

# **Introduction to Programming by MPI for Parallel FEM Report S1 & S2 in C (2/2)**

Kengo Nakajima  
Information Technology Center  
The University of Tokyo

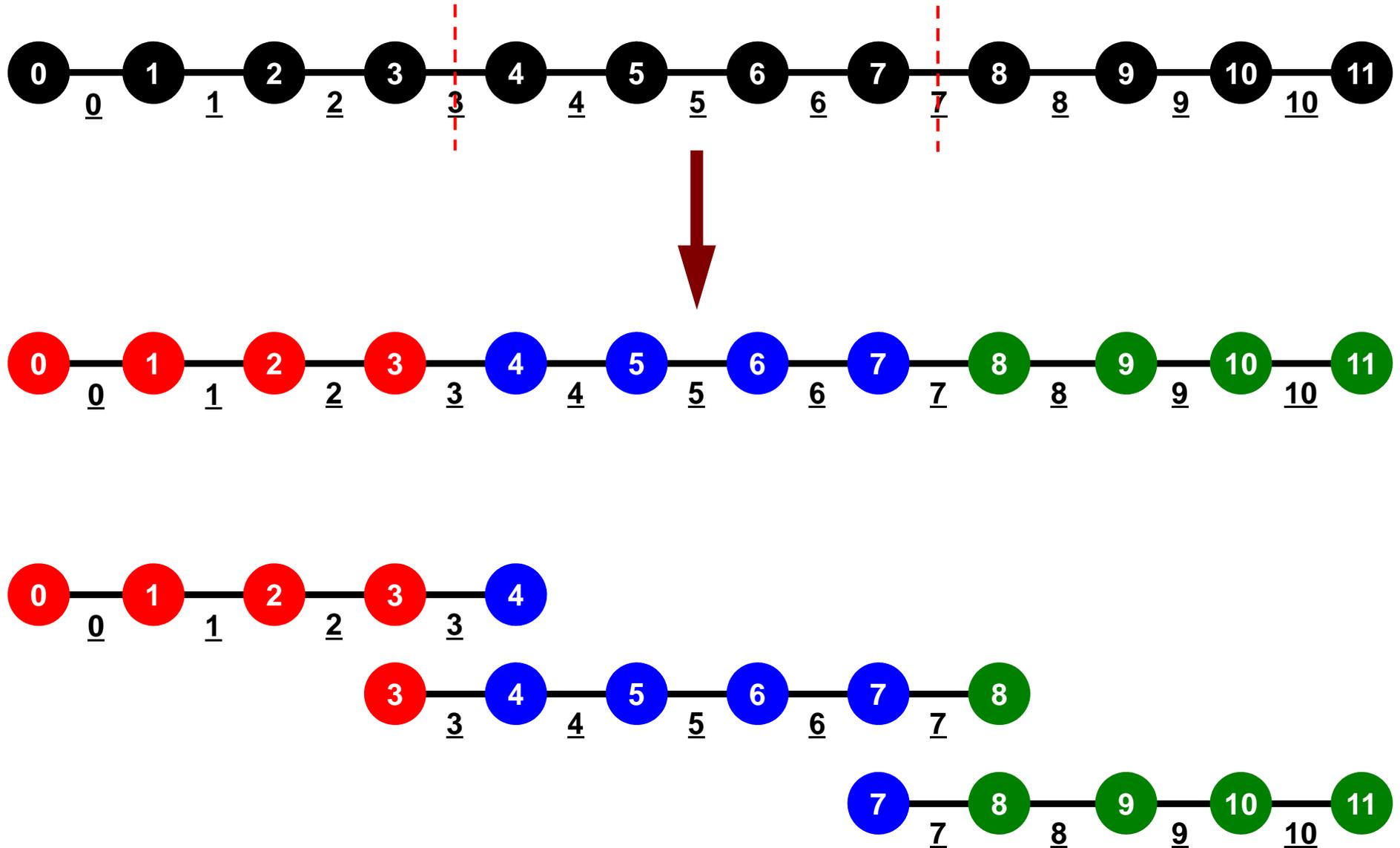
- What is MPI ?
- Your First MPI Program: Hello World
- Collective Communication
- **Point-to-Point Communication**

# Point-to-Point Communication

## 1対1通信

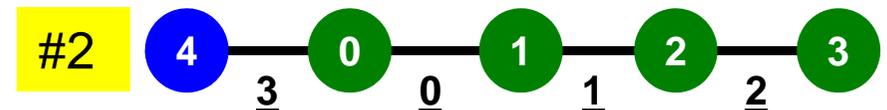
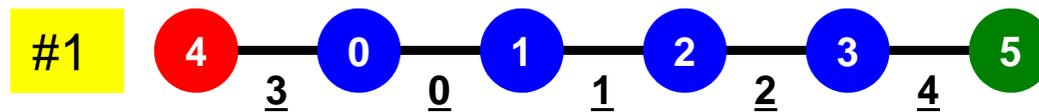
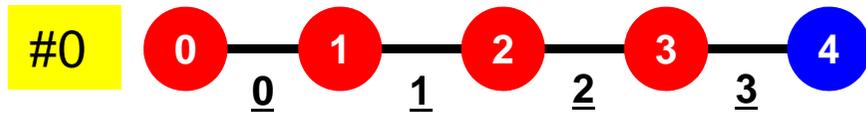
- What is PtoP Communication ?
- 2D Problem, Generalized Communication Table
- Report S2

# 1D FEM: 12 nodes/11 elem's/3 domains



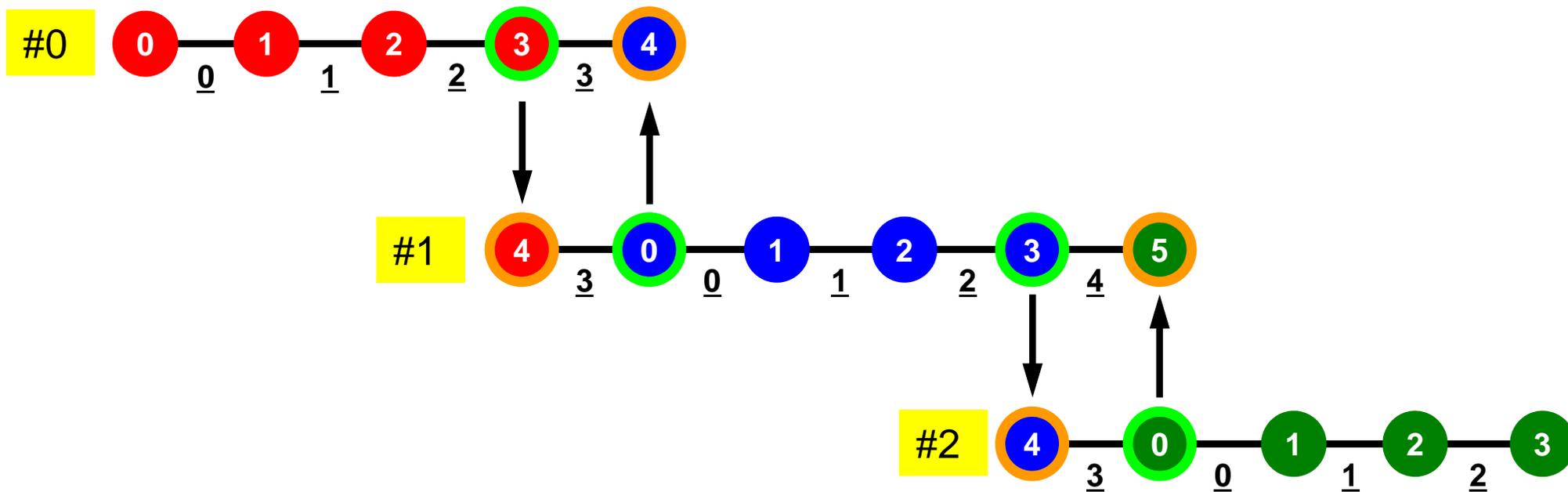
# 1D FEM: 12 nodes/11 elem's/3 domains

Local ID: Starting from 0 for node and elem at each domain



# 1D FEM: 12 nodes/11 elem's/3 domains

Internal/External Nodes



# Preconditioned Conjugate Gradient Method (CG)

```

Compute  $\mathbf{r}^{(0)} = \mathbf{b} - [\mathbf{A}]\mathbf{x}^{(0)}$ 
for  $i = 1, 2, \dots$ 
  solve  $[\mathbf{M}]\mathbf{z}^{(i-1)} = \mathbf{r}^{(i-1)}$ 
   $\rho_{i-1} = \mathbf{r}^{(i-1)} \cdot \mathbf{z}^{(i-1)}$ 
  if  $i=1$ 
     $\mathbf{p}^{(1)} = \mathbf{z}^{(0)}$ 
  else
     $\beta_{i-1} = \rho_{i-1} / \rho_{i-2}$ 
     $\mathbf{p}^{(i)} = \mathbf{z}^{(i-1)} + \beta_{i-1} \mathbf{p}^{(i-1)}$ 
  endif
   $\mathbf{q}^{(i)} = [\mathbf{A}]\mathbf{p}^{(i)}$ 
   $\alpha_i = \rho_{i-1} / \mathbf{p}^{(i)} \cdot \mathbf{q}^{(i)}$ 
   $\mathbf{x}^{(i)} = \mathbf{x}^{(i-1)} + \alpha_i \mathbf{p}^{(i)}$ 
   $\mathbf{r}^{(i)} = \mathbf{r}^{(i-1)} - \alpha_i \mathbf{q}^{(i)}$ 
  check convergence  $|\mathbf{r}|$ 
end

```

Preconditioner:

Diagonal Scaling

Point-Jacobi Preconditioning

$$[\mathbf{M}] = \begin{bmatrix} D_1 & 0 & \dots & 0 & 0 \\ 0 & D_2 & & 0 & 0 \\ \dots & & \dots & & \dots \\ 0 & 0 & & D_{N-1} & 0 \\ 0 & 0 & \dots & 0 & D_N \end{bmatrix}$$

# Preconditioning, DAXPY

Local Operations by Only Internal Points: Parallel Processing is possible

```
/*  
//-- {z} = [Minv]{r}  
*/  
    for (i=0; i<N; i++) {  
        W[Z][i] = W[DD][i] * W[R][i];  
    }
```

```
/*  
//-- {x} = {x} + ALPHA*{p}      DAXPY: double a{x} plus {y}  
// {r} = {r} - ALPHA*{q}  
*/  
    for (i=0; i<N; i++) {  
        PHI[i] += Alpha * W[P][i];  
        W[R][i] -= Alpha * W[Q][i];  
    }
```

0
1
2
3
4
5
6
7
8
9
10
11

# Dot Products

Global Summation needed: Communication ?

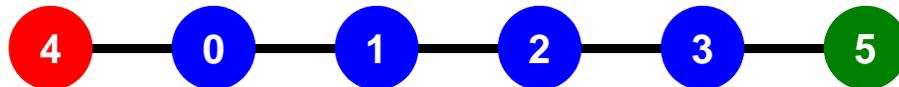
```
/*  
//-- ALPHA= RHO / {p} {q}  
*/  
C1 = 0.0;  
for (i=0; i<N; i++) {  
    C1 += W[P][i] * W[Q][i];  
}  
  
Alpha = Rho / C1;
```

0
1
2
3
4
5
6
7
8
9
10
11

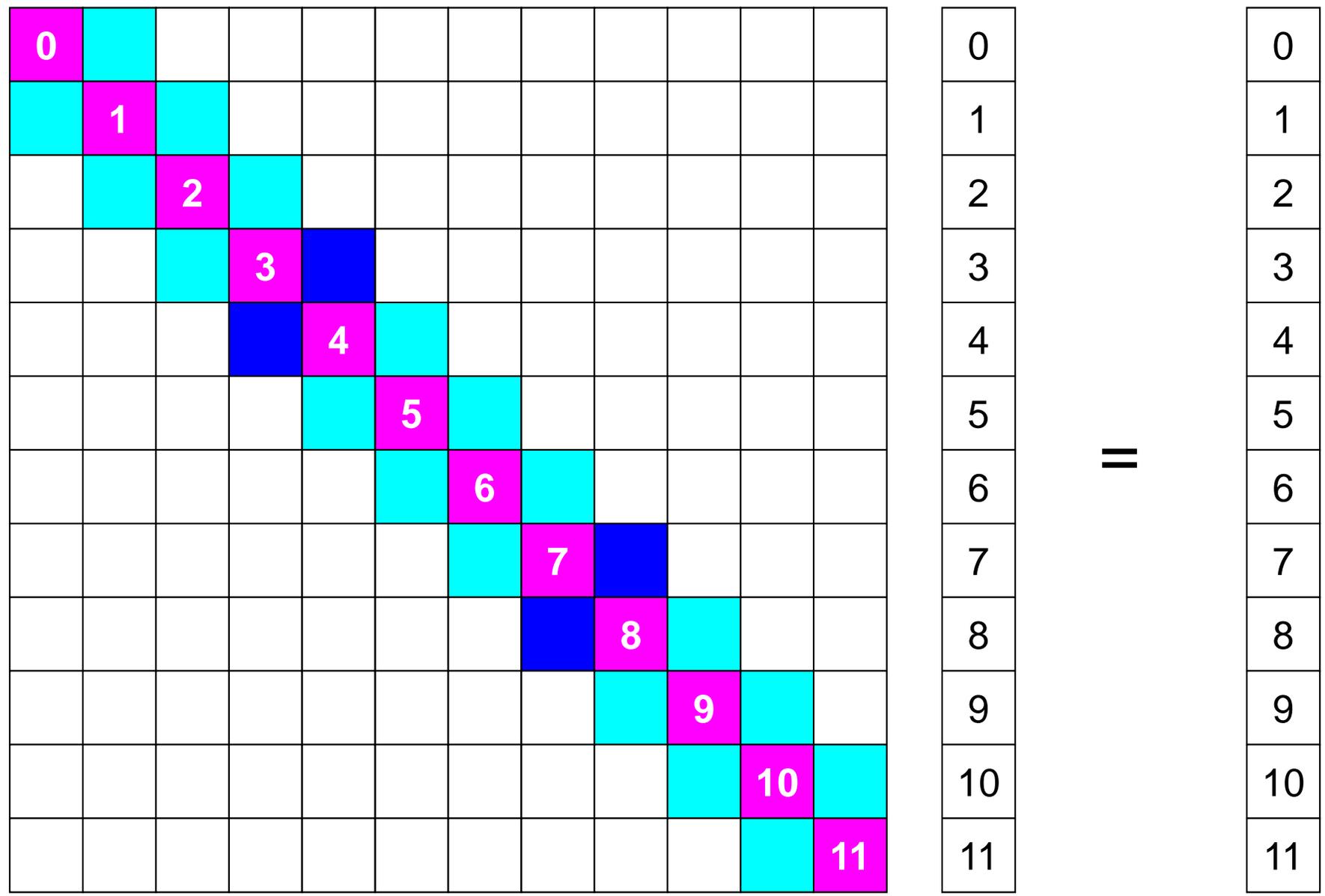
# Matrix-Vector Products

## Values at External Points: P-to-P Communication

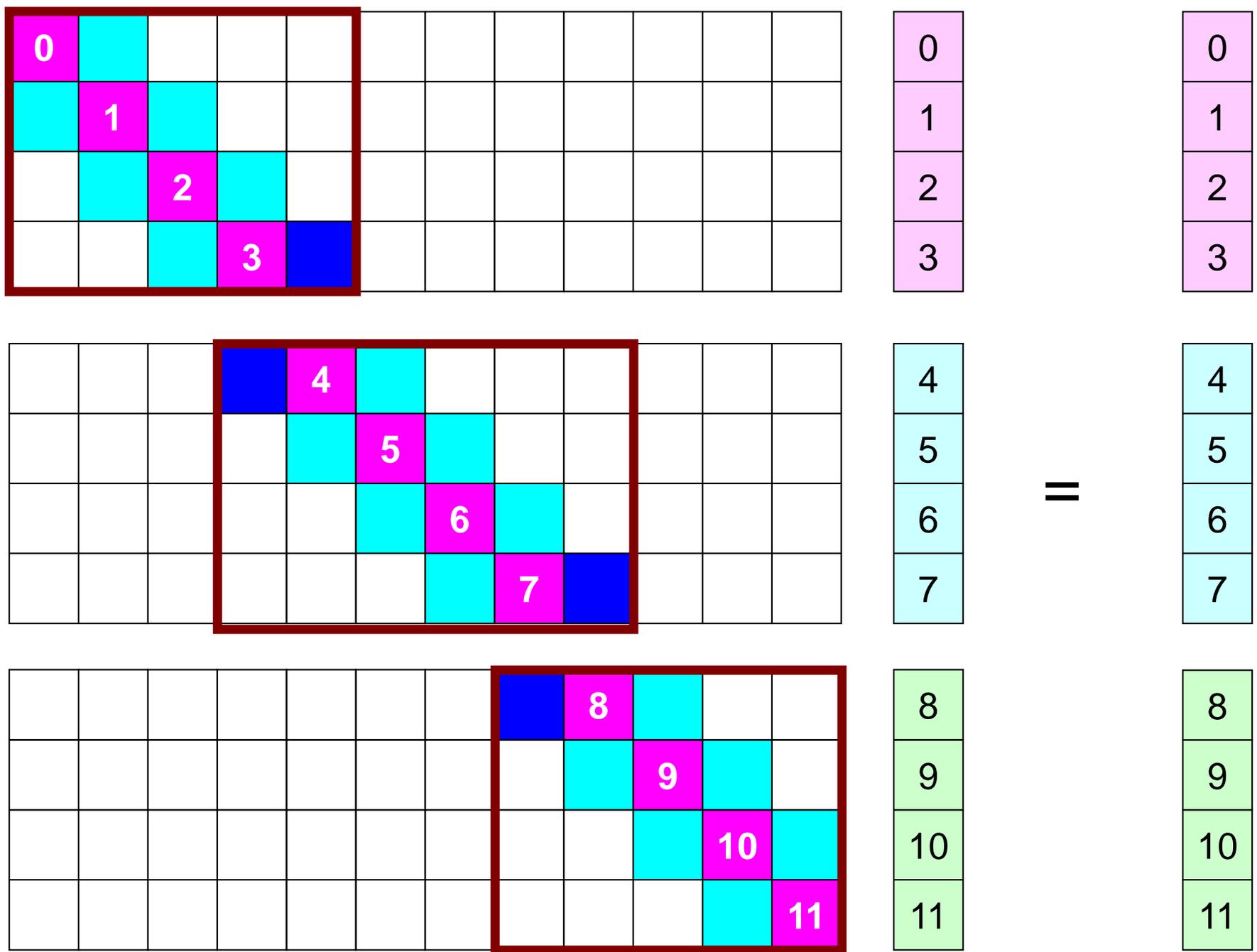
```
/*  
//-- {q} = [A] {p}  
*/  
for (i=0; i<N; i++) {  
    W[Q][i] = Diag[i] * W[P][i];  
    for (j=Index[i]; j<Index[i+1]; j++) {  
        W[Q][i] += AMat[j]*W[P][Item[j]];  
    }  
}
```



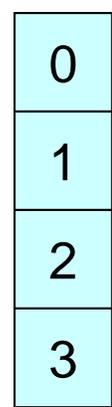
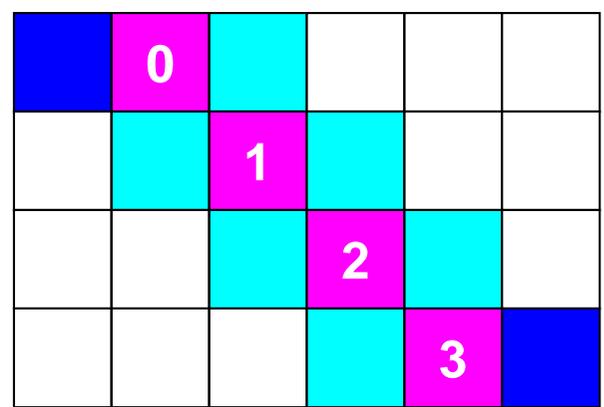
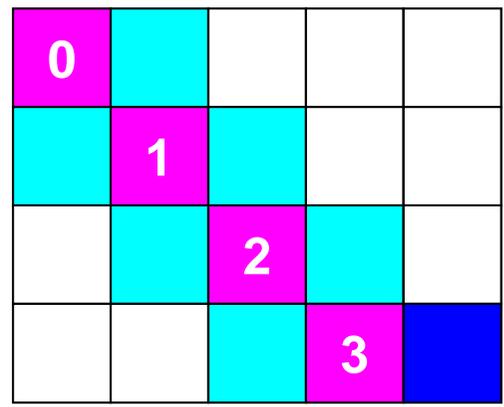
# Mat-Vec Products: Local Op. Possible



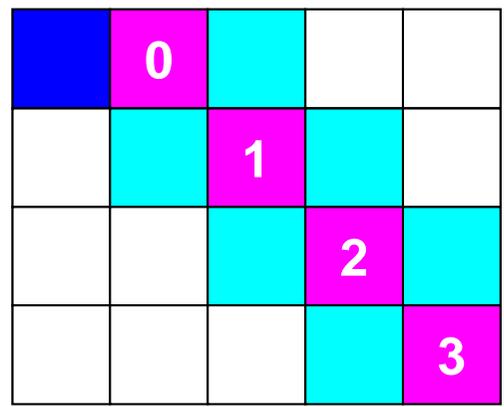
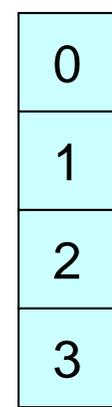
# Mat-Vec Products: Local Op. Possible



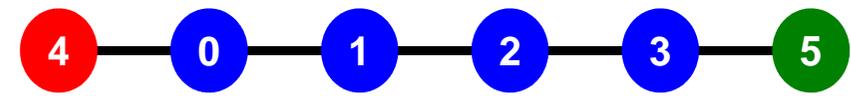
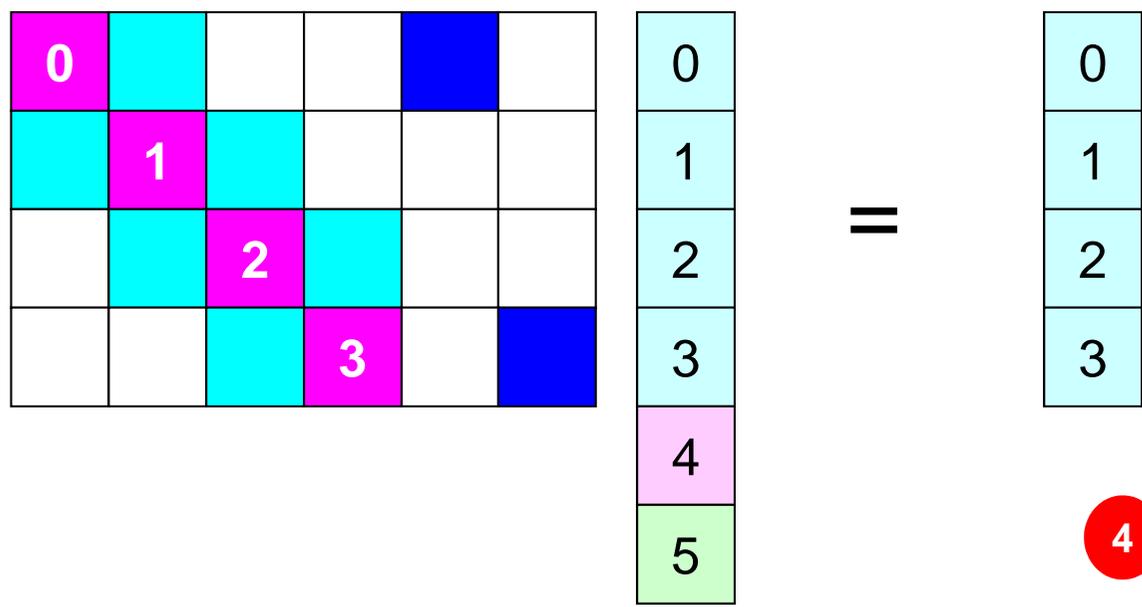
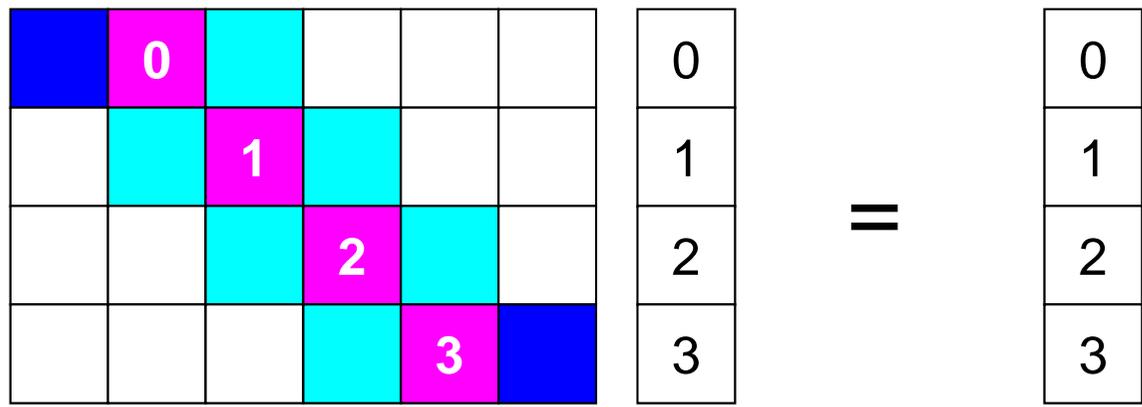
# Mat-Vec Products: Local Op. Possible



=

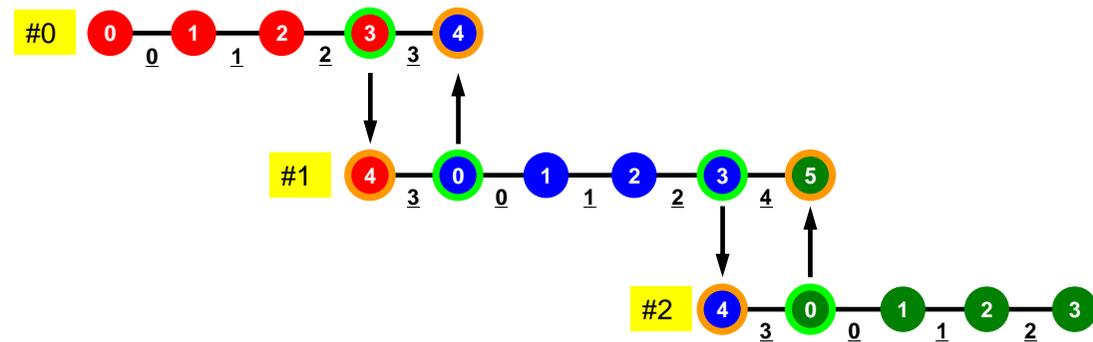


# Mat-Vec Products: Local Op. #1



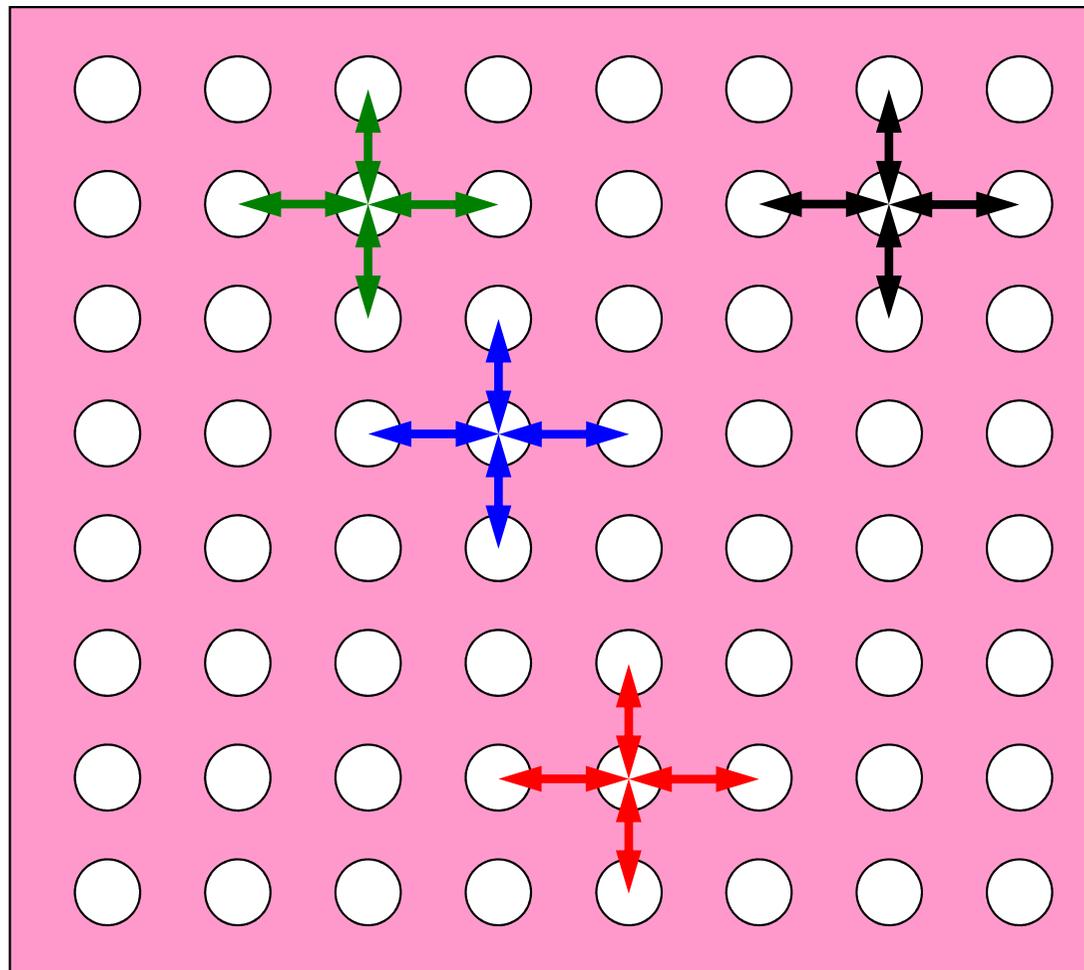
# What is Point-to-Point Comm. ?

