# XcalableMP and XcalableACC for Productivity and Performance in HPC Challenge Award Competition

Masahiro Nakao,[†‡] Hitoshi Murai,[†] Hidetoshi Iwashita[†]

Takenori Shimosaka,[†] Akihiro Tabuchi,[*] Taisuke Boku[‡*]

Mitsuhisa Sato[†‡*]

† RIKEN Advanced Institute for Computational Science, Japan
‡ Center for Computational Sciences, University of Tsukuba
* Graduate School of Systems and Information Engineering, University of Tsukuba
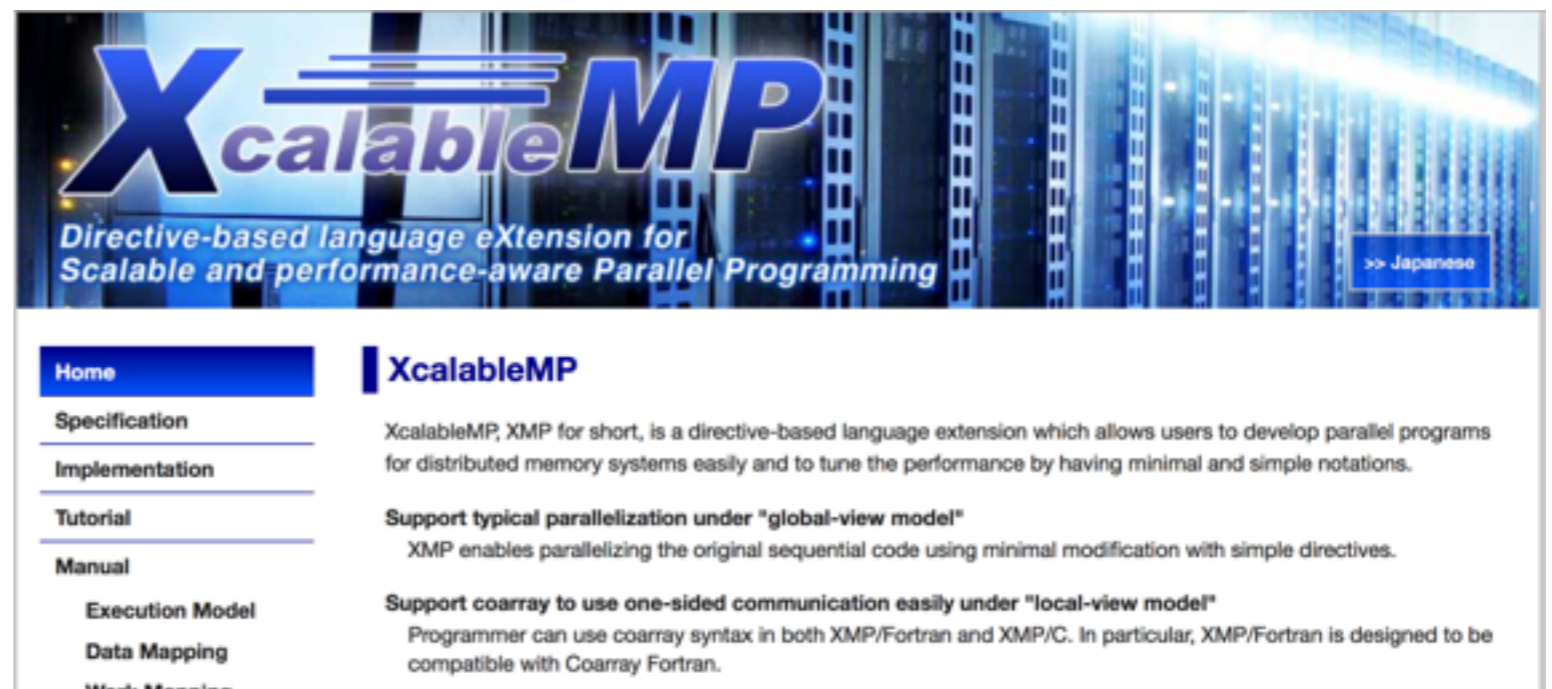
# Outline

1. XcalableMP (XMP) for cluster systems (14min.)

2. XcalableACC (XACC) for accelerator cluster systems (6min.)
   Extension of XMP using OpenACC

Sorry !!, work-in-progress

The submission report is available at *http://xcalablemp.org*

# What is XcalableMP (XMP) ?

**Directive-based language extensions of Fortran and C**

- By XMP specification working group of PC cluster consortium (SC Booth#2924)

- Version 1.2.1 specification available (*http://xcalablemp.org*)

**Support two memory models**

- Global-view (HPF-like data/work mapping directives)

- Local-view (coarray)

**Implementation of Compiler**

- Omni XMP Compiler version 0.9 (*http://omni-compiler.org*)

- Platforms: Fujitsu the K computer and FX10, Cray XT/XE, IBM BlueGene, NEC SX, Hitachi SR, Linux clusters, etc.

# Code example (Global-view)

```
int a[MAX];
#pragma xmp nodes p(4)
#pragma xmp template t(0:MAX-1)
#pragma xmp distribute t(block) on p
#pragma xmp align a[i] with t(i)


main(){
  int i, j, res = 0;



#pragma xmp loop on t(i) reduction(+:res)
  for(i = 0; i <MAX; i++){
    a[i] = func(i);
    res += array[i];
  }
}
```

Data distribution

add to the serial code : incremental parallelization

Work mapping and data synchronization

# Code example (Local-view)

```
double a[100]:[*], b[100]:[*];
int me = xmp_node_num();
```

Define Coarrays

```
if(me == 2)
  a[:]:[1] = b[:];
```

Put Operation

```
if(me == 1)
  a[0:50] = b[0:50]:[2];
```

Get Operation

**Coarray synax in XMP/C**

**array_name[start:length]:[node_number];**

**XMP/Fortran is upward compatible with Fortran 2008**

# Results and Machine

## Summary

| Benchmark | | # Nodes | Performance (/peak) | SLOC |
|-----------|-----|---------|---------------------|------|
| HPL | Ver. 1 | 16,384 | 971 TFlops (46.3%) | 313 |
| | Ver. 2 | 4,096 | 423 TFlops (**80.7%**) | 426 |
| FFT | | **82,944** | 212 TFlops (2.0%) | 205 |
| STREAM | | **82,994** | 3,583 TB/s (67.5%) | 69 |
| RandomAccess | | 16,384 | 254 GUPs | 253 |

