

# Optimised MPI for HPEC applications



- **Middleware Libraries and Application Programming Interfaces**
- **Software Architectures, Reusability, Scalability, and Standards**



## Systems used for Dataflow applications

- Computing power requirements no evenly spread
- Various transport medium may coexist
- Need for QoS type behaviour
- Performance requirement for I/O between nodes

## Requirements

- Need to map process to computing node
- Need to select specific link between process
- Need to implement zero-copy feature



## ■ PROs

- ➔ Available on almost every parallel/cluster machine
- ➔ Ensures application code portability

## ■ CONs

- ➔ Made for collective parallel apps, not distributed apps.
- ➔ No choice of communication interface (only know receiver)
- ➔ Does not care about transport medium
- ➔ No control on timeouts
- ➔ Not a communication library  
(no dynamic connection, no select feature)



## Zero-copy means memory management

- Same memory buffer used by application and I/O system
- At any given time, buffer must belong to application OR I/O

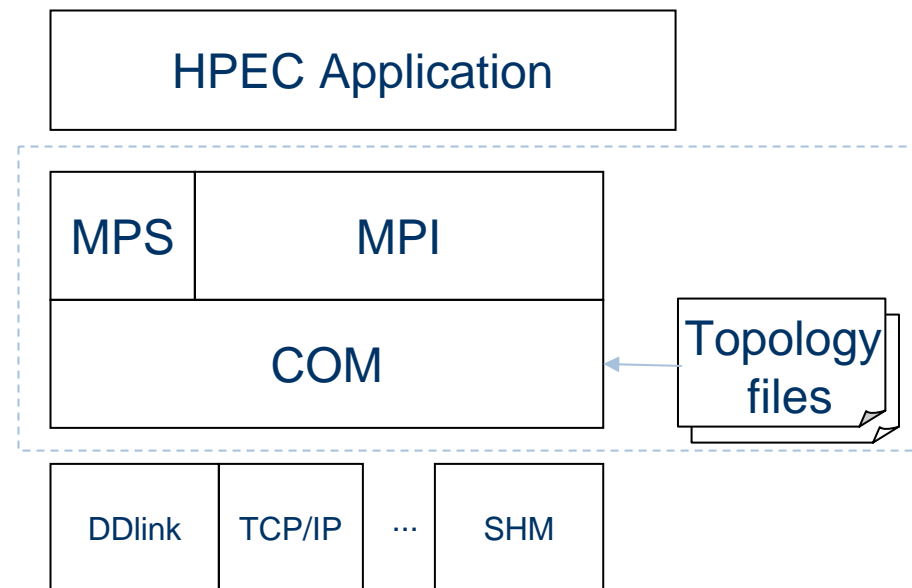
## Zero-copy API

- Buffer Get
  - ➔ Data buffer now part of application data
  - ➔ Can be used as any private memory
- Buffer Release
  - ➔ Data buffer is not to be modified by application any more
  - ➔ Can be used by I/O system (likely hardware DMA)



## ■ MPI Services (MPS) side to side with MPI

- ➔ MPI application source portability
- ➔ Links/Connector relationship
- ➔ Real-Time support
  - Links to select communication channels (~ QoS)
  - Requests timeout support
- ➔ Real zero-copy transfer
  - Buffer Management API (MPS)
- ➔ Heterogeneous machine support
  - Topology files outside application