- Collective Communication
  - MPI\_Reduce, MPI\_Scatter/Gather etc.
  - Communications with all processes in the communicator
  - Application Area
    - BEM, Spectral Method, MD: global interactions are considered
    - Dot products, MAX/MIN: Global Summation & Comparison
- Point-to-Point
  - MPI\_Send, MPI\_Recv
  - Communication with limited processes
    - Neighbors
  - Application Area
    - FEM, FDM: Localized Method



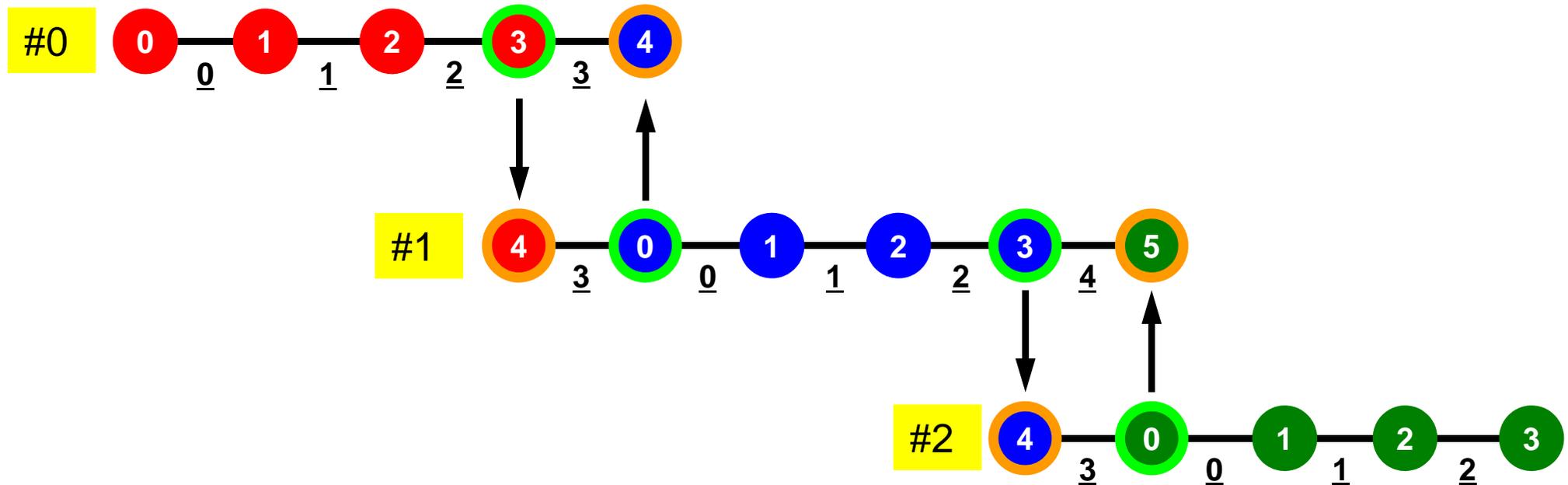
# Collective/PtoP Communications

Interactions with only Neighboring Processes/Element  
Finite Difference Method (FDM), Finite Element Method  
(FEM)



# When do we need PtoP comm.: 1D-FEM

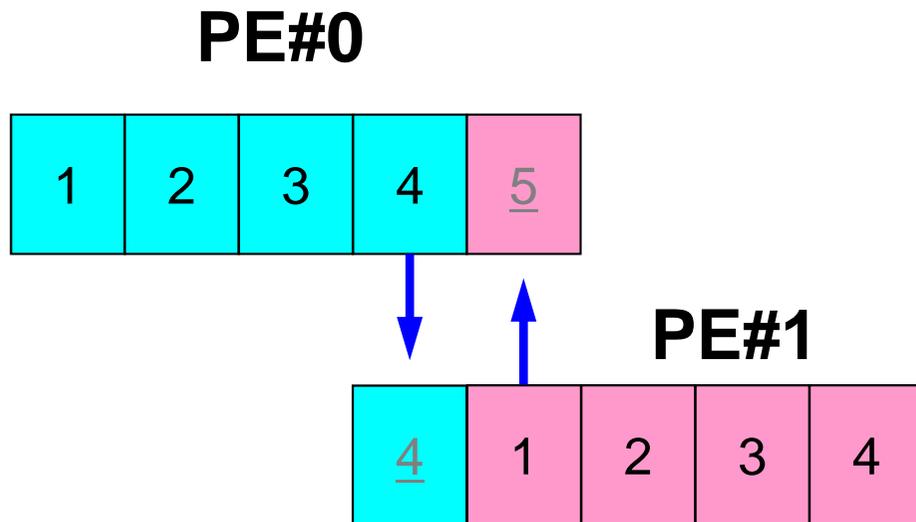
Info in neighboring domains is required for FEM operations  
Matrix assembling, Iterative Method



# Method for PtoP Comm.

- **MPI\_Send, MPI\_Recv**
- These are “blocking” functions. “Dead lock” occurs for these “blocking” functions.
- A “blocking” MPI call means that the program execution will be suspended until the message buffer is safe to use.
- The MPI standards specify that a blocking SEND or RECV does not return until the send buffer is safe to reuse (for MPI\_Send), or the receive buffer is ready to use (for MPI\_Recv).
  - Blocking comm. confirms “secure” communication, but it is very inconvenient.
- Please just remember that “there are such functions”.

# MPI\_Send/MPI\_Recv

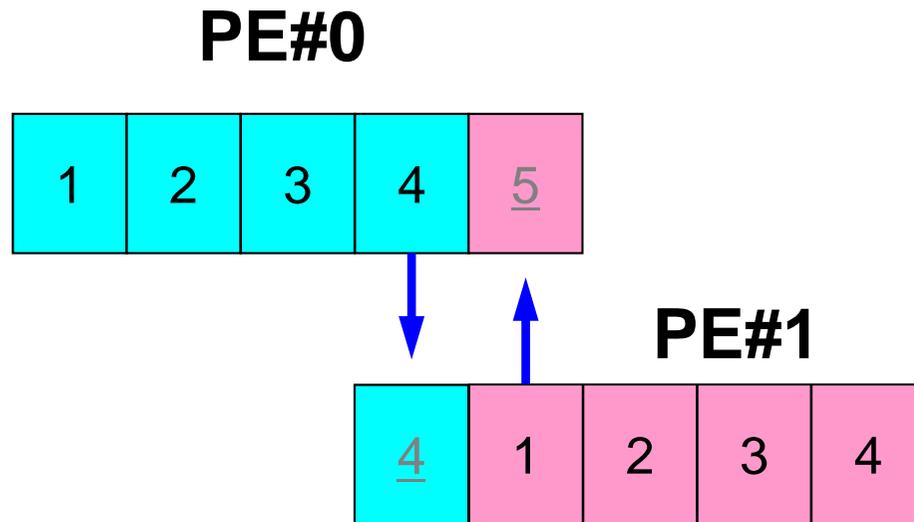


```
if (my_rank.eq.0) NEIB_ID=1
if (my_rank.eq.1) NEIB_ID=0

...
call MPI_SEND (NEIB_ID, arg's)
call MPI_RECV (NEIB_ID, arg's)
...
```

- This seems reasonable, but it stops at MPI\_Send/MPI\_Recv.
  - Sometimes it works (according to implementation).

# MPI\_Send/MPI\_Recv (cont.)



```
if (my_rank.eq.0) NEIB_ID=1
if (my_rank.eq.1) NEIB_ID=0

...
if (my_rank.eq.0) then
  call MPI_SEND (NEIB_ID, arg's)
  call MPI_RECV (NEIB_ID, arg's)
endif

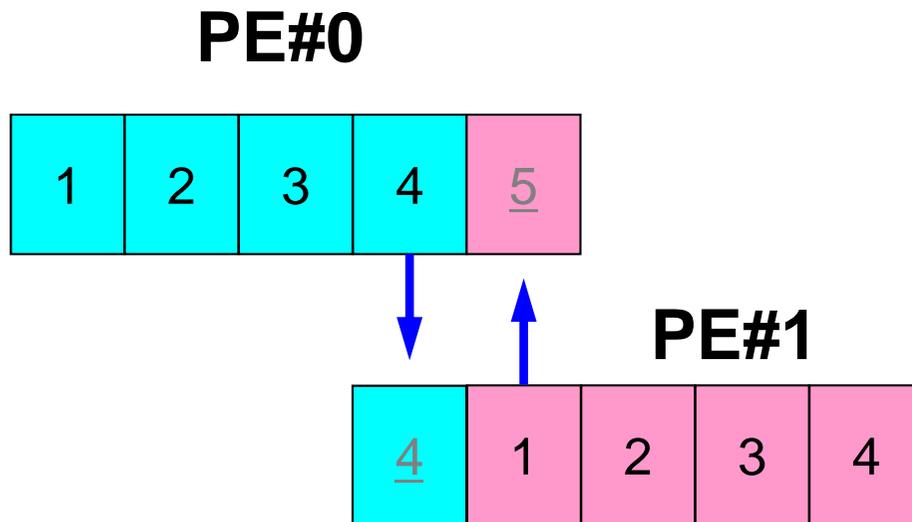
if (my_rank.eq.1) then
  call MPI_RECV (NEIB_ID, arg's)
  call MPI_SEND (NEIB_ID, arg's)
endif

...
```

- It works.
- It is OK for Structured Meshes for FDM (Finite Difference Method).

# How to do PtoP Comm. ?

- Using “non-blocking” functions **MPI\_Isend** & **MPI\_Irecv** together with **MPI\_Waitall** for synchronization
- **MPI\_Sendrecv** is also available.



```

if (my_rank.eq.0) NEIB_ID=1
if (my_rank.eq.1) NEIB_ID=0

...
call MPI_Isend (NEIB_ID, arg's)
call MPI_Irecv (NEIB_ID, arg's)
...
call MPI_Waitall (for Irecv)
...
call MPI_Waitall (for Isend)

```

**MPI\_Waitall** for both of  
**MPI\_Isend/MPI\_Irecv** is possible

# MPI\_Isend

- Begins a non-blocking send
  - Send the contents of sending buffer (starting from **sendbuf**, number of messages: **count**) to **dest** with **tag** .
  - Contents of sending buffer cannot be modified before calling corresponding **MPI\_Waitall**.

- **MPI\_Isend**

**(sendbuf, count, datatype, dest, tag, comm, request)**

- **sendbuf** choice I starting address of sending buffer
- **count** int I number of elements in sending buffer
- **datatype** MPI\_Datatype I datatype of each sending buffer element
- **dest** int I rank of destination
- **tag** int I message tag  
This integer can be used by the application to distinguish messages. Communication occurs if tag's of MPI\_Isend and MPI\_Irecv are matched. Usually tag is set to be "0" (in this class),
- **comm** MPI\_Comm I communicator
- **request** MPI\_Request O **communication request array used in MPI\_Waitall**

# Communication Request: request

## 通信識別子

- **MPI\_Isend**

(**sendbuf**, **count**, **datatype**, **dest**, **tag**, **comm**, **request**)

- **sendbuf** choice I starting address of sending buffer
  - **count** int I number of elements in sending buffer
  - **datatype** MPI\_Datatype I datatype of each sending buffer element
  - **dest** int I rank of destination
  - **tag** int I message tag
- This integer can be used by the application to distinguish messages. Communication occurs if tag's of MPI\_Isend and MPI\_Irecv are matched. Usually tag is set to be "0" (in this class),
- **comm** MPI\_Comm I communicator
  - **request** MPI\_Request O communication request used in MPI\_Waitall

Size of the array is total number of neighboring processes

- Just define the array

# MPI\_Irecv

- Begins a non-blocking receive
  - Receiving the contents of receiving buffer (starting from `recvbuf`, number of messages: `count`) from `source` with `tag` .
  - Contents of receiving buffer cannot be used before calling corresponding `MPI_Waitall`.

- **MPI\_Irecv**

**(recvbuf, count, datatype, source, tag, comm, request)**

- recvbuf choice I starting address of receiving buffer
- count int I number of elements in receiving buffer
- datatype MPI\_Datatype I datatype of each receiving buffer element
- source int I rank of source
- tag int I message tag  
This integer can be used by the application to distinguish messages. Communication occurs if `tag`'s of `MPI_Isend` and `MPI_Irecv` are matched. Usually tag is set to be "0" (in this class),
- comm MPI\_Comm I communicator
- request MPI\_Request O **communication request array used in `MPI_Waitall`**

# MPI\_Waitall

- **MPI\_Waitall** blocks until all comm's, associated with request in the array, complete. It is used for synchronizing MPI\_Isend and MPI\_Irecv in this class.
- At sending phase, contents of sending buffer cannot be modified before calling corresponding **MPI\_Waitall**. At receiving phase, contents of receiving buffer cannot be used before calling corresponding **MPI\_Waitall**.
- MPI\_Isend and MPI\_Irecv can be synchronized simultaneously with a single **MPI\_Waitall** if it is consistent.
  - Same request should be used in MPI\_Isend and MPI\_Irecv.
- Its operation is similar to that of **MPI\_Barrier** but, **MPI\_Waitall** can not be replaced by **MPI\_Barrier**.
  - Possible troubles using **MPI\_Barrier** instead of **MPI\_Waitall**: Contents of **request** and **status** are not updated properly, very slow operations etc.
- **MPI\_Waitall** (count, request, status)
  - count     int        I        number of processes to be synchronized
  - request   MPI\_Request I/O     comm. request used in MPI\_Isend/Irecv (array size: count)
  - status    MPI\_Status  O        **array of status objects**  
                                   MPI\_STATUS\_SIZE: defined in 'mpif.h', 'mpi.h'

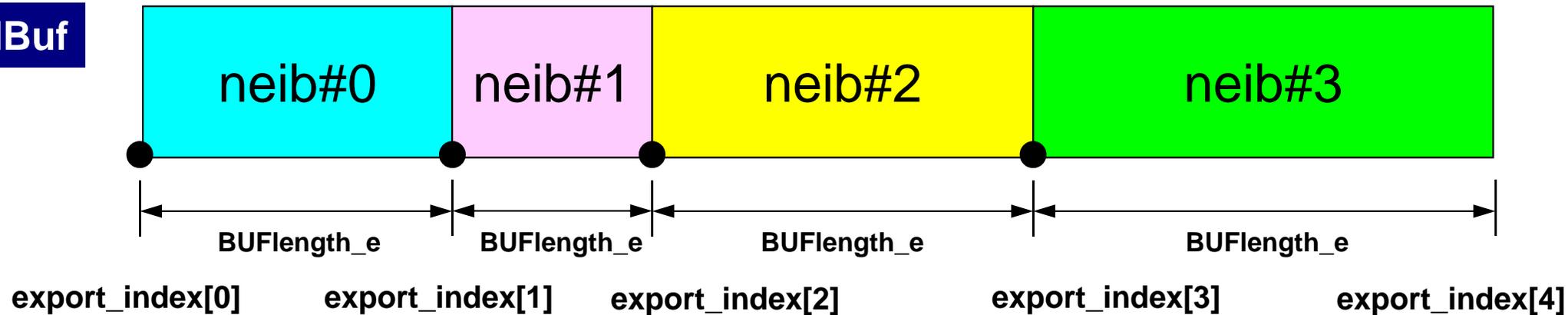
# Array of status object: status 状況オブジェクト配列

- **MPI\_Waitall (count, request, status)**
  - **count** int I number of processes to be synchronized
  - **request** MPI\_Request I/O comm. request used in MPI\_Isend/Irecv (array size: count)
  - **status** MPI\_Status O **array of status objects**  
MPI\_STATUS\_SIZE: defined in 'mpif.h', 'mpi.h'
- Just define the array

# SEND: MPI\_Isend/Irecv/Waitall

C

SendBuf



```

for (neib=0; neib<NeibPETot; neib++){
    tag= 0;
    iS_e= export_index[neib];
    iE_e= export_index[neib+1];
    BUFlength_e= iE_e - iS_e

    ierr= MPI_Isend
        (&SendBuf[iS_e], BUFlength_e, MPI_DOUBLE,
         NeibPE[neib], 0,
         MPI_COMM_WORLD, &ReqSend[neib])
}

MPI_Waitall(NeibPETot, ReqSend, StatSend);

```

# RECV: MPI\_Irecv/Irecv/Waitall

C

```

for (neib=0; neib<NeibPETot; neib++){
  tag= 0;
  iS_i= import_index[neib];
  iE_i= import_index[neib+1];
  BUFlength_i= iE_i - iS_i

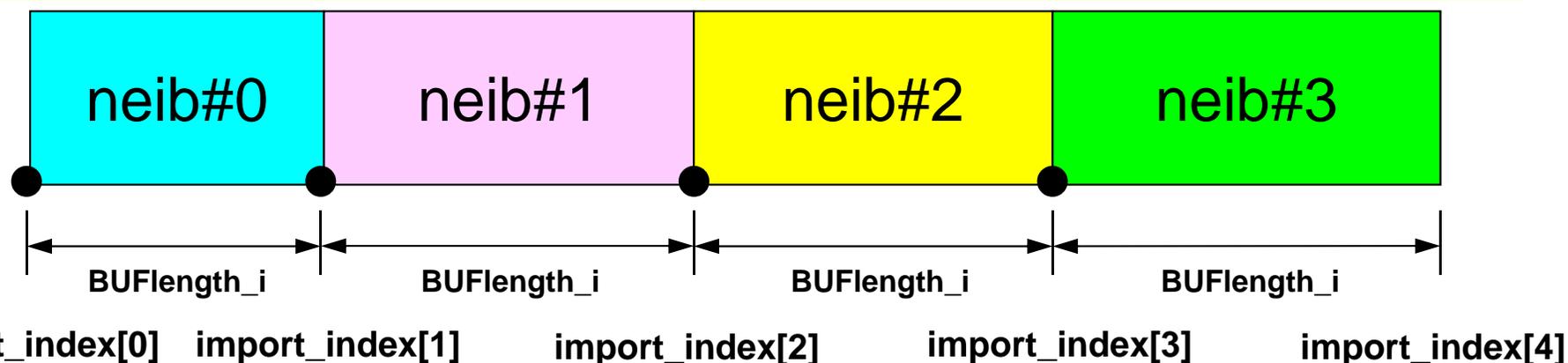
  ierr= MPI_Irecv
        (&RecvBuf[iS_i], BUFlength_i, MPI_DOUBLE,
         NeibPE[neib], 0,
         MPI_COMM_WORLD, &ReqRecv[neib])
}

```

```
MPI_Waitall(NeibPETot, ReqRecv, StatRecv);
```

import\_item (import\_index[neib]:import\_index[neib+1]-1) are received from neib-th neighbor

RecvBuf



# MPI\_Sendrecv

- MPI\_Send+MPI\_Recv: not recommended, many restrictions

- **MPI\_Sendrecv**

(**sendbuf, sendcount, sendtype, dest, sendtag, recvbuf, recvcount, recvtype, source, recvtag, comm, status**)

–	<b><u>sendbuf</u></b>	choice	I	starting address of sending buffer
–	<b><u>sendcount</u></b>	int	I	number of elements in sending buffer
–	<b><u>sendtype</u></b>	MPI_Datatype	I	datatype of each sending buffer element
–	<b><u>dest</u></b>	int	I	rank of destination
–	<b><u>sendtag</u></b>	int	I	message tag for sending
–	<b><u>comm</u></b>	MPI_Comm	I	communicator
–	<b><u>recvbuf</u></b>	choice	I	starting address of receiving buffer
–	<b><u>recvcount</u></b>	int	I	number of elements in receiving buffer
–	<b><u>recvtype</u></b>	MPI_Datatype	I	datatype of each receiving buffer element
–	<b><u>source</u></b>	int	I	rank of source
–	<b><u>recvtag</u></b>	int	I	message tag for receiving
–	<b><u>comm</u></b>	MPI_Comm	I	communicator
–	<b><u>status</u></b>	MPI_Status	O	array of status objects

MPI\_STATUS\_SIZE: defined in 'mpif.h', 'mpi.h'

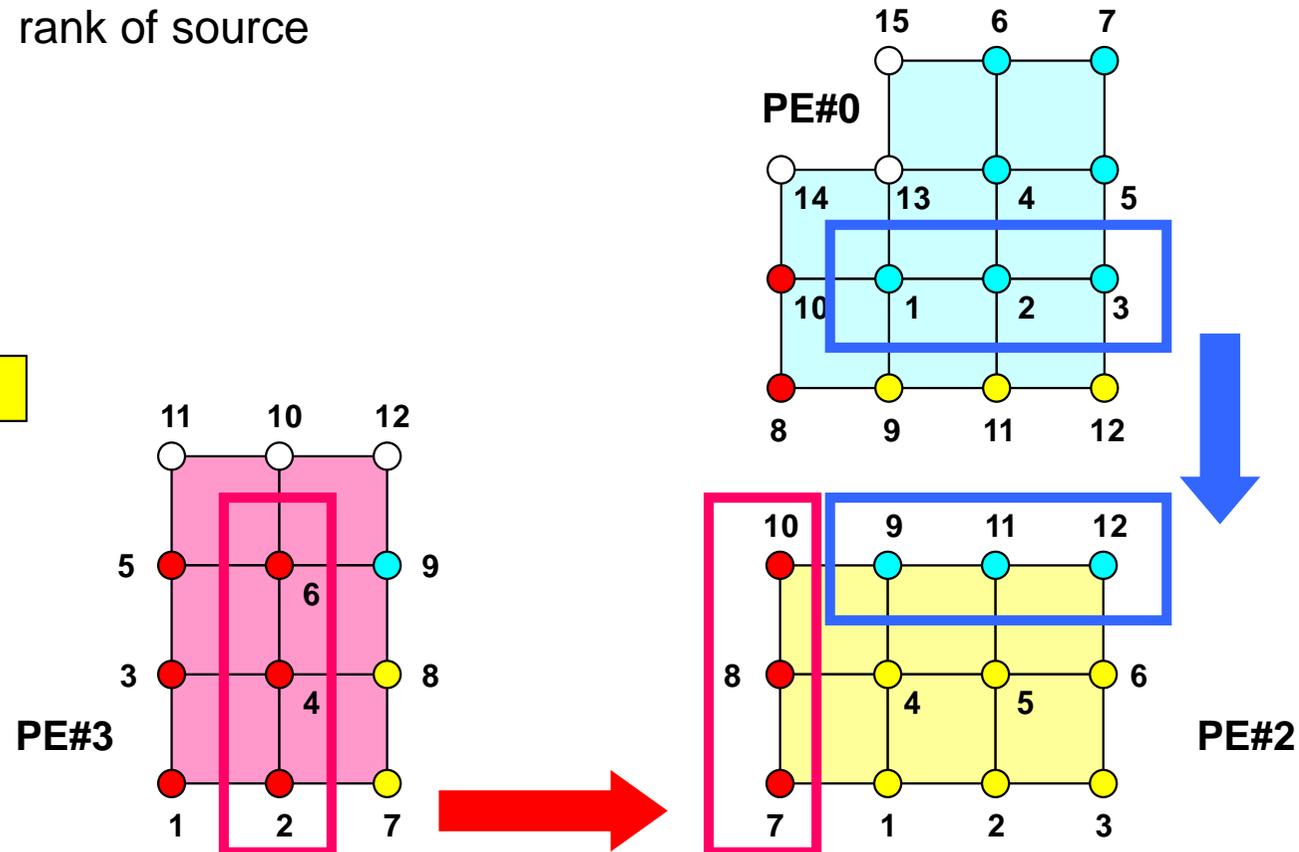
# RECV: receiving to external nodes

Recv. continuous data to recv. buffer from neighbors

- MPI\_Irecv**

(**recvbuf**, **count**, **datatype**, **source**, **tag**, **comm**, **request**)

<b>recvbuf</b>	choice	I	starting address of receiving buffer
<b>count</b>	int	I	number of elements in receiving buffer
<b>datatype</b>	MPI_Datatype	I	datatype of each receiving buffer element
<b>source</b>	int	I	rank of source



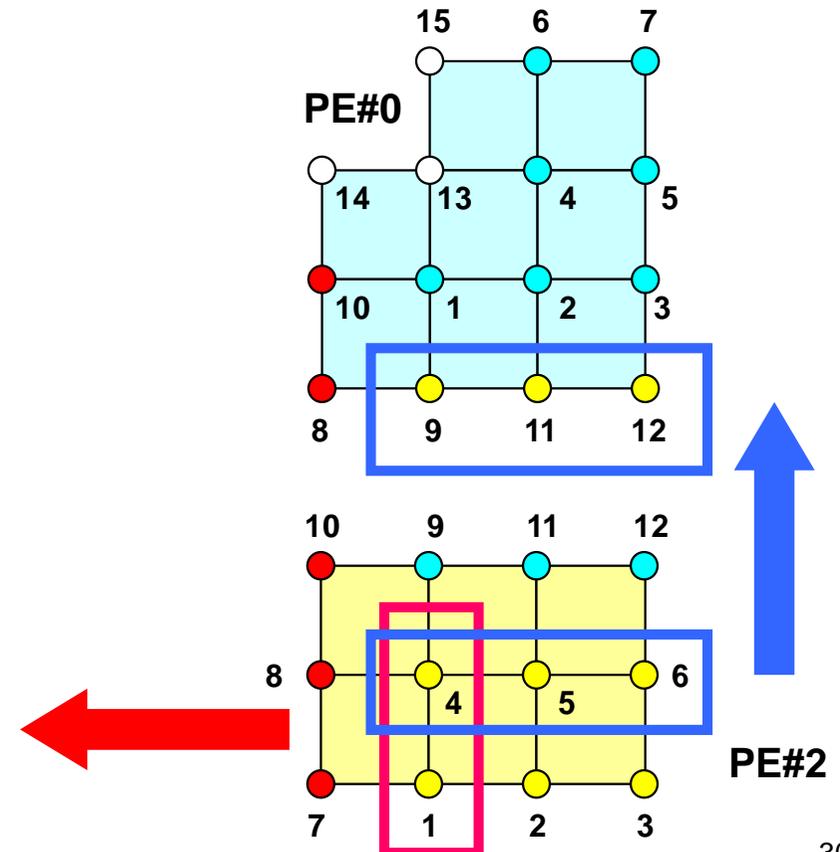
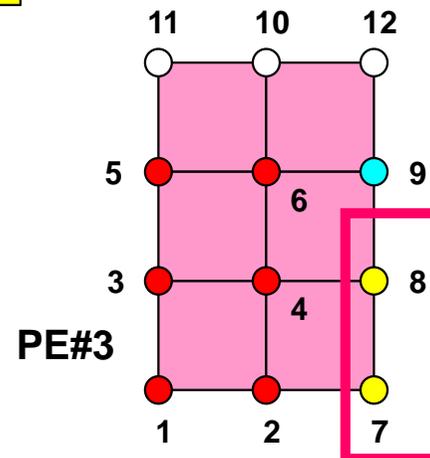
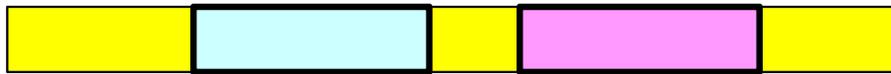
# SEND: sending from boundary nodes

## Send continuous data to send buffer of neighbors

- MPI\_Isend**

(**sendbuf**, **count**, **datatype**, **dest**, **tag**, **comm**, **request**)

<b>sendbuf</b>	choice	I	starting address of sending buffer
<b>count</b>	int	I	number of elements in sending buffer
<b>datatype</b>	MPI_Datatype	I	datatype of each sending buffer element
<b>dest</b>	int	I	rank of destination



# Request, Status in C Language

## Special TYPE of Arrays

- `MPI_Isend`: request
- `MPI_Irecv`: request
- `MPI_Waitall`: request, status

```
MPI_Status *StatSend, *StatRecv;
MPI_Request *RequestSend, *RequestRecv;
...
StatSend = malloc(sizeof(MPI_Status) * NEIBpetot);
StatRecv = malloc(sizeof(MPI_Status) * NEIBpetot);
RequestSend = malloc(sizeof(MPI_Request) * NEIBpetot);
RequestRecv = malloc(sizeof(MPI_Request) * NEIBpetot);
```

- `MPI_Sendrecv`: status

```
MPI_Status *Status;
...
Status = malloc(sizeof(MPI_Status));
```

# Files on OBCX

## Fotran

```
>$ cd /work/gt37/t37XXX/pFEM
>$ cp /work/gt00/z30088/pFEM/F/s2-f.tar .
>$ tar xvf s2-f.tar
```

## C

```
>$ cd /work/gt37/t37XXX/pFEM
>$ cp /work/gt00/z30088/pFEM/C/s2-c.tar .
>$ tar xvf s2-c.tar
```

## Confirm Directory

```
>$ ls
  mpi

>$ cd mpi/S2
```

This directory is called as  $\langle \$O-S2 \rangle$  in this course.

