

Advanced Parallel Programming

Exercise 2

Tim Beringer



TECHNISCHE
UNIVERSITÄT
DARMSTADT

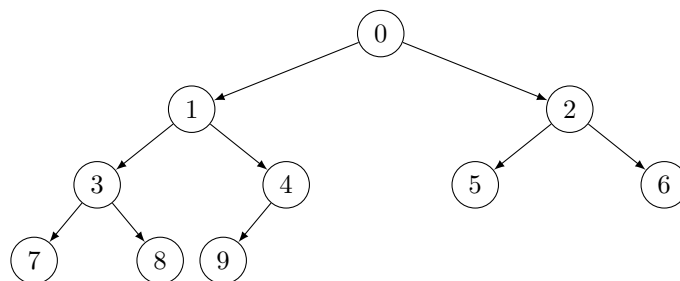
Summer term 2025
22.05.2025

Please solve the following tasks by 22.05.2025. The results are not graded, but a solution is discussed on 22.05.2025.

Task 1: Implement MPI Broadcast

A “broadcast” is a collective communication operation that sends a given message from one process to all other processes. In this task, we will re-implement the same functionality in different ways to get familiar with some important communication patterns.

- Implement a function `Bcast_simple` with the four arguments `buffer`, `count`, `datatype` and `comm` as defined for `MPI_Bcast` (no `root` argument), where rank #0 sends the given message in a simple for-loop to all other ranks.
- A “complete binary tree” is a directed acyclic graph in which each node, except possibly for the last level, which has two child nodes. Set up three formulas that calculate the number for (1) the parent of a child, (2) the left child of a parent, and (3) the right child of a parent, according to the node numbering scheme given in this example (hint: $\lfloor x \rfloor$ is the largest integer not greater than x):



- Implement the above formulas as three separate functions, and add checks for the border cases (i.e., no parent, no child) where `-1` should be returned. Use the function prototypes seen below:

```
1 int get_parent(int child, int num_nodes);  
2 int get_left_child(int parent, int num_nodes);  
3 int get_right_child(int parent, int num_nodes);  
4
```

- The directed edges can be viewed as communication between the MPI processes of the linked node numbers. Use the above functions to implement a more efficient broadcast `Bcast_tree` (same arguments as for `Bcast_simple`) where all ranks except rank #0 first receive a message from their parent using `MPI_Recv`, and then forward this message to their left as well as right children (if available!) using `MPI_Send`.

Solution:

- See Listing 1.

```

1 #include <mpi.h>
2
3 void Bcast_simple(void *buffer, int count, MPI_Datatype datatype, MPI_Comm comm)
4 {
5     int i, my_rank, num_procs, tag=5;
6     MPI_Status status;
7
8     MPI_Comm_size(comm, &num_procs);
9     MPI_Comm_rank(comm, &my_rank);
10
11     if (0 == my_rank) {
12         for (i = 1; i < num_procs; i++) {
13             MPI_Send(buffer, count, datatype, i, tag, comm);
14         }
15     }
16     else {
17         MPI_Recv(buffer, count, datatype, 0, tag, comm, &status);
18     }
19 }

```

Listing 1: Broadcast (simple)

-

$$\text{parent}(x) = \left\lfloor \frac{x-1}{2} \right\rfloor$$

$$\text{leftc}(x) = (2 \cdot x) + 1$$

$$\text{rightc}(x) = (2 \cdot x) + 2$$

- See Listing 2.

```

1 int get_parent(int child, int num_nodes)
2 {
3     int parent;
4
5     /* note: num_nodes is not used int this implementation */
6     if (child > 0) {
7         parent = (child - 1) >> 1;
8     }
9     else parent = -1; /* no parent */
10    return parent;
11 }
12
13 int get_left_child(int parent, int num_nodes)
14 {
15     int child;
16
17     child = (parent << 1) + 1;
18     if (child >= num_nodes) child = -1;
19     return child;
20 }
21
22 int get_right_child(int parent, int num_nodes)
23 {
24     int child;
25
26     child = (parent << 1) + 2;
27     if (child >= num_nodes) child = -1;
28     return child;
29 }

```

Listing 2: Tree formulas

- See Listing 3.

```

1 #include <mpi.h>
2
3 void Bcast_tree(void *buffer, int count, MPI_Datatype datatype, MPI_Comm comm)
4 {
5     int my_rank, num_procs, src, dest, tag=7;
6     MPI_Status status;
7
8     MPI_Comm_size(comm, &num_procs);
9     MPI_Comm_rank(comm, &my_rank);
10
11     src = get_parent(my_rank, num_procs);
12     if (src != -1) MPI_Recv(buffer, count, datatype, src, tag, comm, &status);
13     dest = get_left_child(my_rank, num_procs);
14     if (dest != -1) MPI_Send(buffer, count, datatype, dest, tag, comm);
15     dest = get_right_child(my_rank, num_procs);
16     if (dest != -1) MPI_Send(buffer, count, datatype, dest, tag, comm);
17 }

```

Listing 3: Broadcast (Tree)

Task 2: Benchmark and Analyze Performance of MPI Broadcast

In this task, we will use MPI performance analysis tools to profile and investigate the communication patterns of both broadcast implementations from the previous task. By collecting trace data with Score-P and visualizing it in Vampir, we will examine how each version utilizes the network and identify performance bottlenecks. This will help us better understand the efficiency of different broadcast strategies in terms of communication overhead and process idle times.

