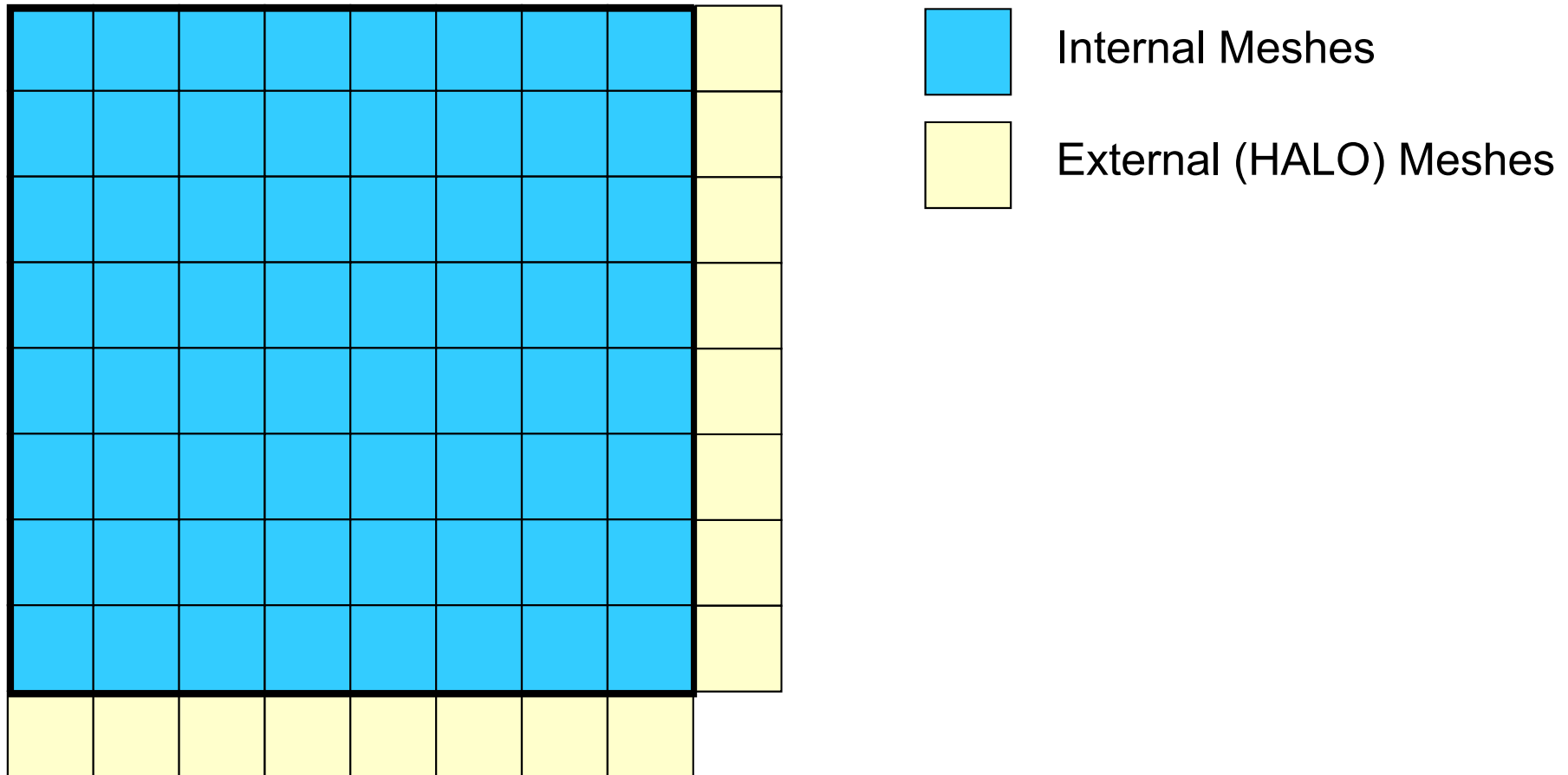


Introduction to Parallel Programming for Multicore/Manycore Clusters

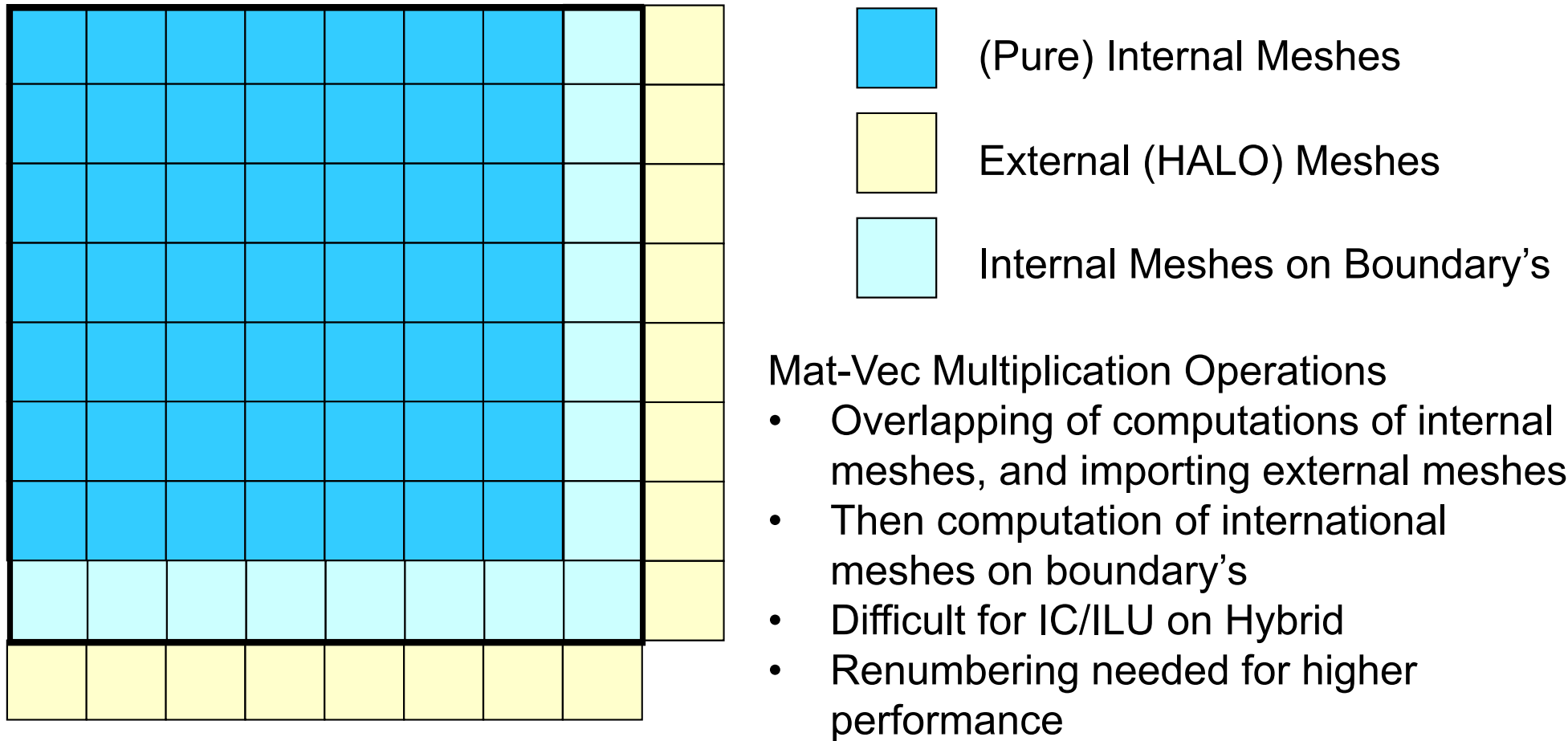
Part II-5: Communication-Computation Overlapping and Some Other Issues

Kengo Nakajima
Information Technology Center
The University of Tokyo

Communication-Computation Overlapping



Communication-Computation Overlapping



Comm.-Comp. Overlapping

<\$O-fvm/src1>

Without Reordering

```
call MPI_Isend
call MPI_Irecv
```

```

█ do i= 1, Nall
    if (BOUNDARY(i).eq.0) then
        (calculations)
    endif
enddo
```

```
call MPI_Waitall
```

```

█ do i= 1, Nall
    if (BOUNDARY(i).eq.1) then
        (calculations)
    endif
enddo
```

