# Introduction to Parallel FEM in Fortran Parallel Data Structure

Kengo Nakajima
Information Technology Center
The University of Tokyo

#### **Parallel Computing**

Faster, Larger & More Complicated

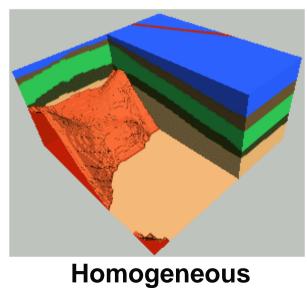
- Scalability
  - Solving N<sup>x</sup> scale problem using N<sup>x</sup> computational resources during same computation time
    - for large-scale problems: Weak Scaling
    - e.g. CG solver: more iterations needed for larger problems
  - Solving a problem using N<sup>x</sup> computational resources during 1/N computation time
    - for faster computation: <u>Strong Scaling</u>

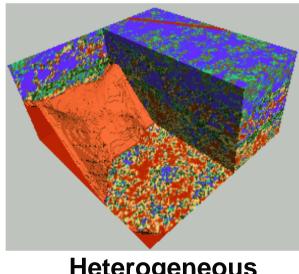
#### What is Parallel Computing? (1/2)

to solve larger problems faster

#### Homogeneous/Heterogeneous **Porous Media**

**Lawrence Livermore National Laboratory** 



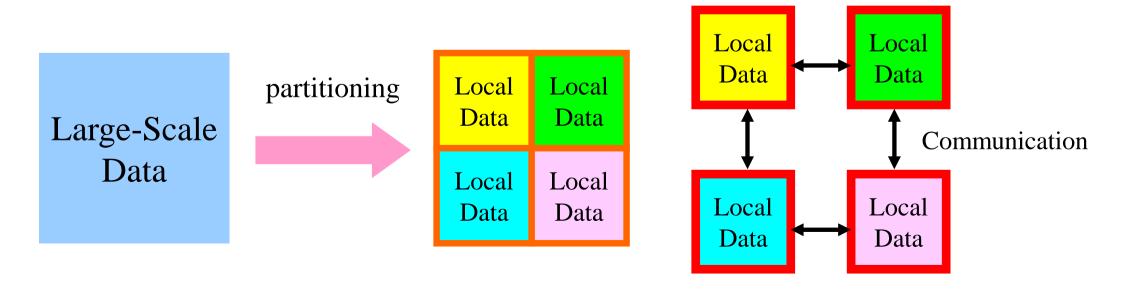


Heterogeneous

very fine meshes are required for simulations of heterogeneous field.

#### What is Parallel Computing? (2/2)

- PC with 1GB memory: 1M meshes are the limit for FEM
  - Southwest Japan with 1,000km x 1,000km x 100km in 1km mesh
     108 meshes
- Large Data -> Domain Decomposition -> Local Operation
- Inter-Domain Communication for Global Operation



#### What is Communication?

Parallel Computing -> Local Operations

 Communications are required in Global Operations for Consistency.

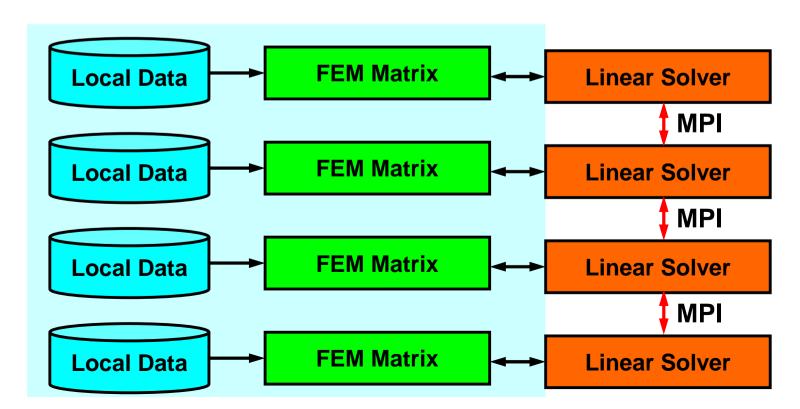
#### **Operations in Parallel FEM**

SPMD: Single-Program Multiple-Data

Large Scale Data -> partitioned into Distributed Local Data Sets.

FEM code can assemble coefficient matrix for each local data set: this part could be completely local, same as serial operations

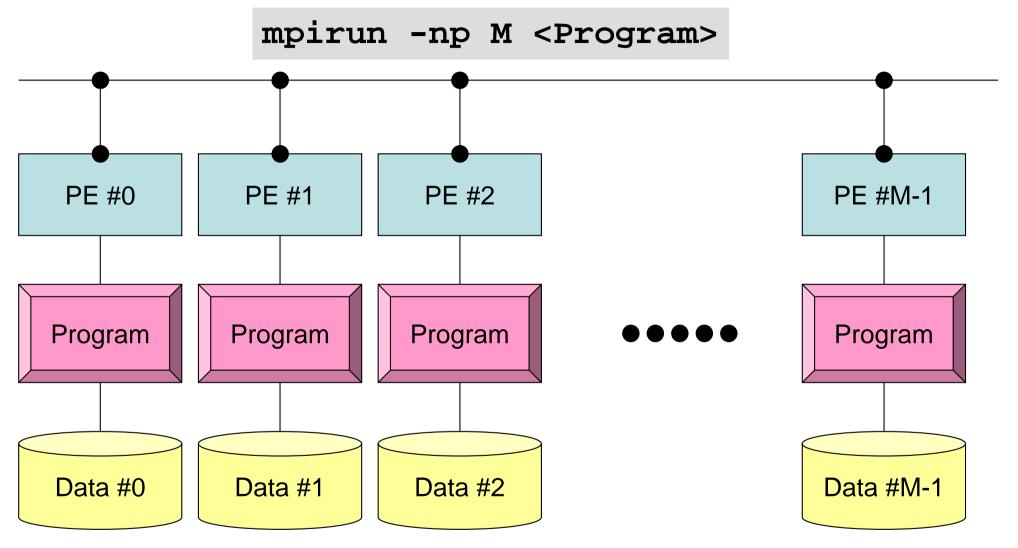
Global Operations & Communications happen only in Linear Solvers dot products, matrix-vector multiply, preconditioning



PE: Processing Element Processor, Domain, Process

#### **SPMD**

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



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.

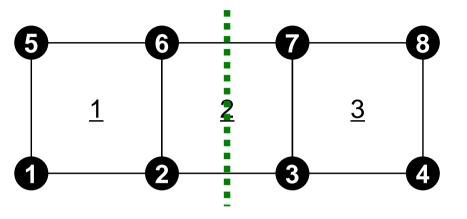
#### Parallel FEM Procedures

- Design on "Local Data Structure" is important
  - for SPMD-type operations in the previous page
- Matrix Generation
- Preconditioned Iterative Solvers for Linear Equations

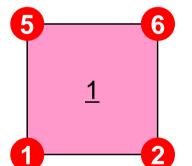
#### **Bi-Linear Square Elements**

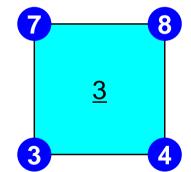
#### Values are defined on each node



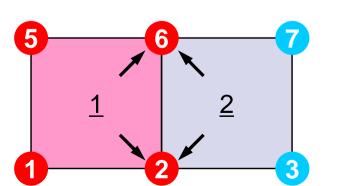


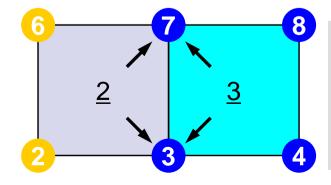
divide into two domains by "node-based" manner, where number of "nodes (vertices)" are balanced.





Local information is not enough for matrix assembling.





Information of overlapped elements and connected nodes are required for matrix assembling on boundary nodes.

