

# MPI

---

Jiarui Guo @ ZJUSCT

July 7, 2024

1. Introduction
2. Basic Concepts
3. Point-to-Point Communication
4. Collective Communication
5. Example
6. Miscellaneous

# Introduction

---

- Before 1990's  
Many libraries.  
Writing code was a difficult and tedious task.

---

## Models commonly adopted: Message Passing Model

An application passes messages among processes in order to perform a task.  
e.g. job assignment, results of sub-problems...

- Supercomputing 1992 conference  
A standard interface was defined.
- 1994  
MPI-1
- 2023  
MPI-4.1 by MPI Forum.

# What is MPI

MPI, a Message Passing Interface.

There exists many implementations:

- OpenMPI
- Intel-MPI
- HMPI (Hyper-MPI)
- MPICH
- .....

- OpenMPI  
Lab0
- Intel-MPI  
Intel-oneAPI(Click Me)
- HMPI  
Huawei (Click Me)

# Hello World!

```
1 #include <mpi.h>
2 #include <stdio.h>
3 int main(int argc, char** argv) {
4     MPI_Init(&argc, &argv);
5     int world_size;
6     MPI_Comm_size(MPI_COMM_WORLD, &world_size);
7     int world_rank;
8     MPI_Comm_rank(MPI_COMM_WORLD, &world_rank);
9     char processor_name[MPI_MAX_PROCESSOR_NAME];
10    int name_len;
11    MPI_Get_processor_name(processor_name, &name_len);
12    printf("Hello, world, from processor %s, rank %d out of %d processors\n",
13           processor_name, world_rank, world_size);
14    MPI_Finalize();
15    return 0;
16 }
```

# Basic Concepts

---



A communicator defines a group of processes that have the ability to communicate with one another.

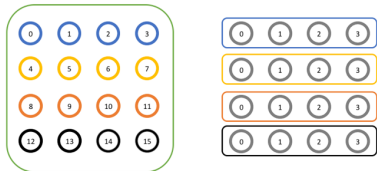
Each process has a **unique** rank.

- MPI\_COMM\_WORLD



- `MPI_COMM_SPLIT(comm, color, key, new_comm)`
  - `comm`: The communicator that will be used as the basis for the new communicators.
  - `color`: Which new communicator each processes will belong.
  - `key`: The ordering (rank) within each new communicator.
  - `new_comm`: [OUT]

Split a Large Communicator Into Smaller Communicators



## Blocking

It does not return until the message data and envelope have been safely stored away so that the sender is free to modify the send buffer.

The message might be copied directly into the matching receive buffer, or it might be copied into a temporary system buffer.

## Non-blocking

A nonblocking call initiates the operation, but does not complete it.

They will return almost immediately.

## Messages are non-overtaking

Order is preserved.(Only under single thread)

If a sender sends two messages in succession to the same destination, and both match the same receive, then this operation cannot receive the second message if the first one is still pending.

If a receiver posts two receives in succession, and both match the same message, then the second receive operation cannot be satisfied by this message, if the first one is still pending.

MPI makes no guarantee of fairness in the handling of communication.

There may be starvation.

## Example

Rank1  $\rightarrow^{send}$  Rank0

Rank2  $\rightarrow^{send}$  Rank0

Rank0  $\leftarrow^{receive}$  from any source.

# Point-to-Point Communication

---

```
1 int MPI_Send(const void* buffer,  
2           int count,  
3           MPI_Datatype datatype,  
4           int recipient ,  
5           int tag,  
6           MPI_Comm communicator);
```

## Parameters:

- **buffer** The buffer to send.
- **count** The number of elements to send.
- **datatype** The type of one buffer element.
- **recipient** The rank of the recipient MPI process.
- **tag** The tag to assign to the message.
- **communicator** The communicator in which the standard send takes place.



```
1 int MPI_Recv(void* buffer,  
2     int count,  
3     MPI_Datatype datatype,  
4     int sender,  
5     int tag,  
6     MPI_Comm communicator,  
7     MPI_Status* status);
```

## Parameters:

- **buffer** The buffer to receive.
- **count** The number of elements to receive.
- **datatype** The type of one buffer element.
- **sender** The rank of the sender MPI process.
- **tag** The tag to assign to the message.
- **communicator** The communicator in which the standard receive takes place.
- **status** The variable in which store the status of the receive operation. Pass MPI\_STATUS\_IGNORE if unused.