## The K computer: 82,944 nodes



*http://www.aics.riken.jp/jp/outreach/photogallery.html*
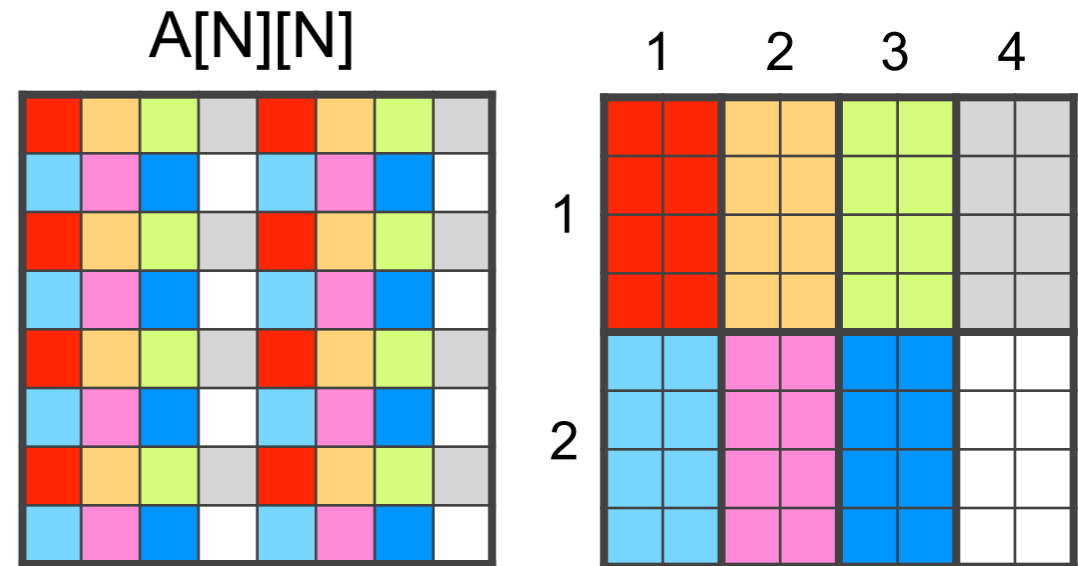
- SPARC64 VIIIfx Chip, 128 GFlops
- DDR3 SDRAM 16GB, 64GB/s
- Tofu Interconnect
  - 6D mesh/torus network
  - 5GB/s x 4links x 2

# HPL version 1

- Source lines of Code (SLOC) is **313**, written in XMP/C

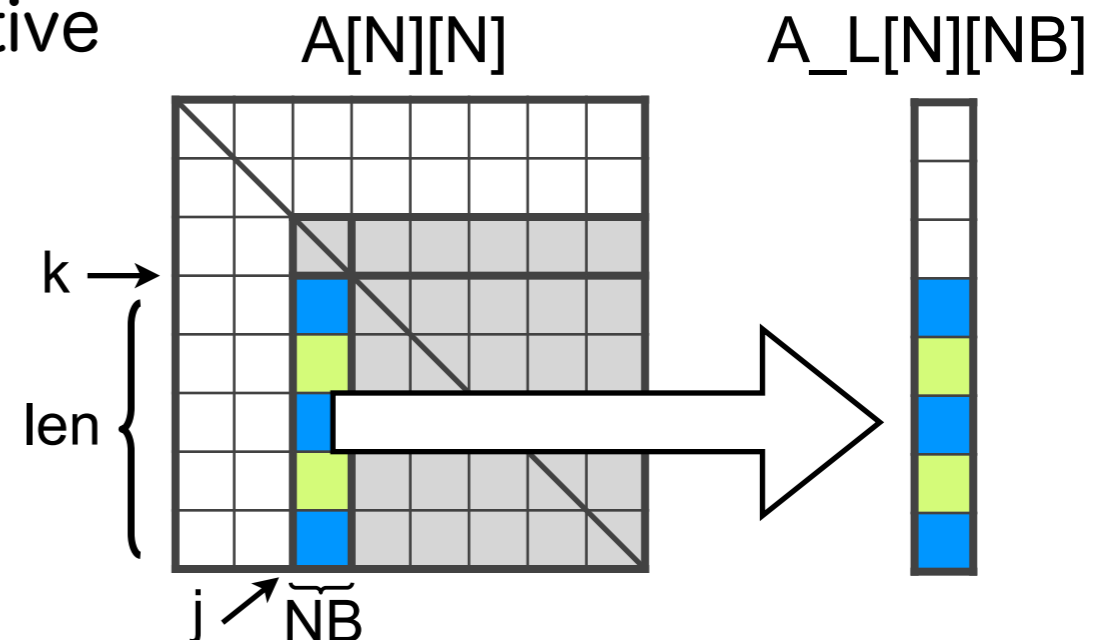- Block-Cyclic Distribution

A[N][N]



```
double A[N][N];
#pragma xmp nodes p(P,Q)
#pragma xmp template t(0:N-1, 0:N-1)
#pragma xmp distribute t(cyclic(NB), \
                cyclic(NB)) onto p
#pragma xmp align A[i][j] with t(j,i)
```

**Programmer can use BLAS for distributed array.**

- Panel Broadcast by using **gmove** directive

A[N][N]          A_L[N][NB]

```
double A_L[N][NB];
#pragma xmp align A_L[i][*] with t(*,i)
    :
#pragma xmp gmove
A_L[k:len][0:NB] = A[k:len][j:NB];
```
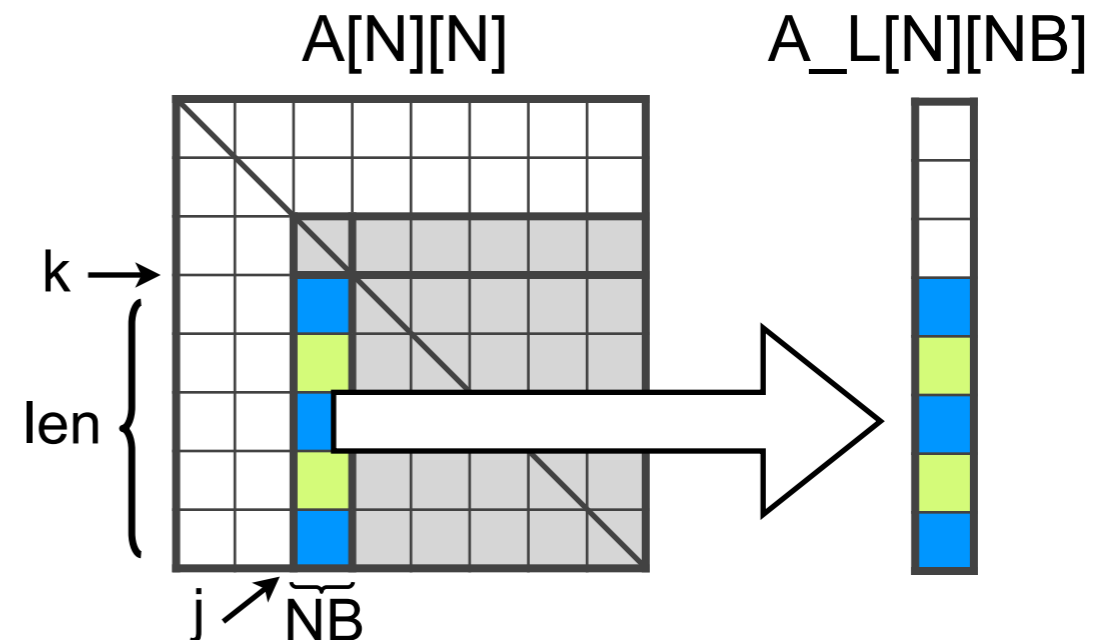
# HPL version 2

