



EU-ASEAN HPC School 2022

5-10 December 2022
Kasetsart University, Bangkok, Thailand



*Introduction of Parallel Programming
using Fugaku*

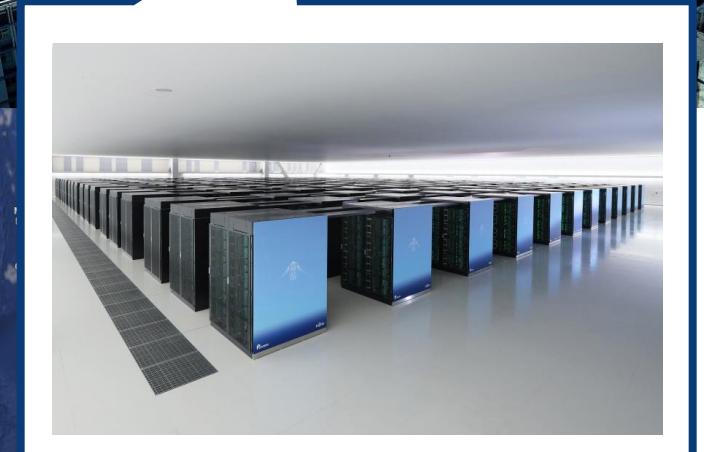
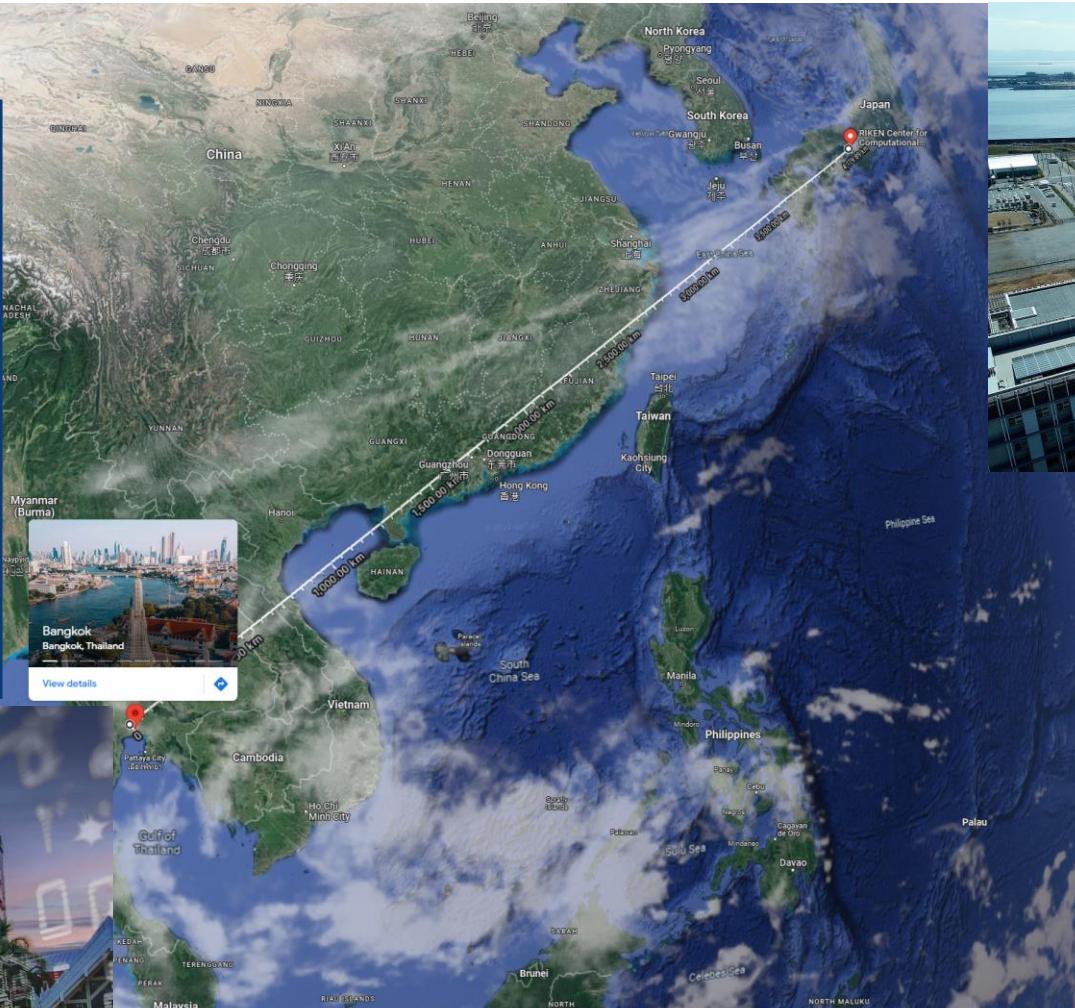
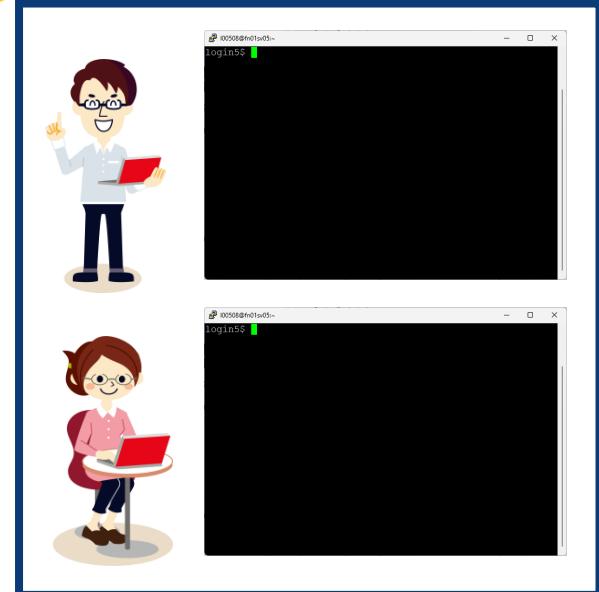
Jorji Nonaka (RIKEN R-CCS)

The EU-ASEAN HPC School is organised by the ASEAN HPC Task Force and carried out in the framework of the Enhanced Regional EU-ASEAN Dialogue Instrument (E-READI)

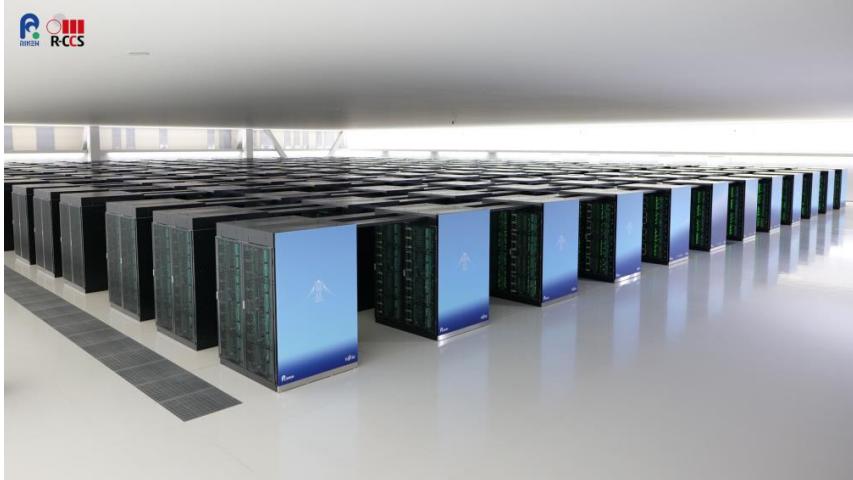
EU-ASEAN HPC School 2022



 RIKEN
Center for
Computational Science



富岳
Fugaku



Users

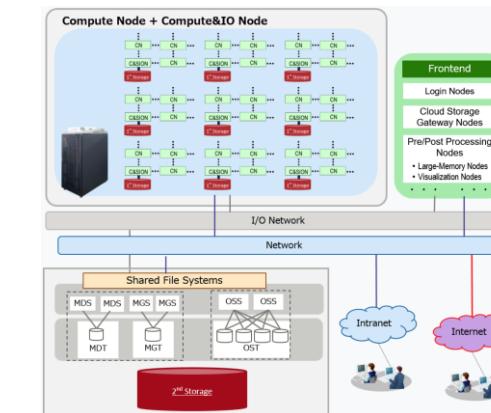
- Scientists
- Engineers
- ...



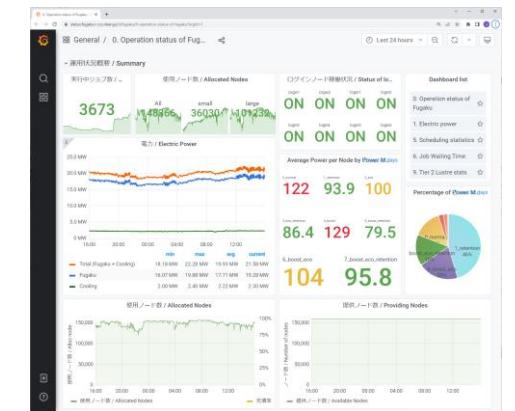
Operations and Computer Technologies Div.



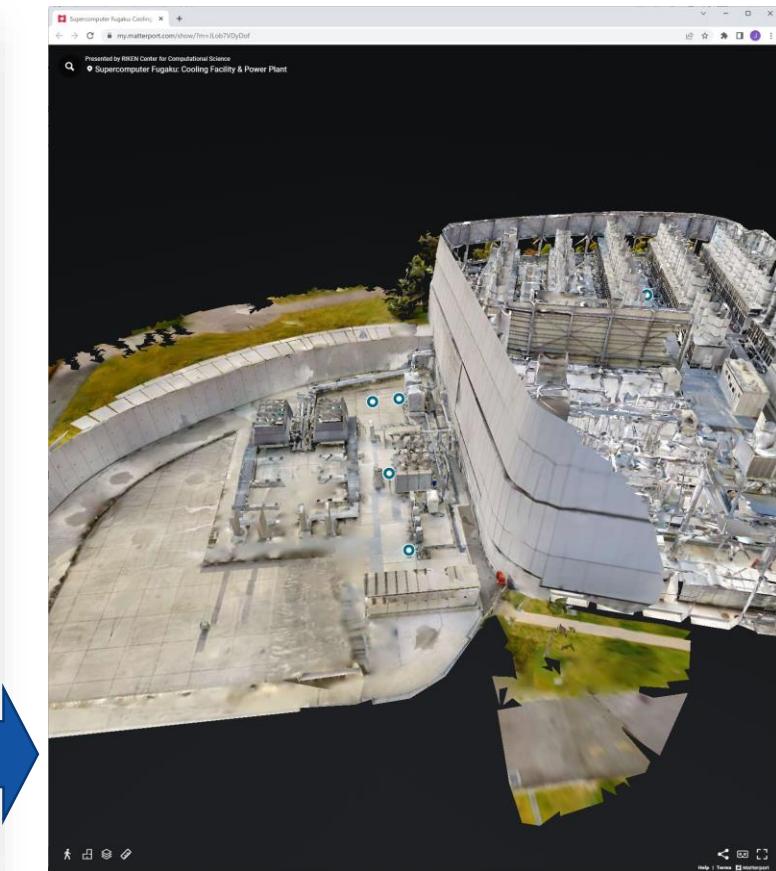
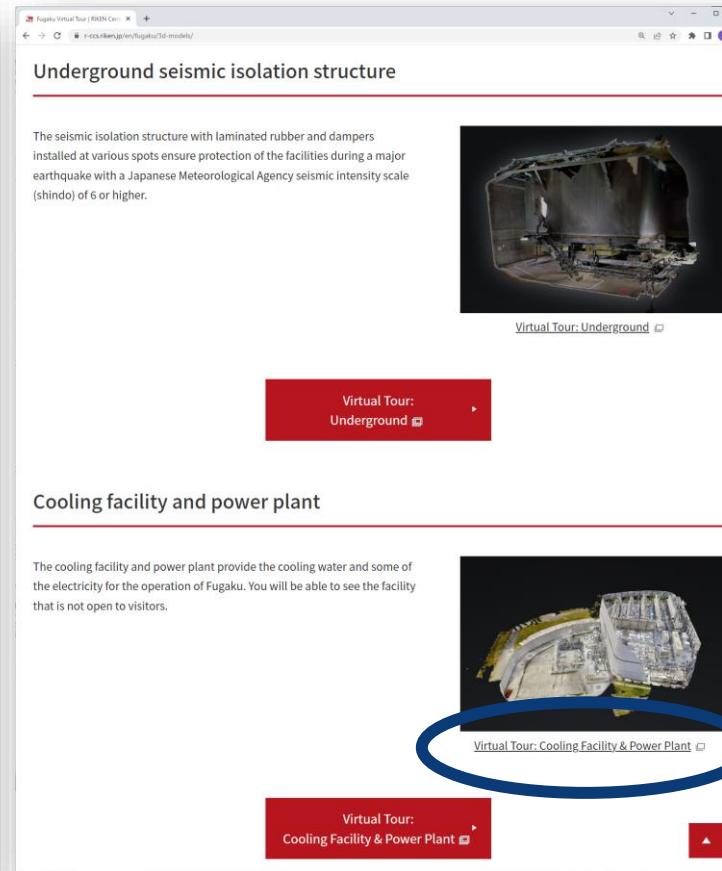
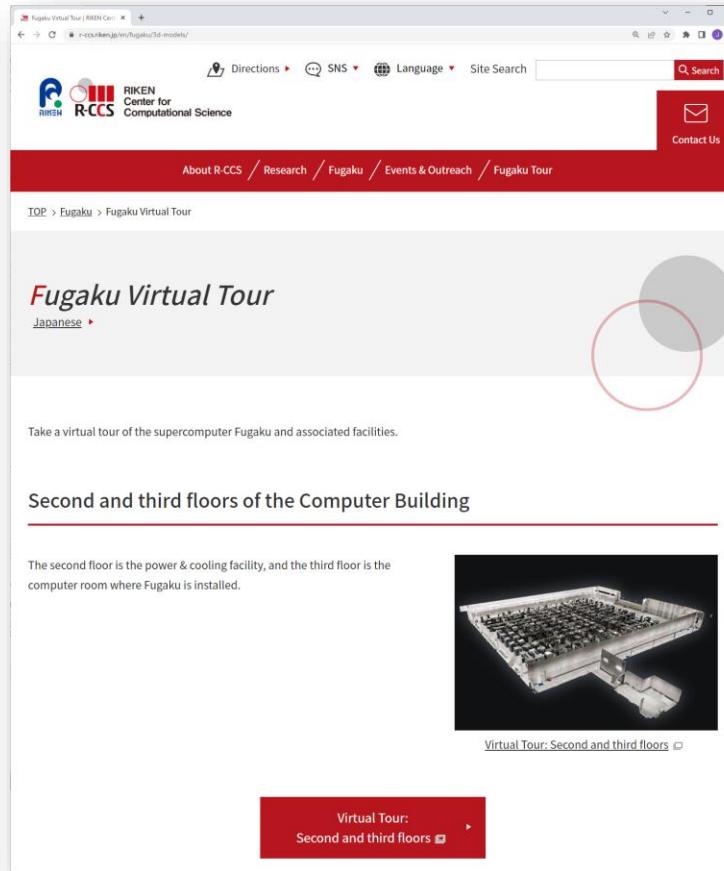
- Power System
 - Substation
 - Power Generator
- Cooling System
 - Air
 - Water



