# Messages and Point-to-Point Communication

1. MPI Overview

2. Process model and language bindings

   MPI_Init()
   MPI_Comm_rank()

3. **Messages and point-to-point communication**
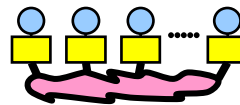   – **the MPI processes can communicate**

4. Non-blocking communication

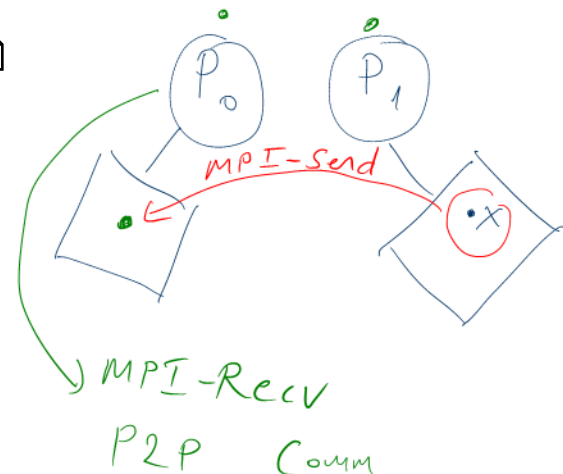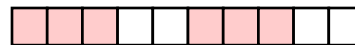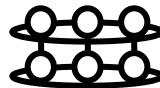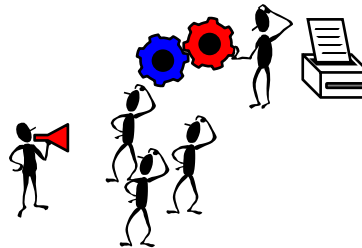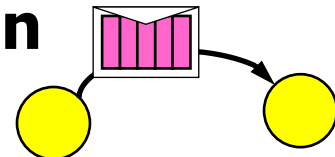5. Collective communication

6. Virtual topologies

7. Derived datatypes

8. Case study

Sub Raitme

MPI-Send

MPI-Recv

P2P Comm

# Messages

- A message contains a number of elements of some particular datatype.

  *Int, float, char, struct, ---*

- MPI datatypes:

  → *MPI_INT , MPI_FLOAT,*

  *MPI_CHAR , ---*

  ✓ – Basic datatype.
  - Derived datatypes

- C types are different from Fortran types.

- Datatype handles are used to describe the type of the data in the memory.

  *float d[7];*

  *float *d*

  *= malloc(7 × sizeof*

  *(float *)*   *(float))*

Example: message with 5 integers

*Int *array =*

| 2345 | 654 | 96574 | -12 | 7676 |
|------|-----|-------|-----|------|

*array ⇒ &array[0]*      *&array[1]*      *&array[4]*

## C

Row oriented
(based)



a

4 × 6

## F

column
oriented
(based)



a

$a[0]$

$a[3]$

$a[0]$

$a[3]$

$i\ k\ j \longrightarrow$ ???

| MPI Datatype | C datatype |
|---|---|
| MPI_CHAR | signed char |
| MPI_SHORT | signed short int |
| MPI_INT | signed int |
| MPI_LONG | signed long int |
| MPI_UNSIGNED_CHAR | unsigned char |
| MPI_UNSIGNED_SHORT | unsigned short int |
| MPI_UNSIGNED | unsigned int |
| MPI_UNSIGNED_LONG | unsigned long int |
| MPI_FLOAT | float |
| MPI_DOUBLE | double |
| MPI_LONG_DOUBLE | long double |
| MPI_BYTE | |
| MPI_PACKED | (Mixed) |

| MPI Datatype | Fortran datatype |
|---|---|
| MPI_INTEGER | INTEGER |
| MPI_REAL | REAL |
| MPI_DOUBLE_PRECISION | DOUBLE PRECISION |
| MPI_COMPLEX | COMPLEX |
| MPI_ LOGICAL | LOGICAL |
| MPI_CHARACTER | CHARACTER(1) |
| MPI_BYTE | |
| MPI_PACKED | |

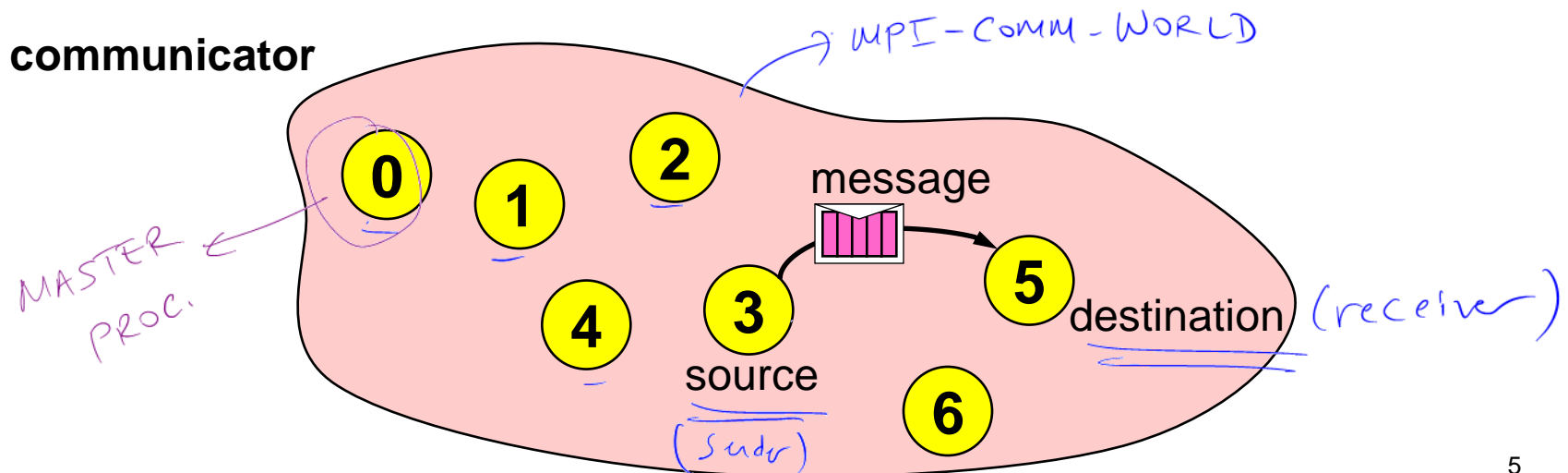| 2345 | 654 | 96574 | -12 | 7676 |
|---|---|---|---|---|