- SLOC is **426**, written in XMP/C

- "Lookahead algorithm" by using **gmove** directive with **async** clause

   Overlap communication and calculation

```
double A_L[N][NB];
#pragma xmp align A_L[i][*] with t(*,i)
    :
#pragma xmp gmove async(1)
A_L[k:len][0:NB] = A[k:len][j:NB];
    :
for(m=j+NB;m<N;m+=NB){
  for(n=j+NB;n<N;n+=NB){
    cblas_dgemm(&A[m][n], ..);
    if(xmp_test_async(1)){
      // receive A[k:len][j:NB];
        :
```
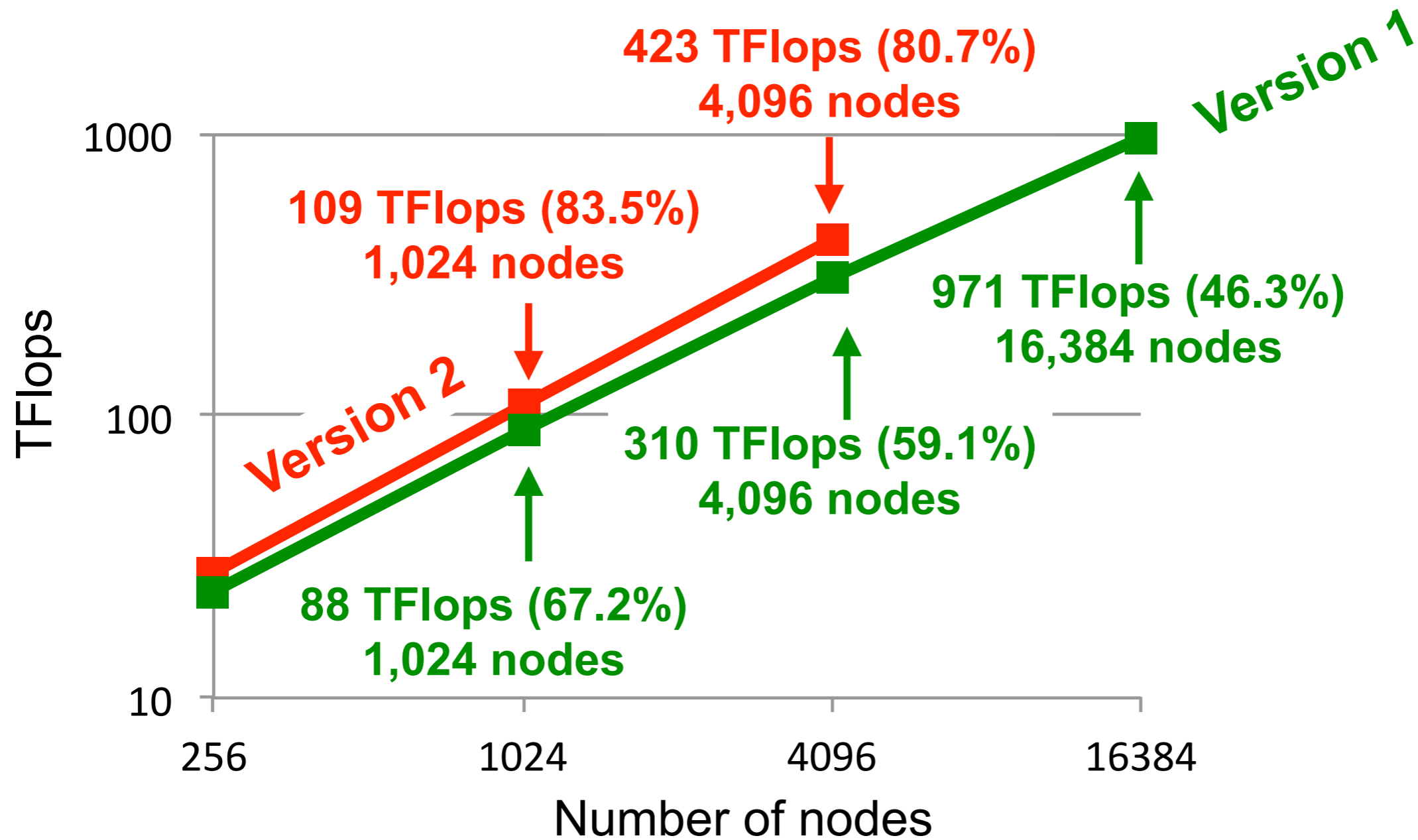
A[N][N]            A_L[N][NB]

k →

len {

j ↗ NB

asynchronous broadcast
communication

**Confirm whether data with async clause comes or not.**

# Performance of HPL



XMP-HPL Version 2 has a good scalability.

Sorry, the measurement in 16,384 nodes is late for this BoF.