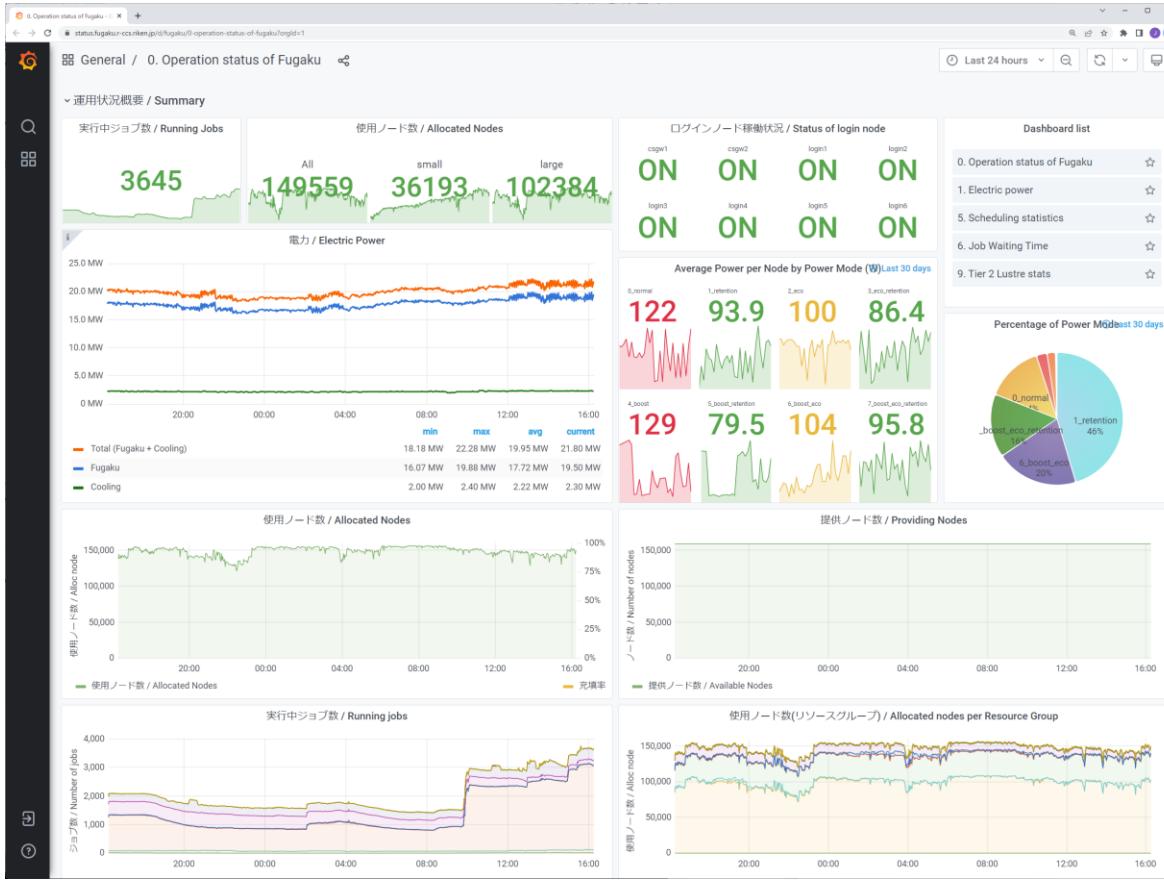
- Management
 - Account
 - Hardware
 - Software
 - Storage
 - ...



<https://www.r-ccs.riken.jp/en/fugaku/3d-models/>



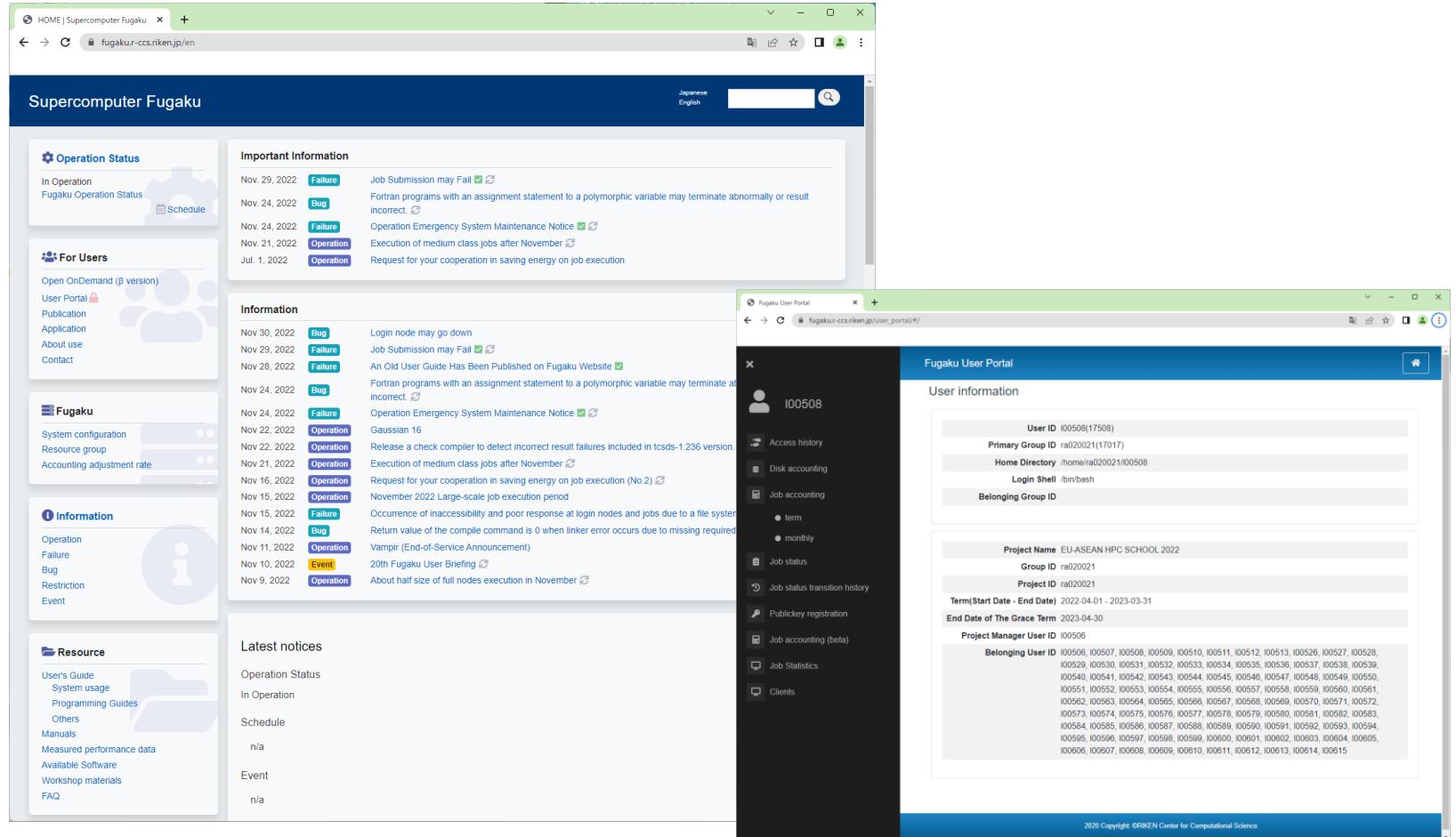
<https://status.fugaku.r-ccs.riken.jp/>



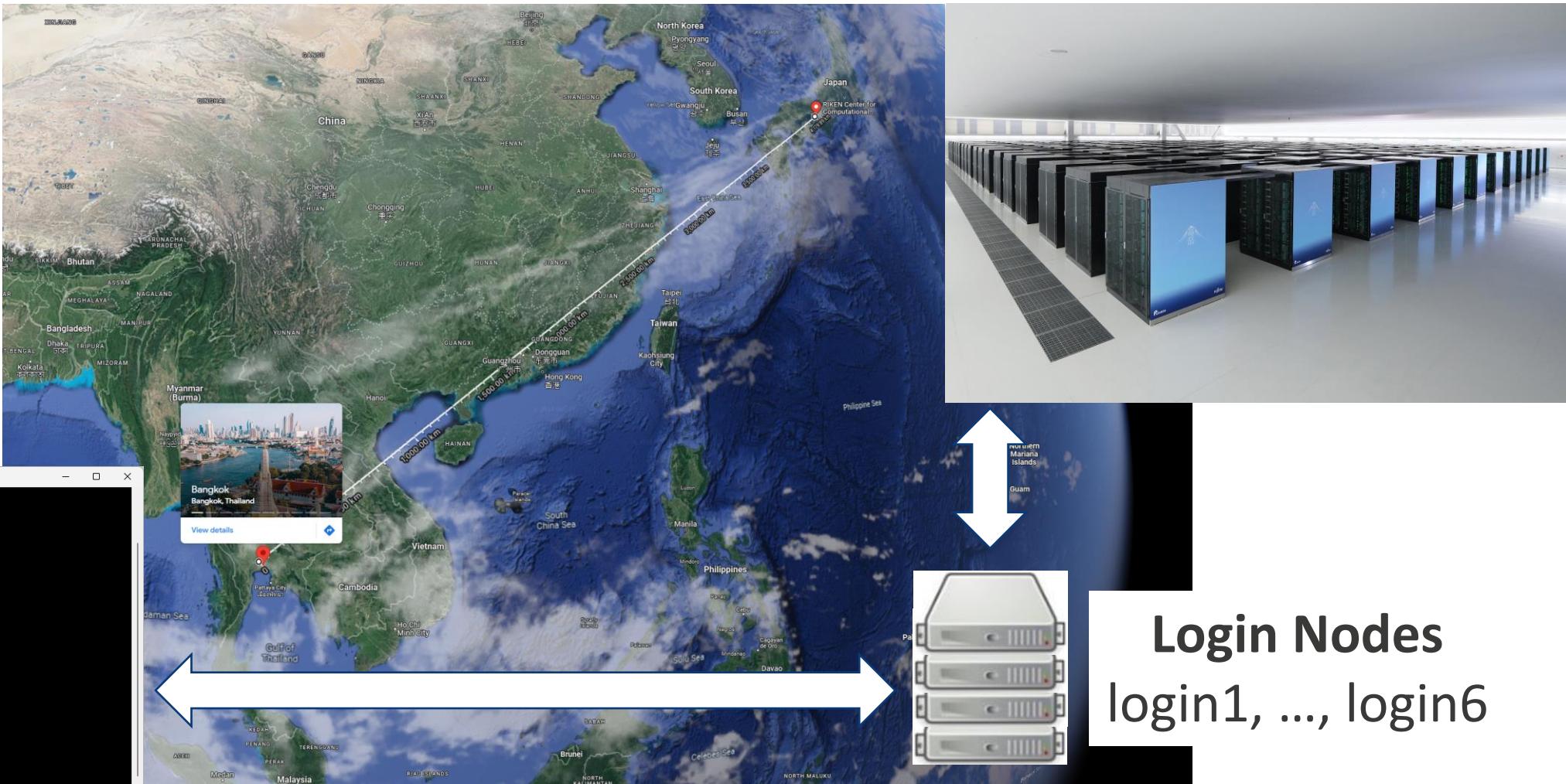
<https://fugaku.r-ccs.riken.jp/>

Fugaku Web Page

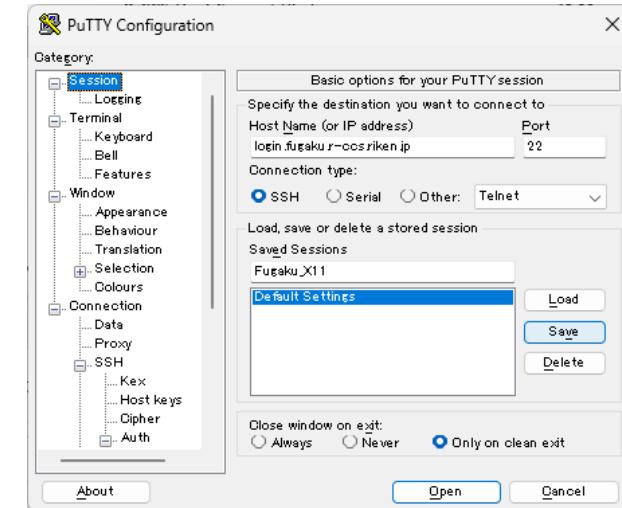
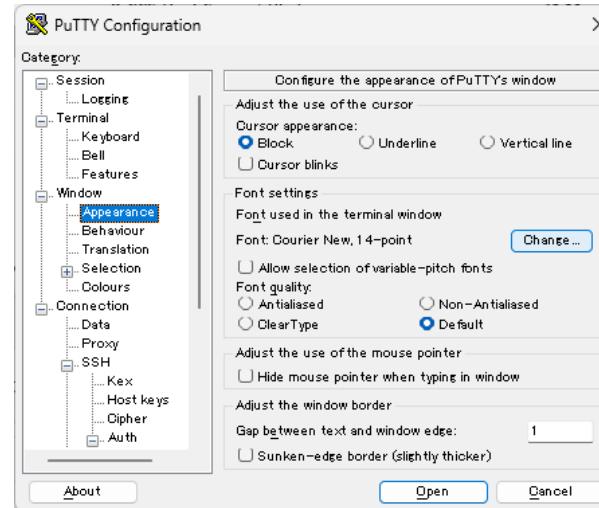
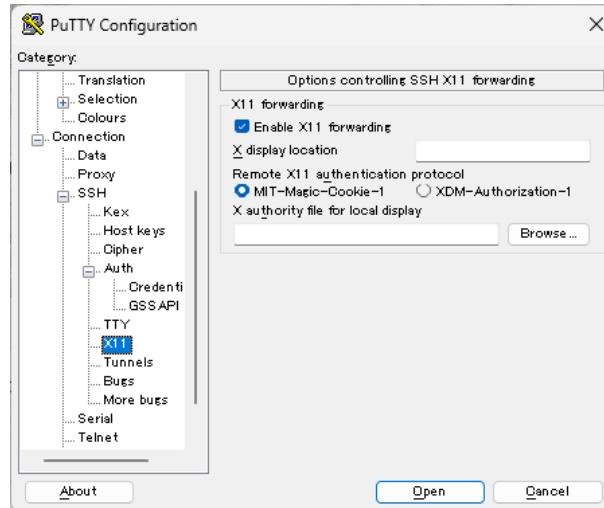
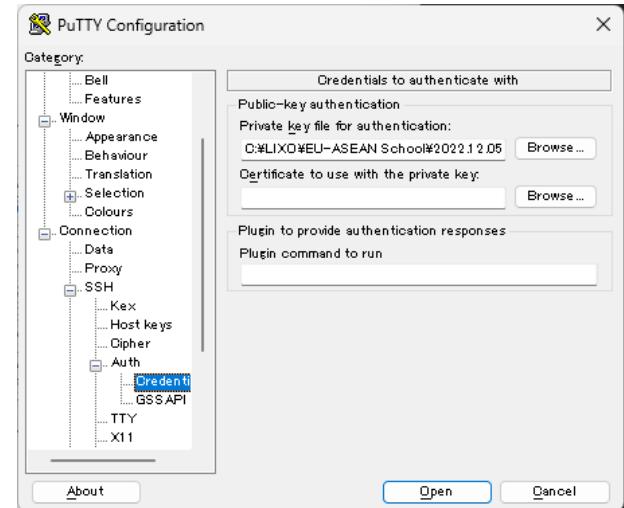
- User Portal
- User's Guide
- Operation Status
- Information
- ...



The screenshot displays two browser windows. The main window shows the 'Supercomputer Fugaku' homepage with sections for 'Operation Status', 'For Users', 'Fugaku', 'Information', and 'Resource'. It includes a sidebar for 'Latest notices' and a 'Important Information' box. The second window shows the 'Fugaku User Portal' for user 'I00508', displaying user information, project details, and job accounting statistics.