MPI\_Status represents the status of a reception operation.

At least 3 attributes:

- MPI\_SOURCE
- MPI\_TAG
- MPI\_ERROR

There may be additional attributes that are implementation-specific.

In addition to the data part, messages carry information that can be used to distinguish messages and selectively receive them.

- source
- destination
- tag
- communicator

- **Buffer Mode**

Can be started whether or not a matching receive was posted. Completion does not depend on the occurrence of a matching receive.

- **Synchronous Mode**

Can be started whether or not a matching receive was posted.

The send will be completed successfully only if a matching receive is posted.

- **Ready Mode**

May be started only if the matching receive is already posted.

- **Standard Mode**

Depends.

Communication mode	Start time	Completion time
Buffer mode	Immediately	Message has gone to buffer
Synchronous mode	Immediately	Matching receive has posted
Ready mode	Matching receive has posted	When the send buffer can be reused
Standard mode	Depends	Depends

## Example

```
1 // n = 2
2 MPI_Comm_rank(comm, &my_rank);
3 MPI_Ssend(sendbuf, count, MPI_INT, my_rank ^ 1, tag, comm);
4 MPI_Recv(recvbuf, count, MPI_INT, my_rank ^ 1, tag, comm, &status);
```

## Example

```
1 // n = 2
2 MPI_Comm_rank(comm, &my_rank);
3 MPI_Ssend(sendbuf, count, MPI_INT, my_rank ^ 1, tag, comm);
4 MPI_Recv(recvbuf, count, MPI_INT, my_rank ^ 1, tag, comm, &status);
```

Deadlock! Any solutions?

## Example

```
1 // n = 2
2 MPI_Comm_rank(comm, &my_rank);
3 if (my_rank == 0) {
4     MPI_Ssend(sendbuf, count, MPI_INT, 1, tag, comm);
5     MPI_Recv(recvbuf, count, MPI_INT, 1, tag, comm, &status);
6 } else if (my_rank == 1) {
7     MPI_Recv(recvbuf, count, MPI_INT, 0, tag, comm, &status);
8     MPI_Ssend(sendbuf, count, MPI_INT, 0, tag, comm);
9 }
```

Any other solutions?



# Blocking Send and Receive

```
1 int MPI_Sendrecv(const void* buffer_send,
2     int count_send,
3     MPI_Datatype datatype_send,
4     int recipient ,
5     int tag_send,
6     void* buffer_recv ,
7     int count_recv ,
8     MPI_Datatype datatype_recv,
9     int sender,
10    int tag_recv ,
11    MPI_Comm communicator,
12    MPI_Status* status);
```

## Notice

The buffers used for send and receive must be different.

Any other solutions?

# Non-Blocking Send and Receive

## Recall:

A nonblocking call initiates the operation, but does not complete it.

They will return almost immediately.

```
1 int MPI_Isend(const void* buffer,  
2      int count,  
3      MPI_Datatype datatype,  
4      int recipient ,  
5      int tag,  
6      MPI_Comm communicator,  
7      MPI_Request* request);
```

- **MPI\_Test**

MPI\_TEST(request, flag, status)

Checks if a non-blocking operation is complete at a given time.

flag=true if completes.

- **MPI\_Wait**

MPI\_WAIT(request, status)

Waits for a non-blocking operation to complete. That is, unlike MPI\_Test, MPI\_Wait will block until the underlying non-blocking operation completes.