# RandomAccess

- SLOC is **253**, written in XMP/C

- Local-view programming with XMP/C coarray syntax

- The XMP RandomAccess is iterated over sets of CHUNK updates on each node

```
u64Int recv[LOGPROCS][RCHUNK+1]:[*];          ←——— Define coarray
...
for (j = 0; j < logNumProcs; j++) {
   recv[j][0:num]:[i_partner] = send[i][0:num];  ←——— Put operation

#pragma xmp sync_memory
#pragma xmp post(p(i_partner), 0)
   :
#pragma xmp wait(p(j_partner))
}
```
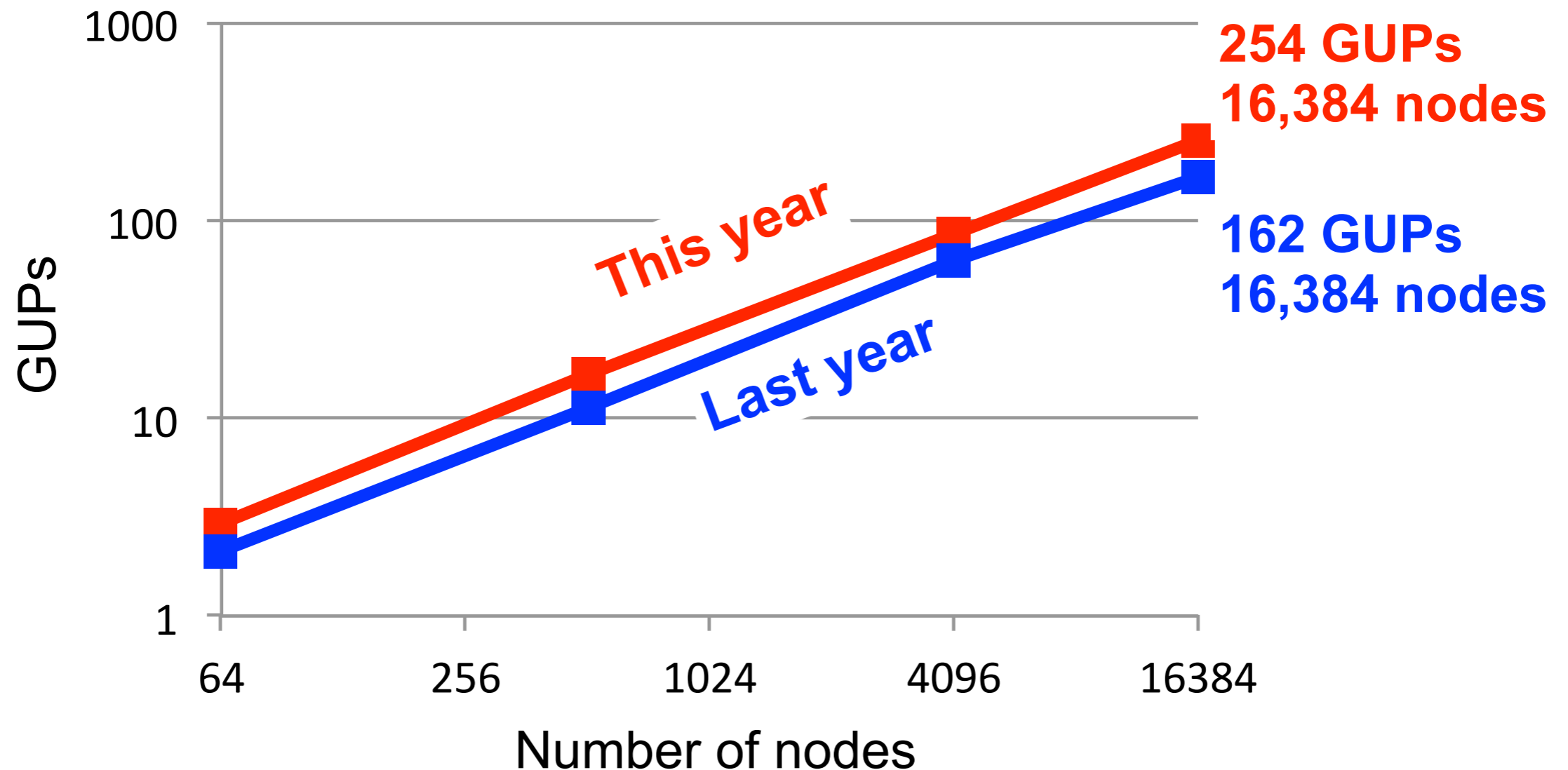
A point-to-point synchronization is specified with the XMP's post and wait directives to realize asynchronous behavior of this algorithm

# Performance of RandomAccess

Last year, to implement the post/wait directives, XMP uses MPI_Send/Recv.
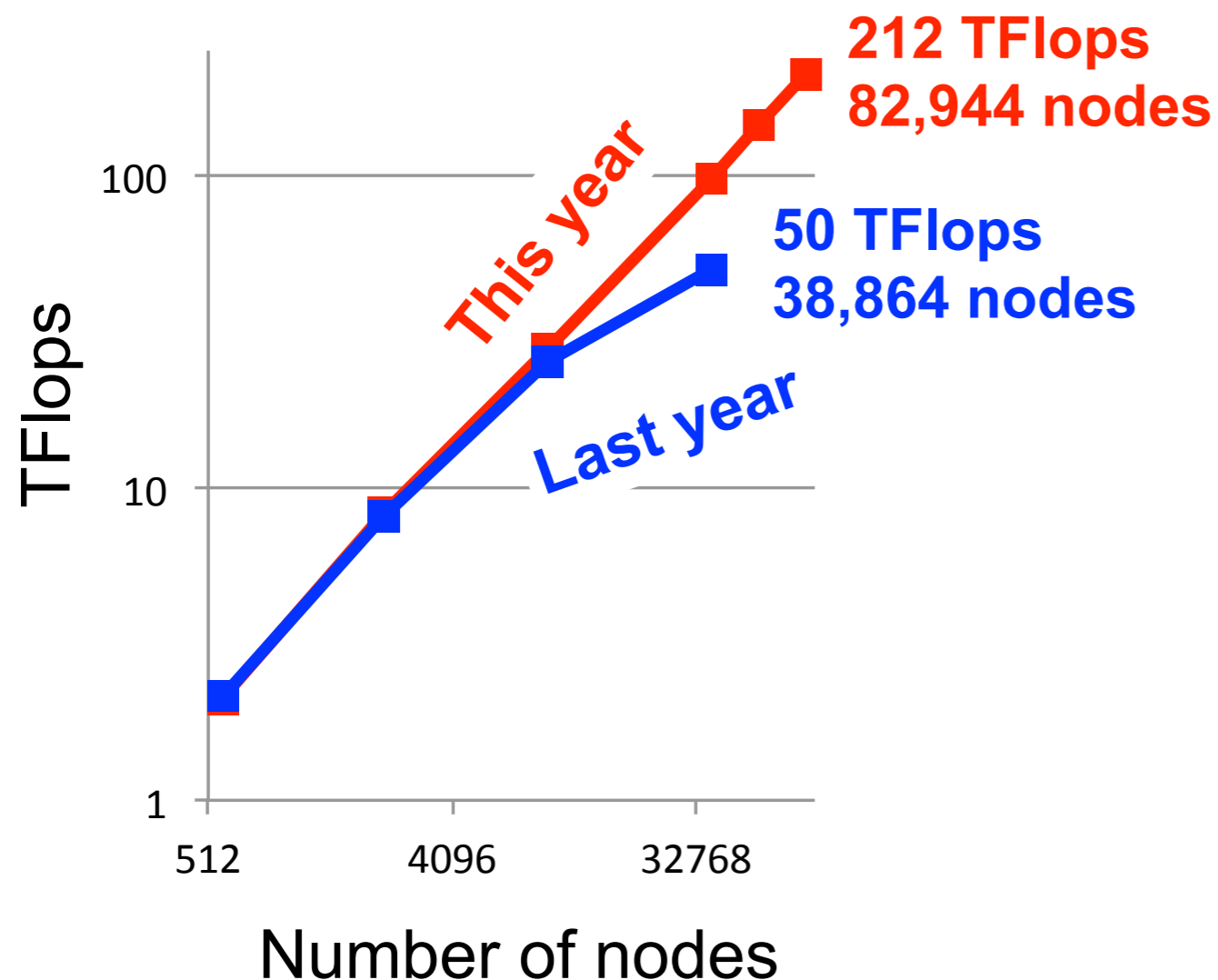This year, to implement them, XMP uses RDMA of the K computer.