$\langle \$O-S2 \rangle = \langle \$O-TOP \rangle / \text{mpi} / S2$

# Ex.1: Send-Recv a Scalar

- Exchange VAL (real, 8-byte) between PE#0 & PE#1

```

if (my_rank=0) NEIB= 1;
if (my_rank=1) NEIB= 0;

MPI_Isend (&VAL      , 1, MPI_DOUBLE, NEIB, ..., req_send[0], ...);
MPI_Irecv (&VALtemp, 1, MPI_DOUBLE, NEIB, ..., req_recv[0], ...);
MPI_Waitall (... , req_recv, stat_recv, ...); Recv.buf VALtemp can be used
MPI_Waitall (... , req_send, stat_send, ...); Send buf VAL can be modified
VAL= VALtemp;

```

```

if (my_rank.eq.0) NEIB= 1
if (my_rank.eq.1) NEIB= 0

call MPI_Sendrecv (VAL      , 1, MPI_DOUBLE_PRECISION, NEIB, ...           &
                  VALtemp, 1, MPI_DOUBLE_PRECISION, NEIB, ...,  status, ...)
VAL= VALtemp

```

Name of recv. buffer could be "VAL", but not recommended.

# Ex.1: Send-Recv a Scalar

## Isend/Irecv/Waitall

```
$> cd /work/gt37/t37XXX/pFEM/mpi/S2
$> mpiicc -O3 ex1-1.c
$> pjsub go2.sh
```

```
#include <stdio.h>
#include <stdlib.h>
#include "mpi.h"
int main(int argc, char **argv){
    int neib, MyRank, PeTot;
    double VAL, VALx;
    MPI_Status *StatSend, *StatRecv;
    MPI_Request *RequestSend, *RequestRecv;

    MPI_Init(&argc, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &PeTot);
    MPI_Comm_rank(MPI_COMM_WORLD, &MyRank);
    StatSend = malloc(sizeof(MPI_Status) * 1);
    StatRecv = malloc(sizeof(MPI_Status) * 1);
    RequestSend = malloc(sizeof(MPI_Request) * 1);
    RequestRecv = malloc(sizeof(MPI_Request) * 1);

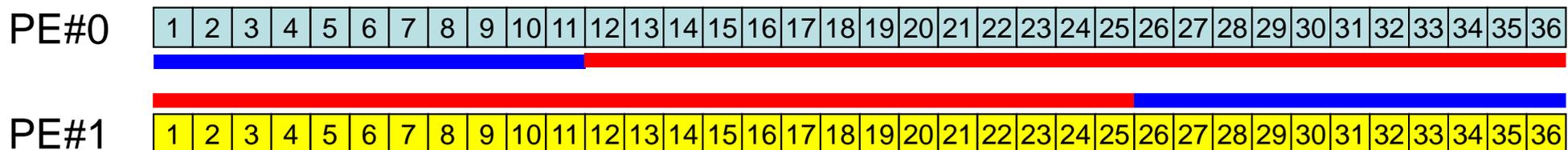
    if(MyRank == 0) {neib= 1; VAL= 10.0;}
    if(MyRank == 1) {neib= 0; VAL= 11.0;}

    MPI_Isend(&VAL , 1, MPI_DOUBLE, neib, 0, MPI_COMM_WORLD, &RequestSend[0]);
    MPI_Irecv(&VALx, 1, MPI_DOUBLE, neib, 0, MPI_COMM_WORLD, &RequestRecv[0]);
    MPI_Waitall(1, RequestRecv, StatRecv);
    MPI_Waitall(1, RequestSend, StatSend);

    VAL=VALx;
    MPI_Finalize();
    return 0; }
```

# Ex.2: Send-Recv an Array (1/3)

- Exchange VEC (real, 8-byte) between PE#0 & PE#1
- PE#0 to PE#1
  - PE#0: send VEC(1)-VEC(11) (length=11)
  - PE#1: recv. as VEC(26)-VEC(36) (length=11)
- PE#1 to PE#0
  - PE#1: send VEC(1)-VEC(25) (length=25)
  - PE#0: recv. as VEC(12)-VEC(36) (length=25)
- Practice: Develop a program for this operation.



# Practice: t1

- Initial status of VEC[:]:
  - PE#0 VEC[0-35]= 101,102,103,~,135,136
  - PE#1 VEC[0-35]= 201,202,203,~,235,236
- Confirm the results in the next page
- MPI\_Isend/Irecv/Waitall

# Estimated Results

**t1**

```

0 #BEFORE# 1      101.
0 #BEFORE# 2      102.
0 #BEFORE# 3      103.
0 #BEFORE# 4      104.
0 #BEFORE# 5      105.
0 #BEFORE# 6      106.
0 #BEFORE# 7      107.
0 #BEFORE# 8      108.
0 #BEFORE# 9      109.
0 #BEFORE# 10     110.
0 #BEFORE# 11     111.
0 #BEFORE# 12     112.
0 #BEFORE# 13     113.
0 #BEFORE# 14     114.
0 #BEFORE# 15     115.
0 #BEFORE# 16     116.
0 #BEFORE# 17     117.
0 #BEFORE# 18     118.
0 #BEFORE# 19     119.
0 #BEFORE# 20     120.
0 #BEFORE# 21     121.
0 #BEFORE# 22     122.
0 #BEFORE# 23     123.
0 #BEFORE# 24     124.
0 #BEFORE# 25     125.
0 #BEFORE# 26     126.
0 #BEFORE# 27     127.
0 #BEFORE# 28     128.
0 #BEFORE# 29     129.
0 #BEFORE# 30     130.
0 #BEFORE# 31     131.
0 #BEFORE# 32     132.
0 #BEFORE# 33     133.
0 #BEFORE# 34     134.
0 #BEFORE# 35     135.
0 #BEFORE# 36     136.

```

```

0 #AFTER # 1      101.
0 #AFTER # 2      102.
0 #AFTER # 3      103.
0 #AFTER # 4      104.
0 #AFTER # 5      105.
0 #AFTER # 6      106.
0 #AFTER # 7      107.
0 #AFTER # 8      108.
0 #AFTER # 9      109.
0 #AFTER # 10     110.
0 #AFTER # 11     111.
0 #AFTER # 12     201.
0 #AFTER # 13     202.
0 #AFTER # 14     203.
0 #AFTER # 15     204.
0 #AFTER # 16     205.
0 #AFTER # 17     206.
0 #AFTER # 18     207.
0 #AFTER # 19     208.
0 #AFTER # 20     209.
0 #AFTER # 21     210.
0 #AFTER # 22     211.
0 #AFTER # 23     212.
0 #AFTER # 24     213.
0 #AFTER # 25     214.
0 #AFTER # 26     215.
0 #AFTER # 27     216.
0 #AFTER # 28     217.
0 #AFTER # 29     218.
0 #AFTER # 30     219.
0 #AFTER # 31     220.
0 #AFTER # 32     221.
0 #AFTER # 33     222.
0 #AFTER # 34     223.
0 #AFTER # 35     224.
0 #AFTER # 36     225.

```

```

1 #BEFORE# 1      201.
1 #BEFORE# 2      202.
1 #BEFORE# 3      203.
1 #BEFORE# 4      204.
1 #BEFORE# 5      205.
1 #BEFORE# 6      206.
1 #BEFORE# 7      207.
1 #BEFORE# 8      208.
1 #BEFORE# 9      209.
1 #BEFORE# 10     210.
1 #BEFORE# 11     211.
1 #BEFORE# 12     212.
1 #BEFORE# 13     213.
1 #BEFORE# 14     214.
1 #BEFORE# 15     215.
1 #BEFORE# 16     216.
1 #BEFORE# 17     217.
1 #BEFORE# 18     218.
1 #BEFORE# 19     219.
1 #BEFORE# 20     220.
1 #BEFORE# 21     221.
1 #BEFORE# 22     222.
1 #BEFORE# 23     223.
1 #BEFORE# 24     224.
1 #BEFORE# 25     225.
1 #BEFORE# 26     226.
1 #BEFORE# 27     227.
1 #BEFORE# 28     228.
1 #BEFORE# 29     229.
1 #BEFORE# 30     230.
1 #BEFORE# 31     231.
1 #BEFORE# 32     232.
1 #BEFORE# 33     233.
1 #BEFORE# 34     234.
1 #BEFORE# 35     235.
1 #BEFORE# 36     236.

```

```

1 #AFTER # 1      201.
1 #AFTER # 2      202.
1 #AFTER # 3      203.
1 #AFTER # 4      204.
1 #AFTER # 5      205.
1 #AFTER # 6      206.
1 #AFTER # 7      207.
1 #AFTER # 8      208.
1 #AFTER # 9      209.
1 #AFTER # 10     210.
1 #AFTER # 11     211.
1 #AFTER # 12     212.
1 #AFTER # 13     213.
1 #AFTER # 14     214.
1 #AFTER # 15     215.
1 #AFTER # 16     216.
1 #AFTER # 17     217.
1 #AFTER # 18     218.
1 #AFTER # 19     219.
1 #AFTER # 20     220.
1 #AFTER # 21     221.
1 #AFTER # 22     222.
1 #AFTER # 23     223.
1 #AFTER # 24     224.
1 #AFTER # 25     225.
1 #AFTER # 26     101.
1 #AFTER # 27     102.
1 #AFTER # 28     103.
1 #AFTER # 29     104.
1 #AFTER # 30     105.
1 #AFTER # 31     106.
1 #AFTER # 32     107.
1 #AFTER # 33     108.
1 #AFTER # 34     109.
1 #AFTER # 35     110.
1 #AFTER # 36     111.

```

## Ex.2: Send-Recv an Array (2/3)

```
if (my_rank=0) then
  MPI_Isend (&VEC[ 0], 11, MPI_DOUBLE, 1, ..., req_send[0], ...)
  MPI_Irecv (&VEC[11], 25, MPI_DOUBLE, 1, ..., req_recv[0], ...)
endif

if (my_rank=1) then
  MPI_Isend (&VEC[ 0], 25, 0, ..., req_send[0], ...)
  MPI_Irecv (&VEC[25], 11, 0, ..., req_recv[0], ...)
endif

MPI_Waitall (... , req_recv, stat_recv, ...)
MPI_Waitall (... , req_send, stat_send, ...)
```

It works, but complicated operations.

Not looks like SPMD.

Not portable.

# Ex.2: Send-Recv an Array (3/3)

```
if (my_rank=0) then
  NEIB= 1
  start_send= 1
  length_send= 11
  start_recv= length_send + 1
  length_recv= 25
endif

if (my_rank=1) then
  NEIB= 0
  start_send= 1
  length_send= 25
  start_recv= length_send + 1
  length_recv= 11
endif

MPI_Isend
(&VEC[start_send-1], length_send, MPI_DOUBLE, NEIB, ..., req_send[0], ...)
MPI_Irecv
(&VEC[start_recv-1], length_recv, MPI_DOUBLE, NEIB, ..., req_recv[0], ...)

MPI_Waitall (... , req_recv, stat_recv, ...)
MPI_Waitall (... , req_send, stat_send, ...)
```

This is “SMPD” !!

t1

# Notice: Send/Recv Arrays

```
#PE0  
send:  
  VEC (start_send) ~  
  VEC (start_send+length_send-1)
```

```
#PE1  
send:  
  VEC (start_send) ~  
  VEC (start_send+length_send-1)
```

```
#PE0  
recv:  
  VEC (start_recv) ~  
  VEC (start_recv+length_recv-1)
```

```
#PE1  
recv:  
  VEC (start_recv) ~  
  VEC (start_recv+length_recv-1)
```

- “length\_send” of sending process must be equal to “length\_recv” of receiving process.
  - PE#0 to PE#1, PE#1 to PE#0
- “sendbuf” and “recvbuf”: different address

# Point-to-Point Communication

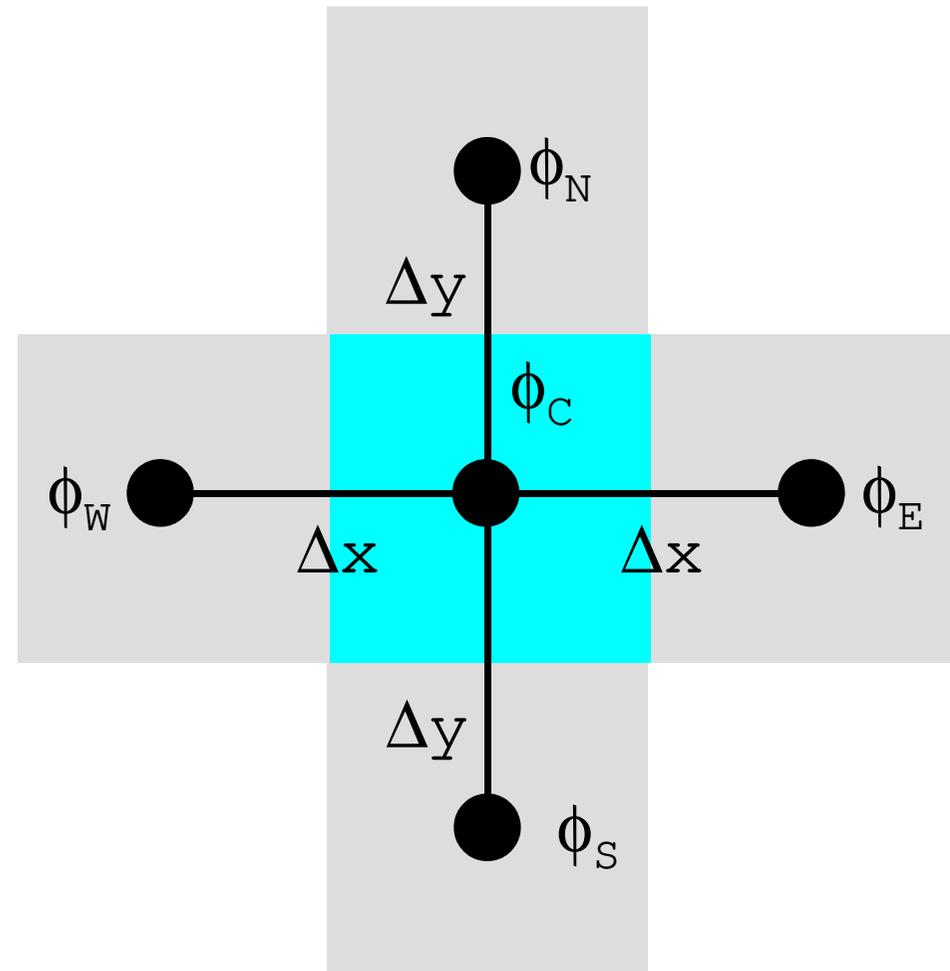
- What is PtoP Communication ?
- **2D Problem, Generalized Communication Table**
  - 2D FDM
  - Problem Setting
  - Distributed Local Data and Communication Table
  - Implementation
- Report S2



# 2D FDM (5-point, central difference)

$$\frac{\partial^2 \phi}{\partial x^2} + \frac{\partial^2 \phi}{\partial y^2} = f$$

$$\left( \frac{\phi_E - 2\phi_C + \phi_W}{\Delta x^2} \right) + \left( \frac{\phi_N - 2\phi_C + \phi_S}{\Delta y^2} \right) = f_C$$



# Decompose into 4 domains

<u>57</u>	<u>58</u>	<u>59</u>	<u>60</u>	<u>61</u>	<u>62</u>	<u>63</u>	<u>64</u>
<u>49</u>	<u>50</u>	<u>51</u>	<u>52</u>	<u>53</u>	<u>54</u>	<u>55</u>	<u>56</u>
<u>41</u>	<u>42</u>	<u>43</u>	<u>44</u>	<u>45</u>	<u>46</u>	<u>47</u>	<u>48</u>
<u>33</u>	<u>34</u>	<u>35</u>	<u>36</u>	<u>37</u>	<u>38</u>	<u>39</u>	<u>40</u>
<u>25</u>	<u>26</u>	<u>27</u>	<u>28</u>	<u>29</u>	<u>30</u>	<u>31</u>	<u>32</u>
<u>17</u>	<u>18</u>	<u>19</u>	<u>20</u>	<u>21</u>	<u>22</u>	<u>23</u>	<u>24</u>
<u>9</u>	<u>10</u>	<u>11</u>	<u>12</u>	<u>13</u>	<u>14</u>	<u>15</u>	<u>16</u>
<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>	<u>5</u>	<u>6</u>	<u>7</u>	<u>8</u>

# 4 domains: Global ID

PE#2

<u>57</u>	<u>58</u>	<u>59</u>	<u>60</u>
<u>49</u>	<u>50</u>	<u>51</u>	<u>52</u>
<u>41</u>	<u>42</u>	<u>43</u>	<u>44</u>
<u>33</u>	<u>34</u>	<u>35</u>	<u>36</u>

PE#3

<u>61</u>	<u>62</u>	<u>63</u>	<u>64</u>
<u>53</u>	<u>54</u>	<u>55</u>	<u>56</u>
<u>45</u>	<u>46</u>	<u>47</u>	<u>48</u>
<u>37</u>	<u>38</u>	<u>39</u>	<u>40</u>

PE#0

<u>25</u>	<u>26</u>	<u>27</u>	<u>28</u>
<u>17</u>	<u>18</u>	<u>19</u>	<u>20</u>
<u>9</u>	<u>10</u>	<u>11</u>	<u>12</u>
<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>

PE#1

<u>29</u>	<u>30</u>	<u>31</u>	<u>32</u>
<u>21</u>	<u>22</u>	<u>23</u>	<u>24</u>
<u>13</u>	<u>14</u>	<u>15</u>	<u>16</u>
<u>5</u>	<u>6</u>	<u>7</u>	<u>8</u>

# 4 domains: Local ID

PE#2

13	14	15	16
9	10	11	12
5	6	7	8
1	2	3	4

PE#3

13	14	15	16
9	10	11	12
5	6	7	8
1	2	3	4

PE#0

13	14	15	16
9	10	11	12
5	6	7	8
1	2	3	4

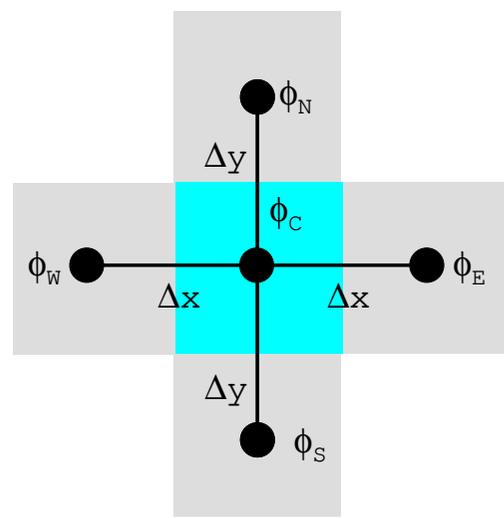
PE#1

13	14	15	16
9	10	11	12
5	6	7	8
1	2	3	4

# External Points: Overlapped Region

PE#2

PE#3



13	14	15	16	13	14	15	16
9	10	11	12	9	10	11	12
5	6	7	8	5	6	7	8
1	2	3	4	1	2	3	4
13	14	15	16	13	14	15	16
9	10	11	12	9	10	11	12
5	6	7	8	5	6	7	8
1	2	3	4	1	2	3	4

PE#0

PE#1

# External Points: Overlapped Region

PE#2

13	14	15	16
9	10	11	12
5	6	7	8
1	2	3	4

PE#3

13	14	15	16
9	10	11	12
5	6	7	8
1	2	3	4

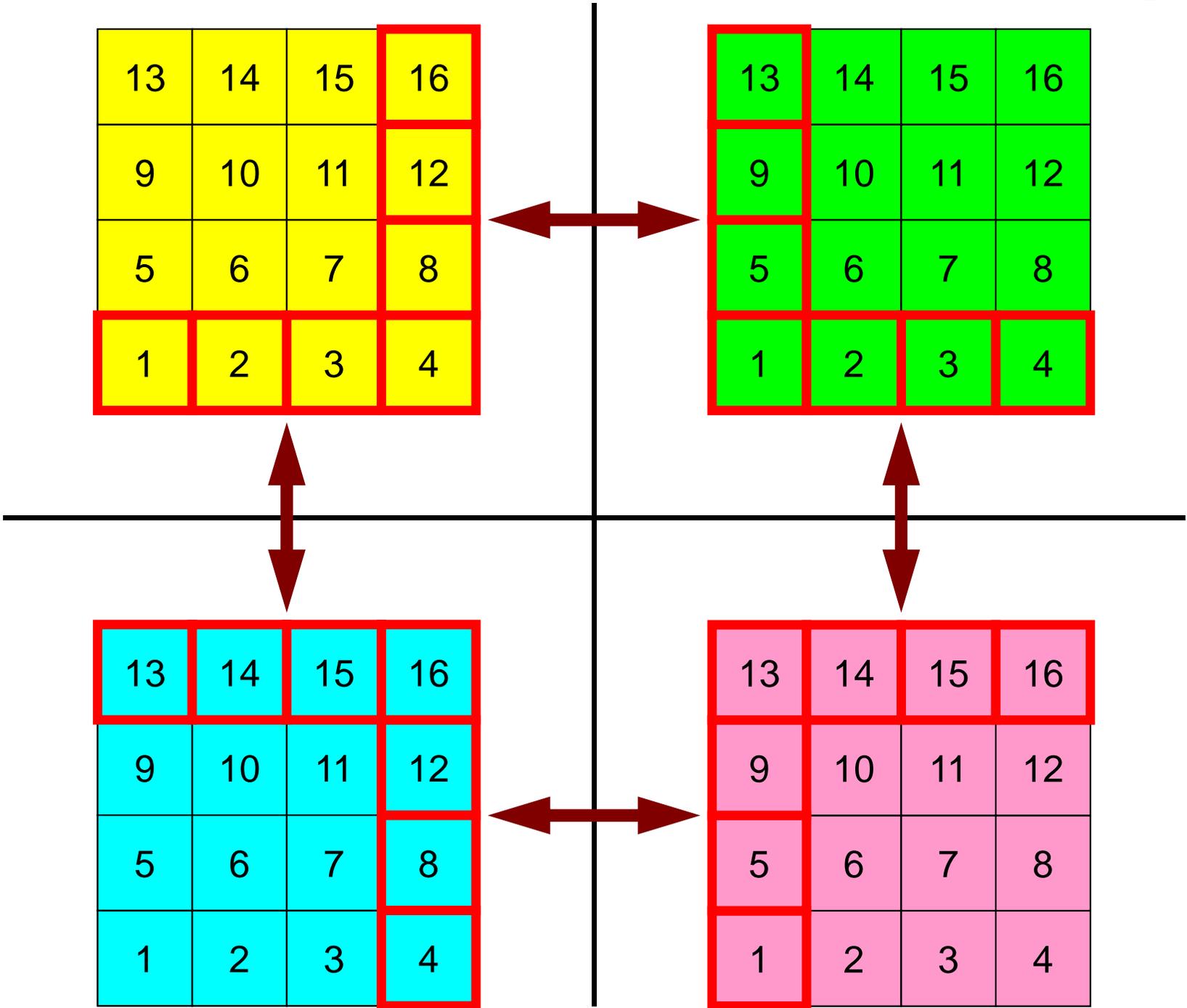


PE#0

13	14	15	16
9	10	11	12
5	6	7	8
1	2	3	4

PE#1

13	14	15	16
9	10	11	12
5	6	7	8
1	2	3	4



# Local ID of External Points ?

PE#2

13	14	15	16	?
9	10	11	12	?
5	6	7	8	?
1	2	3	4	?
?	?	?	?	

PE#3

?	13	14	15	16
?	9	10	11	12
?	5	6	7	8
?	1	2	3	4
	?	?	?	?

PE#0

?	?	?	?	
13	14	15	16	?
9	10	11	12	?
5	6	7	8	?
1	2	3	4	?

PE#1

	?	?	?	?
?	13	14	15	16
?	9	10	11	12
?	5	6	7	8
?	1	2	3	4

# Overlapped Region

PE#2

13	14	15	16	?
9	10	11	12	?
5	6	7	8	?
1	2	3	4	?