count=5                                    INTEGER arr(5)
datatype=MPI_INTEGER

# Point-to-Point Communication

- Communication between two processes.
- Source process sends message to destination process.
- Communication takes place within a communicator, e.g., MPI_COMM_WORLD.
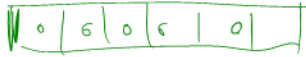- Processes are identified by their ranks in the communicator.

# Sending a Message

- C: int MPI_Send(void *buf, int count, MPI_Datatype datatype, int dest, int tag, MPI_Comm comm)

- Fortran: MPI_SEND(BUF, COUNT, DATATYPE, DEST, TAG, COMM, *IERROR*)

    <type> BUF(*)
    INTEGER COUNT, DATATYPE, DEST, TAG, COMM, IERROR

- <u>buf</u> is the starting point of the message with <u>count</u> elements, each described with <u>datatype</u>.

- <u>dest</u> is the rank of the destination process within the communicator <u>comm</u>.

- <u>tag</u> is an additional nonnegative integer piggyback information, additionally transferred with the message.

- The tag can be used by the program to distinguish different types of messages.

# Receiving a Message

- C: int MPI_Recv(void *_buf_, int count, MPI_Datatype datatype,
  int source, int tag, MPI_Comm comm,
  MPI_Status *_status_)

- Fortran: MPI_RECV(_BUF_, COUNT, DATATYPE, SOURCE, TAG,
  COMM, _STATUS_, _IERROR_)
  <type> BUF(*)
  INTEGER COUNT, DATATYPE, SOURCE, TAG, COMM
  INTEGER STATUS(MPI_STATUS_SIZE), IERROR

- buf/count/datatype describe the receive buffer.

- Receiving the message sent by process with rank <u>source</u> in <u>comm</u>.

- Envelope information is returned in _status_.

- Output arguments are printed _blue-cursive_.

- Only messages with matching <u>tag</u> are received.

# Requirements for Point-to-Point Communications

For a communication to succeed:

- Sender must specify a valid destination rank.

- Receiver must specify a valid source rank.

- The communicator must be the same.

- Tags must match.

- Message datatypes must match.

- Receiver's buffer must be large enough.

# **Wildcards**

- Receiver can wildcard.

- To receive from any source — <u>source</u> = MPI_ANY_SOURCE

- To receive from any tag — <u>tag</u> = MPI_ANY_TAG

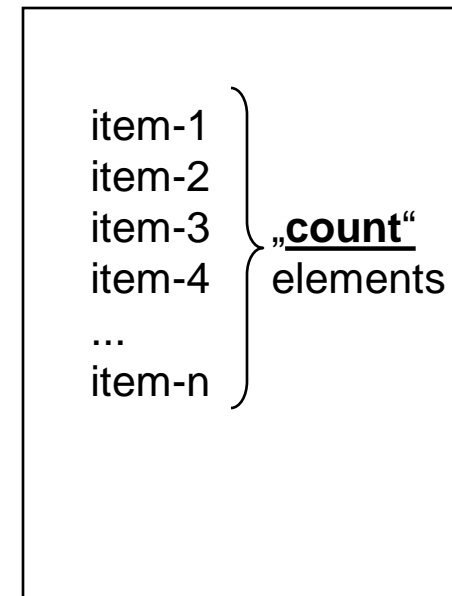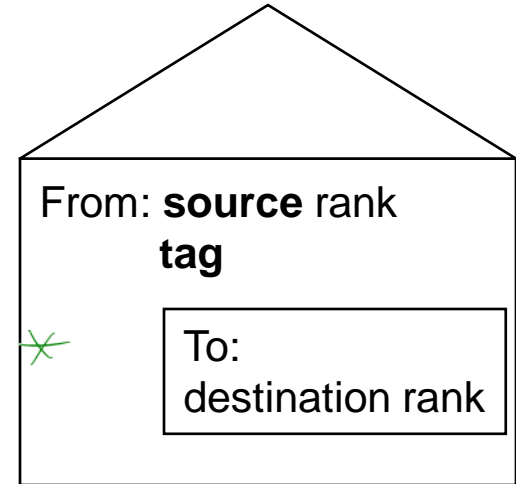- Actual source and tag are returned in the receiver's *status* parameter.

# Communication Envelope

- Envelope information is returned from MPI_RECV in *status*.