<\$O-fvm/src2>

With Reordering

```
call MPI_Isend
call MPI_Irecv
```

```

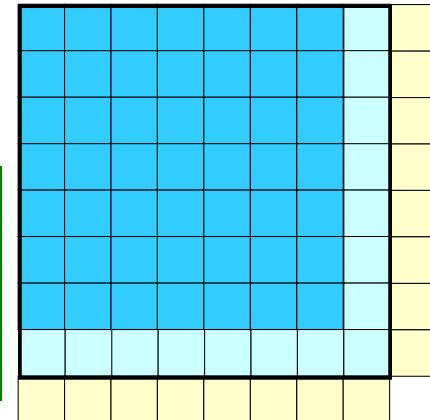
█ do i= 1, Ninn
    (calculations)
enddo
```

```
call MPI_Waitall
```

```

█ do i= Ninn+1, Nall
    (calculations)
enddo
```

Communications and computation for pure internal meshes are overlapped. Therefore, problem size should be large enough for *hiding* communication overhead



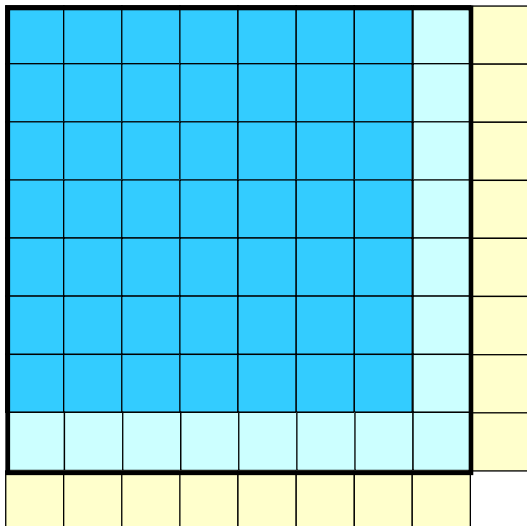
poi_gen (<\$src1>)



```

!$omp parallel do
  do i= 1, NP
    BOUNDARY(i)= 0
  enddo

!$omp parallel do
  do i= 1, N
    do k= indexLU(i-1)+1, indexLU(i)
      jj= itemLU(k)
      if (jj.gt.N) BOUNDARY(i)= 1
    enddo
  enddo

```



If the mesh “i” is connected to “external” meshes  , the mesh is internal mesh on the boundary  .

Finally “BOUNDARY(i)” is set to “1”.

[A]{p}={q} (<\$src1>)

```
!C
!C-- SEND/RECV
```

```
(...) do neib= 1, NEIBPETOT
      call MPI_ISEND (...)
    enddo

(...) do neib= 1, NEIBPETOT
      call MPI_IRECV (...)
    enddo
```

```
!$omp parallel do private (i,k,VAL)
do i= 1, N
  if (BOUNDARY(i).eq.0) then
    VAL= D(i)*W(i,P)
    do k= indexLU(i-1)+1, indexLU(i)
      VAL= VAL + AMAT(k)*W(itemLU(k),P)
    enddo
    W(i,Q)= VAL
  endif
enddo
```

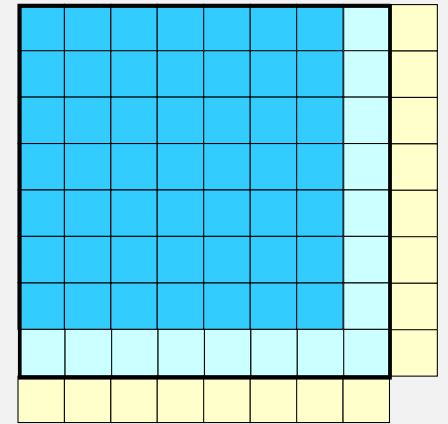
```
call MPI_WAITALL (2*NEIBPETOT, req1, sta1, ierr)
```

```
!$omp parallel do private (i,k,VAL)
do i= 1, N
  if (BOUNDARY(i).eq.1) then
    VAL= D(i)*W(i,P)
    do k= indexLU(i-1)+1, indexLU(i)
      VAL= VAL + AMAT(k)*W(itemLU(k),P)
    enddo
    W(i,Q)= VAL
  endif
enddo
```