PuTTY (SSH Client)



Private Key Setting

- Connection
- SSH
- Credentials

X11 Forwarding

- Connection
- SSH
- Auth
- X11

Appearance (Font settings)

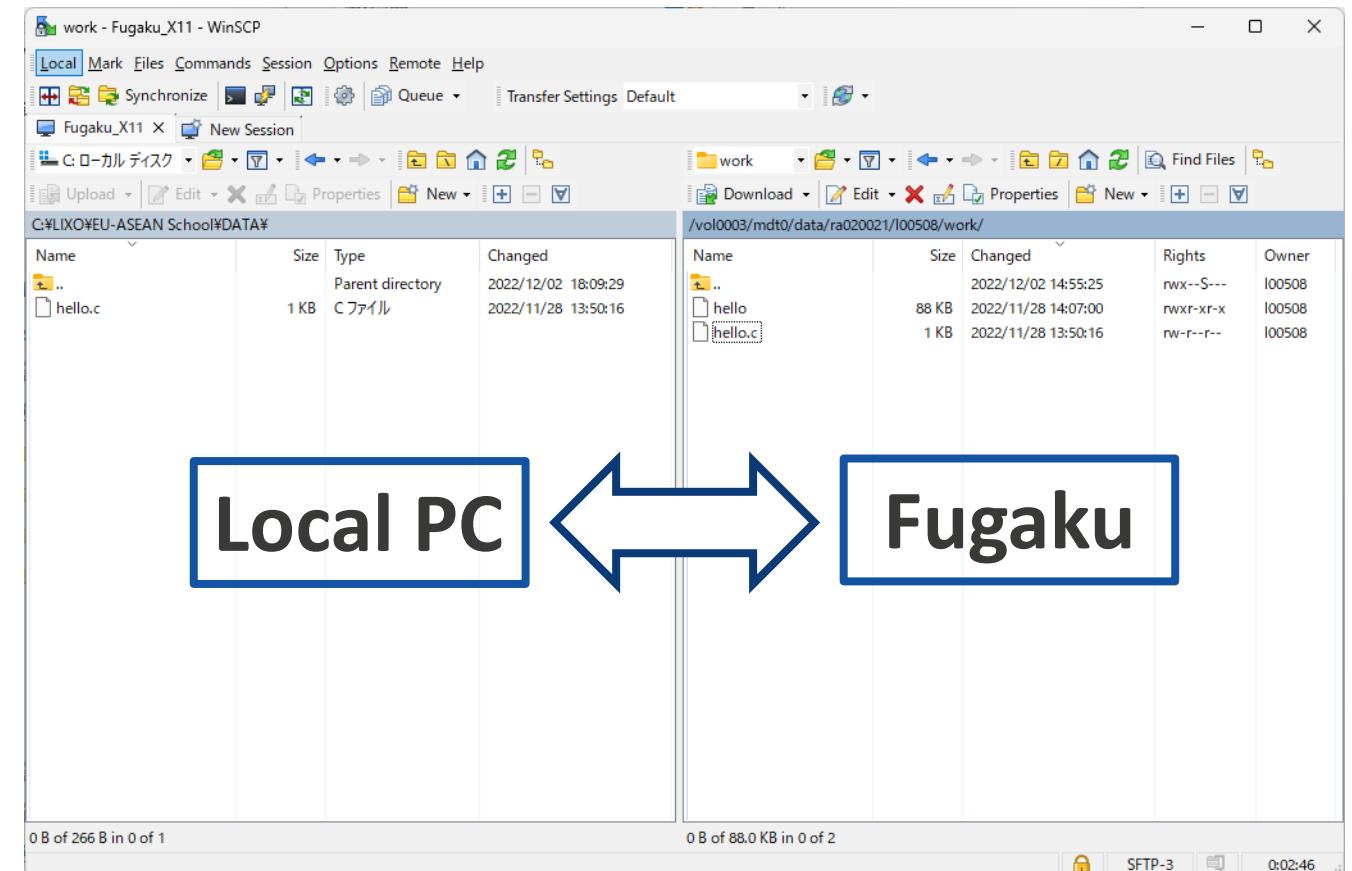
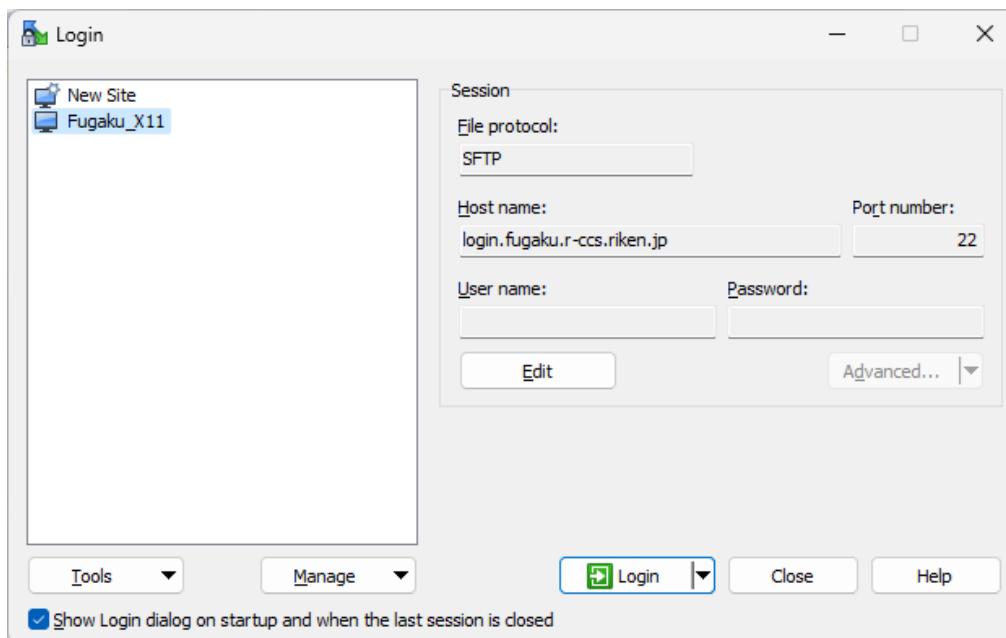
- Window
- Appearance

Save Session (Host Name)

- Session

WinSCP (FTP Client)

Supercomputer Fugaku Startup Guide Version 1.06



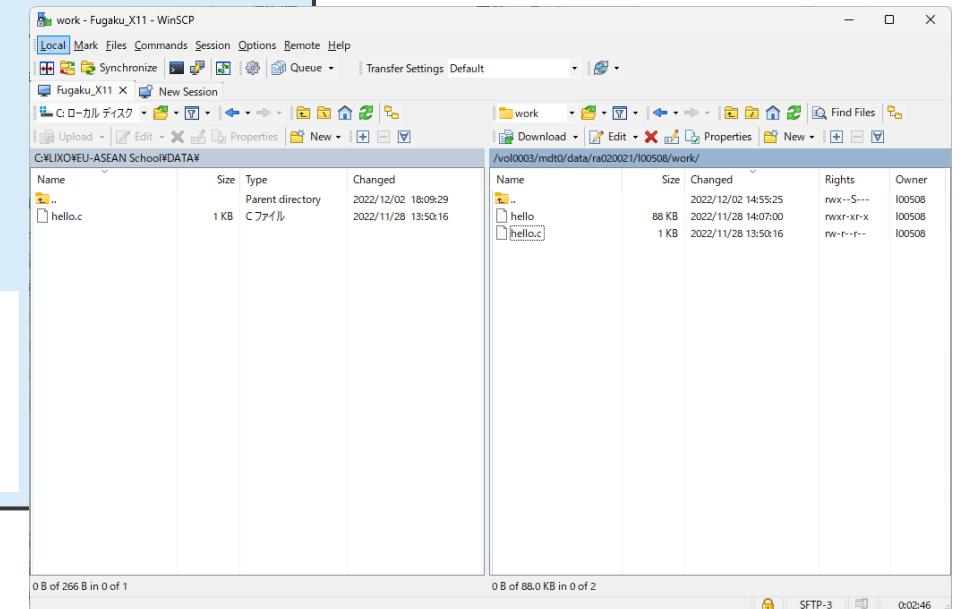
Lecture Slides and Hands-on Data

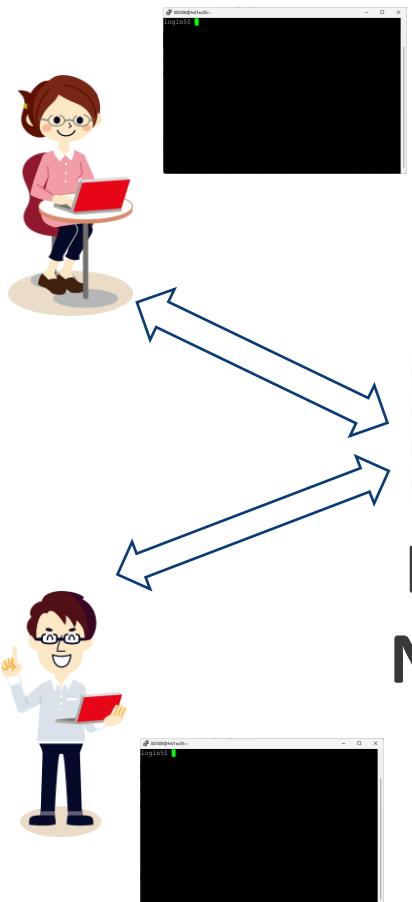
```
$ cd /vol0601/share/ra020021/  
$ cd ComputerScience/20221205_Nonaka  
$ ls  
examples slides
```

Your PC

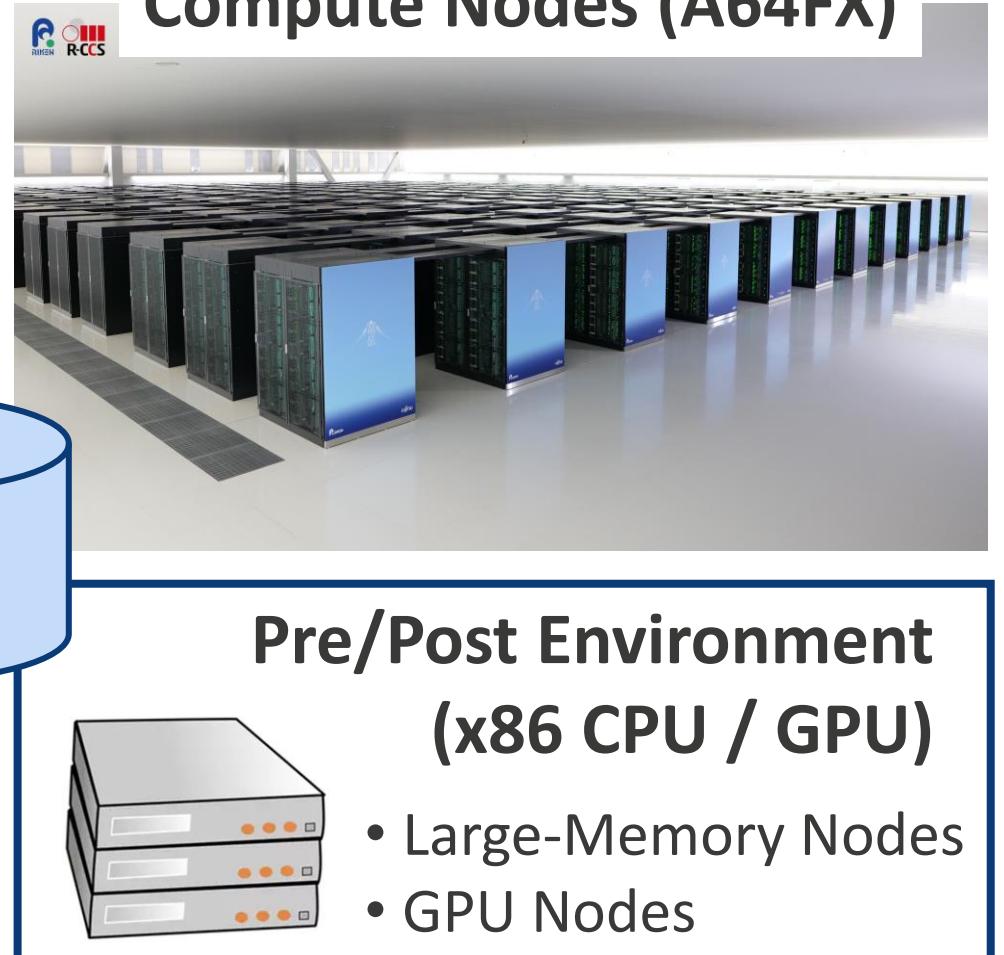
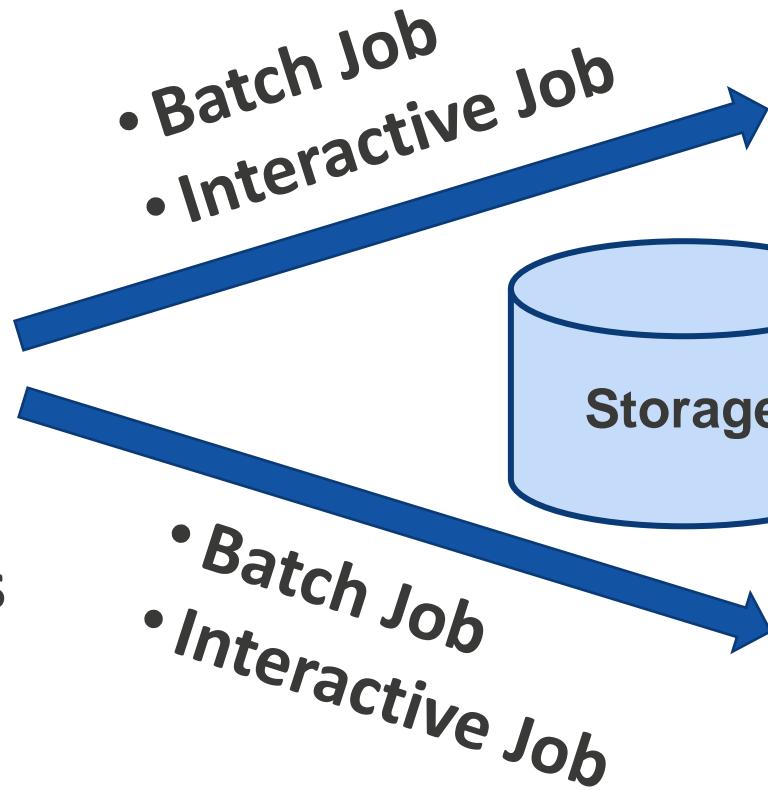
Home Directory
(Fugaku)

```
$ cp -r examples ~/  
$ cp -r slides ~/
```

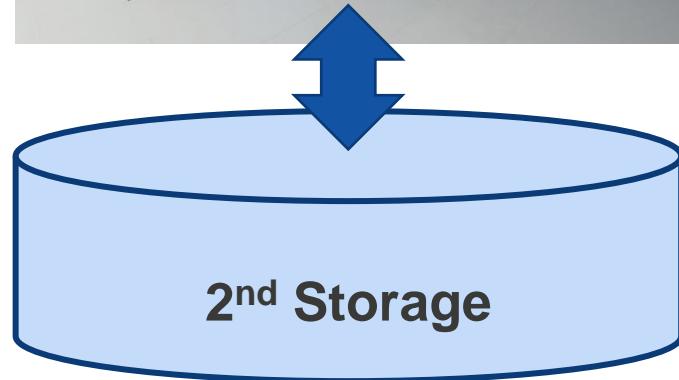




Login
Nodes



Pre/Post Environment Users Guide



\$ accountd -E

```
i00508@fn01sv05:~
login5$ accountd -E
COLLECTDATE : 2022/12/02 13:01:24      unit[GiB]
USER : 100508
*-- [GROUP] --
GROUP      VOLUME          LIMIT      USAGE      AVAILABLE      FILES      USE_RATE
*ra020021  vol0001        5,120       1         5,119        15        0.0%
*ra020021  vol0601        5,120      15         5,105      232,237     0.3%
   /vol0601/data/ra020021
   /vol0601/share/ra020021
*-- [USER] --
USER      VOLUME          LIMIT      USAGE      AVAILABLE      FILES      USE_RATE
100508    vol0300        20          1         19        331        0.1%
   /vol0300/data/ra020021/100508
login5$
```

- Home Directory (20 GB limit)
 - **/vol0003/mdt0/data/ra020021/lxxxxx**
- Data and Share Directory (5TB limit each)
 - **/vol0006/mdt1/data/ra020021**
 - **/vol0006/mdt1/share/ra020021**



Users Guide – Use and job execution

Interactive Job

```
$ pbsub -x PJM_LLIO_GFSCACHE=/vol0004:/vol0006 --interact ...
```

Batch Job

```
$ pbsub batch_job.sh
```



Batch Script: batch_job.sh

```
#PJM -g ra020021
#PJM -L "rscgrp=excl_HPCS_2212-1"
#PJM -x PJM_LLIO_GFSCACHE=/vol0004:/vol0006
...
```

Fugaku Spack User Guide

```
$ ./vol0004/apps/oss/spack/share/spack/setup-env.sh
```