# Non-Blocking Send and Receive(Deadlock revisit)

```
1 MPI_Request req;  
2 MPI_Isend(sendbuf, 0x100, MPI_INT, my_rank^1, 0, MPI_COMM_WORLD, &req);  
3 MPI_Recv(recvbuf, 0x100, MPI_INT, my_rank^1, 0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
```

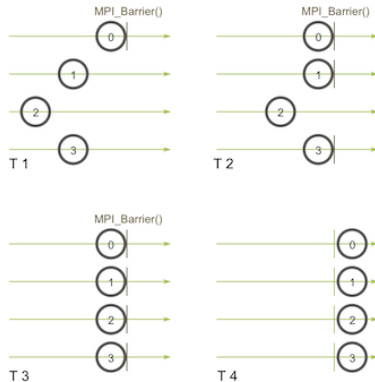
# Collective Communication

---

- **MPI\_Barrier**

MPI\_Barrier(COMM)

Blocks all MPI processes in the given communicator until they all call this routine.

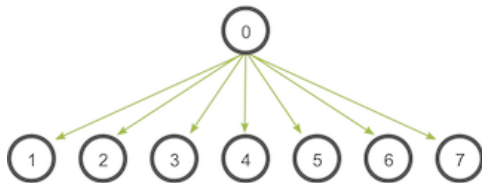


MPI\_Barrier

# BroadCast(One to All)

```
1 int MPI_Bcast(void* buffer,  
2   int count,  
3   MPI_Datatype datatype,  
4   int emitter_rank ,  
5   MPI_Comm communicator);
```

- **emitter\_rank** The rank of the MPI process that broadcasts the data, all other processes receive the data broadcasted.



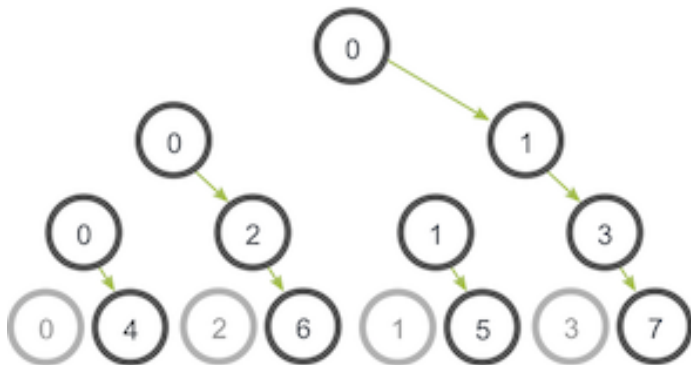
Bcast

## Why not Send and Receive?

```
1 double start = MPI_Wtime();
2
3 if (my_rank == 0){
4     for (int i=1; i<=31; i++)
5         MPI_Send(sendbuf, 0x10000, MPI_INT, i, 0, MPI_COMM_WORLD);
6 } else{
7     MPI_Recv(recvbuf, 0x10000, MPI_INT, 0, 0, MPI_COMM_WORLD, MPI_STATUS_IGNORE);
8 }
9
10 double end = MPI_Wtime();
11
12 if (my_rank == 0) printf(" [Send-Recv]-Finished -in-%f-seconds\n", my_rank, end-start);
13
14 start = MPI_Wtime();
15 MPI_Bcast(&sendbuf, 0x10000, MPI_INT, 0, MPI_COMM_WORLD);
16 end = MPI_Wtime();
17
18 if (my_rank == 0) printf(" [Bcast]-Finished -in-%f-seconds\n", my_rank, end-start);
```



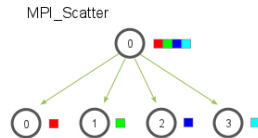
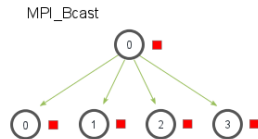
## BroadCast(Tree based algorithm)



# Scatter(One to All)

```
1 int MPI_Scatter(const void* buffer_send ,  
2   int count_send,  
3   MPI_Datatype datatype_send,  
4   void* buffer_recv ,  
5   int count_recv,  
6   MPI_Datatype datatype_recv,  
7   int root,  
8   MPI_Comm communicator);
```

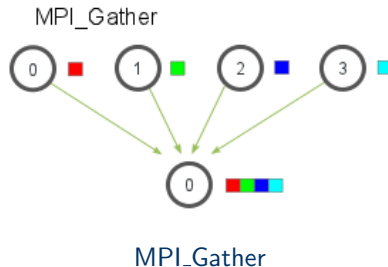
- **count\_send** The number of elements to send to each process, not the total number of elements in the send buffer. For non-root processes, the send parameters like this one are ignored.
- **count\_receive** The number of elements in the receive buffer.