Pure Internal
Meshes

Overlapping

Boundary
Meshes



Communications and computation for pure internal meshes are overlapped. Therefore, problem size should be large enough for *hiding* communication overhead

```
!$omp parallel do private (icel, icN1, icN2, icN3, icN4, icN5, icN6)
do icel= 1, ICELTOT
  icN1= NEIBcell(icel, 1) ...
  if (icN1.gt. ICELTOT) BNODE(icel)= 1
  if (icN2.gt. ICELTOT) BNODE(icel)= 1
  if (icN3.gt. ICELTOT) BNODE(icel)= 1 ...
enddo
```

BNODE(icel)=1
if "icel" is on boundary's

poi_gen (<\$src2>)

```
icou= 0
do icel= 1, ICELTOT
  if (BNODE(icel).eq.0) then
    icou= icou + 1
    OLDtoNEW(icel)= icou
    NEWtoOLD(icou)= icel
  endif
enddo
```

Renumbering
Pure Internal Meshes
ICELTOTinn

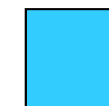
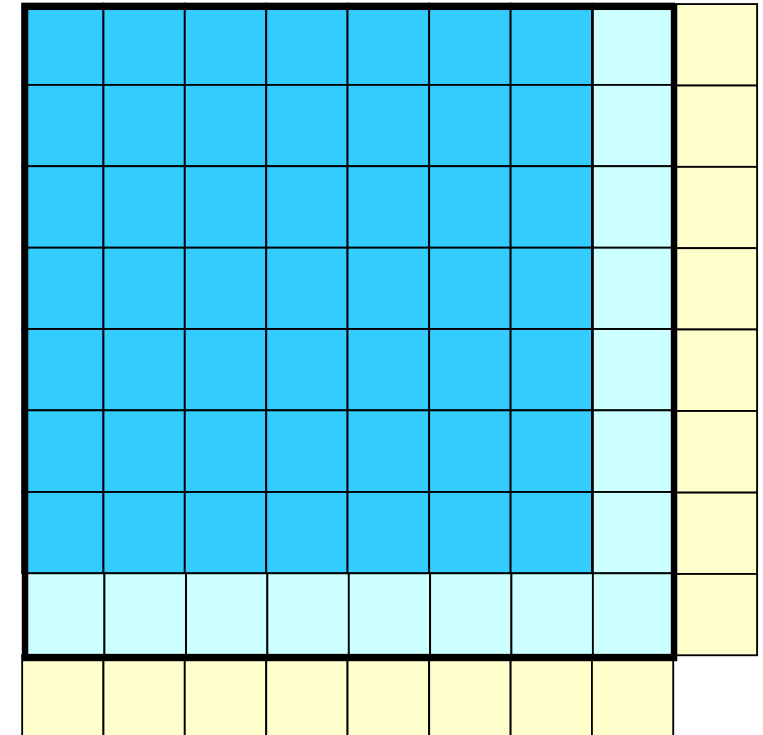
```
ICELTOTinn= icou
do icel= 1, ICELTOT
  if (BNODE(icel).eq.1) then
    icou= icou + 1
    OLDtoNEW(icel)= icou
    NEWtoOLD(icou)= icel
  endif
enddo
```

Renumbering
Boundary Meshes

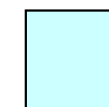
```
do neib= 1, NEIBPETOT
  do k= EXPORT_INDEX(neib-1)+1, EXPORT_INDEX(neib)
    BNODE(k)= EXPORT_ITEM(k)
  enddo
enddo
```

Renumbering
EXPORT_ITEM

```
do neib= 1, NEIBPETOT
  do k= EXPORT_INDEX(neib-1)+1, EXPORT_INDEX(neib)
    EXPORT_ITEM(k)= OLDtoNEW(BNODE(k))
  enddo
enddo
```