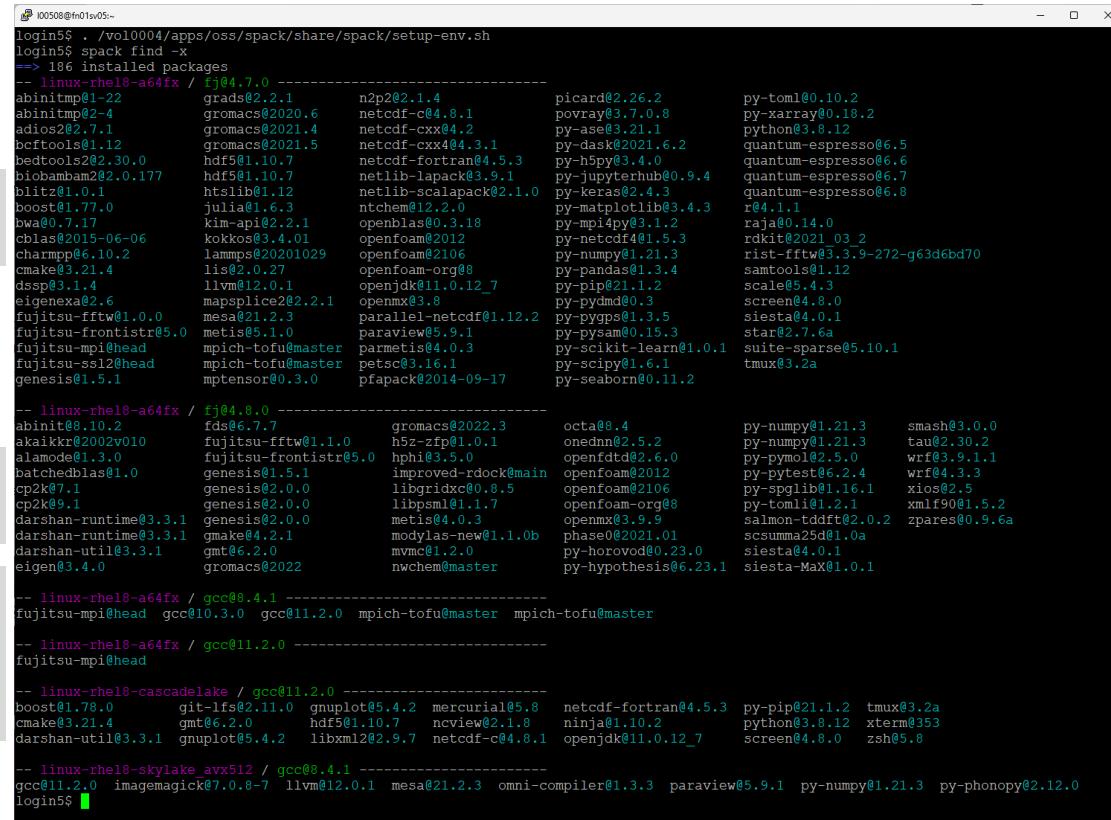
Pre-built packages

```
$ spack find -x
```

Loading packages

\$ spack load NAME

```
$ spack load NAME arch=xxx  
$ spack load /HASH
```



Compute Node

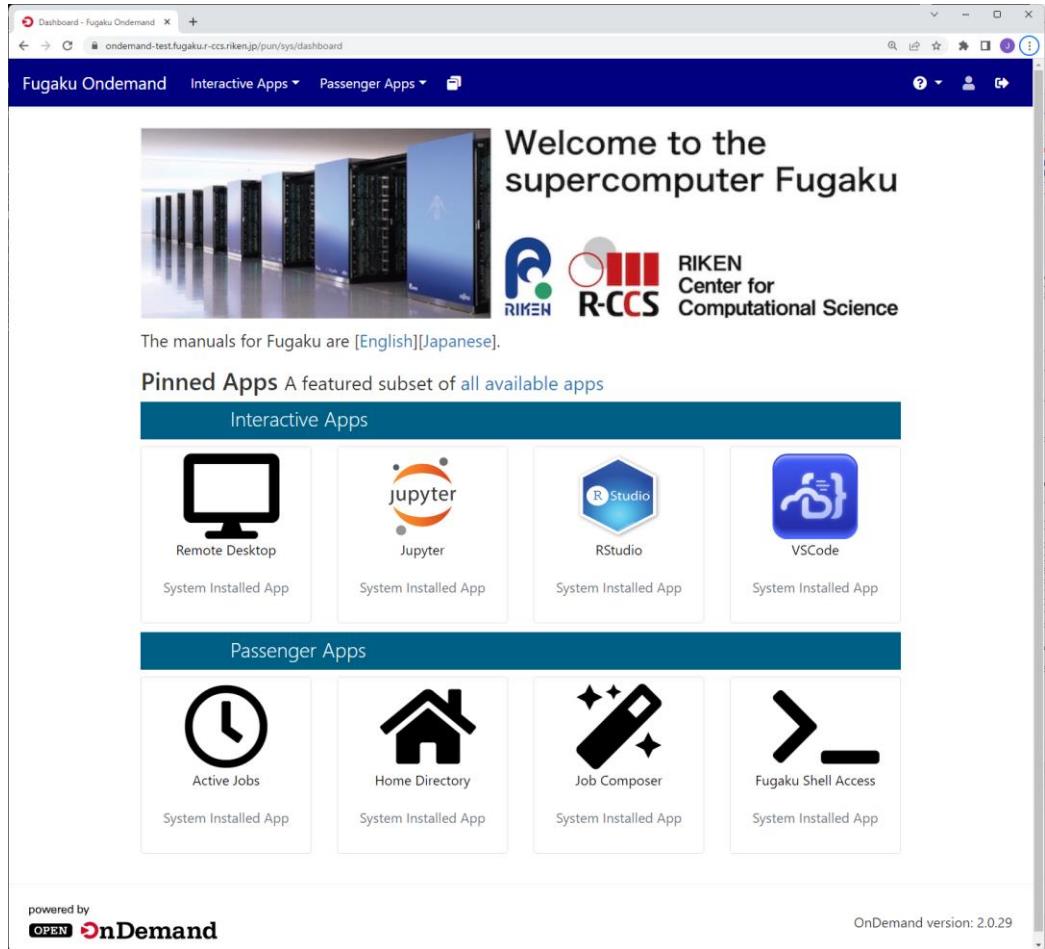
- a64fx

Pre/Post Env.

- cascadelake
 - skylake

Fugaku Open OnDemand Guide

<https://ondemand-test.fugaku.r-ccs.riken.jp>



Welcome to the supercomputer Fugaku

The manuals for Fugaku are [English][Japanese].

Pinned Apps A featured subset of all available apps

Interactive Apps

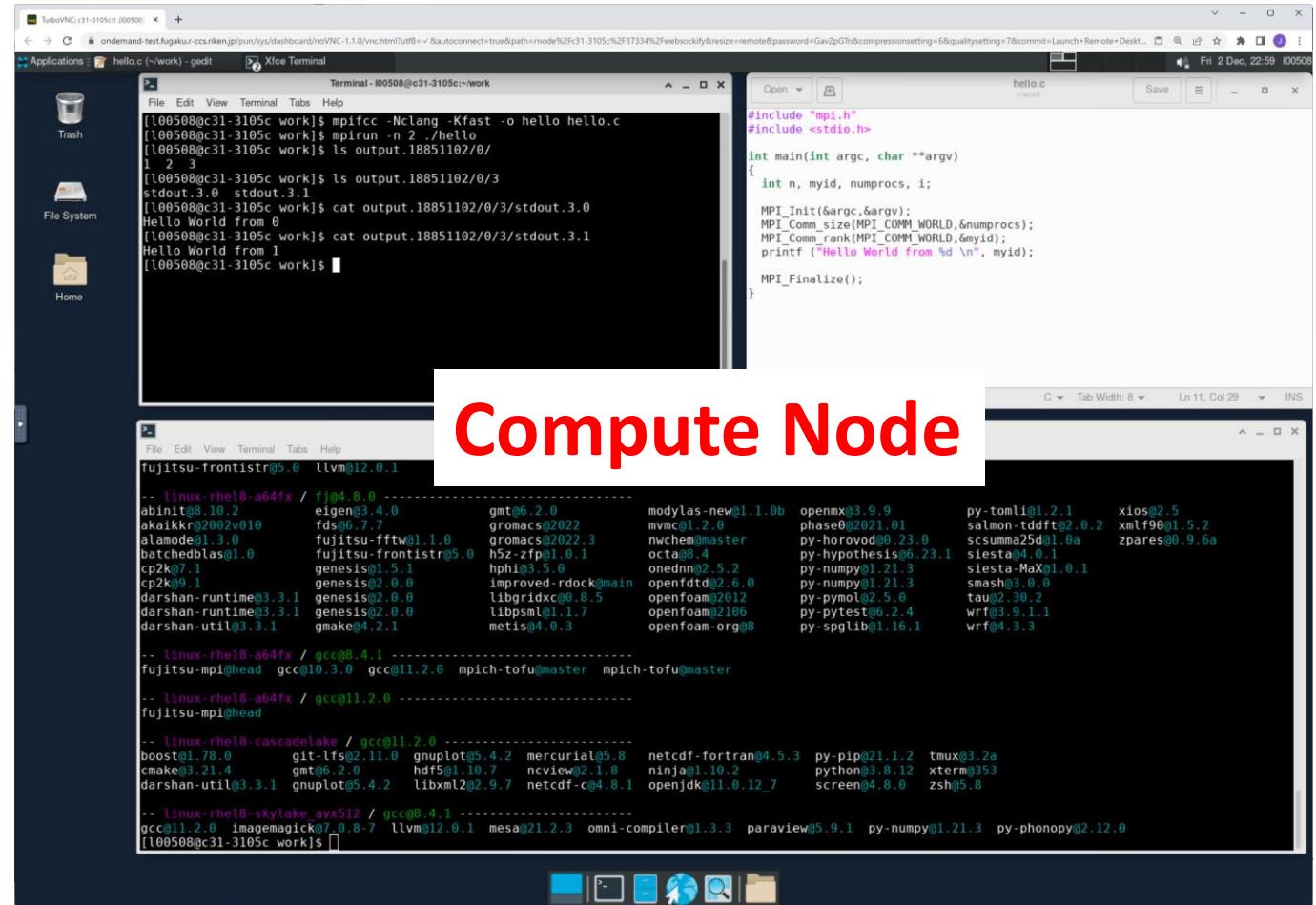
- Remote Desktop
- Jupyter
- RStudio
- VSCode

Passenger Apps

- Active Jobs
- Home Directory
- Job Composer
- Fugaku Shell Access

powered by  OnDemand

OnDemand version: 2.0.29



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

int main(int argc, char **argv)
{
 int n, myid, numprocs, i;

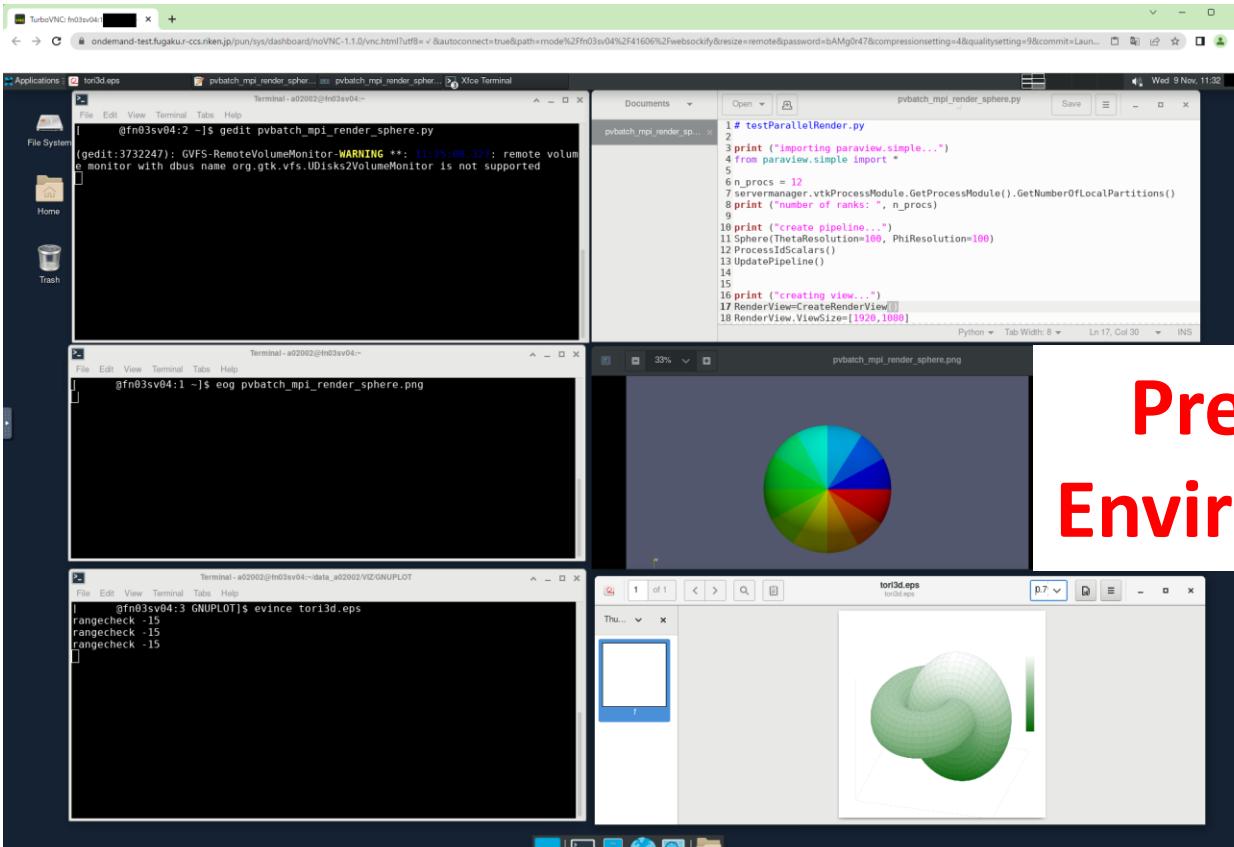
 MPI_Init(&argc,&argv);
 MPI_Comm_size(MPI_COMM_WORLD,&numprocs);
 MPI_Comm_rank(MPI_COMM_WORLD,&myid);
 printf ("Hello World from %d \n", myid);