Scatter

# Gather(All to One)

```
1 int MPI_Gather(const void* buffer_send,  
2 int count_send,  
3 MPI_Datatype datatype_send,  
4 void* buffer_recv ,  
5 int count_recv ,  
6 MPI_Datatype datatype_recv,  
7 int root ,  
8 MPI_Comm communicator);
```



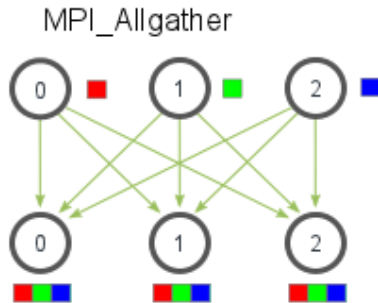
## Example

### Compute average

```
1 MPI_Scatter(buffer, 0x1000000/4, MPI_DOUBLE, local_buffer, 0x1000000/4, MPI_DOUBLE, 0,  
    MPI_COMM_WORLD);  
2 double local_avg = 0;  
3 for(int i=0; i<0x1000000/4; i++){  
4     local_avg += local_buffer[i];  
5 }  
6 local_avg /= 0x1000000/4;  
7 double avgs [4];  
8 MPI_Gather(&local_avg, 1, MPI_DOUBLE, avgs, 1, MPI_DOUBLE, 0, MPI_COMM_WORLD);
```

# Allgather(All to All)

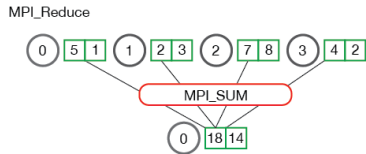
```
1 int MPI_Allgather(const void* buffer_send ,  
2   int count_send,  
3   MPI_Datatype datatype_send,  
4   void* buffer_recv ,  
5   int count_recv ,  
6   MPI_Datatype datatype_recv,  
7   MPI_Comm communicator);  
Actually MPI_Gather + MPI_Bcast.
```



MPI\_Allgather

# Reduce

```
1 int MPI_Reduce(const void* send_buffer,  
2 void* receive_buffer ,  
3 int count,  
4 MPI_Datatype datatype,  
5 MPI_Op operation,  
6 int root,  
7 MPI_Comm communicator);
```



Reduce

## Example

### Compute average revisit

```
1 MPI_Scatter(buffer, 0x1000000/4, MPI_DOUBLE, local_buffer, 0x1000000/4, MPI_DOUBLE, 0,  
    MPI_COMM_WORLD);  
2 double local_avg = 0;  
3 for(int i=0; i<0x1000000/4; i++){  
4     local_avg += local_buffer[i];  
5 }  
6 local_avg /= 0x1000000/4;  
7 double global_avg;  
8 MPI_Reduce(&local_avg, &global_avg, 1, MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD);
```

## Example

---

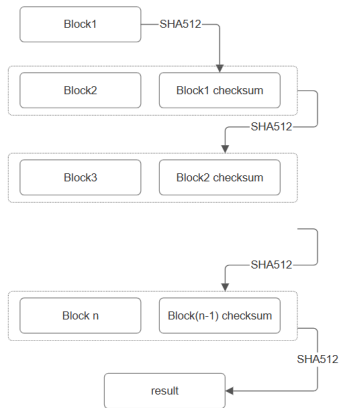


# Task

Implement a data validation algorithm using SHA512.