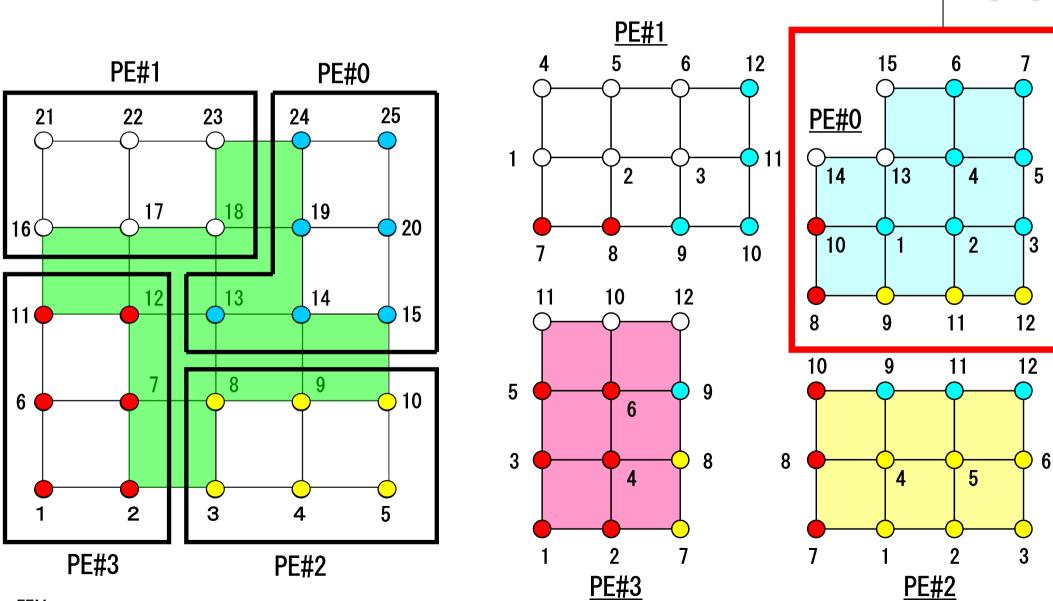
#### **Local Data of Parallel FEM**



- Node-based partitioning for preconditioned iterative solvers
- Local data includes information for :
  - Nodes originally assigned to the partition/PE
  - Elements which include the nodes (originally assigned to the Partition/PE)
  - All nodes which form the elements but out of the partition
- Nodes are classified into the following 3 categories from the viewpoint of the message passing
  - Internal nodes originally assigned nodes
  - External nodes in the overlapped elements but out of the partition
  - Boundary nodes external nodes of other partition (part of internal nodes)
- Communication table between partitions
- NO global information required except partition-to-partition connectivity

#### **Node-based Partitioning**

internal nodes - elements - external nodes

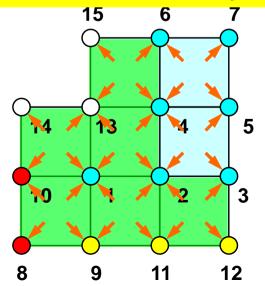


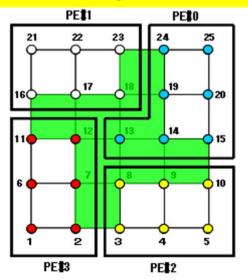
#### **Node-based Partitioning**

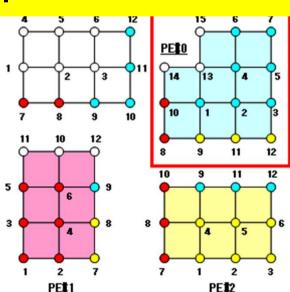
internal nodes - elements - external nodes



- ●Partitioned nodes themselves (Internal Nodes) 内点
- ●Elements which include Internal Nodes 内点を含む要素
- External Nodes included in the Elements 外点 in overlapped region among partitions.
- Info of External Nodes are required for completely local element—based operations on each processor.



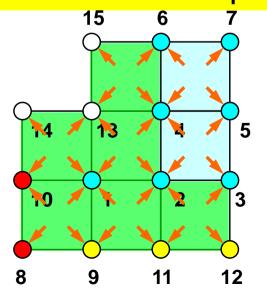


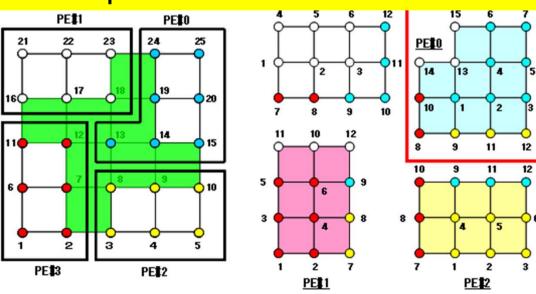


Intr<mark>o pFEM</mark>

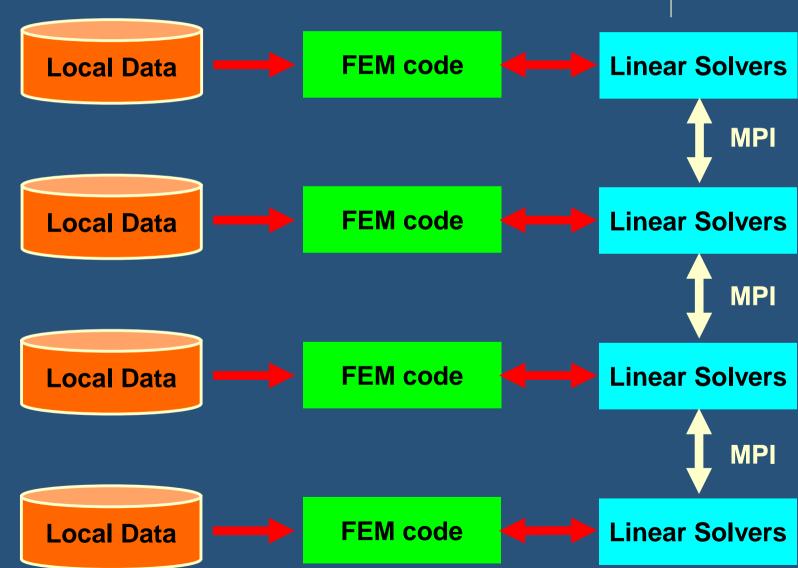
# We do not need communication during matrix assemble!!

- Partitioned nodes themselves (Internal Nodes)
- Elements which include Internal Nodes
- External Nodes included in the Elements in overlapped region among partitions.
- Info of External Nodes are required for completely local element—based operations on each processor.



