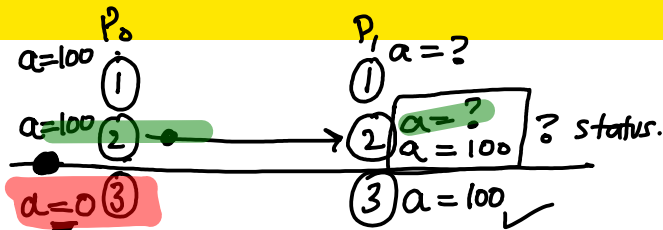


Important MPI Calls

Blocking

Non-Blocking ↓

`MPI_Init(int*, char**);``MPI_Finalize(void);``MPI_Comm_size(MPI_COMM, int);``MPI_Comm_rank(MPI_COMM, int);``MPI_Get_processor_name(char*, int);``MPI_Wtime(void);``MPI_Abort(MPI_COMM);`

// Initiate an MPI Computation

// Terminate an MPI Computation

// How many processes

// Who am I?

// What is the hostname?

// Elapsed time in seconds

// Terminate all processes

① → `int a = 100;`② → `send(&a, P1)`③ → `a = 0`① → `int a`② → `receive(&a, P0)`③ → `printf("%d", a);`

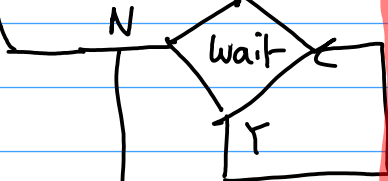
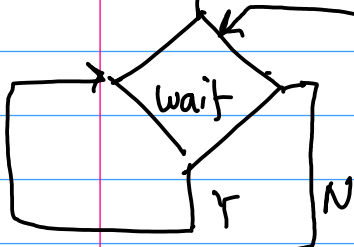
int a=100

int a

idle
wait

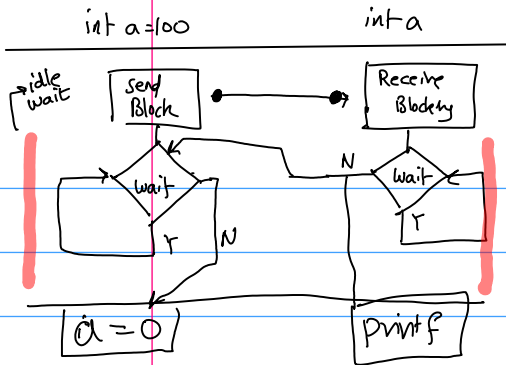
send
Block

Receive
Block



a=0

printf



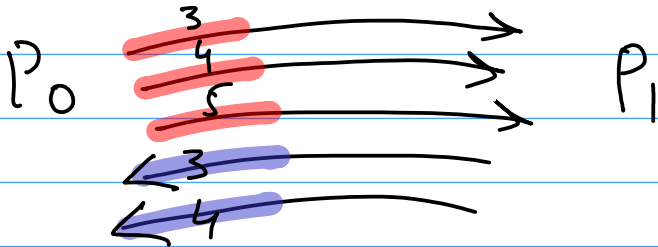
```

int a=100
if (my rank == 0)
{
    send(a, 1)
    ↓ ← Receive
}
a = 0
  
```

```

int a;
if (my rank == 1)
{
    send(a, 0)
    ↓      ↑ recv.
}
printf
  
```

tag



int

date

ack

Approaches to Send/Receive

Blocking (Non-Buffered) Send/Receive

- Follow some form of “handshaking” protocol
 - Request to Send → Clear to Send → Send Data → Acknowledgement
 - Problem 1: Idling Overhead (both sender/receiver side)
 - Problem 2: Deadlock (sending at same time)

Blocking (Buffered) Send/Receive

- Copy send-data to designated buffer, and returns after “copy” operation is completed

- Problem 1: Buffer Size

```
for (i = 0; i < 1000; i++) {
    produce_data(&a);
    send(&a, P1);
}
```

```
for (i = 0; i < 1000; i++) {
    receive(&a, P0);
    consume_data(&a);
}
```

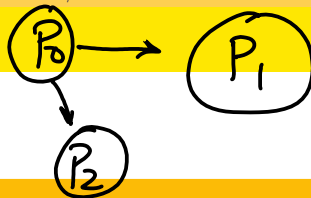
- Problem 2: Deadlock (sending at same time)

Approaches to Send/Receive (cont.)