Pure Internal Meshes



Internal Meshes on
Boundary's

[A]{p}={q} (<\$src2>)

```
!C
!C-- SEND/RECV
```

```
(...) do neib= 1, NEIBPETOT
      call MPI_ISEND (...)
    enddo

(...) do neib= 1, NEIBPETOT
      call MPI_IRecv (...)
    enddo
```

```
!$omp parallel do private (i,k,VAL)
do i= 1, Ninn
  VAL= D(i)*W(i,P)
  do k= indexLU(i-1)+1, indexLU(i)
    VAL= VAL + AMAT(k)*W(itemLU(k),P)
  enddo
  W(i,Q)= VAL
enddo
```

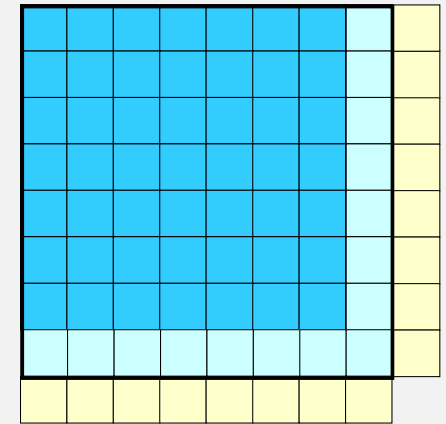
Pure Internal
Meshes

Overlapping

```
call MPI_WAITALL (2*NEIBPETOT, req1, sta1, ierr)
```

```
!$omp parallel do private (i,k,VAL)
do i= Ninn+1, N
  VAL= D(i)*W(i,P)
  do k= indexLU(i-1)+1, indexLU(i)
    VAL= VAL + AMAT(k)*W(itemLU(k),P)
  enddo
  W(i,Q)= VAL
enddo
```

Boundary
Meshes



Communications and computation for pure internal meshes are overlapped. Therefore, problem size should be large enough for *hiding* communication overhead

Compile & Run (Hybrid Only)

<\$src1>

```
>$ cd
>$ cd hybrid/fvm/src1
>$ make clean
>$ make
>$ ls ../run/sol1-mpih
    sol1-mpi
>$ cd ../run
(modify go1.sh, INPUT.DAT)
>$ pjsub go1.sh (or mpiexec ...)
```

<\$src2>

```
>$ cd
>$ cd hybrid/fvm/src2
>$ make clean
>$ make
>$ ls ../run/sol2-mpih
    sol2-mpi
>$ cd ../run
(modify go2.sh, INPUT.DAT)
>$ pjsub go2.sh (or mpiexec ...)
```

Results on Oakleaf-FX

64 nodes, 64 processes, HB 16x1

Time for CG Solvers (sec.)

DOF (N _A x)	<\$O-fvm>/src Original	<\$O-fvm>/src0 Continuous Access	<\$O-fvm>/src1 Comm-Comp Overlapping (if-then-branch)	<\$O-fvm>/src2 Comm-Comp Overlapping (reordered)
508 ³	33.0	33.0	38.7	34.0
512 ³	78.1	78.0	87.3	36.3
516 ³	35.7	35.5	41.7	46.8
636 ³	82.4	82.0	93.9	84.4
640 ³	109.6	109.4	127.4	88.8
644 ³	85.0	85.1	97.3	87.5
744 ³	160.0	159.8	182.5	154.9

