

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

Kengo Nakajima

Technical & Scientific Computing II (4820-1028)

Seminar on Computer Science II (4810-1205)

Parallel FEM

# Target of this Class

- Large-scale parallel computer enables fast computing in large-scale scientific simulations with detailed models. Computational science develops new frontiers of science and engineering.
- Why parallel computing ?
  - faster
  - larger
  - “larger” is more important from the view point of “new frontiers of science & engineering”, but “faster” is also important.
  - + more complicated
  - Ideal: Scalable
    - Solving  $N^x$  scale problem using  $N^x$  computational resources during same computation time.

# Overview

- What is MPI ?
- Your First MPI Program: Hello World
- Global/Local Data
- Collective Communication
- Peer-to-Peer Communication

# What is MPI ? (1/2)

- Message Passing Interface
- “Specification” of message passing API for distributed memory environment
  - Not a program, Not a library
    - <http://phase.hpcc.jp/phase/mpi-j/ml/mpi-j-html/contents.html>
- History
  - 1992 MPI Forum
  - 1994 MPI-1
  - 1997 MPI-2, MPI-3 is now available
- Implementation
  - mpich ANL (Argonne National Laboratory)
  - OpenMPI, MVAPICH
  - H/W vendors
  - C/C++, FOTRAN, Java ; Unix, Linux, Windows, Mac OS

# What is MPI ? (2/2)

- “mpich” (free) is widely used
  - supports MPI-2 spec. (partially)
  - MPICH2 after Nov. 2005.
  - <http://www-unix.mcs.anl.gov/mpi/>
- Why MPI is widely used as *de facto standard* ?
  - Uniform interface through MPI forum
    - Portable, can work on any types of computers
    - Can be called from Fortran, C, etc.
  - mpich
    - free, supports every architecture
- PVM (Parallel Virtual Machine) was also proposed in early 90's but not so widely used as MPI

# References

- W.Gropp et al., Using MPI second edition, MIT Press, 1999.
- M.J.Quinn, Parallel Programming in C with MPI and OpenMP, McGrawhill, 2003.
- W.Gropp et al., MPI: The Complete Reference Vol.I, II, MIT Press, 1998.
- <http://www-unix.mcs.anl.gov/mpi/www/>
  - API (Application Interface) of MPI

# How to learn MPI (1/2)

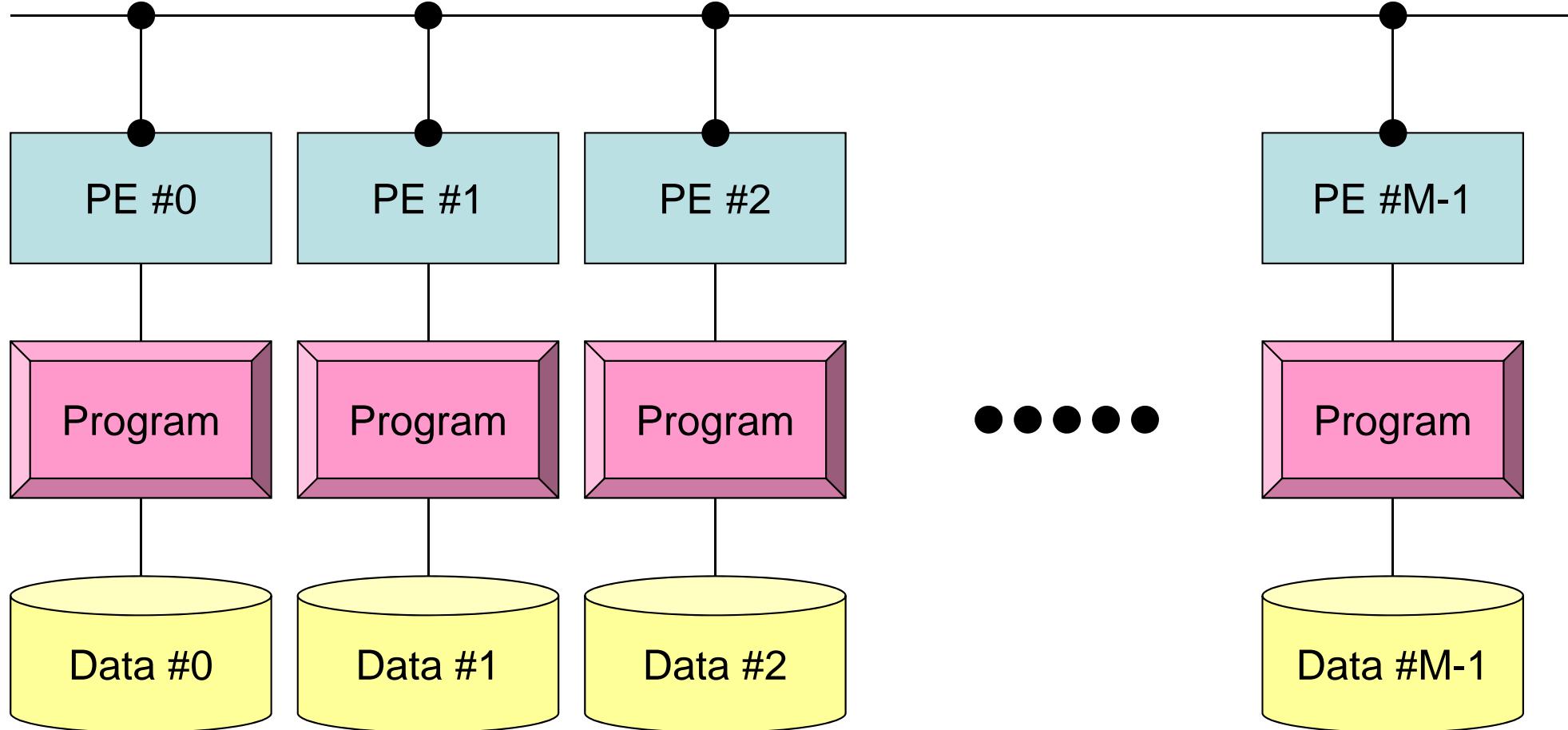
- Grammar
  - 10-20 functions of MPI-1 will be taught in the class
    - although there are many convenient capabilities in MPI-2
  - If you need further information, you can find information from web, books, and MPI experts.
- Practice is important
  - Programming
  - “Running the codes” is the most important
- Be familiar with or “grab” the idea of SPMD/SIMD op’s
  - Single Program/Instruction Multiple Data
  - Each process does same operation for different data
    - Large-scale data is decomposed, and each part is computed by each process
  - Global/Local Data, Global/Local Numbering

PE: Processing Element  
Processor, Domain, Process

# SPMD

You understand 90% MPI, if you understand this figure.

```
mpirun -np M <Program>
```



Each process does same operation for different data

Large-scale data is decomposed, and each part is computed by each process

It is ideal that parallel program is not different from serial one except communication.

# Some Technical Terms

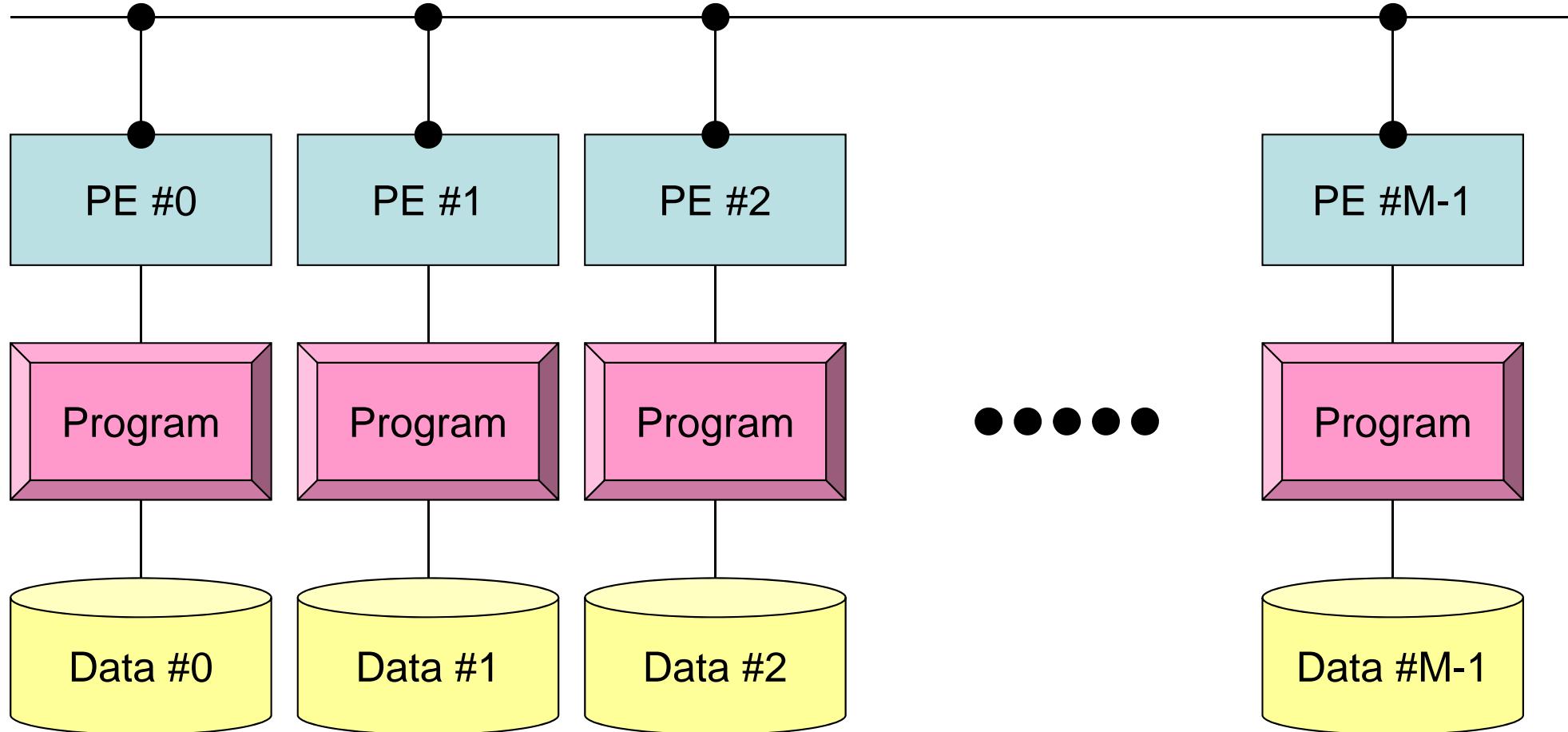
- Processor, Core
  - Processing Unit (H/W), Processor=Core for single-core proc's
- Process
  - Unit for MPI computation, nearly equal to "core"
  - Each core (or processor) can host multiple processes (but not efficient)
- PE (Processing Element)
  - PE originally mean "processor", but it is sometimes used as "process" in this class. Moreover it means "domain" (next)
    - In multicore proc's: PE generally means "core"
- Domain
  - domain=process (=PE), each of "MD" in "SPMD", each data set
- **Process ID of MPI (ID of PE, ID of domain) starts from "0"**
  - if you have 8 processes (PE's, domains), ID is 0~7

PE: Processing Element  
Processor, Domain, Process

# SPMD

You understand 90% MPI, if you understand this figure.

```
mpirun -np M <Program>
```



Each process does same operation for different data

Large-scale data is decomposed, and each part is computed by each process

It is ideal that parallel program is not different from serial one except communication.

# How to learn MPI (2/2)

- NOT so difficult.
- Therefore, 2-3 lectures are enough for just learning grammar of MPI.
- Grab the idea of SPMD !

# Schedule

- MPI
  - Basic Functions
  - Collective Communication
  - Point-to-Point (or Peer-to-Peer) Communication
- 90 min. x 4-5 lectures
  - Collective Communication
    - Report S1
  - Point-to-Point/Peer-to-Peer Communication
    - Report S2: Parallelization of 1D code
  - At this point, you are almost an expert of MPI programming.

- What is MPI ?
- **Your First MPI Program: Hello World**
- Global/Local Data
- Collective Communication
- Peer-to-Peer Communication

# Login to Oakleaf-FX

```
ssh t74**@oakleaf-fx.cc.u-tokyo.ac.jp
```

## Create directory

```
>$ cd  
>$ mkdir fem2 (your favorite name)  
>$ cd fem2
```

In this class this top-directory is called **<\$O-fem2>**.  
Files are copied to this directory.

Under this directory, s1, s2, s1-ref are created:

**<\$O-S1> = <\$O-fem2>/mpi/S1**  
**<\$O-S2> = <\$O-fem2>/mpi/S2**

Oakleaf-FX

ECCS2012

# Copying files on Oakleaf-FX

## Copy

```
>$ cd <$O-fem2>
>$ cp /home/z30088/fem2/C/s1.tar .
>$ tar xvf s1.tar
```

## Confirm directory

```
>$ ls
mpi

>$ cd mpi/S1
```

This directory is called as <\$O-S1>.

**<\$O-S1> = <\$O-fem2>/mpi/S1**

# First Example

**hello.f**

```
implicit REAL*8 (A-H,O-Z)
include 'mpif.h'
integer :: PETOT, my_rank, ierr

call MPI_INIT      (ierr)
call MPI_COMM_SIZE (MPI_COMM_WORLD, PETOT, ierr )
call MPI_COMM_RANK (MPI_COMM_WORLD, my_rank, ierr )

write (* , '(a,2i8)') 'Hello World FORTRAN', my_rank, PETOT

call MPI_FINALIZE (ierr)

stop
end
```

**hello.c**

```
#include "mpi.h"
#include <stdio.h>
int main(int argc, char **argv)
{
    int n, myid, numprocs, i;

    MPI_Init(&argc,&argv);
    MPI_Comm_size(MPI_COMM_WORLD,&numprocs);
    MPI_Comm_rank(MPI_COMM_WORLD,&myid);
    printf ("Hello World %d\n", myid);
    MPI_Finalize();
}
```

# Compiling hello.f/c

```
>$ cd <$O-S1>
>$ mpifrtpx -Kfast hello.f
>$ mpifccpx -Kfast hello.c
```

## FORTRAN

\$> mpifrtpx -Kfast hello.f  
“mpifrtpx”:

required compiler & libraries are included for  
FORTRAN90+MPI

## C

\$> mpifccpx -Kfast hello.c  
“mpifccpx”:

required compiler & libraries are included for C+MPI

# Running Job

- Batch Jobs
  - Only batch jobs are allowed.
  - Interactive executions of jobs are not allowed.
- How to run
  - writing job script
  - submitting job
  - checking job status
  - checking results
- Utilization of computational resources
  - 1-node (16 cores) is occupied by each job.
  - Your node is not shared by other jobs.

# Job Script

- <\$0-\$1>/hello.sh
- Scheduling + Shell Script

```
#!/bin/sh
#PJM -L "node=1"
#PJM -L "elapse=00:10:00"
#PJM -L "rscgrp=lecture4"
#PJM -g "gt74"
#PJM -j
#PJM -o "hello.lst"
#PJM --mpi "proc=4"
```

**mpiexec ./a.out**

Number of Nodes  
Computation Time  
Name of "QUEUE"  
Group Name (Wallet)

Standard Output  
MPI Process #

Execs

8 proc's  
"node=1"  
"proc=8"

16 proc's  
"node=1"  
"proc=16"

32 proc's  
"node=2"  
"proc=32"

64 proc's  
"node=4"  
"proc=64"

192 proc's  
"node=12"  
"proc=192"

# Submitting Jobs

```
>$ cd <$O-S1>
>$ pbsub hello.sh
```

```
>$ cat hello.lst
```

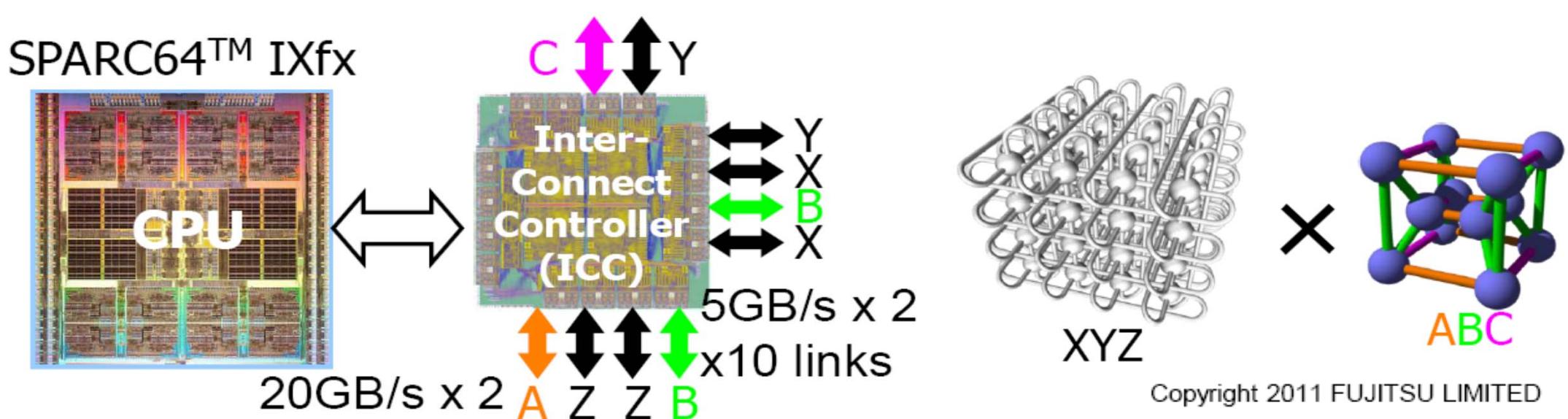
```
Hello World 0
Hello World 3
Hello World 2
Hello World 1
```

# Available QUEUE's

- Following 2 queues are available.
- 1 Tofu (12 nodes) can be used
  - **lecture**
    - 12 nodes (192 cores), 15 min., valid until the end of March, 2014
    - Shared by all “educational” users
  - **lecture4**
    - 12 nodes (192 cores), 15 min., active during class time (Tuesday, 0840-1010)
    - More jobs (compared to **lecture**) can be processed up on availability.

# Tofu Interconnect

- Node Group
  - 12 nodes
  - A-/C- axis: 4 nodes in system board, B-axis: 3 boards
- 6D: (X,Y,Z,A,B,C)
  - ABC 3D Mesh: in each node group:  $2 \times 2 \times 3$
  - XYZ 3D Mesh: connection of node groups:  $10 \times 5 \times 8$
- Job submission according to network topology is possible:
  - Information about used “XYZ” is available after execution.



# Submitting & Checking Jobs

- Submitting Jobs
- Checking status of jobs
- Deleting/aborting
- Checking status of queues
- Detailed info. of queues
- Number of running jobs
- Limitation of submission

pbsub SCRIPT NAME  
pjstat  
pjdel JOB ID  
pjstat --rsc  
pjstat --rsc -x  
pjstat --rsc -b  
pjstat --limit

```
[z30088@oakleaf-fx-6 S2-ref]$ pjstat
```

Oakleaf-FX scheduled stop time: 2012/09/28(Fri) 09:00:00 (Remain: 31days 20:01:46)

JOB_ID	JOB_NAME	STATUS	PROJECT	RSCGROUP	START_DATE	ELAPSE	TOKEN	NODE:COORD
334730	go.sh	RUNNING	gt61	lecture	08/27 12:58:08	00:00:05	0.0	1

# Basic/Essential Functions

```

implicit REAL*8 (A-H,O-Z)
include 'mpif.h'
integer :: PETOT, my_rank, ierr

call MPI_INIT      (ierr)
call MPI_COMM_SIZE (MPI_COMM_WORLD, PETOT, ierr )
call MPI_COMM_RANK (MPI_COMM_WORLD, my_rank, ierr )

write (*,'(a,2i8)') 'Hello World FORTRAN', my_rank, PETOT

call MPI_FINALIZE (ierr)

stop
end

```

```

#include "mpi.h"
#include <stdio.h>
int main(int argc, char **argv)
{
    int n, myid, numprocs, i;

    MPI_Init(&argc,&argv);
    MPI_Comm_size(MPI_COMM_WORLD,&numprocs);
    MPI_Comm_rank(MPI_COMM_WORLD,&myid);

    printf ("Hello World %d\n", myid);
    MPI_Finalize();
}

```

**'mpif.h', "mpi.h"**

Essential Include file

"use mpi" is possible in F90

**MPI\_Init**

Initialization

**MPI\_Comm\_size**

Number of MPI Processes

mpirun -np XX <prog>

**MPI\_Comm\_rank**

Process ID starting from 0

**MPI\_Finalize**

Termination of MPI processes

# Difference between FORTRAN/C

- (Basically) same interface
  - In C, UPPER/lower cases are considered as different
    - e.g.: **MPI\_Comm\_size**
      - MPI: UPPER case
      - First character of the function except “MPI\_” is in UPPER case.
      - Other characters are in lower case.
- In Fortran, return value `ierr` has to be added at the end of the argument list.
- C needs special types for variables:
  - `MPI_Comm`, `MPI_Datatype`, `MPI_Op` etc.
- **MPI\_INIT** is different:
  - `call MPI_INIT ( ierr )`
  - `MPI_Init ( int *argc, char ***argv )`

# What's are going on ?

```
#include "mpi.h"
#include <stdio.h>
int main(int argc, char **argv)
{
    int n, myid, numprocs, i;
    MPI_Init(&argc,&argv);
    MPI_Comm_size(MPI_COMM_WORLD,&numprocs);
    MPI_Comm_rank(MPI_COMM_WORLD,&myid);
    printf ("Hello World %d\n", myid);
    MPI_Finalize();
}
```

<code>#!/bin/sh</code>	
<code>#PJM -L "node=1"</code>	Number of Nodes
<code>#PJM -L "elapse=00:10:00"</code>	Computation Time
<code>#PJM -L "rscgrp=lecture"</code>	Name of "QUEUE"
<code>#PJM -g "gt64"</code>	Group Name (Wallet)
<code>#PJM -j</code>	
<code>#PJM -o "hello.lst"</code>	Standard Output
<code>#PJM --mpi "proc=4"</code>	MPI Process #
<code>mpieexec ./a.out</code>	Execs

- **mpieexec** starts up 4 MPI processes ("proc=4")
  - A single program runs on four processes.
  - each process writes a value of `myid`
- Four processes do same operations, but values of `myid` are different.
- Output of each process is different.
- **That is SPMD !**

# mpi.h, mpif.h

```
implicit REAL*8 (A-H,O-Z)
include 'mpif.h'
integer :: PETOT, my_rank, ierr

call MPI_INIT      (ierr)
call MPI_COMM_SIZE (MPI_COMM_WORLD, PETOT, ierr )
call MPI_COMM_RANK (MPI_COMM_WORLD, my_rank, ierr )

write (*,'(a,2i8)') 'Hello World FORTRAN', my_rank, PETOT

call MPI_FINALIZE (ierr)

stop
end
```

```
#include "mpi.h"
#include <stdio.h>
int main(int argc, char **argv)
{
    int n, myid, numprocs, i;

    MPI_Init(&argc,&argv);
    MPI_Comm_size(MPI_COMM_WORLD,&numprocs);
    MPI_Comm_rank(MPI_COMM_WORLD,&myid);

    printf ("Hello World %d\n", myid);
    MPI_Finalize();
}
```

- Various types of parameters and variables for MPI & their initial values.
- Name of each var. starts from “MPI\_”
- Values of these parameters and variables cannot be changed by users.
- Users do not specify variables starting from “MPI\_” in users’ programs.

# MPI\_Init

- Initialize the MPI execution environment (required)
- It is recommended to put this BEFORE all statements in the program.
- **MPI\_Init (argc, argv)**

```
#include "mpi.h"
#include <stdio.h>
int main(int argc, char **argv)
{
    int n, myid, numprocs, i;

    MPI_Init(&argc,&argv);
    MPI_Comm_size(MPI_COMM_WORLD,&numprocs);
    MPI_Comm_rank(MPI_COMM_WORLD,&myid);

    printf ("Hello World %d\n", myid);
    MPI_Finalize();
}
```

# MPI\_Finalize

- Terminates MPI execution environment (required)
  - It is recommended to put this AFTER all statements in the program.
  - Please do not forget this.
- 
- **MPI\_Finalize ( )**

```
#include "mpi.h"
#include <stdio.h>
int main(int argc, char **argv)
{
    int n, myid, numprocs, i;

    MPI_Init(&argc,&argv);
    MPI_Comm_size(MPI_COMM_WORLD,&numprocs);
    MPI_Comm_rank(MPI_COMM_WORLD,&myid);

    printf ("Hello World %d\n", myid);
    MPI_Finalize();
}
```

# MPI\_Comm\_size

- Determines the size of the group associated with a communicator
- not required, but very convenient function
- **MPI\_Comm\_size (comm, size)**
  - **comm** MPI\_Comm I communicator
  - **size** int O number of processes in the group of communicator

```
#include "mpi.h"
#include <stdio.h>
int main(int argc, char **argv)
{
    int n, myid, numprocs, i;

    MPI_Init(&argc,&argv);
    MPI_Comm_size(MPI_COMM_WORLD,&numprocs);
    MPI_Comm_rank(MPI_COMM_WORLD,&myid);

    printf ("Hello World %d\n", myid);
    MPI_Finalize();
}
```

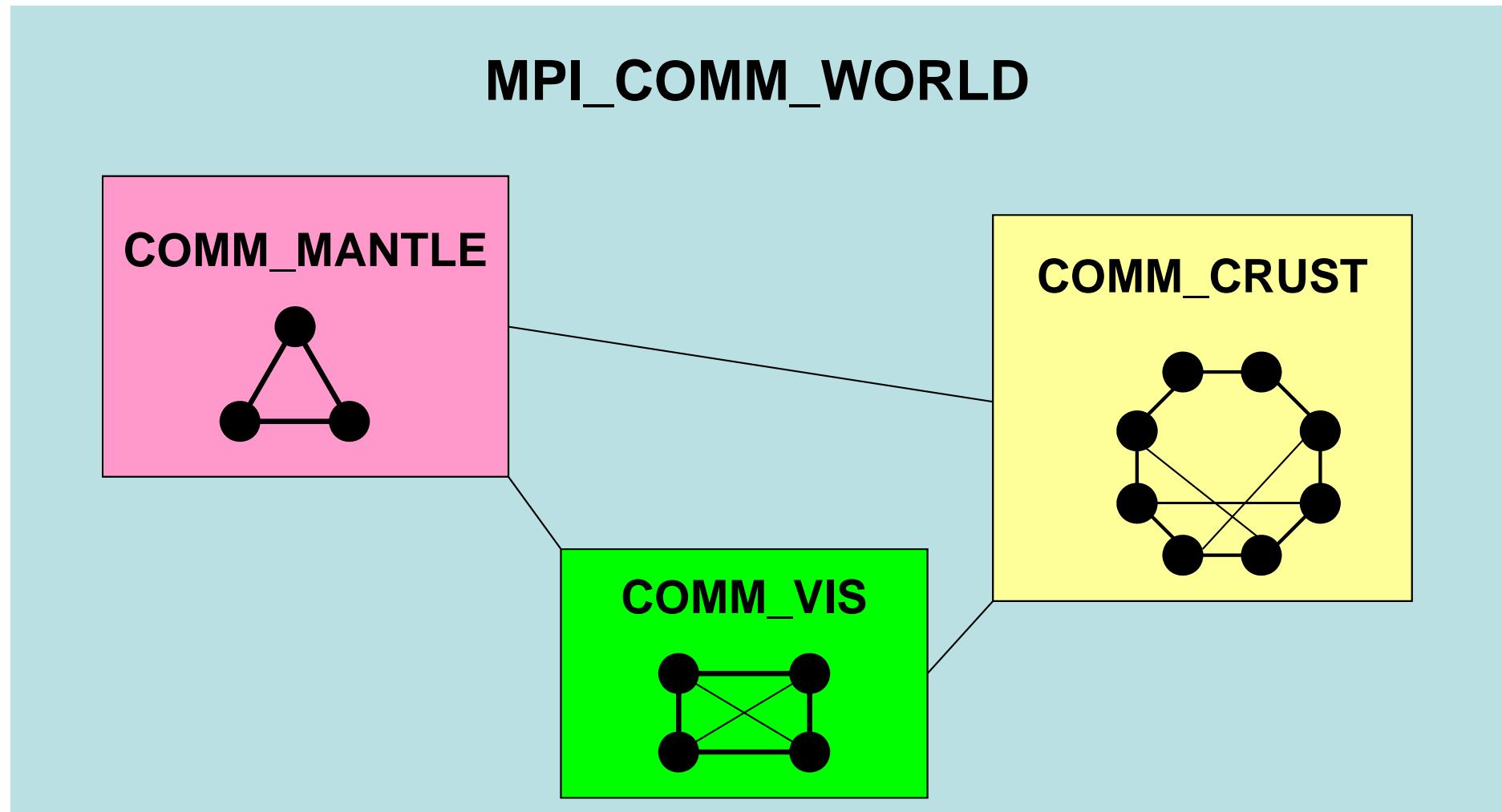
# What is Communicator ?

`MPI_Comm_Size (MPI_COMM_WORLD, PETOT)`

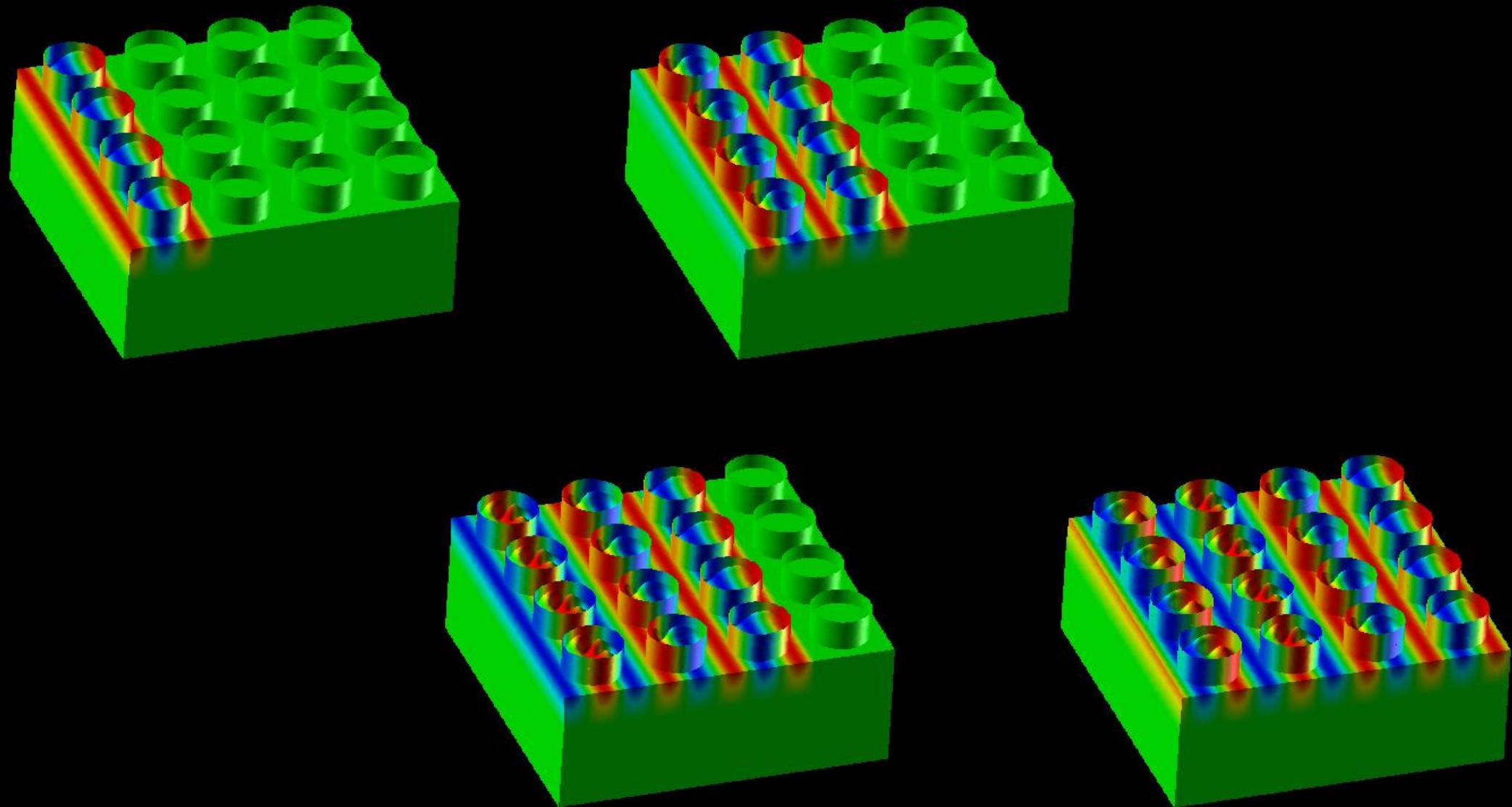
- Group of processes for communication
- Communicator must be specified in MPI program as a unit of communication
- All processes belong to a group, named “**MPI\_COMM\_WORLD**” (default)
- Multiple communicators can be created, and complicated operations are possible.
  - Computation, Visualization
- Only “**MPI\_COMM\_WORLD**” is needed in this class.