- C:        status.MPI_SOURCE
           status.MPI_TAG
           count via MPI_Get_count()

- Fortran:   status(MPI_SOURCE)
            status(MPI_TAG)
            count via MPI_GET_COUNT()

From: **source** rank
**tag**

To:
destination rank

item-1
item-2
item-3         „**count**"
item-4         elements
...
item-n

# Receive Message Count

- C: int MPI_Get_count(MPI_Status *status, MPI_Datatype datatype,
                        int  *_count_)

  → Output parameter.

- Fortran: MPI_GET_COUNT(STATUS, DATATYPE, _COUNT, IERROR_)

  INTEGER STATUS(MPI_STATUS_SIZE)
  INTEGER DATATYPE, COUNT, IERROR

# COMMUNICATION IN MPI

**P2P**

MPI-Send
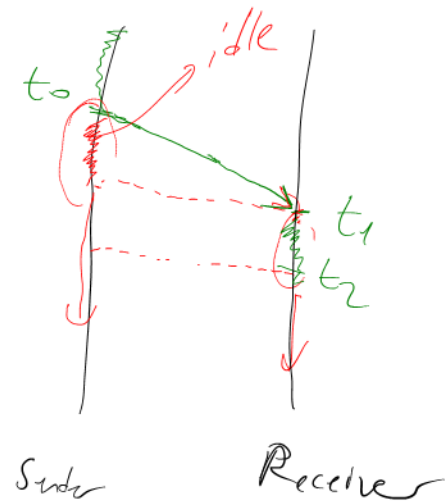MPI-Recv

Blocking

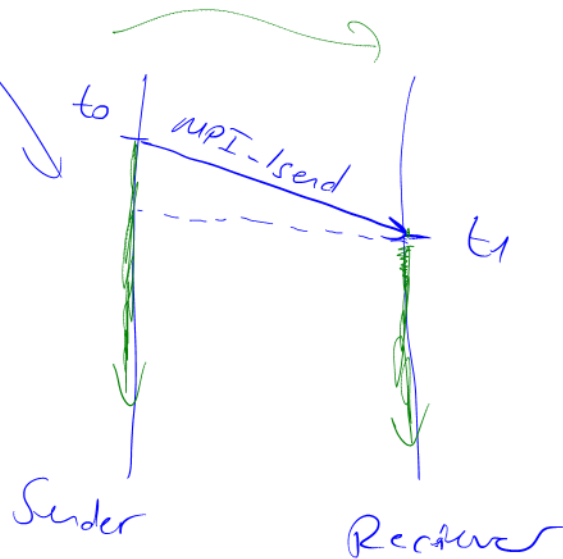Non-blocking

MPI-Isend → Immediate
MPI-IRecv

**Collective** → Next Week !!!

Blocking: Safer
Non-blocking: Faster



idle

$t_0$

$t_1$

$t_2$

↓ time

Sender    Receiver



$t_0$    MPI-Isend    $t_1$

Sender    Receiver

$P_0$     $P_1$     $P_2$

$\longrightarrow$ MPI-Init( )

$t_0$

IDLE

$t_1$

$\longrightarrow$ IDLE

$t_2$

$\longrightarrow$ Blocking

IDLE

$t_3$

$\longrightarrow$ IDLE

$t_4$

$P_0$     $P_1$

MPI-Isend

MPI-Isend

$\longrightarrow$ MPI_IRecv

NoT
IDLE

MPI_IRecv

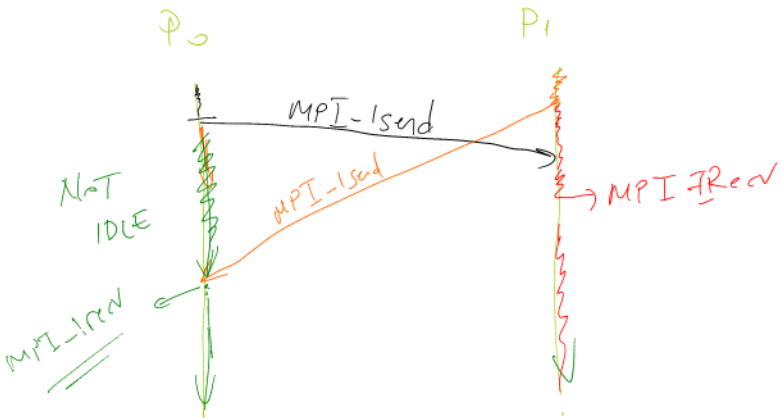$\longrightarrow$ Non-blocking

(Avoiding deadlock)

# Communication Modes

- Send communication modes:
  - synchronous send   → MPI_**S**SEND
  - buffered [asynchronous] send→ MPI_**B**SEND
  - standard send   → MPI_**SEND**
  - Ready send   → MPI_**R**SEND