PE#3

?	13	14	15	16
?	9	10	11	12
?	5	6	7	8
?	1	2	3	4

?	?	?	?
---	---	---	---

?	?	?	?
---	---	---	---

?	?	?	?
---	---	---	---

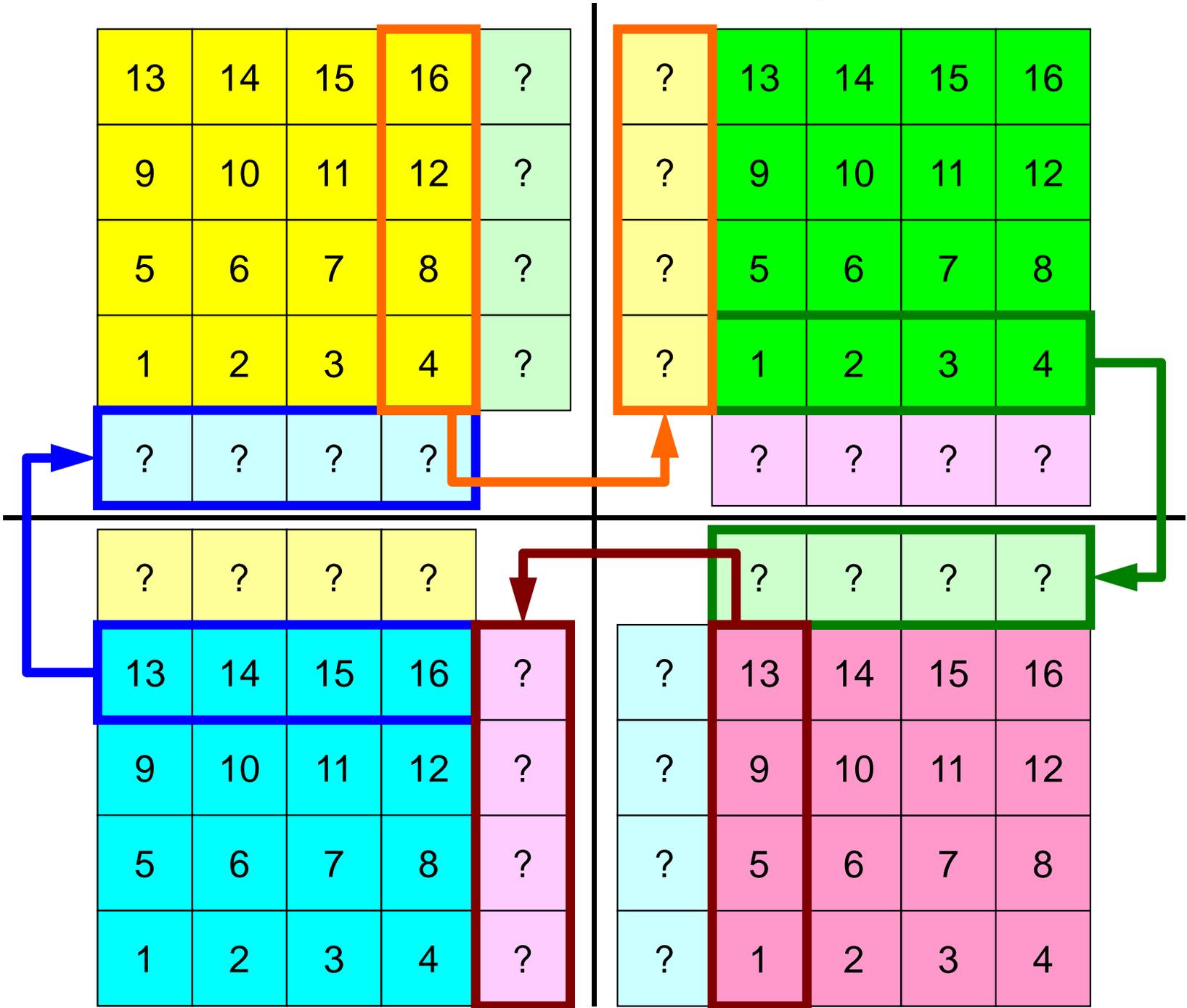
?	?	?	?
---	---	---	---

13	14	15	16	?
9	10	11	12	?
5	6	7	8	?
1	2	3	4	?

?	13	14	15	16
?	9	10	11	12
?	5	6	7	8
?	1	2	3	4

PE#0

PE#1



# Overlapped Region

PE#2

13	14	15	16	?
9	10	11	12	?
5	6	7	8	?
1	2	3	4	?
?	?	?	?	

PE#3

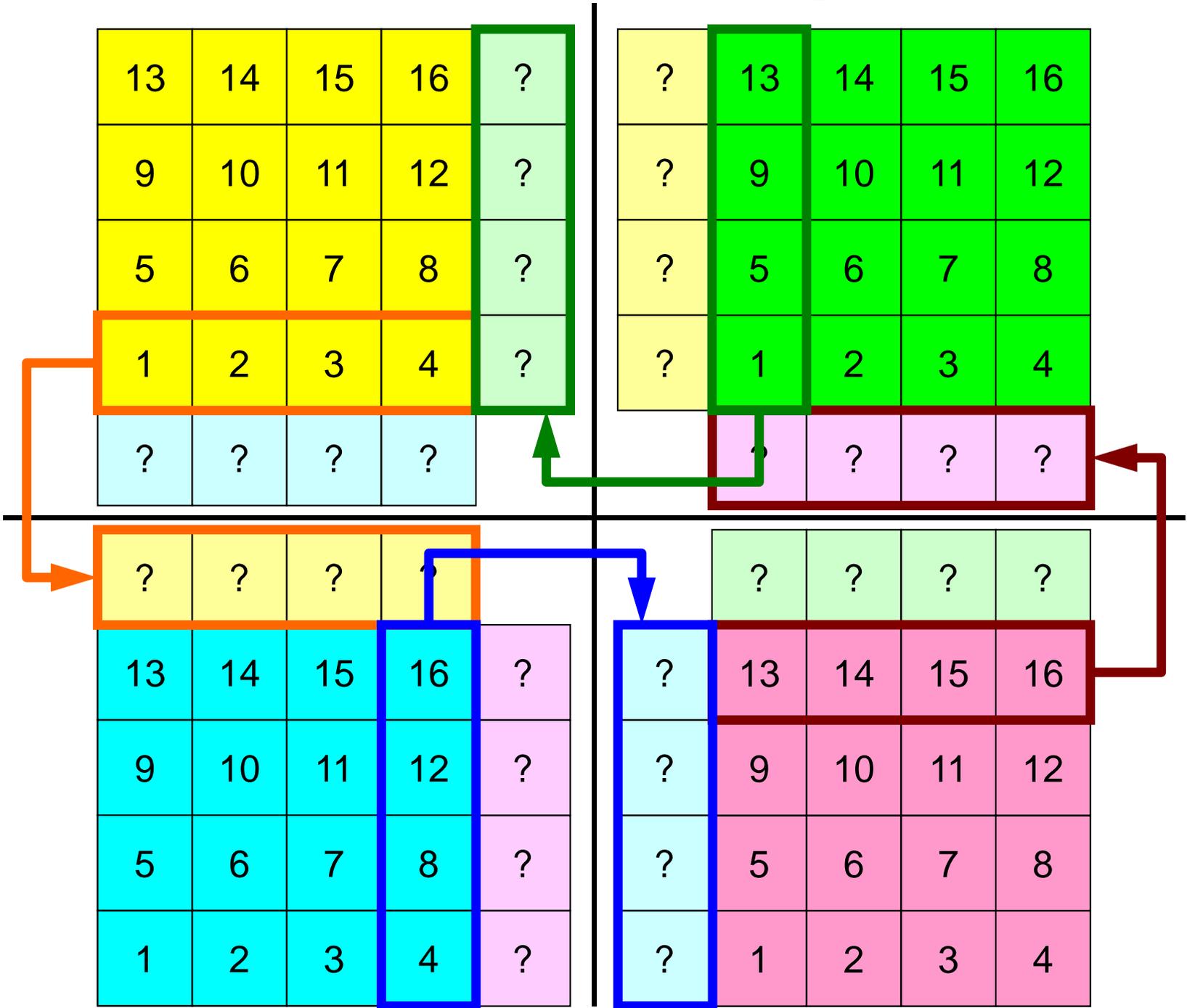
?	13	14	15	16
?	9	10	11	12
?	5	6	7	8
?	1	2	3	4
	?	?	?	?

PE#0

?	?	?	?	
13	14	15	16	?
9	10	11	12	?
5	6	7	8	?
1	2	3	4	?

PE#1

	?	?	?	?
?	13	14	15	16
?	9	10	11	12
?	5	6	7	8
?	1	2	3	4



# Point-to-Point Communication

- What is PtoP Communication ?
- 2D Problem, Generalized Communication Table
  - 2D FDM
  - Problem Setting
  - Distributed Local Data and Communication Table
  - Implementation
- Report S2

# Problem Setting: 2D FDM

<u>57</u>	<u>58</u>	<u>59</u>	<u>60</u>	<u>61</u>	<u>62</u>	<u>63</u>	<u>64</u>
<u>49</u>	<u>50</u>	<u>51</u>	<u>52</u>	<u>53</u>	<u>54</u>	<u>55</u>	<u>56</u>
<u>41</u>	<u>42</u>	<u>43</u>	<u>44</u>	<u>45</u>	<u>46</u>	<u>47</u>	<u>48</u>
<u>33</u>	<u>34</u>	<u>35</u>	<u>36</u>	<u>37</u>	<u>38</u>	<u>39</u>	<u>40</u>
<u>25</u>	<u>26</u>	<u>27</u>	<u>28</u>	<u>29</u>	<u>30</u>	<u>31</u>	<u>32</u>
<u>17</u>	<u>18</u>	<u>19</u>	<u>20</u>	<u>21</u>	<u>22</u>	<u>23</u>	<u>24</u>
<u>9</u>	<u>10</u>	<u>11</u>	<u>12</u>	<u>13</u>	<u>14</u>	<u>15</u>	<u>16</u>
<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>	<u>5</u>	<u>6</u>	<u>7</u>	<u>8</u>

- 2D region with 64 meshes (8x8)
- Each mesh has global ID from 1 to 64
  - In this example, this global ID is considered as dependent variable, such as temperature, pressure etc.
  - Something like computed results

# Problem Setting: Distributed Local Data

**PE#2**

57	58	59	60
49	50	51	52
41	42	43	44
33	34	35	36

**PE#3**

61	62	63	64
53	54	55	56
45	46	47	48
37	38	39	40

- 4 sub-domains.
- Info. of external points (global ID of mesh) is received from neighbors.
  - PE#0 receives

25	26	27	28
17	18	19	20
9	10	11	12
1	2	3	4

29	30	31	32
21	22	23	24
13	14	15	16
5	6	7	8

**PE#2**

57	58	59	60	
49	50	51	52	
41	42	43	44	
33	34	35	36	

**PE#3**

	61	62	63	64
	53	54	55	56
	45	46	47	48
	37	38	39	40

25	26	27	28	
17	18	19	20	
9	10	11	12	
1	2	3	4	

	29	30	31	32
	21	22	23	24
	13	14	15	16
	5	6	7	8

**PE#0**

**PE#1**

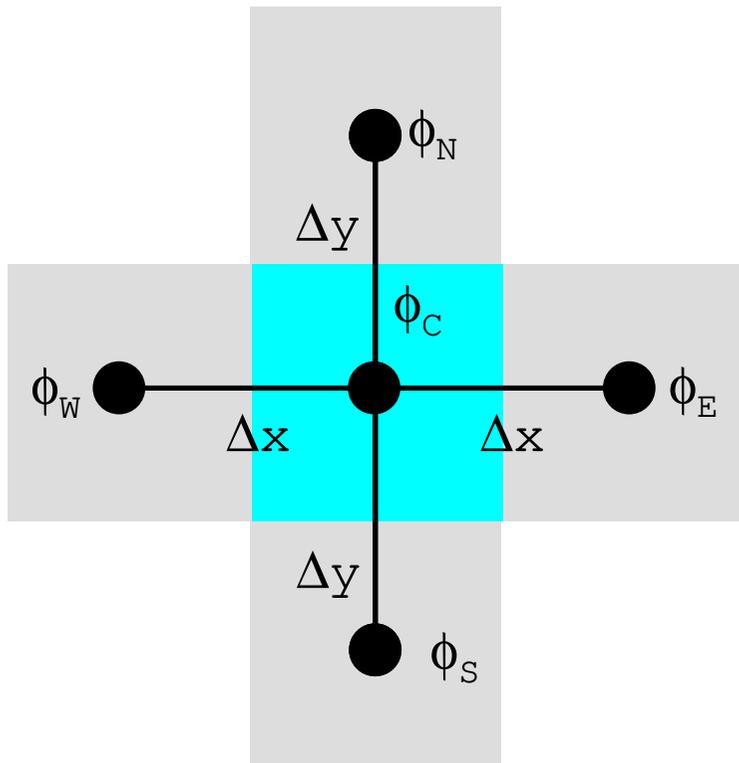
**PE#0**

**PE#1**

# Operations of 2D FDM

$$\frac{\partial^2 \phi}{\partial x^2} + \frac{\partial^2 \phi}{\partial y^2} = f$$

$$\left( \frac{\phi_E - 2\phi_C + \phi_W}{\Delta x^2} \right) + \left( \frac{\phi_N - 2\phi_C + \phi_S}{\Delta y^2} \right) = f_C$$

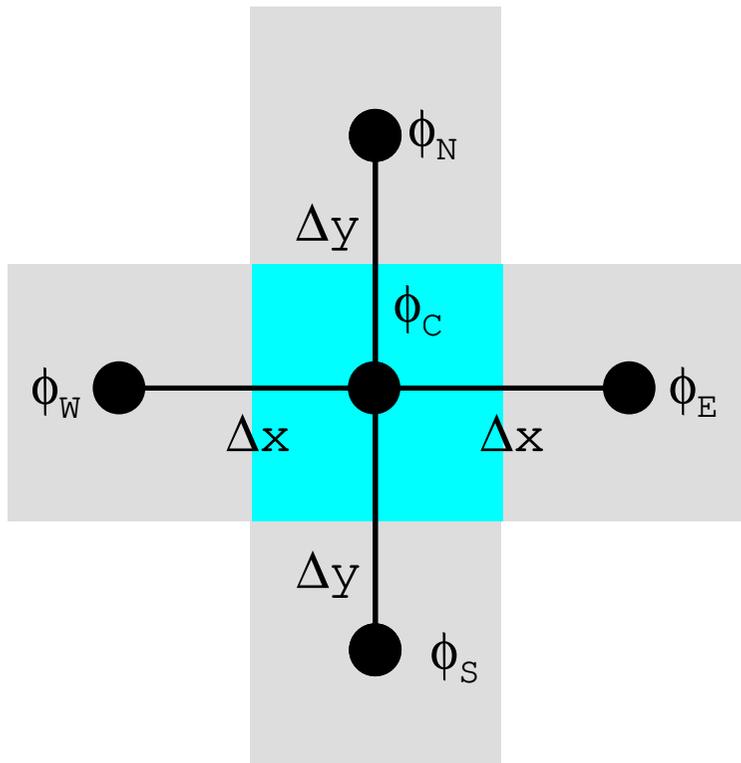


<u>57</u>	<u>58</u>	<u>59</u>	<u>60</u>	<u>61</u>	<u>62</u>	<u>63</u>	<u>64</u>
<u>49</u>	<u>50</u>	<u>51</u>	<u>52</u>	<u>53</u>	<u>54</u>	<u>55</u>	<u>56</u>
<u>41</u>	<u>42</u>	<u>43</u>	<u>44</u>	<u>45</u>	<u>46</u>	<u>47</u>	<u>48</u>
<u>33</u>	<u>34</u>	<u>35</u>	<u>36</u>	<u>37</u>	<u>38</u>	<u>39</u>	<u>40</u>
<u>25</u>	<u>26</u>	<u>27</u>	<u>28</u>	<u>29</u>	<u>30</u>	<u>31</u>	<u>32</u>
<u>17</u>	<u>18</u>	<u>19</u>	<u>20</u>	<u>21</u>	<u>22</u>	<u>23</u>	<u>24</u>
<u>9</u>	<u>10</u>	<u>11</u>	<u>12</u>	<u>13</u>	<u>14</u>	<u>15</u>	<u>16</u>
<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>	<u>5</u>	<u>6</u>	<u>7</u>	<u>8</u>

# Operations of 2D FDM

$$\frac{\partial^2 \phi}{\partial x^2} + \frac{\partial^2 \phi}{\partial y^2} = f$$

$$\left( \frac{\phi_E - 2\phi_C + \phi_W}{\Delta x^2} \right) + \left( \frac{\phi_N - 2\phi_C + \phi_S}{\Delta y^2} \right) = f_C$$



57	58	59	60	61	62	63	64
49	50	51	52	53	54	55	56
41	42	43	44	45	46	47	48
33	34	35	36	37	38	39	40
25	26	27	28	29	30	31	32
17	18	19	20	21	22	23	24
9	10	11	12	13	14	15	16
1	2	3	4	5	6	7	8

# Computation (1/3)

<u>PE#2</u>	<u>57</u>	<u>58</u>	<u>59</u>	<u>60</u>	<u>61</u>	<u>62</u>	<u>63</u>	<u>64</u>	<u>PE#3</u>
	<u>49</u>	<u>50</u>	<u>51</u>	<u>52</u>	<u>53</u>	<u>54</u>	<u>55</u>	<u>56</u>	
	<u>41</u>	<u>42</u>	<u>43</u>	<u>44</u>	<u>45</u>	<u>46</u>	<u>47</u>	<u>48</u>	
	<u>33</u>	<u>34</u>	<u>35</u>	<u>36</u>	<u>37</u>	<u>38</u>	<u>39</u>	<u>40</u>	
	<u>25</u>	<u>26</u>	<u>27</u>	<u>28</u>	<u>29</u>	<u>30</u>	<u>31</u>	<u>32</u>	
	<u>17</u>	<u>18</u>	<u>19</u>	<u>20</u>	<u>21</u>	<u>22</u>	<u>23</u>	<u>24</u>	
	<u>9</u>	<u>10</u>	<u>11</u>	<u>12</u>	<u>13</u>	<u>14</u>	<u>15</u>	<u>16</u>	
<u>PE#0</u>	<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>	<u>5</u>	<u>6</u>	<u>7</u>	<u>8</u>	<u>PE#1</u>

- On each PE, info. of internal pts ( $i=1-N(=16)$ ) are read from distributed local data, info. of boundary pts are sent to neighbors, and they are received as info. of external pts.

# Computation (2/3): Before Send/Recv

**PE#2**

<u>57</u>	<u>58</u>	<u>59</u>	<u>60</u>	
<u>49</u>	<u>50</u>	<u>51</u>	<u>52</u>	
<u>41</u>	<u>42</u>	<u>43</u>	<u>44</u>	
<u>33</u>	<u>34</u>	<u>35</u>	<u>36</u>	

**PE#3**

	<u>61</u>	<u>62</u>	<u>63</u>	<u>64</u>
	<u>53</u>	<u>54</u>	<u>55</u>	<u>56</u>
	<u>45</u>	<u>46</u>	<u>47</u>	<u>48</u>
	<u>37</u>	<u>38</u>	<u>39</u>	<u>40</u>

1: <u>33</u>	9: <u>49</u>	17: <u>?</u>
2: <u>34</u>	10: <u>50</u>	18: <u>?</u>
3: <u>35</u>	11: <u>51</u>	19: <u>?</u>
4: <u>36</u>	12: <u>52</u>	20: <u>?</u>
5: <u>41</u>	13: <u>57</u>	21: <u>?</u>
6: <u>42</u>	14: <u>58</u>	22: <u>?</u>
7: <u>43</u>	15: <u>59</u>	23: <u>?</u>
8: <u>44</u>	16: <u>60</u>	24: <u>?</u>

1: <u>37</u>	9: <u>53</u>	17: <u>?</u>
2: <u>38</u>	10: <u>54</u>	18: <u>?</u>
3: <u>39</u>	11: <u>55</u>	19: <u>?</u>
4: <u>40</u>	12: <u>56</u>	20: <u>?</u>
5: <u>45</u>	13: <u>61</u>	21: <u>?</u>
6: <u>46</u>	14: <u>62</u>	22: <u>?</u>
7: <u>47</u>	15: <u>63</u>	23: <u>?</u>
8: <u>48</u>	16: <u>64</u>	24: <u>?</u>

**PE#0**

<u>25</u>	<u>26</u>	<u>27</u>	<u>28</u>	
<u>17</u>	<u>18</u>	<u>19</u>	<u>20</u>	
<u>9</u>	<u>10</u>	<u>11</u>	<u>12</u>	
<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>	

**PE#1**

	<u>29</u>	<u>30</u>	<u>31</u>	<u>32</u>
	<u>21</u>	<u>22</u>	<u>23</u>	<u>24</u>
	<u>13</u>	<u>14</u>	<u>15</u>	<u>16</u>
	<u>5</u>	<u>6</u>	<u>7</u>	<u>8</u>

1: <u>1</u>	9: <u>17</u>	17: <u>?</u>
2: <u>2</u>	10: <u>18</u>	18: <u>?</u>
3: <u>3</u>	11: <u>19</u>	19: <u>?</u>
4: <u>4</u>	12: <u>20</u>	20: <u>?</u>
5: <u>9</u>	13: <u>25</u>	21: <u>?</u>
6: <u>10</u>	14: <u>26</u>	22: <u>?</u>
7: <u>11</u>	15: <u>27</u>	23: <u>?</u>
8: <u>12</u>	16: <u>28</u>	24: <u>?</u>

1: <u>5</u>	9: <u>21</u>	17: <u>?</u>
2: <u>6</u>	10: <u>22</u>	18: <u>?</u>
3: <u>7</u>	11: <u>23</u>	19: <u>?</u>
4: <u>8</u>	12: <u>24</u>	20: <u>?</u>
5: <u>13</u>	13: <u>29</u>	21: <u>?</u>
6: <u>14</u>	14: <u>30</u>	22: <u>?</u>
7: <u>15</u>	15: <u>31</u>	23: <u>?</u>
8: <u>16</u>	16: <u>32</u>	24: <u>?</u>

# Computation (2/3): Before Send/Recv

1: <u>33</u>	9: <u>49</u>	17: <u>?</u>
2: <u>34</u>	10: <u>50</u>	18: <u>?</u>
3: <u>35</u>	11: <u>51</u>	19: <u>?</u>
4: <u>36</u>	12: <u>52</u>	20: <u>?</u>
5: <u>41</u>	13: <u>57</u>	21: <u>?</u>
6: <u>42</u>	14: <u>58</u>	22: <u>?</u>
7: <u>43</u>	15: <u>59</u>	23: <u>?</u>
8: <u>44</u>	16: <u>60</u>	24: <u>?</u>

**PE#2**

<u>57</u>	<u>58</u>	<u>59</u>	<u>60</u>	
<u>49</u>	<u>50</u>	<u>51</u>	<u>52</u>	
<u>41</u>	<u>42</u>	<u>43</u>	<u>44</u>	
<u>33</u>	<u>34</u>	<u>35</u>	<u>36</u>	

**PE#3**

	<u>61</u>	<u>62</u>	<u>63</u>	<u>64</u>
	<u>53</u>	<u>54</u>	<u>55</u>	<u>56</u>
	<u>45</u>	<u>46</u>	<u>47</u>	<u>48</u>
	<u>37</u>	<u>38</u>	<u>39</u>	<u>40</u>

1: <u>37</u>	9: <u>53</u>	17: <u>?</u>
2: <u>38</u>	10: <u>54</u>	18: <u>?</u>
3: <u>39</u>	11: <u>55</u>	19: <u>?</u>
4: <u>40</u>	12: <u>56</u>	20: <u>?</u>
5: <u>45</u>	13: <u>61</u>	21: <u>?</u>
6: <u>46</u>	14: <u>62</u>	22: <u>?</u>
7: <u>47</u>	15: <u>63</u>	23: <u>?</u>
8: <u>48</u>	16: <u>64</u>	24: <u>?</u>

1: <u>1</u>	9: <u>17</u>	17: <u>?</u>
2: <u>2</u>	10: <u>18</u>	18: <u>?</u>
3: <u>3</u>	11: <u>19</u>	19: <u>?</u>
4: <u>4</u>	12: <u>20</u>	20: <u>?</u>
5: <u>9</u>	13: <u>25</u>	21: <u>?</u>
6: <u>10</u>	14: <u>26</u>	22: <u>?</u>
7: <u>11</u>	15: <u>27</u>	23: <u>?</u>
8: <u>12</u>	16: <u>28</u>	24: <u>?</u>

**PE#0**

<u>25</u>	<u>26</u>	<u>27</u>	<u>28</u>	
<u>17</u>	<u>18</u>	<u>19</u>	<u>20</u>	
<u>9</u>	<u>10</u>	<u>11</u>	<u>12</u>	
<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>	

**PE#1**

	<u>29</u>	<u>30</u>	<u>31</u>	<u>32</u>
	<u>21</u>	<u>22</u>	<u>23</u>	<u>24</u>
	<u>13</u>	<u>14</u>	<u>15</u>	<u>16</u>
	<u>5</u>	<u>6</u>	<u>7</u>	<u>8</u>

1: <u>5</u>	9: <u>21</u>	17: <u>?</u>
2: <u>6</u>	10: <u>22</u>	18: <u>?</u>
3: <u>7</u>	11: <u>23</u>	19: <u>?</u>
4: <u>8</u>	12: <u>24</u>	20: <u>?</u>
5: <u>13</u>	13: <u>29</u>	21: <u>?</u>
6: <u>14</u>	14: <u>30</u>	22: <u>?</u>
7: <u>15</u>	15: <u>31</u>	23: <u>?</u>
8: <u>16</u>	16: <u>32</u>	24: <u>?</u>

**PE#0**

**PE#1**

# Computation (3/3): After Send/Recv

**PE#2**

<u>57</u>	<u>58</u>	<u>59</u>	<u>60</u>	<u>61</u>
<u>49</u>	<u>50</u>	<u>51</u>	<u>52</u>	<u>53</u>
<u>41</u>	<u>42</u>	<u>43</u>	<u>44</u>	<u>45</u>
<u>33</u>	<u>34</u>	<u>35</u>	<u>36</u>	<u>37</u>
<u>25</u>	<u>26</u>	<u>27</u>	<u>28</u>	

**PE#3**

<u>60</u>	<u>61</u>	<u>62</u>	<u>63</u>	<u>64</u>
<u>52</u>	<u>53</u>	<u>54</u>	<u>55</u>	<u>56</u>
<u>44</u>	<u>45</u>	<u>46</u>	<u>47</u>	<u>48</u>
<u>36</u>	<u>37</u>	<u>38</u>	<u>39</u>	<u>40</u>
	<u>29</u>	<u>30</u>	<u>31</u>	<u>32</u>

1: <u>33</u>	9: <u>49</u>	17: <u>37</u>
2: <u>34</u>	10: <u>50</u>	18: <u>45</u>
3: <u>35</u>	11: <u>51</u>	19: <u>53</u>
4: <u>36</u>	12: <u>52</u>	20: <u>61</u>
5: <u>41</u>	13: <u>57</u>	21: <u>25</u>
6: <u>42</u>	14: <u>58</u>	22: <u>26</u>
7: <u>43</u>	15: <u>59</u>	23: <u>27</u>
8: <u>44</u>	16: <u>60</u>	24: <u>28</u>

1: <u>37</u>	9: <u>53</u>	17: <u>36</u>
2: <u>38</u>	10: <u>54</u>	18: <u>44</u>
3: <u>39</u>	11: <u>55</u>	19: <u>52</u>
4: <u>40</u>	12: <u>56</u>	20: <u>60</u>
5: <u>45</u>	13: <u>61</u>	21: <u>29</u>
6: <u>46</u>	14: <u>62</u>	22: <u>30</u>
7: <u>47</u>	15: <u>63</u>	23: <u>31</u>
8: <u>48</u>	16: <u>64</u>	24: <u>32</u>

1: <u>1</u>	9: <u>17</u>	17: <u>5</u>
2: <u>2</u>	10: <u>18</u>	18: <u>14</u>
3: <u>3</u>	11: <u>19</u>	19: <u>21</u>
4: <u>4</u>	12: <u>20</u>	20: <u>29</u>
5: <u>9</u>	13: <u>25</u>	21: <u>33</u>
6: <u>10</u>	14: <u>26</u>	22: <u>34</u>
7: <u>11</u>	15: <u>27</u>	23: <u>35</u>
8: <u>12</u>	16: <u>28</u>	24: <u>36</u>

