- Wait for all the threads to complete the first step. A barrier is used to achieve this.

- Distribute all the black points among the threads. Since the calculation of black points are dependent on red points, which are calculated in the first step, the values of the same iteration will be used for them.

The distribution of points at each step is done by statically allocating rows to the threads. A pseudocode for red-black method is given in Listing 4.

```
Listing 4: Pseudocode for updating values using Red-Black Method
1   thrId = ThreadId.get();
2
3   for (iter = 0; iter < max_iter; iter++) {
4       for (row = thrId; row <= nx; row += num_threads) {
5           for (col = 1; col <= ny; col++) {
6               if (row+col) is even {
7                   // Update T[row][col]
8               }
9           }
10      }
11
12      // Wait for all threads to complete step 1.
13      Barrier.wait();
14
15      for (row = thrId; row <= nx; row += num_threads) {
16          for (col = 1; col <= ny; col++) {
17              if (row+col) is odd {
18                  // Update T[row][col]
19              }
20          }
21      }
22  }
```

### 3.3 Static Row Allocation

A grid of size $(nx, ny)$ has $nx$ rows. Distribution of rows among threads follows a specific pattern based on thread id and number of threads available. Thread $i$ is allocated with rows $\{i, i + k, i + 2k, ...\}$, where k is number of threads.

A pseudocode for wave method is given in Listing 5.

Listing 5: Pseudocode for updating values using Static Row Allocation
```
1   thrId = ThreadId.get();
2
3   for (iter = 0; iter < max_iter; iter++) {
4       for (i = thrId; i < nx ; i += num_threads) {
5           for (j = 0; j < ny ; j++) {
6               // update T[i+1][j+1]
7           }
8       }
9       // wait for all threads to complete
10      barrier.wait();
11  }
```

### 3.4 Static Chunk Allocation

A grid of size $(nx, ny)$ has $nx$ rows. Distribution of rows among threads follows a specific pattern based on thread id, nx and number of threads available. Thread $i$ is allocated with rows $\{i * c, i * c + 1, ..., (i + 1) * c - 1, (i + k) * c, (i + k) * c + 1, ..., (i + k + 1) * c - 1, ...\}$, where k is number of threads and c is chunk size.

A pseudocode for wave method is given in Listing 6.

Listing 6: Pseudocode for updating values using Static Chunk Allocation
```
1   thrId = ThreadId.get();
2
3   for (iter = 0; iter < max_iter; iter++) {
4       id = thrId;
5       while (true) {
6           for (i = id*chunkSize; i < nx && i < (id+1)*chunkSize ; i += num_threads
                ) {
7               for (j = 0; j < ny ; j++) {
8                   // update T[i+1][j+1]
9               }
10          }
11          if (i>=nx) break;
12          id += num_threads;
13      }
14      // wait for all threads to complete
15      barrier.wait();
16  }
```

### 3.5 Dynamic Row Allocation

An atomic counter *row* is initialized as 1. Each thread gets the current row id and increments *row* by 1. If a thread obtains a row id $> nx$ , it breaks out of the loop. After each iteration, row is set to 1 by the error calculation thread.

The pseudocode for a thread which updates the values is given in Listing 7.

```
                        Listing 7: Pseudocode for iteration in batches
1   void dynamicRow ( ) {
2       // Thread local variables
3
4       for ( iter = 0;  iter < max_iter;  iter++) {
5           while ( true ) {
6               // get row no. to operate on
7               currRow = row.fetchAdd(1);
8
9               // iteration has completed
10              if (currRow > nx)  break;
11
12              for ( col = 1;  col <= ny;  col++) {
13                  // update T[currRow][col]
14              }
15          }
16
17          // wait for all threads to complete current iteration
18          barrier.wait();
19      }
20  }
```

## 3.6   Dynamic Chunk Allocation

An atomic counter *row* is initialized as 1. Each thread gets the starting row id and increments *row* by *chunkSize*. The thread then updates the rows $\{startRow, startRow+1, ..., startRow+chunkSize-1\}$. If a thread obtains a row id $> nx$ , it breaks out of the loop. After each iteration, row is set to 1 by the error calculation thread. The optimal size of chunk can depend on cache size of the machine.