- Receiving all modes → MPI_**RECV**

# Communication Modes — Definitions

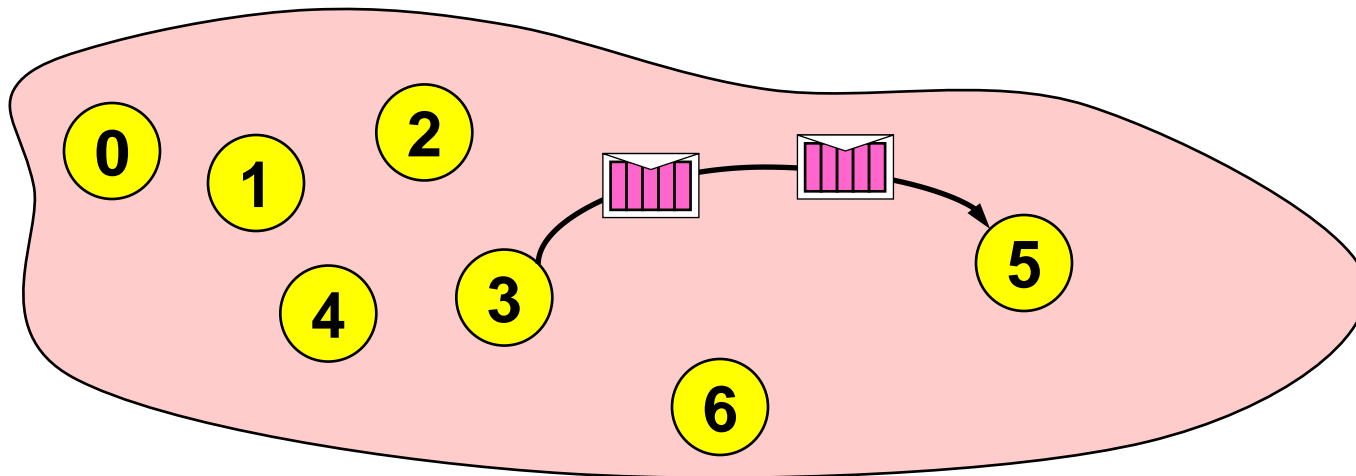| Sender modes | Definition | Notes |
|---|---|---|
| Synchronous send<br>**MPI_SSEND** | Only completes when the receive has started | |
| Buffered send<br>**MPI_BSEND** | Always completes<br>(unless an error occurs), irrespective of receiver | needs application-defined buffer to be declared with<br>MPI_BUFFER_ATTACH |
| Synchronous<br>**MPI_SEND** | Standard send. Either uses an internal buffer or buffered | |
| Ready send<br>**MPI_RSEND** | May be started **only** if the matching receive is already posted! | highly dangerous! |
| Receive<br>**MPI_RECV** | Completes when a the message (data) has arrived | |

# Rules for the communication modes

- Standard send  (**MPI_SEND**)
  - minimal transfer time
  - may block due  to synchronous mode
  - —> risks with synchronous send
- Synchronous send  (**MPI_SSEND**)
  - risk of deadlock
  - risk of serialization
  - risk of waiting —> idle time
  - high latency  /  best bandwidth
- Buffered send  (**MPI_BSEND**)
  - low latency  /  bad bandwidth
- Ready send  (**MPI_RSEND**)
  - use **never**, except you have a 200% guarantee that Recv is already called in the current version and all future versions of your code

# Message Order Preservation

- Rule for messages on the same connection, i.e., same communicator, source, and destination rank:
- **Messages do not overtake each other.**
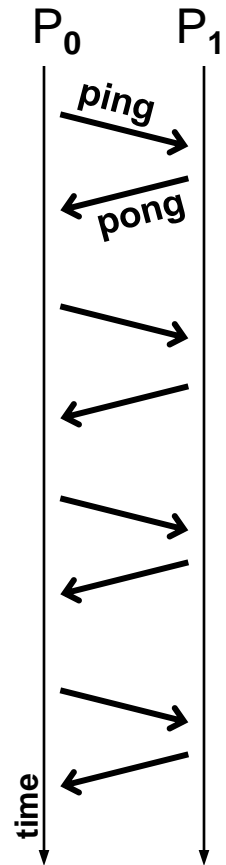- This is true even for non-synchronous sends.



- If both receives match both messages, then the order is preserved.

# Exercise — Ping pong

- Write a program according to the time-line diagram:
  - process 0 sends a message to process 1 (ping)
  - after receiving this message,
    process 1 sends a message back to process 0 (pong)
- Repeat this ping-pong with a loop of length 50
- Add timing calls before and after the loop:
- C: *double MPI_Wtime*(void);
- Fortran: *DOUBLE PRECISION FUNCTION MPI_WTIME*()
- MPI_WTIME returns a wall-clock time in seconds.
- At process 0, print out the transfer time of **one** message
  - in seconds
  - in μs.

$P_0$    $P_1$

ping

pong

time

Ping - Pong Example :   Wimbledon Final      Federer  vs   Nadal

(0.88)                    (0.91)
Federer                   Nadal

1 ball ++;

ball ++; 2

ball ++;

int  ball = 0;

ball % 2 == 1 ⟹ Federer wins
          else
             ⟹ Nadal  wins

# Exercise — Ping pong

rank=0                          rank=1

Send **(dest=1)**

            **(tag=17)**
                              Recv **(source=0)**
                              Send **(dest=0)**

            **(tag=23)**

Recv **(source=1)**

---

```
if (my_rank==0)           /* i.e., emulated multiple program */
    MPI_Send( ... dest=1 ...)
    MPI_Recv( ... source=1 ...)
else
    MPI_Recv( ... source=0 ...)
    MPI_Send( ... dest=0 ...)
fi
```

# Advanced Exercise - Measure latency and bandwidth

- latency = transfer time for zero length messages
- bandwidth = message size (in bytes) / transfer time