 MPI_Finalize();
}

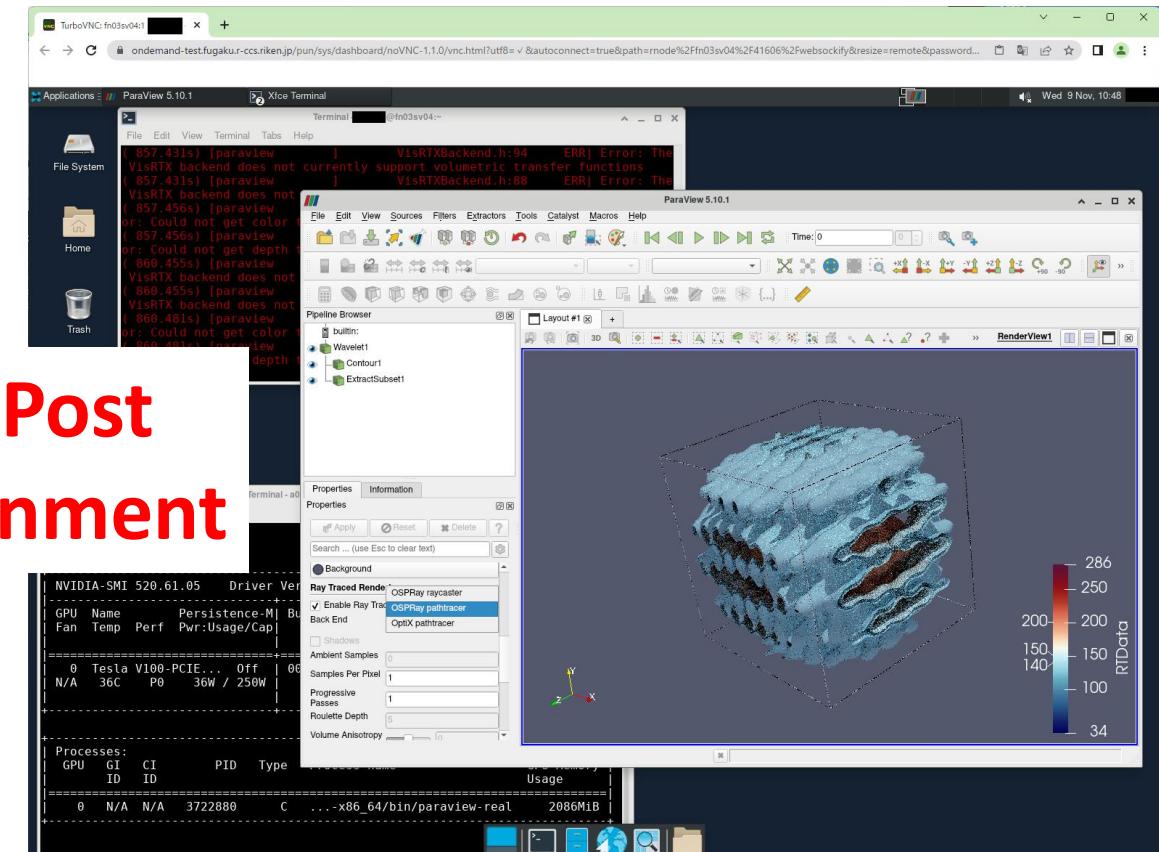
Compute Node

```
-- linux-rhel8-a64fx / fj@4.8.0
abinit@8.10.2 eigen@3.4.0 gmt@0.2.0 modylas-new@1.1.0b openmx@3.9.9 py-tomli@1.2.1 xios@2.5
akalikr@2002v010 fdfs@6.7.7 gromacs@2022 phase@0.2021.01 salmon-tddft@2.0.2 xmff0@1.5.2
alamode@1.3.0 fujitsu-fft@1.1.0 gromacs@2022.3 nwchem@master py-horovod@0.23.0 zcsumma25d@1.0a py-zpages@0.9.6a
batchedblas@1.0 fujitsu-frontistr@5.0 h5z-zfp@1.0.1 octa@8.4 py-hypothesis@6.23.1 siesta@4.0.1
cp2k@7.1 genesis@1.5.1 hphi@3.5.0 onedina@2.5.2 py-numpy@1.21.3 siesta-Max31.0.1
cp2k@9.1 genesis@2.0.0 improved-rdock@main openfdtd@2.6.0 py-numpy@1.21.3 smash@3.0.0
darshan-runtime@3.3.1 genesis@2.0.0 libridx@0.8.5 openfoam@2012 py-pymol@2.5.0 tau@2.30.2
darshan-runtime@3.3.1 genesis@2.0.0 libpsml@1.1.7 openfoam@2106 py-pytest@6.2.4 wrf@3.9.1.1
darshan-util@3.3.1 gmake@4.2.1 metis@4.0.3 openfoam-org@8 py-spglib@1.16.1 wrf@4.3.3
-- linux-rhel8-a64fx / gcc@8.4.1
fujitsu-frontistr@5.0 llvm@12.0.1
-- linux-rhel8-a64fx / gcc@11.2.0
fujitsu-mpi@head gcc@10.3.0 gcc@11.2.0 mpich-tofu@master mpich-tofu@master
-- linux-rhel8-a64fx / gcc@11.2.0
fujitsu-mpi@head
-- linux-rhel8-cascade@lake / gcc@11.2.0
boost@1.78.0 git-lfs@2.11.0 gnuplot@5.4.2 mercurial@5.8 netcdf-fortran@4.5.3 py-pip@21.1.2 tmux@3.2a
cmake@3.1.4 gmt@6.2.0 hdf5@1.10.7 ncview@2.1.8 ninja@1.10.2 python@3.8.12 xterm@353
darshan-util@3.3.1 gnuplot@5.4.2 libxml2@2.9.7 netcdf-c@4.8.1 openjdk@11.0.12_7 screen@4.8.0 zsh@5.8
-- linux-rhel8-skyline_avx512 / gcc@8.4.1
gcc@11.2.0 imagemagick@7.0.8-7 llvm@12.0.1 mesa@21.2.3 omni-compiler@1.3.3 paraview@5.9.1 py-numpy@1.21.3 py-phonopy@2.12.0
```

