

# Introduction to Parallel and Distributed Processing

## One-sided communication

Jaroslav 'Jaric' Zola

<http://www.jzola.org/>

# Problem with Point-to-Point Communication

- Point-to-point communication in MPI (send/recv) requires that both ranks participate
- Hence, the delay on one end slows down the rank on the other end
- This becomes problematic, especially when dealing with irregular communication patterns

# One-sided Communication Idea (RMA)

- Each rank exposes part of its memory to other ranks
- Other ranks can directly read/write to this memory
- No (explicit) synchronization required!

# One-sided Communication Advantage

- Certain communication patterns much easier to express (e.g. recall shifting)
- Modern networks (e.g. IB) support RDMA, and hence provide great performance

# Exposing Memory

- Ranks collectively declare window: memory that will be accessible to all ranks in the communicator
- Memory must be properly aligned (helper API provided)

```
1  #include <mpi.h>
2
3  int main(int argc, char* argv[]) {
4      MPI_Init(&argc, &argv);
5
6      int* buf;
7      MPI_Alloc_mem(1024 * sizeof(int), MPI_INFO_NULL, &buf);
8
9      MPI_Win win;
10     MPI_Win_create(buf, 1024 * sizeof(int), sizeof(int),
11                    MPI_INFO_NULL, MPI_COMM_WORLD, &win);
12
13     // ...
14
15     MPI_Win_free(&win);
16     MPI_Free_mem(buf);
17
18     return MPI_Finalize();
19 } // main
```

# Moving Data

- Primitives for putting, getting and modifying data in a window
- We talk about origin and target for data (instead of sender/receiver)
- Some important caveats
  - No guaranteed ordering of basic operations
  - Basic concurrent writes to the same location undefined

# Getting Data

```
1 | #include <iostream>
2 | #include <mpi.h>
3 |
4 | int main(int argc, char* argv[]) {
5 |     MPI_Init(&argc, &argv);
6 |
7 |     int* buf;
8 |     MPI_Alloc_mem(1024 * sizeof(int), MPI_INFO_NULL, &buf);
9 |
10 |    MPI_Win win;
11 |    MPI_Win_create(buf, 1024 * sizeof(int), sizeof(int),
12 |                  MPI_INFO_NULL, MPI_COMM_WORLD, &win);
13 |
14 |    int x[2];
15 |    MPI_Win_lock(MPI_LOCK_EXCLUSIVE, 0, 0, win);
16 |    // get from buf + 1021 at target 0
17 |    MPI_Get(x, 2, MPI_INT, 0, 1021, 2, MPI_INT, win);
18 |    MPI_Win_unlock(0, win);
19 |
20 |    MPI_Win_free(&win);
21 |    MPI_Free_mem(buf);
22 |
23 |    return MPI_Finalize();
24 | } // main
```

# Putting Data

```
1 | #include <iostream>
2 | #include <mpi.h>
3 |
4 | int main(int argc, char* argv[]) {
5 |     MPI_Init(&argc, &argv);
6 |
7 |     int* buf;
8 |     MPI_Alloc_mem(1024 * sizeof(int), MPI_INFO_NULL, &buf);
9 |
10 |    MPI_Win win;
11 |    MPI_Win_create(buf, 1024 * sizeof(int), sizeof(int),
12 |                  MPI_INFO_NULL, MPI_COMM_WORLD, &win);
13 |
14 |    int x[2];
15 |    MPI_Win_lock(MPI_LOCK_EXCLUSIVE, 0, 0, win);
16 |    // put from x into buf + 1021 at target 0
17 |    MPI_Put(x, 2, MPI_INT, 0, 1021, 2, MPI_INT, win);
18 |    MPI_Win_unlock(0, win);
19 |
20 |    MPI_Win_free(&win);
21 |    MPI_Free_mem(buf);
22 |
23 |    return MPI_Finalize();
24 | } // main
```



# Other Goodies

```
1  #include <iostream>
2  #include <mpi.h>
3
4  int main(int argc, char* argv[]) {
5      MPI_Init(&argc, &argv);
6
7      int* buf;
8      MPI_Alloc_mem(1024 * sizeof(int), MPI_INFO_NULL, &buf);
9
10     MPI_Win win;
11     MPI_Win_create(buf, 1024 * sizeof(int), sizeof(int),
12                    MPI_INFO_NULL, MPI_COMM_WORLD, &win);
13
14     MPI_Win_fence(0, win);
15     MPI_Accumulate(buf, 1024, MPI_INT, 0, 0, 1024, MPI_INT, MPI_SUM, win);
16     MPI_Win_fence(0, win);
17
18     MPI_Win_free(&win);
19     MPI_Free_mem(buf);
20
21     return MPI_Finalize();
22 } // main
```

# For Fun

- Implement a simple distributed hash table using one-sided communication.