**254 GUPs
16,384 nodes**

**162 GUPs
16,384 nodes**

This year

Last year

GUPs

Number of nodes

# FFT and STREAM

Code cleanup and performance improvement.

Please refer to the submission report at *http://xcalablemp.org*

**FFT (SLOC 205, XMP/F)**



**212 TFlops**
**82,944 nodes**

**50 TFlops**
**38,864 nodes**

This year

Last year

TFlops

512    4096    32768

Number of nodes

**STREAM (SLOC 69, XMP/C)**



**3,583 TB/s**
**82,944 nodes**

**2,439 TB/s**
**82,944 nodes**

This year

Last year

TB/s

1024    8192    65536

Number of nodes

# Compare to two versions

- Improvement rate (on the same nodes)  **37 - 94% improvement !!**



Ratio

2.0
1.5
1.0
0.5
0.0

1.37  1.56  1.47  1.94

HPL
**(4,096 nodes)**
RandomAccess
**(16,384 nodes)**
STREAM
**(16,384 nodes)**
FFT
**(36,864 nodes)**

Good

- SLOC



Last year, work-in-progress to clean up code

2500
2000
1500
1000
500
0

313  426    250  253    66  69    2416    205

HPL    RandomAccess    STREAM    FFT

Good

# Outline

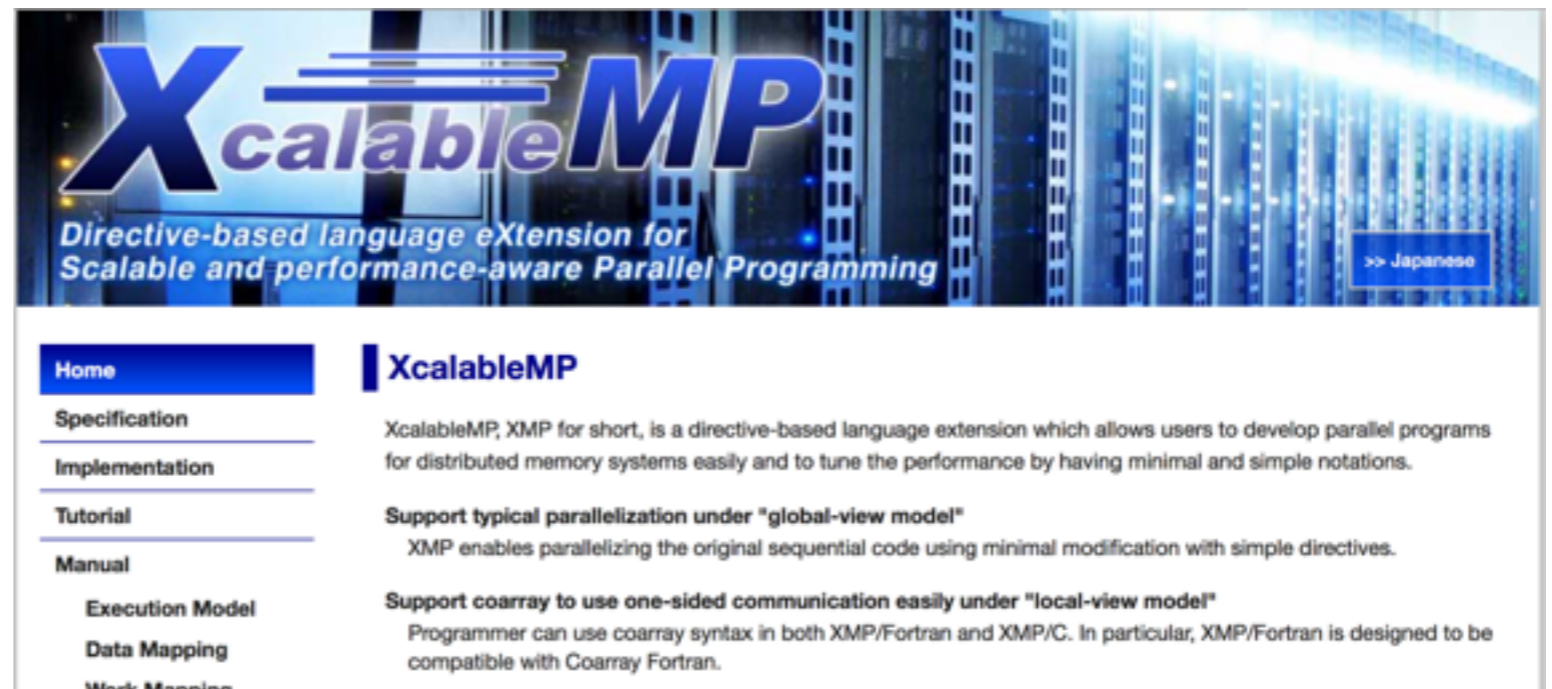1. XcalableMP (XMP) for cluster systems (14min.)