<u>33</u>	<u>34</u>	<u>35</u>	<u>36</u>	
<u>25</u>	<u>26</u>	<u>27</u>	<u>28</u>	<u>29</u>
<u>17</u>	<u>18</u>	<u>19</u>	<u>20</u>	<u>21</u>
<u>9</u>	<u>10</u>	<u>11</u>	<u>12</u>	<u>13</u>
<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>	<u>5</u>

	<u>37</u>	<u>38</u>	<u>39</u>	<u>40</u>
<u>28</u>	<u>29</u>	<u>30</u>	<u>31</u>	<u>32</u>
<u>20</u>	<u>21</u>	<u>22</u>	<u>23</u>	<u>24</u>
<u>12</u>	<u>13</u>	<u>14</u>	<u>15</u>	<u>16</u>
<u>4</u>	<u>5</u>	<u>6</u>	<u>7</u>	<u>8</u>

1: <u>5</u>	9: <u>21</u>	17: <u>4</u>
2: <u>6</u>	10: <u>22</u>	18: <u>12</u>
3: <u>7</u>	11: <u>23</u>	19: <u>20</u>
4: <u>8</u>	12: <u>24</u>	20: <u>28</u>
5: <u>13</u>	13: <u>29</u>	21: <u>37</u>
6: <u>14</u>	14: <u>30</u>	22: <u>38</u>
7: <u>15</u>	15: <u>31</u>	23: <u>39</u>
8: <u>16</u>	16: <u>32</u>	24: <u>40</u>

**PE#0**

**PE#1**

# Point-to-Point Communication

- What is PtoP Communication ?
- 2D Problem, Generalized Communication Table
  - 2D FDM
  - Problem Setting
  - Distributed Local Data and Communication Table
  - Implementation
- Report S2

# Overview of Distributed Local Data

Example on PE#0

PE#2

<u>25</u>	<u>26</u>	<u>27</u>	<u>28</u>	
<u>17</u>	<u>18</u>	<u>19</u>	<u>20</u>	
<u>9</u>	<u>10</u>	<u>11</u>	<u>12</u>	
<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>	

PE#0 PE#1

PE#2

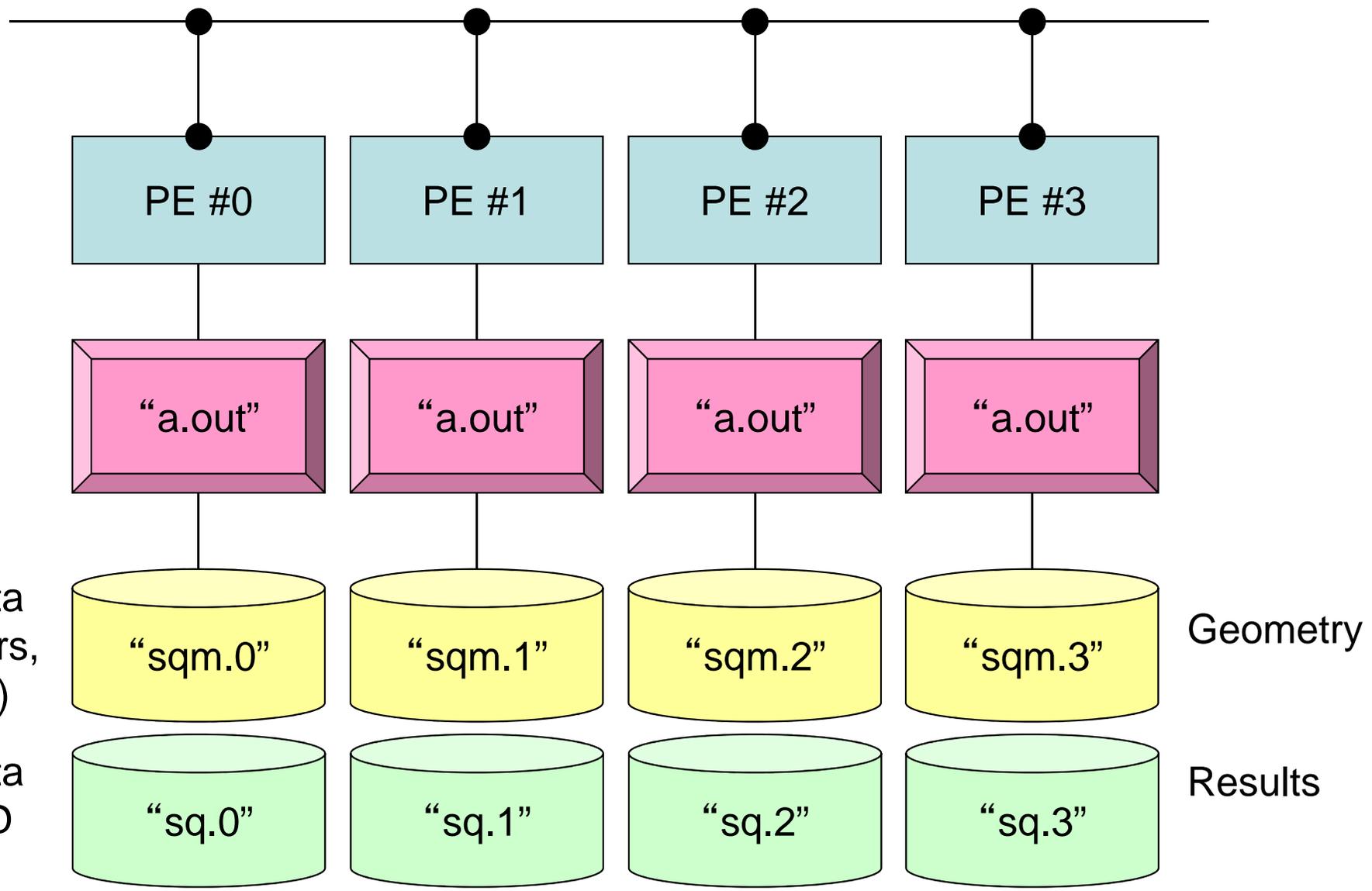
13	14	15	16	
9	10	11	12	
5	6	7	8	
1	2	3	4	

PE#0 PE#1

Value at each mesh (= Global ID)

Local ID

# SPMD...



Dist. Local Data Sets (Neighbors, Comm. Tables)

Dist. Local Data Sets (Global ID of Internal Points)

Geometry

Results

# 2D FDM: PE#0

## Information at each domain (1/4)

### Internal Points

Meshes originally assigned to the domain

13	14	15	16
9	10	11	12
5	6	7	8
1	2	3	4

# 2D FDM: PE#0

## Information at each domain (2/4)

PE#2

●	●	●	●	
13	14	15	16	●
9	10	11	12	●
5	6	7	8	●
1	2	3	4	●

PE#1

### Internal Points

Meshes originally assigned to the domain

### External Points

Meshes originally assigned to different domain, but required for computation of meshes in the domain (meshes in overlapped regions)

- Sleeves
- Halo



# 2D FDM: PE#0

## Information at each domain (3/4)

PE#2

●	●	●	●	
13	14	15	16	●
9	10	11	12	●
5	6	7	8	●
1	2	3	4	●

PE#1

### Internal Points

Meshes originally assigned to the domain

### External Points

Meshes originally assigned to different domain, but required for computation of meshes in the domain (meshes in overlapped regions)

### Boundary Points

Internal points, which are also external points of other domains (used in computations of meshes in other domains)

# 2D FDM: PE#0

## Information at each domain (4/4)

**PE#2**

●	●	●	●	
13	14	15	16	●
9	10	11	12	●
5	6	7	8	●
1	2	3	4	●

**PE#1**

### Internal Points

Meshes originally assigned to the domain

### External Points

Meshes originally assigned to different domain, but required for computation of meshes in the domain (meshes in overlapped regions)

### Boundary Points

Internal points, which are also external points of other domains (used in computations of meshes in other domains)

### Relationships between Domains

Communication Table: External/Boundary Points  
Neighbors

# Description of Distributed Local Data

21	22	23	24	
13	14	15	16	20
9	10	11	12	19
5	6	7	8	18
1	2	3	4	17

- **Internal/External Points**
  - Numbering: Starting from internal pts, then external pts after that
- **Neighbors**
  - Shares overlapped meshes
  - Number and ID of neighbors
- **Import Table (Receive)**
  - From where, how many, and which external points are received/imported ?
- **Export Table (Send)**
  - To where, how many and which boundary points are sent/exported ?

# Overview of Distributed Local Data

Example on PE#0

PE#2

<u>25</u>	<u>26</u>	<u>27</u>	<u>28</u>	
<u>17</u>	<u>18</u>	<u>19</u>	<u>20</u>	
<u>9</u>	<u>10</u>	<u>11</u>	<u>12</u>	
<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>	

PE#0 PE#1

PE#2

21	22	23	24	
13	14	15	16	20
9	10	11	12	19
5	6	7	8	18
1	2	3	4	17

PE#0 PE#1

Value at each mesh (= Global ID)

Local ID

# Generalized Comm. Table: Send

0	1.1 ⊙	2.4 ①,0	3.2 ④,1		
1	3.6 ①	4.3 ⊙,2	2.5 ③,3	3.7 ⑤,4	9.1 ⑦,5
2	5.7 ②	1.5 ④,6	3.1 ⑥,7		
3	9.8 ③	4.1 ①,8	2.5 ④,9	2.7 ⑤,10	
4	11.5 ④	3.1 ⊙,11	9.5 ①,12	10.4 ②,13	4.3 ⑥,14
5	12.4 ⑤	6.5 ②,15	9.5 ⑥,16		
6	23.1 ⑥	6.4 ①,17	2.5 ②,18	1.4 ⑤,19	13.1 ⑦,20
7	51.3 ⑦	9.5 ①,21	1.3 ②,22	9.6 ③,23	3.1 ⑤,24

<b>Diag [i]</b>	Diagonal Components (REAL, i=0~N-1)
<b>Index[i]</b>	Number of Non-Zero Off-Diagonals at Each ROW (INT, i=0~N)
<b>Item[k]</b>	Off-Diagonal Components (Corresponding Column ID) (INT, k=0, index[N])
<b>AMat [k]</b>	Off-Diagonal Components (Value) ( REAL, k=0, index[N] )

{Y}=[A] {X}

```

for (i=0; i<N; i++) {
  Y[i] = Diag[i] * X[i];
  for (k=Index[i]; k<Index[i+1]; k++) {
    Y[i] += AMat[k]*X[Index[k]];
  }
}

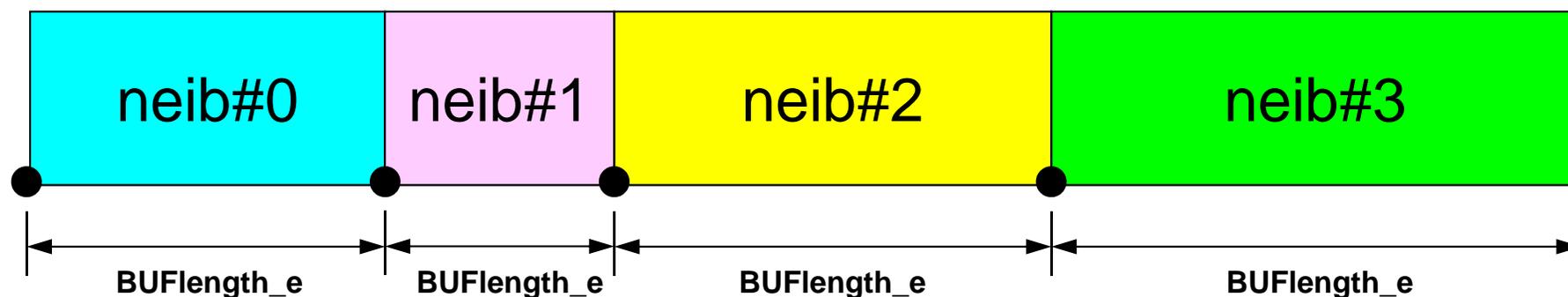
```

- Neighbors
  - NeibPETot, NeibPE[NeibPETot]
- Message size for each neighbor
  - export\_index[NeibPETot+1]
- ID of boundary points
  - export\_item[export\_index[NeibPETot+1]]
- Messages to each neighbor
  - SendBuf[export\_index[NeibPETot+1]]

# SEND: MPI\_Isend/Irecv/Waitall

C

SendBuf



export\_index[0]      export\_index[1]      export\_index[2]      export\_index[3]      export\_index[4]

export\_item (export\_index[neib]:export\_index[neib+1]-1) are sent to neib-th neighbor

```
for (neib=0; neib<NeibPETot;neib++){
  for (k=export_index[neib];k<export_index[neib+1];k++){
    kk= export_item[k];
    SendBuf[k]= VAL[kk];
  }
}
```

Copied to sending buffers

```
for (neib=0; neib<NeibPETot; neib++){
  tag= 0;
  iS_e= export_index[neib];
  iE_e= export_index[neib+1];
  BUFlength_e= iE_e - iS_e

  ierr= MPI_Isend
    (&SendBuf[iS_e], BUFlength_e, MPI_DOUBLE, NeibPE[neib], 0,
     MPI_COMM_WORLD, &ReqSend[neib])
}
```

```
MPI_Waitall(NeibPETot, ReqSend, StatSend);
```

# Generalized Comm. Table: Receive

- Neighbors
  - NeibPETot, NeibPE[NeibPETot]
- Message size for each neighbor
  - import\_index [NeibPETot+1]
- ID of **external** points
  - import\_item [import\_index[NeibPETot+1]]
- Messages from each neighbor
  - RecvBuf [import\_index[NeibPETot+1]]

# RECV: MPI\_Irecv/Irecv/Waitall

C

```

for (neib=0; neib<NeibPETot; neib++){
  tag= 0;
  iS_i= import_index[neib];
  iE_i= import_index[neib+1];
  BUFlength_i= iE_i - iS_i

  ierr= MPI_Irecv
    (&RecvBuf[iS_i], BUFlength_i, MPI_DOUBLE, NeibPE[neib], 0,
     MPI_COMM_WORLD, &ReqRecv[neib])
}

```

```
MPI_Waitall(NeibPETot, ReqRecv, StatRecv);
```

```

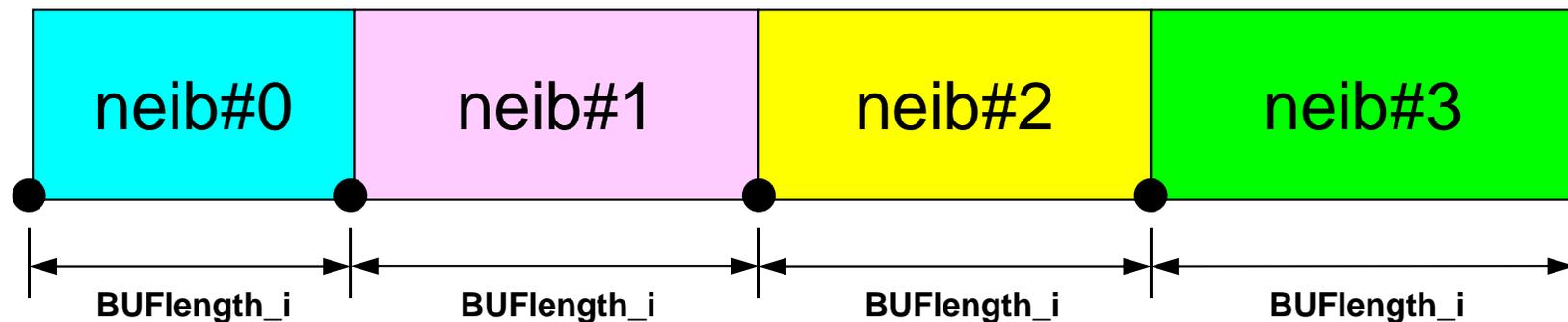
for (neib=0; neib<NeibPETot;neib++){
  for (k=import_index[neib];k<import_index[neib+1];k++){
    kk= import_item[k];
    VAL[kk]= RecvBuf[k];
  }
}

```

Copied from receiving buffer

**import\_item (import\_index[neib]:import\_index[neib+1]-1) are received from neib-th neighbor**

**RecvBuf**



import\_index[0]   import\_index[1]   import\_index[2]   import\_index[3]   import\_index[4]

# Relationship SEND/RECV

```

do neib= 1, NEIBPETOT
  iS_e= export_index(neib-1) + 1
  iE_e= export_index(neib )
  BUFlength_e= iE_e + 1 - iS_e

  call MPI_ISEND
&      (SENDbuf(iS_e), BUFlength_e, MPI_INTEGER, NEIBPE(neib), 0, &
&      MPI_COMM_WORLD, request_send(neib), ierr)
  enddo

```

```

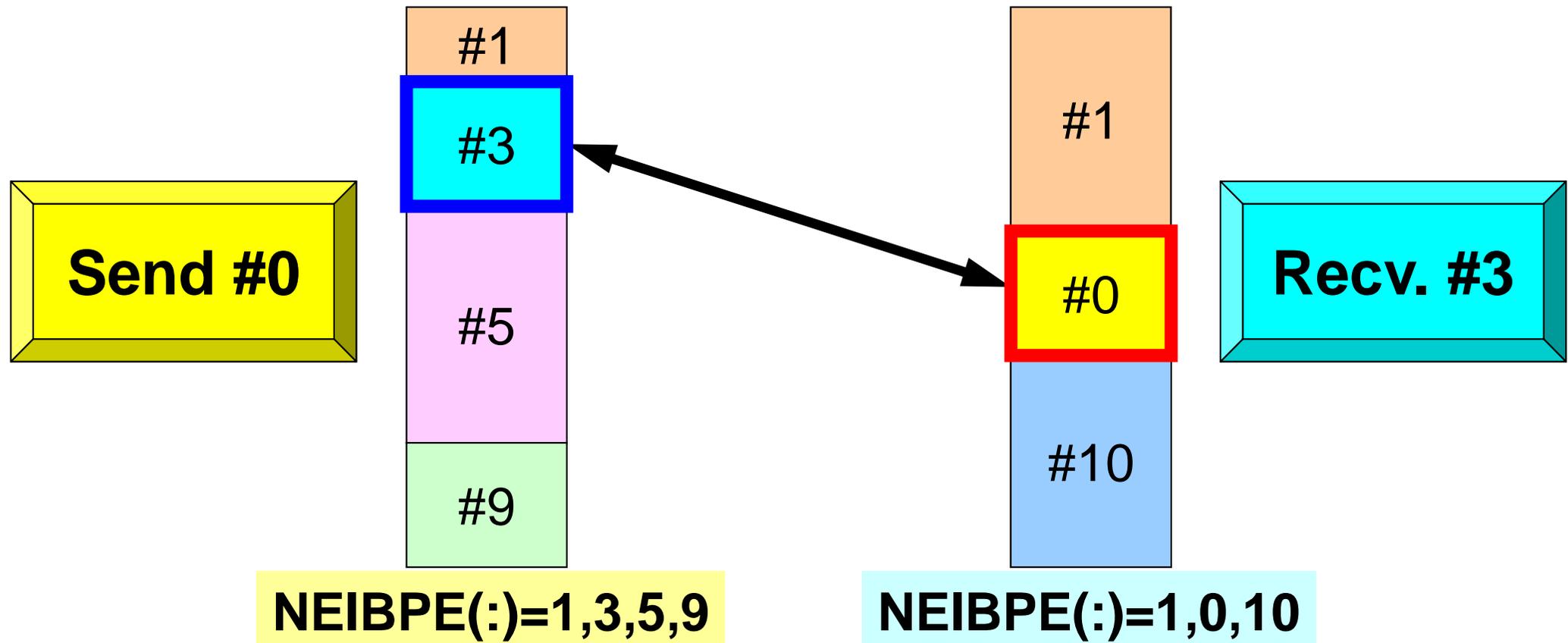
do neib= 1, NEIBPETOT
  iS_i= import_index(neib-1) + 1
  iE_i= import_index(neib )
  BUFlength_i= iE_i + 1 - iS_i

  call MPI_IRECV
&      (RECVbuf(iS_i), BUFlength_i, MPI_INTEGER, NEIBPE(neib), 0, &
&      MPI_COMM_WORLD, request_recv(neib), ierr)
  enddo

```

- Consistency of ID's of sources/destinations, size and contents of messages !
- Communication occurs when NEIBPE(neib) matches

# Relationship SEND/RECV (#0 to #3)



- Consistency of ID's of sources/destinations, size and contents of messages !
- Communication occurs when NEIBPE(neib) matches

# Generalized Comm. Table (1/6)

**PE#2**

<b>21</b>	<b>22</b>	<b>23</b>	<b>24</b>	
<b>13</b>	<b>14</b>	<b>15</b>	<b>16</b>	<b>20</b>
<b>9</b>	<b>10</b>	<b>11</b>	<b>12</b>	<b>19</b>
<b>5</b>	<b>6</b>	<b>7</b>	<b>8</b>	<b>18</b>
<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>17</b>

**PE#1**

```

#NEIBPEtot
2
#NEIBPE
1 2
#NODE
24 16
#IMPORT_index
4 8
#IMPORT_items
17
18
19
20
21
22
23
24
#EXPORT_index
4 8
#EXPORT_items
4
8
12
16
13
14
15
16
    
```

# Generalized Comm. Table (2/6)

**PE#2**

21	22	23	24	
13	14	15	16	20
9	10	11	12	19
5	6	7	8	18
1	2	3	4	17

**PE#1**

```

#NEIBPEtot  Number of neighbors
2
#NEIBPE     ID of neighbors
1 2
#NODE
24 16      Ext/Int Pts, Int Pts
#IMPORT_index
4 8
#IMPORT_items
17
18
19
20
21
22
23
24
#EXPORT_index
4 8
#EXPORT_items
4
8
12
16
13
14
15
16
    
```

# Generalized Comm. Table (3/6)

PE#2

21	22	23	24	
13	14	15	16	20
9	10	11	12	19
5	6	7	8	18
1	2	3	4	17

PE#1

```

#NEIBPEtot
2
#NEIBPE
1 2
#NODE
24 16
#IMPORT_index
4 8
#IMPORT_items
17
18
19
20
21
22
23
24
#EXPORT_index
4 8
#EXPORT_items
4
8
12
16
13
14
15
16
    
```

Four ext pts (1<sup>st</sup>-4<sup>th</sup> items) are imported from 1<sup>st</sup> neighbor (PE#1), and four (5<sup>th</sup>-8<sup>th</sup> items) are from 2<sup>nd</sup> neighbor (PE#2).

# Generalized Comm. Table (4/6)

PE#2

21	22	23	24	
13	14	15	16	20
9	10	11	12	19
5	6	7	8	18
1	2	3	4	17

PE#1

```

#NEIBPEtot
2
#NEIBPE
1 2
#NODE
24 16
#IMPORT_index
4 8
#IMPORT_items
17
18 imported from 1st Neighbor
19 (PE#1) (1st-4th items)
20
21 imported from 2nd Neighbor
22 (PE#2) (5th-8th items)
23
24
#EXPORT_index
4 8
#EXPORT_items
4
8
12
16
13
14
15
16
    
```

# Generalized Comm. Table (5/6)

PE#2

21	22	23	24	
13	14	15	16	20
9	10	11	12	19
5	6	7	8	18
1	2	3	4	17

PE#1

```

#NEIBPEtot
2
#NEIBPE
1 2
#NODE
24 16
#IMPORT_index
4 8
#IMPORT_items
17
18
19
20
21
22
23
24
#EXPORT_index
4 8
#EXPORT_items
4
8
12
16
13
14
15
16

```

Four boundary pts (1<sup>st</sup>-4<sup>th</sup> items) are exported to 1<sup>st</sup> neighbor (PE#1), and four (5<sup>th</sup>-8<sup>th</sup> items) are to 2<sup>nd</sup> neighbor (PE#2).

# Generalized Comm. Table (6/6)

**PE#2**

21	22	23	24	
13	14	15	16	20
9	10	11	12	19
5	6	7	8	18
1	2	3	4	17

**PE#1**

```

#NEIBPEtot
2
#NEIBPE
1 2
#NODE
24 16
#IMPORT_index
4 8
#IMPORT_items
17
18
19
20
21
22
23
24
#EXPORT_index
4 8
#EXPORT_items
4
8
12
16
13
14
15
16
    
```

exported to 1<sup>st</sup> Neighbor (PE#1) (1<sup>st</sup>-4<sup>th</sup> items)

exported to 2<sup>nd</sup> Neighbor (PE#2) (5<sup>th</sup>-8<sup>th</sup> items)

# Generalized Comm. Table (6/6)

**PE#2**

21	22	23	24	
13	14	15	16	20
9	10	11	12	19
5	6	7	8	18
1	2	3	4	17

**PE#1**

An external point is only sent from its original domain.

A boundary point could be referred from more than one domain, and sent to multiple domains (e.g. 16<sup>th</sup> mesh).