- Print out message <u>transfer time</u> and <u>bandwidth</u>
  - for following send modes:
    - for standard send (MPI_Send)
    - for synchronous send (MPI_Ssend)
  - for following message sizes:
    - 8 bytes (e.g., one double or double precision value)
    - 512 B (= 8*64 bytes)
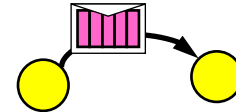    - 32 kB (= 8*64**2 bytes)
    - 2 MB  (= 8*64**3 bytes)
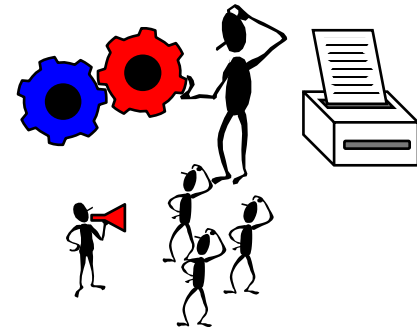
# Chap.4 Non-Blocking Communication

1. MPI Overview

   `MPI_Init()`
   `MPI_Comm_rank()`

2. Process model and language bindings

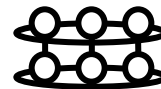3. Messages and point-to-point communication

4. **Non-blocking communication**
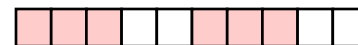   – **to avoid idle time and deadlocks**

5. Collective communication
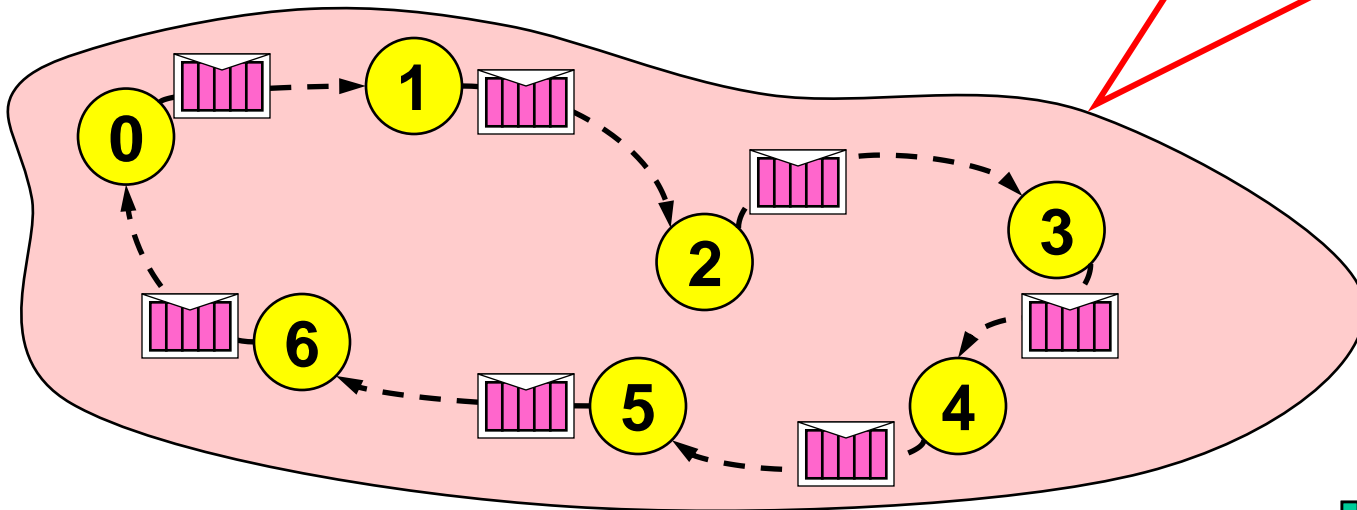
6. Virtual topologies

7. Derived datatypes

8. Case study

# Deadlock

- Code in each MPI process:
  MPI_Ssend(..., right_rank, ...)
  MPI_Recv(   ..., left_rank,    ...)

Will block and never return, because MPI_Recv cannot be called in the right-hand MPI process



- Same problem with standard send mode (MPI_Send), if MPI implementation chooses synchronous protocol
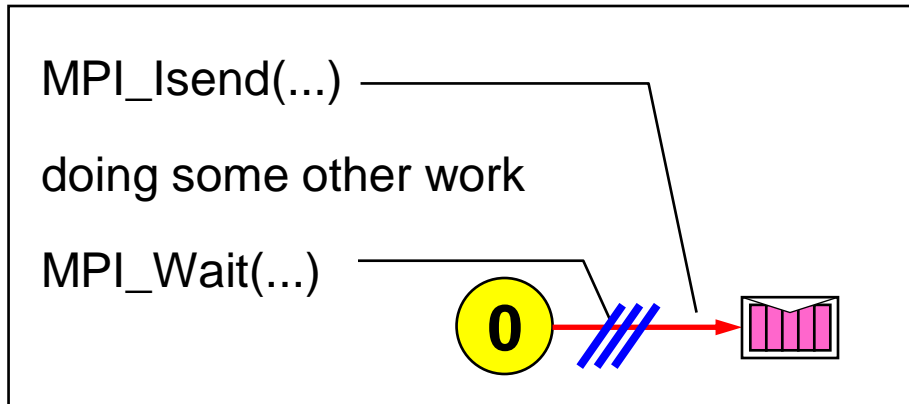
# Non-Blocking Communications

- Separate communication into three phases:
- Initiate non-blocking communication
  - returns **I**mmediately
  - routine name starting with MPI_**I**…
- Do some work
  - "latency hiding"
- Wait for non-blocking communication to complete