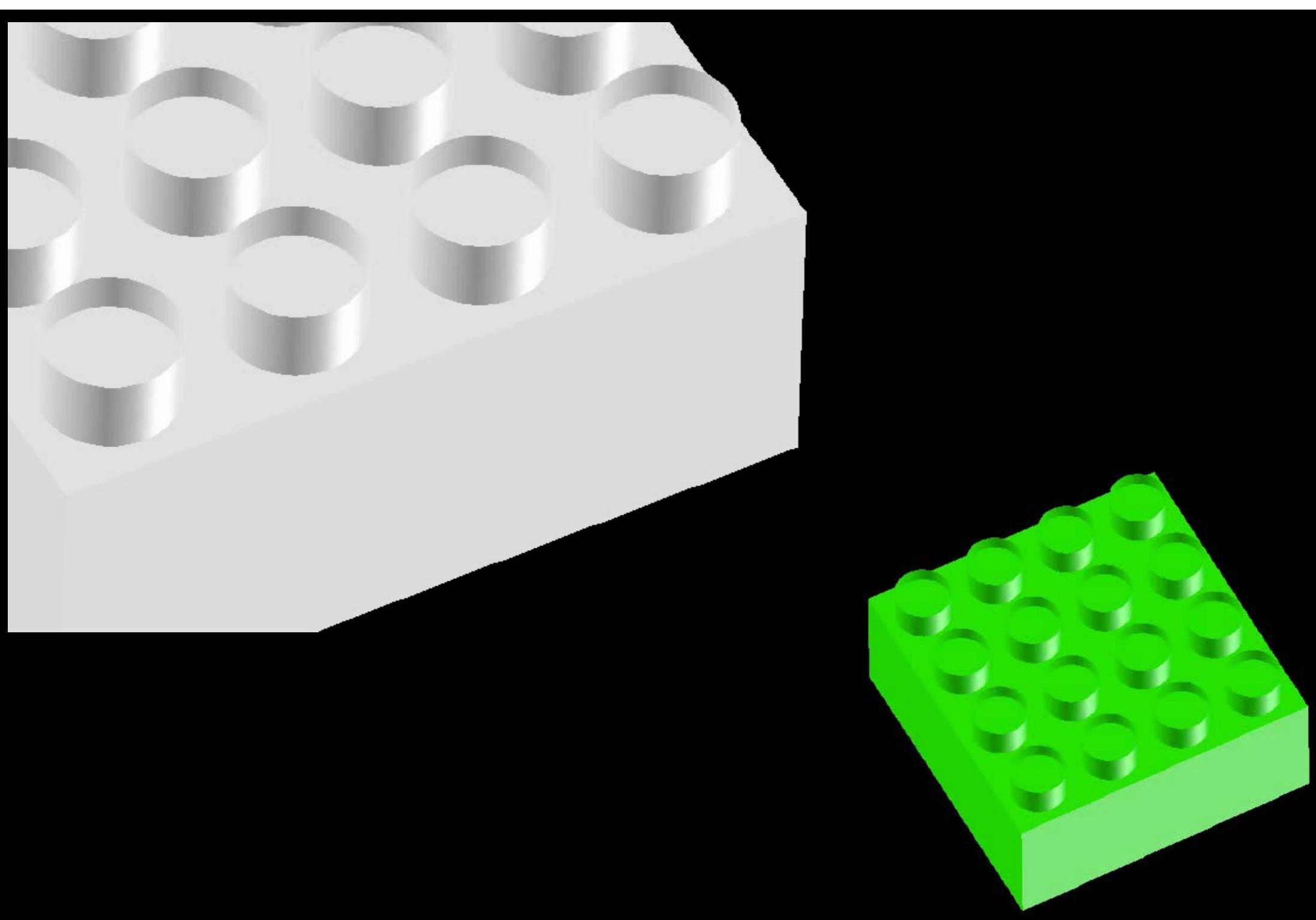
# Communicator in MPI

One process can belong to multiple communicators



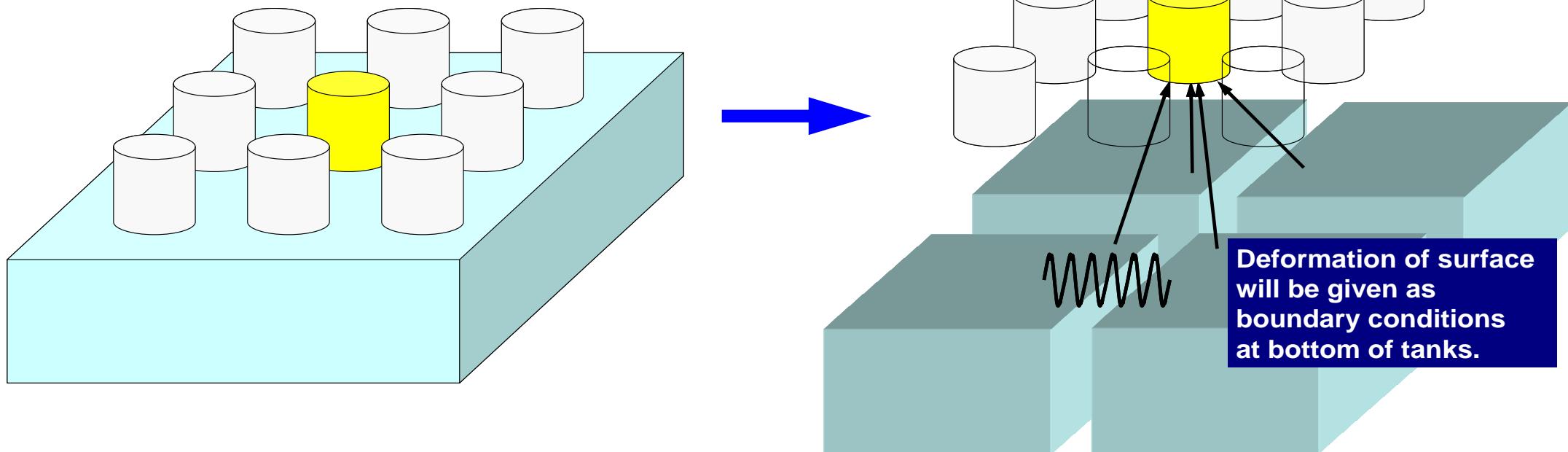
# Coupling between “Ground Motion” and “Sloshing of Tanks for Oil-Storage”





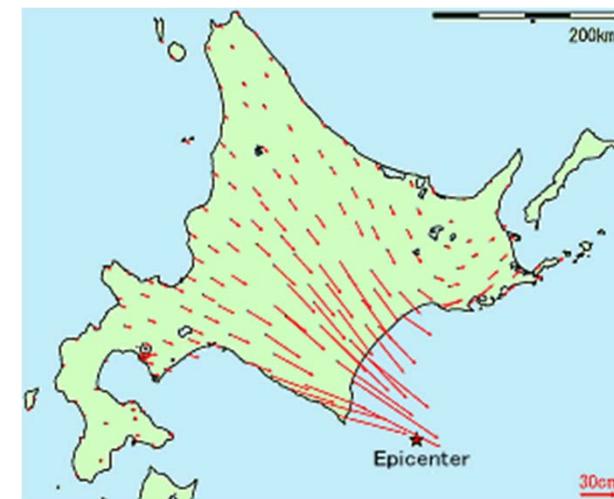
# Target Application

- Coupling between “Ground Motion” and “Sloshing of Tanks for Oil-Storage”
  - “One-way” coupling from “Ground Motion” to “Tanks”.
  - Displacement of ground surface is given as forced displacement of bottom surface of tanks.
  - 1 Tank = 1 PE (serial)

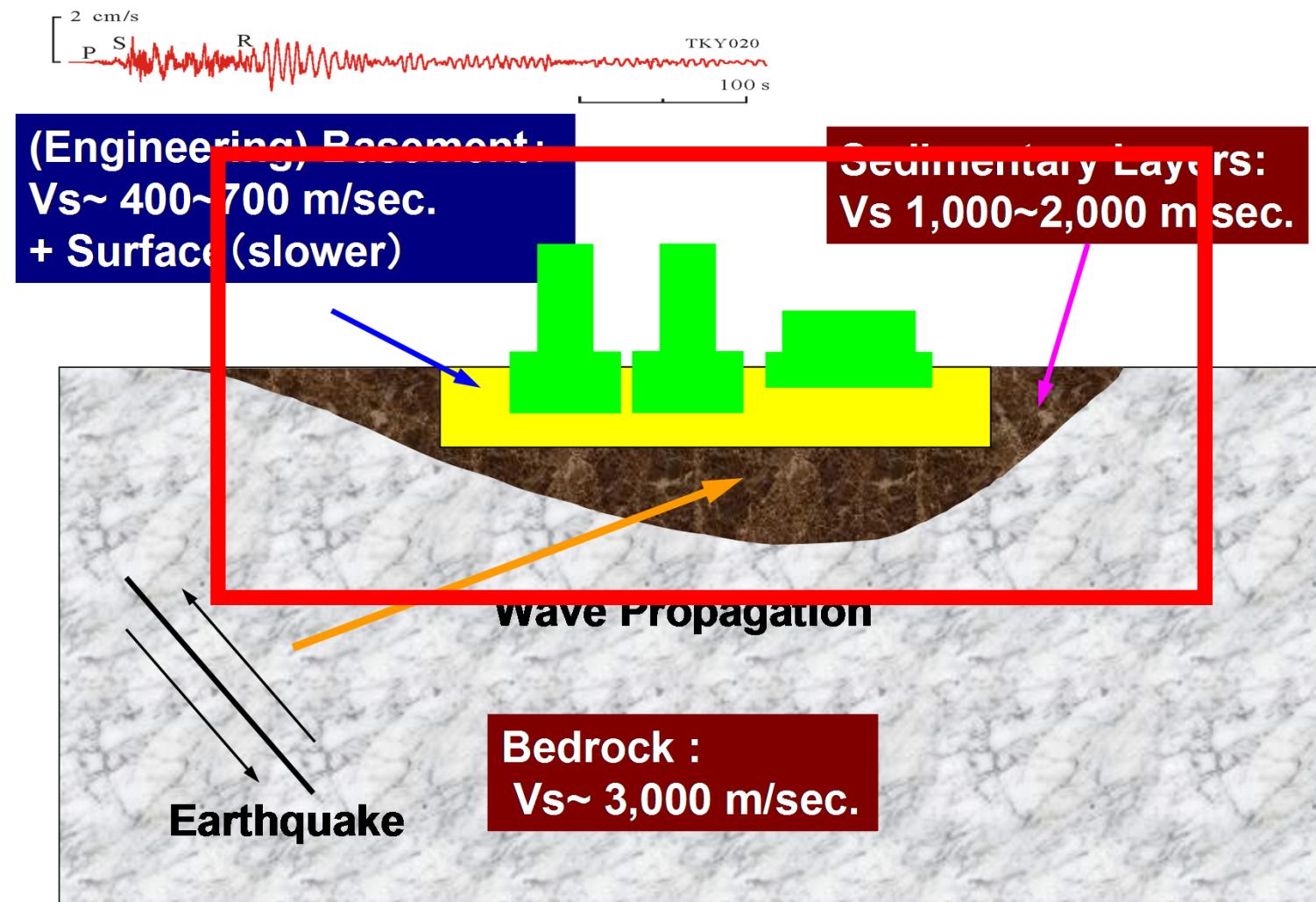


# 2003 Tokachi Earthquake (M8.0)

Fire accident of oil tanks due to long period ground motion (surface waves) developed in the basin of Tomakomai

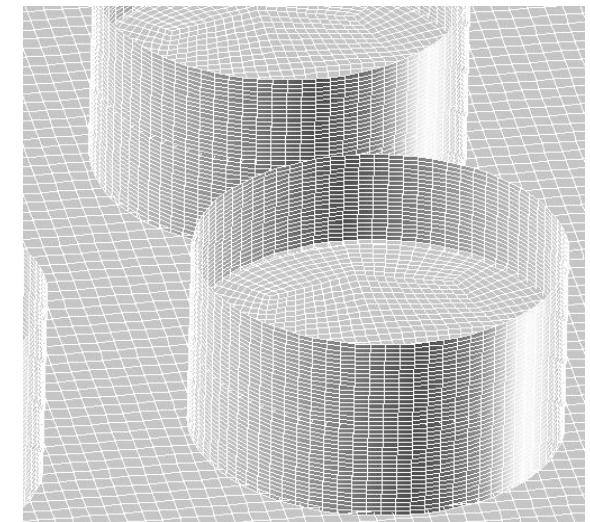
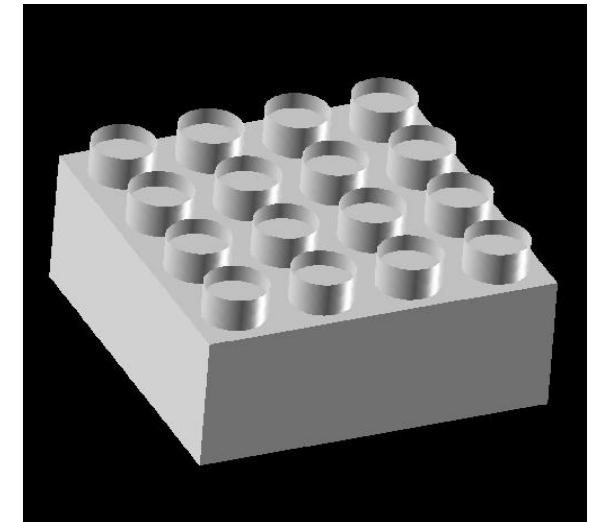


# Seismic Wave Propagation, Underground Structure

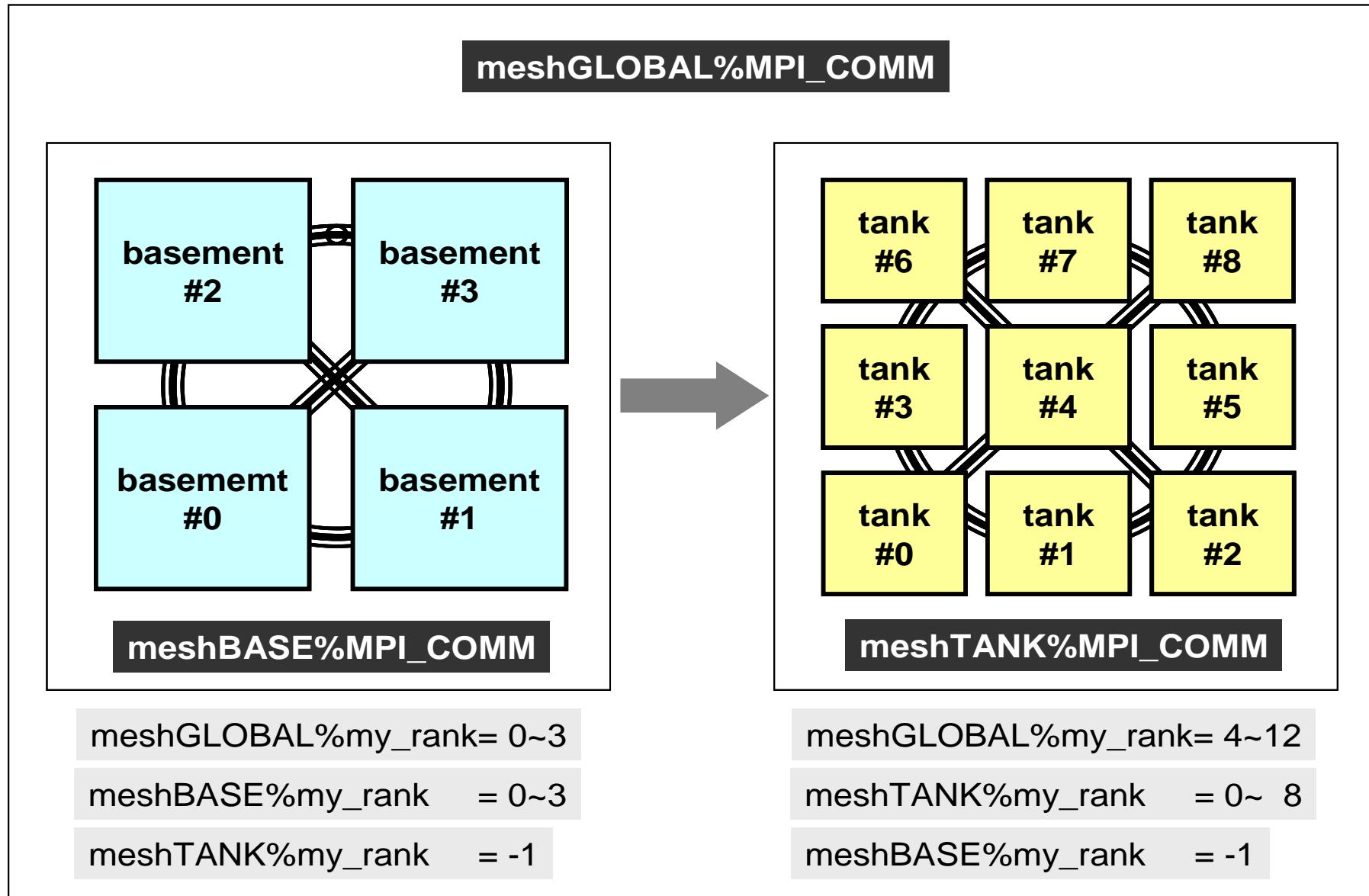


# Simulation Codes

- Ground Motion (Ichimura): Fortran
  - Parallel FEM, 3D Elastic/Dynamic
    - Explicit forward Euler scheme
  - Each element:  $2\text{m} \times 2\text{m} \times 2\text{m}$  cube
  - $240\text{m} \times 240\text{m} \times 100\text{m}$  region
- Sloshing of Tanks (Nagashima): C
  - Serial FEM (Embarrassingly Parallel)
    - Implicit backward Euler, Skyline method
    - Shell elements + Inviscid potential flow
  - D: 42.7m, H: 24.9m, T: 20mm,
  - Frequency: 7.6sec.
  - 80 elements in circ., 0.6m mesh in height
  - Tank-to-Tank: 60m,  $4 \times 4$
- Total number of unknowns: 2,918,169



# Three Communicators



# **MPI\_Comm\_rank**

C

- Determines the rank of the calling process in the communicator
    - “ID of MPI process” is sometimes called “rank”
  - **MPI\_Comm\_rank (comm, rank)**
    - comm      MPI\_Comm I                    communicator
    - rank      int            O                    rank of the calling process in the group of comm  
**Starting from “0”**

```
#include "mpi.h"
#include <stdio.h>
int main(int argc, char **argv)
{
    int n, myid, numprocs, i;

    MPI_Init(&argc,&argv);
    MPI_Comm_size(MPI_COMM_WORLD,&numprocs);
    MPI_Comm_rank(MPI_COMM_WORLD,&myid);

    printf ("Hello World %d\n", myid);
    MPI_Finalize();
}
```

# MPI\_Abort

- Aborts MPI execution environment
- **MPI\_Abort (comm, errcode)**
  - comm      MPI\_Comm    I      communicator
  - errcode    int      O      Error Code

# MPI\_Wtime

- Returns an elapsed time on the calling processor
- **time= MPI\_Wtime ()**
  - **time** double 0 Time in seconds since an arbitrary time in the past.

```
...
double Stime, Etime;

Stime= MPI_Wtime ();

(...)

Etime= MPI_Wtime ();
```

# Example of MPI\_Wtime

```
$> cd <$O-S1>
```

```
$> mpifccpx -O1 time.c
```

```
$> mpifrtpx -O1 time.f
```

(modify go4.sh, 4 processes)

```
$> pbsub go4.sh
```

0	1.113281E+00
3	1.113281E+00
2	1.117188E+00
1	1.117188E+00

Process ID	Time
------------	------

# MPI\_Wtick

- Returns the resolution of MPI\_Wtime
- depends on hardware, and compiler
- **time= MPI\_Wtick ( )**
  - time double 0 Time in seconds of resolution of MPI\_Wtime

```
implicit REAL*8 (A-H,O-Z)
include 'mpif.h'

...
TM= MPI_WTICK ( )
write (*,* ) TM
...
```

```
double Time;

...
Time = MPI_Wtick();
printf( "%5d%16.6E\n", MyRank, Time);
...
```

# Example of MPI\_Wtick

```
$> cd <$O-S1>  
  
$> mpifccpx -O1 wtick.c  
$> mpifrtpx -O1 wtick.f  
  
(modify go1.sh, 1 process)  
$> pbsub go1.sh
```

# MPI\_Barrier

- Blocks until all processes in the communicator have reached this routine.
- Mainly for debugging, huge overhead, not recommended for real code.
- **MPI\_Barrier (comm)**
  - comm MPI\_Comm I communicator

- What is MPI ?
- Your First MPI Program: Hello World
- **Global/Local Data**
- Collective Communication
- Peer-to-Peer Communication

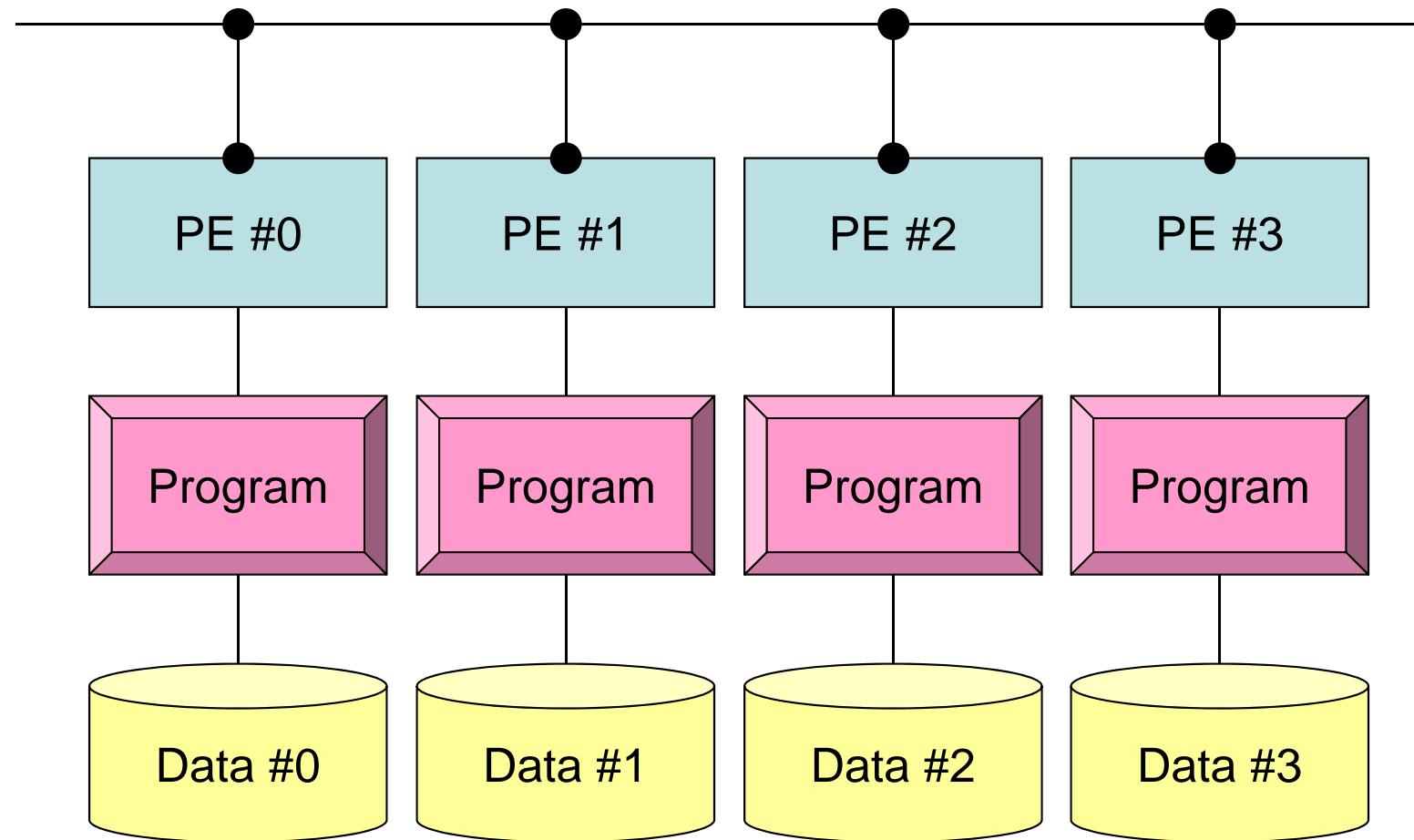
# Data Structures & Algorithms

- Computer program consists of data structures and algorithms.
- They are closely related. In order to implement an algorithm, we need to specify an appropriate data structure for that.
  - We can even say that “Data Structures=Algorithms”
  - Some people may not agree with this, I (KN) think it is true for scientific computations from my experiences.
- Appropriate data structures for parallel computing must be specified before starting parallel computing.

# SPMD: Single Program Multiple Data

- There are various types of “*parallel computing*”, and there are many algorithms.
- Common issue is SPMD (Single Program Multiple Data).
- It is ideal that parallel computing is done in the same way for serial computing (except communications)
  - It is required to specify processes with communications and those without communications.

# What is a data structure which is appropriate for SPMD ?



# Data Structure for SMPD (1/2)

- SPMD: Large data is decomposed into small pieces, and each piece is processed by each processor/process
- Consider the following simple computation for vector **Vg** with length of **Ng** (=20):

```
int main(){
    int i,Ng;
    double Vg[20];
    Ng=20;
    for(i=0;i<Ng;i++){
        Vg[i] = 2.0*Vg[i];}
    return 0;}
```

- If you compute this using four processors, each processor stores and processes 5 (=20/4) components of **Vg**.

# Data Structure for SMPD (2/2)

- i.e.

```
int main(){
    int i,Nl;
    double v1[5];
    Nl=5;
    for(i=0;i<Nl;i++){
        v1[i] = 2.0*v1[i];}
    return 0;}
```

- Thus, a “single program” can execute parallel processing.
  - In each process, components of “V1” are different: Multiple Data
  - Computation using only “V1” (as long as possible) leads to efficient parallel computation.
  - Program is not different from that for serial CPU (in the previous page).

# Global & Local Data

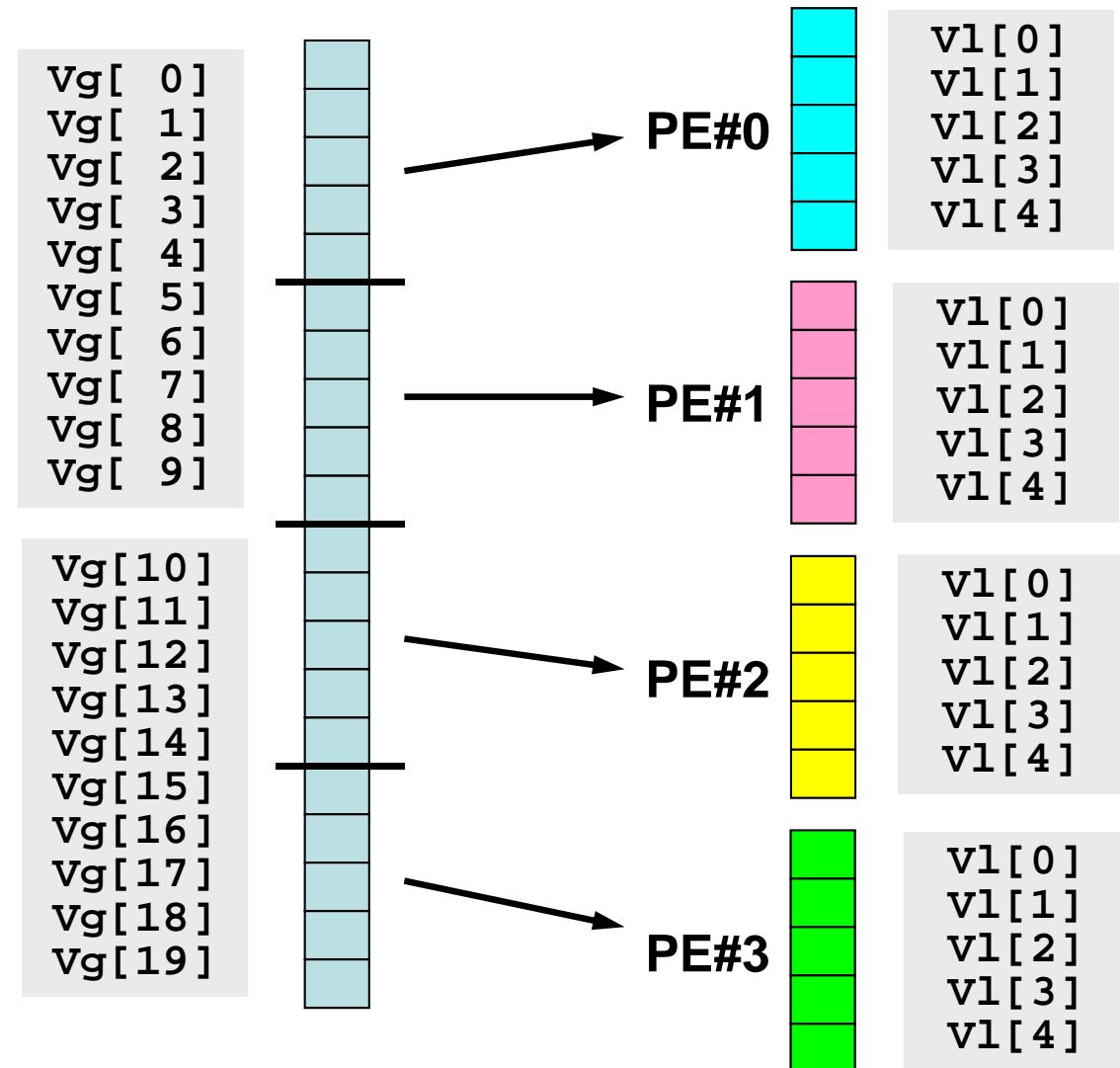
- $V_g$ 
  - Entire Domain
  - “Global Data” with “Global ID” from 1 to 20
- $V_l$ 
  - for Each Process (PE, Processor, Domain)
  - “Local Data” with “Local ID” from 1 to 5
  - **Efficient utilization of local data leads to excellent parallel efficiency.**

# Idea of Local Data in C

**Vg:** Global Data

- 0<sup>th</sup>-4<sup>th</sup> comp. on PE#0
- 5<sup>th</sup>-9<sup>th</sup> comp. on PE#1
- 10<sup>th</sup>-14<sup>th</sup> comp. on PE#2
- 15<sup>th</sup>-19<sup>th</sup> comp. on PE#3

Each of these four sets corresponds to 0<sup>th</sup>-4<sup>th</sup> components of VI (local data) where local ID's are 0-4.



# Global & Local Data

- $V_g$ 
  - Entire Domain
  - “Global Data” with “Global ID” from 1 to 20
- $VI$ 
  - for Each Process (PE, Processor, Domain)
  - “Local Data” with “Local ID” from 1 to 5
- Please keep your attention to the following:
  - How to generate  $VI$  (local data) from  $V_g$  (global data)
  - How to map components, from  $V_g$  to  $VI$ , and from  $VI$  to  $V_g$ .
  - What to do if  $VI$  cannot be calculated on each process in independent manner.
  - Processing as localized as possible leads to excellent parallel efficiency:
    - Data structures & algorithms for that purpose.