Algorithm procedure:

1. Tile the input file into blocks of 1MB. (If the last block is smaller than 1MB, pad it with zeros.)
2. For the  $i^{th}$  block, concatenate it with the validation sum SHA512 of  $(i - 1)^{th}$  block and calculate validation sum of SHA512.
3. The validation sum of the last block is considered as the validation sum of the entire file.



```
1 int num_block = (len + BLOCK_SIZE - 1) / BLOCK_SIZE;
2 uint8_t prev_md[SHA512_DIGEST_LENGTH];
3
4 EVP_MD_CTX *ctx = EVP_MD_CTX_new();
5 EVP_MD *sha512 = EVP_MD_fetch(nullptr, "SHA512", nullptr);
6
7 SHA512(nullptr, 0, prev_md);
```

```
1 for (int i = 0; i < num_block; i++) {
2     uint8_t buffer[BLOCK_SIZE]{};
3     EVP_DigestInit_ex(ctx, sha512, nullptr);
4     std::memcpy(buffer, data + i * BLOCK_SIZE,
5                 std::min(BLOCK_SIZE, len - i * BLOCK_SIZE));
6     EVP_DigestUpdate(ctx, buffer, BLOCK_SIZE);
7     EVP_DigestUpdate(ctx, prev_md, SHA512_DIGEST_LENGTH);
8
9     unsigned int len = 0;
10    EVP_DigestFinal_ex(ctx, prev_md, &len);
```

## Notice

EVP\_DigestUpdate(a); EVP\_DigestUpdate(b);

Equivalent to

EVP\_DigestUpdate(concat(a,b)) !

Computation is dependent on the result of the previous one.

How to exploit MPI?

Computation is dependent on the result of the previous one.

How to exploit MPI?

**Answer:**

File **I/O** accounts! We can **overlap** I/O operations with computation.

**Non-Blocking receives the previous block's checksum.**

```
1 if(i != 0){
2     MPI_Irecv((void *)prev_md,
3             SHA512_DIGEST_LENGTH,
4             MPI_UINT8_T,
5             sender,
6             0,
7             MPI_COMM_WORLD,
8             &request);
9 }
```

**Meanwhile... File I/O and Digest**

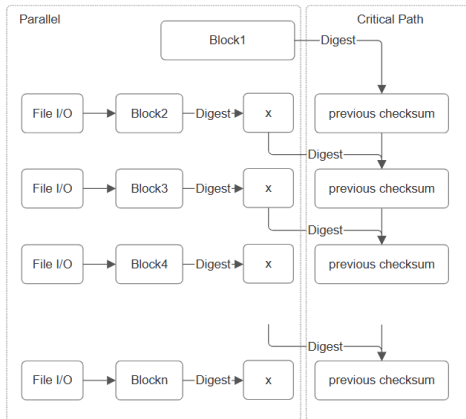
```
1 istrm.seekg(i * BLOCK_SIZE);
2 istrm.read( reinterpret_cast <char *>(data + i * BLOCK_SIZE), std::
               min(BLOCK_SIZE*local_size, file_size - i * BLOCK_SIZE));
3
4 for(int j=i; j<upper_bound; j++){
5     uint8_t buffer2[BLOCK_SIZE]{};
6     EVP_DigestInit_ex(ctx[j-i], sha512, nullptr);
7     std::memcpy(buffer2, data + j * BLOCK_SIZE,
8                 std::min(BLOCK_SIZE, len - j * BLOCK_SIZE));
9     EVP_DigestUpdate(ctx[j-i], buffer2, BLOCK_SIZE);
10 }
11
12 if(i != 0){
13     MPI_Wait(&request,
14             MPI_STATUS_IGNORE);
15 }
```