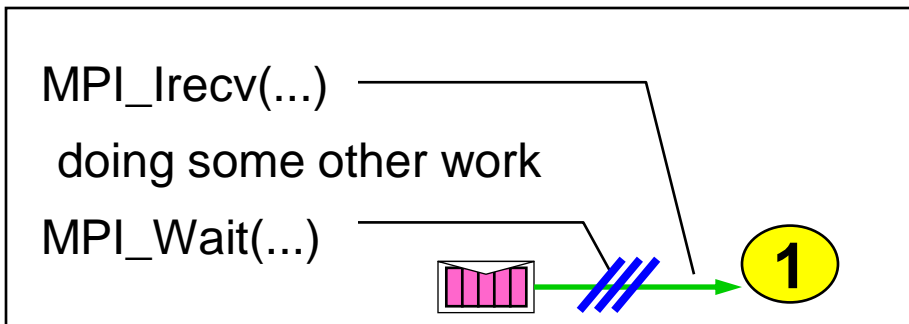# Non-Blocking Examples

- Non-blocking **send**

MPI_Isend(...)

doing some other work

MPI_Wait(...)

**0**

- Non-blocking **receive**

MPI_Irecv(...)

doing some other work
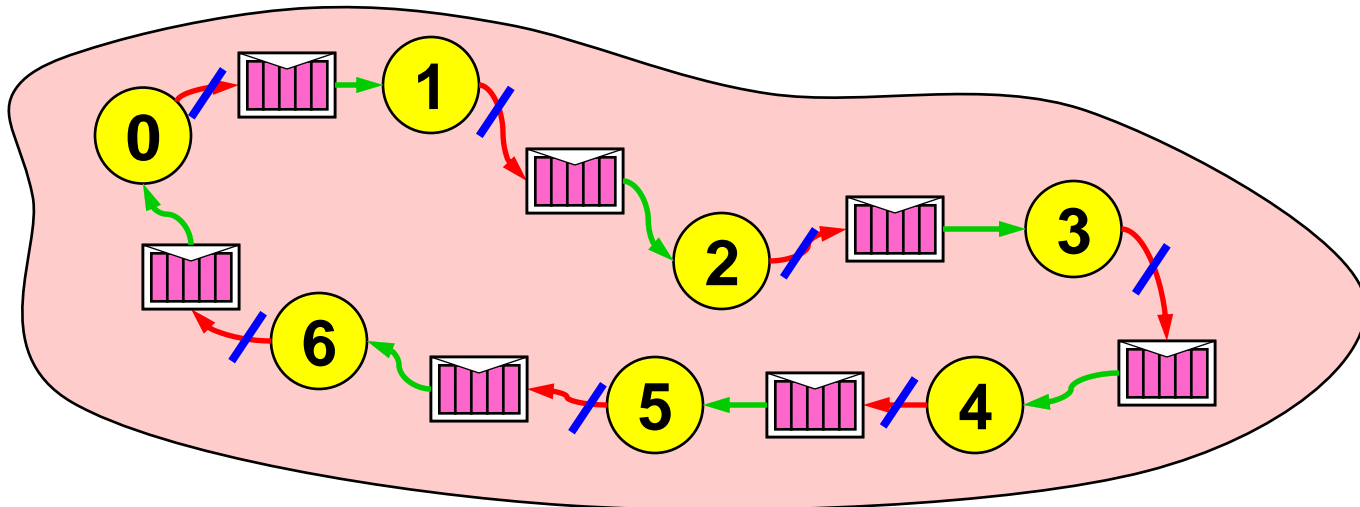
MPI_Wait(...)

**1**

/// = waiting until operation locally completed
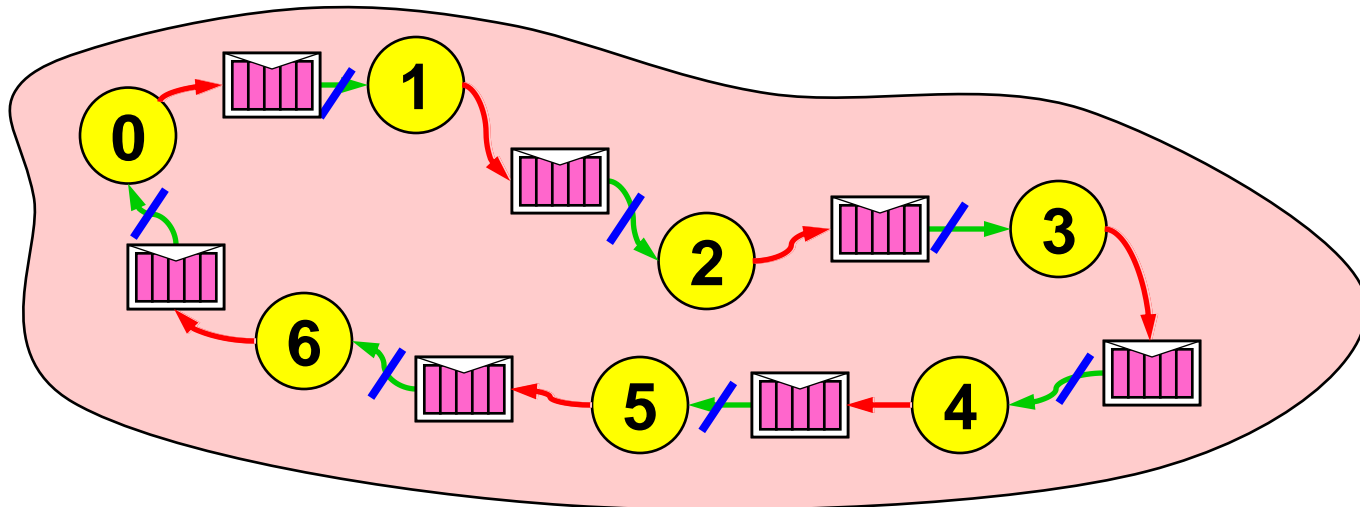
# Non-Blocking Send