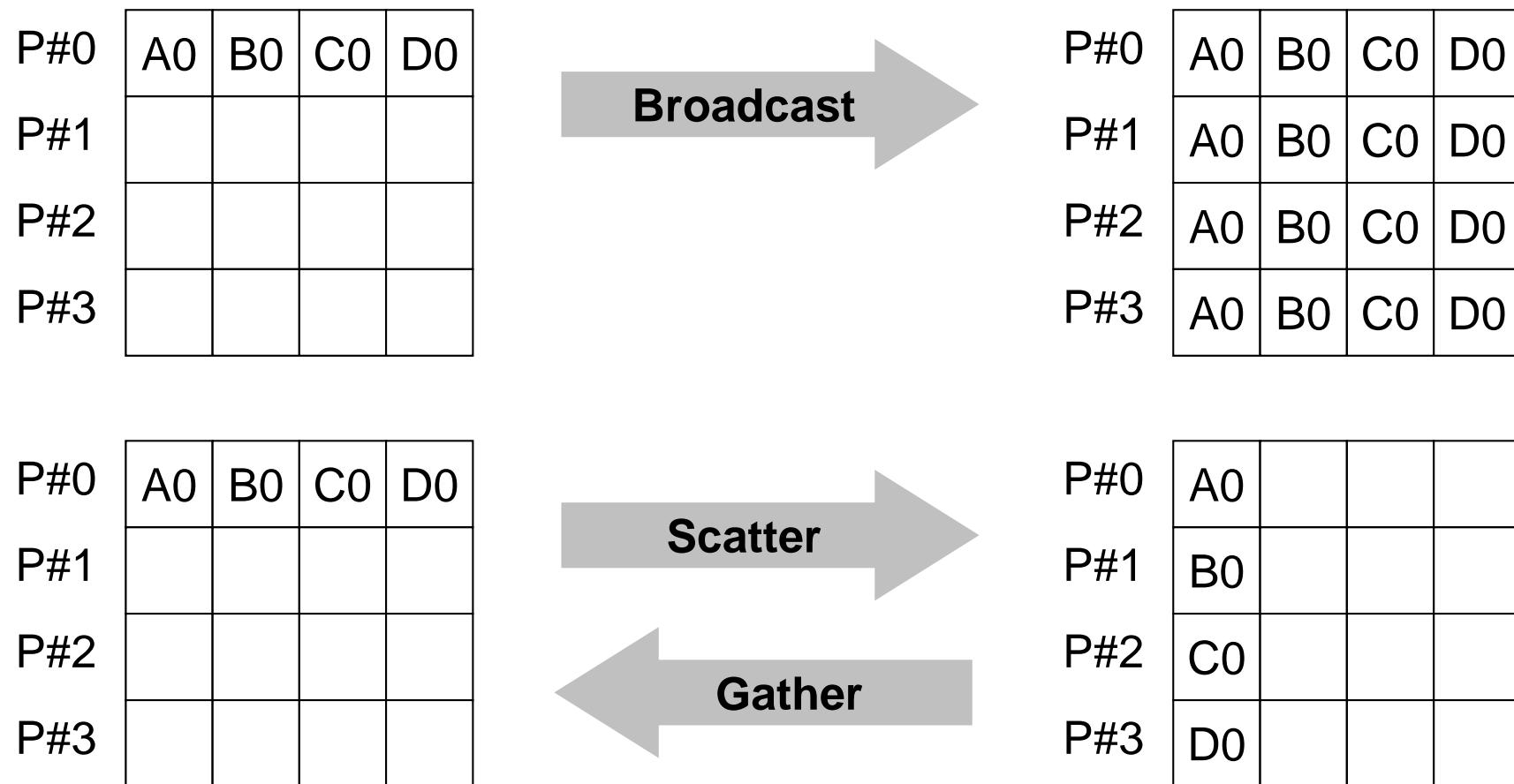
- What is MPI ?
- Your First MPI Program: Hello World
- Global/Local Data
- **Collective Communication**
- Peer-to-Peer Communication

# What is Collective Communication ?

## 集団通信, グループ通信

- Collective communication is the process of exchanging information between multiple MPI processes in the communicator: one-to-all or all-to-all communications.
- Examples
  - Broadcasting control data
  - Max, Min
  - Summation
  - Dot products of vectors
  - Transformation of dense matrices

# Example of Collective Communications (1/4)



# Example of Collective Communications (2/4)

P#0	A0			
P#1	B0			
P#2	C0			
P#3	D0			

All gather

P#0	A0	B0	C0	D0
P#1	A0	B0	C0	D0
P#2	A0	B0	C0	D0
P#3	A0	B0	C0	D0

P#0	A0	A1	A2	A3
P#1	B0	B1	B2	B3
P#2	C0	C1	C2	C3
P#3	D0	D1	D2	D3

All-to-All

P#0	A0	B0	C0	D0
P#1	A1	B1	C1	D1
P#2	A2	B2	C2	D2
P#3	A3	B3	C3	D3

# Example of Collective Communications (3/4)

P#0	A0	B0	C0	D0
P#1	A1	B1	C1	D1
P#2	A2	B2	C2	D2
P#3	A3	B3	C3	D3

Reduce

P#0	op.A0-A3	op.B0-B3	op.C0-C3	op.D0-D3
P#1				
P#2				
P#3				

P#0	A0	B0	C0	D0
P#1	A1	B1	C1	D1
P#2	A2	B2	C2	D2
P#3	A3	B3	C3	D3

All reduce

P#0	op.A0-A3	op.B0-B3	op.C0-C3	op.D0-D3
P#1	op.A0-A3	op.B0-B3	op.C0-C3	op.D0-D3
P#2	op.A0-A3	op.B0-B3	op.C0-C3	op.D0-D3
P#3	op.A0-A3	op.B0-B3	op.C0-C3	op.D0-D3

# Example of Collective Communications (4/4)

P#0	A0	B0	C0	D0
P#1	A1	B1	C1	D1
P#2	A2	B2	C2	D2
P#3	A3	B3	C3	D3

Reduce scatter

P#0	op.A0-A3			
P#1	op.B0-B3			
P#2	op.C0-C3			
P#3	op.D0-D3			

# Examples by Collective Comm.

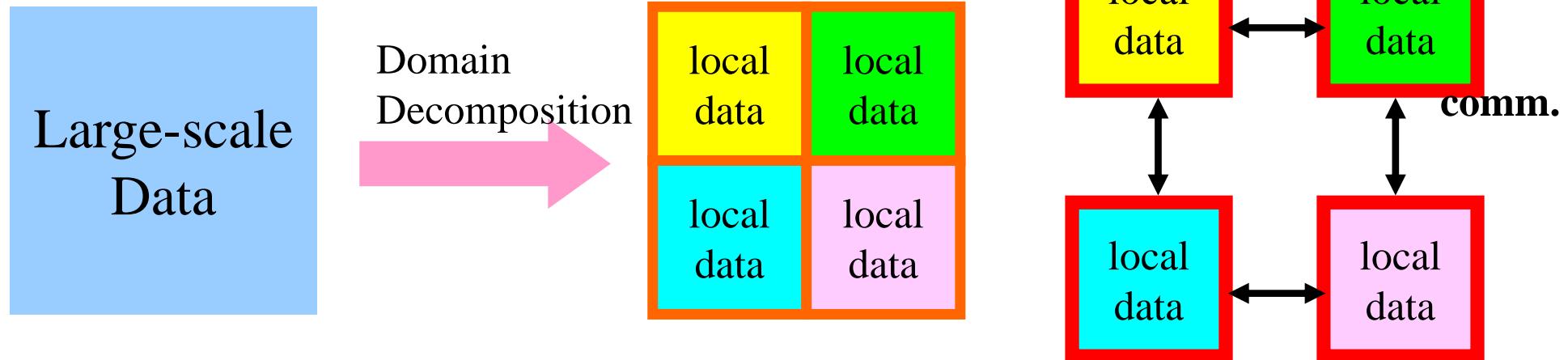
- Dot Products of Vectors
- Scatter/Gather
- Reading Distributed Files

# Global/Local Data

- Data structure of parallel computing based on SPMD, where large scale “global data” is decomposed to small pieces of “local data”.

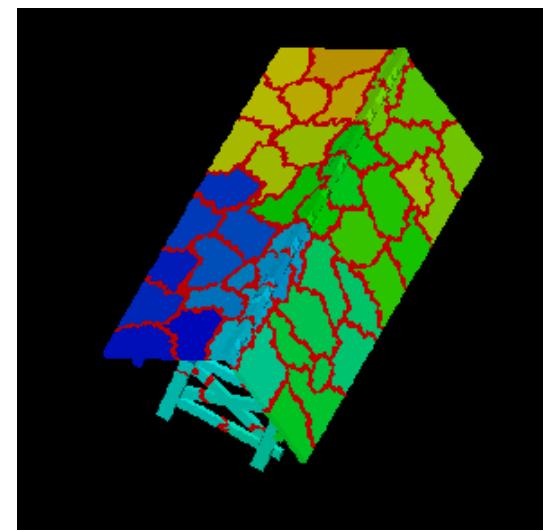
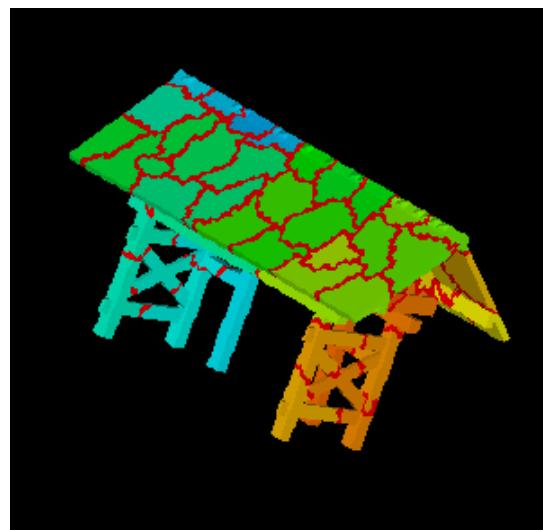
# Domain Decomposition/Partitioning

- PC with 1GB RAM: can execute FEM application with up to  $10^6$  meshes
  - $10^3\text{km} \times 10^3\text{ km} \times 10^2\text{ km}$  (SW Japan):  $10^8$  meshes by 1km cubes
- Large-scale Data: Domain decomposition, parallel & local operations
- Global Computation: Comm. among domains needed



# Local Data Structure

- It is important to define proper local data structure for target computation (and its algorithm)
  - Algorithms= Data Structures
- Main objective of this class !



# Global/Local Data

- Data structure of parallel computing based on SPMD, where large scale “global data” is decomposed to small pieces of “local data”.
- Consider the dot product of following VECp and VECs with length=20 by parallel computation using 4 processors

<b>VECp[ 0 ] =</b>	<b>2</b>
<b>[ 1 ] =</b>	<b>2</b>
<b>[ 2 ] =</b>	<b>2</b>
...	
<b>[ 17 ] =</b>	<b>2</b>
<b>[ 18 ] =</b>	<b>2</b>
<b>[ 19 ] =</b>	<b>2</b>

<b>VECs[ 0 ] =</b>	<b>3</b>
<b>[ 1 ] =</b>	<b>3</b>
<b>[ 2 ] =</b>	<b>3</b>
...	
<b>[ 17 ] =</b>	<b>3</b>
<b>[ 18 ] =</b>	<b>3</b>
<b>[ 19 ] =</b>	<b>3</b>

# <\$O-S1>/dot.c

```
implicit REAL*8 (A-H,O-Z)
real(kind=8),dimension(20):: &
    VECp,   VECs

do i= 1, 20
    VECp(i)= 2.0d0
    VECs(i)= 3.0d0
enddo

sum= 0.d0
do ii= 1, 20
    sum= sum + VECp(ii)*VECs(ii)
enddo

stop
end
```

```
#include <stdio.h>
int main(){
    int i;
    double VECp[20], VECs[20]
    double sum;

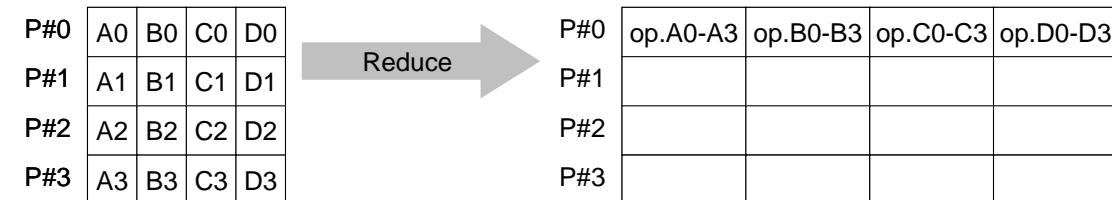
    for(i=0;i<20;i++){
        VECp[i]= 2.0;
        VECs[i]= 3.0;
    }

    sum = 0.0;
    for(i=0;i<20;i++){
        sum += VECp[i] * VECs[i];
    }
    return 0;
}
```

# Execution of <\$O-S1>/dot.f, dot.c (do this on ECCS2012)

```
>$ cd <$T-S1>  
  
>$ cc -O3 dot.c  
>$ f90 -O3 dot.f  
  
>$ ./a.out  
  
1          2.          3.  
2          2.          3.  
3          2.          3.  
...  
18         2.          3.  
19         2.          3.  
20         2.          3.  
  
dot product      120.
```

# MPI\_Reduce



- Reduces values on all processes to a single value
  - Summation, Product, Max, Min etc.
- `MPI_Reduce ( sendbuf, recvbuf, count, datatype, op, root, comm )`**
  - sendbuf** choice I starting address of send buffer
  - recvbuf** choice O starting address receive buffer  
*type is defined by "datatype"*
  - count** int I number of elements in send/receive buffer
  - datatype** MPI\_Datatype I data type of elements of send/receive buffer
    - FORTAN MPI\_INTEGER, MPI\_REAL, MPI\_DOUBLE\_PRECISION, MPI\_CHARACTER etc.
    - C MPI\_INT, MPI\_FLOAT, MPI\_DOUBLE, MPI\_CHAR etc
  - op** MPI\_Op I reduce operation
    - MPI\_MAX, MPI\_MIN, MPI\_SUM, MPI\_PROD, MPI\_LAND, MPI\_BAND etc*Users can define operations by [MPI\\_OP\\_CREATE](#)*
  - root** int I rank of root process
  - comm** MPI\_Comm I communicator

# Send/Receive Buffer (Sending/Receiving)

- Arrays of “send (sending) buffer” and “receive (receiving) buffer” often appear in MPI.
- Addresses of “send (sending) buffer” and “receive (receiving) buffer” must be different.

# Example of MPI\_Reduce (1/2)

**MPI\_Reduce**

(**sendbuf**,**recvbuf**,**count**,**datatype**,**op**,**root**,**comm**)

```
double x0, x1;  
  
MPI_Reduce  
(&x0, &x1, 1, MPI_DOUBLE, MPI_MAX, 0, <comm>);
```

```
double x0[4], xmax[4];  
  
MPI_Reduce  
(x0, xmax, 4, MPI_DOUBLE, MPI_MAX, 0, <comm>);
```

Global Max values of X0[i] go to XMAX[i] on #0 process (i=0~3)

# Example of MPI\_Reduce (2/2)

**MPI\_Reduce**

(**sendbuf**, **recvbuf**, **count**, **datatype**, **op**, **root**, **comm**)

```
double X0, XSUM;  
  
MPI_Reduce  
(&X0, &XSUM, 1, MPI_DOUBLE, MPI_SUM, 0, <comm>)
```

Global summation of X0 goes to XSUM on #0 process.

```
double X0[4];  
  
MPI_Reduce  
(&X0[0], &X0[2], 2, MPI_DOUBLE_PRECISION, MPI_SUM, 0, <comm>)
```

- Global summation of X0[0] goes to X0[2] on #0 process.
- Global summation of X0[1] goes to X0[3] on #0 process.

# MPI\_Bcast

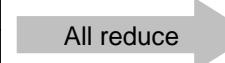
P#0	A0	B0	C0	D0
P#1				
P#2				
P#3				



P#0	A0	B0	C0	D0
P#1	A0	B0	C0	D0
P#2	A0	B0	C0	D0
P#3	A0	B0	C0	D0

- Broadcasts a message from the process with rank "root" to all other processes of the communicator
- **`MPI_Bcast (buffer, count, datatype, root, comm)`**
  - **buffer** choice I/O starting address of buffer  
**type is defined by "datatype"**
  - **count** int I number of elements in send/receive buffer
  - **datatype** MPI\_Datatype I data type of elements of send/receive buffer
    - FORTTRAN MPI\_INTEGER, MPI\_REAL, MPI\_DOUBLE\_PRECISION, MPI\_CHARACTER etc.
    - C MPI\_INT, MPI\_FLOAT, MPI\_DOUBLE, MPI\_CHAR etc.
  - **root** int I **rank of root process**
  - **comm** MPI\_Comm I communicator

# MPI\_Allreduce



P#0	A0	B0	C0	D0		P#0	op.A0-A3	op.B0-B3	op.C0-C3	op.D0-D3
P#1	A1	B1	C1	D1	All reduce	P#1	op.A0-A3	op.B0-B3	op.C0-C3	op.D0-D3
P#2	A2	B2	C2	D2		P#2	op.A0-A3	op.B0-B3	op.C0-C3	op.D0-D3
P#3	A3	B3	C3	D3		P#3	op.A0-A3	op.B0-B3	op.C0-C3	op.D0-D3

- MPI\_Reduce + MPI\_Bcast
- Summation (of dot products) and MAX/MIN values are likely to utilized in each process
- **call MPI\_Allreduce**  
**(sendbuf,recvbuf,count,datatype,op, comm)**
  - **sendbuf** choice I starting address of send buffer
  - **recvbuf** choice O starting address receive buffer  
type is defined by "datatype"
  - **count** int I number of elements in send/receive buffer
  - **datatype** MPI\_Datatype I data type of elements of send/recive buffer
  - **op** MPI\_Op I reduce operation
  - **comm** MPI\_Comm I communicator

C

# “op” of MPI\_Reduce/Allreduce

## **MPI\_Reduce**

```
(sendbuf,recvbuf,count,datatype,op,root,comm)
```

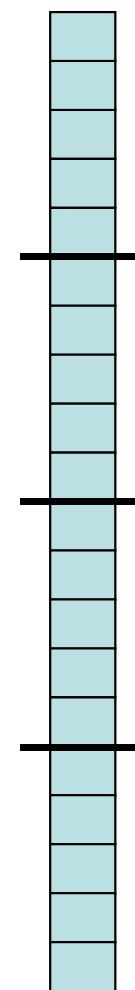
- **MPI\_MAX, MPI\_MIN** Max, Min
- **MPI\_SUM, MPI\_PROD** Summation, Product
- **MPI LAND** Logical AND

# Local Data (1/2)

- Decompose vector with length=20 into 4 domains (processes)
- Each process handles a vector with length= 5

```
VECp[ 0 ]= 2  
[ 1 ]= 2  
[ 2 ]= 2  
...  
[17 ]= 2  
[18 ]= 2  
[19 ]= 2
```

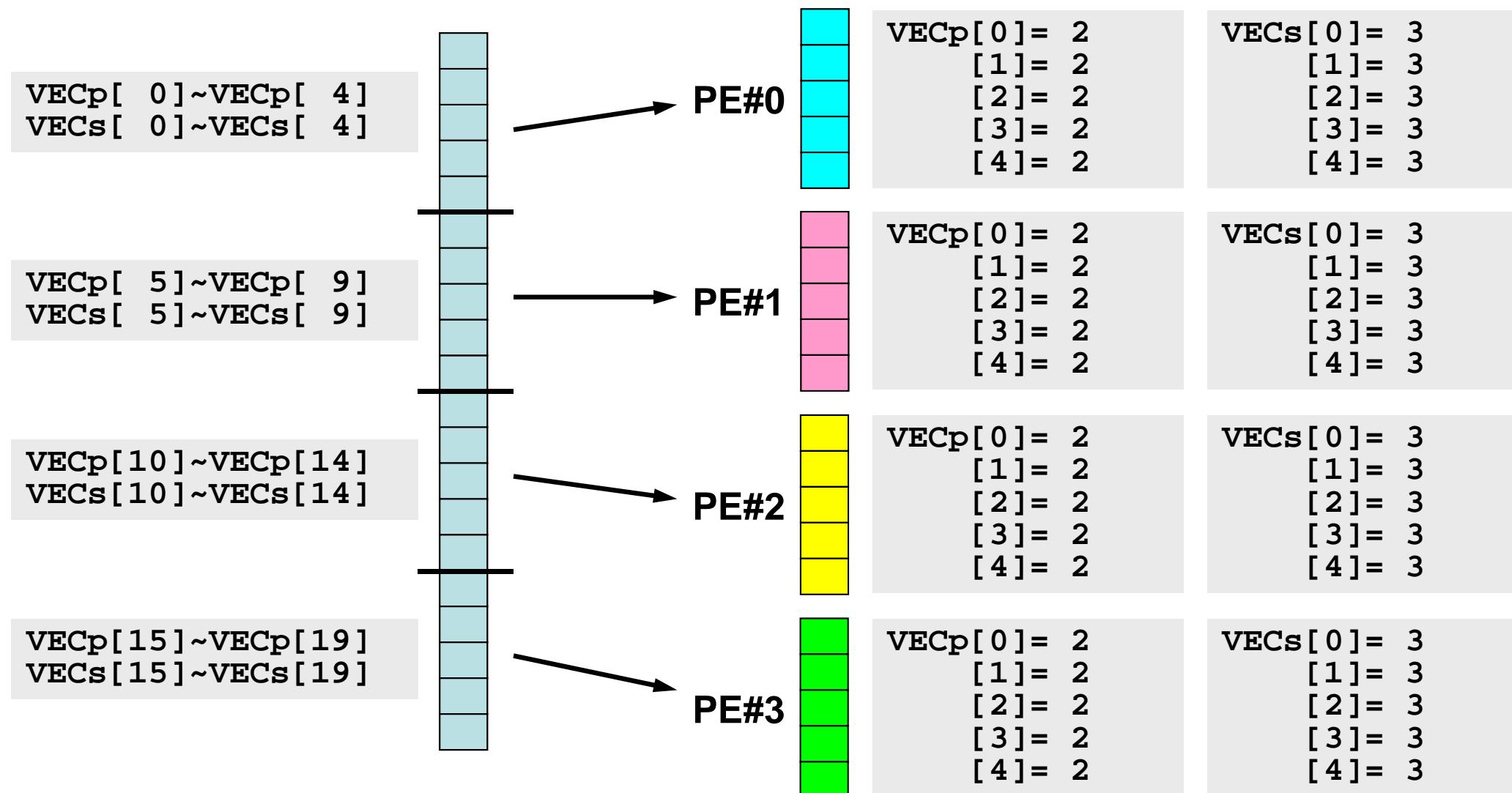
```
VECs[ 0 ]= 3  
[ 1 ]= 3  
[ 2 ]= 3  
...  
[17 ]= 3  
[18 ]= 3  
[19 ]= 3
```



C

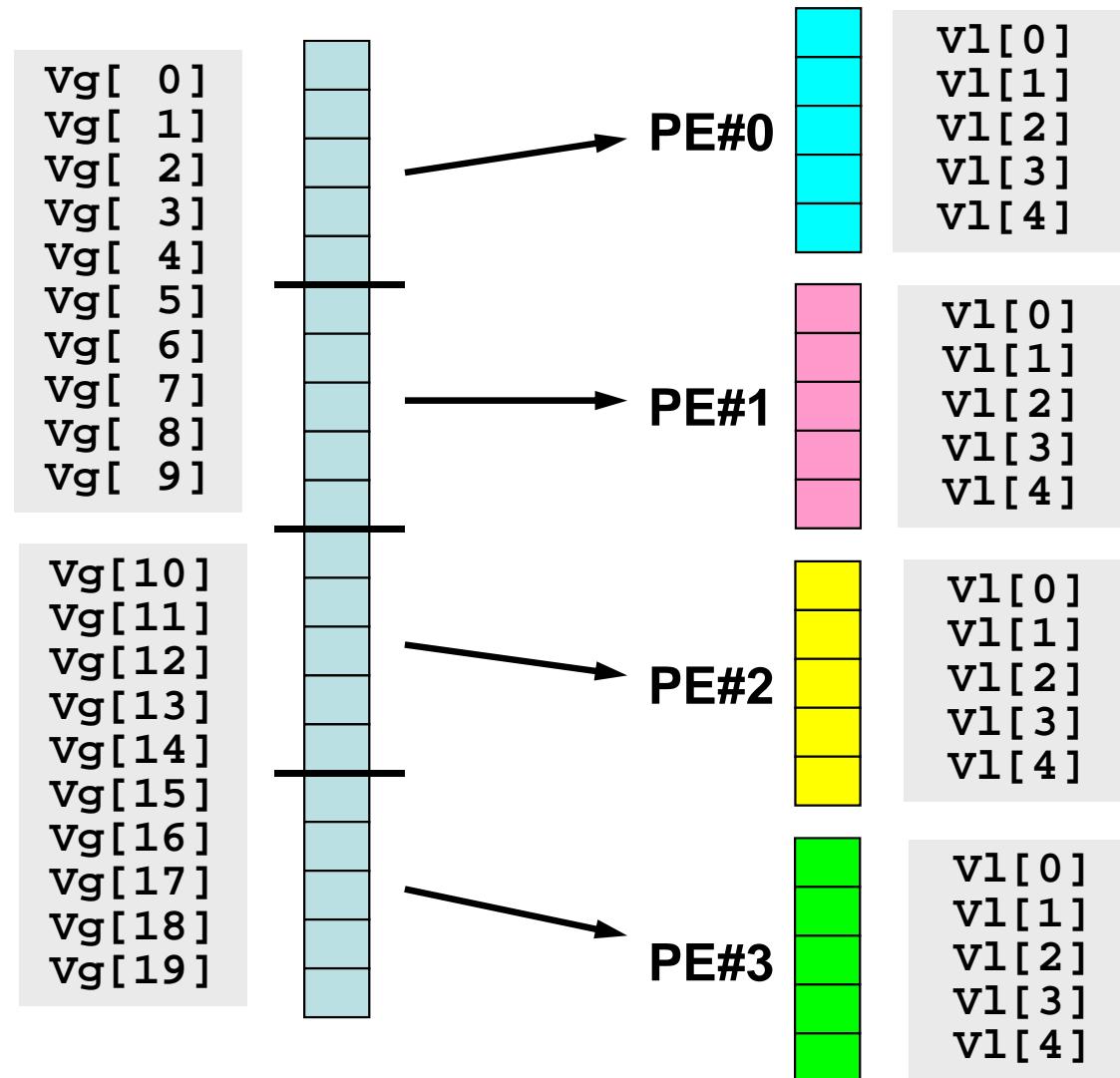
# Local Data (2/2)

- 1<sup>th</sup>-5<sup>th</sup> components of original global vector go to 1<sup>th</sup>-5<sup>th</sup> components of PE#0, 6<sup>th</sup>-10<sup>th</sup> -> PE#1, 11<sup>th</sup>-15<sup>th</sup> -> PE#2, 16<sup>th</sup>-20<sup>th</sup> -> PE#3.



# But ...

- It is too easy !! Just decomposing and renumbering from 1 (or 0).
- Of course, this is not enough. Further examples will be shown in the latter part.



# Example: Dot Product (1/3)

<\$O-S1>/allreduce.c

```
#include <stdio.h>
#include <stdlib.h>
#include "mpi.h"

int main(int argc, char **argv){
    int i,N;
    int PeTot, MyRank;
double VECp[5], VECs[5];
    double sumA, sumR, sum0;

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

    sumA= 0.0;
    sumR= 0.0;

N=5;
for(i=0;i<N;i++){
    VECp[i] = 2.0;
    VECs[i] = 3.0;
}

    sum0 = 0.0;
    for(i=0;i<N;i++){
        sum0 += VECp[i] * VECs[i];
    }
}
```

Local vector is generated at each local process.

# Example: Dot Product (2/3)

<\$O-S1>/allreduce.c

```
MPI_Reduce(&sum0, &sumR, 1, MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD);
MPI_Allreduce(&sum0, &sumA, 1, MPI_DOUBLE, MPI_SUM, MPI_COMM_WORLD);
printf("before BCAST %5d %15.0F %15.0F\n", MyRank, sumA, sumR);

MPI_Bcast(&sumR, 1, MPI_DOUBLE, 0, MPI_COMM_WORLD);
printf("after BCAST %5d %15.0F %15.0F\n", MyRank, sumA, sumR);

MPI_Finalize();

return 0;
}
```

# Example: Dot Product (3/3)

<\$O-S1>/allreduce.c

```
MPI_Reduce(&sum0, &sumR, 1, MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD);  
MPI_Allreduce(&sum0, &sumA, 1, MPI_DOUBLE, MPI_SUM, MPI_COMM_WORLD);
```

## Dot Product

Summation of results of each process (sum0)  
“sumR” has value only on PE#0.

“sumA” has value on all processes by MPI\_Allreduce

```
MPI_Bcast(&sumR, 1, MPI_DOUBLE, 0, MPI_COMM_WORLD);
```

“sumR” has value on PE#1-#3 by MPI\_Bcast

# Execute <\$O-S1>/allreduce.f/c

```
$> mpifccpx -Kfast allreduce.c  
$> mpifrtpx -Kfast allreduce.f  
(modify go4.sh, 4 process)  
$> pbsub go4.sh
```

(my_rank, sumALLREDUCE, sumREDUCE)			
before BCAST	0	1.200000E+02	1.200000E+02
after BCAST	0	1.200000E+02	1.200000E+02
before BCAST	1	1.200000E+02	0.000000E+00
after BCAST	1	1.200000E+02	1.200000E+02
before BCAST	3	1.200000E+02	0.000000E+00
after BCAST	3	1.200000E+02	1.200000E+02
before BCAST	2	1.200000E+02	0.000000E+00
after BCAST	2	1.200000E+02	1.200000E+02

# Examples by Collective Comm.

- Dot Products of Vectors
- Scatter/Gather
- Reading Distributed Files

# Global/Local Data (1/3)

- Parallelization of an easy process where a real number  $\alpha$  is added to each component of real vector **VECg**:

```
do i= 1, NG  
    VECg(i)= VECg(i) + ALPHA  
enddo
```

```
for (i=0; i<NG; i++){  
    VECg[i]= VECg[i] + ALPHA  
}
```

# Global/Local Data (2/3)

- Configuration
  - **NG= 32 (length of the vector)**
  - **ALPHA=1000.**
  - Process # of MPI= 4
- Vector VECg has following 32 components  
( $\langle \$O-S1 \rangle / a1x.all$ ):

(101. 0, 103. 0, 105. 0, 106. 0, 109. 0, 111. 0, 121. 0, 151. 0,
201. 0, 203. 0, 205. 0, 206. 0, 209. 0, 211. 0, 221. 0, 251. 0,
301. 0, 303. 0, 305. 0, 306. 0, 309. 0, 311. 0, 321. 0, 351. 0,
401. 0, 403. 0, 405. 0, 406. 0, 409. 0, 411. 0, 421. 0, 451. 0)

# Global/Local Data (3/3)

- Procedure
  - ① Reading vector **VECg** with length=32 from one process (e.g. 0<sup>th</sup> process)
    - Global Data
  - ② Distributing vector components to 4 MPI processes equally (*i.e.* length= 8 for each processes)
    - Local Data, Local ID/Numbering
  - ③ Adding **ALPHA** to each component of the local vector (with length= 8) on each process.
  - ④ Merging the results to global vector with length= 32.
- Actually, we do not need parallel computers for such a kind of small computation.