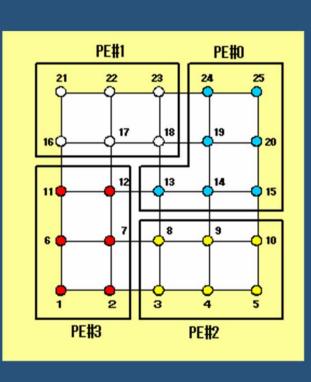


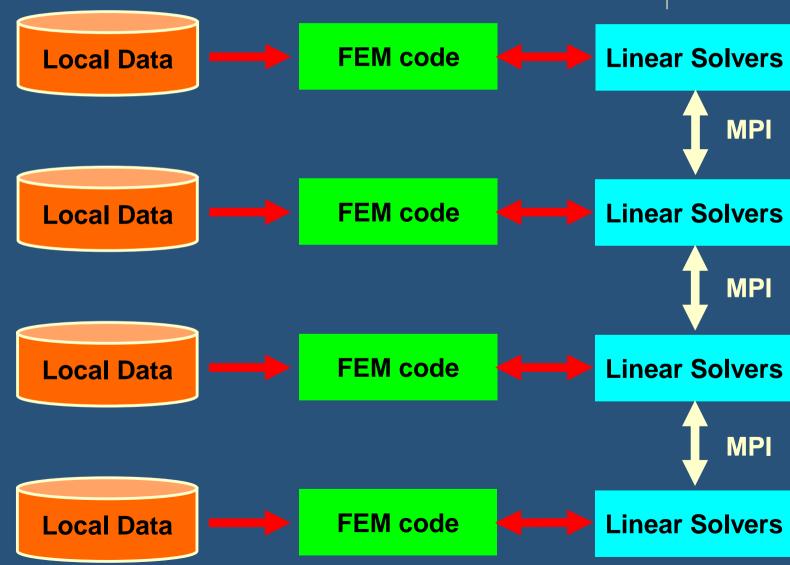




Intro-pFEM 14

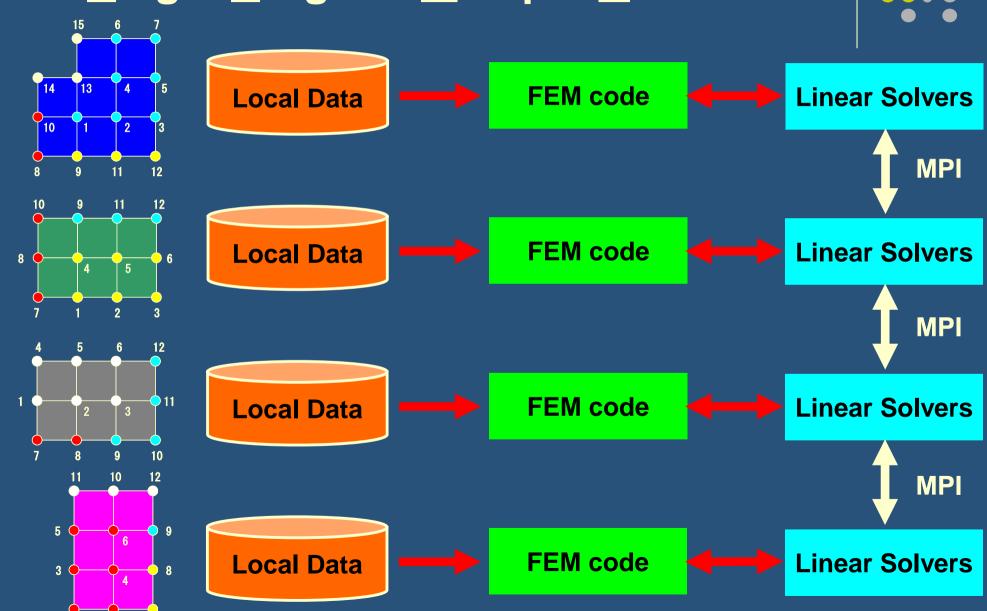




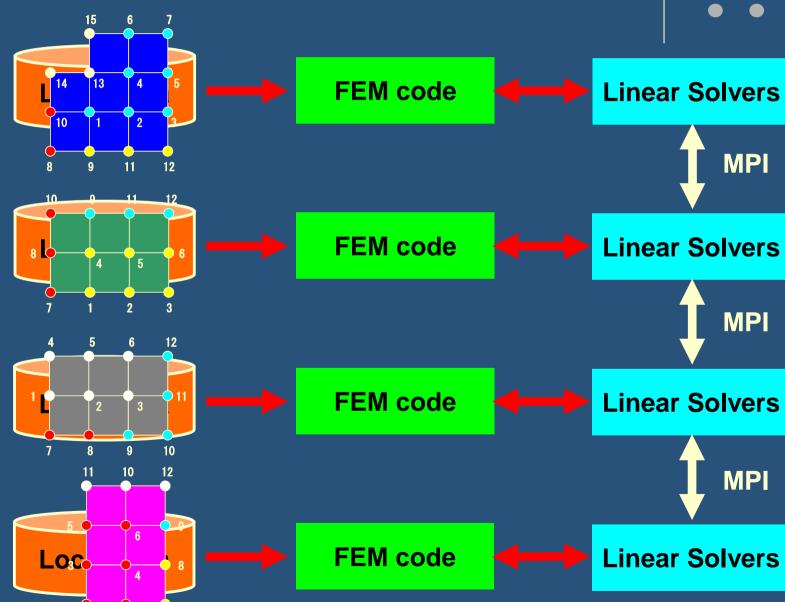


Intro-pFEM 15

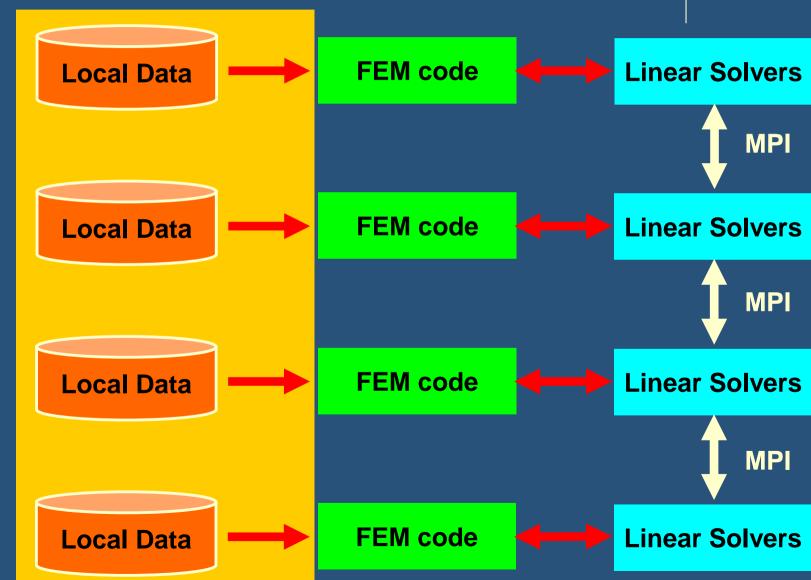












#### What is Communications?

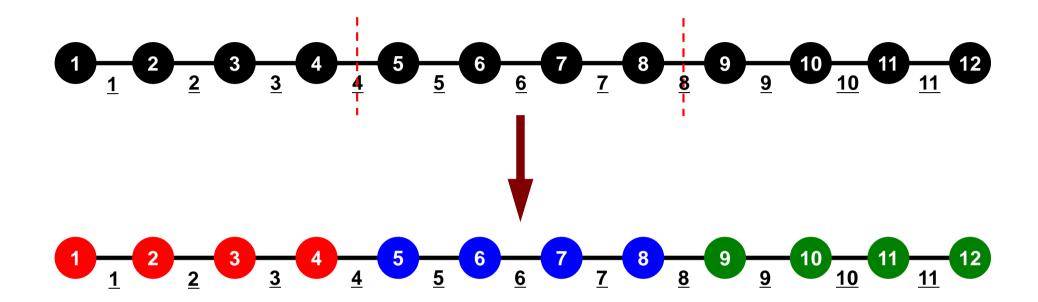


 to get information of "external nodes" from external partitions (local data)