# Notice: Send/Recv Arrays

```
#PE0  
send:  
  VEC (start_send) ~  
  VEC (start_send+length_send-1)
```

```
#PE1  
send:  
  VEC (start_send) ~  
  VEC (start_send+length_send-1)
```

```
#PE0  
recv:  
  VEC (start_recv) ~  
  VEC (start_recv+length_recv-1)
```

```
#PE1  
recv:  
  VEC (start_recv) ~  
  VEC (start_recv+length_recv-1)
```

- “length\_send” of sending process must be equal to “length\_recv” of receiving process.
  - PE#0 to PE#1, PE#1 to PE#0
- “sendbuf” and “recvbuf”: different address

# Point-to-Point Communication

- What is PtoP Communication ?
- **2D Problem, Generalized Communication Table**
  - 2D FDM
  - Problem Setting
  - Distributed Local Data and Communication Table
  - **Implementation**
- Report S2

# Sample Program for 2D FDM

```
$ cd /work/gt37/t37XXX/pFEM/mpi/S2
```

```
$ mpiifort -O3 sq-sr1.f
```

```
$ mpiicc -O3 sq-sr1.c
```

(modify go4.sh for 4 processes)

```
$ pjsub go4.sh
```

# Example: sq-sr1.c (1/6)

C

## Initialization

```
#include <stdio.h>
#include <stdlib.h>
#include <string.h>
#include <assert.h>
#include "mpi.h"
int main(int argc, char **argv){

    int n, np, NeibPeTot, BufLength;
    MPI_Status *StatSend, *StatRecv;
    MPI_Request *RequestSend, *RequestRecv;

    int MyRank, PeTot;
    int *val, *SendBuf, *RecvBuf, *NeibPe;
    int *ImportIndex, *ExportIndex, *ImportItem, *ExportItem;

    char FileName[80], line[80];
    int i, nn, neib;
    int iStart, iEnd;
    FILE *fp;

    /*
    !C +-----+
    !C | INIT. MPI |
    !C +-----+
    !C===*/

    MPI_Init(&argc, &argv);
    MPI_Comm_size(MPI_COMM_WORLD, &PeTot);
    MPI_Comm_rank(MPI_COMM_WORLD, &MyRank);
```

# Example: sq-sr1.c (2/6)

C

## Reading distributed local data files (sqm.\*)

```
/*
!C +-----+
!C | DATA file |
!C +-----+
!C===*/

    sprintf(FileName, "sqm.%d", MyRank);
    fp = fopen(FileName, "r");

    fscanf(fp, "%d", &NeibPeTot);
    NeibPe = calloc(NeibPeTot, sizeof(int));
    ImportIndex = calloc(1+NeibPeTot, sizeof(int));
    ExportIndex = calloc(1+NeibPeTot, sizeof(int));

    for(neib=0;neib<NeibPeTot;neib++){
        fscanf(fp, "%d", &NeibPe[neib]);
    }
    fscanf(fp, "%d %d", &np, &n);

    for(neib=1;neib<NeibPeTot+1;neib++){
        fscanf(fp, "%d", &ImportIndex[neib]);}
    nn = ImportIndex[NeibPeTot];
    ImportItem = malloc(nn * sizeof(int));
    for(i=0;i<nn;i++){
        fscanf(fp, "%d", &ImportItem[i]); ImportItem[i]--;}

    for(neib=1;neib<NeibPeTot+1;neib++){
        fscanf(fp, "%d", &ExportIndex[neib]);}
    nn = ExportIndex[NeibPeTot];
    ExportItem = malloc(nn * sizeof(int));

    for(i=0;i<nn;i++){
        fscanf(fp, "%d", &ExportItem[i]);ExportItem[i]--;}
```

# Example: sq-sr1.c (2/6)

C

## Reading distributed local data files (sqm.\*)

```

/*
!C +-----+
!C | DATA file |
!C +-----+
!C===*/

    sprintf(FileName, "sqm.%d", MyRank);
    fp = fopen(FileName, "r");

    fscanf(fp, "%d", &NeibPeTot);
    NeibPe = calloc(NeibPeTot, sizeof(int));
    ImportIndex = calloc(1+NeibPeTot, sizeof(int));
    ExportIndex = calloc(1+NeibPeTot, sizeof(int));

    for(neib=0;neib<NeibPeTot;neib++){
        fscanf(fp, "%d", &NeibPe[neib]);
    }
    fscanf(fp, "%d %d", &np, &n);

    for(neib=1;neib<NeibPeTot+1;neib++){
        fscanf(fp, "%d", &ImportIndex[neib]);}
    nn = ImportIndex[NeibPeTot];
    ImportItem = malloc(nn * sizeof(int));
    for(i=0;i<nn;i++){
        fscanf(fp, "%d", &ImportItem[i]); ImportItem[i]-

    for(neib=1;neib<NeibPeTot+1;neib++){
        fscanf(fp, "%d", &ExportIndex[neib]);}
    nn = ExportIndex[NeibPeTot];
    ExportItem = malloc(nn * sizeof(int));

    for(i=0;i<nn;i++){
        fscanf(fp, "%d", &ExportItem[i]);ExportItem[i]--;}

```

```

#NEIBPEtot
2
#NEIBPE
1 2
#NODE
24 16
#IMPORTindex
4 8
#IMPORTitems
17
18
19
20
21
22
23
24
#EXPORTindex
4 8
#EXPORTitems
4
8
12
16
13
14
15
16

```

# Example: sq-sr1.c (2/6)

C

## Reading distributed local data files (sqm.\*)

```

/*
!C +-----+
!C | DATA file |
!C +-----+
!C===*/

    sprintf(FileName, "sqm.%d", MyRank);
    fp = fopen(FileName, "r");

    fscanf(fp, "%d", &NeibPeTot);
    NeibPe = calloc(NeibPeTot, sizeof(int));
    np   Number of all meshes (internal + external)
    n    Number of internal meshes
    for (neib=0;neib<NeibPeTot;neib++){
        fscanf(fp, "%d", &NeibPe[neib]);
    }
    fscanf(fp, "%d %d", &np, &n);

    for (neib=1;neib<NeibPeTot+1;neib++){
        fscanf(fp, "%d", &ImportIndex[neib]);}
    nn = ImportIndex[NeibPeTot];
    ImportItem = malloc(nn * sizeof(int));
    for (i=0;i<nn;i++){
        fscanf(fp, "%d", &ImportItem[i]); ImportItem[i]-

    for (neib=1;neib<NeibPeTot+1;neib++){
        fscanf(fp, "%d", &ExportIndex[neib]);}
    nn = ExportIndex[NeibPeTot];
    ExportItem = malloc(nn * sizeof(int));

    for (i=0;i<nn;i++){
        fscanf(fp, "%d", &ExportItem[i]);ExportItem[i]--;}

```

```

#NEIBPEtot
2
#NEIBPE
1 2
#NODE
24 16
#IMPORTindex
4 8
#IMPORTitems
17
18
19
20
21
22
23
24
#EXPORTindex
4 8
#EXPORTitems
4
8
12
16
13
14
15
16

```

# Example: sq-sr1.c (2/6)

C

## Reading distributed local data files (sqm.\*)

```

/*
!C +-----+
!C | DATA file |
!C +-----+
!C===*/

    sprintf(FileName, "sqm.%d", MyRank);
    fp = fopen(FileName, "r");

    fscanf(fp, "%d", &NeibPeTot);
    NeibPe = calloc(NeibPeTot, sizeof(int));
    ImportIndex = calloc(1+NeibPeTot, sizeof(int));
    ExportIndex = calloc(1+NeibPeTot, sizeof(int));

    for(neib=0;neib<NeibPeTot;neib++){
        fscanf(fp, "%d", &NeibPe[neib]);
    }
    fscanf(fp, "%d %d", &np, &n);

    for(neib=1;neib<NeibPeTot+1;neib++){
        fscanf(fp, "%d", &ImportIndex[neib]);}
    nn = ImportIndex[NeibPeTot];
    ImportItem = malloc(nn * sizeof(int));
    for(i=0;i<nn;i++){
        fscanf(fp, "%d", &ImportItem[i]); ImportItem[i]-
    }

    for(neib=1;neib<NeibPeTot+1;neib++){
        fscanf(fp, "%d", &ExportIndex[neib]);}
    nn = ExportIndex[NeibPeTot];
    ExportItem = malloc(nn * sizeof(int));

    for(i=0;i<nn;i++){
        fscanf(fp, "%d", &ExportItem[i]);ExportItem[i]--;}

```

```

#NEIBPEtot
2
#NEIBPE
1 2
#NODE
24 16
#IMPORTindex
4 8
#IMPORTitems
17
18
19
20
21
22
23
24
#EXPORTindex
4 8
#EXPORTitems
4
8
12
16
13
14
15
16

```

# Example: sq-sr1.c (2/6)

C

## Reading distributed local data files (sqm.\*)

```

/*
!C +-----+
!C | DATA file |
!C +-----+
!C===*/

    sprintf(FileName, "sqm.%d", MyRank);
    fp = fopen(FileName, "r");

    fscanf(fp, "%d", &NeibPeTot);
    NeibPe = calloc(NeibPeTot, sizeof(int));
    ImportIndex = calloc(1+NeibPeTot, sizeof(int));
    ExportIndex = calloc(1+NeibPeTot, sizeof(int));

    for(neib=0;neib<NeibPeTot;neib++){
        fscanf(fp, "%d", &NeibPe[neib]);
    }
    fscanf(fp, "%d %d", &np, &n);

    for(neib=1;neib<NeibPeTot+1;neib++){
        fscanf(fp, "%d", &ImportIndex[neib]);}
    nn = ImportIndex[NeibPeTot];
    ImportItem = malloc(nn * sizeof(int));
    for(i=0;i<nn;i++){
        fscanf(fp, "%d", &ImportItem[i]); ImportItem[i]-

    for(neib=1;neib<NeibPeTot+1;neib++){
        fscanf(fp, "%d", &ExportIndex[neib]);}
    nn = ExportIndex[NeibPeTot];
    ExportItem = malloc(nn * sizeof(int));

    for(i=0;i<nn;i++){
        fscanf(fp, "%d", &ExportItem[i]);ExportItem[i]--;}

```

```

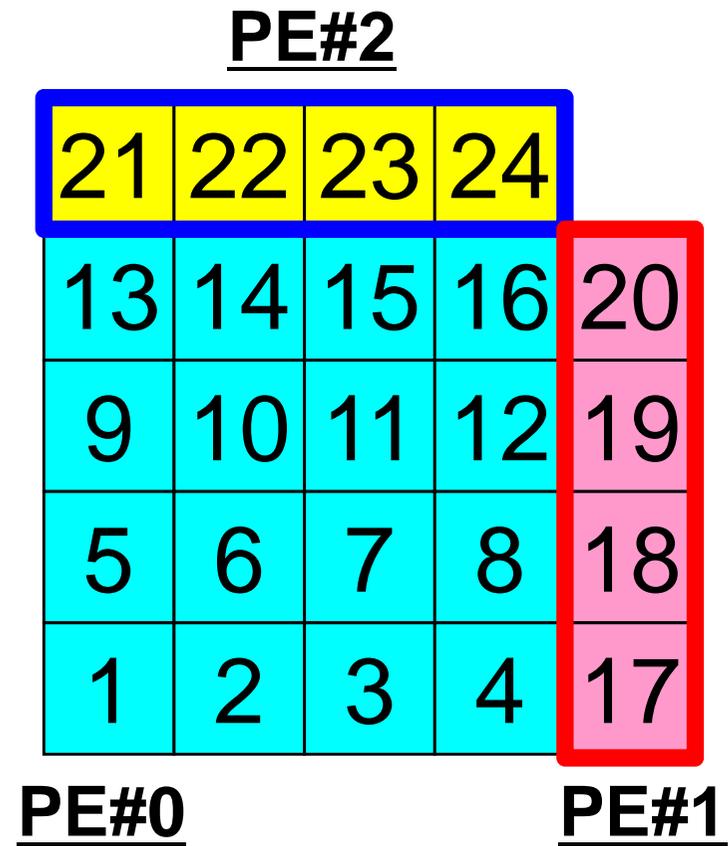
#NEIBPEtot
2
#NEIBPE
1 2
#NODE
24 16
#IMPORTindex
4 8
#IMPORTitems
17
18
19
20
21
22
23
24
#EXPORTindex
4 8
#EXPORTitems
4
8
12
16
13
14
15
16

```

# RECV/Import: PE#0

```

#NEIBPEtot
2
#NEIBPE
1 2
#NODE
24 16
#IMPORTindex
4 8
#IMPORTitems
17
18
19
20
21
22
23
24
#EXPORTindex
4 8
#EXPORTitems
4
8
12
16
13
14
15
16
    
```



# Example: sq-sr1.c (2/6)

C

## Reading distributed local data files (sqm.\*)

```

/*
!C +-----+
!C | DATA file |
!C +-----+
!C===*/

    sprintf(FileName, "sqm.%d", MyRank);
    fp = fopen(FileName, "r");

    fscanf(fp, "%d", &NeibPeTot);
    NeibPe = calloc(NeibPeTot, sizeof(int));
    ImportIndex = calloc(1+NeibPeTot, sizeof(int));
    ExportIndex = calloc(1+NeibPeTot, sizeof(int));

    for(neib=0;neib<NeibPeTot;neib++){
        fscanf(fp, "%d", &NeibPe[neib]);
    }
    fscanf(fp, "%d %d", &np, &n);

    for(neib=1;neib<NeibPeTot+1;neib++){
        fscanf(fp, "%d", &ImportIndex[neib]);}
    nn = ImportIndex[NeibPeTot];
    ImportItem = malloc(nn * sizeof(int));
    for(i=0;i<nn;i++){
        fscanf(fp, "%d", &ImportItem[i]); ImportItem[i]-

for(neib=1;neib<NeibPeTot+1;neib++){
        fscanf(fp, "%d", &ExportIndex[neib]);}
nn = ExportIndex[NeibPeTot];
    ExportItem = malloc(nn * sizeof(int));

    for(i=0;i<nn;i++){
        fscanf(fp, "%d", &ExportItem[i]);ExportItem[i]--;}

```

```

#NEIBPEtot
2
#NEIBPE
1 2
#NODE
24 16
#IMPORTindex
4 8
#IMPORTitems
17
18
19
20
21
22
23
24
#EXPORTindex
4 8
#EXPORTitems
4
8
12
16
13
14
15
16

```

# Example: sq-sr1.c (2/6)

C

## Reading distributed local data files (sqm.\*)

```

/*
!C +-----+
!C | DATA file |
!C +-----+
!C===*/

    sprintf(FileName, "sqm.%d", MyRank);
    fp = fopen(FileName, "r");

    fscanf(fp, "%d", &NeibPeTot);
    NeibPe = calloc(NeibPeTot, sizeof(int));
    ImportIndex = calloc(1+NeibPeTot, sizeof(int));
    ExportIndex = calloc(1+NeibPeTot, sizeof(int));

    for(neib=0;neib<NeibPeTot;neib++){
        fscanf(fp, "%d", &NeibPe[neib]);
    }
    fscanf(fp, "%d %d", &np, &n);

    for(neib=1;neib<NeibPeTot+1;neib++){
        fscanf(fp, "%d", &ImportIndex[neib]);}
    nn = ImportIndex[NeibPeTot];
    ImportItem = malloc(nn * sizeof(int));
    for(i=0;i<nn;i++){
        fscanf(fp, "%d", &ImportItem[i]); ImportItem[i]-

    for(neib=1;neib<NeibPeTot+1;neib++){
        fscanf(fp, "%d", &ExportIndex[neib]);}
    nn = ExportIndex[NeibPeTot];
    ExportItem = malloc(nn * sizeof(int));

    for(i=0;i<nn;i++){
        fscanf(fp, "%d", &ExportItem[i]);ExportItem[i]--;

```

```

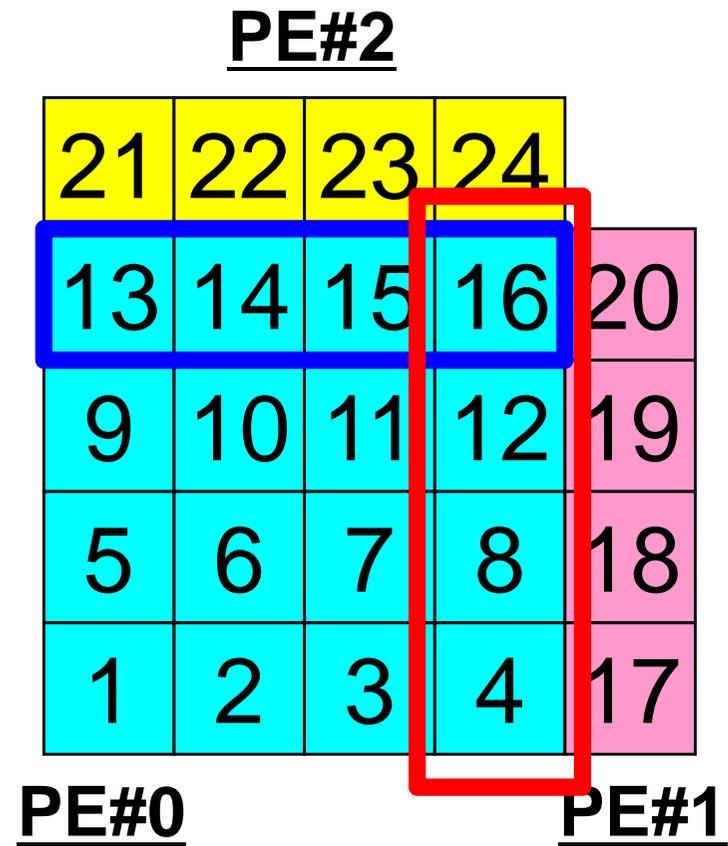
#NEIBPEtot
2
#NEIBPE
1 2
#NODE
24 16
#IMPORTindex
4 8
#IMPORTitems
17
18
19
20
21
22
23
24
#EXPORTindex
4 8
#EXPORTitems
4
8
12
16
13
14
15
16

```

# SEND/Export: PE#0

```

#NEIBPEtot
2
#NEIBPE
1 2
#NODE
24 16
#IMPORTindex
4 8
#IMPORTitems
17
18
19
20
21
22
23
24
#EXPORTindex
4 8
#EXPORTitems
4
8
12
16
13
14
15
16
    
```



# Example: sq-sr1.c (3/6)

C

## Reading distributed local data files (sq.\*)

```

sprintf(FileName, "sq.%d", MyRank);

fp = fopen(FileName, "r");
assert(fp != NULL);

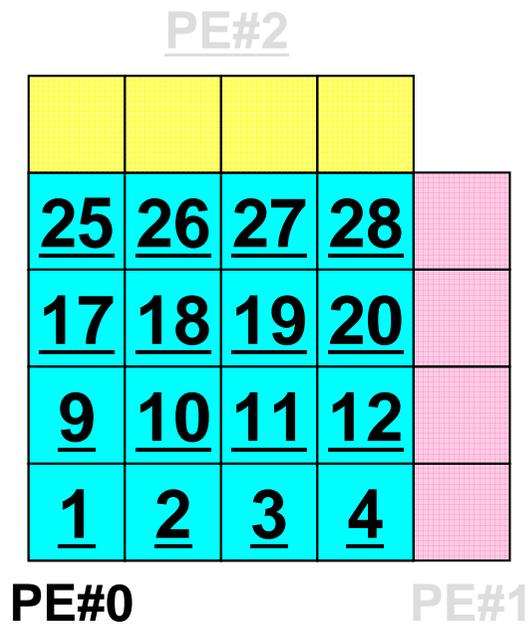
val = calloc(np, sizeof(*val));
for(i=0;i<n;i++){
    fscanf(fp, "%d", &val[i]);
}

```

**n** : Number of internal points

**val** : Global ID of meshes

**val** on external points are unknown at this stage.



1  
2  
3  
4  
9  
10  
11  
12  
17  
18  
19  
20  
25  
26  
27  
28

# Example: sq-sr1.c (4/6)

C

## Preparation of sending/receiving buffers

```

/*
!C
!C +-----+
!C | BUFFER |
!C +-----+
!C===*/
    SendBuf = calloc(ExportIndex[NeibPeTot], sizeof(*SendBuf));
    RecvBuf = calloc(ImportIndex[NeibPeTot], sizeof(*RecvBuf));

    for(neib=0;neib<NeibPeTot;neib++){
        iStart = ExportIndex[neib];
        iEnd   = ExportIndex[neib+1];
        for(i=iStart;i<iEnd;i++){
            SendBuf[i] = val[ExportItem[i]];
        }
    }

```

Info. of boundary points is written into sending buffer (**SendBuf**).

Info. sent to **NeibPe[neib]** is stored in

**SendBuf[ExportIndex[neib] :**

**ExportIndex[neib+1]-1]**

# Sending Buffer is nice ...

C

```

for (neib=0; neib<NeibPETot; neib++){
  tag= 0;
  iS_e= export_index[neib];
  iE_e= export_index[neib+1];
  BUFlength_e= iE_e - iS_e

  ierr= MPI_Isend
        (&SendBuf[iS_e], BUFlength_e, MPI_DOUBLE, NeibPE[neib], 0,
         MPI_COMM_WORLD, &ReqSend[neib])
}

```

**PE#2**

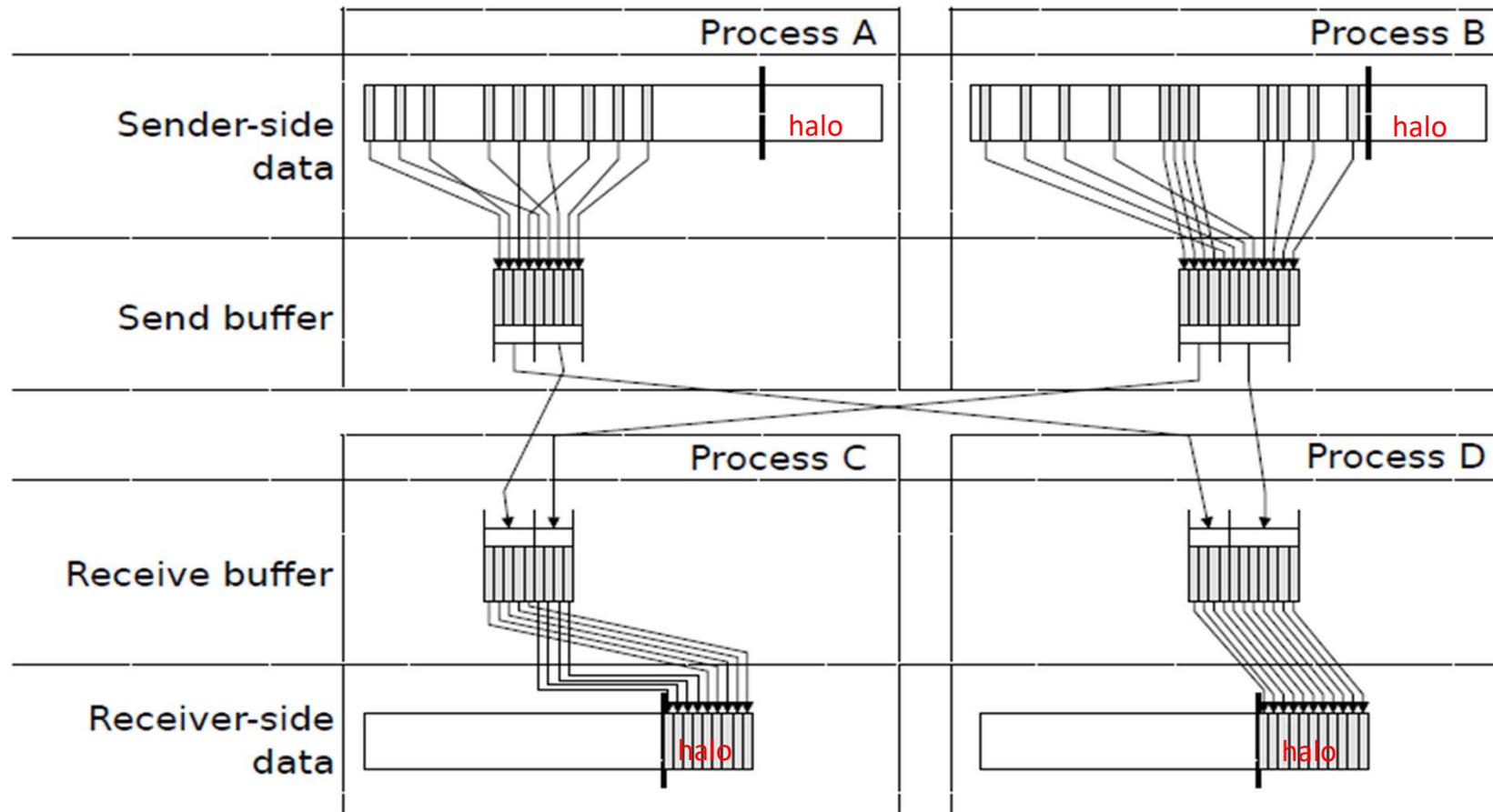
21	22	23	24	
13	14	15	16	20
9	10	11	12	19
5	6	7	8	18
1	2	3	4	17

**PE#0** **PE#1**

Numbering of these boundary nodes is not continuous, therefore the following procedure of MPI\_Isend is not applied directly:

- Starting address of sending buffer
- XX-messages from that address

# Communication Pattern using 1D Structure



# Example: sq-sr1.c (5/6)

C

## SEND/Export: MPI\_Isend

```

/*
!C
!C +-----+
!C | SEND-RECV |
!C +-----+
!C===*/

StatSend = malloc(sizeof(MPI_Status) * NeibPeTot);
StatRecv = malloc(sizeof(MPI_Status) * NeibPeTot);
RequestSend = malloc(sizeof(MPI_Request) * NeibPeTot);
RequestRecv = malloc(sizeof(MPI_Request) * NeibPeTot);

for (neib=0; neib<NeibPeTot; neib++) {
    iStart = ExportIndex[neib];
    iEnd = ExportIndex[neib+1];
    BufLength = iEnd - iStart;
    MPI_Isend(&SendBuf[iStart], BufLength, MPI_INT,
              NeibPe[neib], 0, MPI_COMM_WORLD, &RequestSend[neib]);
}

for (neib=0; neib<NeibPeTot; neib++) {
    iStart = ImportIndex[neib];
    iEnd = ImportIndex[neib+1];
    BufLength = iEnd - iStart;

    MPI_Irecv(&RecvBuf[iStart], BufLength, MPI_INT,
              NeibPe[neib], 0, MPI_COMM_WORLD, &RequestRecv[neib]);
}

```

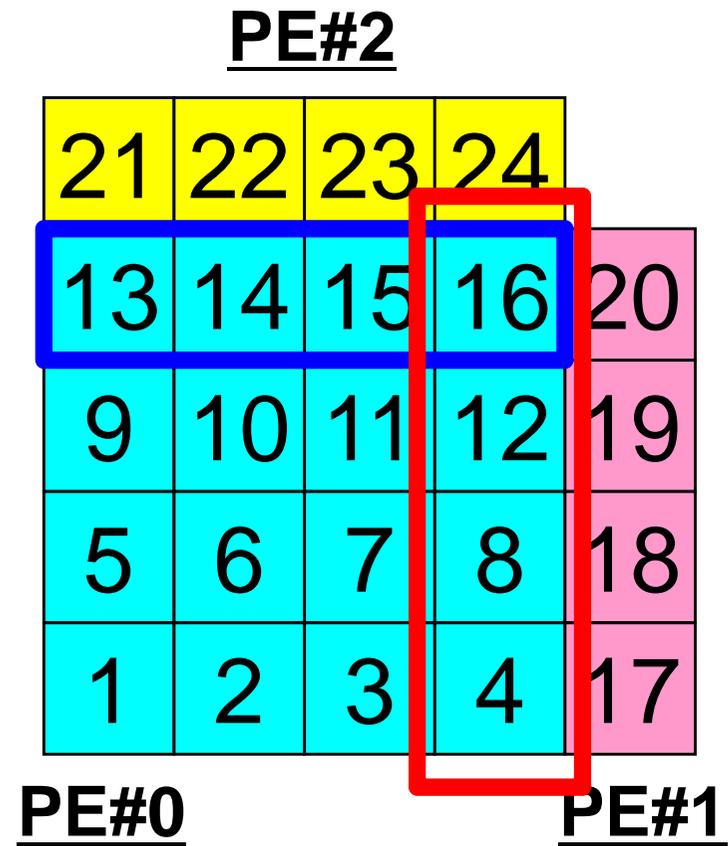
PE#2				PE#3			
57	58	59	60	61	62	63	64
49	50	51	52	53	54	55	56
41	42	43	44	45	46	47	48
33	34	35	36	37	38	39	40
25	26	27	28	29	30	31	32
17	18	19	20	21	22	23	24
9	10	11	12	13	14	15	16
1	2	3	4	5	6	7	8

PE#0 PE#1

# SEND/Export: PE#0

```

#NEIBPEtot
2
#NEIBPE
1 2
#NODE
24 16
#IMPORTindex
4 8
#IMPORTitems
17
18
19
20
21
22
23
24
#EXPORTindex
4 8
#EXPORTitems
4
8
12
16
13
14
15
16
    
```



# SEND: MPI\_Isend/Irecv/Waitall

C

SendBuf



export\_index[0]      export\_index[1]      export\_index[2]      export\_index[3]      export\_index[4]

export\_item (export\_index[neib]:export\_index[neib+1]-1) are sent to neib-th neighbor

```
for (neib=0; neib<NeibPETot;neib++){
  for (k=export_index[neib];k<export_index[neib+1];k++){
    kk= export_item[k];
    SendBuf[k]= VAL[kk];
  }
}
```

Copied to sending buffers

```
for (neib=0; neib<NeibPETot; neib++){
  tag= 0;
  iS_e= export_index[neib];
  iE_e= export_index[neib+1];
  BUFlength_e= iE_e - iS_e

  ierr= MPI_Isend
    (&SendBuf[iS_e], BUFlength_e, MPI_DOUBLE, NeibPE[neib], 0,
     MPI_COMM_WORLD, &ReqSend[neib])
}
```

```
MPI_Waitall(NeibPETot, ReqSend, StatSend);
```

# Notice: Send/Recv Arrays

```
#PE0  
send:  
  VEC (start_send) ~  
  VEC (start_send+length_send-1)
```

```
#PE1  
send:  
  VEC (start_send) ~  
  VEC (start_send+length_send-1)
```

```
#PE0  
recv:  
  VEC (start_recv) ~  
  VEC (start_recv+length_recv-1)
```

```
#PE1  
recv:  
  VEC (start_recv) ~  
  VEC (start_recv+length_recv-1)
```

- “length\_send” of sending process must be equal to “length\_recv” of receiving process.
  - PE#0 to PE#1, PE#1 to PE#0
- “sendbuf” and “recvbuf”: different address

