

# Introduction to Parallel and Distributed Processing

## Introduction to MPI

Jaroslav 'Jaric' Zola

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

# Suggested Reading

- Message Passing Interface forum  
<http://www.mpi-forum.org/>
- MPI tutorial <https://computing.llnl.gov/tutorials/mpi/>

# Limitations of Shared Memory

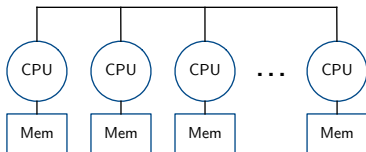
- Shared memory system is not scalable
- As we add more cores/processors, we put more memory pressure
- Memory subsystem is hard to expand (very slow progress)

# Distributed Memory Systems

- Combine many individual “nodes” each with its own memory
- Enable them to communicate via fast interconnect
- When needed add more nodes (hence memory) and improve network

# Distributed Memory Systems

- How to interconnect nodes – network topology?
- How to program such system – coordination, and communication



# Network Topologies

- Overall, you want low-latency, high (bisection) bandwidth
- Static – nodes connect directly point-to-point  
2D/3D/5D mesh or torus, hypercube, etc.
- Dynamic – nodes connected via “switching” element  
often hierarchically, e.g. fat tree

# Programming Distributed Memory

- We are always dealing with MIMD (usually SPMD)
- Overall two general strategies:
  - Explicitly describe communication patterns (e.g. MPI)  
tight control and great expressiveness, but laborious
  - Use higher-level model (e.g. Map/Reduce, Pregel)  
easy to use, lack of flexibility
  - Intermediate solutions (e.g. UPC, Chapel)  
dedicated DSLs for improved productivity

# Message Passing Interface

- Standard for parallel programming with messages, support for C/C++ and Fortran
- Meant to be portable, efficient and practical (but no fault tolerance)
- The workhorse of the HPC industry

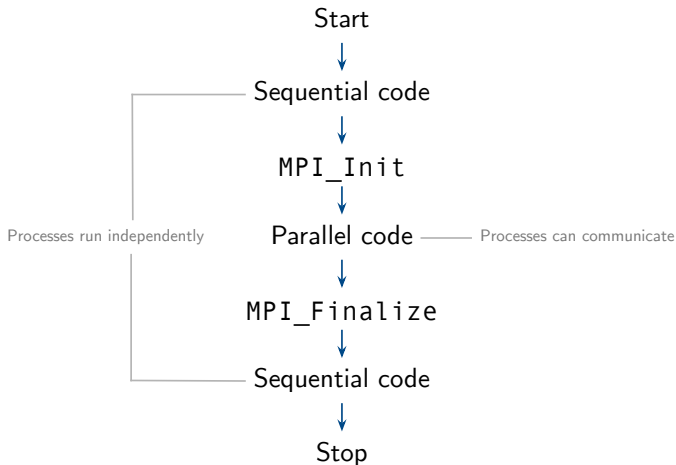


# MPI Support

- Multiple implementations, e.g. MPICH, OpenMPI, each HPC vendor has a version
- MPI standard specifies three components:
  - MPI library with `mpi.h` header
  - Compiler wrapper: `mpicc`, `mpicxx`
  - Runtime system with `mpiexec` command
- Example:

```
1 | mpicxx -std=c++11 -O3 mpiapp.cpp -o mpiapp
2 | mpiexec -machinefile nodeslist.txt -n 65536 mpiapp
```

# MPI Program Flow



# MPI Environment

- Processes are grouped into a communicator
- Each process within a communicator has a unique rank
- Ranks are integers from  $0 \dots p - 1$
- Initially all ranks belong to **MPI\_COMM\_WORLD**

# MPI Environment

```
1 // ex01.cpp
2 #include <iostream>
3 #include <mpi.h>
4
5 int main(int argc, char* argv[]) {
6     MPI_Init(&argc, &argv);
7
8     int rank, size;
9
10    MPI_Comm_size(MPI_COMM_WORLD, &size);
11    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
12
13    std::cout << rank << " " << size << std::endl;
14
15    return MPI_Finalize();
16 } // main
```

# Point to Point Communication

- The most basic form of messaging: one rank sends message to some other
- Different variants possible, e.g. blocking/non-blocking