Bank conflict occurs at 512³ & 640³

In <\$src2> it is avoided due to renumbering

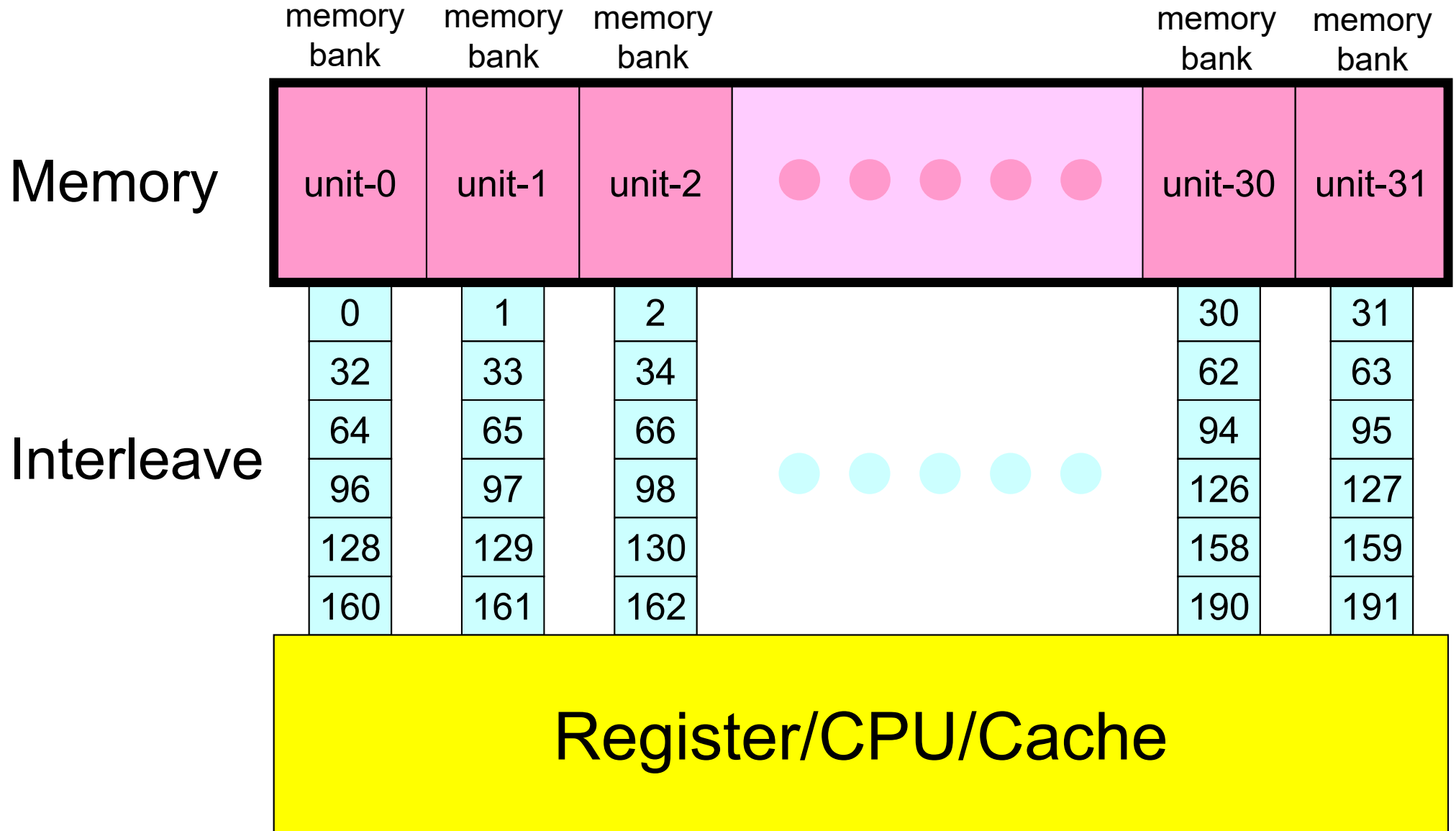
DOF (NAX)	<\$O-fvm>/src Original	<\$O-fvm>/src0 Continuous Access	<\$O-fvm>/src1 Comm-Comp Overlapping (if-then-branch)	<\$O-fvm>/src2 Comm-Comp Overlapping (reordered)
508 ³	33.0	33.0	38.7	34.0
512 ³	78.1	78.0	87.3	36.3
516 ³	35.7	35.5	41.7	46.8
636 ³	82.4	82.0	93.9	84.4
640 ³	109.6	109.4	127.4	88.8
644 ³	85.0	85.1	97.3	87.5
744 ³	160.0	159.8	182.5	154.9

Memory Interleaving/Bank Conflict

- Memory Interleaving
 - Method for fast data transfer to/from memory.
 - Parallel I/O for multiple memory banks.
- Memory Bank
 - Unit for memory management, small pieces of memory
 - Usually, there are 2^n independent modules.
 - Single bank can execute a single reading or writing at one time. Therefore, performance gets worse if data components on same bank are accessed simultaneously.
 - For example, “bank conflict” occurs if off-set of data access is 32 (in next page).
 - Remedy: Change of array size, loop exchange, reordering etc.