- To evaluate the performance of the two broadcast implementations, write a complete MPI program that initializes MPI, allocates a data buffer, and measures the runtime of the broadcast call. Use a buffer of integers with a reasonably large size (e.g., one million elements). On rank #0, initialize the buffer with a fixed value. Repeat the broadcast operation multiple times to reduce measurement noise. Use `MPI_Barrier` before and after the timing region to synchronize processes, and use `MPI_Wtime` to measure the elapsed time. Only rank #0 should print the final result, which should be the average time per broadcast in seconds. Add a command line argument to choose between the two implementations (`Bcast_simple` or `Bcast_tree`). Compile your program with `mpicc++` using optimization flags (e.g., `-O2`).
- To run your code on the Lichtenberg cluster, use a SLURM batch script that requests 32 MPI tasks across two nodes. Run both versions (simple and tree) on 32 processes with the same input size and compare the average broadcast time. Which version is faster for $n = 32$ processes? Why?
Hint: You can inspire from the following bash script:

```

1 #!/bin/bash
2 #SBATCH --job-name=bcast32
3 #SBATCH --nodes=2
4 #SBATCH --ntasks-per-node=16    # total = 32
5 #SBATCH --time=00:05:00
6 #SBATCH --output=bcast32_%j.out
7
8 module load mpi
9
10 echo "Running simple broadcast with 32 procs"
11 srun ./bcast_test 1000000 simple
12
13 echo "Running tree broadcast with 32 procs"
14 srun ./bcast_test 1000000 tree
15

```

More information on cluster usage can be found on the HRZ website.¹

¹https://www.hrz.tu-darmstadt.de/hlr/nutzung_hlr/rechenjobs_und_batchsystem_hlr/index.en.jsp

- Analyze both broadcast implementations using Score-P ² and Vampir ³. On the Lichtenberg cluster, load the required modules for Score-P and Vampir before compiling and running the program. Enable both profiling and tracing in Score-P by setting the appropriate environment variables and compiling the code with Score-P instrumentation. Run both versions of the program with tracing enabled. Then, use the Vampir GUI to analyze the resulting trace files. You can open the GUI remotely using SSH with the -X option for X11 forwarding. In Vampir, navigate to the timeline view to observe the communication behavior.
 - For the simple broadcast version, examine how much time the root rank spends actively sending compared to how long the non-root ranks are idle and waiting to receive. Describe the communication pattern you observe.
 - Compare the simple to the tree-based broadcast version. How does the idle time of non-root ranks differ between the two versions? Based on your analysis, explain which implementation makes better use of parallel communication and why.

Hint: You can enable profiling and tracing with the following variables.

```
1 export SCOREP_ENABLE_PROFILING=true
2 export SCOREP_ENABLE_TRACING=true
3 export SCOREP_EXPERIMENT_DIRECTORY="scorep_simple_32"
4
```

Solution:

- See Listing 4.

```
1 #include <mpi.h>
2 #include <iostream>
3 #include <vector>
4 #include <chrono>
5 #include <cstring>
6
7 int main(int argc, char* argv[]) {
8     int count = 1000000;
9     bool use_tree = false;
10
11     if (argc < 2) {
12         std::cerr << "Usage: " << argv[0] << " <count> [tree]" << std::endl;
13         return EXIT_FAILURE;
14     }
15
16     count = std::stoi(argv[1]);
17     if (argc > 2 && std::string(argv[2]) == "tree") {
18         use_tree = true;
19     }
20
21     MPI_Init(&argc, &argv);
22
23     int rank, size;
24     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
25     MPI_Comm_size(MPI_COMM_WORLD, &size);
26
27     std::vector<int> data(count);
28
29     if (rank == 0) {
30         for (int i = 0; i < count; ++i) {
31             data[i] = i;
32         }
33     }
34
35     MPI_Barrier(MPI_COMM_WORLD);
36     double start = MPI_Wtime();
37
38     for (int i = 0; i < 1000; i++) {
```

²<https://www.vi-hps.org/projects/score-p/>

³<https://vampir.eu/>

```

39     if (use_tree) {
40         Bcast_tree(data.data(), count, MPI_INT, MPI_COMM_WORLD);
41     } else {
42         Bcast_simple(data.data(), count, MPI_INT, MPI_COMM_WORLD);
43     }
44 }
45
46 MPI_Barrier(MPI_COMM_WORLD);
47 double end = MPI_Wtime();
48
49 bool error = false;
50 for (int i = 0; i < count; ++i) {
51     if (data[i] != i) {
52         error = true;
53         break;
54     }
55 }
56
57 if (error) {
58     std::cerr << "Rank " << rank << ": Broadcast error!" << std::endl;
59 }
60
61 if (rank == 0) {
62     std::cout << "Broadcast (" << (use_tree ? "Tree" : "Simple")
63         << "): " << (end - start) << " seconds" << std::endl;
64 }
65
66 MPI_Finalize();
67 return 0;
68 }

```

Listing 4: Broadcast benchmark

- The tree-based broadcast is faster for 32 processes because it reduces the number of communication steps from linear (in the simple version) to logarithmic. This results in less overall waiting time and better parallel utilization, especially as the number of processes increases.
- In the simple broadcast, the root process sends messages sequentially, leading to significant idle time on non-root ranks. In contrast, the tree-based version overlaps communication across processes, reducing idle time and distributing the load more evenly. The tree implementation makes better use of parallel communication due to its logarithmic structure.