The pseudocode for a thread which updates the values is given in Listing 8.

```
                        Listing 8: Pseudocode for iteration in batches
1   void dynamicRow ( ) {
2       // Thread local variables
3
4       for ( iter = 0;  iter < max_iter;  iter++) {
5           while ( true ) {
6               // get row no. to operate on
7               startRow = row.fetchAdd(chunkSize);
8
9               // iteration has completed
10              if (currRow > nx)  break;
11
12              for (currRow = startRow;  currRow < min(startRow+chunkSize, nx+1);
                        currRow++) {
13                  for ( col = 1;  col <= ny;  col++) {
14                      // update T[currRow][col]
15                  }
16              }
17          }
18
19          // wait for all threads to complete current iteration
20          barrier.wait();
21      }
22  }
```

## 3.7 Iteration in Batches

Recognizing the inherent inefficiency introduced by barriers, we devised a modified implementation wherein threads execute iterations in batches rather than synchronizing at every iteration. This method has proven effective in mitigating synchronization overhead and, consequently, reducing the overall computational time.

In particular, the approach involves grouping iterations into batches, allowing threads to operate autonomously within each batch. By adopting this strategy, threads are relieved of the burden of constant synchronization, which becomes especially beneficial when dealing with a large number of iterations. Consequently, the computational workload is distributed more efficiently, resulting in an appreciable improvement in the overall speedup.

This optimization becomes particularly noteworthy when considering the substantial scale of computational iterations—often numbering in the order of tens of thousands. Therefore, it makes perfect sense to perform a few hundred iterations in batches.

The pseudocode for a thread which updates the values is given in Listing 9.

Listing 9: Pseudocode for iteration in batches

```
1   void batchUpdateVals() {
2       // Thread local variables
3       batch_iter_count = 0;
4       temp_err = 0.0;
5
6       for (iter = 0; iter < max_iter; iter++) {
7           // Update T
8           // Calculate error only when (batch_iter_count = iteration_batch_size)
9
10          if (batch_iter_count == iteration_batch_size) {
11              thread_err[ThreadId.get()] = temp_err;
12
13              Barrier_1.wait();
14              // Now calculate error thread will run
15
16              Barrier_2.wait();
17
18              if (err < threshold) break;
19          }
20
21          if (batch_iter_count == iteration_batch_size) batch_iter_count = 0;
22          else batch_iter_count++;
23      }
24  }
```

# 4 Barrier Free Algorithm

We devised an alternative parallel algorithm for the Gauss-Seidel method that eliminates the need for barriers, simplifying thread coordination. In this approach, each thread computes the error during every iteration and records it in a 2D array (size: (`max_iter, num_threads`)). Simultaneously, each thread fetches and increments an atomic integer related to the iteration within a designated array, `count_arr`. When the fetched value equals `num_threads - 1`, the thread calculates the error. If the computed error falls below a specified threshold, the thread signals other threads to terminate through a shared memory boolean variable. This strategy streamlines the parallel computation process without the reliance on barriers, contributing to improved efficiency.

The pseudocode for this approach is given in Listing 10

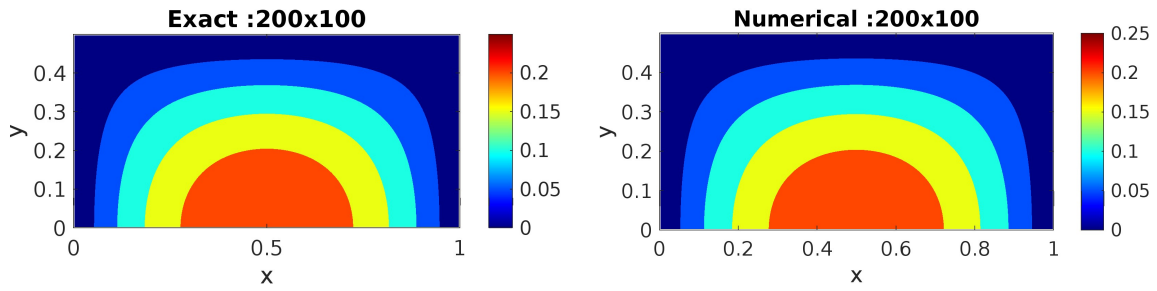Listing 10: Pseudocode for barrier free algorithm

```
1   // global arrays
2   count_arr[max_iter] = {0};
3   thread_err[max_iter][num_threads];
4
5   done_calculation = false; // global variable
6
7   void barrierFreeUpdateVals() {
8       for (iter = 0; iter < max_iter; iter++) {
9           if (done_calculation) break;
10
11          // Update T and calculate temp_err
12          // Keep checking the value of done_calculation
13
14          if (done_calculation) break;
15          thread_err[iter][ThreadId.get()] = temp_err;
16
17          get_count_value = count_arr[iter].fetchAndAdd(1);
18
19          if (get_count_value == num_threads − 1) {
20              // Calculate error
21
22              if (err < threshold) {
23                  done_calculation = true;
24                  break;
25              }
26          }
27      }
28  }
```