Visualization User Guide

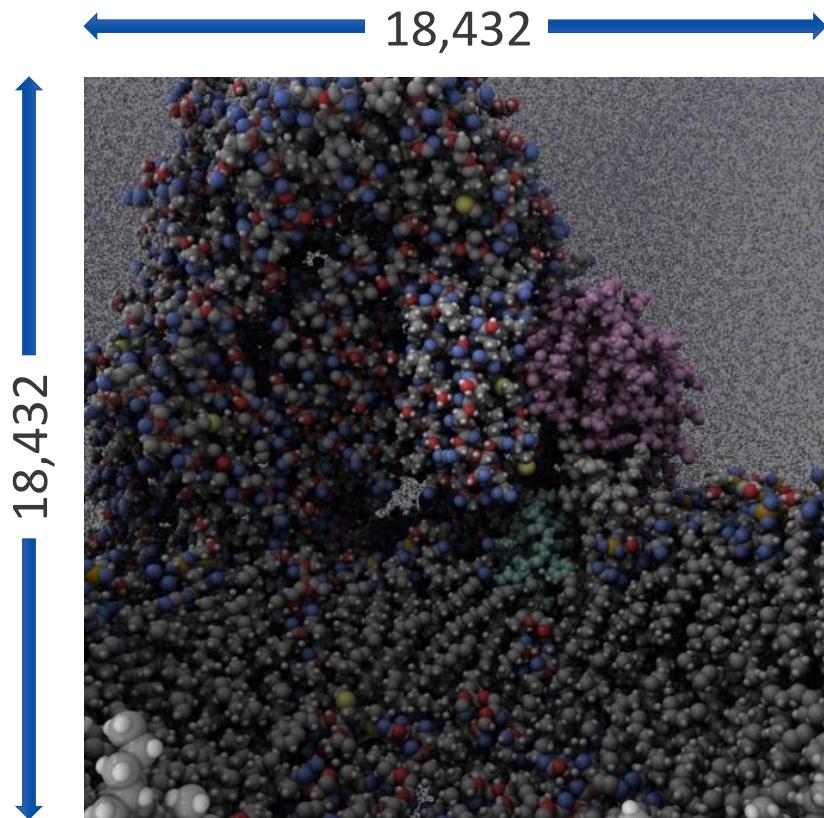


Pre/Post Environment



Parallel Processing

Divide-and-conquer

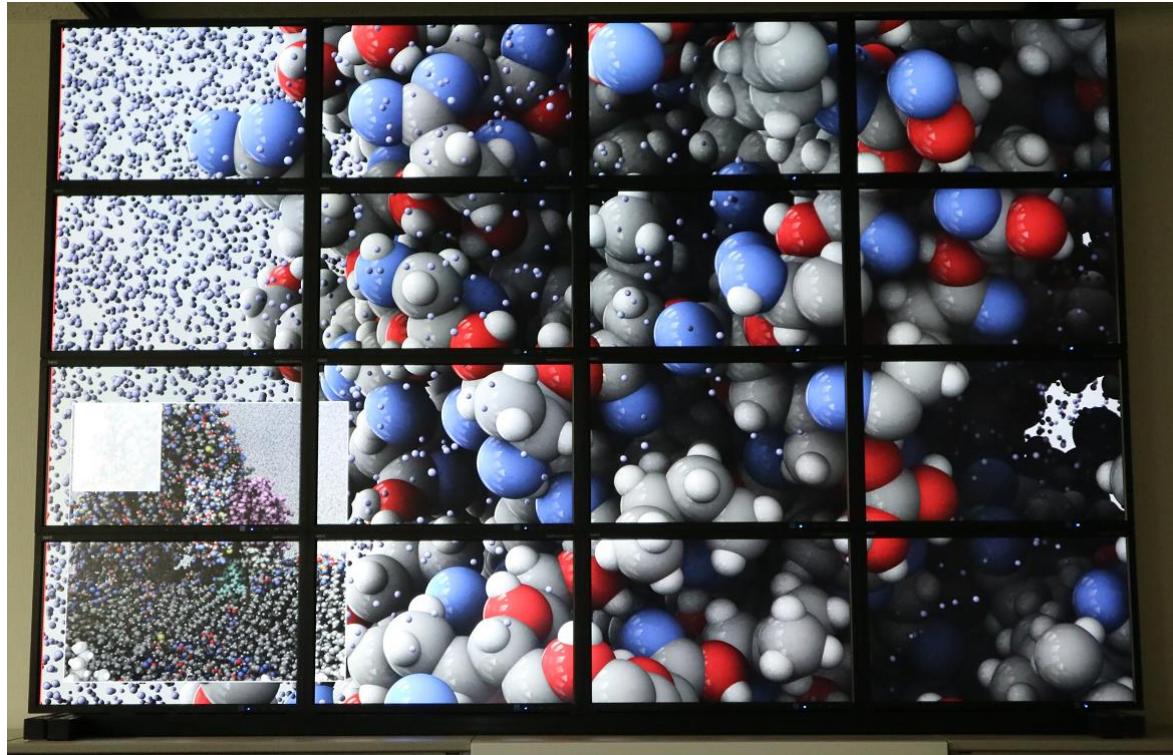


18,432



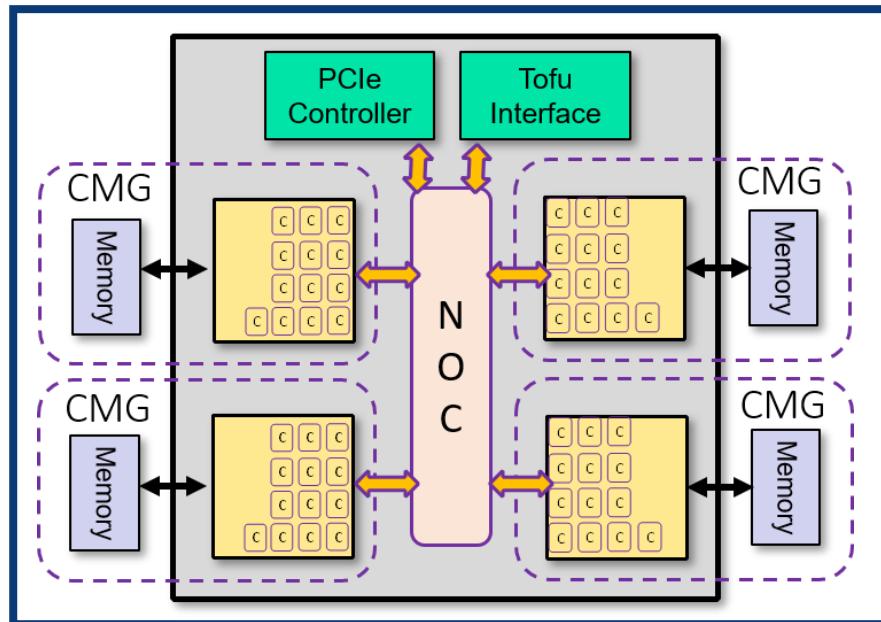
82,944 Nodes

- 288 x 288 (2D subdivision)
- 64 x 64 pixels

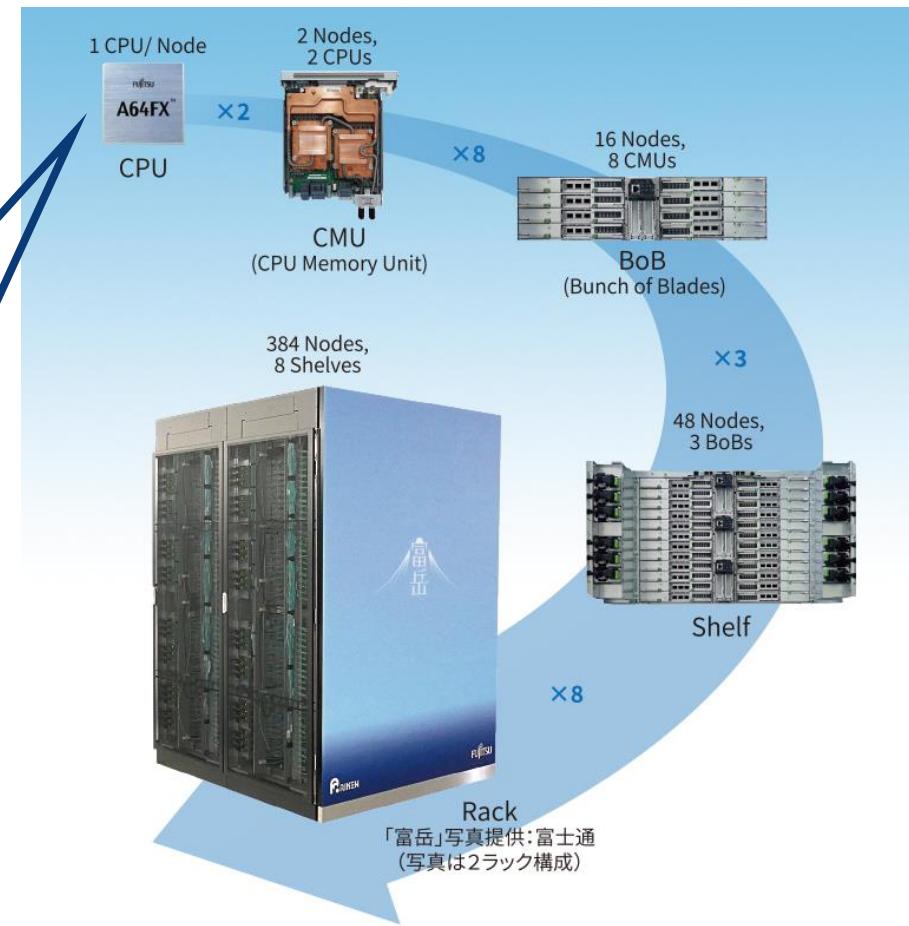


Distributed Memory

Shared Memory



48 Computer Cores per CPU



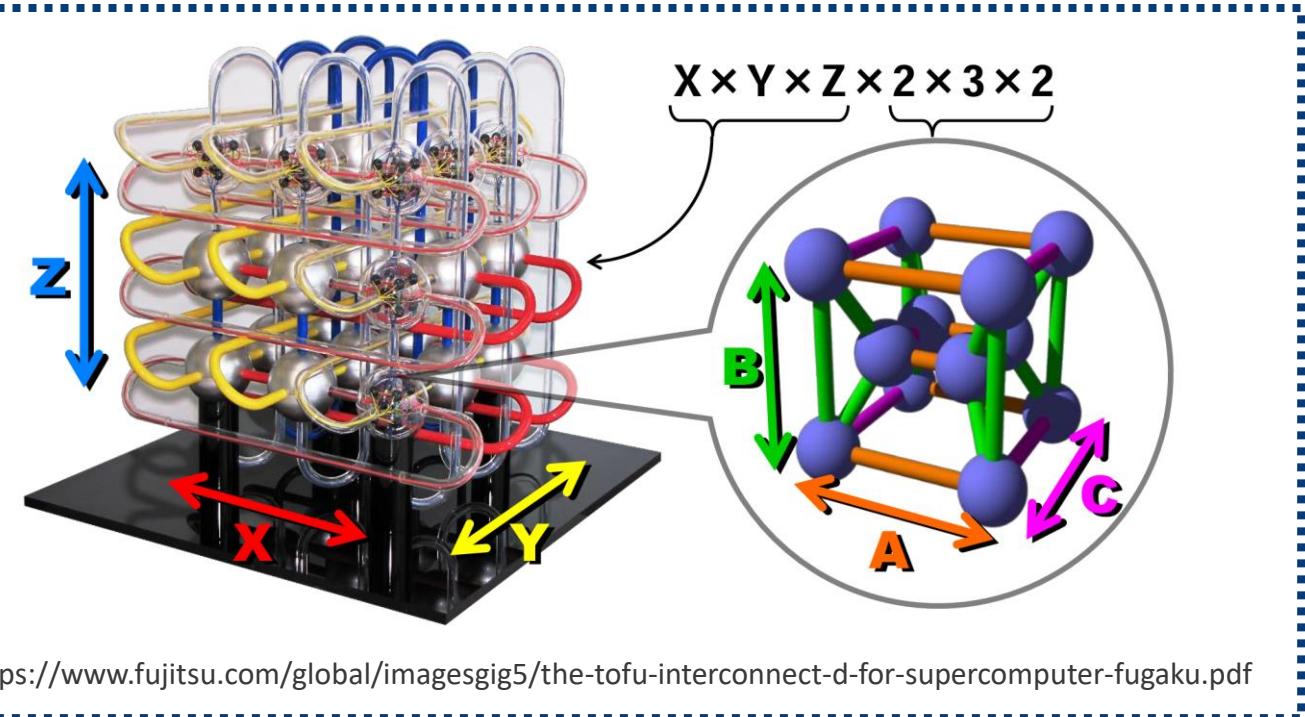
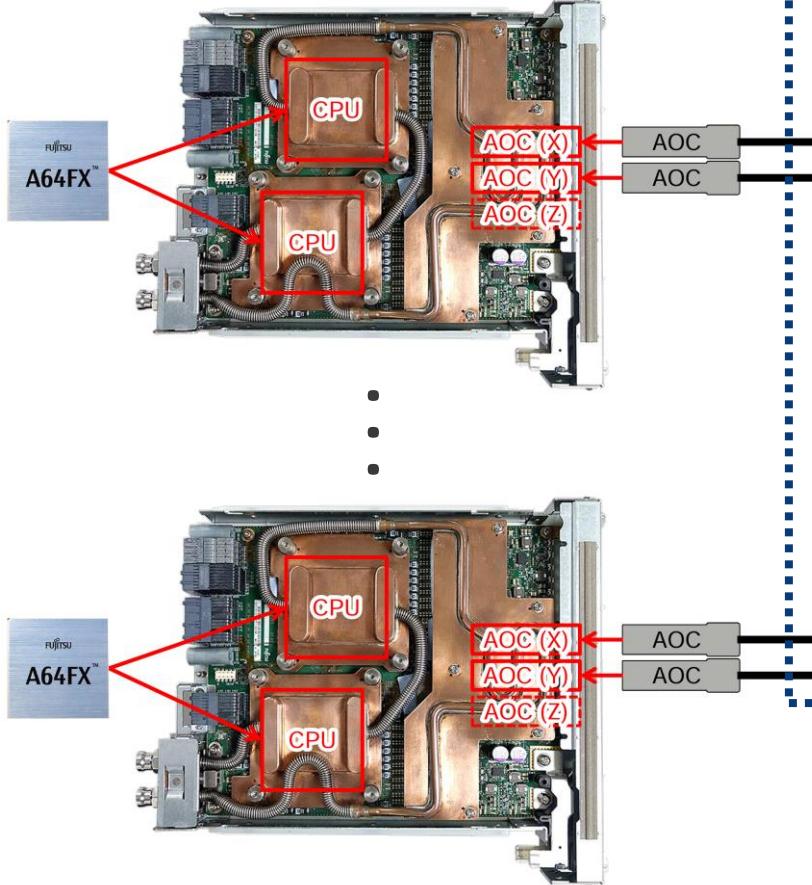
432 Racks

- 396 (Full)
- 36 (Half)

Distributed Memory

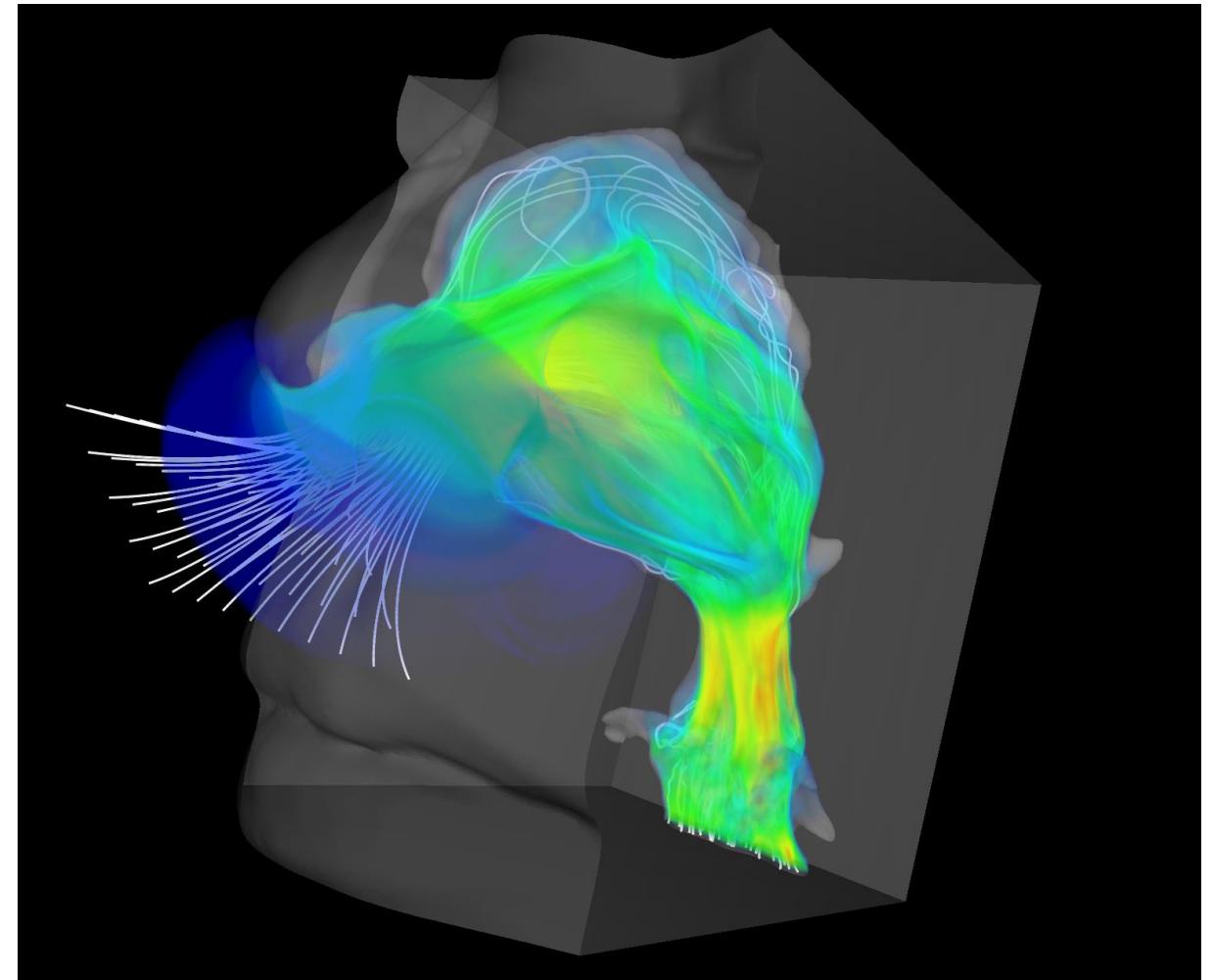
Independent Nodes

CPU Memory Unit (CMU)



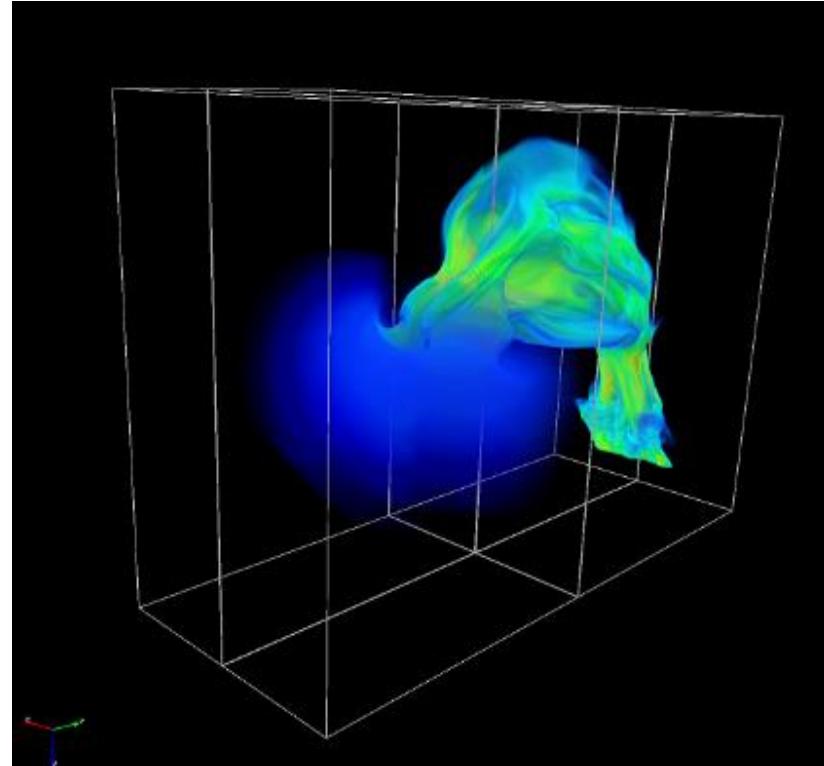
Shape (1D / 2D / 3D)

Rendering

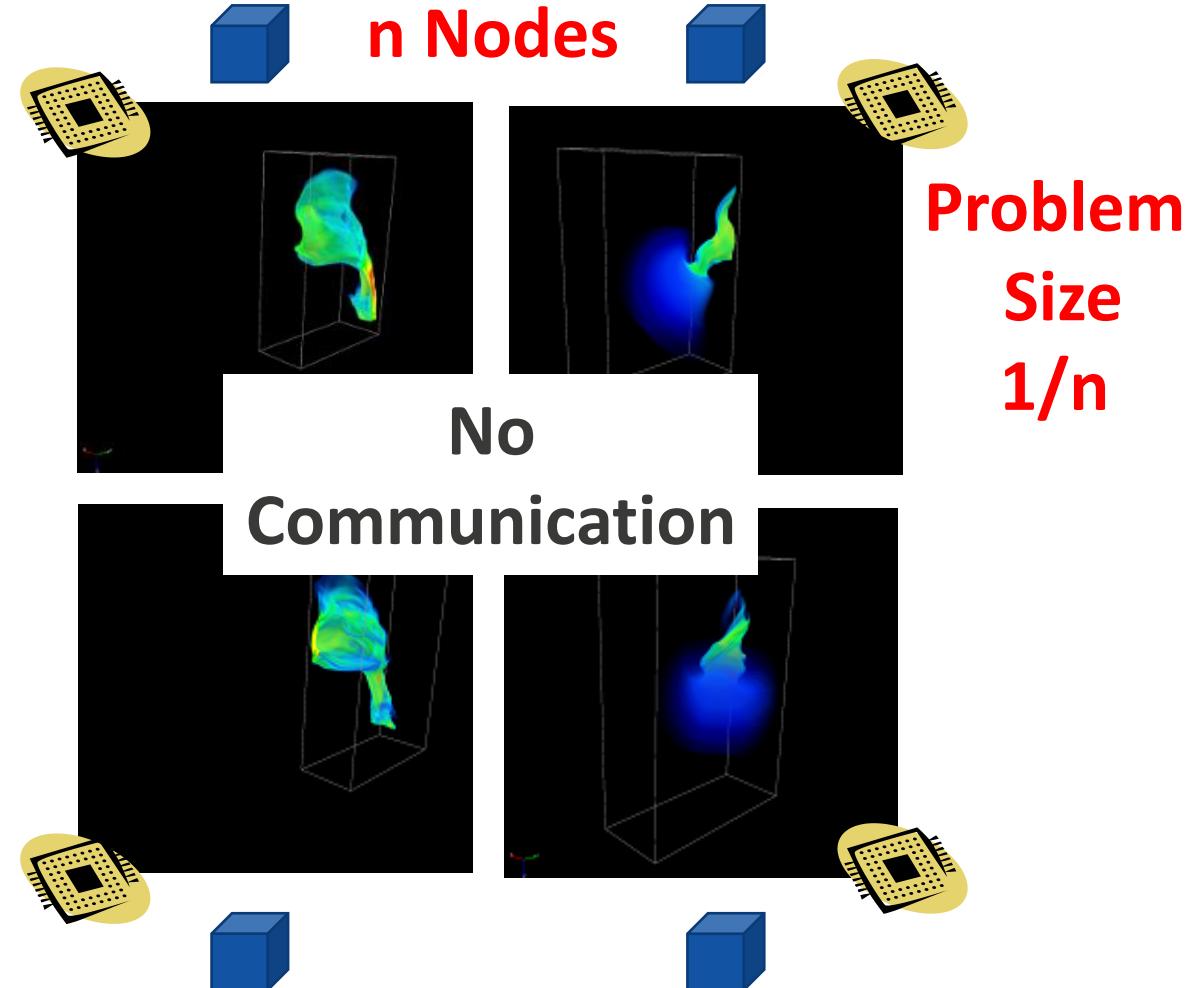


Parallel Rendering

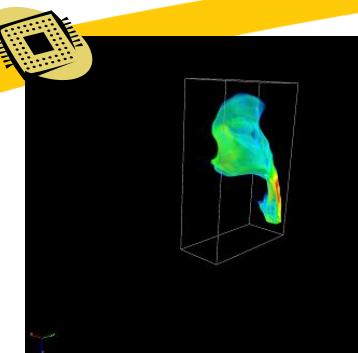
1



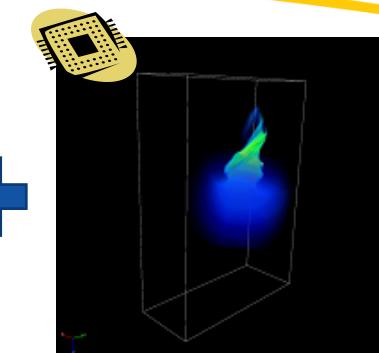
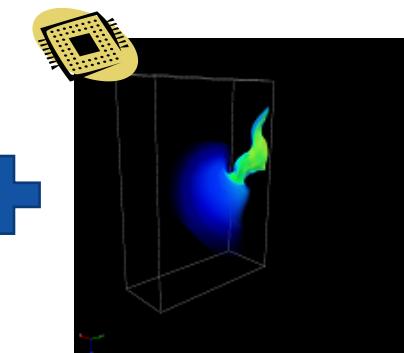
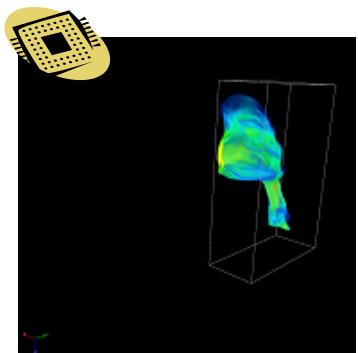
n Nodes



Parallel Image Composition



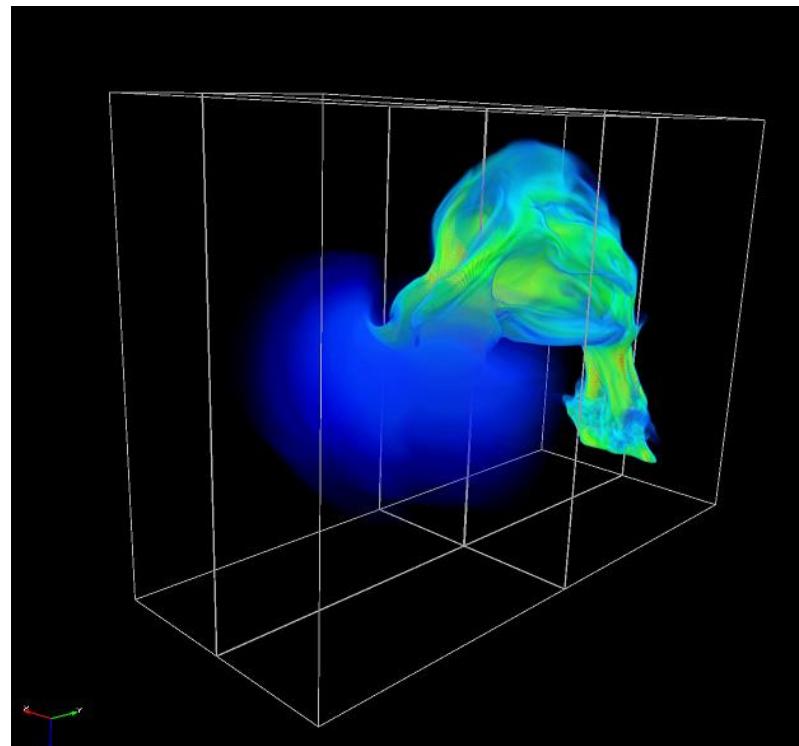
1



Problem
Size
 $1 \times n$

n Nodes

Communication
(Send/Recv
Image Data)