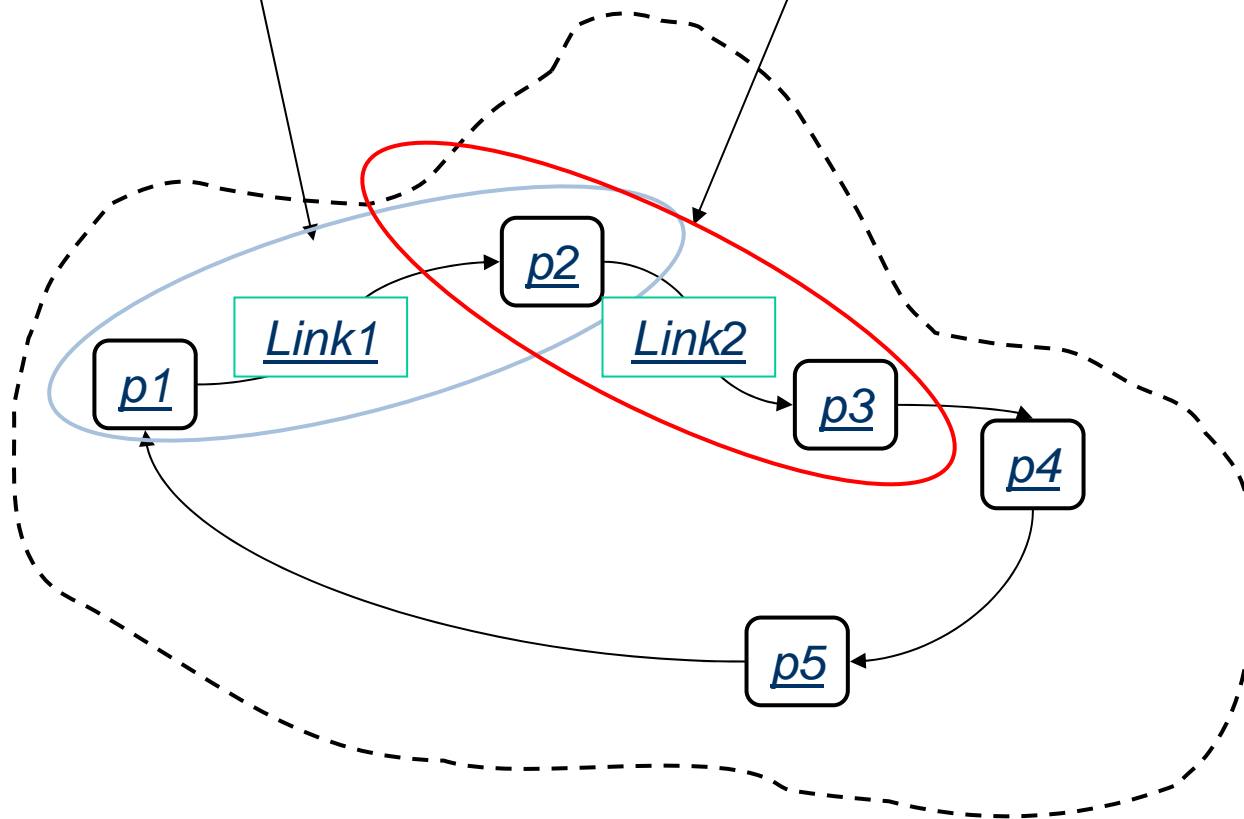




# Dedicated MPI Communicator for Zero-copy Link

com12

com23



MPI\_COMM\_WORLD

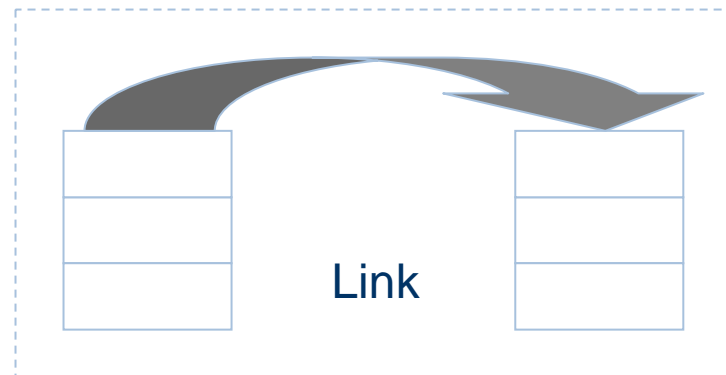
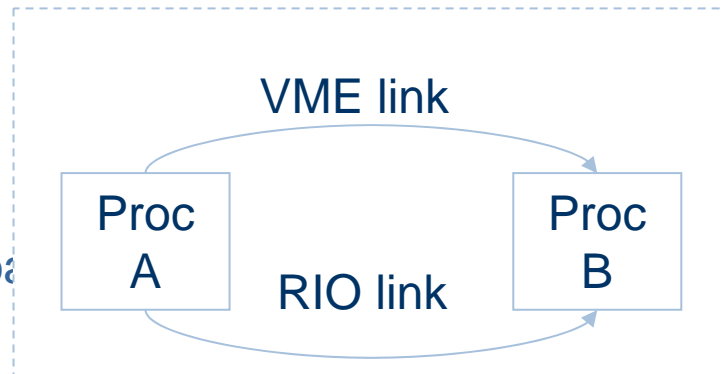
- System topology described outside the application code
- External ASCII files with:

## ➡ Process

- Process name
- Process Hardware location (board)

## ➡ Link

- Link name
- Medium type (+medium-specific parameters)
- Buffer size
- Buffer count



## MPS\_Channel\_create

(*\*chan\_name*, *\*rendpoint*, MPI\_Comm *\*comm*, int *\*lrank*, int *\*rrank*) ;

<i>link name</i>	^	^			
<i>remote end name</i>					
<i>specific communicator for the link</i>	v				
<i>my rank in new communicator</i>	v				
<i>remote end rank in new communicator</i>	v				

## MPS\_Process\_get\_name (int rank, char *\*name*) ;

<i>rank in MPI_COMM_WORLD</i>	^	
<i>my name in link/process file</i>	v	

## MPS\_Process\_get\_rank (char *\*name*, int *\*rank*) ;

<i>name in link/process file</i>	^	
<i>my rank in MPI_COMM_WORLD</i>	v	





## MPS\_Buf\_pool\_init

(MPI\_Comm com, way, \* p\_bufsize, \* p\_bufcount, \*p\_mps\_pool)

<i>MPI communicator</i>	^	^			
<i>Send or Receive</i>			v	v	
			<i>buffer size &amp; count</i>		v
					<i>MPS pool handle</i>

## MPS\_Buf\_get (p\_mps\_pool, void \*\*p\_buffer)

get buffer from pool (may block, or return EEMPTY)

## MPS\_Buf\_release (p\_mps\_pool, void \*buffer)

give buffer to I/O system (compulsory at each use) busy???