"Communication tables" contain the information

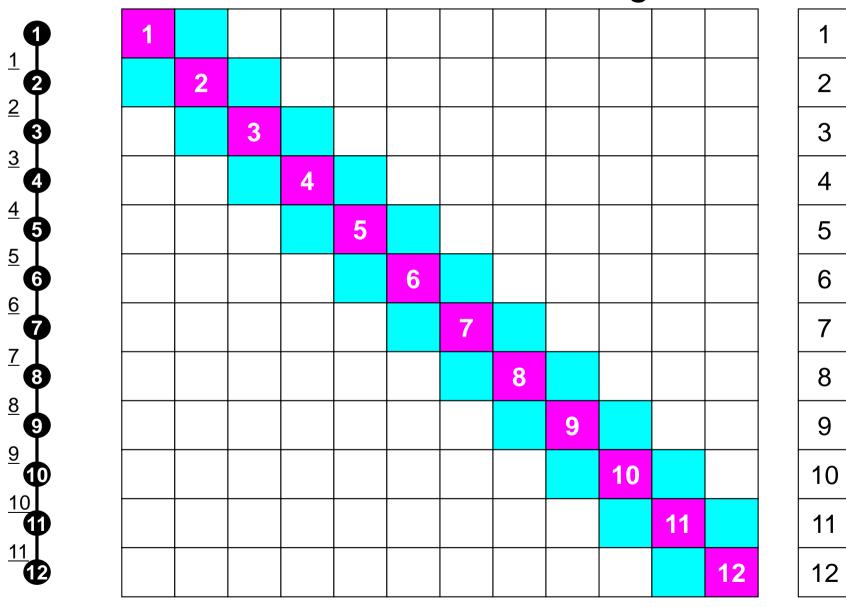
Intro-pFEM 19

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

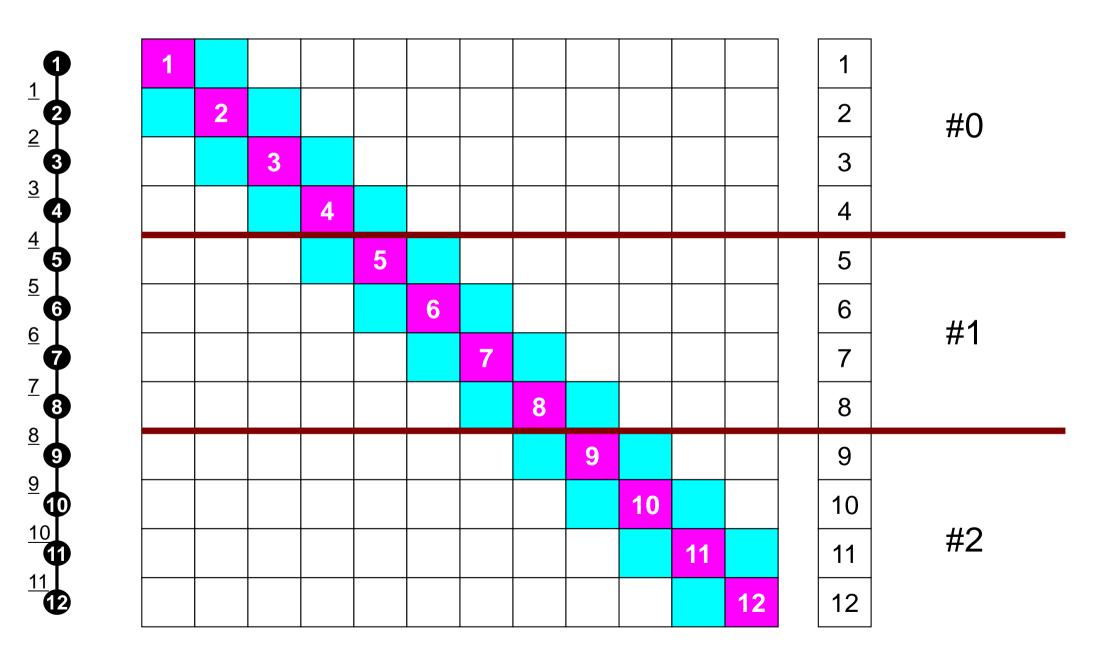


#### 1D FEM: 12 nodes/11 elem's/3 domains

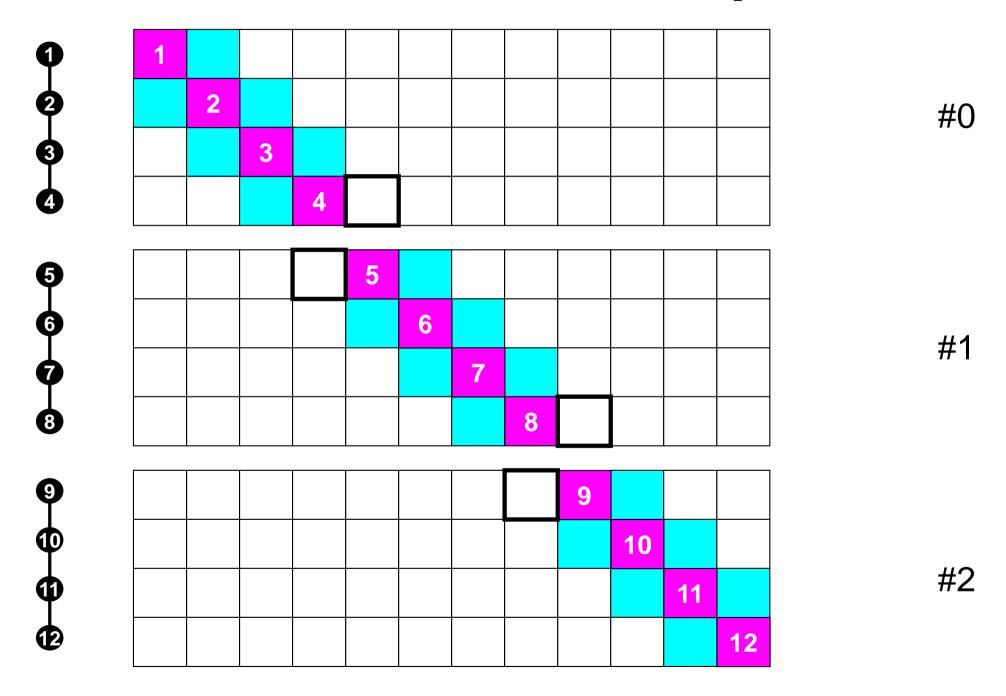
三重対角行列: Tri-Diagonal Matrix



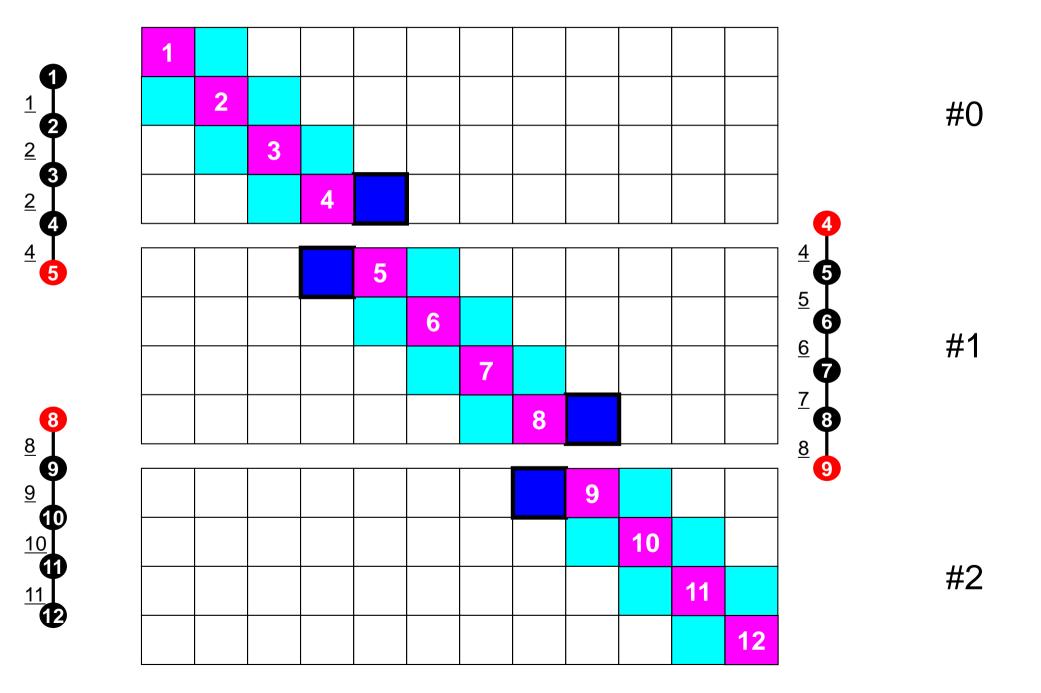
#### # "Internal Nodes" should be balanced