Computation
(Alpha Blending)

Message Passing Interface (MPI)

- “Specification” of message passing API for distributed memory environment
 - <http://www.mcs.anl.gov/mpi/www>
 - <https://www mpi-forum.org/docs>
- History
 - **1992** MPI Forum; **1994** MPI-1.0; **1997** MPI-2.0; **2012** MPI-3.0
 - **2015** MPI-3.1
 - **2021** MPI-4.0

Fugaku Users Guide

- Fujitsu MPI
 - MPI-3.1
 - MPI-4.0 subset

Based on Prof. Kengo Nakajima's Slides (Introduction to Programming by MPI for Parallel FEM)

Message Passing Interface (MPI)

- Implementation
 - MPICH
 - OpenMPI (Fujitsu Compiler)
- Cross-Compilers
 - C language
 - mpifccpx -Nclang -Kfast
 - C++ language
 - mpiFCCpx -Nclang -Kfast
 - FORTRAN language
 - mpifrtpx
- Native Compilers
 - C language
 - mpifcc
 - C++ language
 - mpiFCC
 - FORTRAN language
 - mpifrt

Spack (MPICH and OpenMPI)

```
-- linux-rhel8-a64fx / fj@4.7.0 -----
fujitsu-mpi@head
mpich-tofu@master

-- linux-rhel8-a64fx / gcc@8.4.1 -----
fujitsu-mpi@head
mpich-tofu@master
```

Hello World (C Program)

```
$ cd /vol0601/share/ra020021/ComputerScience/20221205_Nonaka
$ cp -r examples ~/
```

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

int main()
{
    printf("Hello World... \n");

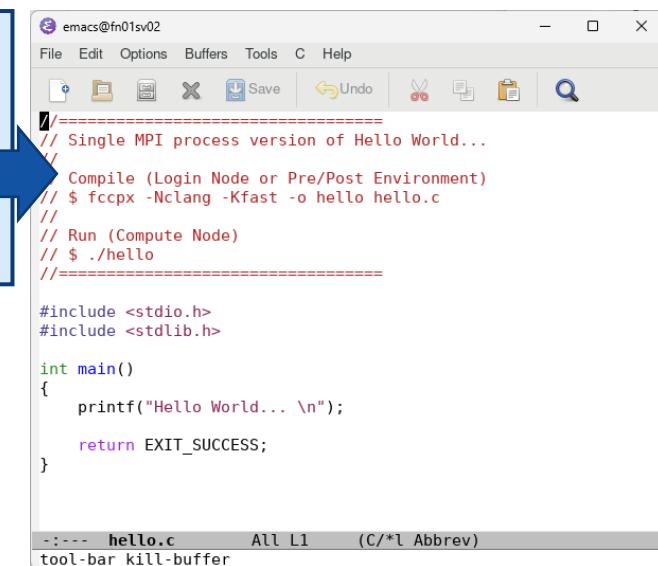
    return EXIT_SUCCESS;
}
```

hello.c

```
$ cd
$ cd examples
$ cat mpi_hello.c
$ fccpx -Nclang -Kfast -o hello hello.c
$ ./hello
```

To Edit

```
$ emacs hello.c (GUI)
$ vi hello.c
```



The screenshot shows the Emacs graphical user interface (GUI) window titled "emacs@fn01sv02". The menu bar includes File, Edit, Options, Buffers, Tools, C, and Help. The toolbar contains icons for file operations like Open, Save, Undo, and Cut/Paste. The main buffer area displays the C code for "hello.c". A blue callout box points from the text "To Edit" to the "emacs" command in the list below.

```
// Single MPI process version of Hello World...
// Compile (Login Node or Pre/Post Environment)
// $ fccpx -Nclang -Kfast -o hello hello.c
// Run (Compute Node)
// $ ./hello
=====

#include <stdio.h>
#include <stdlib.h>

int main()
{
    printf("Hello World... \n");

    return EXIT_SUCCESS;
}
```

MPI Processes (Distributed Memory)

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

int main()
{
    int my_rank;      // My MPI process number (Rank)
    int num_procs; // Total number MPI processes

    // Initialize MPI environment
    MPI_Init(NULL, NULL);
    // Get the total number of MPI processes
    MPI_Comm_size( MPI_COMM_WORLD, &num_procs );
    // Get my rank number (Start from ZERO)
    MPI_Comm_rank( MPI_COMM_WORLD, &my_rank );
```

mpi_hello.c

```
printf("Hello World... from %d of %d \n", my_rank, num_procs);

// Finalize MPI environment
MPI_Finalize();
}
```

\$ cat mpi_hello.c

```
$ mpifccpx -Nclang -Kfast -o mpi_hello mpi_hello.c
$ mpirun -n 48 ./mpi_hello
```

Batch Job (MPI)

```
#!/bin/bash
#PJM -g ra020021 # Group name

#PJM -L "rscgrp=small" # Resource group small (up to 384 nodes)
##PJM -L "rscgrp=excl_HPCS_2212-1" # EU-ASEAN HPC School: 12/5 0:00 to 12/9 24:00

#PJM -L "node=1"           # Number of nodes
#PJM --mpi "max-proc-per-node=48" # Max. MPI processes per node

#PJM -L "elapse=5:00" # Running time
#PJM -L "eco_state=2" # Eco mode setting

##PJM -x PJM_LLIO_GFSCACHE=/vol0004:/vol0006 # Volume names used by the job
#PJM -s # Output statistical information

export PLE_MPI_STD_EMPTYFILE=off
cd ${HOME}/intro_pp
mpiexec -n 48 ./mpi_hello
```

job_mpi_hello.sh

```
// Job submission
$ pjsub job_mpi_hello.sh

// Check the job status
$ pjstat

// Check the job output
$ cat job_mpi_hello.sh.xxxxxxx.err
$ cat job_mpi_hello.sh.xxxxxxx.stats

$ cd output.xxxxxxxx
$ cd x/x
$ cat stdout.x.x
```

OpenMP (Shared Memory)

- “Specification” for a set of compiler directives, library routines, and environment variables that can be used to specify high-level parallelism in Fortran and C/C++ programs
 - <https://www.openmp.org/specifications/>
- History
 - 1998 C/C++ version 1.0; 2002 C/C++ version 2.0
 - 2011 OpenMP-3.1
 - 2013 OpenMP-4.0
 - 2015 OpenMP-4.5
 - 2020 OpenMP-5.0

Fugaku Users Guide

- Fujitsu Compiler
 - OpenMP 3.1
 - OpenMP 4.0
 - OpenMP 4.5
 - OpenMP 5.0 subset

OpenMP Threads

```
#include <omp.h>  
  
#include <stdio.h>  
#include <stdlib.h>  
  
int main(int argc, char* argv[]){  
    int my_thread; // My OpenMP thread  
    int num_threads; // Total number of OpenMP threads  
  
    // Beginning of OpenMP Parallel Region  
    #pragma omp parallel  
    {  
        my_thread = omp_get_thread_num();  
        num_threads = omp_get_num_threads();  
    }
```

openmp_hello.c

```
    printf("Hello World... from thread %d of %d \n", my_thread, num_threads);  
}  
// End of OpenMP parallel region  
  
return EXIT_SUCCESS;  
}
```

\$ cat openmp_hello.c

\$ fccpx -Nclang -Kfast –Kopenmp -o openmp_hello
openmp_hello.c

\$ export OMP_NUM_THREADS=xx
\$./openmp_hello

Batch Job (OpenMP)

```
#!/bin/bash
#PJM -g ra020021 # Group name

#PJM -L "rscgrp=small"      # Resource group small (up to 384 nodes)
# #PJM -L "rscgrp=excl_HPCS_2212-1" # EU-ASEAN HPC School: 12/5 0:00 to 12/9 24:00

#PJM -L "node=1"           # Number of nodes

#PJM -L "elapse=5:00" # Running time
#PJM -L "eco_state=2" # Eco mode setting

# #PJM -x PJM_LLIO_GFSCACHE=/vol0004:/vol0006 # Volume names used by the job
#PJM -s # Output statistical information

export OMP_NUM_THREADS=48
cd ${HOME}/intro_pp
./openmp_hello
```

job_openmp_hello.sh

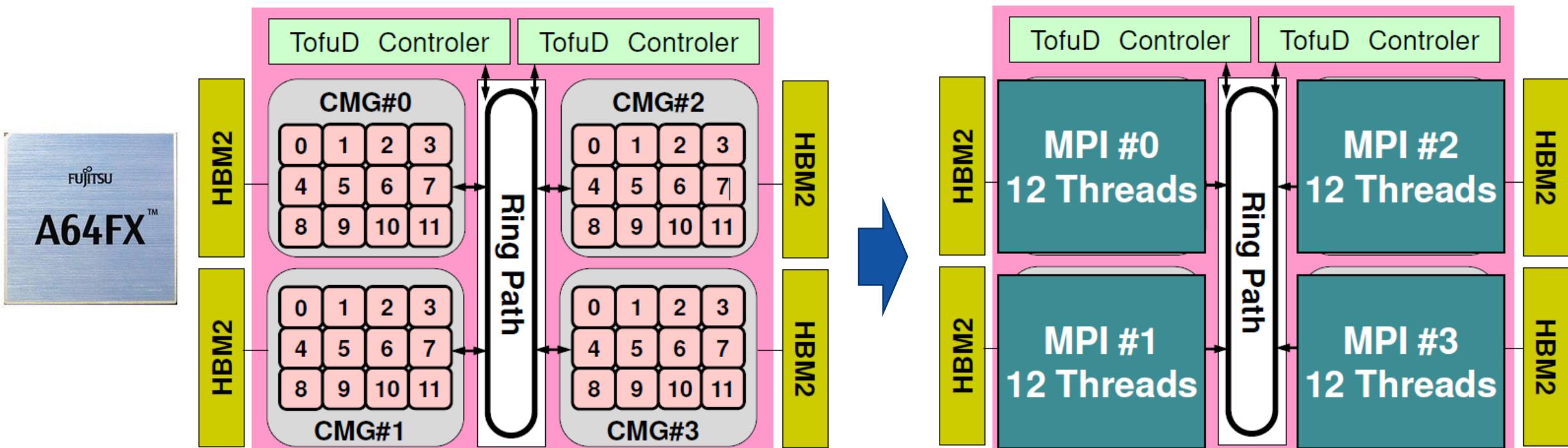
```
// Job submission
$ pbsub job_openmp_hello.sh

// Check the job status
$ pjstat

// Check the job output
$ cat job_openmp_hello.sh.xxxxxxx.err
$ cat job_openmp_hello.sh.xxxxxxx.stats

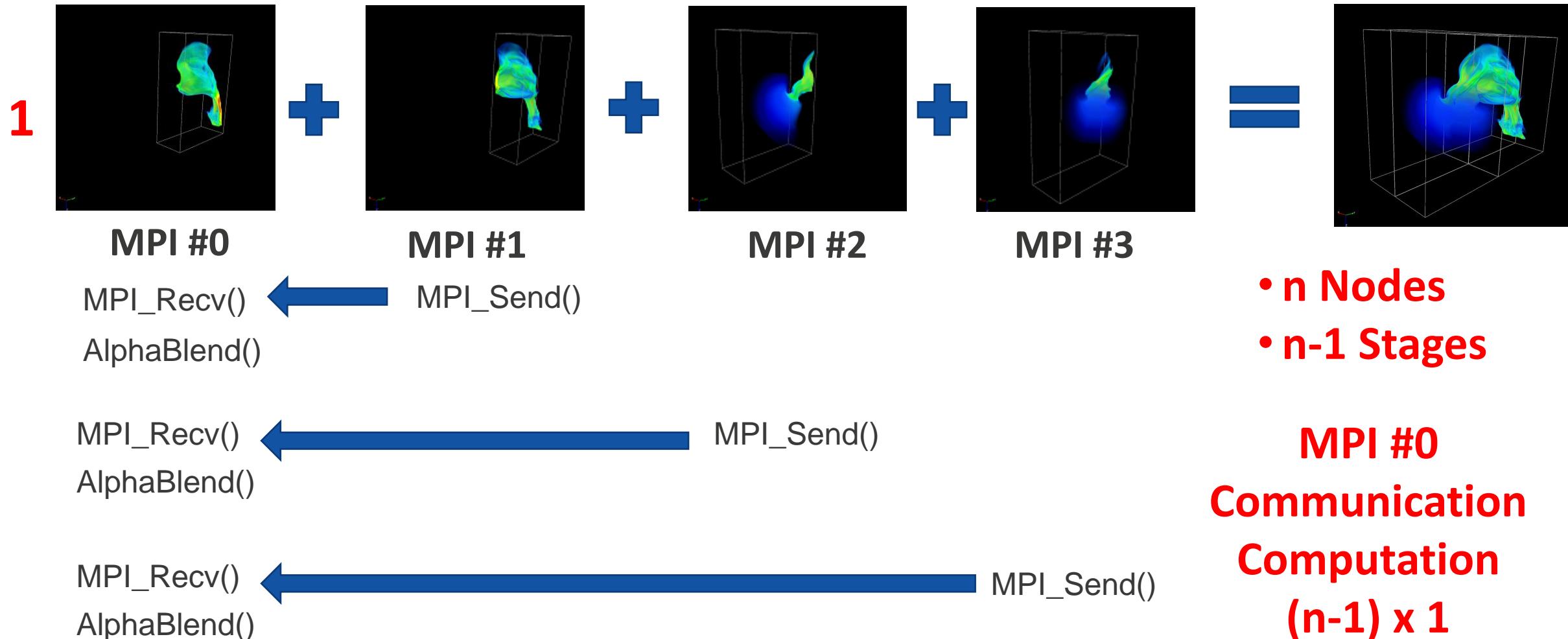
$ cd output.xxxxxxxx
$ cd x/x
$ cat stdout.x.x
```

MPI + OpenMP (Distributed and Shared Memory)

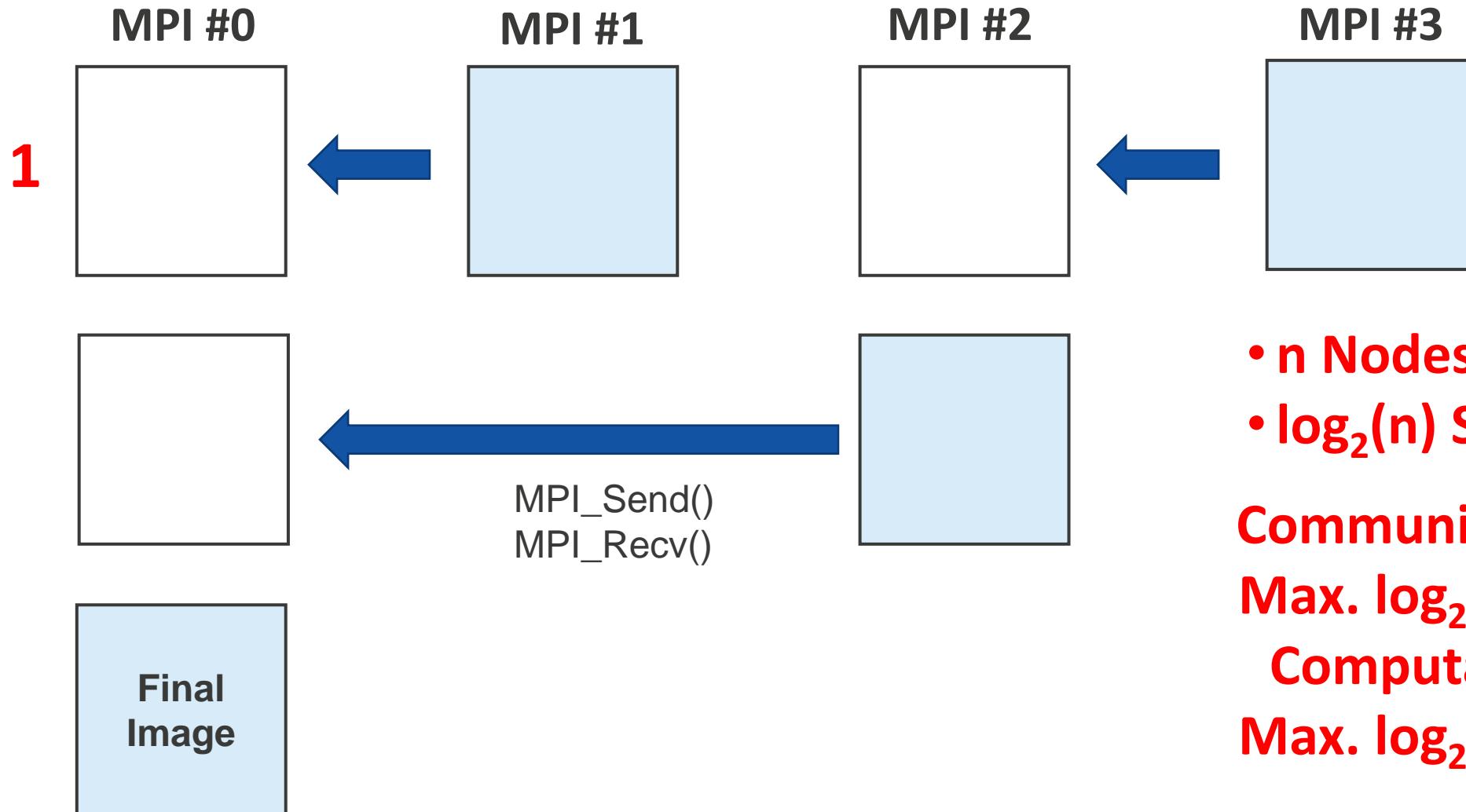


Based on Prof. Kengo Nakajima's Slides (Introduction to Programming by MPI for Parallel FEM)