# Operations of Scatter/Gather (1/8)

Reading VECg (length=32) from a process (e.g. #0)

- Reading global data from #0 process

```
include    'mpif.h'
integer, parameter :: NG= 32
real(kind=8), dimension(NG) :: VECg

call MPI_INIT (ierr)
call MPI_COMM_SIZE (<comm>, PETOT , ierr)
call MPI_COMM_RANK (<comm>, my_rank, ierr)

if (my_rank.eq.0) then
  open (21, file= 'a1x.all', status= 'unknown')
  do i= 1, NG
    read (21,*) VECg(i)
  enddo
  close (21)
endif
```

```
#include <mpi.h>
#include <stdio.h>
#include <math.h>
#include <assert.h>

int main(int argc, char **argv) {
  int i, NG=32;
  int PeTot, MyRank, MPI_Comm;
  double VECg[32];
  char filename[80];
  FILE *fp;

  MPI_Init(&argc, &argv);
  MPI_Comm_size(<comm>, &PeTot);
  MPI_Comm_rank(<comm>, &MyRank);

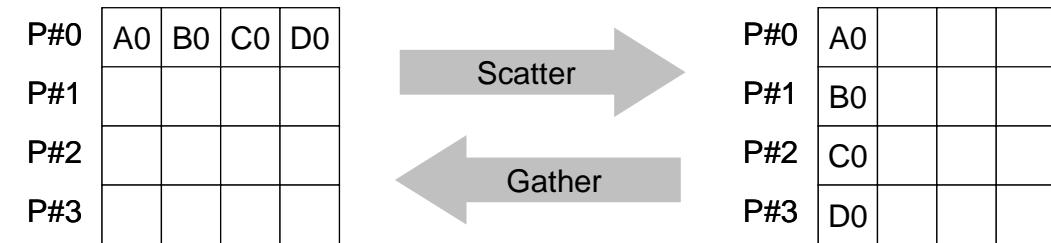
  fp = fopen("a1x.all", "r");
  if (!MyRank) for (i=0; i<NG; i++) {
    fscanf(fp, "%lf", &VECg[i]);
  }
```

# Operations of Scatter/Gather (2/8)

Distributing global data to 4 process equally (*i.e.* length=8 for each process)

- MPI\_Scatter

# MPI\_Scatter



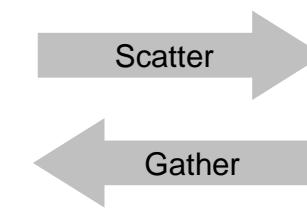
- Sends data from one process to all other processes in a communicator
  - scount-size messages are sent to each process
- `MPI_Scatter (sendbuf, scount, sendtype, recvbuf, rcount, recvtype, root, comm)`**
  - sendbuf** choice I starting address of sending buffer  
**type is defined by "datatype"**
  - scount** int I number of elements sent to each process
  - sendtype** MPI\_Datatype I data type of elements of sending buffer  
FORTRAN: MPI\_INTEGER, MPI\_REAL, MPI\_DOUBLE\_PRECISION, MPI\_CHARACTER etc.  
C: MPI\_INT, MPI\_FLOAT, MPI\_DOUBLE, MPI\_CHAR etc.
  - recvbuf** choice O starting address of receiving buffer
  - rcount** int I number of elements received from the root process
  - recvtype** MPI\_Datatype I data type of elements of receiving buffer  
**rank of root process**
  - root** int I
  - comm** MPI\_Comm I communicator

# MPI\_Scatter (cont.)

- **`MPI_Scatter (sendbuf, scount, sendtype, recvbuf, rcount, recvtype, root, comm)`**

- **`sendbuf`** choice I
  - **`scount`** int I
  - **`sendtype`** MPI\_Datatype I
  - **`recvbuf`** choice O
  - **`rcount`** int I
  - **`recvtype`** MPI\_Datatype I
  - **`root`** int I
  - **`comm`** MPI\_Comm I

P#0	A0	B0	C0	D0
P#1				
P#2				
P#3				



P#0	A0			
P#1	B0			
P#2	C0			
P#3	D0			

starting address of sending buffer  
 number of elements sent to each process  
 data type of elements of sending buffer  
 starting address of receiving buffer  
 number of elements received from the root process  
 data type of elements of receiving buffer  
**rank of root process**  
 communicator

- **Usually**
  - **`scount = rcount`**
  - **`sendtype= recvtype`**
- This function sends **`scount`** components starting from **`sendbuf`** (sending buffer) at process **#root** to each process in **`comm`**. Each process receives **`rcount`** components starting from **`recvbuf`** (receiving buffer).

# Operations of Scatter/Gather (3/8)

Distributing global data to 4 processes equally

- Allocating receiving buffer **VEC** (length=8) at each process.
- 8 components sent from sending buffer **VECg** of process #0 are received at each process #0-#3 as 1<sup>st</sup>-8<sup>th</sup> components of receiving buffer **VEC**.

```
integer, parameter :: N = 8
real(kind=8), dimension(N) :: VEC
...
call MPI_Scatter
    & (VECg, N, MPI_DOUBLE_PRECISION, &
    VEC , N, MPI_DOUBLE_PRECISION, &
    0, <comm>, ierr)
```

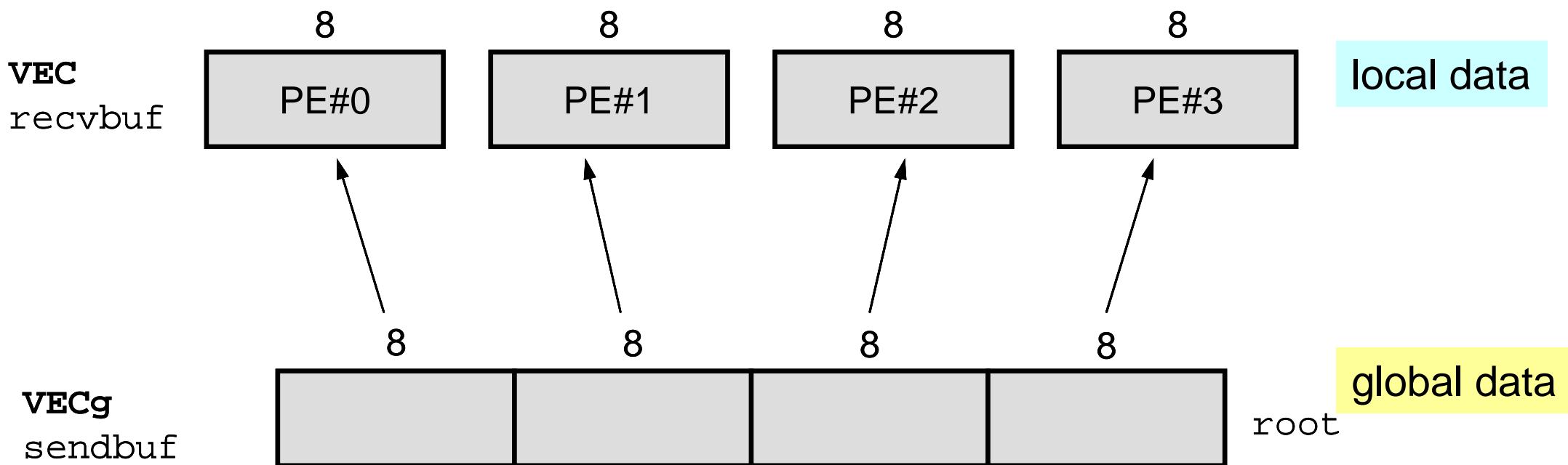
```
int N=8;
double VEC [8];
...
MPI_Scatter (&VECg, N, MPI_DOUBLE, &VEC, N,
MPI_DOUBLE, 0, <comm>);
```

```
call MPI_SCATTER
(sndbuf, scount, sndtype, rcvbuf, rcount,
rcvtype, root, comm, ierr)
```

# Operations of Scatter/Gather (4/8)

Distributing global data to 4 processes equally

- 8 components are scattered to each process from root (#0)
- 1<sup>st</sup>-8<sup>th</sup> components of **VECg** are stored as 1<sup>st</sup>-8<sup>th</sup> ones of **VEC** at **#0**, 9<sup>th</sup>-16<sup>th</sup> components of **VECg** are stored as 1<sup>st</sup>-8<sup>th</sup> ones of **VEC** at **#1**, etc.
  - **VECg**: Global Data, **VEC**: Local Data



# Operations of Scatter/Gather (5/8)

Distributing global data to 4 processes equally

- Global Data: 1<sup>st</sup>-32<sup>nd</sup> components of **VECg** at **#0**
- Local Data: 1<sup>st</sup>-8<sup>th</sup> components of **VEC** at each process
- Each component of **VEC** can be written from each process in the following way:

```
do i= 1, N
    write (*, '(a, 2i8, f10.0)') 'before', my_rank, i, VEC(i)
enddo
```

```
for(i=0;i<N;i++) {
    printf("before %5d %5d %10.0F\n", MyRank, i+1, VEC[i]);}
```

# Operations of Scatter/Gather (5/8)

Distributing global data to 4 processes equally

- Global Data: 1<sup>st</sup>-32<sup>nd</sup> components of **VECg** at **#0**
- Local Data: 1<sup>st</sup>-8<sup>th</sup> components of **VEC** at each process
- Each component of **VEC** can be written from each process in the following way:

PE#0

before 0 1	101.
before 0 2	103.
before 0 3	105.
before 0 4	106.
before 0 5	109.
before 0 6	111.
before 0 7	121.
before 0 8	151.

PE#1

before 1 1	201.
before 1 2	203.
before 1 3	205.
before 1 4	206.
before 1 5	209.
before 1 6	211.
before 1 7	221.
before 1 8	251.

PE#2

before 2 1	301.
before 2 2	303.
before 2 3	305.
before 2 4	306.
before 2 5	309.
before 2 6	311.
before 2 7	321.
before 2 8	351.

PE#3

before 3 1	401.
before 3 2	403.
before 3 3	405.
before 3 4	406.
before 3 5	409.
before 3 6	411.
before 3 7	421.
before 3 8	451.

# Operations of Scatter/Gather (6/8)

On each process, **ALPHA** is added to each of 8 components of **VEC**

- On each process, computation is in the following way

```
real(kind=8), parameter :: ALPHA= 1000.
do i= 1, N
    VEC(i)= VEC(i) + ALPHA
enddo
```

```
double ALPHA=1000. ;
...
for(i=0; i<N; i++) {
    VEC[i]= VEC[i] + ALPHA;}
```

- Results:

## PE#0

after 0 1	1101.
after 0 2	1103.
after 0 3	1105.
after 0 4	1106.
after 0 5	1109.
after 0 6	1111.
after 0 7	1121.
after 0 8	1151.

## PE#1

after 1 1	1201.
after 1 2	1203.
after 1 3	1205.
after 1 4	1206.
after 1 5	1209.
after 1 6	1211.
after 1 7	1221.
after 1 8	1251.

## PE#2

after 2 1	1301.
after 2 2	1303.
after 2 3	1305.
after 2 4	1306.
after 2 5	1309.
after 2 6	1311.
after 2 7	1321.
after 2 8	1351.

## PE#3

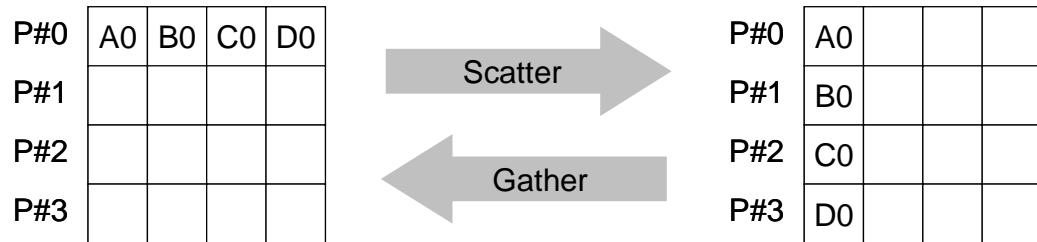
after 3 1	1401.
after 3 2	1403.
after 3 3	1405.
after 3 4	1406.
after 3 5	1409.
after 3 6	1411.
after 3 7	1421.
after 3 8	1451.

# Operations of Scatter/Gather (7/8)

Merging the results to global vector with length= 32

- Using MPI\_Gather (inverse operation to MPI\_Scatter)

# MPI\_Gather



- Gathers together values from a group of processes, inverse operation to **MPI\_Scatter**
- **`MPI_Gather (sendbuf, scount, sendtype, recvbuf, rcount, recvtype, root, comm)`**
  - **`sendbuf`** choice I starting address of sending buffer
  - **`scount`** int I number of elements sent to each process
  - **`sendtype`** MPI\_Datatype I data type of elements of sending buffer
  - **`recvbuf`** choice O starting address of receiving buffer
  - **`rcount`** int I number of elements received from the root process
  - **`recvtype`** MPI\_Datatype I data type of elements of receiving buffer
  - **`root`** int I rank of root process
  - **`comm`** MPI\_Comm I communicator
- **`recvbuf`** is on **root** process.

# Operations of Scatter/Gather (8/8)

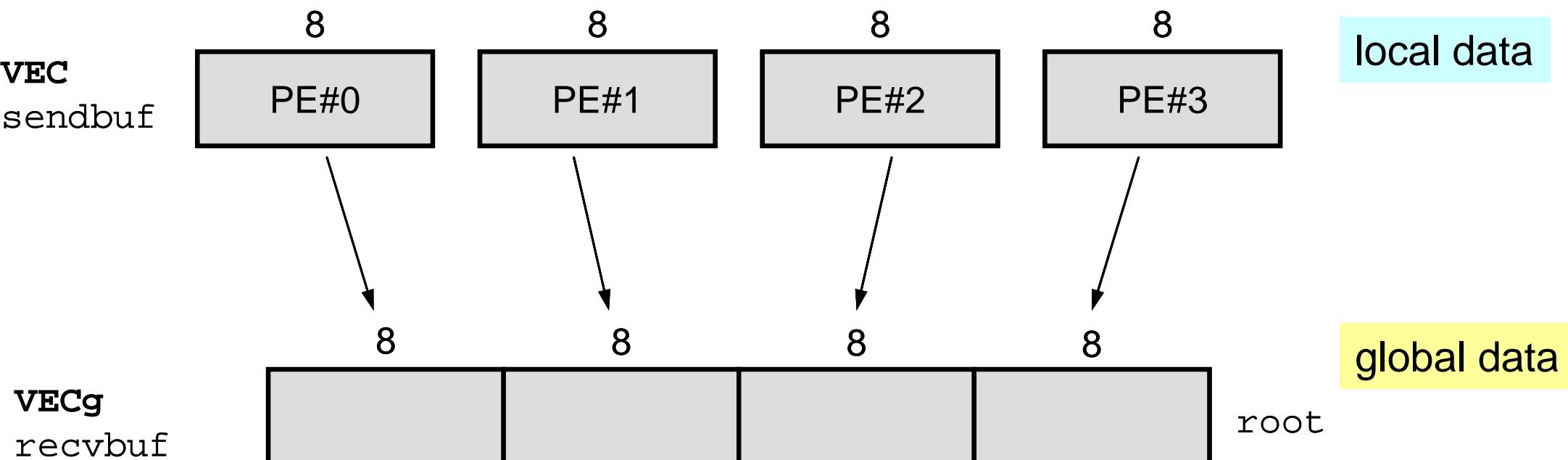
Merging the results to global vector with length= 32

- Each process components of **VEC** to **VECg** on root (#0 in this case).

```
call MPI_Gather          &
  (VEC , N, MPI_DOUBLE_PRECISION, &
   VECg, N, MPI_DOUBLE_PRECISION, &
   0, <comm>, ierr)
```

```
MPI_Gather (&VEC, N, MPI_DOUBLE, &VECg, N,
MPI_DOUBLE, 0, <comm>);
```

- 8 components are gathered from each process to the root process.



# <\$O-S1>/scatter-gather.f/c example

```
$> mpifccpx -Kfast scatter-gather.c  
$> mpifrtpx -Kfast scatter-gather.f  
$> (exec. 4 proc's) go4.sh
```

**PE#0**

before 0 1 101.  
before 0 2 103.  
before 0 3 105.  
before 0 4 106.  
before 0 5 109.  
before 0 6 111.  
before 0 7 121.  
before 0 8 151.

**PE#1**

before 1 1 201.  
before 1 2 203.  
before 1 3 205.  
before 1 4 206.  
before 1 5 209.  
before 1 6 211.  
before 1 7 221.  
before 1 8 251.

**PE#2**

before 2 1 301.  
before 2 2 303.  
before 2 3 305.  
before 2 4 306.  
before 2 5 309.  
before 2 6 311.  
before 2 7 321.  
before 2 8 351.

**PE#3**

before 3 1 401.  
before 3 2 403.  
before 3 3 405.  
before 3 4 406.  
before 3 5 409.  
before 3 6 411.  
before 3 7 421.  
before 3 8 451.

**PE#0**

after 0 1 1101.  
after 0 2 1103.  
after 0 3 1105.  
after 0 4 1106.  
after 0 5 1109.  
after 0 6 1111.  
after 0 7 1121.  
after 0 8 1151.

**PE#1**

after 1 1 1201.  
after 1 2 1203.  
after 1 3 1205.  
after 1 4 1206.  
after 1 5 1209.  
after 1 6 1211.  
after 1 7 1221.  
after 1 8 1251.

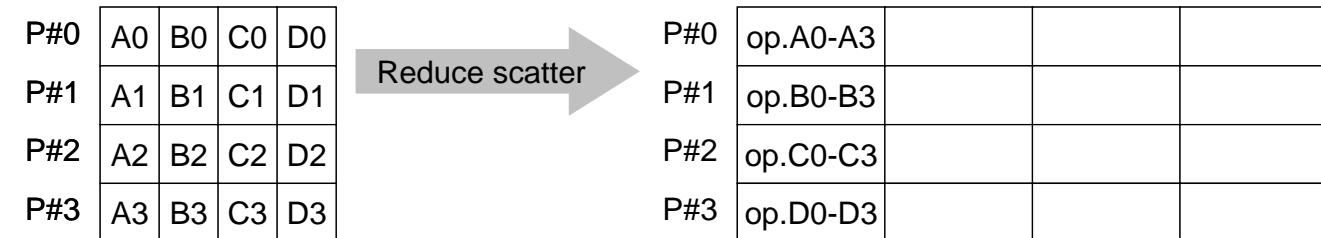
**PE#2**

after 2 1 1301.  
after 2 2 1303.  
after 2 3 1305.  
after 2 4 1306.  
after 2 5 1309.  
after 2 6 1311.  
after 2 7 1321.  
after 2 8 1351.

**PE#3**

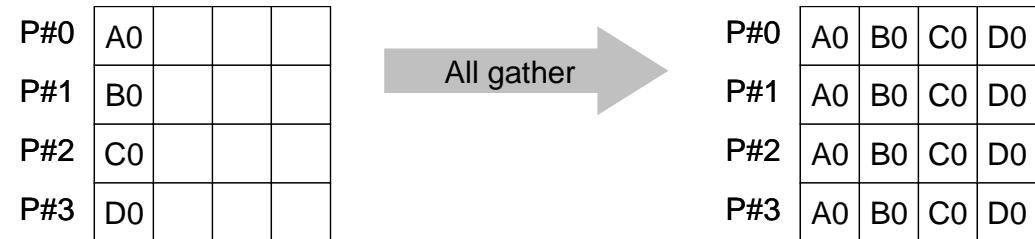
after 3 1 1401.  
after 3 2 1403.  
after 3 3 1405.  
after 3 4 1406.  
after 3 5 1409.  
after 3 6 1411.  
after 3 7 1421.  
after 3 8 1451.

# MPI\_Reduce\_scatter



- MPI\_Reduce + MPI\_Scatter
- **`MPI_Reduce_Scatter (sendbuf, recvbuf, rcount, datatype, op, comm)`**
  - **`sendbuf`** choice I starting address of sending buffer
  - **`recvbuf`** choice O starting address of receiving buffer
  - **`rcount`** int I integer array specifying the number of elements in result distributed to each process. Array must be identical on all calling processes.
  - **`datatype`** MPI\_Datatype I data type of elements of sending/receiving buffer
  - **`op`** MPI\_Op I reduce operation
    - MPI\_MAX, MPI\_MIN, MPI\_SUM, MPI\_PROD, MPI\_LAND, MPI\_BAND etc
  - **`comm`** MPI\_Comm I communicator

# MPI\_Allgather



- MPI\_Gather+MPI\_Bcast
  - Gathers data from all tasks and distribute the combined data to all tasks
- **MPI\_Allgather (sendbuf, scount, sendtype, recvbuf, rcount, recvtype, comm)**
  - **sendbuf** choice I starting address of sending buffer
  - **scount** int I number of elements sent to each process
  - **sendtype** MPI\_Datatype I data type of elements of sending buffer
  - **recvbuf** choice O starting address of receiving buffer
  - **rcount** int I number of elements received from the root process
  - **recvtype** MPI\_Datatype I data type of elements of receiving buffer
  - **comm** MPI\_Comm I communicator

# MPI\_Alltoall

P#0	A0	A1	A2	A3
P#1	B0	B1	B2	B3
P#2	C0	C1	C2	C3
P#3	D0	D1	D2	D3



P#0	A0	B0	C0	D0
P#1	A1	B1	C1	D1
P#2	A2	B2	C2	D2
P#3	A3	B3	C3	D3

- Sends data from all to all processes: transformation of dense matrix
- `MPI_Alltoall (sendbuf, scount, sendtype, recvbuf, rcount, recvtype, comm)`**

- <b><u>sendbuf</u></b>	choice	I	starting address of sending buffer
- <b><u>scount</u></b>	int	I	number of elements sent to each process
- <b><u>sendtype</u></b>	MPI_Datatype	I	data type of elements of sending buffer
- <b><u>recvbuf</u></b>	choice	O	starting address of receiving buffer
- <b><u>rcount</u></b>	int	I	number of elements received from the root process
- <b><u>recvtype</u></b>	MPI_Datatype	I	data type of elements of receiving buffer
- <b><u>comm</u></b>	MPI_Comm	I	communicator

# Examples by Collective Comm.

- Dot Products of Vectors
- Scatter/Gather
- Reading Distributed Files

# Operations of Distributed Local Files

- In Scatter/Gather example, PE#0 reads global data, that is *scattered* to each processer, then parallel operations are done.
- If the problem size is very large, a single processor may not read entire global data.
  - If the entire global data is decomposed to distributed local data sets, each process can read the local data.
  - If global operations are needed to a certain sets of vectors, MPI functions, such as MPI\_Gather etc. are available.

# Reading Distributed Local Files: Uniform Vec. Length (1/2)

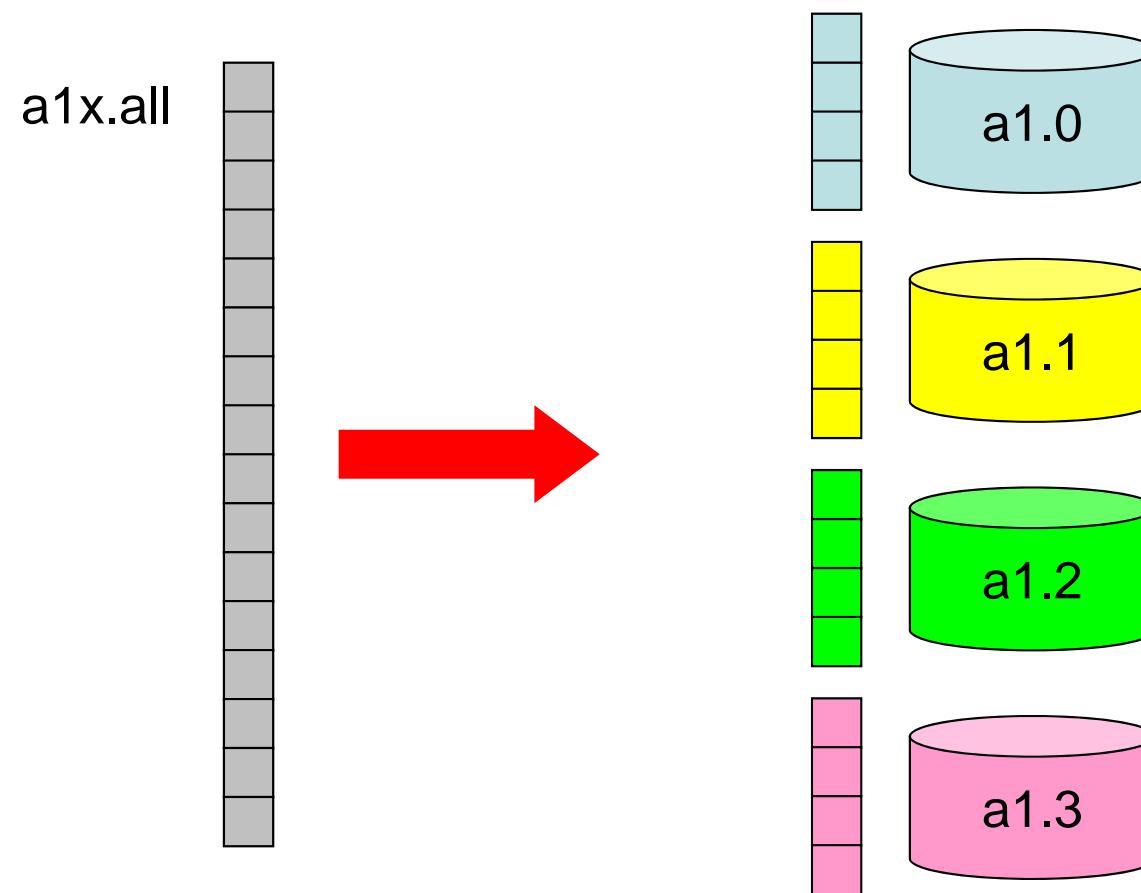
```
>$ cd <$O-S1>
>$ ls a1.*
    a1.0 a1.1 a1.2 a1.3      a1x.all is decomposed to
                                4 files.

>$ mpifccpx -Kfast file.c
>$ mpifrtpx -Kfast file.f

(modify go4.sh for 4 processes)
>$ pbsub go4.sh
```

# Operations of Distributed Local Files

- Local files `a1.0~a1.3` are originally from global file `a1x.all`.



# Reading Distributed Local Files: Uniform Vec. Length (2/2)

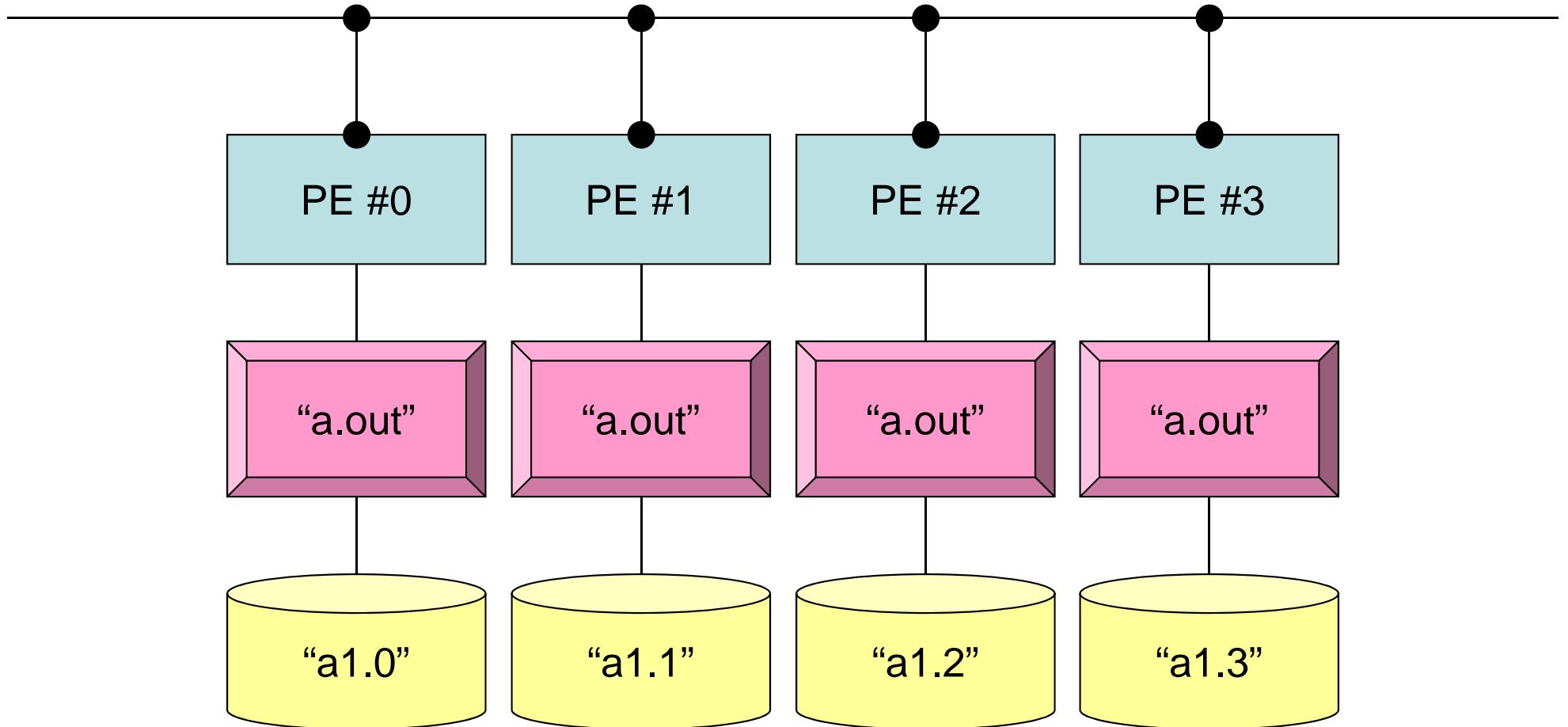
<\$O-S1>/file.c

```
int main(int argc, char **argv){  
    int i;  
    int PeTot, MyRank;  
    MPI_Comm SolverComm;  
    double vec[8];  
    char FileName[80];  
    FILE *fp;  
  
    MPI_Init(&argc, &argv);  
    MPI_Comm_size(MPI_COMM_WORLD, &PeTot);  
    MPI_Comm_rank(MPI_COMM_WORLD, &MyRank);  
  
    sprintf(FileName, "a1.%d", MyRank);  
  
    fp = fopen(FileName, "r");  
    if(fp == NULL) MPI_Abort(MPI_COMM_WORLD, -1) Local ID is 0-7  
    for(i=0;i<8;i++){  
        fscanf(fp, "%lf", &vec[i]); }  
  
    for(i=0;i<8;i++){  
        printf("%5d%5d%10.0f\n", MyRank, i+1, vec[i]);  
    }  
    MPI_Finalize();  
    return 0;  
}
```

Similar to  
“Hello”

Local ID is 0-7

# Typical SPMD Operation



```
mpirun -np 4 a.out
```

# Non-Uniform Vector Length (1/2)

```
>$ cd <$O-S1>
>$ ls a2.*
    a2.0 a2.1 a2.2 a2.3
>$ cat a2.0
    5      Number of Components at each Process
    201.0   Components
    203.0
    205.0
    206.0
    209.0

>$ mpifccpx -Kfast file2.c
>$ mpifrtpx -Kfast file2.f

(modify go4.sh for 4 processes)
>$ pbsub go4.sh
```

# Non-Uniform Vector Length (2/2)

<\$O-S1>/file2.c

```
int main(int argc, char **argv){  
    int i, int PeTot, MyRank;  
    MPI_Comm SolverComm;  
    double *vec, *vec2, *vecg;  
    int num;  
    double sum0, sum;  
    char filename[80];  
    FILE *fp;  
  
    MPI_Init(&argc, &argv);  
    MPI_Comm_size(MPI_COMM_WORLD, &PeTot);  
    MPI_Comm_rank(MPI_COMM_WORLD, &MyRank);  
  
    sprintf(filename, "a2.%d", MyRank);  
    fp = fopen(filename, "r");  
    assert(fp != NULL);  
    "num" is different at each process  
    fscanf(fp, "%d", &num);  
    vec = malloc(num * sizeof(double));  
    for(i=0;i<num;i++){fscanf(fp, "%lf", &vec[i]);}  
  
    for(i=0;i<num;i++){  
        printf(" %5d%5d%5d%10.0f\n", MyRank, i+1, num, vec[i]);}  
  
    MPI_Finalize();  
}
```

# How to generate local data

- Reading global data ( $N=NG$ )
  - Scattering to each process
  - Parallel processing on each process
  - (If needed) reconstruction of global data by gathering local data
- Generating local data ( $N=NL$ ), or reading distributed local data
  - Generating or reading local data on each process
  - Parallel processing on each process
  - (If needed) reconstruction of global data by gathering local data
- In future, latter case is more important, but former case is also introduced in this class for understanding of operations of global/local data.