# Relationship SEND/RECV

```
do neib= 1, NEIBPETOT
  iS_e= export_index(neib-1) + 1
  iE_e= export_index(neib )
  BUFlength_e= iE_e + 1 - iS_e

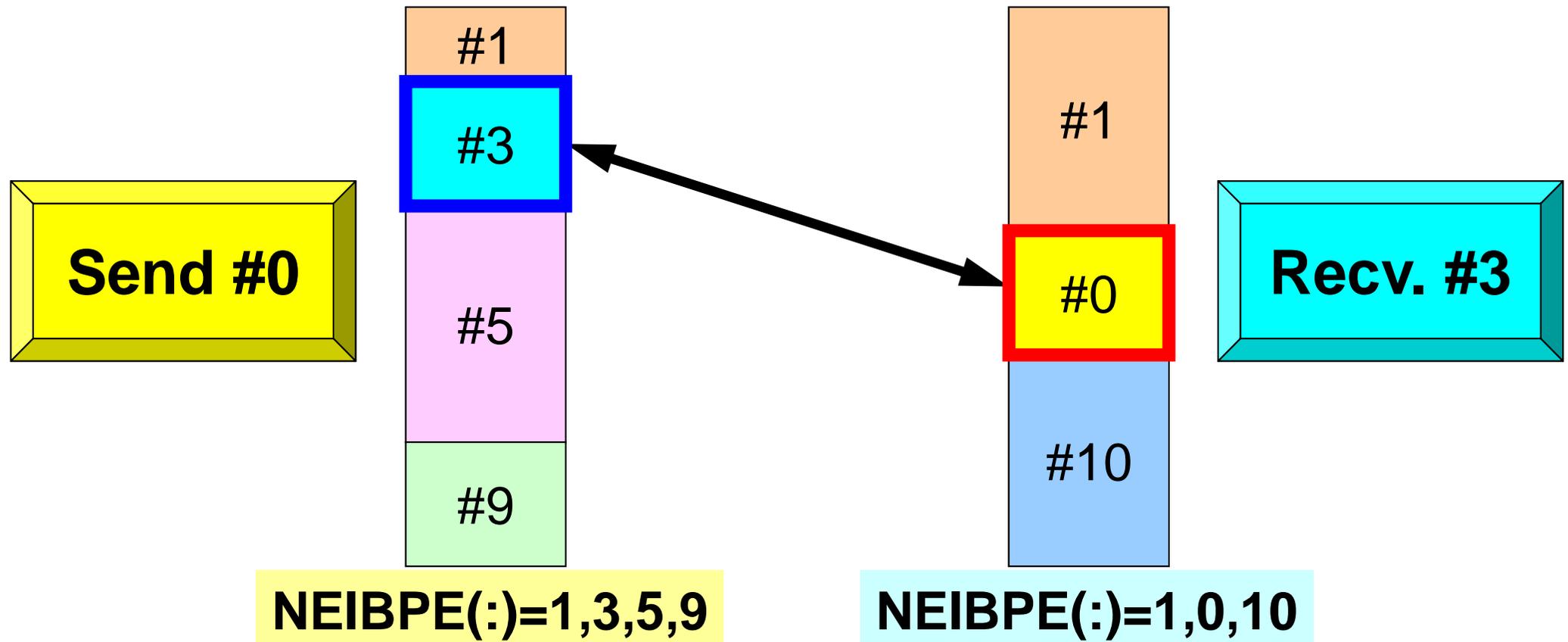
  call MPI_ISEND
&      (SENDbuf(iS_e), BUFlength_e, MPI_INTEGER, NEIBPE(neib), 0, &
&      MPI_COMM_WORLD, request_send(neib), ierr)
enddo
```

```
do neib= 1, NEIBPETOT
  iS_i= import_index(neib-1) + 1
  iE_i= import_index(neib )
  BUFlength_i= iE_i + 1 - iS_i

  call MPI_IRECV
&      (RECVbuf(iS_i), BUFlength_i, MPI_INTEGER, NEIBPE(neib), 0, &
&      MPI_COMM_WORLD, request_recv(neib), ierr)
enddo
```

- Consistency of ID's of sources/destinations, size and contents of messages !
- Communication occurs when NEIBPE(neib) matches

# Relationship SEND/RECV (#0 to #3)



- Consistency of ID's of sources/destinations, size and contents of messages !
- Communication occurs when NEIBPE(neib) matches

# Example: sq-sr1.c (5/6)

C

## RECV/Import: MPI\_Irecv

```

/*
!C
!C +-----+
!C | SEND-RECV |
!C +-----+
!C===*/

StatSend = malloc(sizeof(MPI_Status) * NeibPeTot);
StatRecv = malloc(sizeof(MPI_Status) * NeibPeTot);
RequestSend = malloc(sizeof(MPI_Request) * NeibPeTot);
RequestRecv = malloc(sizeof(MPI_Request) * NeibPeTot);

for(neib=0;neib<NeibPeTot;neib++){
    iStart = ExportIndex[neib];
    iEnd    = ExportIndex[neib+1];
    BufLength = iEnd - iStart;
    MPI_Isend(&SendBuf[iStart], BufLength, MPI_INT,
              NeibPe[neib], 0, MPI_COMM_WORLD, &RequestSend[neib]);
}

for(neib=0;neib<NeibPeTot;neib++){
    iStart = ImportIndex[neib];
    iEnd    = ImportIndex[neib+1];
    BufLength = iEnd - iStart;

    MPI_Irecv(&RecvBuf[iStart], BufLength, MPI_INT,
              NeibPe[neib], 0, MPI_COMM_WORLD, &RequestRecv[neib]);
}

```

PE#2				PE#3			
57	58	59	60	61	62	63	64
49	50	51	52	53	54	55	56
41	42	43	44	45	46	47	48
33	34	35	36	37	38	39	40
25	26	27	28	29	30	31	32
17	18	19	20	21	22	23	24
9	10	11	12	13	14	15	16
1	2	3	4	5	6	7	8

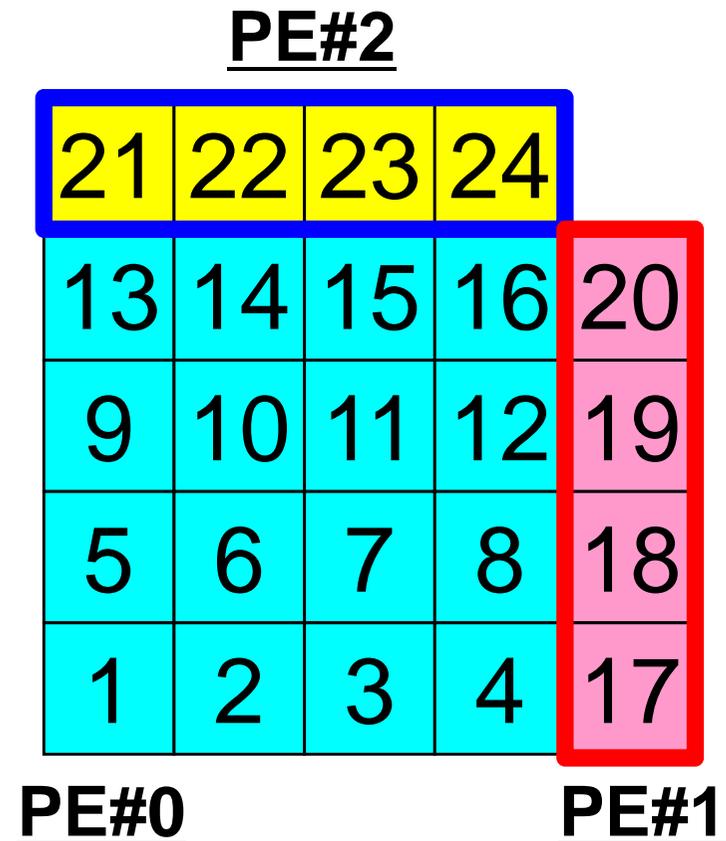
PE#0 PE#1

# RECV/Import: PE#0

```

#NEIBPEtot
2
#NEIBPE
1 2
#NODE
24 16
#IMPORTindex
4 8
#IMPORTitems
17
18
19
20
21
22
23
24
#EXPORTindex
4 8
#EXPORTitems
4
8
12
16
13
14
15
16

```



# RECV: MPI\_Irecv/Irecv/Waitall

C

```

for (neib=0; neib<NeibPETot; neib++){
  tag= 0;
  iS_i= import_index[neib];
  iE_i= import_index[neib+1];
  BUFlength_i= iE_i - iS_i

  ierr= MPI_Irecv
    (&RecvBuf[iS_i], BUFlength_i, MPI_DOUBLE, NeibPE[neib], 0,
     MPI_COMM_WORLD, &ReqRecv[neib])
}

```

```
MPI_Waitall(NeibPETot, ReqRecv, StatRecv);
```

```

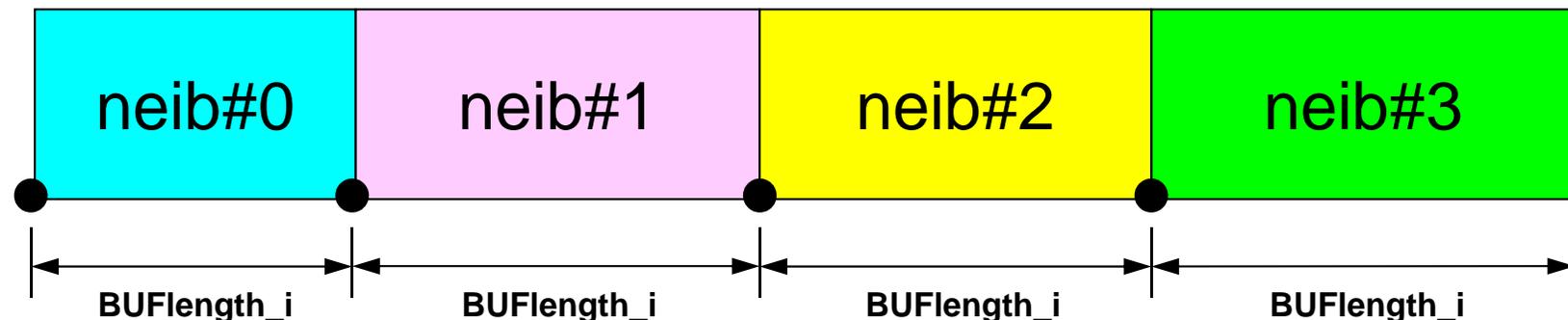
for (neib=0; neib<NeibPETot;neib++){
  for (k=import_index[neib];k<import_index[neib+1];k++){
    kk= import_item[k];
    VAL[kk]= RecvBuf[k];
  }
}

```

Copied from receiving buffer

import\_item (import\_index[neib]:import\_index[neib+1]-1) are received from neib-th neighbor

RecvBuf



import\_index[0]   import\_index[1]   import\_index[2]   import\_index[3]   import\_index[4]

# Example: sq-sr1.c (6/6)

C

## Reading info of ext pts from receiving buffers

```
MPI_Waitall(NeibPeTot, RequestRecv, StatRecv);
```

```
for(neib=0;neib<NeibPeTot;neib++){
    iStart = ImportIndex[neib];
    iEnd    = ImportIndex[neib+1];
    for(i=iStart;i<iEnd;i++){
        val[ImportItem[i]] = RecvBuf[i];
    }
}
```

Contents of RecvBuf are copied to values at external points.

```
MPI_Waitall(NeibPeTot, RequestSend, StatSend); /*
```

```
!C +-----+
!C | OUTPUT |
!C +-----+
!C===*/
```

```
for(neib=0;neib<NeibPeTot;neib++){
    iStart = ImportIndex[neib];
    iEnd    = ImportIndex[neib+1];
    for(i=iStart;i<iEnd;i++){
        int in = ImportItem[i];
        printf("RECVbuf%8d%8d%8d¥n", MyRank, NeibPe[neib], val[in]);
    }
}
MPI_Finalize();

return 0;
```

```
}
```

# Example: sq-sr1.c (6/6)

C

## Writing values at external points

```
MPI_Waitall(NeibPeTot, RequestRecv, StatRecv);

for(neib=0;neib<NeibPeTot;neib++){
    iStart = ImportIndex[neib];
    iEnd   = ImportIndex[neib+1];
    for(i=iStart;i<iEnd;i++){
        val[ImportItem[i]] = RecvBuf[i];
    }
}
MPI_Waitall(NeibPeTot, RequestSend, StatSend); /*
```

```
!C +-----+
!C |  OUTPUT  |
!C +-----+
!C===*/
```

```
    for(neib=0;neib<NeibPeTot;neib++){
        iStart = ImportIndex[neib];
        iEnd   = ImportIndex[neib+1];
        for(i=iStart;i<iEnd;i++){
            int in = ImportItem[i];
            printf("RECVbuf%8d%8d%8d\n", MyRank, NeibPe[neib], val[in]);
        }
    }
MPI_Finalize();

return 0;
}
```

# Results (PE#0)

**PE#2**

<u>57</u>	<u>58</u>	<u>59</u>	<u>60</u>
<u>49</u>	<u>50</u>	<u>51</u>	<u>52</u>
<u>41</u>	<u>42</u>	<u>43</u>	<u>44</u>
<u>33</u>	<u>34</u>	<u>35</u>	<u>36</u>

**PE#3**

<u>61</u>	<u>62</u>	<u>63</u>	<u>64</u>
<u>53</u>	<u>54</u>	<u>55</u>	<u>56</u>
<u>45</u>	<u>46</u>	<u>47</u>	<u>48</u>
<u>37</u>	<u>38</u>	<u>39</u>	<u>40</u>

<u>25</u>	<u>26</u>	<u>27</u>	<u>28</u>
<u>17</u>	<u>18</u>	<u>19</u>	<u>20</u>
<u>9</u>	<u>10</u>	<u>11</u>	<u>12</u>
<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>

<u>29</u>	<u>30</u>	<u>31</u>	<u>32</u>
<u>21</u>	<u>22</u>	<u>23</u>	<u>24</u>
<u>13</u>	<u>14</u>	<u>15</u>	<u>16</u>
<u>5</u>	<u>6</u>	<u>7</u>	<u>8</u>

**PE#0**

**PE#1**

RECVbuf	0	1	5
RECVbuf	0	1	13
RECVbuf	0	1	21
RECVbuf	0	1	29
RECVbuf	0	2	33
RECVbuf	0	2	34
RECVbuf	0	2	35
RECVbuf	0	2	36
RECVbuf	1	0	4
RECVbuf	1	0	12
RECVbuf	1	0	20
RECVbuf	1	0	28
RECVbuf	1	3	37
RECVbuf	1	3	38
RECVbuf	1	3	39
RECVbuf	1	3	40
RECVbuf	2	3	37
RECVbuf	2	3	45
RECVbuf	2	3	53
RECVbuf	2	3	61
RECVbuf	2	0	25
RECVbuf	2	0	26
RECVbuf	2	0	27
RECVbuf	2	0	28
RECVbuf	3	2	36
RECVbuf	3	2	44
RECVbuf	3	2	52
RECVbuf	3	2	60
RECVbuf	3	1	29
RECVbuf	3	1	30
RECVbuf	3	1	31
RECVbuf	3	1	32

# Results (PE#1)

**PE#2**

<u>57</u>	<u>58</u>	<u>59</u>	<u>60</u>
<u>49</u>	<u>50</u>	<u>51</u>	<u>52</u>
<u>41</u>	<u>42</u>	<u>43</u>	<u>44</u>
<u>33</u>	<u>34</u>	<u>35</u>	<u>36</u>

**PE#3**

<u>61</u>	<u>62</u>	<u>63</u>	<u>64</u>
<u>53</u>	<u>54</u>	<u>55</u>	<u>56</u>
<u>45</u>	<u>46</u>	<u>47</u>	<u>48</u>
<u>37</u>	<u>38</u>	<u>39</u>	<u>40</u>

<u>25</u>	<u>26</u>	<u>27</u>	<u>28</u>
<u>17</u>	<u>18</u>	<u>19</u>	<u>20</u>
<u>9</u>	<u>10</u>	<u>11</u>	<u>12</u>
<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>

<u>29</u>	<u>30</u>	<u>31</u>	<u>32</u>
<u>21</u>	<u>22</u>	<u>23</u>	<u>24</u>
<u>13</u>	<u>14</u>	<u>15</u>	<u>16</u>
<u>5</u>	<u>6</u>	<u>7</u>	<u>8</u>

**PE#0**

**PE#1**

RECVbuf	0	1	5
RECVbuf	0	1	13
RECVbuf	0	1	21
RECVbuf	0	1	29
RECVbuf	0	2	33
RECVbuf	0	2	34
RECVbuf	0	2	35
RECVbuf	0	2	36
<b>RECVbuf</b>	<b>1</b>	<b>0</b>	<b>4</b>
<b>RECVbuf</b>	<b>1</b>	<b>0</b>	<b>12</b>
<b>RECVbuf</b>	<b>1</b>	<b>0</b>	<b>20</b>
<b>RECVbuf</b>	<b>1</b>	<b>0</b>	<b>28</b>
<b>RECVbuf</b>	<b>1</b>	<b>3</b>	<b>37</b>
<b>RECVbuf</b>	<b>1</b>	<b>3</b>	<b>38</b>
<b>RECVbuf</b>	<b>1</b>	<b>3</b>	<b>39</b>
<b>RECVbuf</b>	<b>1</b>	<b>3</b>	<b>40</b>
RECVbuf	2	3	37
RECVbuf	2	3	45
RECVbuf	2	3	53
RECVbuf	2	3	61
RECVbuf	2	0	25
RECVbuf	2	0	26
RECVbuf	2	0	27
RECVbuf	2	0	28
RECVbuf	3	2	36
RECVbuf	3	2	44
RECVbuf	3	2	52
RECVbuf	3	2	60
RECVbuf	3	1	29
RECVbuf	3	1	30
RECVbuf	3	1	31
RECVbuf	3	1	32

# Results (PE#2)

**PE#2**

<u>57</u>	<u>58</u>	<u>59</u>	<u>60</u>
<u>49</u>	<u>50</u>	<u>51</u>	<u>52</u>
<u>41</u>	<u>42</u>	<u>43</u>	<u>44</u>
<u>33</u>	<u>34</u>	<u>35</u>	<u>36</u>

**PE#3**

<u>61</u>	<u>62</u>	<u>63</u>	<u>64</u>
<u>53</u>	<u>54</u>	<u>55</u>	<u>56</u>
<u>45</u>	<u>46</u>	<u>47</u>	<u>48</u>
<u>37</u>	<u>38</u>	<u>39</u>	<u>40</u>

<u>25</u>	<u>26</u>	<u>27</u>	<u>28</u>
<u>17</u>	<u>18</u>	<u>19</u>	<u>20</u>
<u>9</u>	<u>10</u>	<u>11</u>	<u>12</u>
<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>

<u>29</u>	<u>30</u>	<u>31</u>	<u>32</u>
<u>21</u>	<u>22</u>	<u>23</u>	<u>24</u>
<u>13</u>	<u>14</u>	<u>15</u>	<u>16</u>
<u>5</u>	<u>6</u>	<u>7</u>	<u>8</u>

**PE#0**

**PE#1**

RECVbuf	0	1	5
RECVbuf	0	1	13
RECVbuf	0	1	21
RECVbuf	0	1	29
RECVbuf	0	2	33
RECVbuf	0	2	34
RECVbuf	0	2	35
RECVbuf	0	2	36
RECVbuf	1	0	4
RECVbuf	1	0	12
RECVbuf	1	0	20
RECVbuf	1	0	28
RECVbuf	1	3	37
RECVbuf	1	3	38
RECVbuf	1	3	39
RECVbuf	1	3	40
<b>RECVbuf</b>	<b>2</b>	<b>3</b>	<b>37</b>
<b>RECVbuf</b>	<b>2</b>	<b>3</b>	<b>45</b>
<b>RECVbuf</b>	<b>2</b>	<b>3</b>	<b>53</b>
<b>RECVbuf</b>	<b>2</b>	<b>3</b>	<b>61</b>
<b>RECVbuf</b>	<b>2</b>	<b>0</b>	<b>25</b>
<b>RECVbuf</b>	<b>2</b>	<b>0</b>	<b>26</b>
<b>RECVbuf</b>	<b>2</b>	<b>0</b>	<b>27</b>
<b>RECVbuf</b>	<b>2</b>	<b>0</b>	<b>28</b>
RECVbuf	3	2	36
RECVbuf	3	2	44
RECVbuf	3	2	52
RECVbuf	3	2	60
RECVbuf	3	1	29
RECVbuf	3	1	30
RECVbuf	3	1	31
RECVbuf	3	1	32

# Results (PE#3)

**PE#2**

<u>57</u>	<u>58</u>	<u>59</u>	<u>60</u>
<u>49</u>	<u>50</u>	<u>51</u>	<u>52</u>
<u>41</u>	<u>42</u>	<u>43</u>	<u>44</u>
<u>33</u>	<u>34</u>	<u>35</u>	<u>36</u>

**PE#3**

<u>61</u>	<u>62</u>	<u>63</u>	<u>64</u>
<u>53</u>	<u>54</u>	<u>55</u>	<u>56</u>
<u>45</u>	<u>46</u>	<u>47</u>	<u>48</u>
<u>37</u>	<u>38</u>	<u>39</u>	<u>40</u>

<u>25</u>	<u>26</u>	<u>27</u>	<u>28</u>
<u>17</u>	<u>18</u>	<u>19</u>	<u>20</u>
<u>9</u>	<u>10</u>	<u>11</u>	<u>12</u>
<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>

<u>29</u>	<u>30</u>	<u>31</u>	<u>32</u>
<u>21</u>	<u>22</u>	<u>23</u>	<u>24</u>
<u>13</u>	<u>14</u>	<u>15</u>	<u>16</u>
<u>5</u>	<u>6</u>	<u>7</u>	<u>8</u>

**PE#0**

**PE#1**

RECVbuf	0	1	5
RECVbuf	0	1	13
RECVbuf	0	1	21
RECVbuf	0	1	29
RECVbuf	0	2	33
RECVbuf	0	2	34
RECVbuf	0	2	35
RECVbuf	0	2	36
RECVbuf	1	0	4
RECVbuf	1	0	12
RECVbuf	1	0	20
RECVbuf	1	0	28
RECVbuf	1	3	37
RECVbuf	1	3	38
RECVbuf	1	3	39
RECVbuf	1	3	40
RECVbuf	2	3	37
RECVbuf	2	3	45
RECVbuf	2	3	53
RECVbuf	2	3	61
RECVbuf	2	0	25
RECVbuf	2	0	26
RECVbuf	2	0	27
RECVbuf	2	0	28
<b>RECVbuf</b>	<b>3</b>	<b>2</b>	<b>36</b>
<b>RECVbuf</b>	<b>3</b>	<b>2</b>	<b>44</b>
<b>RECVbuf</b>	<b>3</b>	<b>2</b>	<b>52</b>
<b>RECVbuf</b>	<b>3</b>	<b>2</b>	<b>60</b>
<b>RECVbuf</b>	<b>3</b>	<b>1</b>	<b>29</b>
<b>RECVbuf</b>	<b>3</b>	<b>1</b>	<b>30</b>
<b>RECVbuf</b>	<b>3</b>	<b>1</b>	<b>31</b>
<b>RECVbuf</b>	<b>3</b>	<b>1</b>	<b>32</b>

# Distributed Local Data Structure for Parallel Computation

- Distributed local data structure for domain-to-domain communications has been introduced, which is appropriate for such applications with sparse coefficient matrices (e.g. FDM, FEM, FVM etc.).
  - SPMD
  - Local Numbering: Internal pts to External pts
  - Generalized communication table
- Everything is easy, if proper data structure is defined:
  - Values at boundary pts are copied into sending buffers
  - Send/Recv
  - Values at external pts are updated through receiving buffers

# Initial Mesh

**t2**

<u>21</u>	<u>22</u>	<u>23</u>	<u>24</u>	<u>25</u>
<u>16</u>	<u>17</u>	<u>18</u>	<u>19</u>	<u>20</u>
<u>11</u>	<u>12</u>	<u>13</u>	<u>14</u>	<u>15</u>
<u>6</u>	<u>7</u>	<u>8</u>	<u>9</u>	<u>10</u>
<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>	<u>5</u>

# Three Domains

**t2**

**#PE2**

<u>21</u>	<u>22</u>	<u>23</u>	<u>24</u>
<u>16</u>	<u>17</u>	<u>18</u>	<u>19</u>
<u>11</u>	<u>12</u>	<u>13</u>	<u>14</u>
<u>6</u>	<u>7</u>	<u>8</u>	

**#PE1**

<u>23</u>	<u>24</u>	<u>25</u>
<u>18</u>	<u>19</u>	<u>20</u>
<u>13</u>	<u>14</u>	<u>15</u>
<u>8</u>	<u>9</u>	<u>10</u>
	<u>4</u>	<u>5</u>

**#PE0**

<u>11</u>	<u>12</u>	<u>13</u>		
<u>6</u>	<u>7</u>	<u>8</u>	<u>9</u>	<u>10</u>
<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>	<u>5</u>

# Three Domains

t2

#PE2

<b>7</b> <u>21</u>	<b>8</b> <u>22</u>	<b>9</b> <u>23</u>	<b>15</b> <u>24</u>
<b>4</b> <u>16</u>	<b>5</b> <u>17</u>	<b>6</b> <u>18</u>	<b>14</b> <u>19</u>
<b>1</b> <u>11</u>	<b>2</b> <u>12</u>	<b>3</b> <u>13</u>	<b>13</b> <u>14</u>
<b>10</b> <u>6</u>	<b>11</b> <u>7</u>	<b>12</b> <u>8</u>	

#PE1

<b>14</b> <u>23</u>	<b>7</b> <u>24</u>	<b>8</b> <u>25</u>
<b>13</b> <u>18</u>	<b>5</b> <u>19</u>	<b>6</b> <u>20</u>
<b>12</b> <u>13</u>	<b>3</b> <u>14</u>	<b>4</b> <u>15</u>
<b>11</b> <u>8</u>	<b>1</b> <u>9</u>	<b>2</b> <u>10</u>
	<b>9</b> <u>4</u>	<b>10</b> <u>5</u>

#PE0

<b>11</b> <u>11</u>	<b>12</b> <u>12</u>	<b>13</b> <u>13</u>			
<b>6</b> <u>6</u>	<b>7</b> <u>7</u>	<b>8</b> <u>8</u>	<b>9</b> <u>9</u>	<b>10</b> <u>10</u>	
<b>1</b> <u>1</u>	<b>2</b> <u>2</u>	<b>3</b> <u>3</u>	<b>4</b> <u>4</u>	<b>5</b> <u>5</u>	

# PE#0: sqm.0: fill O's

#PE2

<b>7</b> 21	<b>8</b> 22	<b>9</b> 23	<b>15</b> 24
<b>4</b> 16	<b>5</b> 17	<b>6</b> 18	<b>14</b> 19
<b>1</b> 11	<b>2</b> 12	<b>3</b> 13	<b>13</b> 14
<b>10</b> 6	<b>11</b> 7	<b>12</b> 8	

#PE1

<b>14</b> 23	<b>7</b> 24	<b>8</b> 25
<b>13</b> 18	<b>5</b> 19	<b>6</b> 20
<b>12</b> 13	<b>3</b> 14	<b>4</b> 15
<b>11</b> 8	<b>1</b> 9	<b>2</b> 10
	<b>9</b> 4	<b>10</b> 5

#PE0

<b>11</b> 11	<b>12</b> 12	<b>13</b> 13		
<b>6</b> 6	<b>7</b> 7	<b>8</b> 8	<b>9</b> 9	<b>10</b> 10
<b>1</b> 1	<b>2</b> 2	<b>3</b> 3	<b>4</b> 4	<b>5</b> 5

```

#NEIBPEtot
    2
#NEIBPE
    1    2
#NODE
    13    8 (int+ext, int pts)
#IMPORTindex
    ○    ○
#IMPORTitems
    ○...
#EXPORTindex
    ○    ○
#EXPORTitems
    ○...