- Initiate non-blocking send
  - ➡ in the ring example: Initiate non-blocking send to the right neighbor
- Do some work:
  - ➡ in the ring example: Receiving the message from left neighbor
- Now, the message transfer can be completed
- Wait for non-blocking send to complete

# Non-Blocking Receive

- Initiate non-blocking receive
  - ⟶ in the ring example: Initiate non-blocking receive from left neighbor
- Do some work:
  - ⟶ in the ring example: Sending the message to the right neighbor
- Now, the message transfer can be completed
- Wait for non-blocking receive to complete

# Handles, already known

- Predefined handles
  - defined in mpi.h / mpif.h
  - communicator, e.g., MPI_COMM_WORLD
  - datatype, e.g., MPI_INT, MPI_INTEGER, …

- Handles **can** also be stored in local variables
  - memory for datatype handles
    - in C: MPI_Datatype
    - in Fortran: INTEGER
  - memory for communicator handles
    - in C: MPI_Comm
    - in Fortran: INTEGER

# Request Handles

## Request handles

- are used for non-blocking communication

- **must** be stored in local variables

  C: MPI_Request

  Fortran: INTEGER

- the value
  - **is generated** by a non-blocking communication routine
  - **is used** (and freed) in the MPI_WAIT routine

# Non-blocking Synchronous Send

- C:

  MPI_Issend(buf, count, datatype, dest, tag, comm,
             OUT &*request_handle*);

  MPI_Wait(INOUT &request_handle, &*status*);

- Fortran:

  CALL MPI_ISSEND(buf, count, datatype, dest, tag, comm,
                  OUT *request_handle*, *ierror*)

  CALL MPI_WAIT(INOUT request_handle, *status*, *ierror*)

- <u>buf</u> must not be used between <u>Issend</u> and <u>Wait</u>  (in all progr. languages)
  MPI 1.1, page 40, lines 44-45

- "<u>Issend</u> + <u>Wait</u> directly after Issend" is equivalent to blocking call (<u>Ssend</u>)

- <u>status</u> is not used in <u>Issend</u>, but in <u>Wait</u> (with send: nothing returned)

- Fortran problems, see MPI-2, Chap. 10.2.2, pp 284-290

# Non-blocking Receive

- C:

  MPI_Irecv(*buf*, count, datatype, source, tag, comm,
  OUT &*request_handle*);

  MPI_Wait(INOUT &request_handle, &*status*);

- Fortran:

  CALL MPI_IRECV (*buf*, count, datatype, source, tag, comm,
  OUT *request_handle*, *ierror*)

  CALL MPI_WAIT(  INOUT request_handle, *status*, *ierror*)

- <u>buf</u> must not be used between <u>Irecv</u> and <u>Wait</u> (in all progr. languages)

# Blocking and Non-Blocking

- Send and receive can be blocking or non-blocking.

- A blocking send can be used with a non-blocking receive,
and vice-versa.

- Non-blocking sends can use any mode
  - standard        – MPI_ISEND
  - synchronous    – MPI_ISSEND
  - buffered        – MPI_IBSEND
  - ready          – MPI_IRSEND

# Completion

- C:

    MPI_Wait( &request_handle, &*status*);

    MPI_Test( &request_handle, &*flag*, &*status*);

- Fortran:

    CALL MPI_WAIT( request_handle, *status*, *ierror*)

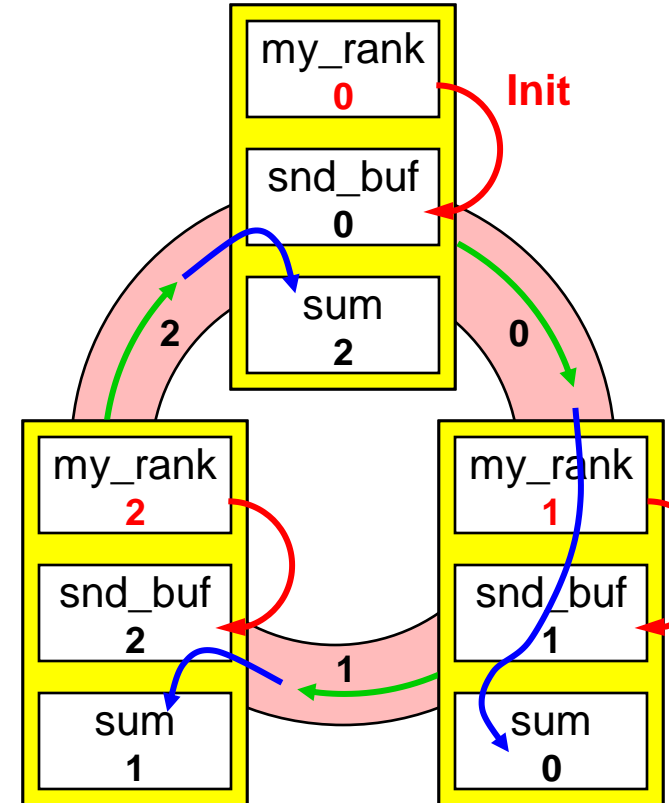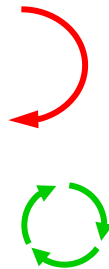    CALL MPI_TEST( request_handle, *flag*, *status*, *ierror*)

- one must
    - WAIT or
    - loop with TEST until request is completed, i.e., flag == 1 or .TRUE.