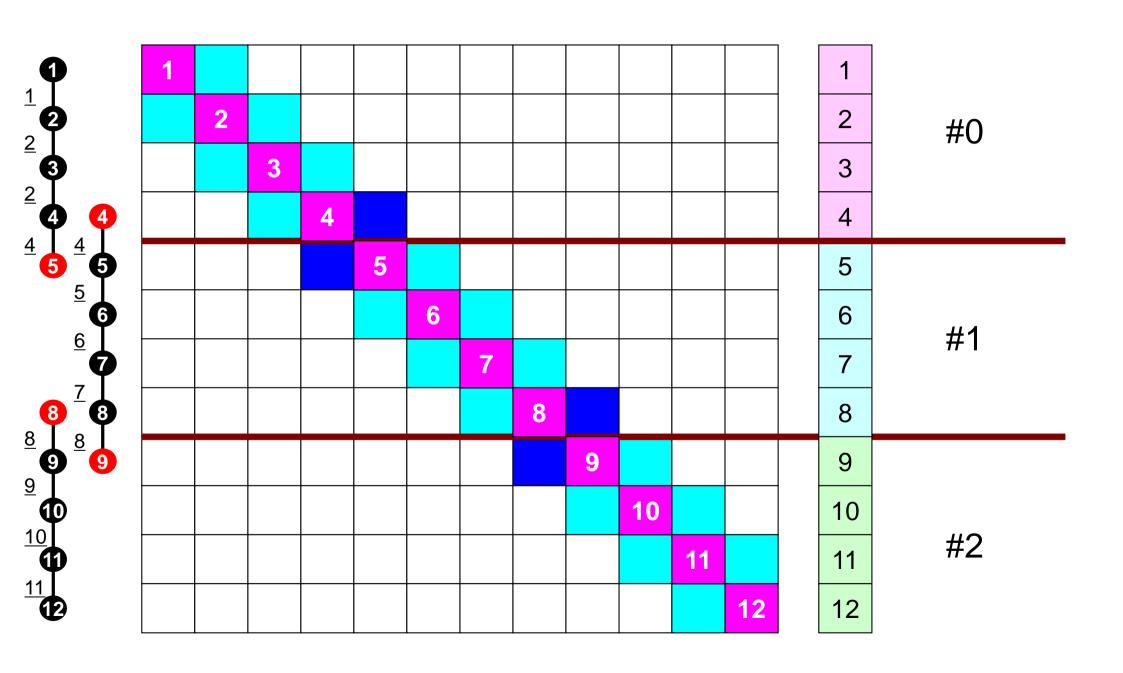
#### Matrices are incomplete!



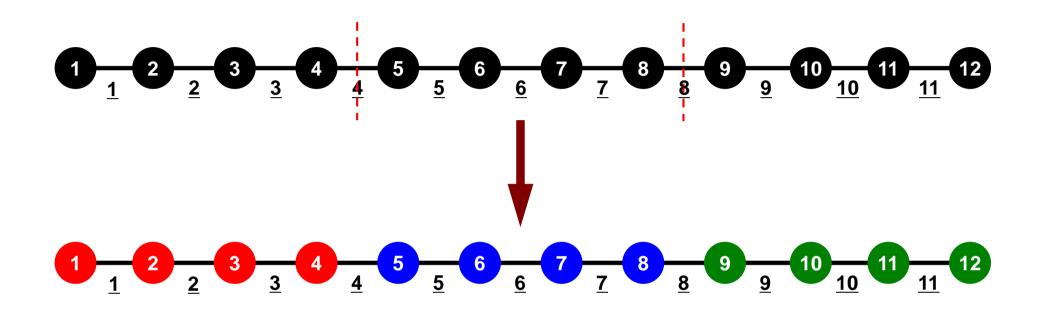
#### **Connected Elements + External Nodes**

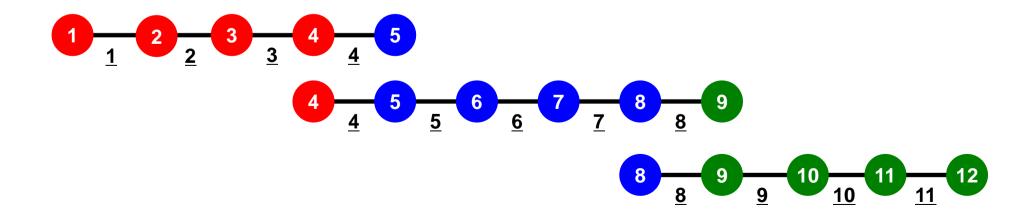


#### 1D FEM: 12 nodes/11 elem's/3 domains



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

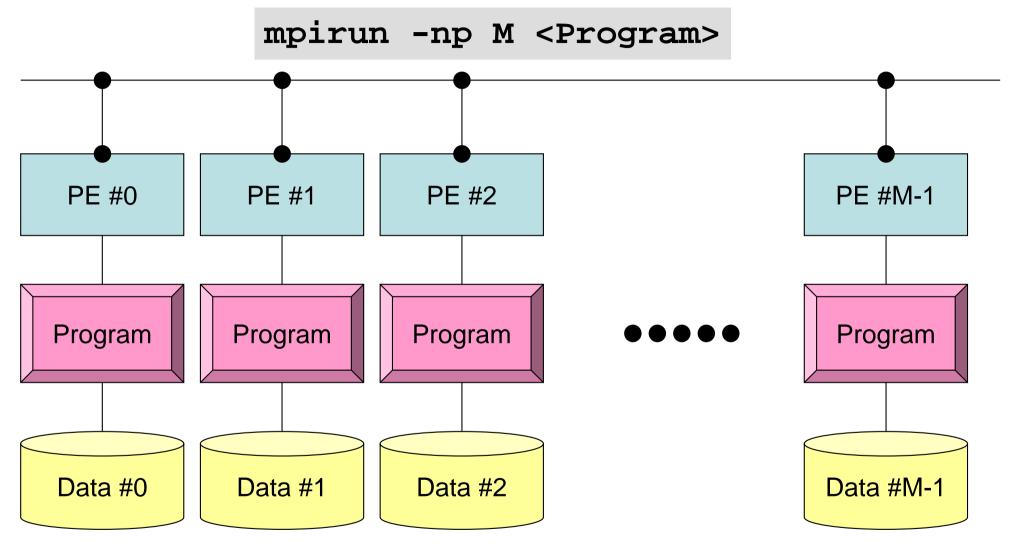




PE: Processing Element Processor, Domain, Process

#### **SPMD**

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

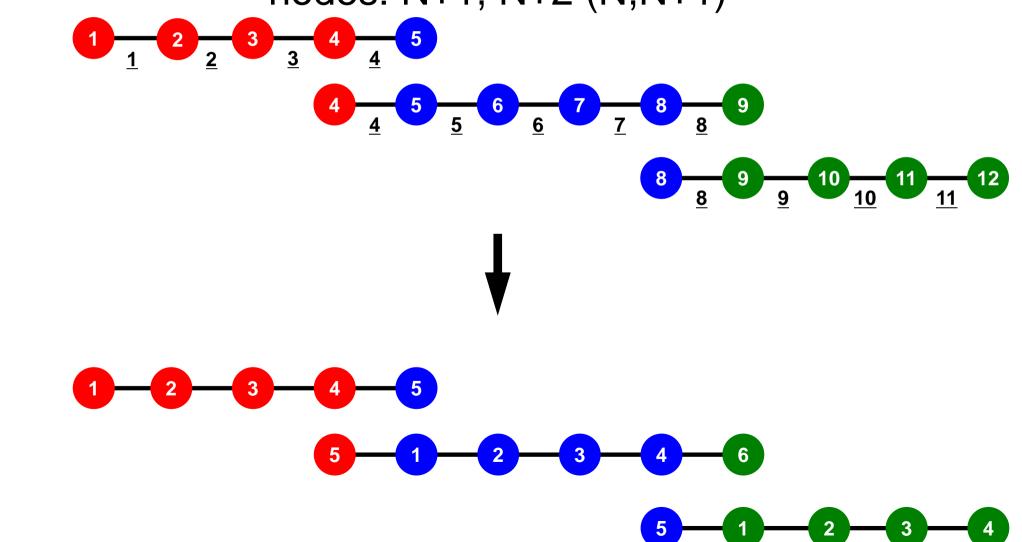


Each process does same operation for different data

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.