Bank Conflict

If off-set of data access is 32, only a single bank is utilized



Avoiding Bank Conflict

X

```
REAL*8 A(32,10000)
```

```
k= N
```

```
do i= 1, 10000
```

```
    A(k,i)= 0.d0
```

```
enddo
```

O

```
REAL*8 A(33,10000)
```

```
k= N
```

```
do i= 1, 10000
```

```
    A(k,i)= 0.d0
```

```
enddo
```

- Arrays with size of 2^n should be avoided.

If $N=2^m$, e.g. 8x8

- Bank conflict always occurs when non-zero off-diagonals are accessed, 16 threads
- Remedy for Bank Conflict
 - Padding by compiler, Reordering

73	74	75	76	77	78	79	80	
57	58	59	60	61	62	63	64	72
49	50	51	52	53	54	55	56	71
41	42	43	44	45	46	47	48	70
33	34	35	36	37	38	39	40	69
25	26	27	28	29	30	31	32	68
17	18	19	20	21	22	23	24	67
9	10	11	12	13	14	15	16	66
1	2	3	4	5	6	7	8	65



73	74	75	76	77	78	79	80	
57	58	59	60	61	62	63	64	72
43	44	45	46	47	48	49	56	71
36	37	38	39	40	41	42	55	70
29	30	31	32	33	34	35	54	69
22	23	24	25	26	27	28	53	68
15	16	17	18	19	20	21	52	67
8	9	10	11	12	13	14	51	66
1	2	3	4	5	6	7	50	65

If you want to do these on Oakleaf-FX

by 12 nodes, 12 MPI processes, HB 16x1
128³ for each node (MPI process)

mg.sh

```
#!/bin/sh
#PJM -L "node=1"
#PJM -L "elapse=00:05:00"
#PJM -j
#PJM -L "rscgrp=lecture7"
#PJM -g "gt17"
#PJM -o "test.lst"
#PJM --mpi "proc=12"

mpiexec ./pmesh
```

mesh.inp

```
384 256 256
  3   2   2
```

goh.sh

```
#!/bin/sh
#PJM -L "node=12"
#PJM -L "elapse=00:05:00"
#PJM -j
#PJM -L "rscgrp=lecture7"
#PJM -g "gt17"
#PJM -o "test.lst"
#PJM --mpi "proc=12"

export OMP_NUM_THREADS=16
mpiexec ./sol-mpih
```


Cache Thrashing

- FX10: L1D cache with 32KB for each core, 2-way
 - n-way set associative (n群連想記憶式)
 - Cache is divided into “n” banks
 - Each bank is divided into “cache lines”
 - Number of Cache Lines, Size of Cache Line (128 bytes for FX10) $\Rightarrow 2^m$
- This “2-way” cache is very harmful
 - If “ $N=2^m$ ”, memory addresses of $W(i, P)$, $W(i, Q)$, $W(i, R)$ map to the same cache address in CG computation.
 - Cache Thrashing: Lower Performance
 - $R=1$, $P=2$, $Q=3$
 - $X(i)$ is not affected

```
!$omp parallel do private(i)
  do i= 1, N
    X(i) = X(i) + ALPHA * W(i, P)
    W(i, R) = W(i, R) - ALPHA * W(i, Q)
  enddo
```

Remedy

- If the loop is split into 2 loops, up to 2 cache lines of W are referred. Therefore, cache thrashing does not occur (Remedy-1).

```
!$omp parallel do private(i)
  do i= 1, N
    X(i) = X (i)  + ALPHA * W(i,P)
  enddo

!$omp parallel do private(i)
  do i= 1, N
    W(i,R) = W(i,R) - ALPHA * W(i,Q)
  enddo
```

- If “ $N=2^m$ “, certain numbers (e.g. 64, 128 ...) can be added to N . Thus, size of the array is not equal to 2^m , and cache thrashing is avoided (Remedy-2).
 - No such operation is needed for x

```
N2=128
allocate (W(N+N2, 4))
```

Remedy-2: Try by yourself $\langle \\$src2 \rangle$

solver_PCG: C

```

N3= N+N2; N2=128
W = (double **)malloc(sizeof(double *)*4);
...
for (i=0; i<4; i++) {
    W[i] = (double *)malloc(sizeof(double)*N3);
    ...

