

The MPI Message-passing Standard Practical use and implementation (II)

SPD Course

03/10/2010

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)
```

Synchronous Send

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

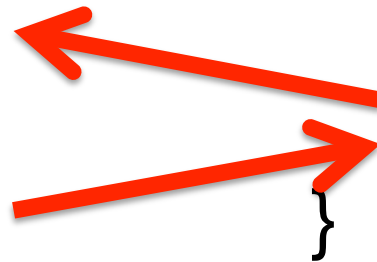
Ready Send and SendRecv

```
// 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

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_WAITSSOME(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_TESTSSOME(incount, array_of _requests, outcount, array_of _indices, array_of _statuses)

TEST/WAIT comments

- 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)

- MPI standard Relevant Material for 3rd 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