## Non-blocking send my checksum

```
1 unsigned int len = 0;
2 for (int j=i; j<upper_bound; j++){
3     EVP_DigestUpdate(ctx[j-i], prev_md,
4         SHA512_DIGEST_LENGTH);
5     EVP_DigestFinal_ex(ctx[j-i], prev_md, &len);
6 }
7 if (upper_bound != num_block) {
8     MPI_Isend(prev_md,
9         SHA512_DIGEST_LENGTH,
10        MPI_UINT8_T,
11        0,
12        MPI_COMM_WORLD,
13        &request);
14 }
```

```
(hpc101) jrguo@669:~/hpc101/sha512$ mpirun -n 1 baseline 2G1.bin baseline.out
2G1.bin size: 2147483648
checksum: 7224ec372ae2480f6609c35fe3ed6e1c7ea37f87a9d5e3e79b90838647eab9299db7d818a21lacb7fb53884993869d5e16604e69ae5e459e5fbb3bf12c8
0b270
checksum time cost: 4238ms
total time cost: 22635ms
(hpc101) jrguo@669:~/hpc101/sha512$ cp 2G1.bin 2G2.bin
(hpc101) jrguo@669:~/hpc101/sha512$ mpirun -n 8 mycode 2G2.bin mycode.out
2G2.bin size: 2147483648
checksum: 7224ec372ae2480f6609c35fe3ed6e1c7ea37f87a9d5e3e79b90838647eab9299db7d818a21lacb7fb53884993869d5e16604e69ae5e459e5fbb3bf12c8
0b270
checksum time cost: 18401ms
total time cost: 18403ms
```

# Wrap up



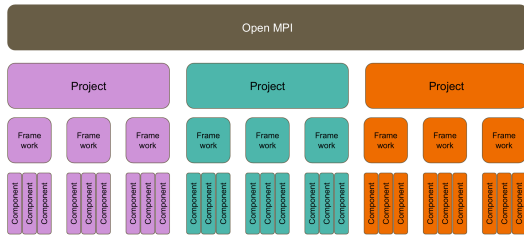
## Miscellaneous

---



## Modular Component Architecture(MCA)

- MCA framework
- MCA component
- MCA module



Open MPI Overall Architecture Terminology

## 3 Types of Open MPI Framework

- In the MPI layer (OMPI)
- In the run-time layer (ORTE)
- In the operating system/platform layer (OPAL)

You might think of these frameworks as ways to group MCA parameters by function. (e.g. btl in OMPI)

```
~$ ompi_info --param btl all
MCA btl: vader (MCA v2.1.0, API v3.1.0, Component v4.1.6)
MCA btl: self (MCA v2.1.0, API v3.1.0, Component v4.1.6)
MCA btl: tcp (MCA v2.1.0, API v3.1.0, Component v4.1.6)
MCA btl tcp: -----
MCA btl tcp: parameter "btl_tcp_if_include" (current value: "",
data source: default, level: 1 user/basic, type:
string)
Comma-delimited list of devices and/or CIDR
notation of networks to use for MPI communication
(e.g., "eth0,192.168.0.0/16"). Mutually exclusive
with btl_tcp_if_exclude.
MCA btl tcp: parameter "btl_tcp_if_exclude" (current value:
"127.0.0.1/8,sppp", data source: default, level: 1
user/basic, type: string)
Comma-delimited list of devices and/or CIDR
notation of networks to NOT use for MPI
communication -- all devices not matching these
specifications will be used (e.g.,
"eth0,192.168.0.0/16"). If set to a non-default
value, it is mutually exclusive with
btl_tcp_if_include.
MCA btl tcp: parameter "btl_tcp_progress_thread" (current value:
"0", data source: default, level: 1 user/basic,
type: int)
```

ompi\_info

## Specify Compilers

`./configure CC=/path/to/clang`

`CXX=/path/to/clang++ FC=/path/to/gfortran ...`

## Static or Shared ?

- `-enable-static / -disable-static` (default)  
libmpi.a
- `-enable-shared / -disable-shared`  
libmpi.so

## Communication Library

UCX (Unified Communication X)

`-with-ucx[=UCX_INSTALL_DIR]`

## With CUDA support

`./configure --with-cuda[=/path/to/cuda]`

- -x [env]  
Passes environment variables to remote nodes.
- -bind-to core
- -hostfile [hostfile]
- ...

- `$ module load openmpi/5.0.3-pe46zvn`
- `$ module load intel-oneapi-mpi/2021.13.0-hpxfbao`

Later lectures. (7.12)

# Thank You

## Questions?

<https://rookiehpc.org/mpi/docs/index.html>

<https://docs.open-mpi.org/en/v5.0.x/index.html>

<https://www.intel.com/content/www/us/en/docs/mpi-library/developer-reference-linux/2021-13/overview.html>