```

solver_PCG: Fortran

```

allocate (W(N+N2, 4)) N2=128

```

Results Fortran

	NX=NY=NZ=128 2,097,152 meshes Load/Store= 8.28×10^{10}	NX=NY=NZ=129 2,146,689 meshes Load/Store= 8.53×10^{10}
sol10 (original)	19.50 sec. 24.11 GB/sec 13.59 %	9.15 sec. 52.98 GB/sec 3.97 %
sol20 (CM reordering)	10.15 sec. 45.60 GB/sec 5.65 %	9.44 sec. 50.64 GB/sec 4.11 %
sol2x (CM+ Remedy-2)	9.69 sec. 47.77 GB/sec 4.20 %	9.54 sec. 50.12 GB/sec 4.11 %

C (N=128³)

- sol10: 20.03 sec.
- sol20: 13.72
- sol2x: 10.05

- **Comp. Time**
- **Memory Throughput**
- **L1D Miss Ratio**
(to Load/Store)

Cache Thrashing, Bank Conflict

- Problem is that the cache is 2-way on FX10
 - This is the reason for both of bank conflict and cache thrashing
- Most of modern architectures based on 4-way or 8-way
 - e.g. Intel CPU
- FX-100 (successor of FX10) has 4-way cache

Back to “Communication-Computation Overlapping”

- Problem size should be large enough for hiding communications by overlapping.
- Number of compute nodes should be large enough (too small in this case)
- Generally speaking, effect of communication/computation overlapping is rather smaller for Krylov iterative solvers
 - Explicit time-marching method by FDM

[A]{p}={q} (<\$src2>)

```
!C
!C-- SEND/RECV
```

```
do neib= 1, NEIBPETOT
(...)
  call MPI_ISEND (...)
enddo

do neib= 1, NEIBPETOT
(...)
  call MPI_Irecv (...)
enddo
```

```
!$omp parallel do private (i,k,VAL)
do i= 1, Ninn
  VAL= D(i)*W(i,P)
  do k= indexLU(i-1)+1, indexLU(i)
    VAL= VAL + AMAT(k)*W(itemLU(k),P)
  enddo
  W(i,Q)= VAL
enddo
```

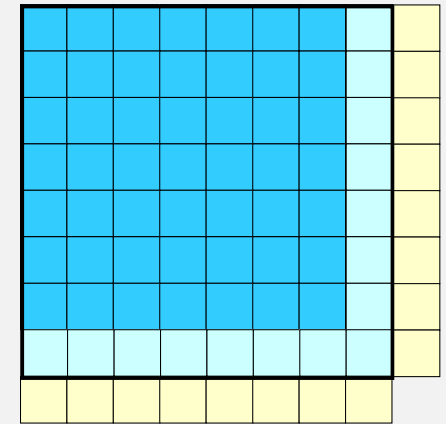
Pure Internal
Meshes

Overlapping

```
call MPI_WAITALL (2*NEIBPETOT, req1, sta1, ierr)
```

```
!$omp parallel do private (i,k,VAL)
do i= Ninn+1, N
  VAL= D(i)*W(i,P)
  do k= indexLU(i-1)+1, indexLU(i)
    VAL= VAL + AMAT(k)*W(itemLU(k),P)
  enddo
  W(i,Q)= VAL
enddo
```

Boundary
Meshes



Communications and computation for pure internal meshes are overlapped. Therefore, problem size should be large enough for *hiding* communication overhead

Bank conflict occurs at 512³ & 640³

In <\$src2> it is avoided due to renumbering

DOF (NAX)	<\$O-fvm>/src Original	<\$O-fvm>/src0 Continuous Access	<\$O-fvm>/src1 Comm-Comp Overlapping (if-then-branch)	<\$O-fvm>/src2 Comm-Comp Overlapping (reordered)
508 ³	33.0	33.0	38.7	34.0
512 ³	78.1	78.0	87.3	36.3
516 ³	35.7	35.5	41.7	46.8
636 ³	82.4	82.0	93.9	84.4
640 ³	109.6	109.4	127.4	88.8
644 ³	85.0	85.1	97.3	87.5
744 ³	160.0	159.8	182.5	154.9