# Examples by Collective Comm.

- Dot Products of Vectors
- Scatter/Gather
- Reading Distributed Files
- **MPI\_Allgatherv**

# **MPI\_Gatherv, MPI\_Scatterv**

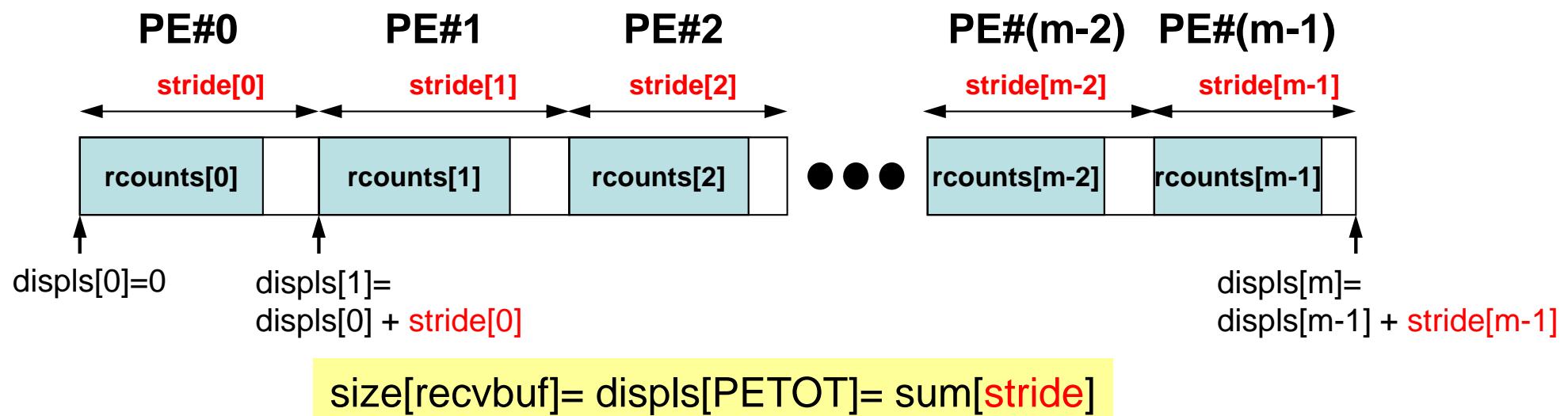
- **MPI\_Gather, MPI\_Scatter**
  - Length of message from/to each process is uniform
- **MPI\_XXXv** extends functionality of **MPI\_XXX** by allowing a varying count of data from each process:
  - **MPI\_Gatherv**
  - **MPI\_Scatterv**
  - **MPI\_Allgatherv**
  - **MPI\_Alltoallv**

# MPI\_Allgatherv

- Variable count version of MPI\_Allgather
  - creates “global data” from “local data”
- **`MPI_Allgatherv (sendbuf, scount, sendtype, recvbuf, rcounts, displs, recvtype, comm)`**
  - **`sendbuf`** choice I starting address of sending buffer
  - **`scount`** int I number of elements sent to each process
  - **`sendtype`** MPI\_Datatype I data type of elements of sending buffer
  - **`recvbuf`** choice O starting address of receiving buffer
  - **`rcounts`** int I integer array (of length group size) containing the number of elements that are to be received from each process  
(array: size= PETOT)
  - **`displs`** int I integer array (of length group size). Entry *i* specifies the displacement (relative to recvbuf ) at which to place the incoming data from process *i* (array: size= PETOT+1)
  - **`recvtype`** MPI\_Datatype I data type of elements of receiving buffer
  - **`comm`** MPI\_Comm I communicator

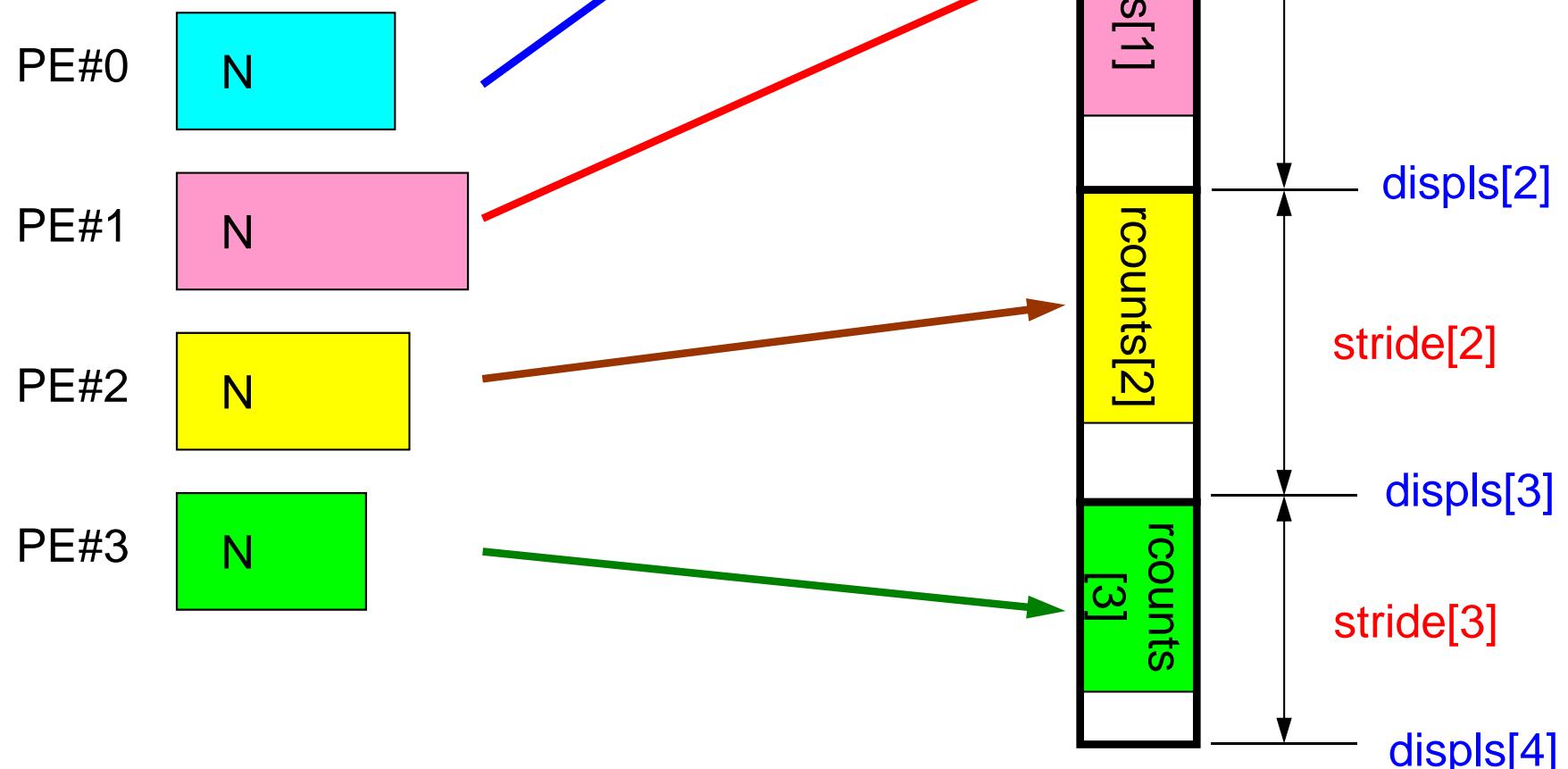
# MPI\_Allgatherv (cont.)

- **`MPI_Allgatherv (sendbuf, scount, sendtype, recvbuf, rcounts, displs, recvtype, comm)`**
  - **`rcounts`** int I integer array (of length group size) containing the number of elements that are to be received from each process (array: size= PETOT)
  - **`displs`** int I integer array (of length group size). Entry  $i$  specifies the displacement (relative to `recvbuf`) at which to place the incoming data from process  $i$  (array: size= PETOT+1)
  - These two arrays are related to size of final “global data”, therefore each process requires information of these arrays (`rcounts`, `displs`)
    - Each process must have same values for all components of both vectors
  - Usually, `stride(i)=rcounts(i)`



# What MPI\_Allgatherv is doing

Generating global data from local data

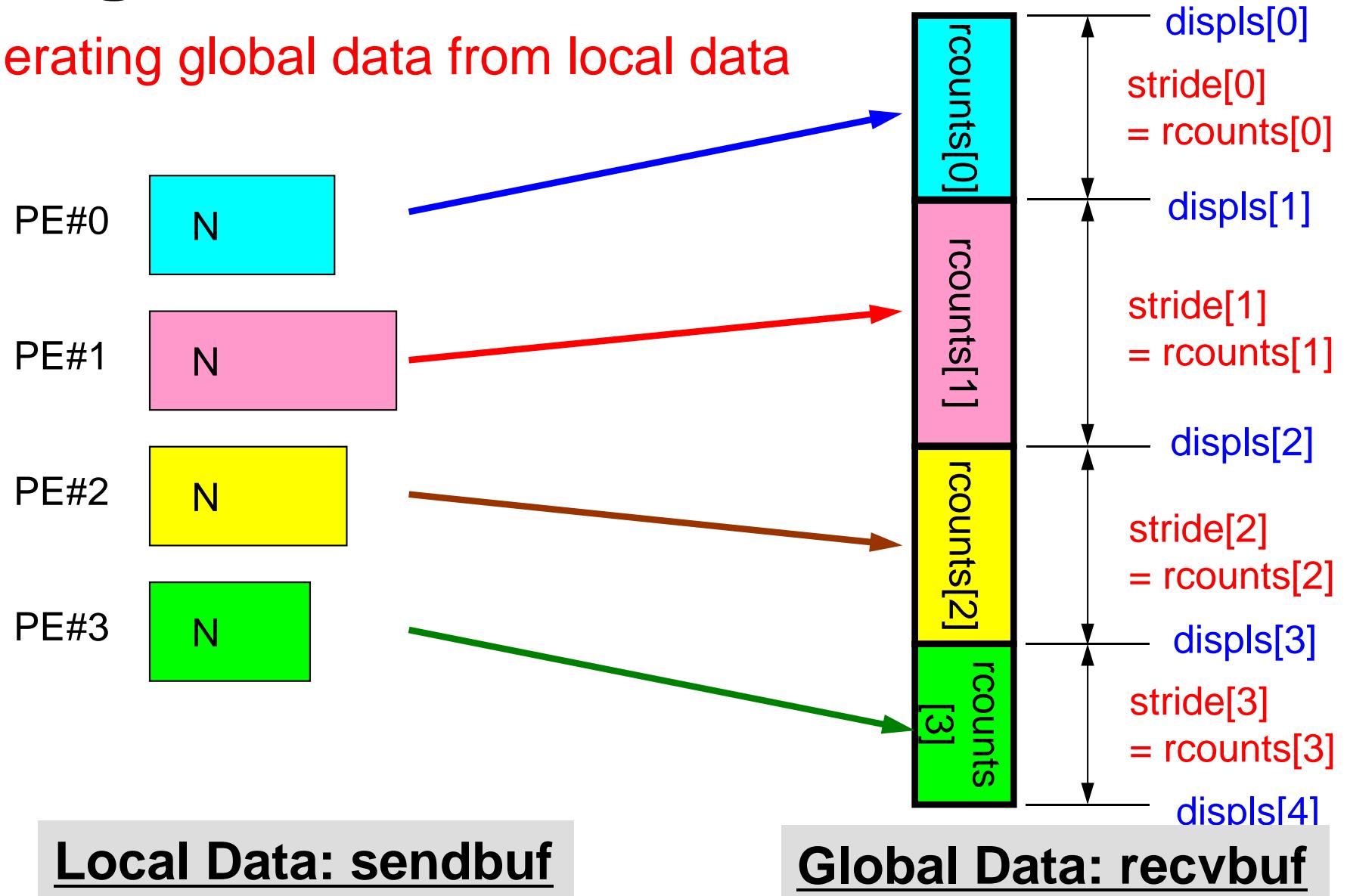


Local Data: sendbuf

Global Data: recvbuf

# What MPI\_Allgatherv is doing

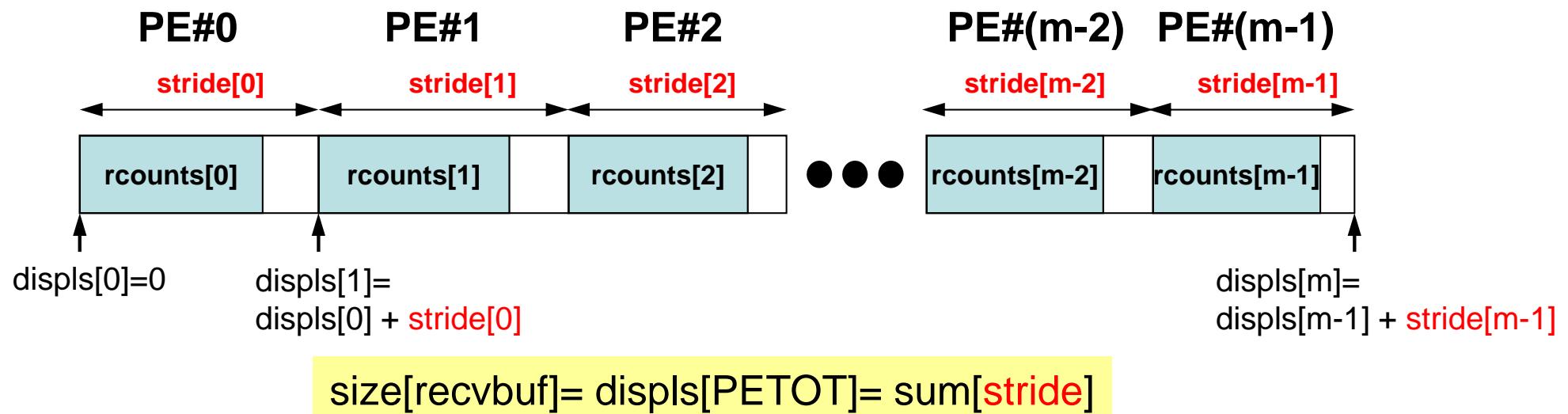
Generating global data from local data



# MPI\_Allgatherv in detail (1/2)

C

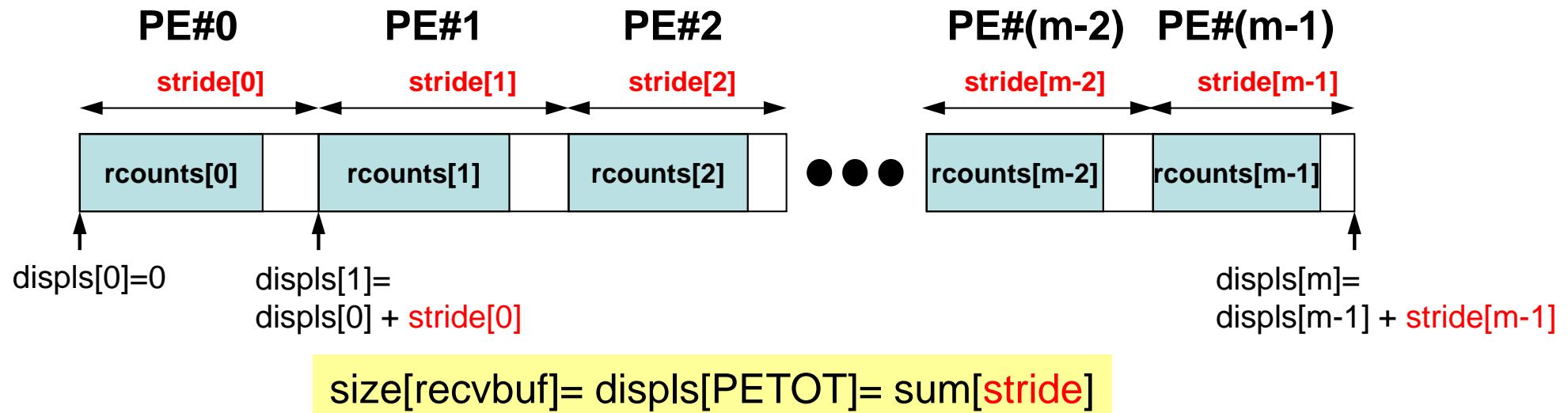
- **`MPI_Allgatherv`** (`sendbuf`, `scount`, `sendtype`, `recvbuf`, `rcounts`, `displs`, `recvtype`, `comm`)
- **`rcounts`**
  - Size of message from each PE: Size of Local Data (Length of Local Vector)
- **`displs`**
  - Address/index of each local data in the vector of global data
  - `displs(PETOT+1)`= Size of Entire Global Data (Global Vector)



# MPI\_Allgatherv in detail (2/2)

C

- Each process needs information of **rcounts** & **displs**
  - “**rcounts**” can be created by gathering local vector length “**N**” from each process.
  - On each process, “**displs**” can be generated from “**rcounts**” on each process.
    - `stride[i] = rcounts[i]`
  - Size of “**recvbuf**” is calculated by summation of “**rcounts**”.



# Preparation for MPI\_Allgatherv

## <\$O-S1>/agv.c

- Generating global vector from “a2.0”~”a2.3”.
- Length of the each vector is 8, 5, 7, and 3, respectively. Therefore, size of final global vector is 23 (= 8+5+7+3).

# a2.0~a2.3

## PE#0

8  
101.0  
103.0  
105.0  
106.0  
109.0  
111.0  
121.0  
151.0

## PE#1

5  
201.0  
203.0  
205.0  
206.0  
209.0

## PE#2

7  
301.0  
303.0  
305.0  
306.0  
311.0  
321.0  
351.0

## PE#3

3  
401.0  
403.0  
405.0

# Preparation: MPI\_Allgatherv (1/4)

<\$O-S1>/agv.c

```

int main(int argc, char **argv){
    int i;
    int PeTot, MyRank;
    MPI_Comm SolverComm;
    double *vec, *vec2, *vecg;
    int *Rcounts, *Displs;
    int n;
    double sum0, sum;
    char filename[80];
    FILE *fp;

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

    sprintf(filename, "a2.%d", MyRank);
    fp = fopen(filename, "r");
    assert(fp != NULL);

    fscanf(fp, "%d", &n);
    vec = malloc(n * sizeof(double));
    for(i=0;i<n;i++){
        fscanf(fp, "%lf", &vec[i]);
    }
}

```

`n`(NL) is different at each process

# Preparation: MPI\_Allgatherv (2/4)

<\$O-S1>/agv.c

```
Rcounts= calloc(PeTot, sizeof(int));
Displs = calloc(PeTot+1, sizeof(int));

printf("before %d %d", MyRank, n);
for(i=0;i<PeTot;i++){printf(" %d", Rcounts[i]);}

MPI_Allgather(&n, 1, MPI_INT, Rcounts, 1, MPI_INT, MPI_COMM_WORLD);

printf("after  %d %d", MyRank, n);
for(i=0;i<PeTot;i++){printf(" %d", Rcounts[i]);}          Rcounts on each PE

Displs[0] = 0;
```

PE#0 N=8

PE#1 N=5

PE#2 N=7

PE#3 N=3



**MPI\_Allgather**

Rcounts[0:3]= {8, 5, 7, 3}

Rcounts[0:3]={8, 5, 7, 3}

Rcounts[0:3]={8, 5, 7, 3}

Rcounts[0:3]={8, 5, 7, 3}

# Preparation: MPI\_Allgatherv (2/4)

<\$O-S1>/agv.c

```
Rcounts= calloc(PeTot, sizeof(int));
Displs = calloc(PeTot+1, sizeof(int));

printf("before %d %d", MyRank, n);
for(i=0;i<PeTot;i++){printf(" %d", Rcounts[i]);}

MPI_Allgather(&n, 1, MPI_INT, Rcounts, 1, MPI_INT, MPI_COMM_WORLD);

printf("after  %d %d", MyRank, n);
for(i=0;i<PeTot;i++){printf(" %d", Rcounts[i]);}      Rcounts on each PE

Displs[0] = 0;
for(i=0;i<PeTot;i++){
        Displs[i+1] = Displs[i] + Rcounts[i];}

    printf("CoundIndex  %d ", MyRank);                  Displs on each PE
    for(i=0;i<PeTot+1;i++){
        printf(" %d", Displs[i]);
    }
    MPI_Finalize();
    return 0;
}
```

# Preparation: MPI\_Allgatherv (3/4)

```
> cd <$O-S1>
> mpifccpx -Kfast agv.c

(modify go4.sh for 4 processes)
> pjsub go4.sh
```

before	0	8	0	0	0	0
after	0	8	8	5	7	3
Displs	0	0	8	13	20	23
before	1	5	0	0	0	0
after	1	5	8	5	7	3
Displs	1	0	8	13	20	23
before	3	3	0	0	0	0
after	3	3	8	5	7	3
Displs	3	0	8	13	20	23
before	2	7	0	0	0	0
after	2	7	8	5	7	3
Displs	2	0	8	13	20	23

# Preparation: MPI\_Allgatherv (4/4)

- Only "recvbuf" is not defined yet.
- Size of "recvbuf" = "Displs[ PETOT ] "

```
MPI_Allgatherv
  ( VEC , N, MPI_DOUBLE,
    recvbuf, rcounts, displs, MPI_DOUBLE,
    MPI_COMM_WORLD );
```

# Report S1 (1/2)

- Deadline: 17:00 February 15<sup>th</sup> (Sa), 2014.
  - Send files via e-mail at `nakajima(at)cc.u-tokyo.ac.jp`
- Problem S1-1
  - Read local files `<$O-S1>/a1.0~a1.3`, `<$O-S1>/a2.0~a2.3`.
  - Develop codes which calculate norm  $\|x\|$  of global vector for each case.
    - `<$O-S1>file.c`, `<$O-S1>file2.c`
- Problem S1-2
  - Read local files `<$O-S1>/a2.0~a2.3`.
  - Develop a code which constructs “global vector” using `MPI_Allgatherv`.

# Report S1 (2/2)

- Problem S1-3
  - Develop parallel program which calculates the following numerical integration using “trapezoidal rule” by MPI\_Reduce, MPI\_Bcast etc.
  - Measure computation time, and parallel performance

$$\int_0^1 \frac{4}{1+x^2} dx$$

- Report
  - Cover Page: Name, ID, and Problem ID (S1) must be written.
  - Less than two pages including figures and tables (A4) for each of three sub-problems
    - Strategy, Structure of the Program, Remarks
  - Source list of the program (if you have bugs)
  - Output list (as small as possible)

- What is MPI ?
- Your First MPI Program: Hello World
- Global/Local Data
- Collective Communication
- **Peer-to-Peer Communication**

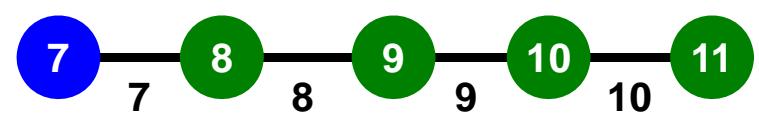
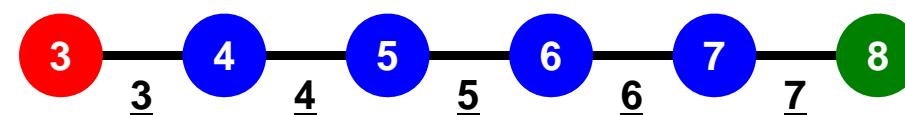
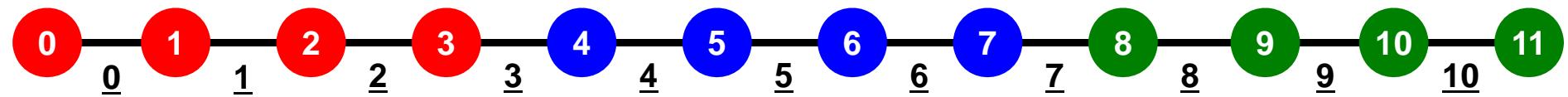
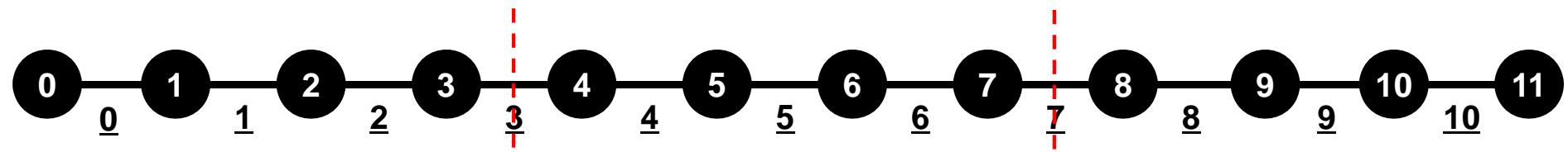
# Peer-to-Peer Communication

## Point-to-Point Communication

### 1対1通信

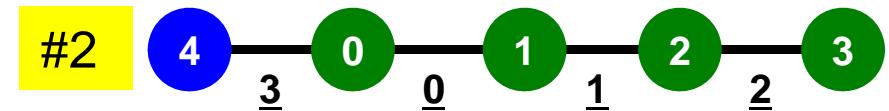
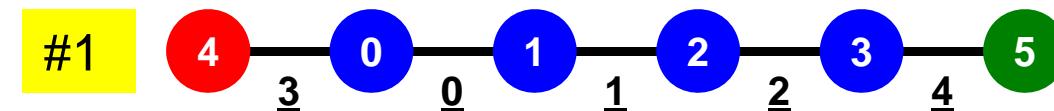
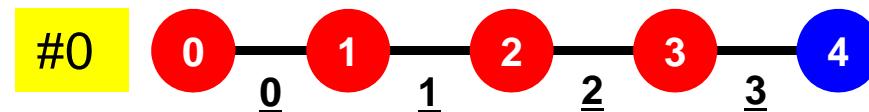
- What is P2P Communication ?
- 2D Problem, Generalized Communication Table
- Report S2

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



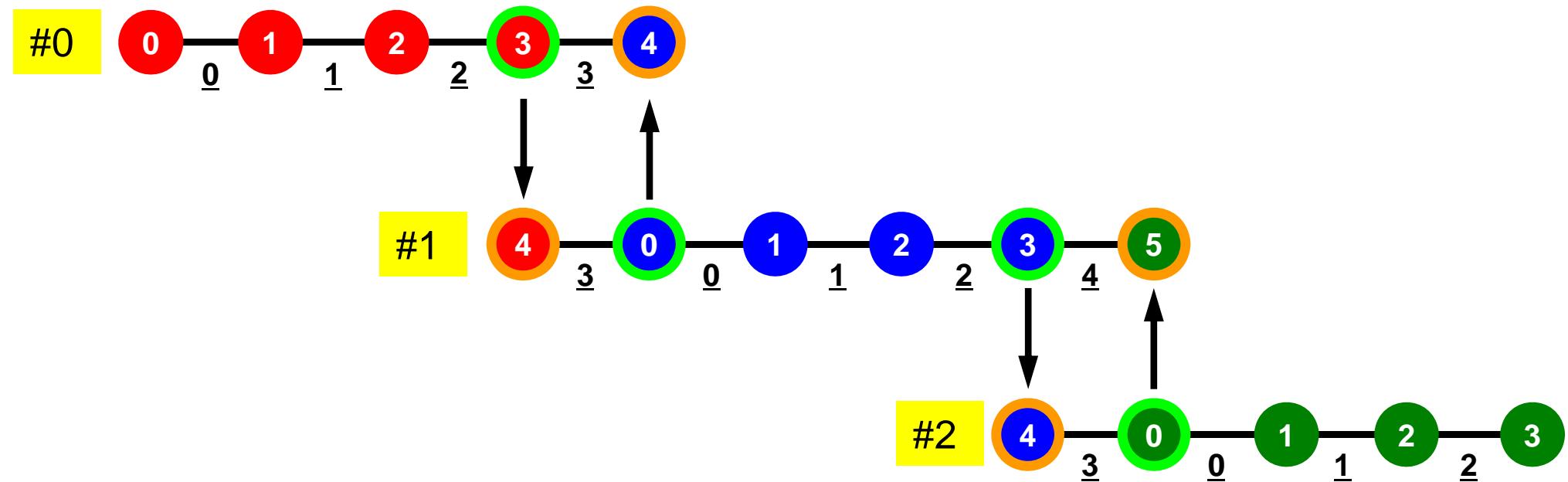
# 1D FEM: 12 nodes/11 elem's/3 domains

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, ...
    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

# Preconditioning, DAXPY

Local Operations by Only Internal Points: Parallel Processing  
is possible

```
/*
//-- {z} = [M-1] {r}
*/
    for (i=0; i<N; i++) {
        W[Z][i] = W[DD][i] * W[R][i];
    }
```