Non-Blocking Send/Receive

- Return from Send/Receive operation before it is “safe” to return.
 - Programmer responsibility to ensure that “sending data” is not altered immediately
-
- Blocking Operations: Safe and Easy Programming (at cost of overhead and risk of deadlocks)
 - Non-Blocking Operations: Useful for Performance optimization, and breaking deadlocks (but brings in plenty of race-conditions if programmer not careful)

Point to Point Communication



Types of Point-to-Point Send/Receive Calls

- **Synchronous Transfer:** Send/Receive routines return only when the message transfer is completed. Not only does this transfer data, but it also synchronizes processes

`MPI_Send()` *// Blocking Send*

`MPI_Recv()` *// Blocking Receive*

- **Asynchronous Transfers:** Send/Receive do not wait for transfer data and proceeds with execution next line of instruction. (Precaution: Do not modify the send/receive buffers)

`MPI_Isend()` *// Non-Blocking Send*

`MPI_Irecv()` *// Non-Blocking Receive*

Point to Point Communication (cont.)

① MPI-Send(&a, 1 - - - -)
 ② MPI-Send(a, 100 - - - -)

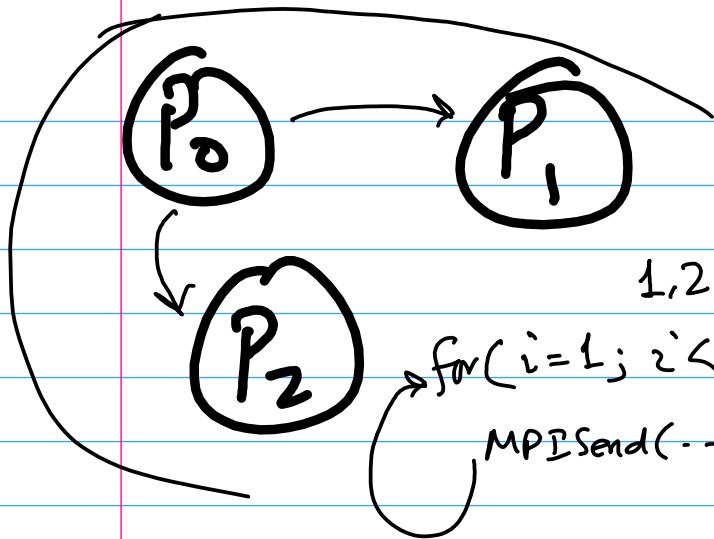
Sending

```
int MPI_Send(void *buffer,    int count,    MPI_DATATYPE datatype,
              int destination, int tag,      MPI_Comm comm);
```

- Send the data stored in **buffer** •
- **Count** is the number of entries in the buffer
- What is the **datatype** of the buffer (MPI_CHAR, MPI_INT, MPI_FLOAT, MPI_DOUBLE, MPI_LONG_DOUBLE, MPI_LONG, MPI_SHORT, MPI_UNSIGNED_CHAR, etc.)
- **Destination** is the rank of process, to whom buffer is to be sent to, residing in communication universe **comm**
- The **tag** of the message (to distinguish between different types of messages)

① int a
 ② int a[100]

get-size MPI Comm-world



Point to Point Communication (cont.)

Receiving

```
int MPI_Recv(void *buffer,    int count,    MPI_DATATYPE datatype,
             int source,      int tag,      MPI_Comm comm,
             MPI_Status *status);
```

- Store the received message in **buffer**
- **Count** is the number of entries to be received in the buffer. If number of entries is larger than the capacity of buffer, an overflow error `MPI_ERR_TRUNCATE` is returned.
- **Datatype** is the type of data that has been received
- **Source** is the rank of process, residing in communication domain **comm**, from whom buffer is received. Source can be hard-set, or a wild-card `MPI_ANY_SOURCE`.
- To retrieve message of certain type, set the **tag** argument. If there are many messages of same tag from same process, any one of them may be retrieved. If message of any tag is to be retrieved, use the wild-card `MPI_ANY_TAG`.
- Store status of received message in **status** (next slide). If not needed, use `MPI_STATUS_IGNORE`