**Note:** Link for Github repository.

# 5 Results



(a) Exact Solution                    (b) Solution using BarrierWithBatch method

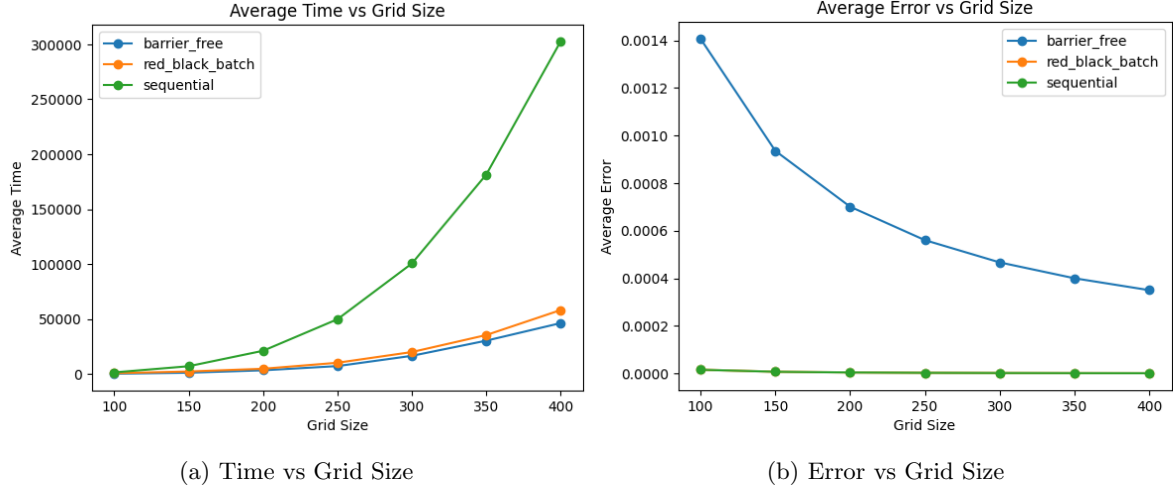Figure 3: Temperature contours for exact and numerical solution for a grid size of $(200, 100)$

(a) Time vs Grid Size

(b) Error vs Grid Size

Figure 4: Sequential vs Red-Black-Batch vs Barrier-Free (# threads = 16)



(a) Time vs Grid Size

(b) Error vs Grid Size

Figure 5: Barrier vs BarrierWithBatch vs BarrierFree using Static Allocation (# threads = 16)

(a) Time vs Grid Size

(b) Error vs Grid Size

Figure 6: Comparison of Barrier and BarrierWithBatch using dynamic allocation (# threads = 16)



(a) Time vs Threads

(b) Error vs Threads

Figure 7: Barrier vs BarrierWithBatch vs BarrierFree using Static Allocation (GridSize: 300 x 150)

(a) Time vs Threads          (b) Error vs Threads

Figure 8: Barrier vs BarrierWithBatch using dynamic allocation (GridSize: 300 x 150)