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



# 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];
}

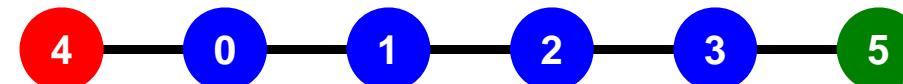
Alpha = Rho / C1;
```



# 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

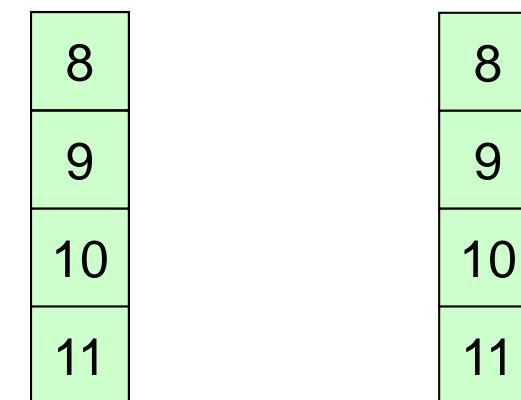
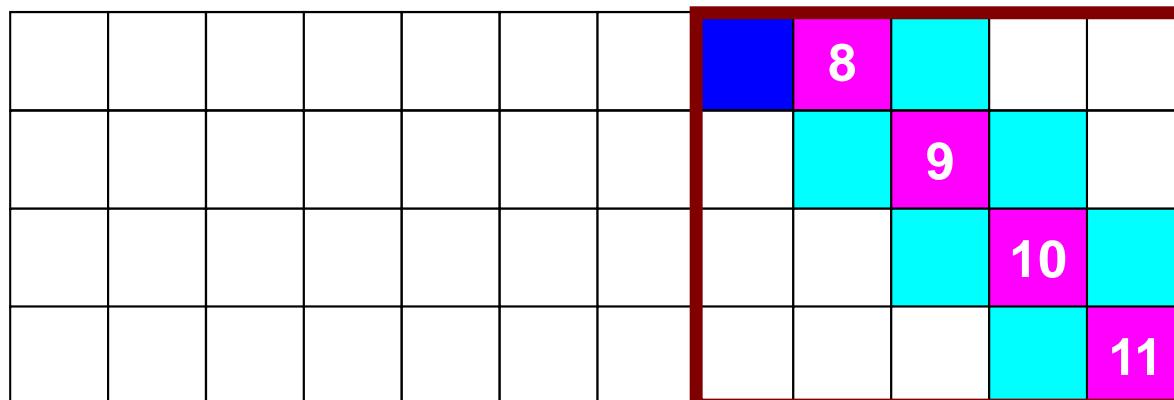
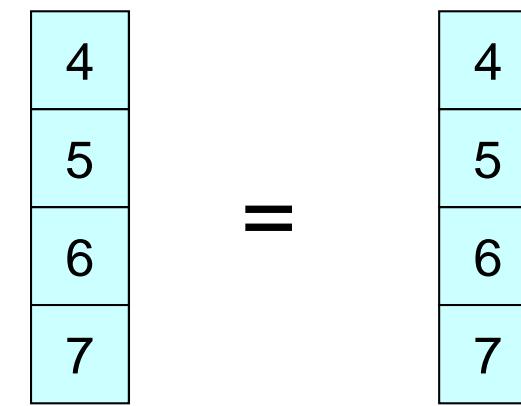
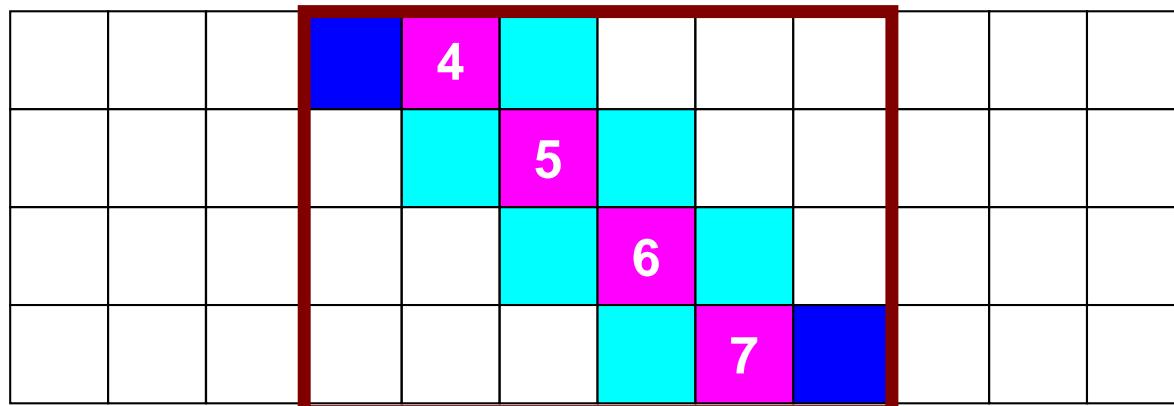
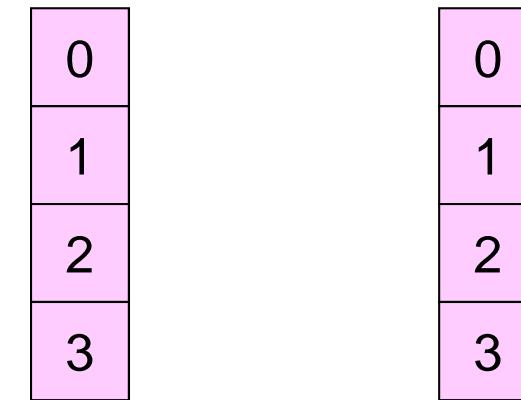
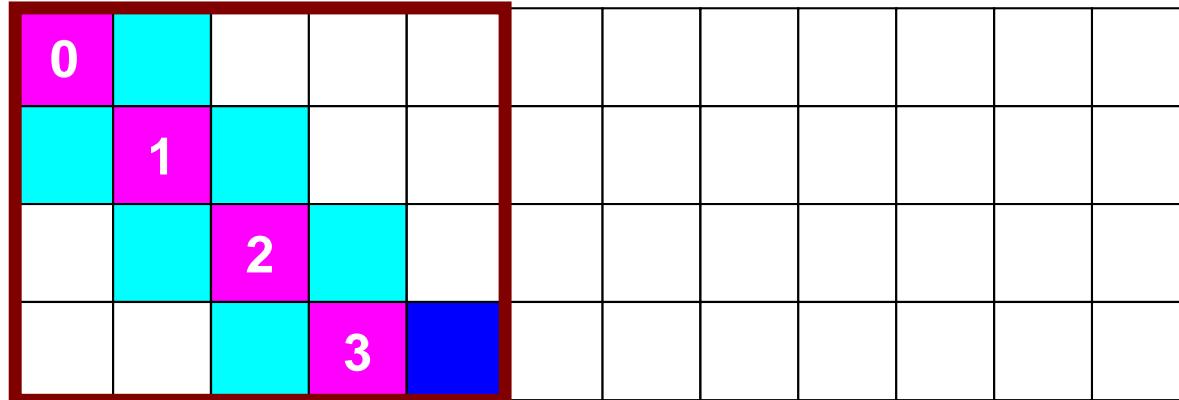
0												
	1											
		2										
			3									
				4								
					5							
						6						
							7					
								8				
									9			
										10		
											11	

=

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

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

# Mat-Vec Products: Local Op. Possible



# Mat-Vec Products: Local Op. Possible

0				
	1			
		2		
			3	

0
1
2
3

0
1
2
3

	0					
		1				
			2			
				3		

0
1
2
3

0
1
2
3

=

	0				
		1			
			2		
				3	

0
1
2
3

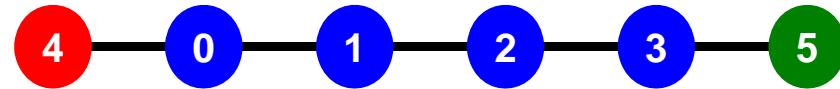
0
1
2
3

# Mat-Vec Products: Local Op. #1

$$\begin{array}{|c|c|c|c|c|c|} \hline & \textcolor{blue}{\boxed{}} & \textcolor{magenta}{\boxed{0}} & \textcolor{cyan}{\boxed{}} & \textcolor{white}{\boxed{}} & \textcolor{white}{\boxed{}} & \textcolor{white}{\boxed{}} \\ \hline \textcolor{white}{\boxed{}} & \textcolor{cyan}{\boxed{}} & \textcolor{magenta}{\boxed{1}} & \textcolor{cyan}{\boxed{}} & \textcolor{white}{\boxed{}} & \textcolor{white}{\boxed{}} & \textcolor{white}{\boxed{}} \\ \hline \textcolor{white}{\boxed{}} & \textcolor{cyan}{\boxed{}} & \textcolor{magenta}{\boxed{2}} & \textcolor{cyan}{\boxed{}} & \textcolor{cyan}{\boxed{}} & \textcolor{white}{\boxed{}} & \textcolor{white}{\boxed{}} \\ \hline \textcolor{white}{\boxed{}} & \textcolor{white}{\boxed{}} & \textcolor{cyan}{\boxed{}} & \textcolor{magenta}{\boxed{3}} & \textcolor{cyan}{\boxed{}} & \textcolor{blue}{\boxed{}} & \textcolor{white}{\boxed{}} \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline 0 \\ \hline 1 \\ \hline 2 \\ \hline 3 \\ \hline \end{array}$$

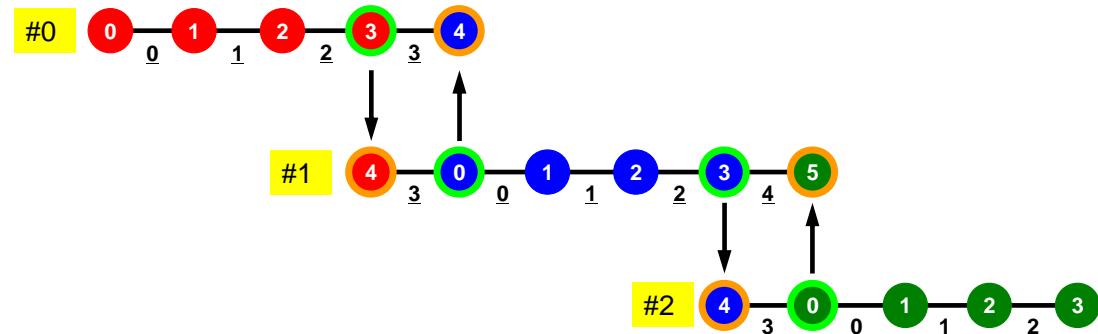


$$\begin{array}{|c|c|c|c|c|c|} \hline \textcolor{magenta}{\boxed{0}} & \textcolor{cyan}{\boxed{}} & \textcolor{white}{\boxed{}} & \textcolor{white}{\boxed{}} & \textcolor{blue}{\boxed{}} & \textcolor{white}{\boxed{}} \\ \hline \textcolor{cyan}{\boxed{}} & \textcolor{magenta}{\boxed{1}} & \textcolor{cyan}{\boxed{}} & \textcolor{white}{\boxed{}} & \textcolor{white}{\boxed{}} & \textcolor{white}{\boxed{}} \\ \hline \textcolor{white}{\boxed{}} & \textcolor{cyan}{\boxed{}} & \textcolor{magenta}{\boxed{2}} & \textcolor{cyan}{\boxed{}} & \textcolor{white}{\boxed{}} & \textcolor{white}{\boxed{}} \\ \hline \textcolor{white}{\boxed{}} & \textcolor{white}{\boxed{}} & \textcolor{cyan}{\boxed{3}} & \textcolor{magenta}{\boxed{}} & \textcolor{white}{\boxed{}} & \textcolor{blue}{\boxed{}} \\ \hline \end{array} = \begin{array}{|c|c|c|c|c|} \hline 0 \\ \hline 1 \\ \hline 2 \\ \hline 3 \\ \hline 4 \\ \hline 5 \\ \hline \end{array}$$



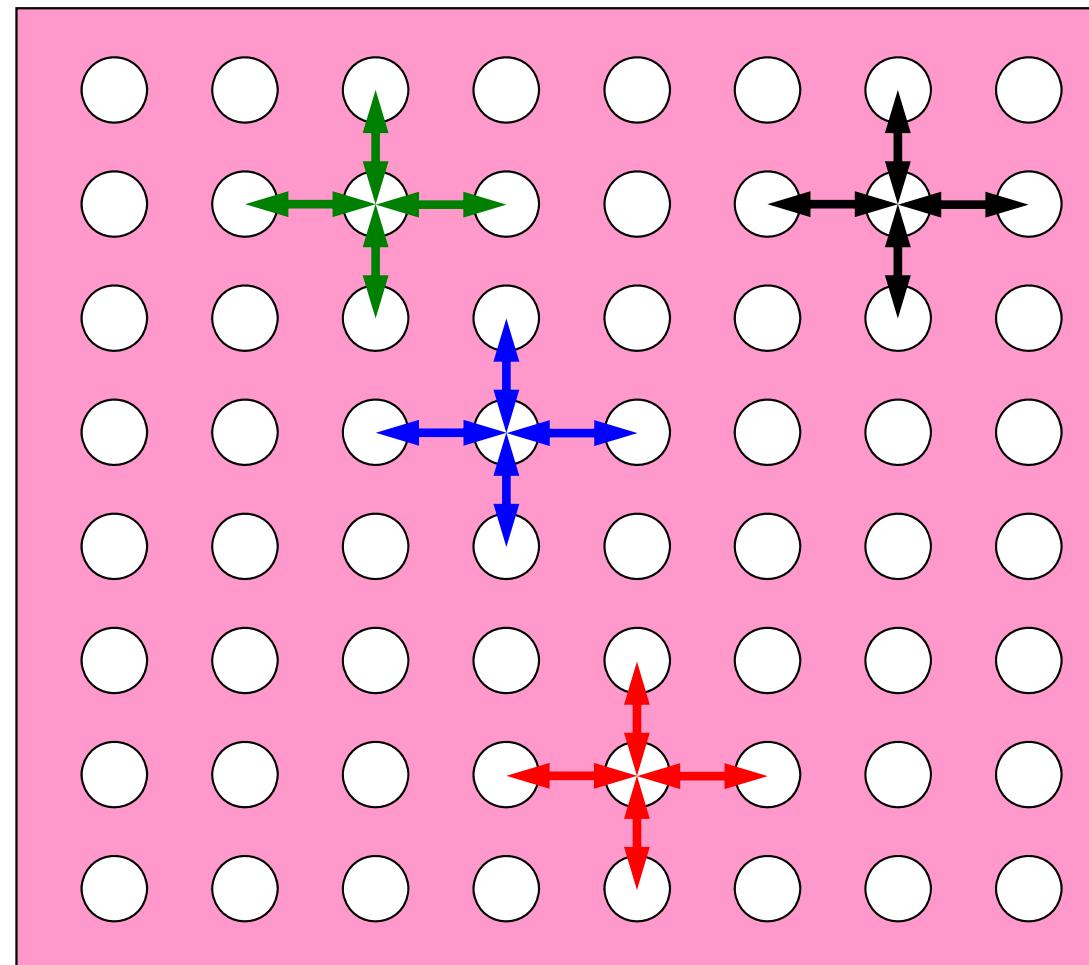
# What is Peer-to-Peer Communication ?

- 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
- Peer-toPeer/Point-to-Point
  - MPI\_Send, MPI\_Receive
  - Communication with limited processes
    - Neighbors
  - Application Area
    - FEM, FDM: Localized Method



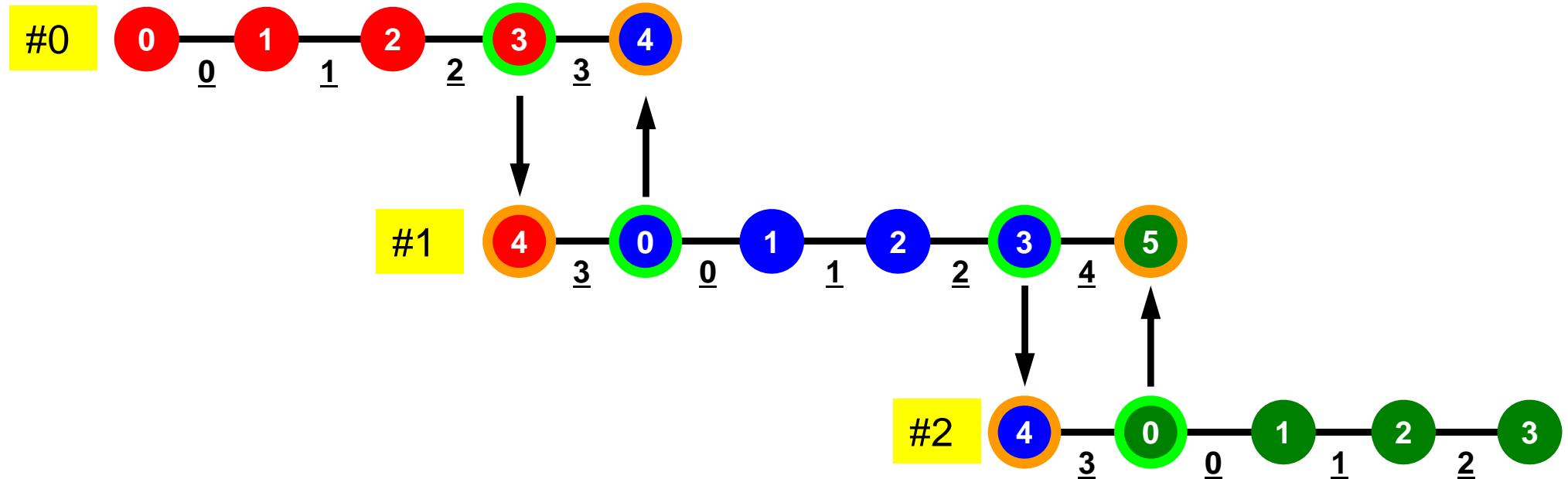
# Collective/P2P Communications

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



# When do we need P2P comm.: 1D-FEM

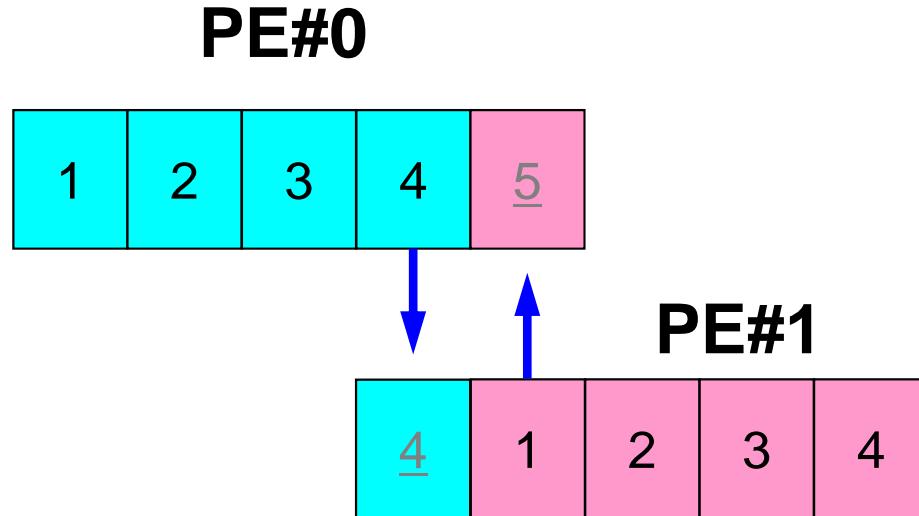
Info in neighboring domains is required for FEM operations  
Matrix assembling, Iterative Method



# Method for P2P 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”.

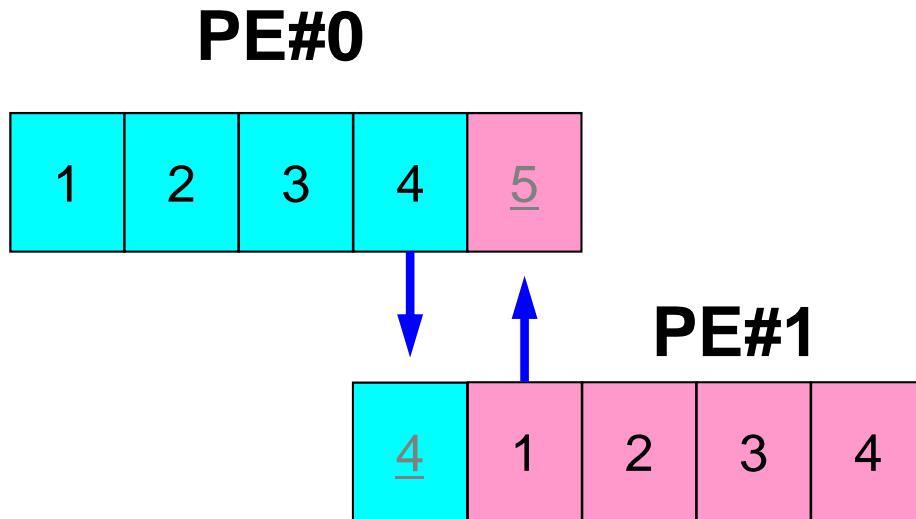
# MPI\_Send/MPI\_Recv



```
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.)

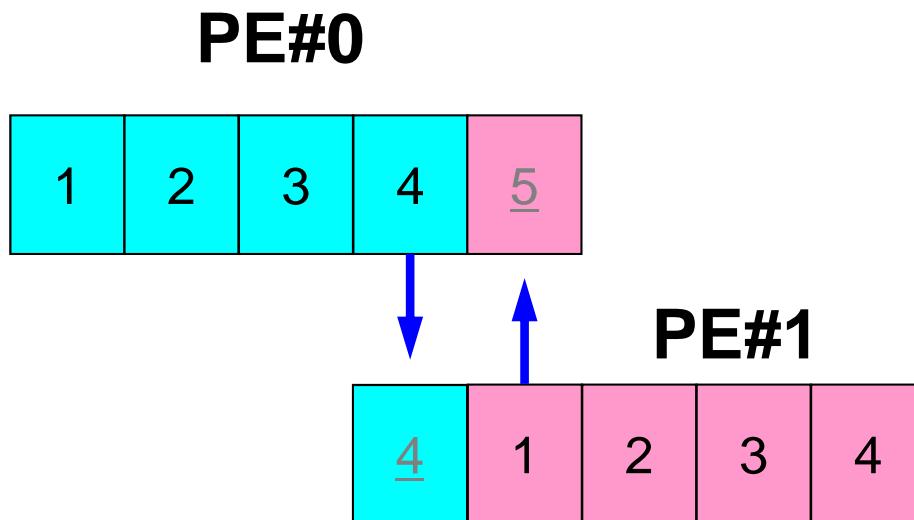


```
if (my_rank.eq.0) NEIB_ID=1  
if (my_rank.eq.1) NEIB_ID=0  
...  
if (my_rank.eq.0) then  
  call MPI_SEND (NEIB_ID, arg's)  
  call MPI_RECV (NEIB_ID, arg's)  
endif  
  
if (my_rank.eq.1) then  
  call MPI_RECV (NEIB_ID, arg's)  
  call MPI_SEND (NEIB_ID, arg's)  
endif  
...
```

- It works ... but

# How to do P2P 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 )**

– <b><u>sendbuf</u></b>	choice	I	starting address of sending buffer
– <b><u>count</u></b>	int	I	number of elements in sending buffer
– <b><u>datatype</u></b>	MPI_Datatype	I	datatype of each sending buffer element
– <b><u>dest</u></b>	int	I	rank of destination
– <b><u>tag</u></b>	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),
– <b><u>comm</u></b>	MPI_Comm	I	communicator
– <b><u>request</u></b>	MPI_Request	O	communication request array used in MPI_Waitall

# Communication Request: request

- **MPI\_Isend**

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

- <u>sendbuf</u>	choice	I	starting address of sending buffer
- <u>count</u>	int	I	number of elements in sending buffer
- <u>datatype</u>	MPI_Datatype	I	datatype of each sending buffer element
- <u>dest</u>	int	I	rank of destination
- <u>tag</u>	int	I	message tag  This integer can be used by the application to distinguish messages. Communication occurs if tag's of MPI_Irecv and MPI_Irecv are matched. Usually tag is set to be "0" (in this class),
- <u>comm</u>	MPI_Comm	I	communicator
- <u>request</u>	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`)**

– <u><code>recvbuf</code></u>	choice	I	starting address of receiving buffer
– <u><code>count</code></u>	int	I	number of elements in receiving buffer
– <u><code>datatype</code></u>	MPI_Datatype	I	datatype of each receiving buffer element
– <u><code>source</code></u>	int	I	rank of source
– <u><code>tag</code></u>	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),
– <u><code>comm</code></u>	MPI_Comm	I	communicator
– <u><code>request</code></u>	MPI_Request	O	communication request array used in MPI_Waitall

# MPI\_Waitall

C

- **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** (request[2\*count])
- 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\_Waitall** (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`

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

# 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'

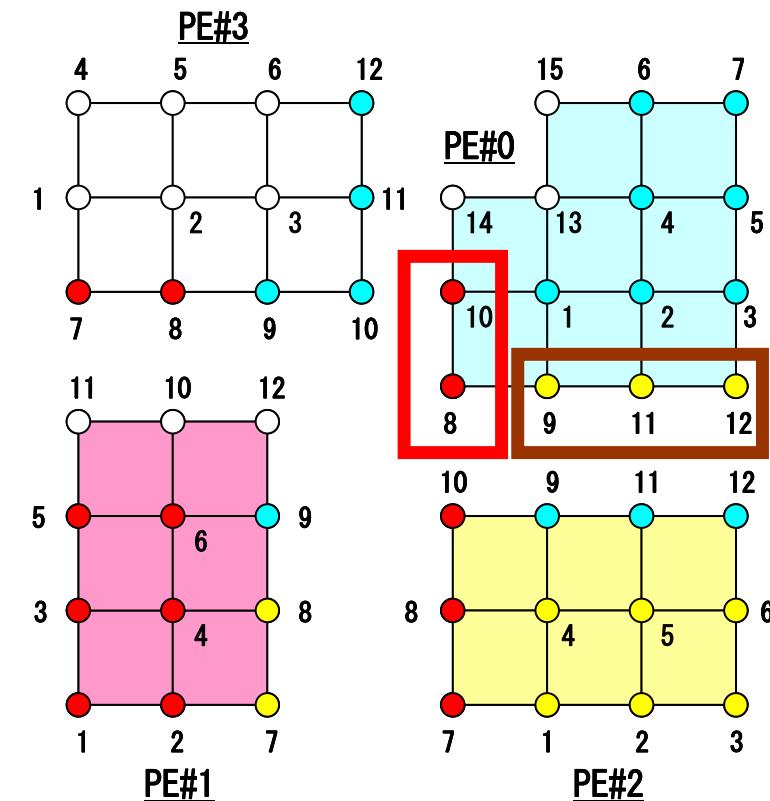
# RECV: receiving to external nodes

Recv. continuous data to recv. buffer from neighbors

- **`MPI_Irecv`**

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

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



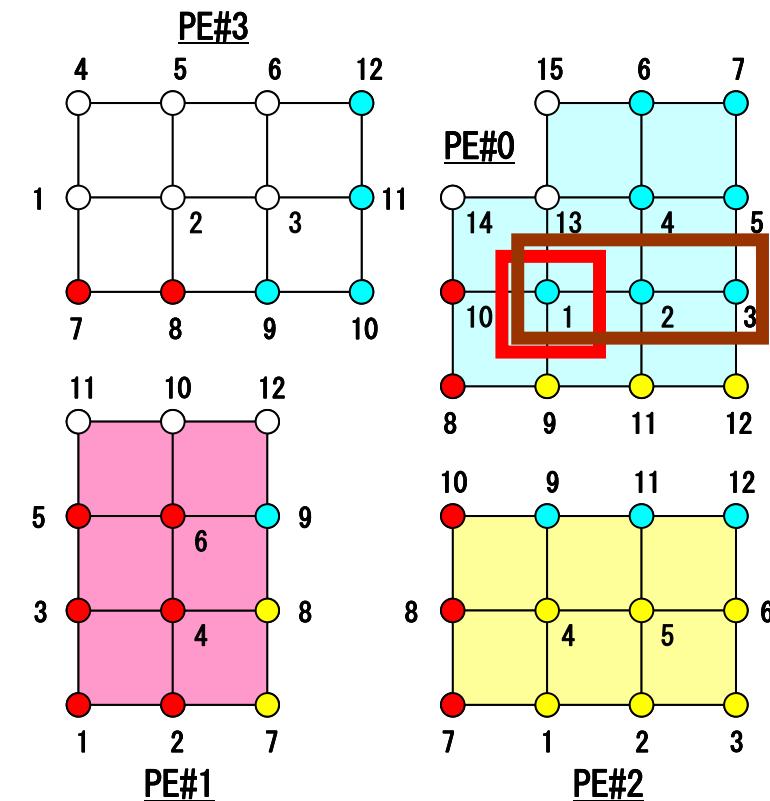
# SEND: sending from boundary nodes

Send continuous data to send buffer of neighbors

- **`MPI_Isend`**

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

<u><code>sendbuf</code></u>	choice	I	starting address of sending buffer
<u><code>count</code></u>	int	I	number of elements in sending buffer
<u><code>datatype</code></u>	MPI_Datatype	I	datatype of each sending buffer element
<u><code>dest</code></u>	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 Oakleaf-FX

## Copy

```
>$ cd <$O-TOP>
>$ cp /home/z30088/fem2/C/s2.tar .
>$ tar xvf s2.tar
```

## Confirm Directory

```
>$ ls
mpi

>$ cd mpi/S2
```

This directory is called as <\$O-S2> in this course.

<\$O-S2> = <\$O-fem2>/mpi/S2

# Ex.1: Send-Recv a Scalar

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

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

call MPI_Isend (VAL ,1,MPI_DOUBLE_PRECISION,NEIB,...,req_send,...)
call MPI_Irecv (VALtemp,1,MPI_DOUBLE_PRECISION,NEIB,...,req_recv,...)
call MPI_Waitall (...,req_recv,stat_recv,...): Recv.buf VALtemp can be used
call 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

```
$> cd <$O-S2>
$> mpifccpx -Kfast ex1-1.c
$> pbsub 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.1: Send-Recv a Scalar

## SendRecv

```
$> cd <$O-S2>
$> mpifccpx -Kfast ex1-2.c
$> pbsub go2.sh
```

```
#include <stdio.h>
#include <stdlib.h>
#include "mpi.h"
int main(int argc, char **argv){
    int neib;
    int MyRank, PeTot;
    double VAL, VALtemp;
    MPI_Status *StatSR;

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

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

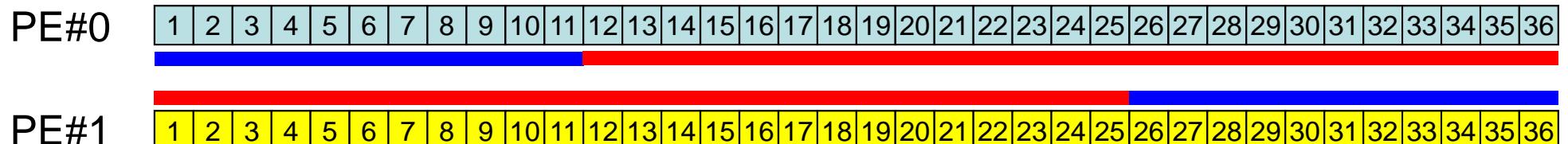
    StatSR = malloc(sizeof(MPI_Status));

    MPI_Sendrecv(&VAL      , 1, MPI_DOUBLE, neib, 0,
                &VALtemp, 1, MPI_DOUBLE, neib, 0, MPI_COMM_WORLD, StatSR);
    VAL=VALtemp;

    MPI_Finalize();
    return 0;
}
```

# Ex.2: Send-Recv an Array (1/4)

- 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.



t1

# Practice: t1

- Initial status of VEC(:):
  - PE#0 VEC(1-36)= 101,102,103,~,135,136
  - PE#1 VEC(1-36)= 201,202,203,~,235,236
- Confirm the results in the next page
- Using following two functions:
  - MPI\_Isend/Irecv/Waitall
  - MPI\_Sendrecv

# 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 (2/4)

t1

```
if (my_rank.eq.0) then
    call MPI_Isend (VEC( 1),11,MPI_DOUBLE_PRECISION,1,...,req_send,...)
    call MPI_Irecv (VEC(12),25,MPI_DOUBLE_PRECISION,1,...,req_recv,...)
endif

if (my_rank.eq.1) then
    call MPI_Isend (VEC( 1),25,MPI_DOUBLE_PRECISION,0,...,req_send,...)
    call MPI_Irecv (VEC(26),11,MPI_DOUBLE_PRECISION,0,...,req_recv,...)
endif

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

It works, but complicated operations.  
Not looks like SPMD.  
Not portable.

# Ex.2: Send-Recv an Array (3/4)

t1

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

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

call MPI_Isend
(&
(VEC(start_send),length_send,MPI_DOUBLE_PRECISION,NEIB,...,req_send,...)
call MPI_Irecv
(&
(VEC(start_recv),length_recv,MPI_DOUBLE_PRECISION,NEIB,...,req_recv,...)

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

This is “SMPD” !!