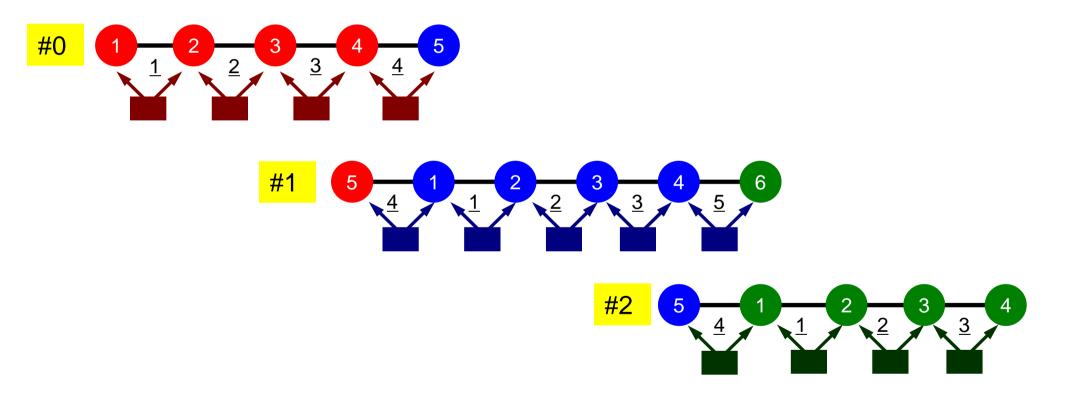
#### **Local Numbering for SPMD**

Numbering of internal nodes is 1-N (0-N-1), same operations in serial program can be applied. Numbering of external nodes: N+1, N+2 (N,N+1)

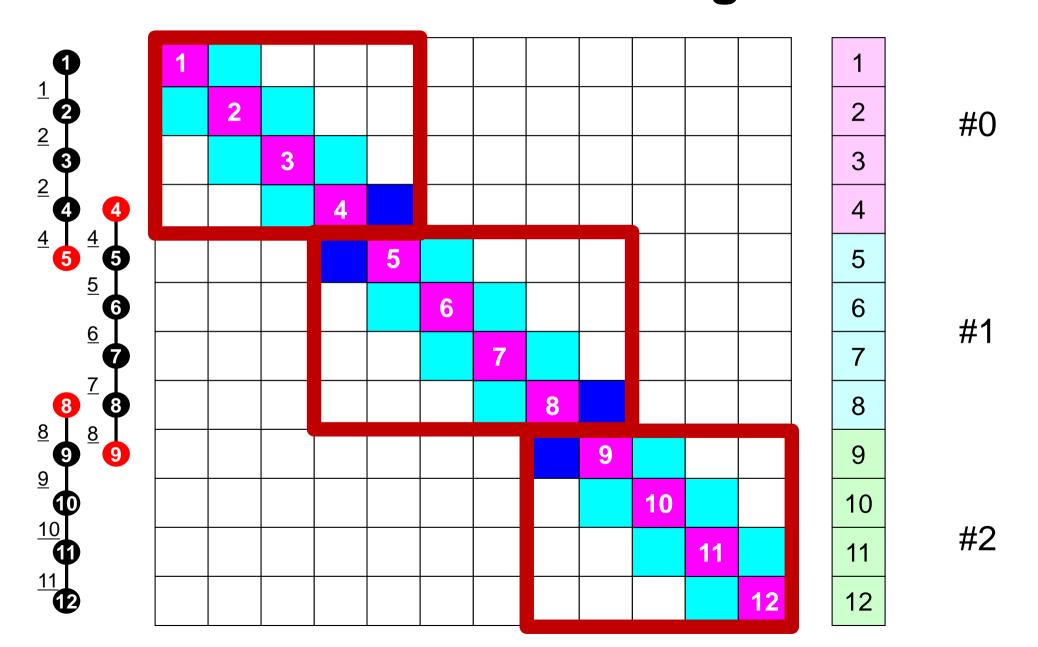


#### 1D FEM: 12 nodes/11 elem's/3 domains

Integration on each element, element matrix -> global matrix Operations can be done by info. of internal/external nodes and elements which include these nodes



#### Because the matrix is sparse, the union of the local matrices forms the global matrix!



#### **Finite Element Procedures**

- Initialization
  - Control Data
  - Node, Connectivity of Elements (N: Node#, NE: Elem#)
  - Initialization of Arrays (Global/Element Matrices)
  - Element-Global Matrix Mapping (Index, Item)
- Generation of Matrix
  - Element-by-Element Operations (do icel= 1, NE)
    - Element matrices
    - Accumulation to global matrix
  - Boundary Conditions
- Linear Solver
  - Conjugate Gradient Method

#### **Preconditioned CG Solver**

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

- Preconditioning
  - Diagonal Scaling/Point Jacobi
- Parallel operations are required in
  - Dot Products
  - Mat-Vec. Multiplication
    - SpMV: Sparse Mat-Vec. Mult.

#### Preconditioning, DAXPY

Local Operations by Only Internal Points: Parallel

Processing is possible

```
!C
!C-- {z}= [Minv] {r}

do i= 1, N
    W(i, Z)= W(i, DD) * W(i, R)
enddo
```

```
!C
!C-- {x}= {x} + ALPHA*{p}
!C {r}= {r} - ALPHA*{q}

do i= 1, N
    PHI(i)= PHI(i) + ALPHA * W(i, P)
    W(i, R)= W(i, R) - ALPHA * W(i, Q)
    enddo
DAXPY: double a{x} plus {y}

DAXPY: double a{x} plus {y}

LAPPHA*{q}

DAXPY: double a{x} plus {y}

LAPPHA*{q}

DAXPY: double a{x} plus {y}

LAPPHA*{q}

IN Company to the plus {y}

PHI(i) = PHI(i) + ALPHA * W(i, P)

W(i, R) = W(i, R) - ALPHA * W(i, Q)

Enddo
```

#### **Dot Products**

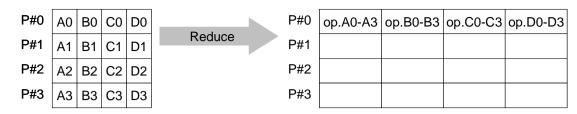
Global Summation needed: Communication 2

```
!C
!C-- ALPHA= RHO / {p} {q}

C1= 0. d0
    do i= 1, N
        C1= C1 + W(i, P)*W(i, Q)
    enddo
    ALPHA= RHO / C1
```

MPI Programming 35

#### MPI\_REDUCE



- Reduces values on all processes to a single value
  - Summation, Product, Max, Min etc.
- call MPI\_REDUCE

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

```
- <u>sendbuf</u>
             choice
                               starting address of send buffer
                               starting address receive buffer
recvbuf choice
                               type is defined by "datatype"
                               number of elements in send/receive buffer
  count
                               data type of elements of send/recive buffer
                      Τ
datatype I
    FORTRAN MPI_INTEGER, MPI_REAL, MPI_DOUBLE PRECISION, MPI CHARACTER etc.
             MPI_INT, MPI_FLOAT, MPI_DOUBLE, MPI_CHAR etc
                               reduce operation
  op
    MPI_MAX, MPI_MIN, MPI_SUM, MPI_PROD, MPI_LAND, MPI_BAND etc
    Users can define operations by MPI OP CREATE
```

_	<u>root</u>	I	I	rank of root process
_	comm	I	I	communicator
_	<u>ierr</u>	I	Ο	completion code

Fortran

#### **Preconditioned CG Solver**

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

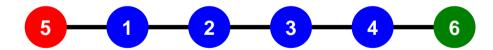
- Preconditioning
  - Diagonal Scaling/Point Jacobi
- Parallel operations are required in
  - Dot Products
  - Mat-Vec. Multiplication
    - SpMV: Sparse Mat-Vec. Mult.