```

# PE#1: sqm.1: fill O's

#PE2

<b>7</b> 21	<b>8</b> 22	<b>9</b> 23	<b>15</b> 24
<b>4</b> 16	<b>5</b> 17	<b>6</b> 18	<b>14</b> 19
<b>1</b> 11	<b>2</b> 12	<b>3</b> 13	<b>13</b> 14
<b>10</b> 6	<b>11</b> 7	<b>12</b> 8	

#PE1

<b>14</b> 23	<b>7</b> 24	<b>8</b> 25
<b>13</b> 18	<b>5</b> 19	<b>6</b> 20
<b>12</b> 13	<b>3</b> 14	<b>4</b> 15
<b>11</b> 8	<b>1</b> 9	<b>2</b> 10
	<b>9</b> 4	<b>10</b> 5

#PE0

<b>11</b> 11	<b>12</b> 12	<b>13</b> 13		
<b>6</b> 6	<b>7</b> 7	<b>8</b> 8	<b>9</b> 9	<b>10</b> 10
<b>1</b> 1	<b>2</b> 2	<b>3</b> 3	<b>4</b> 4	<b>5</b> 5

```

#NEIBPEtot
    2
#NEIBPE
    0    2
#NODE
    14    8 (int+ext, int pts)
#IMPORTindex
    ○    ○
#IMPORTitems
    ○...
#EXPORTindex
    ○    ○
#EXPORTitems
    ○...

```

# PE#2: sqm.2: fill O's

#PE2

<b>7</b> 21	<b>8</b> 22	<b>9</b> 23	<b>15</b> 24
<b>4</b> 16	<b>5</b> 17	<b>6</b> 18	<b>14</b> 19
<b>1</b> 11	<b>2</b> 12	<b>3</b> 13	<b>13</b> 14
<b>10</b> 6	<b>11</b> 7	<b>12</b> 8	

#PE1

<b>14</b> 23	<b>7</b> 24	<b>8</b> 25
<b>13</b> 18	<b>5</b> 19	<b>6</b> 20
<b>12</b> 13	<b>3</b> 14	<b>4</b> 15
<b>11</b> 8	<b>1</b> 9	<b>2</b> 10
	<b>9</b> 4	<b>10</b> 5

#PE0

<b>11</b> 11	<b>12</b> 12	<b>13</b> 13		
<b>6</b> 6	<b>7</b> 7	<b>8</b> 8	<b>9</b> 9	<b>10</b> 10
<b>1</b> 1	<b>2</b> 2	<b>3</b> 3	<b>4</b> 4	<b>5</b> 5

```
#NEIBPEtot
    2
#NEIBPE
    1    0
#NODE
    15    9 (int+ext, int pts)
#IMPORTindex
    ○    ○
#IMPORTitems
    ○...
#EXPORTindex
    ○    ○
#EXPORTitems
    ○...
```

t2

#PE2

<b>7</b> <u>21</u>	<b>8</b> <u>22</u>	<b>9</b> <u>23</u>	<b>15</b> <u>24</u>
<b>4</b> <u>16</u>	<b>5</b> <u>17</u>	<b>6</b> <u>18</u>	<b>14</b> <u>19</u>
<b>1</b> <u>11</u>	<b>2</b> <u>12</u>	<b>3</b> <u>13</u>	<b>13</b> <u>14</u>
<b>10</b> <u>6</u>	<b>11</b> <u>7</u>	<b>12</b> <u>8</u>	

#PE1

<b>14</b> <u>23</u>	<b>7</b> <u>24</u>	<b>8</b> <u>25</u>
<b>13</b> <u>18</u>	<b>5</b> <u>19</u>	<b>6</b> <u>20</u>
<b>12</b> <u>13</u>	<b>3</b> <u>14</u>	<b>4</b> <u>15</u>
<b>11</b> <u>8</u>	<b>1</b> <u>9</u>	<b>2</b> <u>10</u>
	<b>9</b> <u>4</u>	<b>10</b> <u>5</u>

#PE0

<b>11</b> <u>11</u>	<b>12</b> <u>12</u>	<b>13</b> <u>13</u>		
<b>6</b> <u>6</u>	<b>7</b> <u>7</u>	<b>8</b> <u>8</u>	<b>9</b> <u>9</u>	<b>10</b> <u>10</u>
<b>1</b> <u>1</u>	<b>2</b> <u>2</u>	<b>3</b> <u>3</u>	<b>4</b> <u>4</u>	<b>5</b> <u>5</u>

# Procedures

- Number of Internal/External Points
- Where do External Pts come from ?
  - **IMPORTindex, IMPORTitems**
  - Sequence of **NEIBPE**
- Then check destinations of Boundary Pts.
  - **EXPORTindex, EXPORTitems**
  - Sequence of **NEIBPE**
- “sq.\*” are in <\$O-S2>/ex
- Create “sqm.\*” by yourself
- copy <\$O-S2>/a.out (by sq-sr1.f) to <\$O-S2>/ex
- pjsub go3.sh

# Report S2 (1/2)

- Parallelize 1D code (1d.c) using MPI
- Read entire element number, and decompose into sub-domains in your program
- **Validate the results**
  - Answer of Original Code = Answer of Parallel Code
  - Explain why number of iterations does not change, as number of MPI processes changes.
- **Measure parallel performance**

# Report S2 (2/2)

- **Deadline: January 28<sup>th</sup> (Tue), 2020, 17:00**
  - Send files via e-mail at **`nakajima (at) cc.u-tokyo.ac.jp`**
- **Problem**
  - Apply “Generalized Communication Table”
  - Read entire elem. #, decompose into sub-domains in your program
  - **Evaluate parallel performance**
    - You need huge number of elements, to get excellent performance.
    - Fix number of iterations (e.g. 100), if computations cannot be completed.
- **Report**
  - Cover Page: Name, ID, and Problem ID (S2) must be written.
  - Less than eight pages including figures and tables (A4).
    - Strategy, Structure of the Program, Remarks
  - Source list of the program (if you have bugs)
  - Output list (as small as possible)

# Options for Optimization

## General Option (-O3)

```
$ mpiifort -O3 test.f
```

```
$ mpiicc -O3 test.c
```

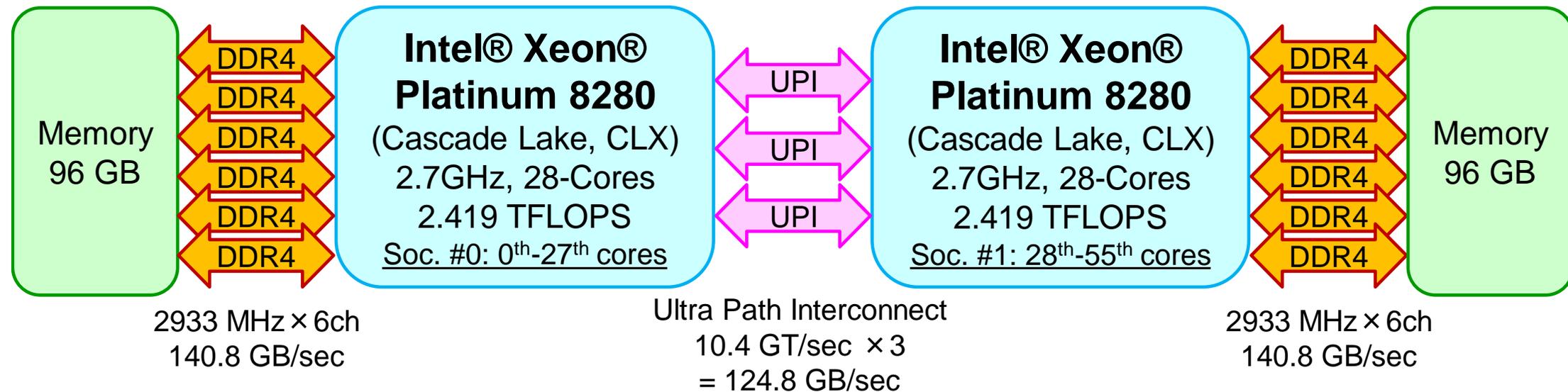
## Special Options for AVX512, NOT Necessarily Fast ...

```
$ mpiifort -align array64byte -O3 -axCORE-AVX512 test.f
```

```
$ mpiicc -align -O3 -axCORE-AVX512 test.c
```

# Process Number

#PJM -L node=1; #PJM --mpi proc= 1	1-node, 1-proc, 1-proc/n
#PJM -L node=1; #PJM --mpi proc= 4	1-node, 4-proc, 4-proc/n
#PJM -L node=1; #PJM --mpi proc=16	1-node, 16-proc, 16-proc/n
#PJM -L node=1; #PJM --mpi proc=28	1-node, 28-proc, 28-proc/n
#PJM -L node=1; #PJM --mpi proc=56	1-node, 56-proc, 56-proc/n
#PJM -L node=4; #PJM --mpi proc=128	4-node, 128-proc, 32-proc/n
#PJM -L node=8; #PJM --mpi proc=256	8-node, 256-proc, 32-proc/n
#PJM -L node=8; #PJM --mpi proc=448	8-node, 448-proc, 56-proc/n

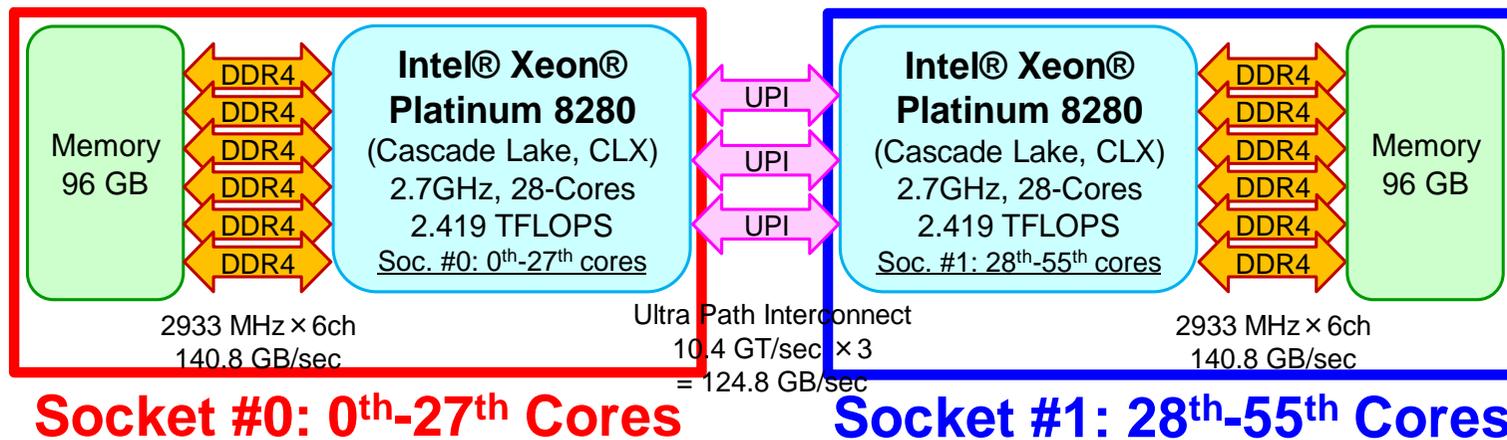


# <\$O-S1>/a01.sh: Use 1-core (0<sup>th</sup>)

```
#!/bin/sh
#PJM -N "test"
#PJM -L rscgrp=lecture
#PJM -L node=1
#PJM --mpi proc=1
#PJM -L elapse=00:15:00
#PJM -g gt37
#PJM -j
#PJM -e err
#PJM -o test.lst
```

```
export I_MPI_PIN_PROCESSOR_LIST=0
```

```
mpiexec.hydra -n ${PJM_MPI_PROC} ./a.out
```

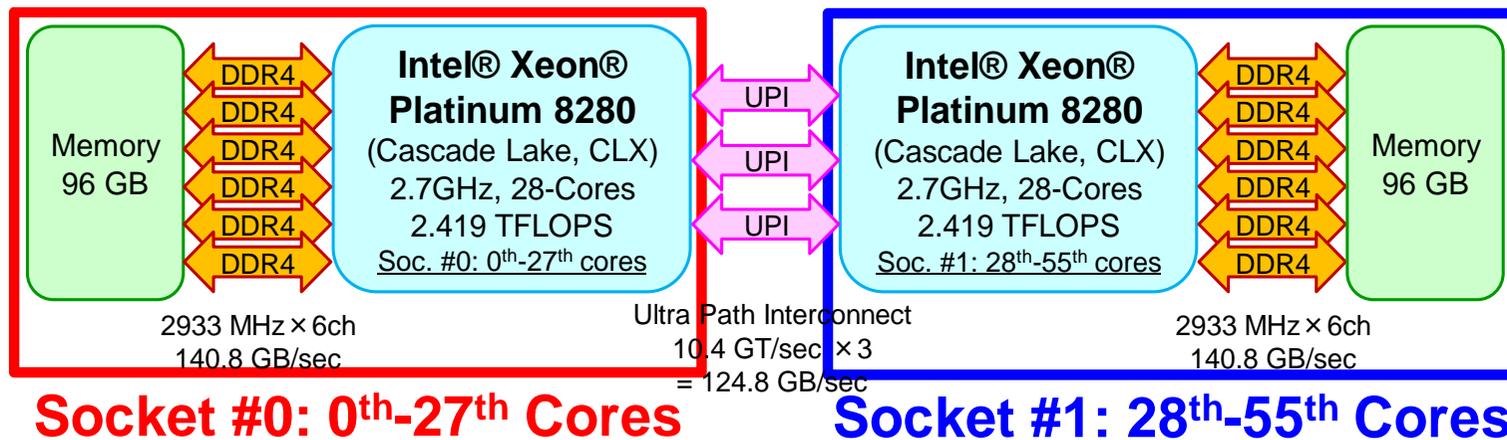


# <\$O-S1>/a24.sh: Use 24-cores (0<sup>th</sup>-23<sup>rd</sup>)

```
#!/bin/sh
#PJM -N "test"
#PJM -L rscgrp=lecture
#PJM -L node=1
#PJM --mpi proc=24
#PJM -L elapse=00:15:00
#PJM -g gt37
#PJM -j
#PJM -e err
#PJM -o test.lst
```

```
export I_MPI_PIN_PROCESSOR_LIST=0-23
```

```
mpiexec.hydra -n ${PJM_MPI_PROC} ./a.out
```

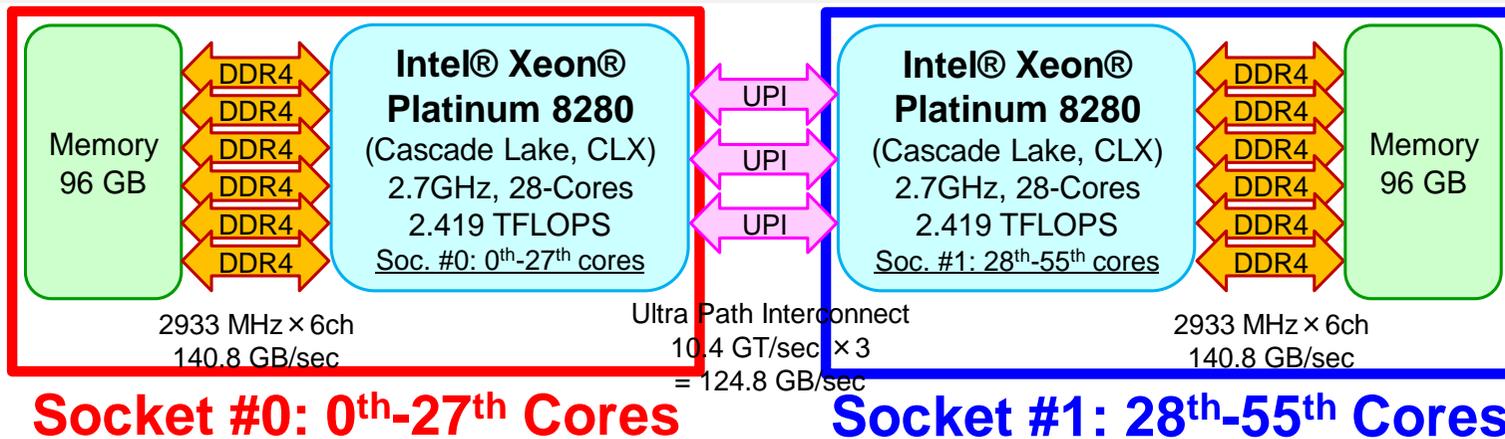


# <\$O-S1>/a48.sh: Use 48-cores (0<sup>th</sup>-23<sup>rd</sup>, 28<sup>th</sup>-51<sup>st</sup>)

```
#!/bin/sh
#PJM -N "test"
#PJM -L rscgrp=lecture
#PJM -L node=1
#PJM --mpi proc=48
#PJM -L elapse=00:15:00
#PJM -g gt37
#PJM -j
#PJM -e err
#PJM -o test.lst
```

```
export I_MPI_PIN_PROCESSOR_LIST=0-23,28-51
```

```
mpiexec.hydra -n ${PJM_MPI_PROC} ./a.out
```



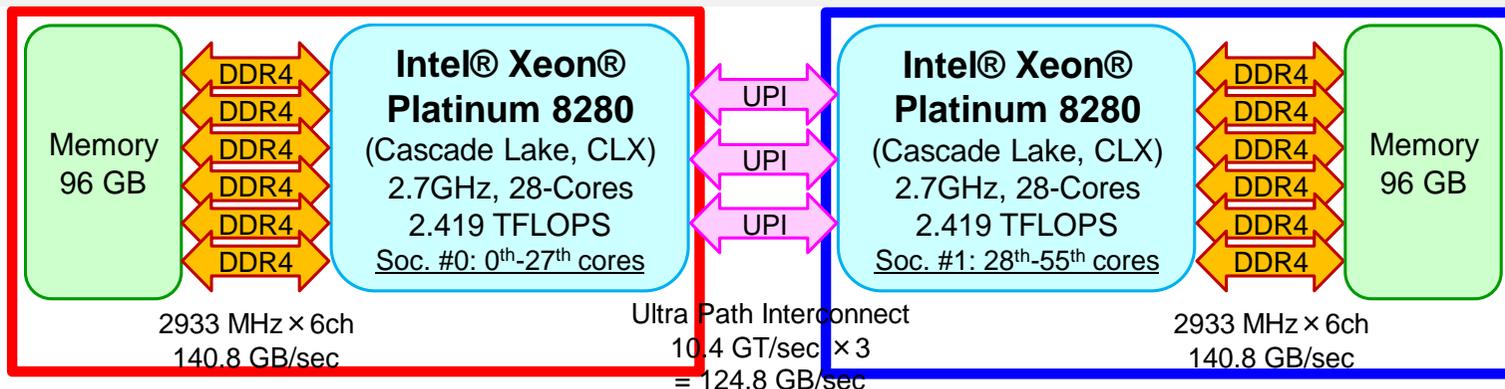
# <\$O-S1>/b48.sh: Use 8x48-cores (0<sup>th</sup>-23<sup>rd</sup>, 28<sup>th</sup>-51<sup>st</sup>)

```
#!/bin/sh
#PJM -N "test"
#PJM -L rscgrp=lecture
#PJM -L node=8
#PJM --mpi proc=384
#PJM -L elapse=00:15:00
#PJM -g gt37
#PJM -j
#PJM -e err
#PJM -o test.lst
```

$384 \div 8 = 48\text{-cores/node}$

```
export I_MPI_PIN_PROCESSOR_LIST=0-23,28-51
```

```
mpiexec.hydra -n ${PJM_MPI_PROC} ./a.out
```



**Socket #0: 0<sup>th</sup>-27<sup>th</sup> Cores**

**Socket #1: 28<sup>th</sup>-55<sup>th</sup> Cores**

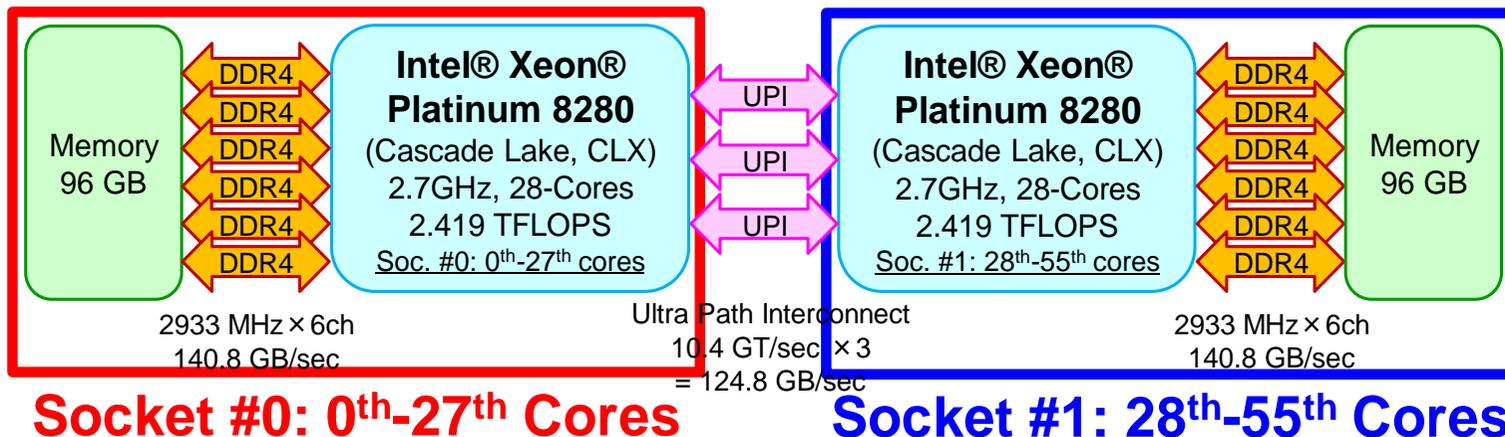
# <\$O-S1>/b48b.sh: Use 8x48-cores (0<sup>th</sup>-23<sup>rd</sup>, 28<sup>th</sup>-51<sup>st</sup>)

```
#!/bin/sh
#PJM -N "test"
#PJM -L rscgrp=lecture
#PJM -L node=8
#PJM --mpi proc=384
#PJM -L elapse=00:15:00
#PJM -g gt37
#PJM -j
#PJM -e err
#PJM -o test.lst
```

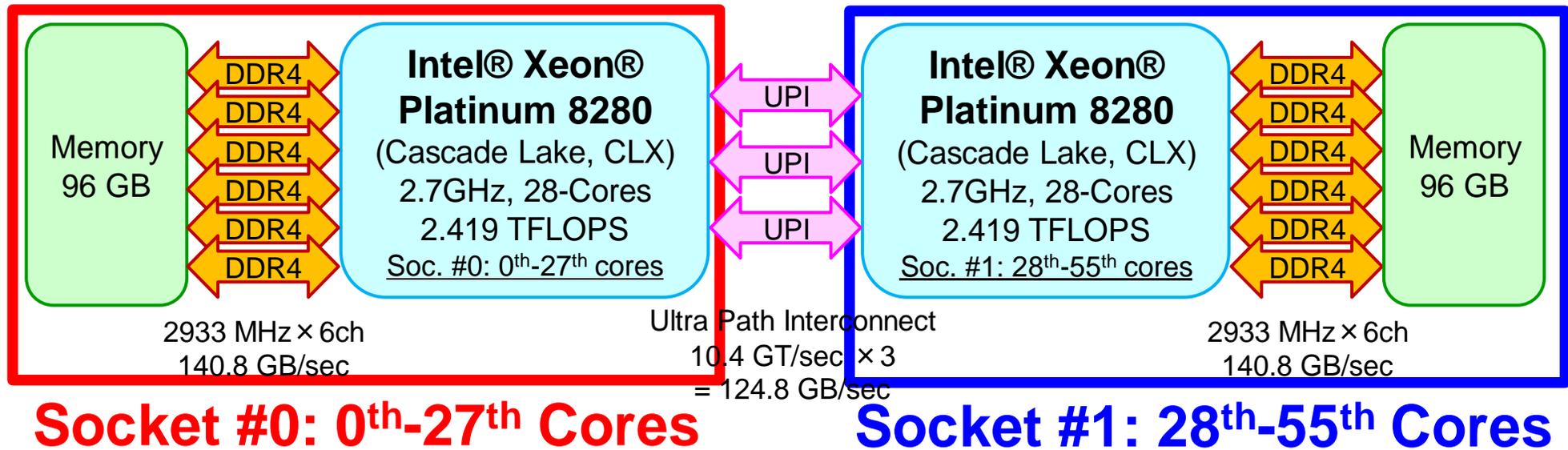
$$384 \div 8 = 48\text{-cores/node}$$

```
export I_MPI_PIN_PROCESSOR_LIST=0-23,28-51
```

```
mpiexec.hydra -n ${PJM_MPI_PROC} numactl -l ./a.out
```



# NUMA Architecture



- Each Node of Oakbridge-CX (OBCX)
  - 2 Sockets (CPU's) of Intel CLX
  - Each socket has 28 cores
- Each core of a socket can access to the memory on the other socket : NUMA (Non-Uniform Memory Access)
  - `numactl -l` : local memory to be used
  - Sometimes (not always), more stable with this option

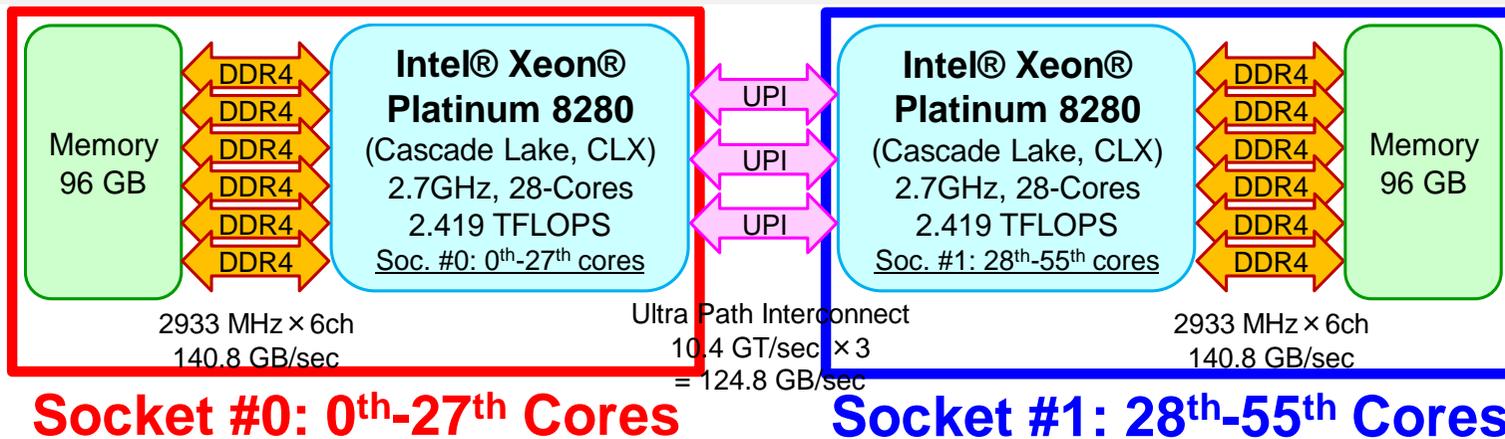
# Use 8x48-cores, 48-cores are randomly selected from 56-cores

```
#!/bin/sh
#PJM -N "test"
#PJM -L rscgrp=lecture
#PJM -L node=8
#PJM --mpi proc=384
#PJM -L elapse=00:15:00
#PJM -g gt37
#PJM -j
#PJM -e err
#PJM -o test.lst
```

$$384 \div 8 = 48\text{-cores/node}$$

```
export I_MPI_PIN_PROCESSOR_LIST=0-23,28-51
```

```
mpiexec.hydra -n ${PJM_MPI_PROC} ./a.out
```



# Use 8x56-cores

```
#!/bin/sh
#PJM -N "test"
#PJM -L rscgrp=lecture
#PJM -L node=8
#PJM --mpi proc=448
#PJM -L elapse=00:15:00
#PJM -g gt37
#PJM -j
#PJM -e err
#PJM -o test.lst
```

$$448 \div 8 = 56\text{-cores/node}$$

```
mpiexec.hydra -n ${PJM_MPI_PROC} ./a.out
```

