

# The MPI Message-passing Standard

## Practical use and implementation (III)

SPD Course

05/03/2014

Massimo Coppola

# POINT-TO-POINT COMMUNICATION MODES

# Buffered Send

## **MPI\_BSEND (buf, count, datatype, dest, tag, comm)**

```
MPI_Bsend(void* buf, int count, MPI_Datatype datatype,  
int dest, int tag, MPI_Comm comm)
```

- Same parameters as the standard send
- Explicitly relies on buffering
- Programmer has to allocate enough buffers for the process needs, and pass them to the MPI implementation

```
int MPI_Buffer_attach(void* buffer, int size)
```

```
int MPI_Buffer_detach(void* buffer_addr, int*  
size)
```

## MPI\_SSEND (buf, count, datatype, dest, tag, comm)

```
MPI_Ssend(void* buf, int count, MPI_Datatype  
          datatype, int dest, int tag, MPI_Comm comm)
```

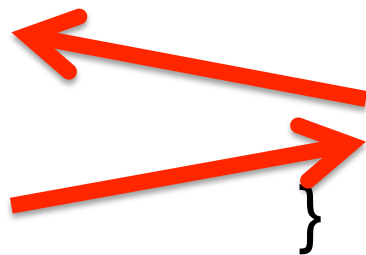
- Same parameters as the standard send
- Enforces synchronous send operation
  - A program is **safe** if all its sends are Synchronous

## **MPI\_RSEND (buf, count, datatype, dest, tag, comm)**

- Again same parameters
  - Optimizes implementation assuming a matching receive has been already posted
    - Used with **permanent** requests
    - When program semantics ensures the precondition
    - Together With SendRecv primitives
    - Note that:
      - Permanent requests and SendRecv are used solely as example cases
- SendRec a single primitive for send and receive combined

# Ready Send and SendRecv

```
// Process A
while (true) {
    recv ( ... B ...)
    do_compute()
    Rsend ( ...B... )
}
```



```
//Process B
while (true) {
    do_compute()
    sendRecv( ...A...)
}
```

# BLOCKING AND NON-BLOCKING POINT-TO-POINT

# Incomplete operations

- Separate communication **start** from its **completion**
- Available for **both** send and receive
- Primitive calls can return before completion
- Resources are NOT free
- Separate primitives for checking communication completion/status
- Useful if actual communication is offloaded to DMA, coprocessors etc.



# Incomplete Send / Recv

**MPI\_ISEND(buf, count, datatype, dest, tag, comm, request)**  
**MPI\_IRecv (buf, count, datatype, source, tag, comm, request)**

```
int MPI_Isend(void* buf, int count,  
             MPI_Datatype datatype, int dest, int tag, MPI_Comm  
             comm, MPI_Request *request)
```

```
int MPI_Irecv(void* buf, int count,  
             MPI_Datatype datatype, int source,  
             int tag, MPI_Comm comm, MPI_Request *request)
```

- MPI\_ISEND Also combines with all modes
- MPI\_IBSEND
- MPI\_ISSEND
- MPI\_IRSEND

# Request objects

- Opaque objects
- Fully identify the communication operation
  - One to one match with communications
  - Requests are allocated by MPI when they become **active** (communication started, but not completed)
  - Requests are active until completion is not checked
- Can provide status and completion information
- The MPI\_request type is the object handle
  - Uninitialized/**inactive** handle value:  
MPI\_REQUEST\_NULL
  - MPI does this whenever a request object is no longer needed (it becomes inactive) and it is freed

# Waiting and Testing

MPI\_WAIT(request, status)

- INOUT request request (handle)
- OUT status status object (Status)
- Waits until the operation is complete
  - Returns the operations status
  - Clears the request handle

MPI\_TEST(request, flag, status)

- Returns immediately
  - flag=true if the operation is complete
  - In this case, behaves as a completed WAIT
- Wait is a non-local operation, Test is a local one
- MPI\_REQUEST\_NULL handles are silently ignored

# Multiple Wait / Test

- MPI\_WAITANY (count, array\_of \_requests, index, status)
  - Wait for one request from an array to complete (nondeterministic behaviour, no fairness)
  - index=MPI\_UNDEFINED if no request is active
- MPI\_WAITALL (count, array\_of \_requests, array\_of \_statuses)
  - Wait for all requests to complete
- MPI\_WAITSOME (incount, array\_of \_requests, outcount, array\_of \_indices, array\_of \_statuses)
  - Wait for at least one request to complete, possibly several ones
  - You can implement your own preferred nondeterministic behaviour
  - outcount=MPI\_UNDEFINED if no request is active
- MPI\_TESTANY (count, array\_of \_requests, index, flag, status)
- MPI\_TESTALL (count, array\_of \_requests, flag, array\_of \_statuses)
- MPI\_TESTSOME (incount, array\_of \_requests, outcount, array\_of \_indices, array\_of \_statuses)

- It is safe to call again and again the same primitive: eventually, all requests become inactive
- MPI\_requests are handles
  - can be copied
  - it's programmer's responsibility not to use more than one copy (better invalidate them!)
- Null handle is not the same as inactive
  - MPI\_REQUEST\_NULL is also inactive ofc

# MORE DERIVED DATATYPE CONSTRUCTORS

## **MPI\_TYPE\_CREATE\_STRUCT (count, array\_of\_blocklengths, array\_of\_displacements, array\_of\_types, newtype)**

IN count      number of blocks (non-negative integer)

- also number of entries in arrays array\_of\_types,  
array\_of\_displacements and array\_of\_blocklengths

IN array\_of\_blocklength      elements in each block  
(array of non-negative integer)

IN array\_of\_displacements      byte displacement of  
each block (array of integer)

IN array\_of\_types      type of elements in each block  
(array of handles to datatype objects)

OUT newtype      new datatype (handle)

# Reference Texts

- MPI standard Relevant Material for 3<sup>rd</sup> lesson
  - Chapter 2:  
sec.
  - Chapter 3:  
sec. 3.5, 3.6 (3.6.1 can be skipped), 3.7, 3.11  
persistent comm.s and sendRecv are 3.9, 3.10
  - Chapter 4:  
sec. 4.1 – to 4.1.2, (skip 4.1.3, 4.1.4), 4.1.9 – 4.1.11