2. XcalableACC (XACC) for accelerator cluster systems (6min.)
   Extension of XMP using OpenACC

Sorry !!, work-in-progress

The submission report is available at *http://xcalablemp.org*

# What is XcalableACC?

**Extension of XMP using OpenACC for accelerator clusters**

**Feature:**

- Mix XMP and OpenACC directives seamlessly

- Support transferring data among accelerators directly

# Difference XMP and XACC memory models

- **XMP** memory model

Global Indexing
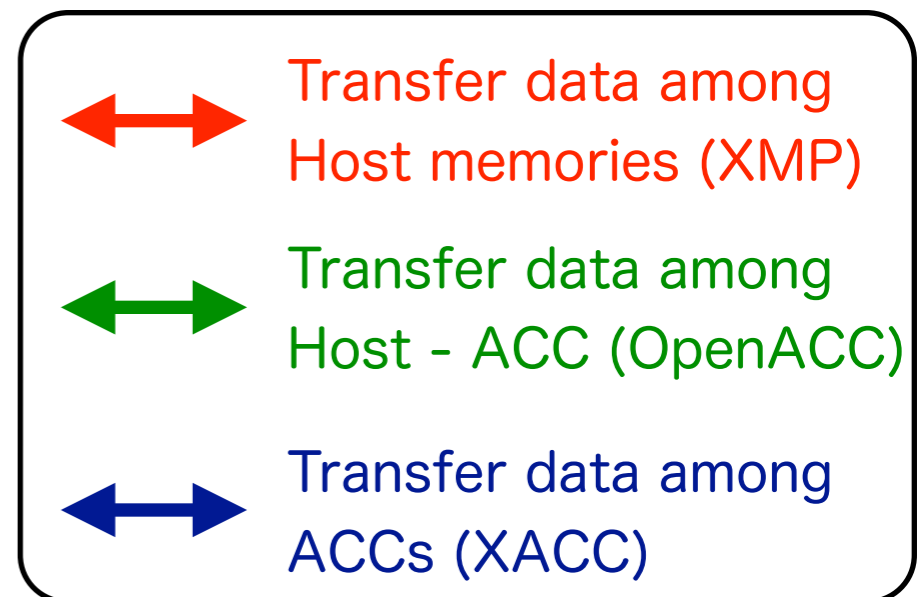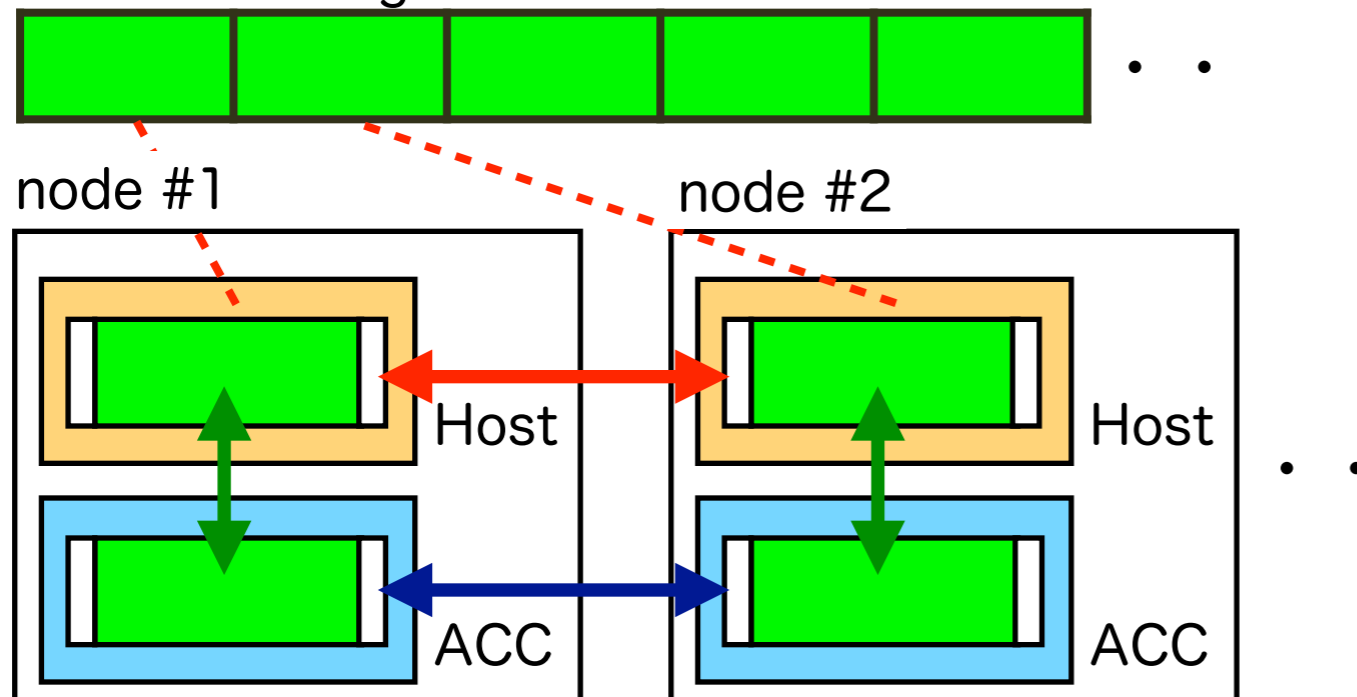


Transfer data among Host memories (XMP)

- **XACC** memory model

Map **"global Indexing"** to accelerators

Global Indexing



Transfer data among Host memories (XMP)

Transfer data among Host - ACC (OpenACC)

Transfer data among ACCs (XACC)

# XACC code example

```
double u[XSIZE][YSIZE], uu[XSIZE][YSIZE];
#pragma xmp nodes p(x, y)
#pragma xmp template t(0:YSIZE-1, 0:XSIZE-1)
#pragma xmp distribute t(block, block) onto p
#pragma xmp align [j][i] with t(i,j) :: u, uu
#pragma xmp shadow uu[1:1][1:1]

…
#pragma acc data copy(u) copyin(uu)
{
  for(k=0; k<MAX_ITER; k++){
#pragma xmp loop (y,x) on t(y,x)
#pragma acc parallel loop collapse(2)
    for(x=1; x<XSIZE-1; x++)
      for(y=1; y<YSIZE-1; y++)
        uu[x][y] = u[x][y];

#pragma xmp reflect (uu) acc

#pragma xmp loop (y,x) on t(y,x)
#pragma acc parallel loop collapse(2)
    for(x=1; x<XSIZE-1; x++)
      for(y=1; y<YSIZE-1; y++)
        u[x][y] = (uu[x-1][y]+uu[x+1][y]+
                   uu[x][y-1]+uu[x][y+1])/4.0;
  } // end k
} // end data
```

**Laplace's equation**

Data Distribution

Transfer XMP distributed arrays to accelerator

OpenACC directive parallelizes the loop statement parallelized by XMP directive

Exchange halo region of uu[][]

When "acc" clause is specified in XMP communication directive, data on accelerator is transferred.