# Ex.2: Send-Recv an Array (4/4)

t1

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

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

call MPI_Sendrecv
(VEC(start_send),length_send,MPI_DOUBLE_PRECISION,NEIB,... &
 VEC(start_recv),length_recv,MPI_DOUBLE_PRECISION,NEIB,..., status,...)
```

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

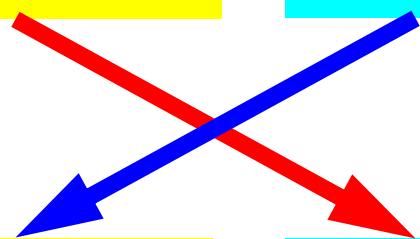
recv:

```
VEC(start_recv)~  
VEC(start_recv+length_recv-1)
```

#PE1

recv:

```
VEC(start_recv)~  
VEC(start_recv+length_recv-1)
```



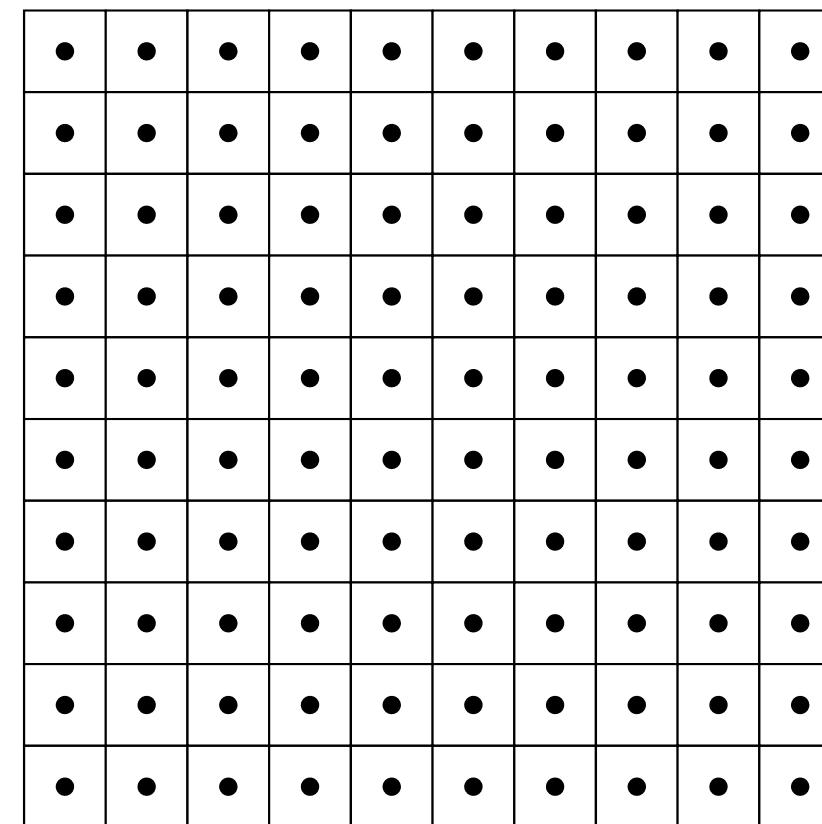
- “length\_send” of sending process must be equal to “length\_recv” of receiving process.
  - PE#0 to PE#1, PE#1 to PE#0
- “sendbuf” and “recvbuf”: different address

# Peer-to-Peer Communication

- What is P2P Communication ?
- 2D Problem, Generalized Communication Table
  - 2D FDM
  - Problem Setting
  - Distributed Local Data and Communication Table
  - Implementation
- Report S2

# 2D FDM (1/5)

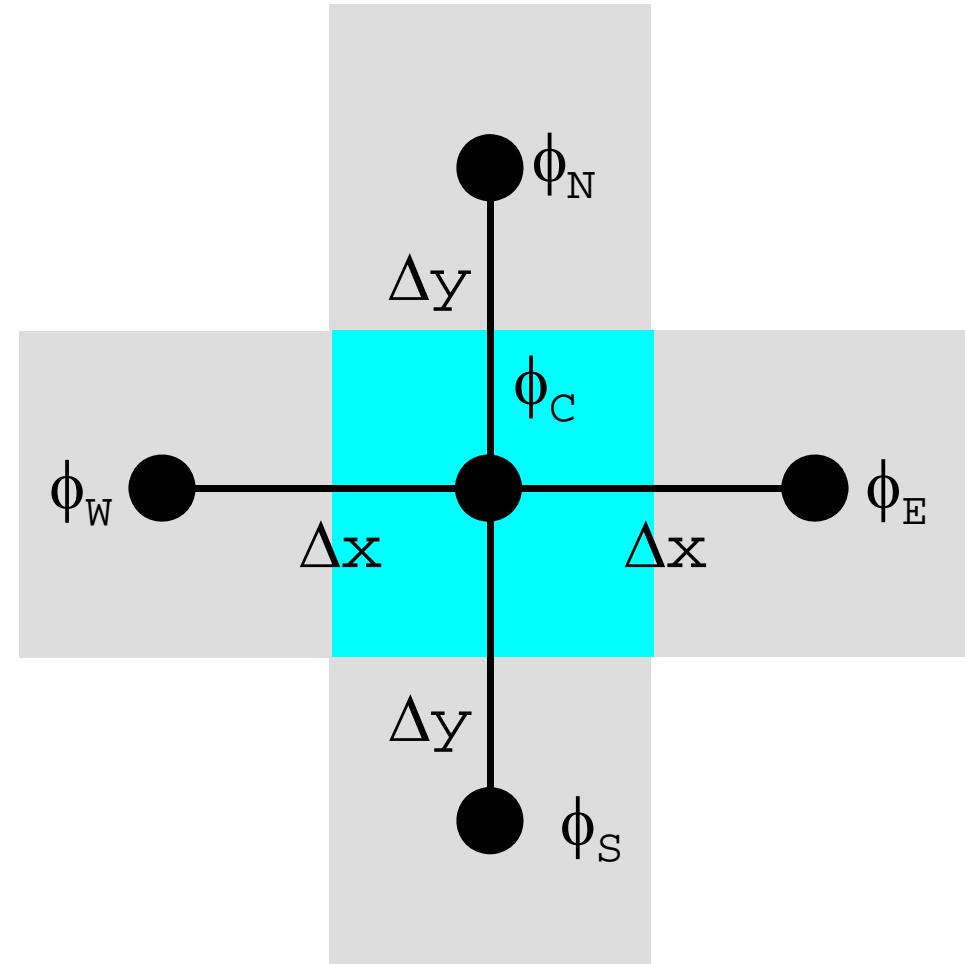
## Entire Mesh



# 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#3

<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#2

<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#3

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

PE#2

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

PE#0

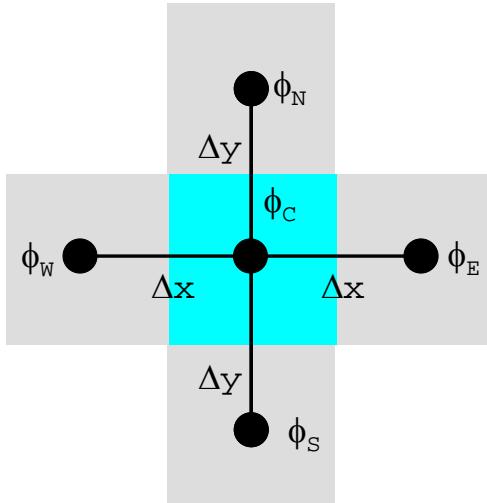
13	14	15	16
9	10	11	12
5	6	7	8
1	2	3	4

PE#1

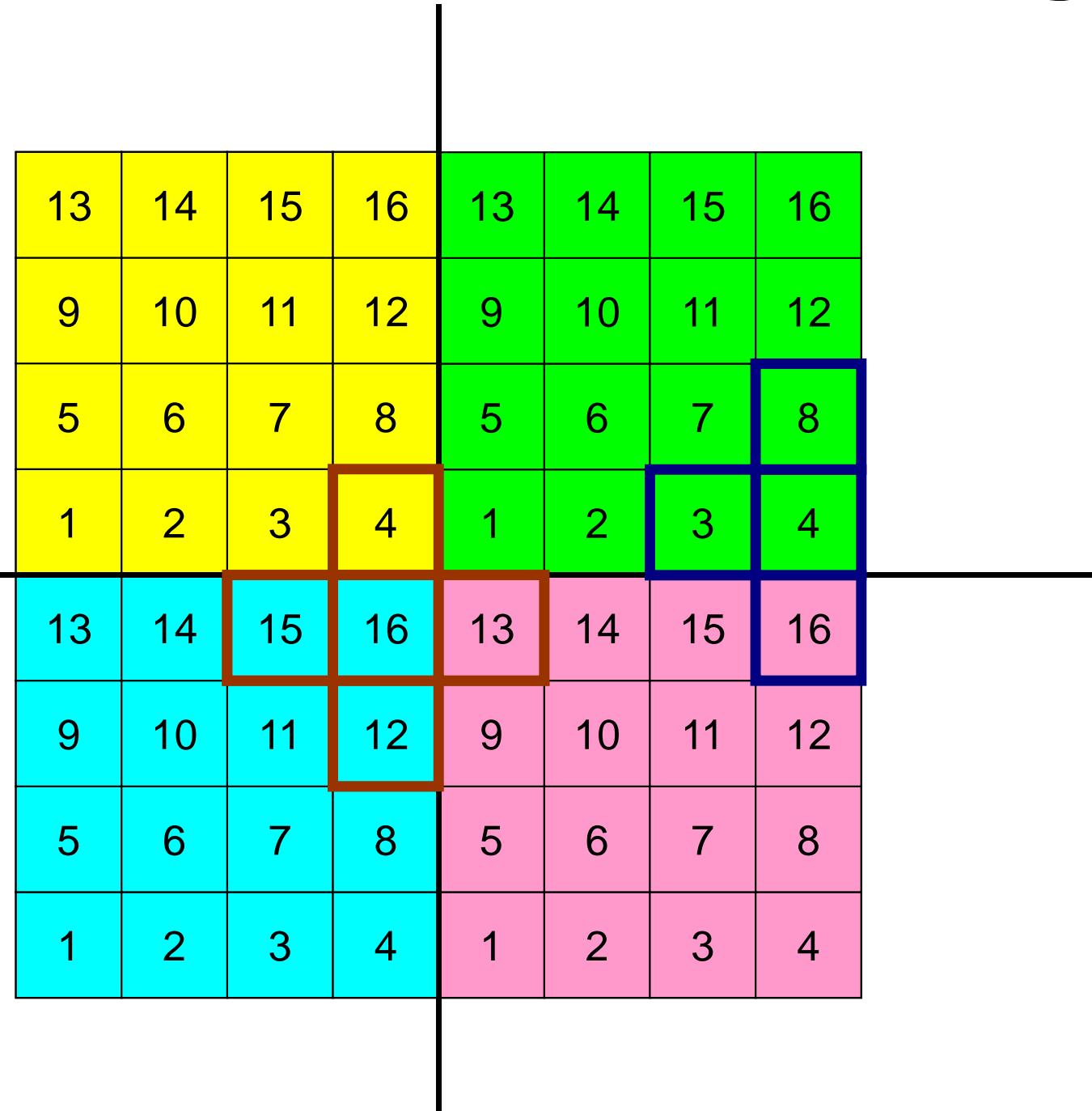
13	14	15	16
9	10	11	12
5	6	7	8
1	2	3	4

# External Points: Overlapped Region

PE#3



PE#2



PE#0

PE#1

# External Points: Overlapped Region

PE#3

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

PE#2

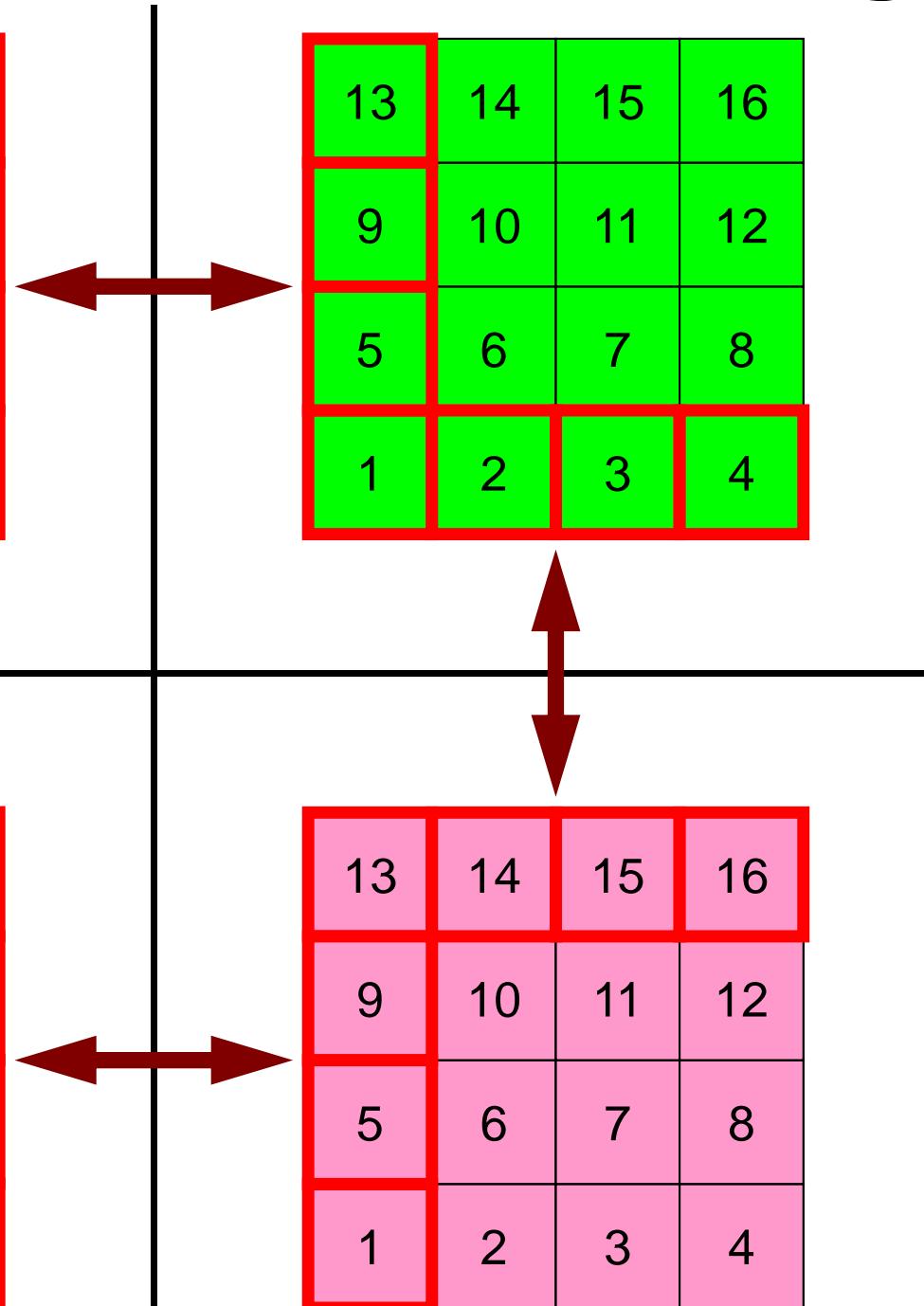
13	14	15	16
9	10	11	12
5	6	7	8
1	2	3	4

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#3

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

PE#2

?	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#3

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

PE#2

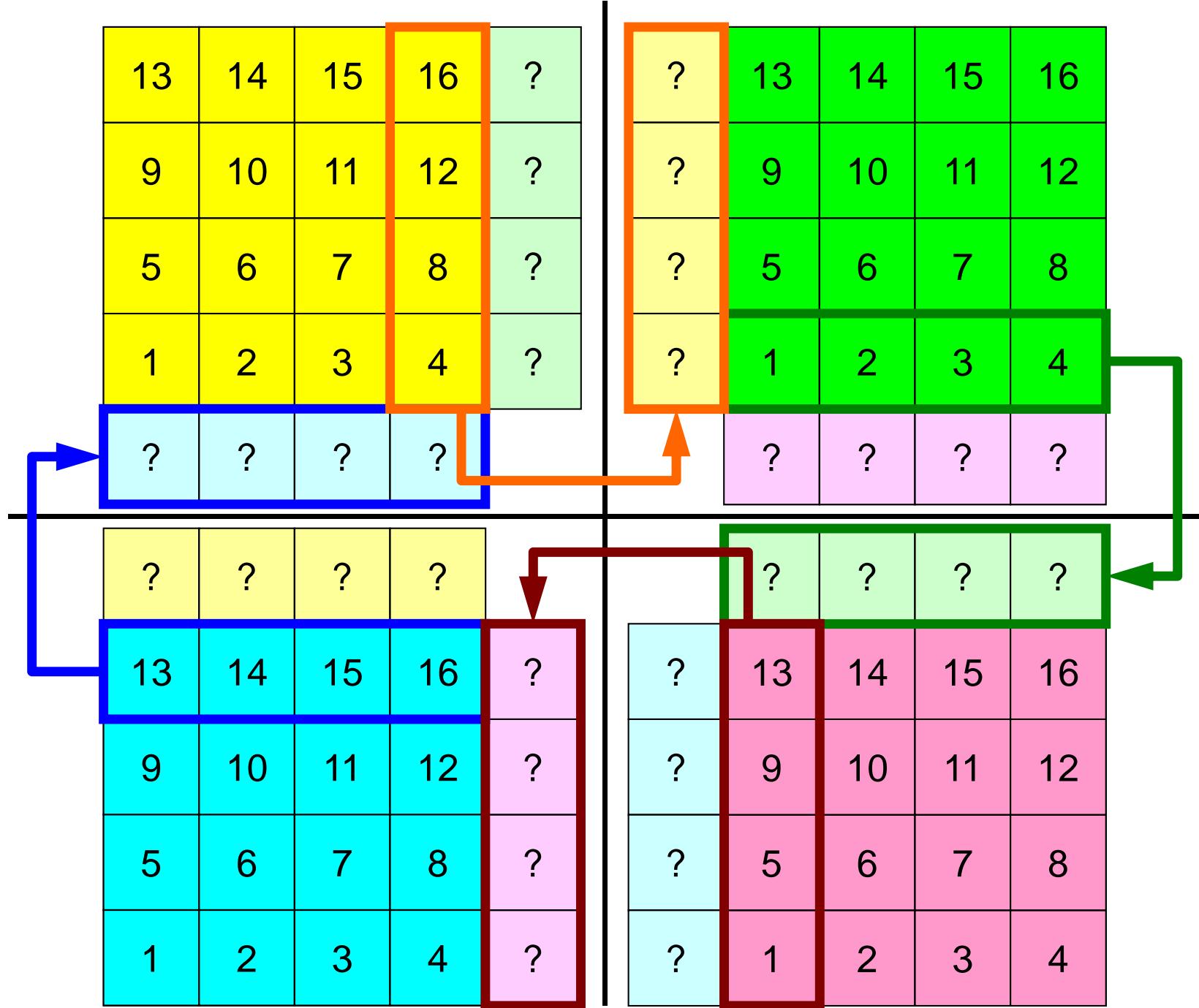
?	13	14	15	16
?	9	10	11	12
?	5	6	7	8
?	1	2	3	4

PE#0

?	?	?	?	?
13	14	15	16	?
9	10	11	12	?
5	6	7	8	?

PE#1

?	?	?	?	?
?	13	14	15	16
?	9	10	11	12
?	5	6	7	8



# Overlapped Region

PE#3

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

PE#2

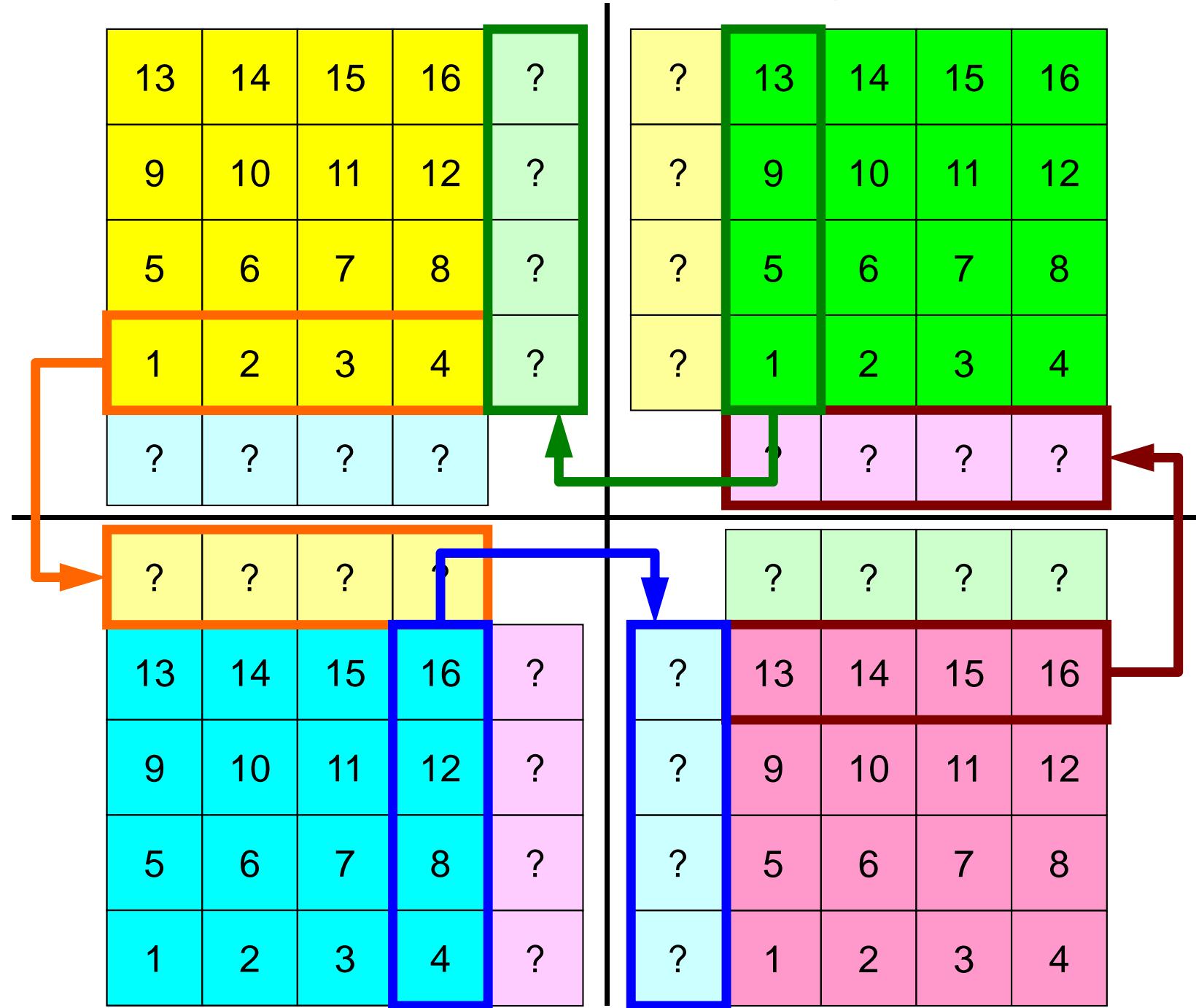
?	13	14	15	16
?	9	10	11	12
?	5	6	7	8
?	1	2	3	4
?	?	?	?	?

PE#0

?	?	?	1	?
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
?	?	?	?	?



# Peer-to-Peer Communication

- What is P2P Communication ?
- 2D Problem, Generalized Communication Table
  - 2D FDM
  - Problem Setting
  - Distributed Local Data and Communication Table
  - Implementation
- Report S2

# Problem Setting: 2D FDM

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

- 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

**PE#0**

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#1****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	

**PE#0**

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

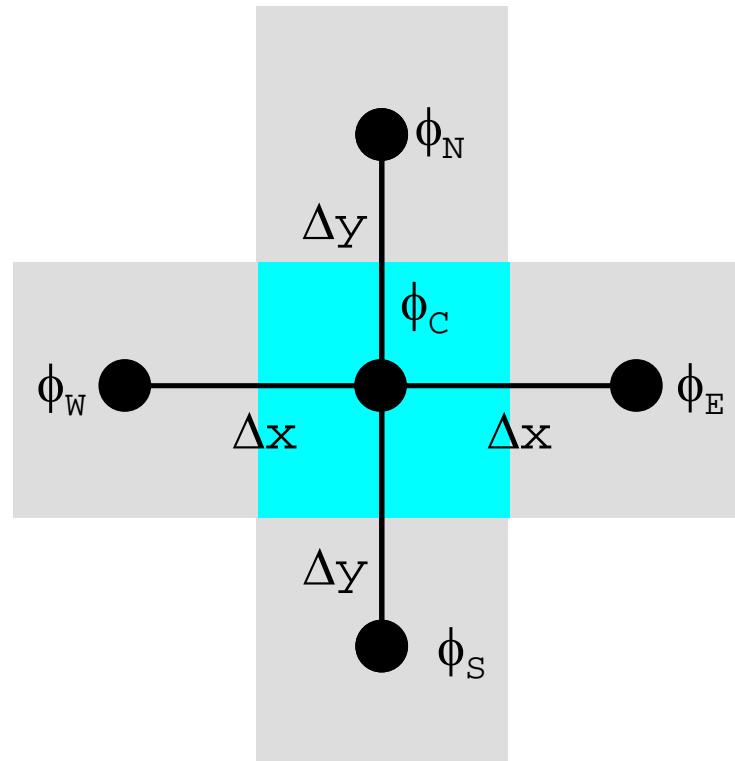
**PE#1**

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

# Operations of 2D FDM

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

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

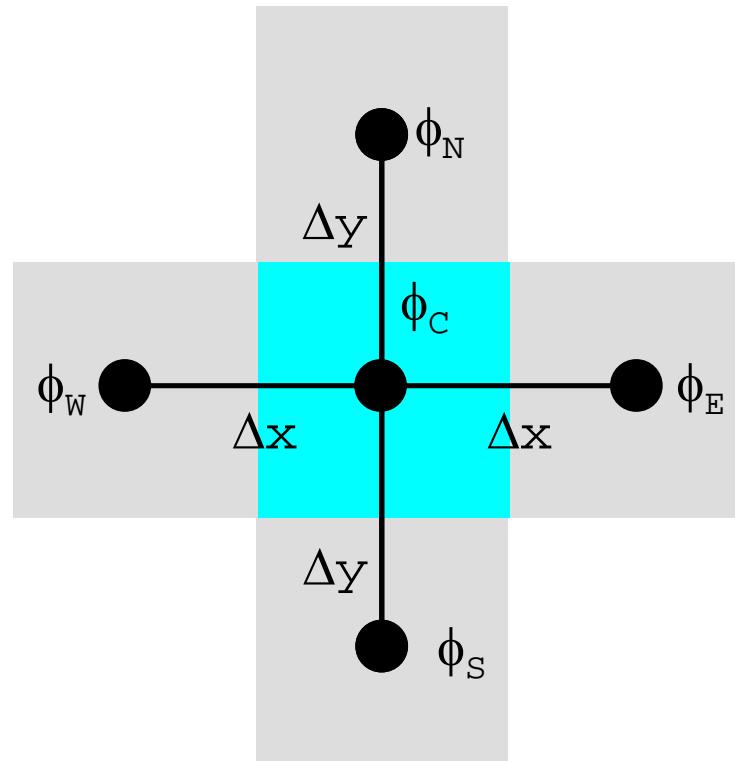


57	58	59	60	61	62	63	64
49	50	51	52	53	54	55	56
41	42	43	44	45	46	47	48
33	34	35	36	37	38	39	40
25	26	27	28	29	30	31	32
17	18	19	20	21	22	23	24
9	10	11	12	13	14	15	16
1	2	3	4	5	6	7	8

# 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>	57	58	59	60	61	62	63	64	<u>PE#3</u>
<u>PE#0</u>	49	50	51	52	53	54	55	56	<u>PE#1</u>
	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	

- 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

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

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

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

<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**

	<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#1**

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

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

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

**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

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

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

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

**PE#0**

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

**PE#1**

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

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

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>

**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>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>

**PE#0**

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

**PE#1**

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>

# Peer-to-Peer Communication

- What is P2P 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	

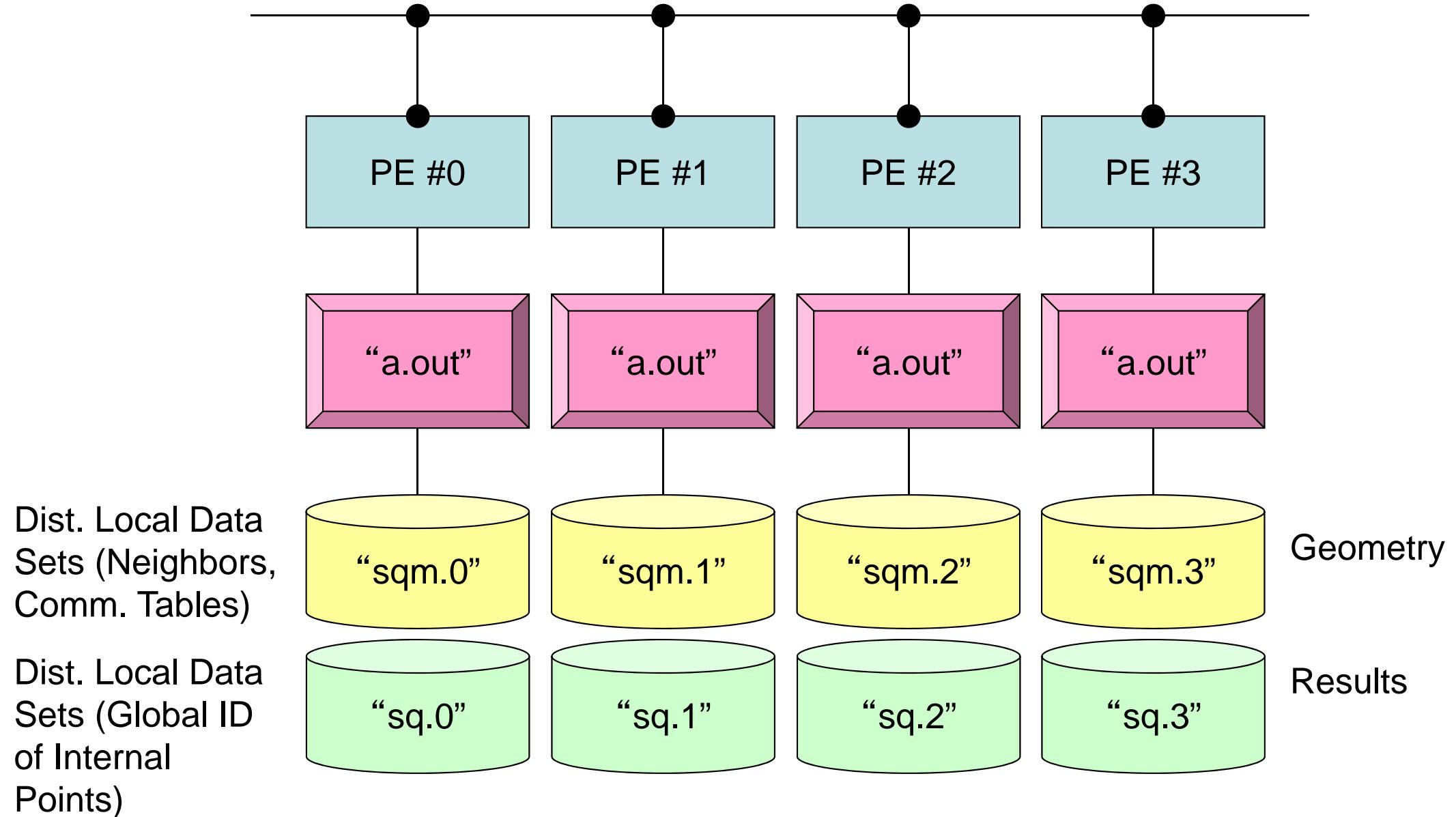
PE#0

PE#1

Value at each mesh (= Global ID)

Local ID

# SPMD . . .



# 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#3**

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)

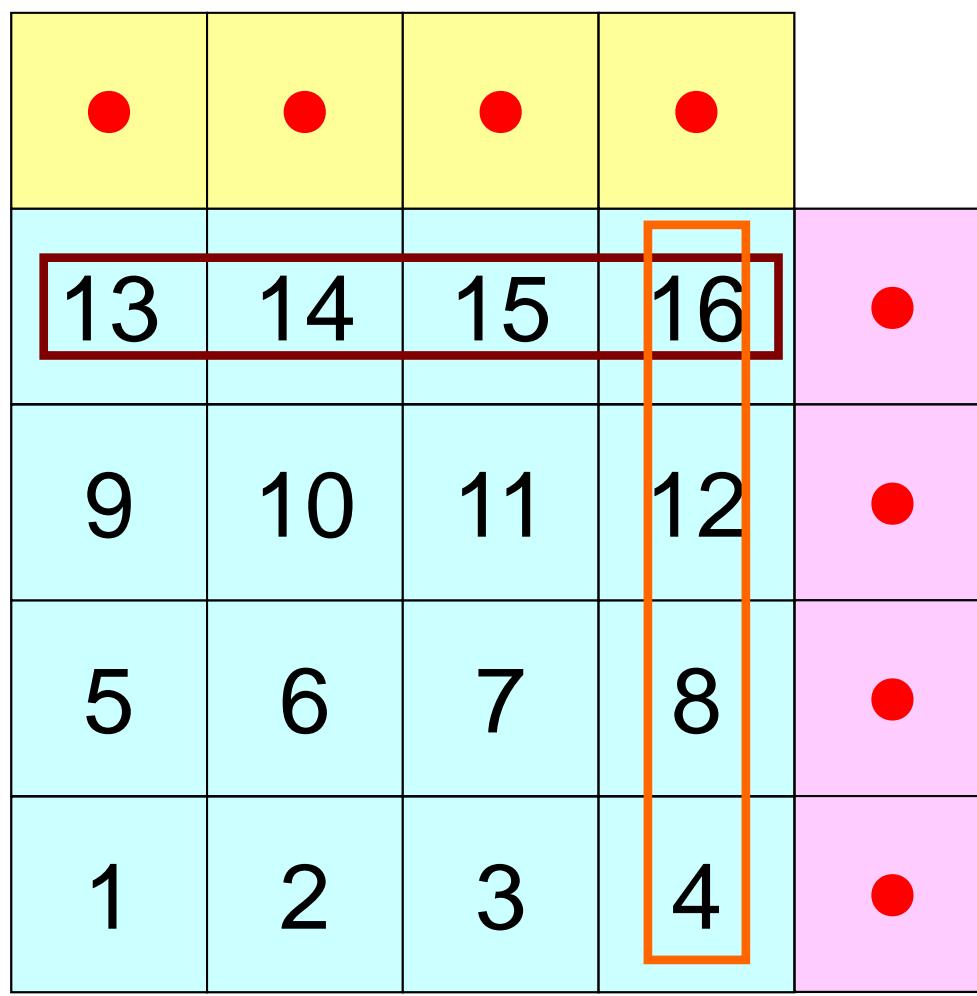
- Sleeves
- Halo



# 2D FDM: PE#0

## Information at each domain (3/4)

PE#3



### 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)

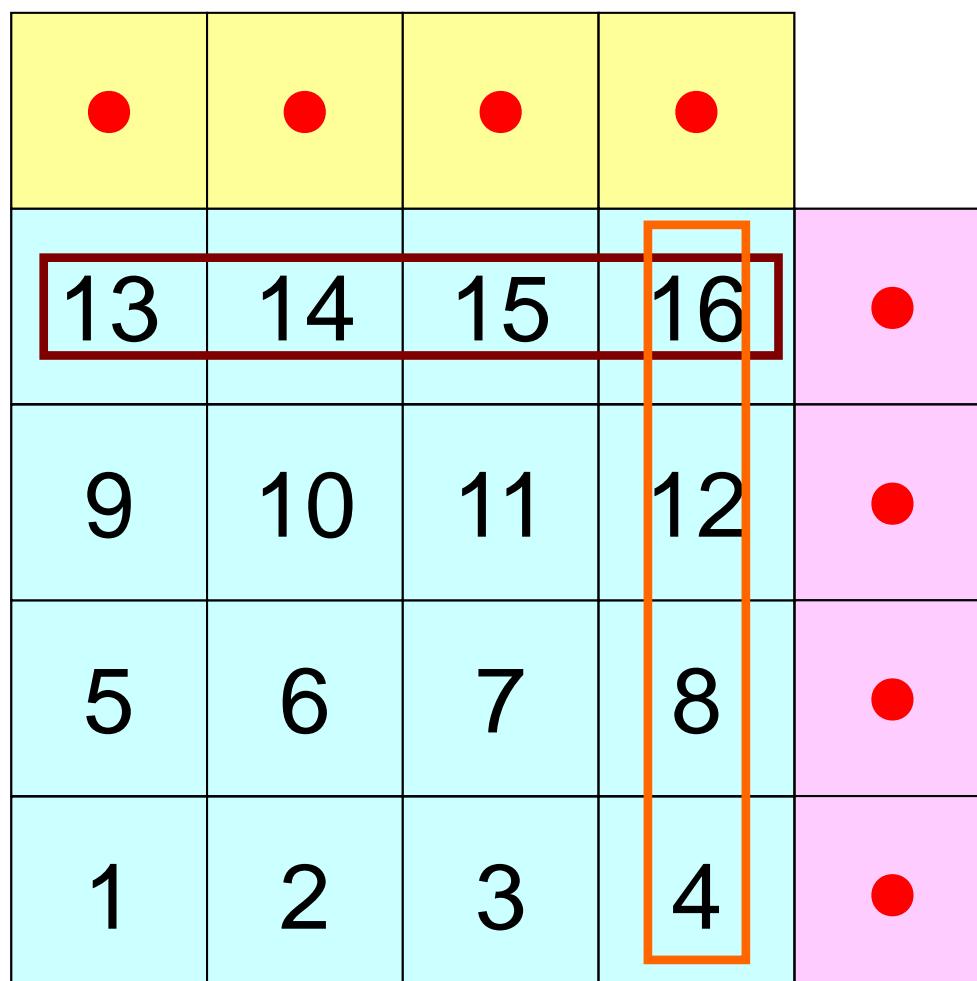
### Boundary Points

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#3



**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
- External Points
  - From where, how many, and which external points are received/imported ?
- Boundary Points
  - 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

- Neighbors
  - NeibPETot, NeibPE[neib]
- Message size for each neighbor
  - export\_index[neib], neib= 0, NeibPETot-1
- ID of boundary points
  - export\_item[k], k= 0, export\_index[NeibPETot]-1
- Messages to each neighbor
  - SendBuf[k], k= 0, 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[neib]
- Message size for each neighbor
  - import\_index[neib], neib= 0, NeibPETot-1
- ID of external points
  - import\_item[k], k= 0, import\_index[NeibPETot]-1
- Messages from each neighbor
  - RecvBuf[k], k= 0, import\_index[NeibPETot]-1

# RECV: MPI\_Isend/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);

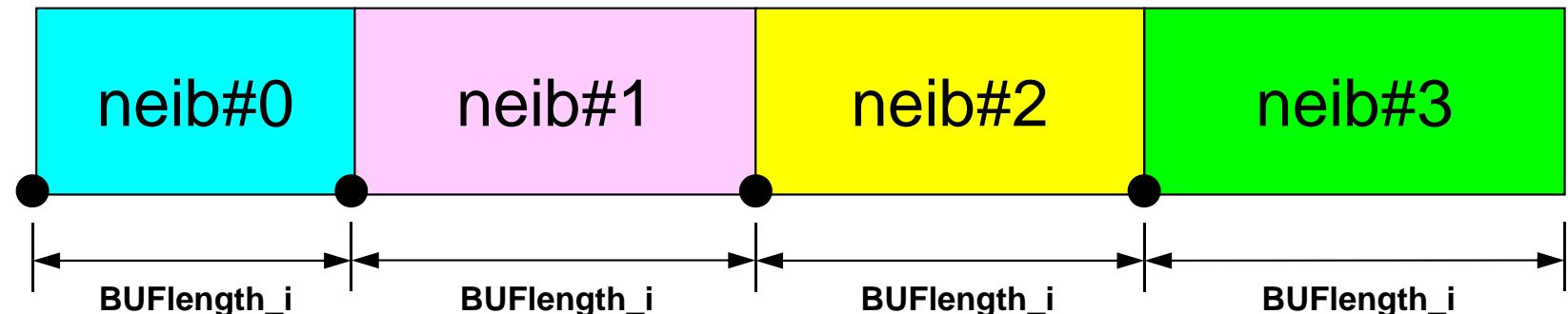
for (neib=0; neib<NeibPETot; neib++){
    for (k=import_index[neib];k<import_index[neib+1];k++){
        kk= import_item[k];
        VAL[kk]= RecvBuf[k];
    }
}