# Point to Point Communication

```
1 // ex02.cpp
2 #include <mpi.h>
3
4 int main(int argc, char* argv[]) {
5     MPI_Init(&argc, &argv);
6
7     int rank;
8     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
9
10    const int N = 32;
11    int tab[N];
12
13    if (rank == 0) {
14        tab[N - 1] = 13;
15        MPI_Send(tab, N, MPI_INT, 1, 111, MPI_COMM_WORLD);
16    } else if (rank == 1) {
17        MPI_Status stat;
18        MPI_Recv(tab, N, MPI_INT, 0, 111, MPI_COMM_WORLD, &stat);
19    }
20
21    return MPI_Finalize();
22 } // main
```

# Message

- Data part: buffer, count, type
- Message envelope: source, destination, tag, communicator

```
1 | MPI_Send(tab, N, MPI_INT, 1, 111, MPI_COMM_WORLD);  
2 | MPI_Recv(tab, N, MPI_INT, 0, 111, MPI_COMM_WORLD, &stat);
```

- Basic data types `MPI_CHAR`, `MPI_INT`, `MPI_DOUBLE`,  
new types can be created when needed

# Point to Point Mechanics

- `MPI_Send` and `MPI_Recv` are blocking!  
Will not finish until message is sent/received
- Each send operation must be matched by receive operation
- Messages are selected for receiving based on their envelopes
- Messages with the same envelope follow the order in which they were sent



# Point to Point

## Incorrect:

```
1 | // ex04.cpp
2 | #include <mpi.h>
3 |
4 | int main(int argc, char* argv[]) {
5 |     int rank;
6 |
7 |     MPI_Init(&argc, &argv);
8 |     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
9 |
10 |    int x;
11 |
12 |    if (rank == 0) {
13 |        x = 13;
14 |        MPI_Send(&x, 1, MPI_INT, 1, 111, MPI_COMM_WORLD);
15 |    } else if (rank == 1) {
16 |        MPI_Status stat;
17 |        MPI_Recv(&x, 1, MPI_INT, 0, 222, MPI_COMM_WORLD, &stat);
18 |    }
19 |
20 |    return MPI_Finalize();
21 | } // main
```

# Point to Point

## Incorrect:

```
1 // ex05.cpp
2 #include <mpi.h>
3
4 int main(int argc, char* argv[]) {
5     int rank;
6
7     MPI_Init(&argc, &argv);
8     MPI_Comm_rank(MPI_COMM_WORLD, &rank);
9
10    int x[65536];
11    float y[65536];
12
13    if (rank == 0) {
14        x[13] = 13;
15        y[13] = 0.13;
16        MPI_Send(x, 65536, MPI_INT, 1, 111, MPI_COMM_WORLD);
17        MPI_Send(y, 65536, MPI_FLOAT, 1, 222, MPI_COMM_WORLD);
18    } else if (rank == 1) {
19        MPI_Status stat;
20        MPI_Recv(y, 65536, MPI_FLOAT, 0, 222, MPI_COMM_WORLD, &stat);
21        MPI_Recv(x, 65536, MPI_INT, 0, 111, MPI_COMM_WORLD, &stat);
22    }
23
24    return MPI_Finalize();
25 } // main
```

# Useful Trick

- `MPI_ANY_TAG` and `MPI_ANY_SOURCE` in `MPI_Recv` can receive any message from any processor
- `MPI_Status` stores envelope of the received message

```
1 | int x;  
2 | MPI_Status stat;  
3 | MPI_Recv(&x, 1, MPI_INT, MPI_ANY_SOURCE, 111,  
4 |         MPI_COMM_WORLD, &stat);  
5 | int source = stat.MPI_SOURCE;  
6 | int tag = stat.MPI_TAG;
```

# Blocking Send-Recv

- `MPI_Sendrecv` – handy for cyclic communication like shifts

```
1 // ex06.cpp
2 #include <mpi.h>
3
4 int main(int argc, char* argv[]) {
5     int size, rank;
6
7     MPI_Init(&argc, &argv);
8
9     MPI_Comm_size(MPI_COMM_WORLD, &size);
10    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
11
12    int x = rank;
13    int y = 0;
14
15    MPI_Status stat;
16    for (int i = 0; i < size - 1; ++i) {
17        MPI_Sendrecv(&x, 1, MPI_INT, (size + rank - 1) % size, 111,
18                    &y, 1, MPI_INT, (rank + 1) % size, 111,
19                    MPI_COMM_WORLD, &stat);
20        std::swap(x, y);
21    }
22
23    return MPI_Finalize();
24 } // main
```

# For Fun

- Implement passing a message in a ring: processor with rank 0 sends to 1, which forwards to 2, which forwards..., and finally  $p-1$  forwards to 0