# Results and Machine

## Summary

Three HPCC Benchmarks and HIMENO Benchmark

| Benchmark | #Nodes | #CPUs | #GPUs | Performance (/peak) | SLOC |
|-----------|-------:|------:|------:|:-------------------:|-----:|
| HPL       | 32     | 64    | 128   | 7 TFlops (4.2%)     | 343  |
| FFT       | 32     | 64    | -     | 257 GFlops (0.1%)   | 205  |
| STREAM    | 64     | 128   | 256   | 15 TB/s (20.4%)     | 84   |
| HIMENO    | 64     | 128   | 256   | 14 TFlops (1.4%)    | 253  |

## HA-PACS/TCA: 64 nodes



- Ivy Bridge E5-2680v2, 224GFlops x 2 Sockets
- DDR3 SDRAM 128GB, 59.7GB/s x 2
- Infiniband 4xQDR x 2 rails : 8GB/s
- NVIDIA K20X (4GPUs / Node)
  - 1.31 TFlops/GPU(SP), 3.95 TFlops/GPU(DP)
  - 250GB/s/GPU

*http://www.ccs.tsukuba.ac.jp/CCS/eng/research-activities/projects/ha-pacs*

# STREAM

The XACC STREAM uses both CPUs and GPUs together,
**XMP, OpenACC**, and **OpenMP** directives are used.

```
#pragma xmp nodes p(*)
#pragma acc data copy(a[0:GPU_SIZE], b[0:GPU_SIZE], c[0:GPU_SIZE])
  {
    for(k=0; k<NTIMES; k++) {
#pragma xmp barrier
    times[k] = -xmp_wtime();

#pragma acc parallel loop async
      for (j=0; j<GPU_SIZE; j++)          on GPU
        a[j] = b[j] + scalar*c[j];

#pragma omp parallel for
      for (j=GPU_SIZE; j<MAX_SIZE; j++)    on CPU
        a[j] = b[j] + scalar*c[j];

#pragma acc wait                          Wait until GPU task completes

#pragma xmp barrier
    times[k] += xmp_wtime();
    }
} // acc data
```

# Performance of STREAM



reasonable performance

- 15 TB/s 64 nodes (256GPUs)
- 5 TB/s 64 nodes

XACC (CPU + GPU) SLOC:84

XMP (Only CPU) SLOC: 69

GB/s vs Number of nodes

# HIMENO Benchmark

- Stencil application of incompressible fluid analysis code

- Solving the Poisson's equation

- Sequential and MPI Version HIMENO Benchmark is available at

  *http://accc.riken.jp/2444.htm*

```
float p[MIMAX][MJMAX][MKMAX];
// Define distributed array and halo

#pragma acc data copy(p) ..
{
..
#pragma xmp reflect (p) acc
..
#pragma xmp loop (k,j,i) on t(k,j,i)
#pragma acc parallel loop ..
for(i=1; i<MIMAX; ++i)
  for(j=1; j<MJMAX; ++j){
#pragma acc loop vector ..
    for(k=1; k<MKMAX; ++k){
      S0 = p[i+1][j][k] * ..;
```
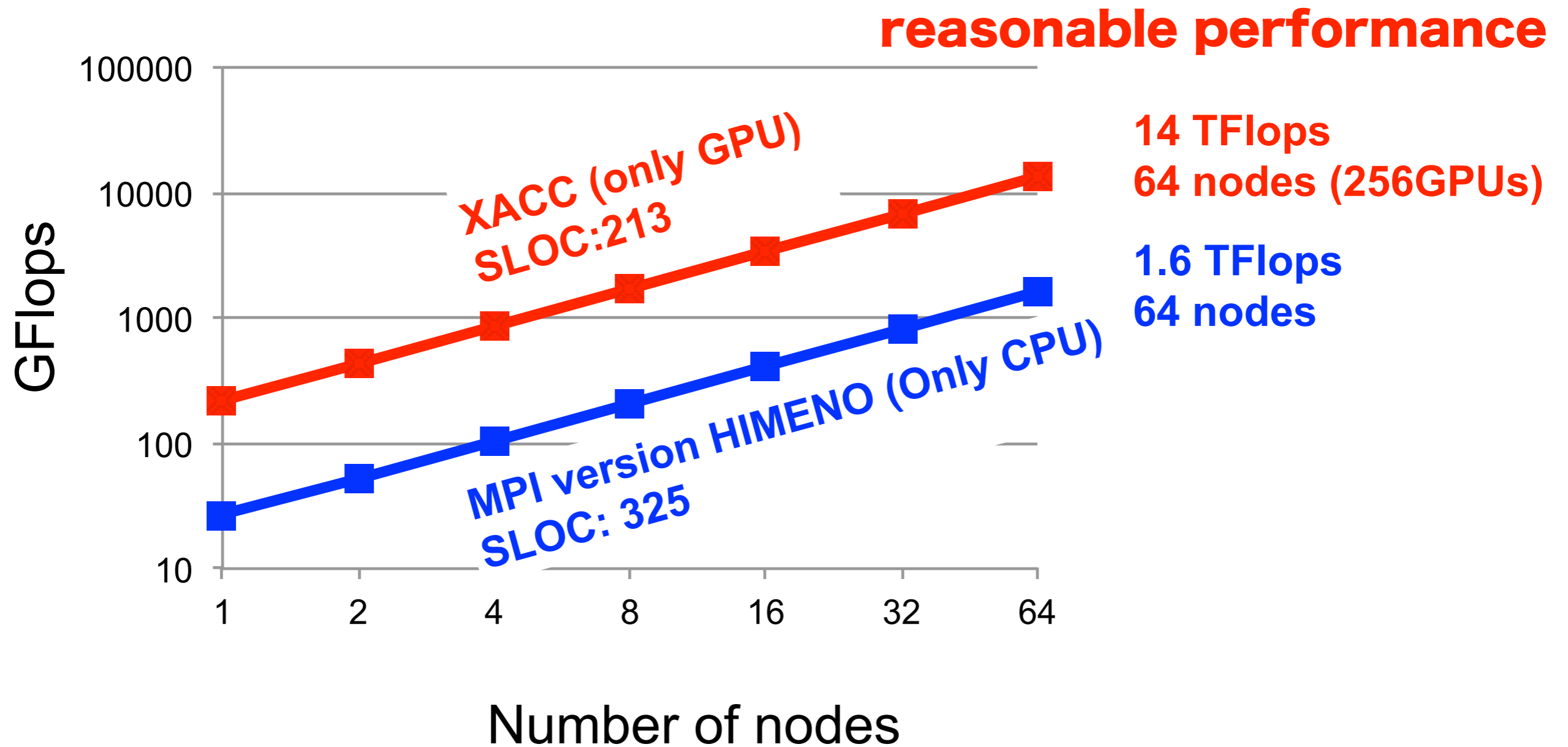
Transfer distributed array to accelerator

Exchange halo region

Parallelize loop statement

Only add XMP and OpenACC directives into the sequential Himeno benchmark.