```

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

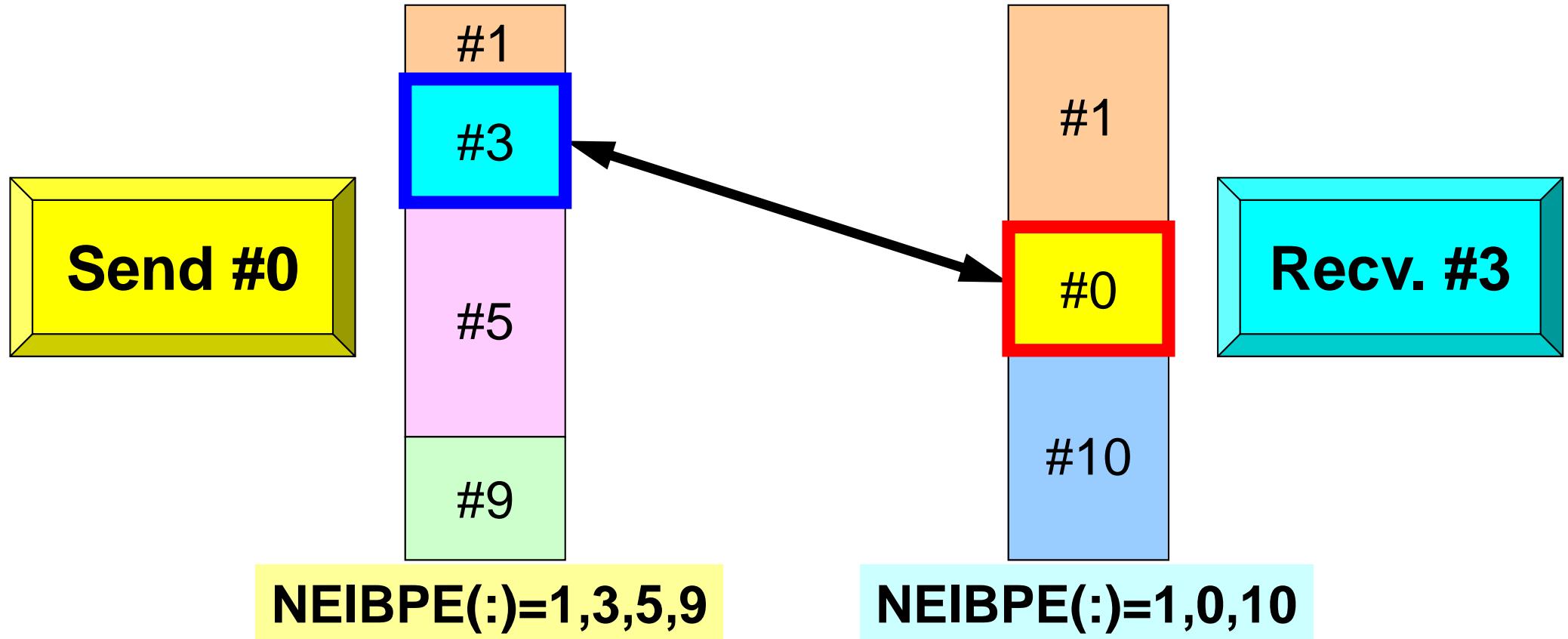
    call MPI_ISEND
&          (SENDbuf(iS_e), BUFlength_e, MPI_INTEGER, NEIBPE(neib), 0,&
&           MPI_COMM_WORLD, request_send(neib), ierr)
enddo
```

```
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#3

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 3  
#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#3

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 3
#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#3

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 3
```

```
#NODE
```

```
24 16
```

```
#IMPORT_index
```

```
4 8
```

```
#IMPORT_items
```

```
17
```

```
18
```

```
19
```

```
20
```

```
21
```

```
22
```

```
23
```

```
24
```

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#3).

```
#EXPORT_index
```

```
4 8
```

```
#EXPORT_items
```

```
4
```

```
8
```

```
12
```

```
16
```

```
13
```

```
14
```

```
15
```

```
16
```

# Generalized Comm. Table (4/6)

PE#3

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 3
#NODE
24 16
#IMPORT_index
4 8
#IMPORT_items
17
18 imported from 1st Neighbor
19 (PE#1) (1st-4th items)
20
21
22 imported from 2nd Neighbor
23 (PE#3) (5th-8th items)
24
#EXPORT_index
4 8
#EXPORT_items
4
8
12
16
13
14
15
16

```

# Generalized Comm. Table (5/6)

PE#3

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 3
#NODE
24 16
#IMPORT_index
4 8
#IMPORT_items
```

17  
18  
19  
20  
21  
22  
23  
24

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#3).

```
#EXPORT_index
4 8
#EXPORT_items
4
8
12
16
13
14
15
16
```

# Generalized Comm. Table (6/6)

PE#3

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 3
#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          exported to 1st Neighbor
12          (PE#1) (1st-4th items)
16
13
14          exported to 2nd Neighbor
15          (PE#3) (5th-8th items)
16

```

# Generalized Comm. Table (6/6)

PE#3

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

# Peer-to-Peer Communication

- What is P2P 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 <$O-S2>

$ mpifrtpx -Kfast sq-sr1.f
$ mpifccpx -Kfast sq-sr1.c

(modify go4.sh for 4 processes)
$ pbsub 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==*/
```

**np** Number of all meshes (internal + external)  
**n** Number of internal meshes

```

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

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

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] = 1;

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 = malloc(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));

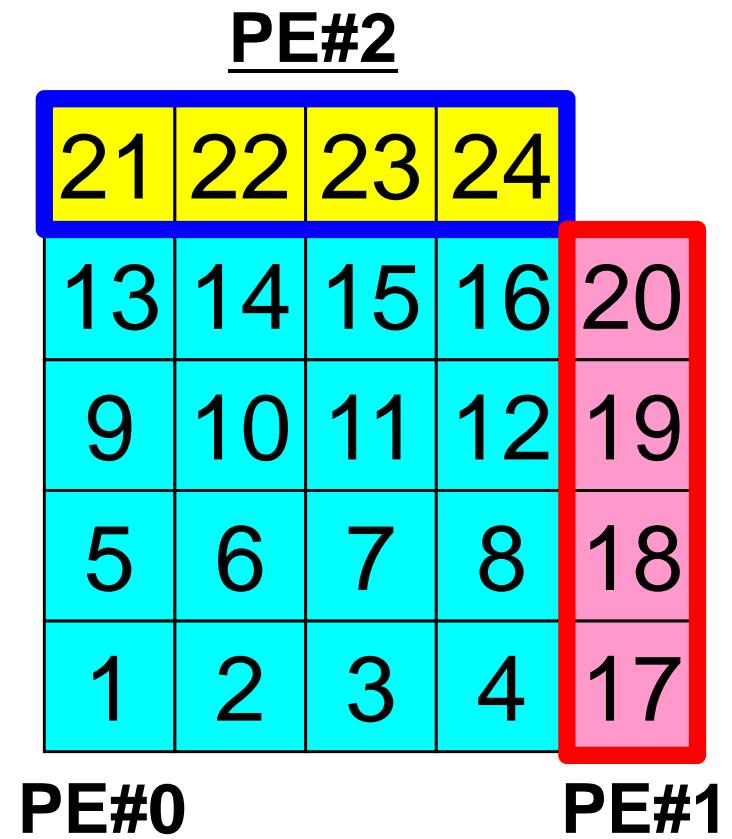
        for(i=0;i<nn;i++){
            fscanf(fp, "%d", &ExportItem[i]); ExportItem[i]--;
}

```

#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] = 0;

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 = malloc(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));

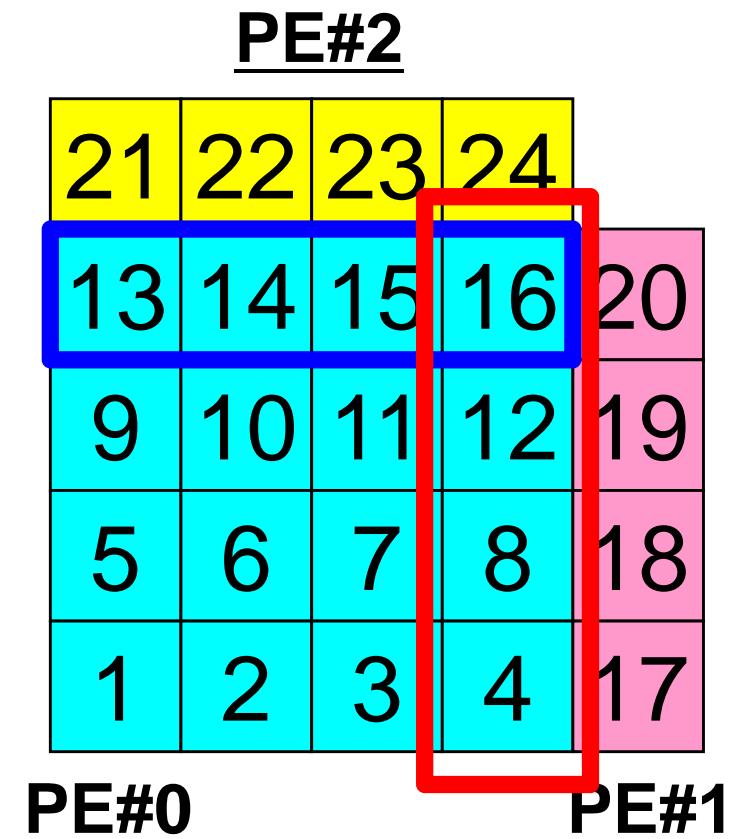
        for(i=0;i<nn;i++){
            fscanf(fp, "%d", &ExportItem[i]); ExportItem[i]--;

```

#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.

PE#2

<b>25</b>	<b>26</b>	<b>27</b>	<b>28</b>
<b>17</b>	<b>18</b>	<b>19</b>	<b>20</b>
<b>9</b>	<b>10</b>	<b>11</b>	<b>12</b>

PE#0

PE#1

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]];
        }
    }
}
```

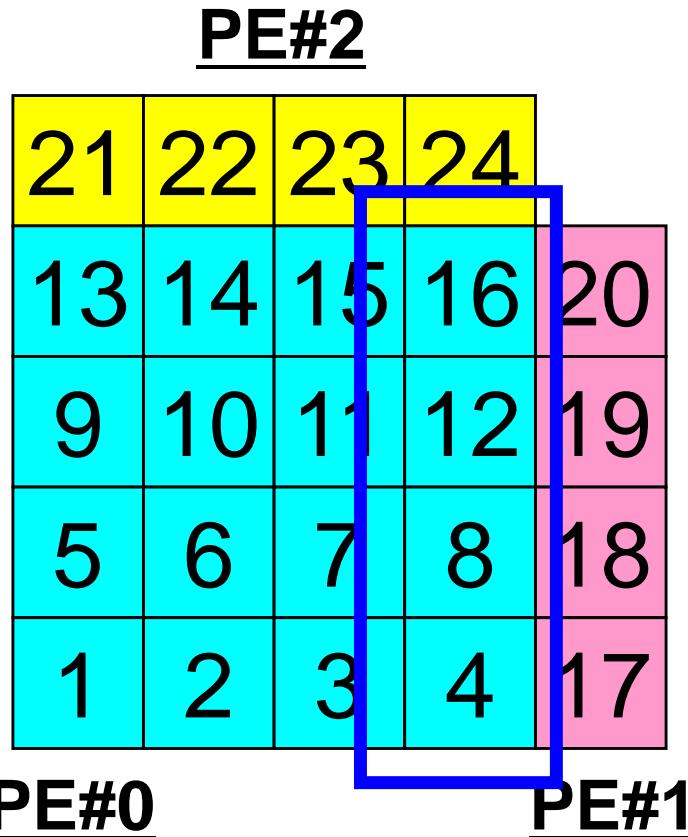
Info. of boundary points is written into sending buffer (`SendBuf`).

Info. sent to `NeibPe[neib]` is stored in `SendBuf[ExportIndex[neib]:ExportIndex[neib+1]-1]`:

# Sending Buffer is nice ...

```
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

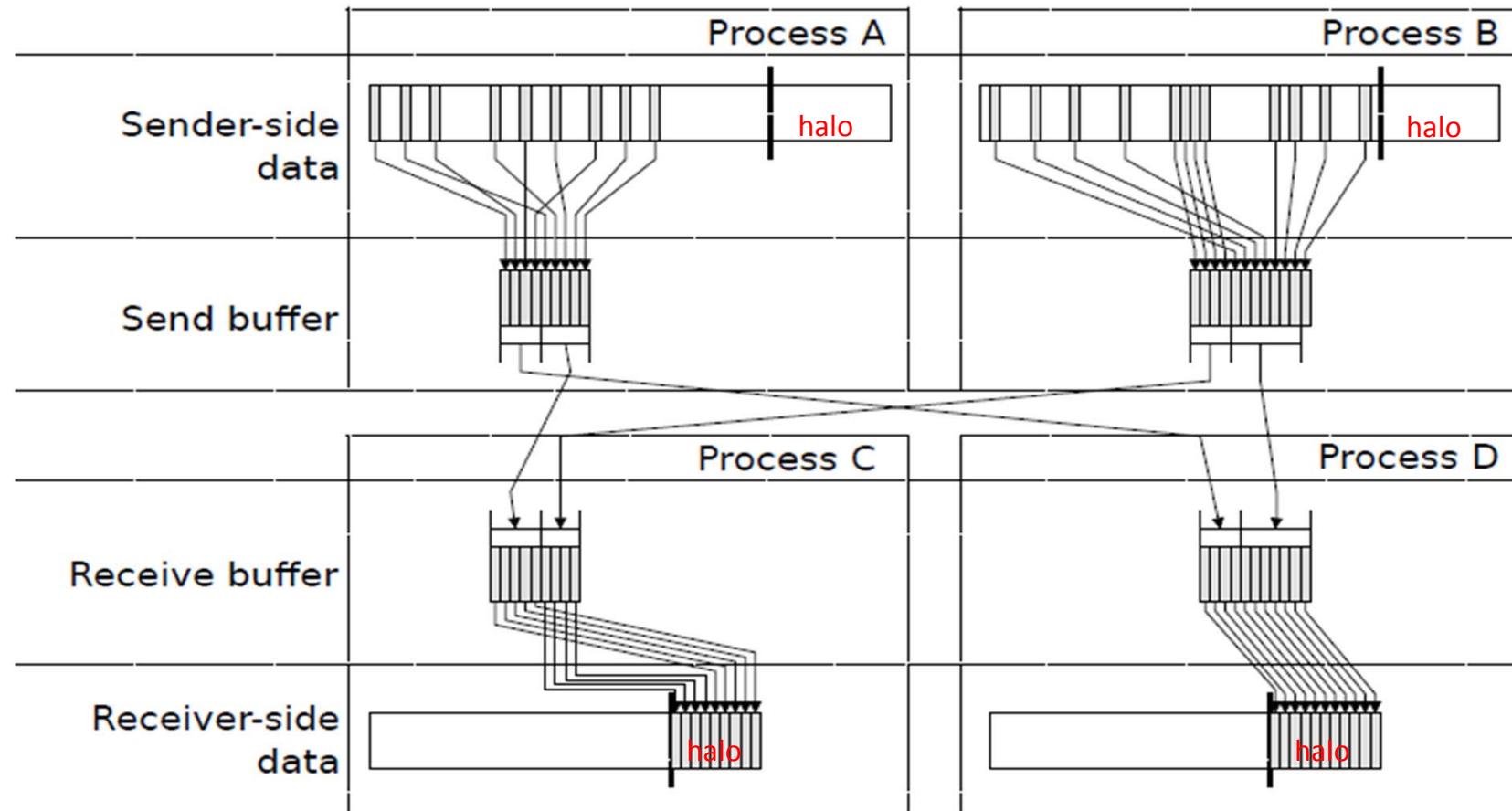
    ierr= MPI_Isend
        (&SendBuf[iS_e], BUFlength_e, MPI_DOUBLE, NeibPE[neib], 0,
         MPI_COMM_WORLD, &ReqSend[neib])
}
```



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



Dr. Osni Marques  
(Lawrence Berkeley National Laboratory)

# 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]);
}
```

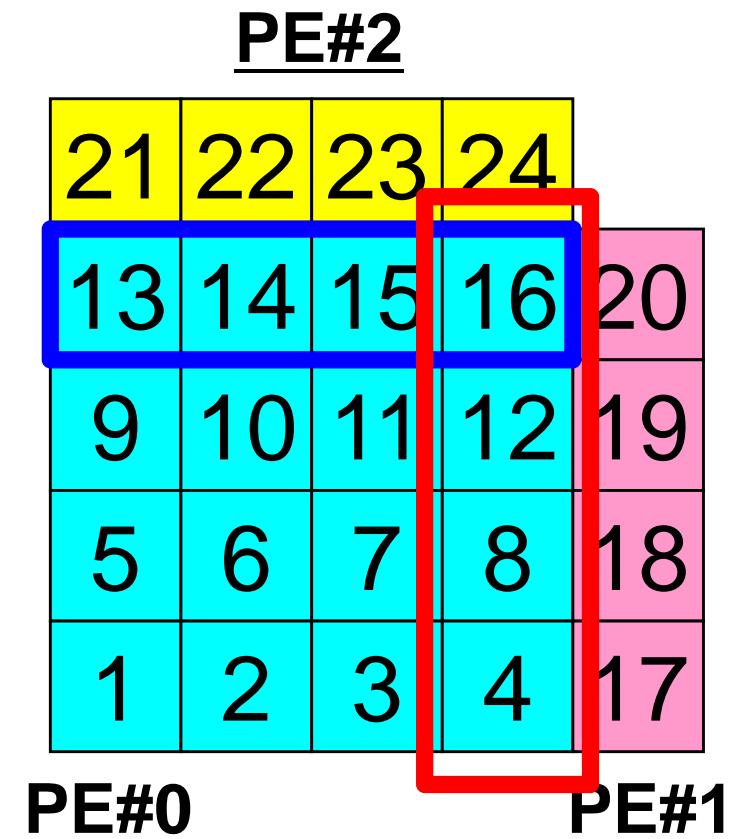
```
for(neib=0;neib<NeibPeTot;neib++){
    iStart = ImportIndex[neib];
    iEnd   = ImportIndex[neib+1];
    BufLength = iEnd - iStart;

    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
PE#0	PE#1
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

# 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];
    }
}

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);

```

Copied to sending buffers

# 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

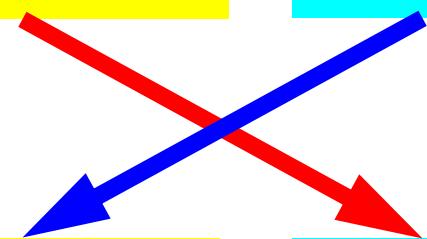
recv:

```
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

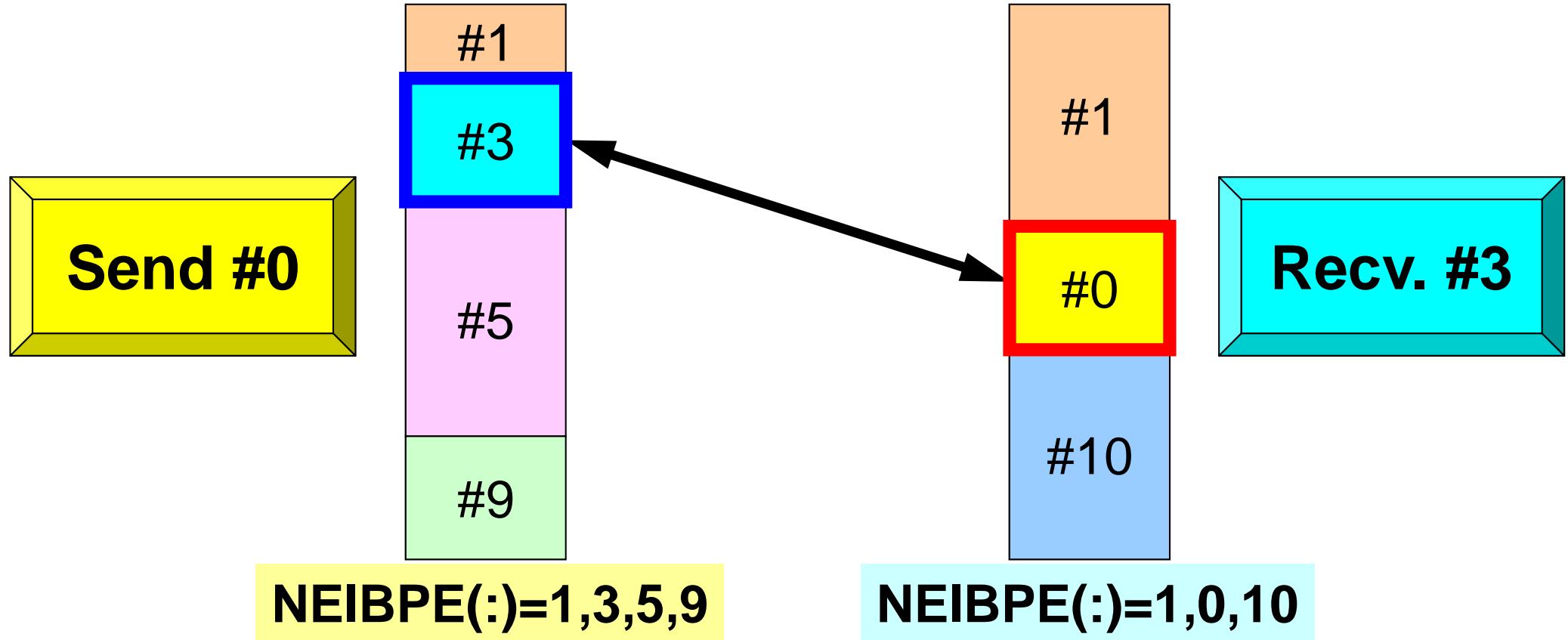
    call MPI_ISEND
&          (SENDbuf(iS_e), BUFlength_e, MPI_INTEGER, NEIBPE(neib), 0,&
&           MPI_COMM_WORLD, request_send(neib), ierr)
enddo
```

```
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]);
}
```

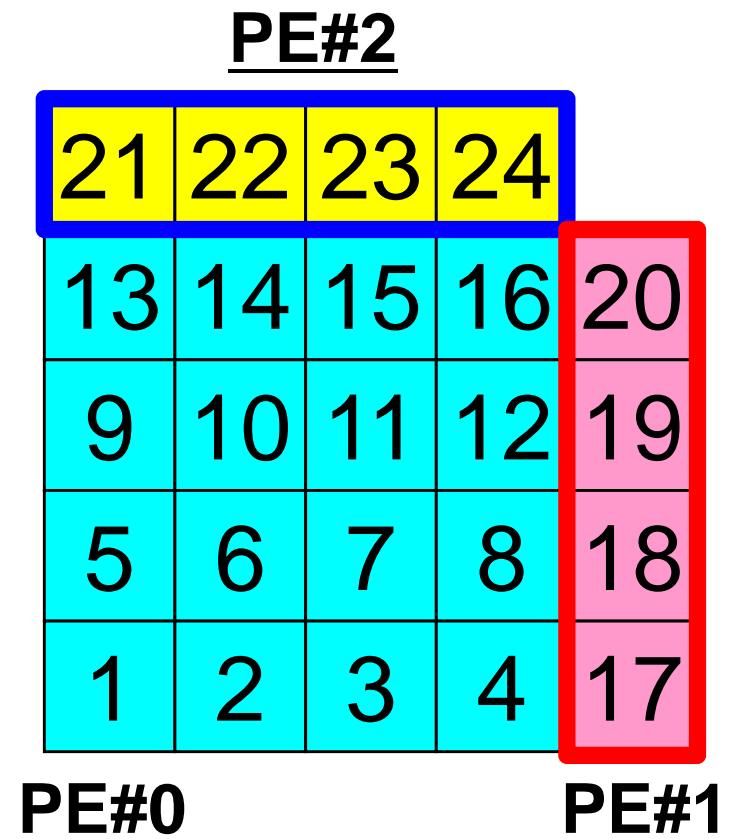
```
for(neib=0;neib<NeibPeTot;neib++){
    iStart = ImportIndex[neib];
    iEnd   = ImportIndex[neib+1];
    BufLength = iEnd - iStart;

    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
PE#0	PE#1
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

# 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\_Isend/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])
}

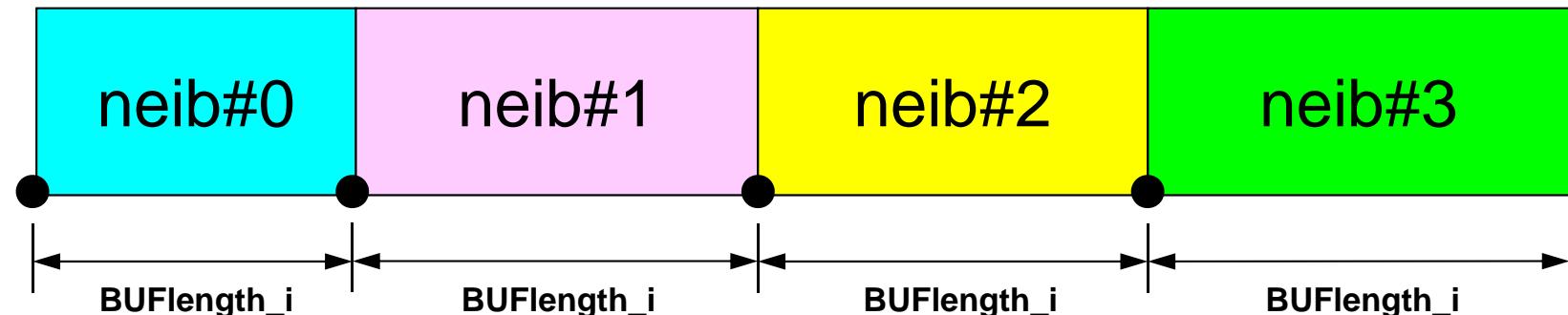
MPI_Waitall(NeibPETot, ReqRecv, StatRecv);

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];
    }
}
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;
}
```

Contents of RecvBuf are copied to values at external points.

# 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**

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

**PE#0**

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

**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**

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

**PE#0**

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

**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#2)

**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

29	30	31	32
21	22	23	24
13	14	15	16

**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#3)

**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

**PE#0**

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

**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

# 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>

#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>	

# Three Domains

#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>

#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

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

#PE0

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

#PE1

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

#NEIBPETot

2

#NEIBPE

1 2

#NODE

13 8 (int+ext, int pts)

#IMPORTindex

O O

#IMPORTitems

O...

#EXPORTindex

O O

#EXPORTitems

O...

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

#PE2

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

#PE0

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

#PE1

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

#NEIBPETot

2

#NEIBPE

0 2

#NODE

8 14 (int+ext, int pts)

#IMPORTindex

O O

#IMPORTitems

O...

#EXPORTindex

O O

#EXPORTitems

O...

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

#PE2

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

#PE0

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

#PE1

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

#NEIBPETot

2

#NEIBPE

1 0

#NODE

9 15 (int+ext, int pts)

#IMPORTindex

O O

#IMPORTitems

O...

#EXPORTindex

O O

#EXPORTitems

O...

#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

t2

- 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.c) to <\$O-S2>/ex
- pbsub 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
- Measure parallel performance

# Report S2 (2/2)

- Deadline: 17:00 February 15th (Sa), 2014.
  - 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)