## MPS\_Buf\_pool\_finalize (p\_mps\_pool)

free all buffers, all coms must have completed first

```
MPI_Init(&argc, &argv);
```

*Create Dedicated Link  
Get Specific connector  
Initialize memory pool*

```
MPS_Channel_create("link1", "proc2", &com, &lrnk, &rrnk);
```

```
MPS_buf_pool_init(com, (sender) ? MPS_SND : MPS_RCV, &bufsize, &bufcount, &pool);
```

```
if (sender) {
```

```
    MPS_Buf_get(pool, &buf);
```

*Fill in with data*

```
    MPI_Isend(buf, size/sizeof(int), MPI_INT, rrank, 99, com, &req);
```

```
    MPI_Wait(req, &status);
```

```
    MPS_Buf_release(pool, buf);
```

*Take buffer ownership*

*Send on connector*

*Release buffer*

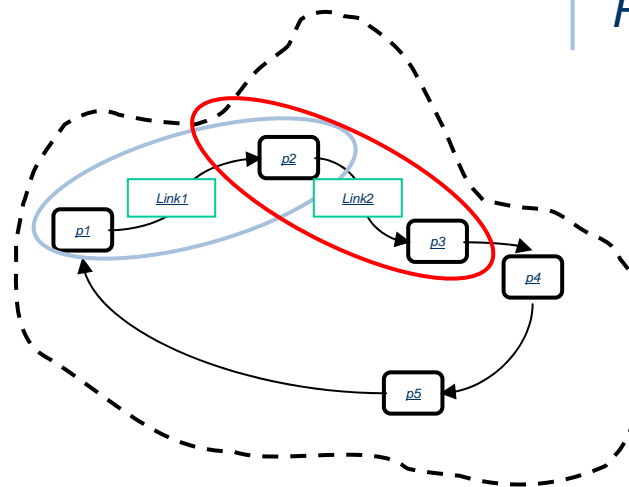
```
} else {
```

```
    ...
```

```
}
```

```
MPS_Buf_pool_finalize(pool);
```

```
MPI_Finalize();
```



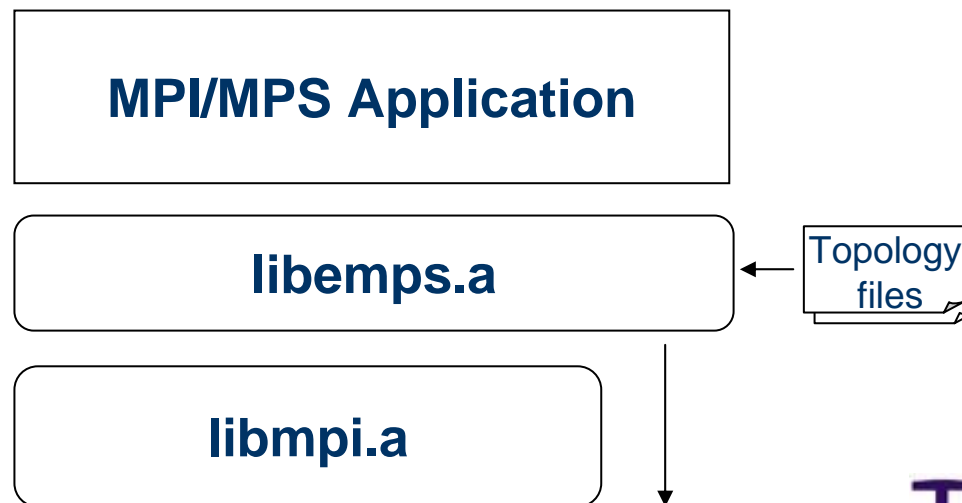


## MPI application easily ported to MPI/MPS API

- See example

## MPI/MPS application can run on any platform: EMPS

- EMPS is MPS emulation on top of standard MPI com
- Allow to run MPI/MPS code unmodified
  - ➡ Includes buffer and link management





## Based on MICH ?? Version etc...

### Software

- IA32 Red Hat, PowerPC LynxOS 4.0

### HW Targets

- PC, Thales multiprocessor VME boards

### Multi-protocol support in COM layer

- DDlink : Direct Deposit zero copy layer
  - ➔ Fibre Channel RDMA, Shared Memory, VME 2eSST, RapidIO
- Standard unix/posix I/O
  - ➔ Shared Memory, TCP/IP



## **Finalize process mapping**

- MPI\_RUN and HPEC compatible process mapping

## **Towards automatic code generation**

- Create MPS / MPI code from HPEC application tools

## **More support for MPI-aware debug tools**

- Like TotalView™

- Thank you  
vincent.chuffart@thalescomputers.fr