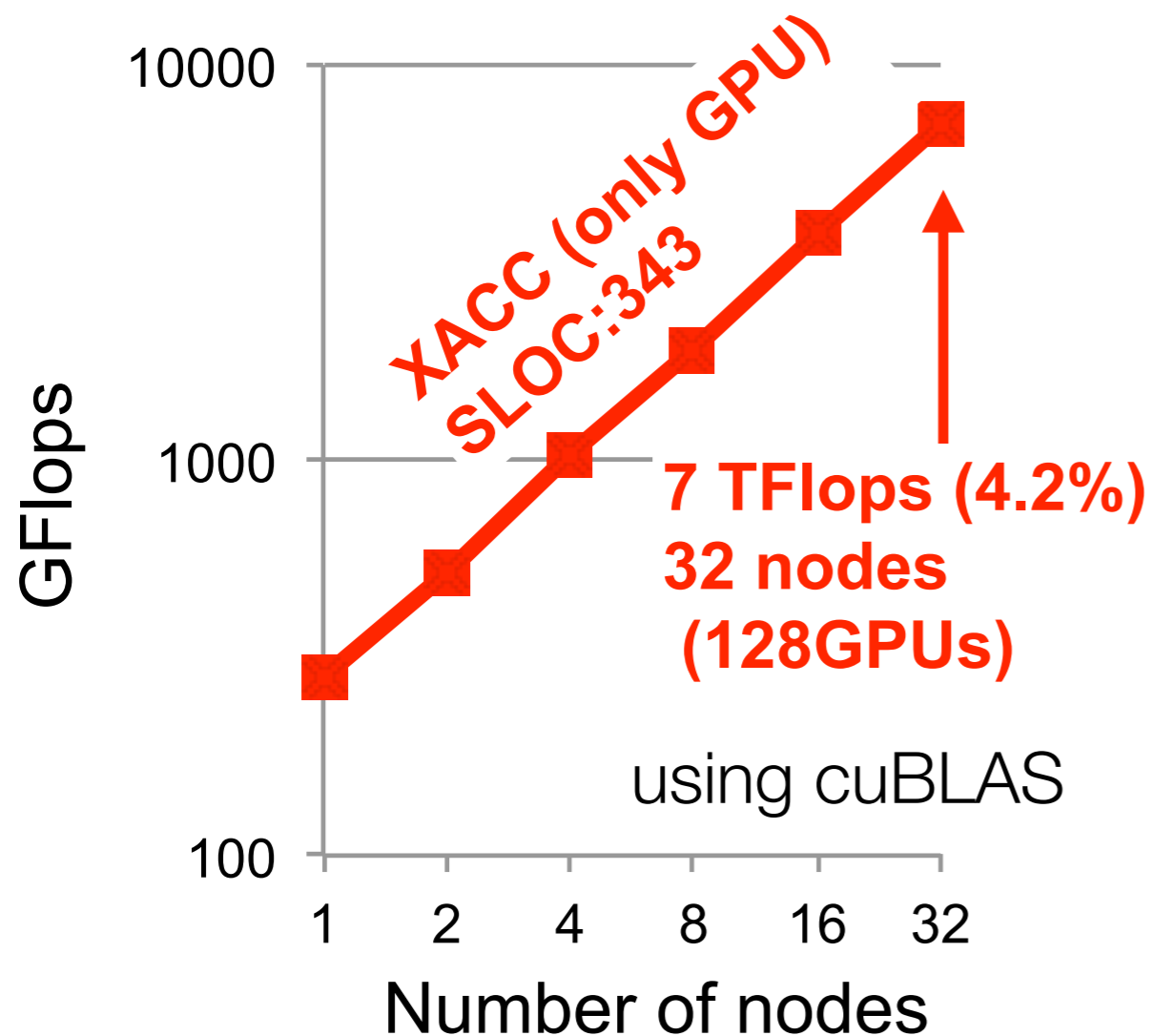# Performance of HIMENO



reasonable performance

XACC (only GPU)
SLOC:213

14 TFlops
64 nodes (256GPUs)

1.6 TFlops
64 nodes

MPI version HIMENO (Only CPU)
SLOC: 325

GFlops

Number of nodes

# HPL and FFT

**Sorry !! work-in-progress for implementing and tuning.**



**HPL**

GFlops vs Number of nodes

XACC (only GPU) SLOC:343

7 TFlops (4.2%)
32 nodes
(128GPUs)

using cuBLAS

**FFT**

GFlops vs Number of nodes

XMP (only CPU) SLOC:205

257 GFlops (0.1%)
32 nodes

will use FFTE-CUDA

Time of transfer data between CPU and host memory dominates the total computation time

# Conclusion

## XMP on the K computer

**Good productivity and performance !!**

| Benchmark | | # Nodes | Performance (/peak) | SLOC |
|---|---|---|---|---|
| HPL | Ver. 1 | 16,384 | 971 TFlops (46.3%) | 313 |
| | Ver. 2 | 4,096 | 423 TFlops (**80.7**%) | 426 |
| FFT | | **82,944** | 212 TFlops (2.0%) | 205 |
| STREAM | | **82,994** | 3,583 TB/s (67.5%) | 69 |
| RandomAccess | | 16,384 | 254 GUPs | 253 |

## XACC on HA-PACS/TCA

**We will improve HPL and FFT next year.**

| Benchmark | #Nodes | #CPUs | #GPUs | Performance (/peak) | SLOC |
|---|---|---|---|---|---|
| HPL | 32 | 64 | 128 | 7 TFlops (4.2%) | 343 |
| FFT | 32 | 64 | - | 257 GFlops (0.1%) | 205 |
| STREAM | 64 | 128 | 256 | 15 TB/s (20.4%) | 84 |
| HIMENO | 64 | 128 | 256 | 14 TFlops (1.4%) | 253 |

# For more information

**Please visit our booth !!**

- RIKEN AICS (Advanced Institute for Computational Science) #2413

- Center for Computational Sciences, University of Tsukuba #3215