#### Matrix-Vector Products

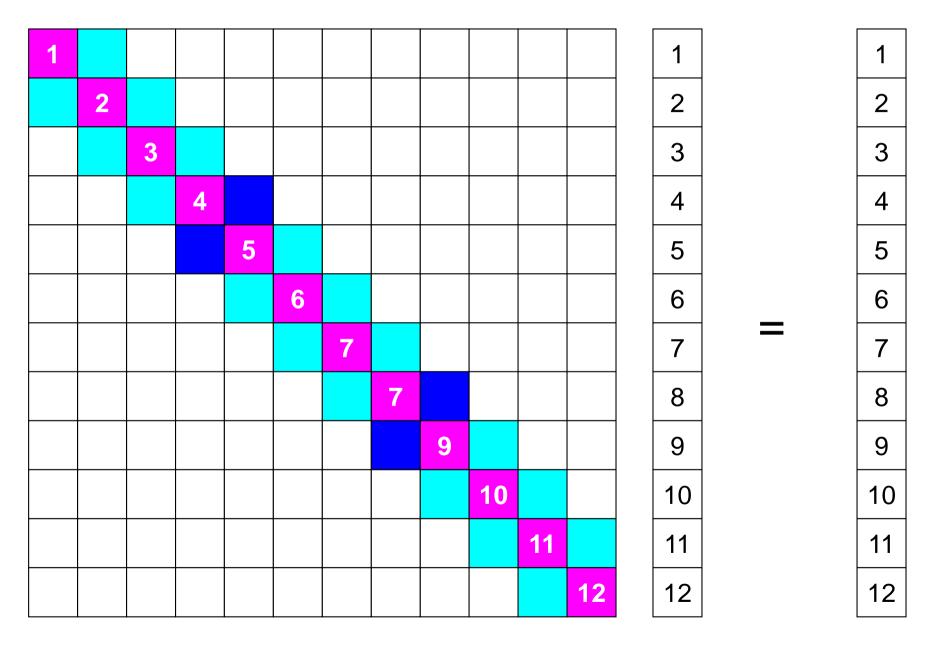
Values at External Points: P-to-P Communication

```
!C
!C-- {q}= [A] {p}

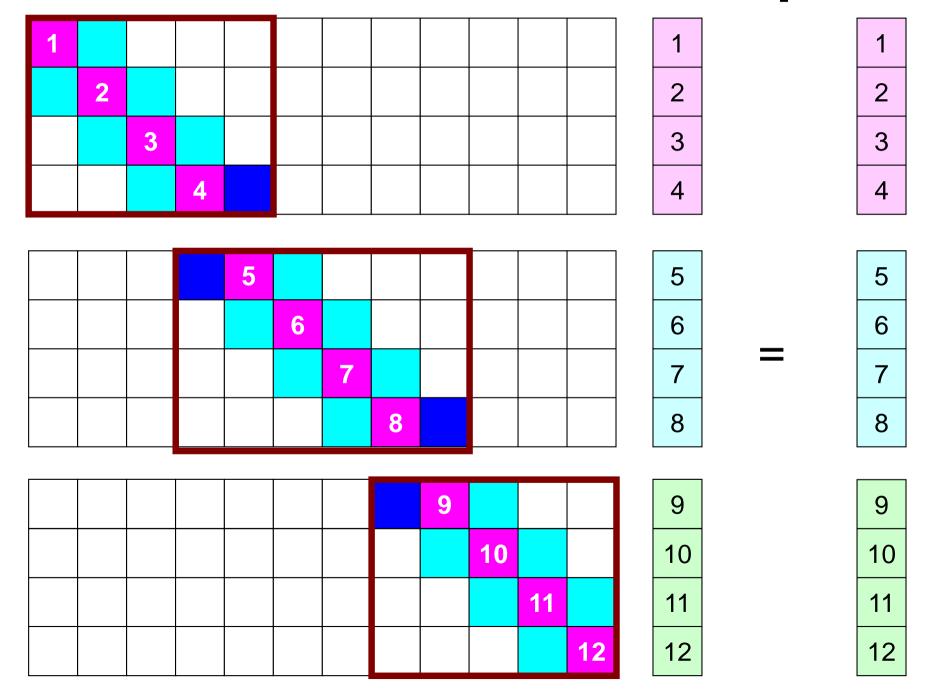
do i= 1, N
    W(i,Q) = DIAG(i)*W(i,P)
    do j= INDEX(i-1)+1, INDEX(i)
        W(i,Q) = W(i,Q) + AMAT(j)*W(ITEM(j),P)
    enddo
enddo
```