# Exercise — Rotating information around a ring

- A set of processes are arranged in a ring.
- Each process stores its rank in MPI_COMM_WORLD into an integer variable *snd_buf*.
- Each process passes this on to its neighbor on the right.
- Each processor calculates the sum of all values.
- Keep passing it around the ring until the value is back where it started, i.e.
- each process calculates sum of all ranks.
- Use non-blocking MPI_Issend
  - to avoid deadlocks
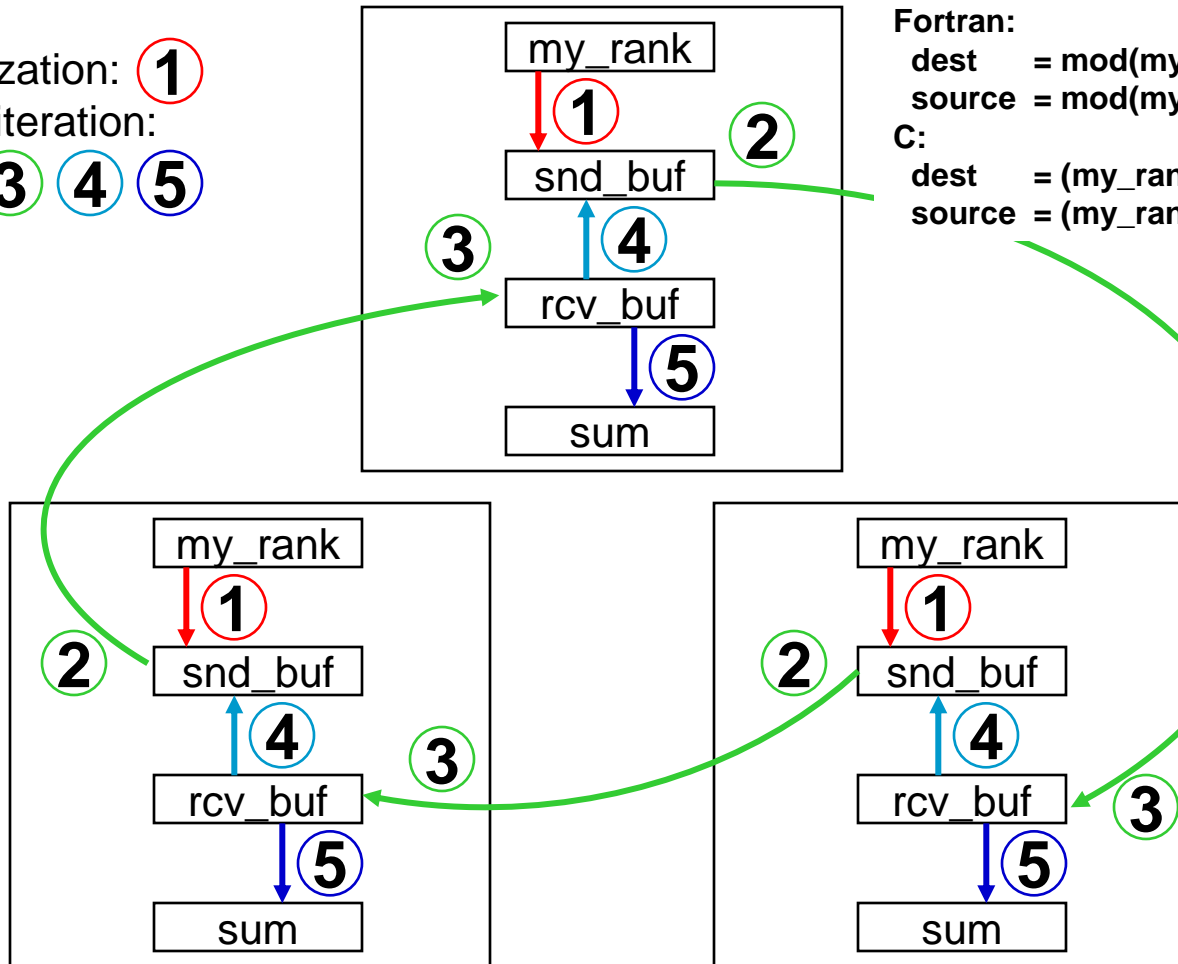  - to verify the correctness, because blocking synchronous send will cause a deadlock

Initialization: ①
Each iteration:
② ③ ④ ⑤

my_rank
①
snd_buf ②
③ ④
rcv_buf
⑤
sum

my_rank
①
② snd_buf
④
③
rcv_buf
⑤
sum

my_rank
①
② snd_buf
④
rcv_buf ③
⑤
sum

**Fortran:**
 **dest      = mod(my_rank+1,size)**
 **source  = mod(my_rank–1+size,size)**
**C:**
 **dest      = (my_rank+1) % size;**
 **source  = (my_rank–1+size) % size;**

**Single
Program !!!**