(a) Time vs Grid Size          (b) Error vs Grid Size

Figure 9: Sequential vs RedBlack vs BarrierWithBatch vs BarrierFree (# threads = 16)
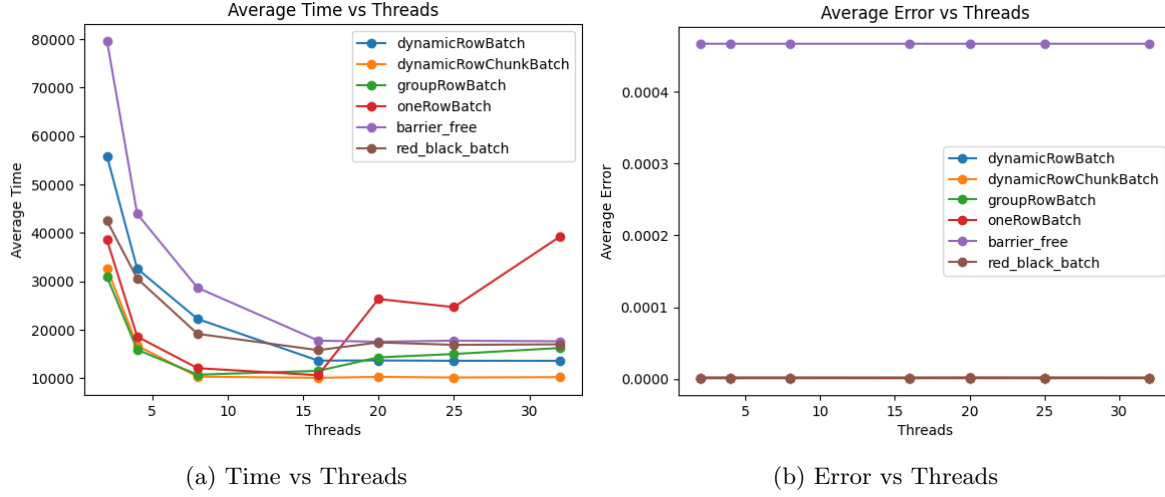
(a) Time vs Threads



(b) Error vs Threads

Figure 10: Sequential vs RedBlack vs BarrierWithBatch vs BarrierFree (GridSize: 300 x 150)

# 6 Observations and Conclusions

Following are the observations obtained from various executions of all the methods:

- When grid size is increased, the error decreases because the values of gradients calculated are more accurate and a finer grid allows for a more detailed representation of the solution that might be missed on a coarser grid. This can be seen in Figure 4, 5 and 6.

- With the increase in threads, the error remains constant as grid size is kept constant. The error is calculated with respect to the exact solution, and the error changes with a change in grid size. This can be seen in Figure 7, 8 and 10.

- While using barriers, allocating chunks of rows to each thread is more beneficial than allocating single rows (for both static as well as dynamic allocation).

- Performing iterations in batches proves to be really useful as can be seen in Figure 7 and 8.

- The speedup obtained by methods using barriers reaches convergence after employing a limited number of threads.

- The barrier free algorithm scales really well in the sense that when the number of threads are by a factor of $f$, then the time taken is decreased by a factor of approx. $f$. However, because the threads are not synchronizing anymore, the average error increases by a significant factor. Although the error obtained using the barrier-free algorithm ($\approx 1e - 4$) is significantly higher than the sequential algorithm ($\approx 1e - 6$), it remains well within an acceptable range. Users, prioritizing speedup and a low overall error may find this solution satisfactory.

# 7 System Specification, Input File and Naming

- All performance tests are conducted on Intel(R) Core(TM) i7-10700 CPU @ 2.90GHz, having 8 cores 16 thread socket.

- Domain size is kept 0.0 to 1.0 in x-direction and 0.0 to 0.5 in y-direction.

- For a grid of dimensions $(nx, ny)$, grid size is changed by changing $nx$ and keeping $ny = nx/2$.

13

- Time is measured in milliseconds in all the experiments and plots.

- oneRow refers to static row allocation, while groupRow refers to static chunk allocation.

- dynamicRow refers to dynamic row allocation, while dynamicRowChunk refers to dynamic chunk allocation.

- If Batch is used, then BarrierWithBatch method is used, else Barrier method unless BarrierFree is not specified.

# 8   References

- `https://www.sciencedirect.com/science/article/pii/S089812210900042X`

- `https://github.com/rpandya1990/Gauss-seidel-Parallel-Implementation`

- `https://gist.github.com/shelterz/4b13459668eec743f15be6c200aa91b2`

**Note:** Link for Github repository.