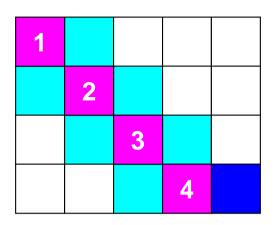
## Mat-Vec Products: Local Op. Possible

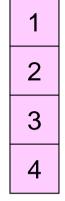


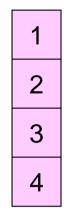
# Mat-Vec Products: Local Op. Possible

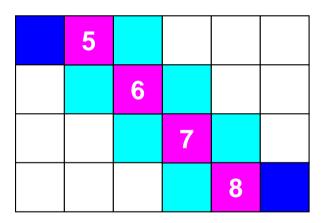


## Mat-Vec Products: Local Op. Possible

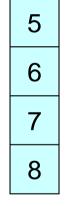




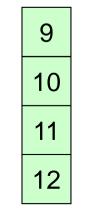


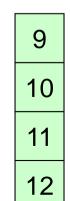


5	
6	
7	
8	

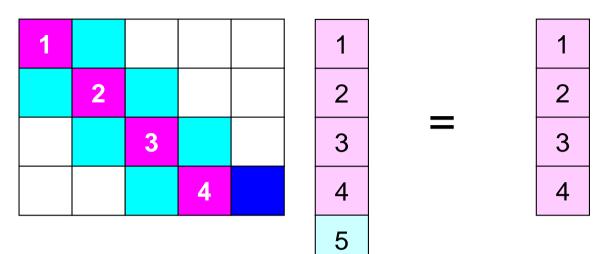


9			
	10		
		11	
			12



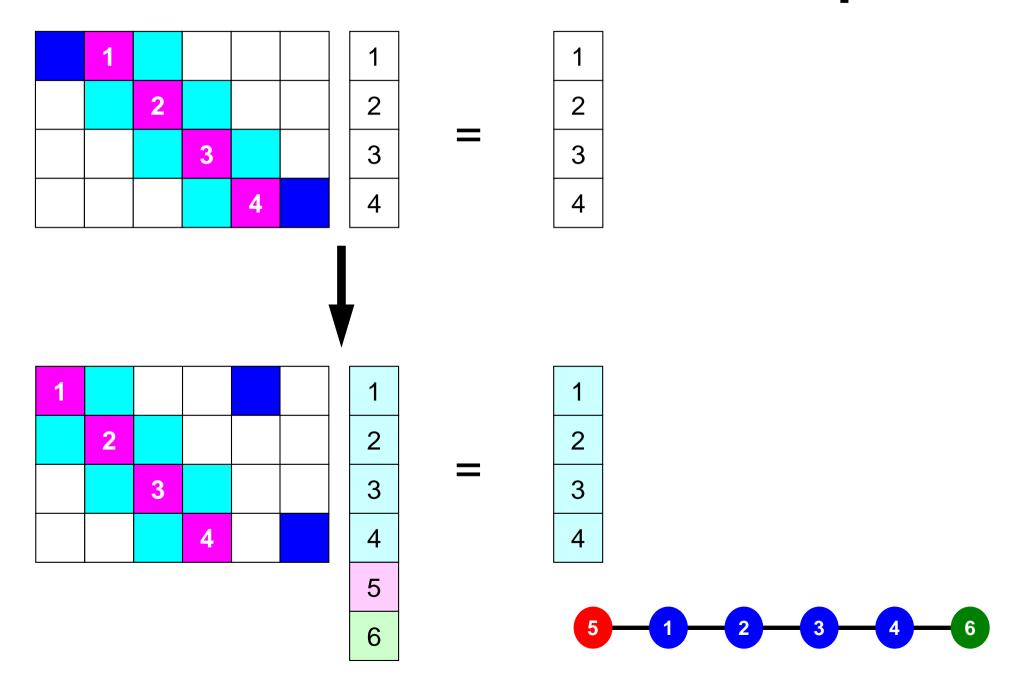


### Mat-Vec Products: Local Op. #0

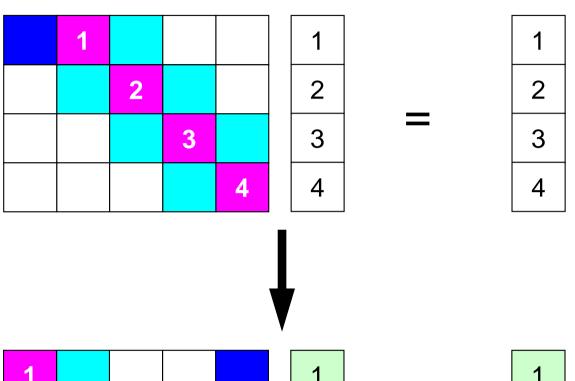




# Mat-Vec Products: Local Op. #1



# Mat-Vec Products: Local Op. #2

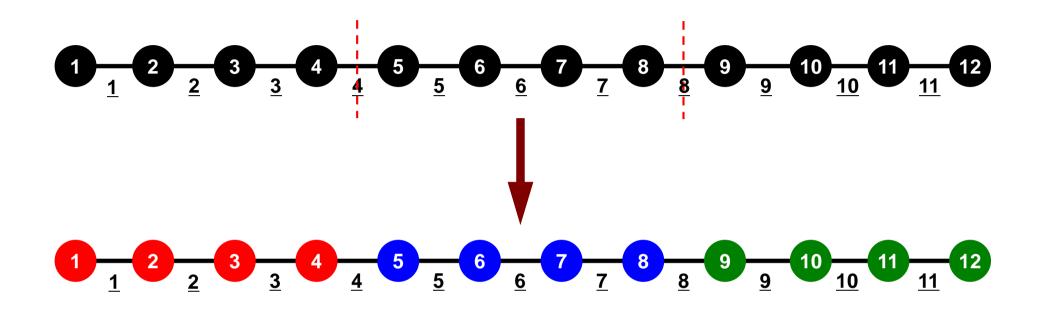


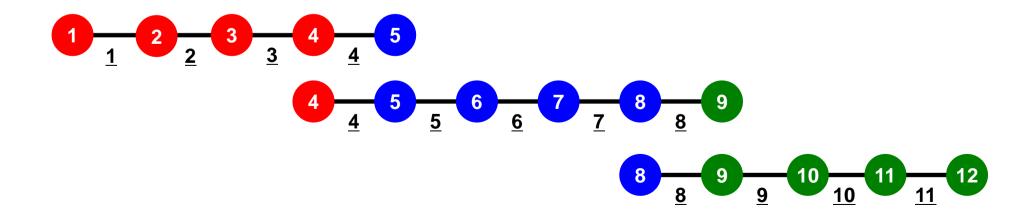
1					1
	2				2
		3			3
			4		4
					5





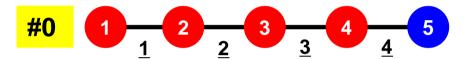
#### 1D FEM: 12 nodes/11 elem's/3 domains

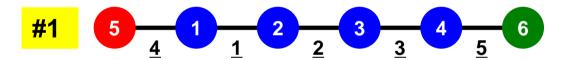


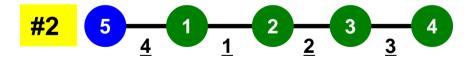


#### 1D FEM: 12 nodes/11 elem's/3 domains

Local ID: Starting from 1 for node and elem at each domain



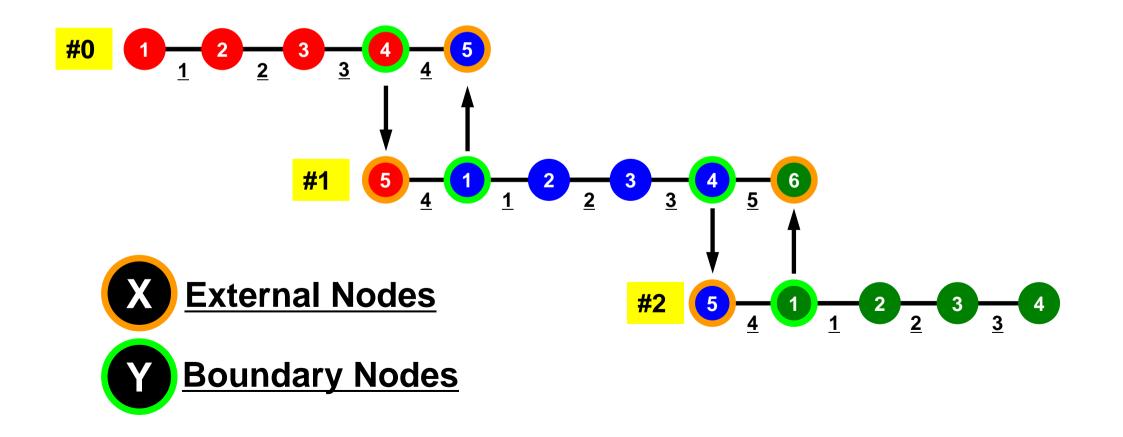




#### 1D FEM: 12 nodes/11 elem's/3 domains

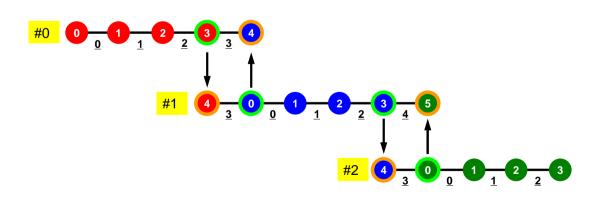
Internal/External/Boundary Nodes

Boundary Nodes: Part of Internal Nodes, and External Nodes of Other Domains



#### 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\_Recv
  - Communication with limited processes
    - Neighbors
  - Application Area
    - FEM, FDM: Localized Method





### 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.

# call MPI\_ISEND (sendbuf,count,datatype,dest,tag,comm,request, ierr)

_	<u>sendbuf</u>	choice	I	starting address of sending buffer
_	<u>count</u>	I	I	number of elements sent to each process
_	<u>datatype</u>	I	I	data type of elements of sending buffer
_	<u>dest</u>	I	I	rank of destination
_	<u>tag</u>	I	I	message tag
				This integer can be used by the application to distinguish
				messages. Communication occurs if tag's of
				MPI_Isend and MPI_Irecv are matched.
				Usually tag is set to be "0" (in this class),
_	comm	I	I	communicator
_	<u>request</u>	I	0	communication request array used in MPI_Waitall
_	<u>ierr</u>	I	0	completion code





- 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.

# call MPI\_IRECV (recvbuf,count,datatype,dest,tag,comm,request, ierr)

_	<u>recvbuf</u>	choice	I	starting address of receiving buffer
_	<u>count</u>	I	I	number of elements in receiving buffer
_	<u>datatype</u>	I	I	data type of elements of receiving buffer
_	<u>source</u>	I	I	rank of source
_	<u>tag</u>	I	I	message tag
				This integer can be used by the application to distinguish messages. Communication occurs if tag's of MPI_Isend and MPI_Irecv are matched. Usually tag is set to be "0" (in this class),
_	COMM	I	I	communicator
_	<u>request</u>	I	0	communication request used in MPI_Waitall
_	<u>ierr</u>	I	0	completion code

### MPI\_WAITALL



- MPI\_Waitall blocks until all comm's, associated with <u>request</u> in the array, complete. It is used for synchronizing <u>MPI\_Isend</u> and <u>MPI\_Irecv</u> 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 consitent.
  - Same <u>request</u> should be used in <u>MPI\_Isend</u> and <u>MPI\_Irecv</u>.
- 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.
- call MPI\_WAITALL (count, request, status, ierr)