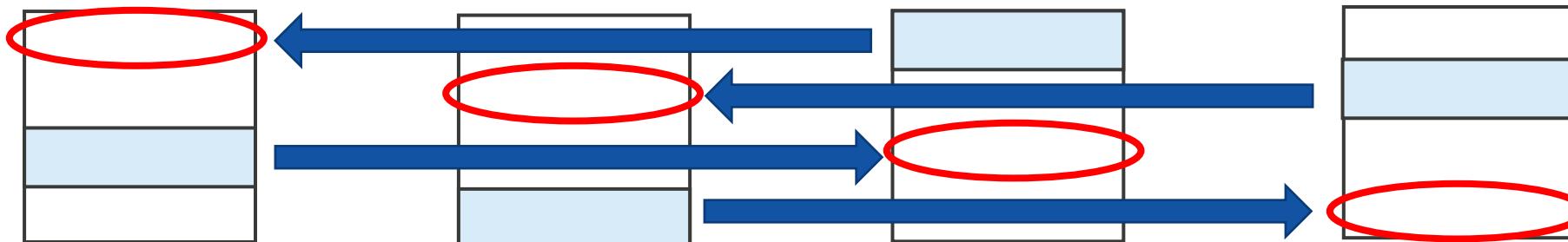
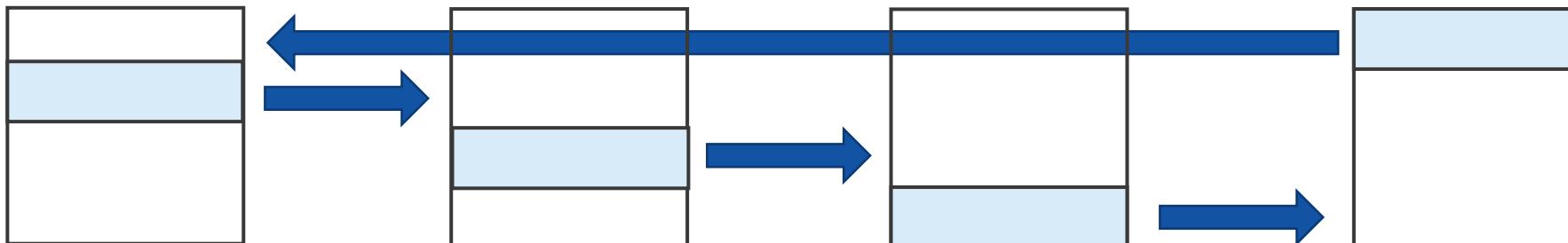
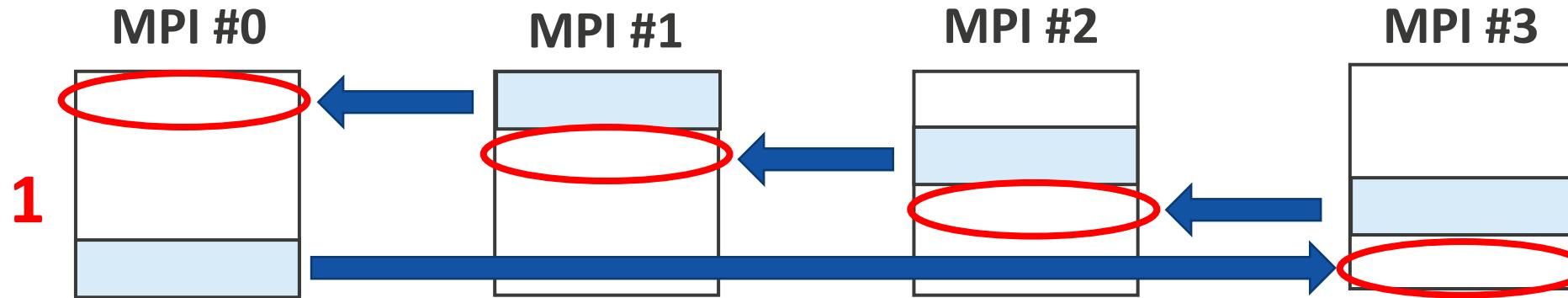
Communication & Computation



Binary Tree (2^n)



Direct Send



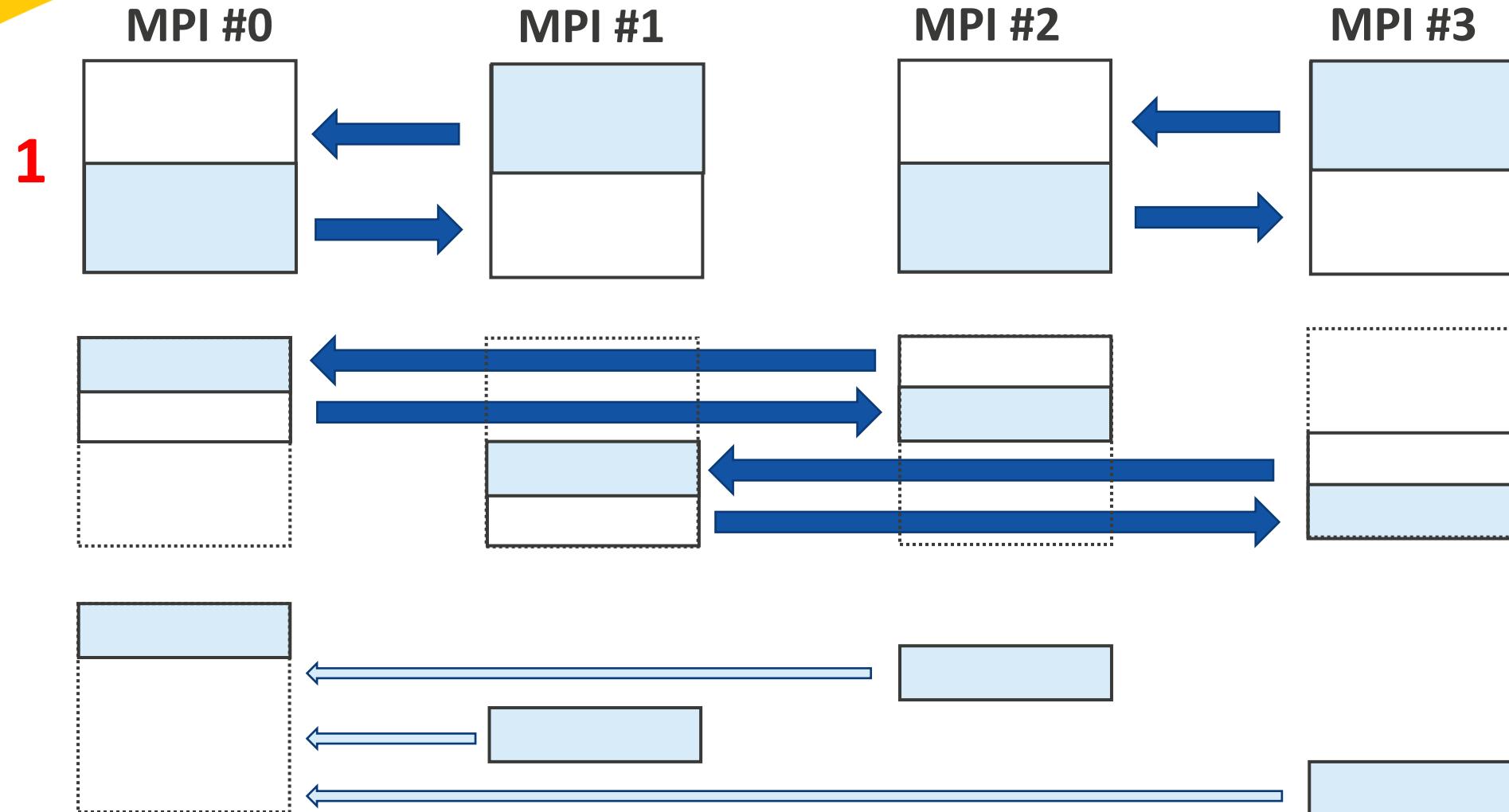
1

- n Nodes
- n-1 Stages

Communication
Computation
 $(n-1) \times 1/n$

Final Gather
Communication
 $1 - (1/n)$

Binary-Swap



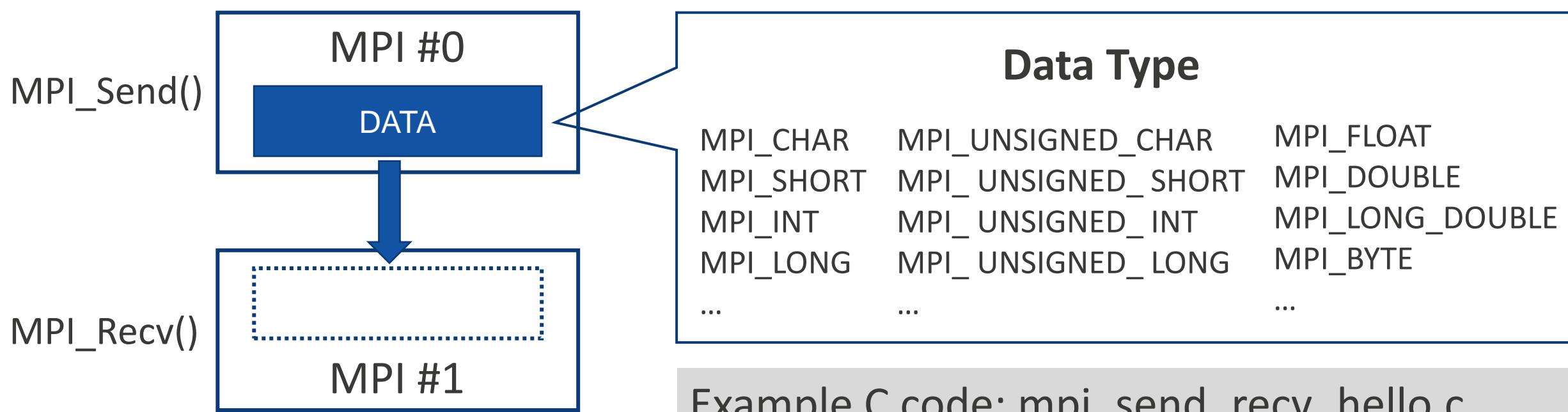
- n Nodes
- $\log_2(n)$ Stages

Communication
Computation

$$\sum_{i=1}^{\log_2 n} \left(\frac{1}{2^i}\right)$$

Final Gather
Communication
 $1 - (1/n)$

Point-to-Point Communication



Example Code (Point-to-Point Communication)

MPI_Send / MPI_Recv

```
$ cat mpi_send_recv_hello.c
```

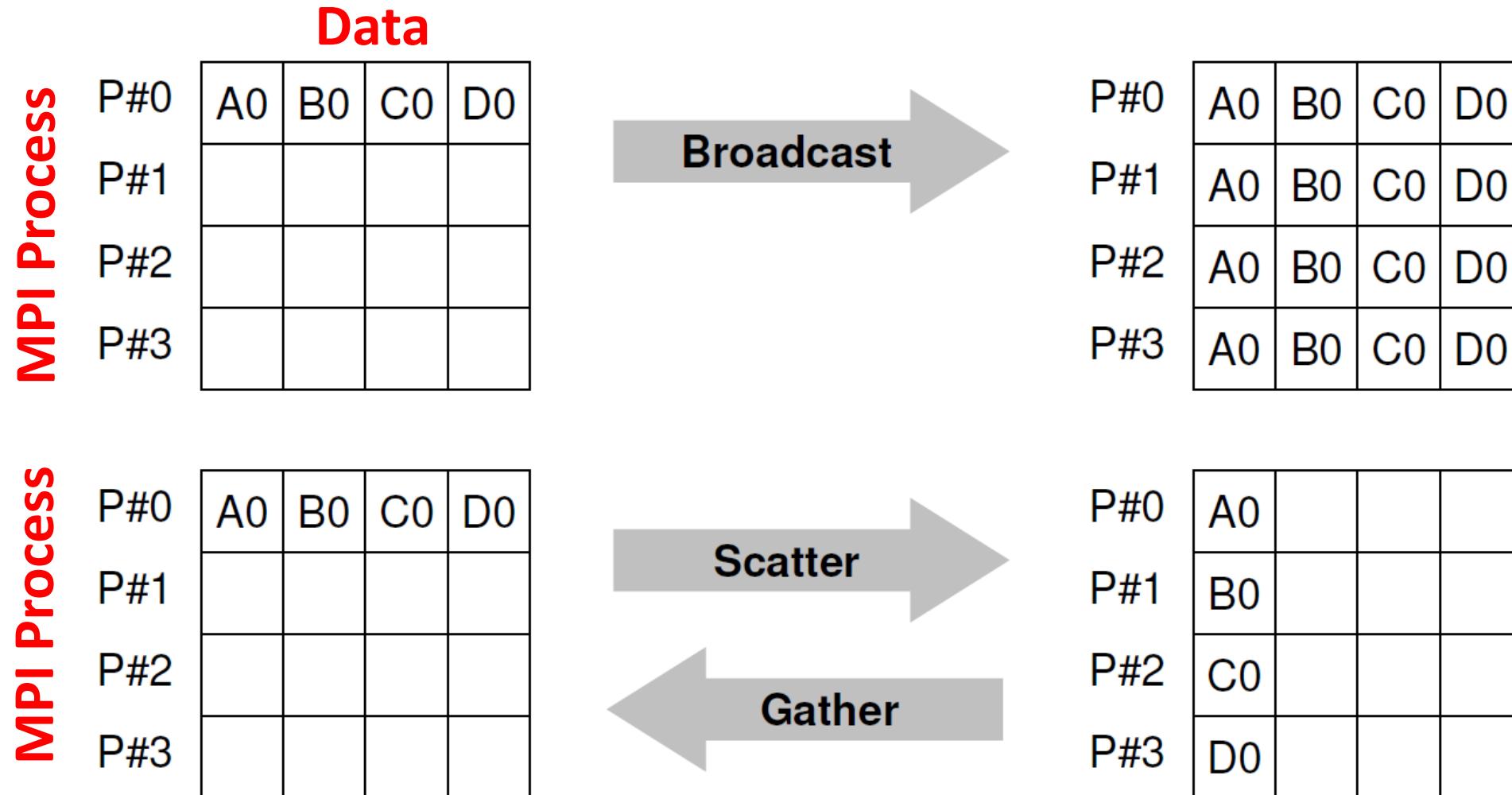
```
$ mpifccpx -Nclang -Kfast -o mpi_send_recv_hello mpi_send_recv_hello.c
```

```
$ pbsub mpi_send_recv_hello.c
```

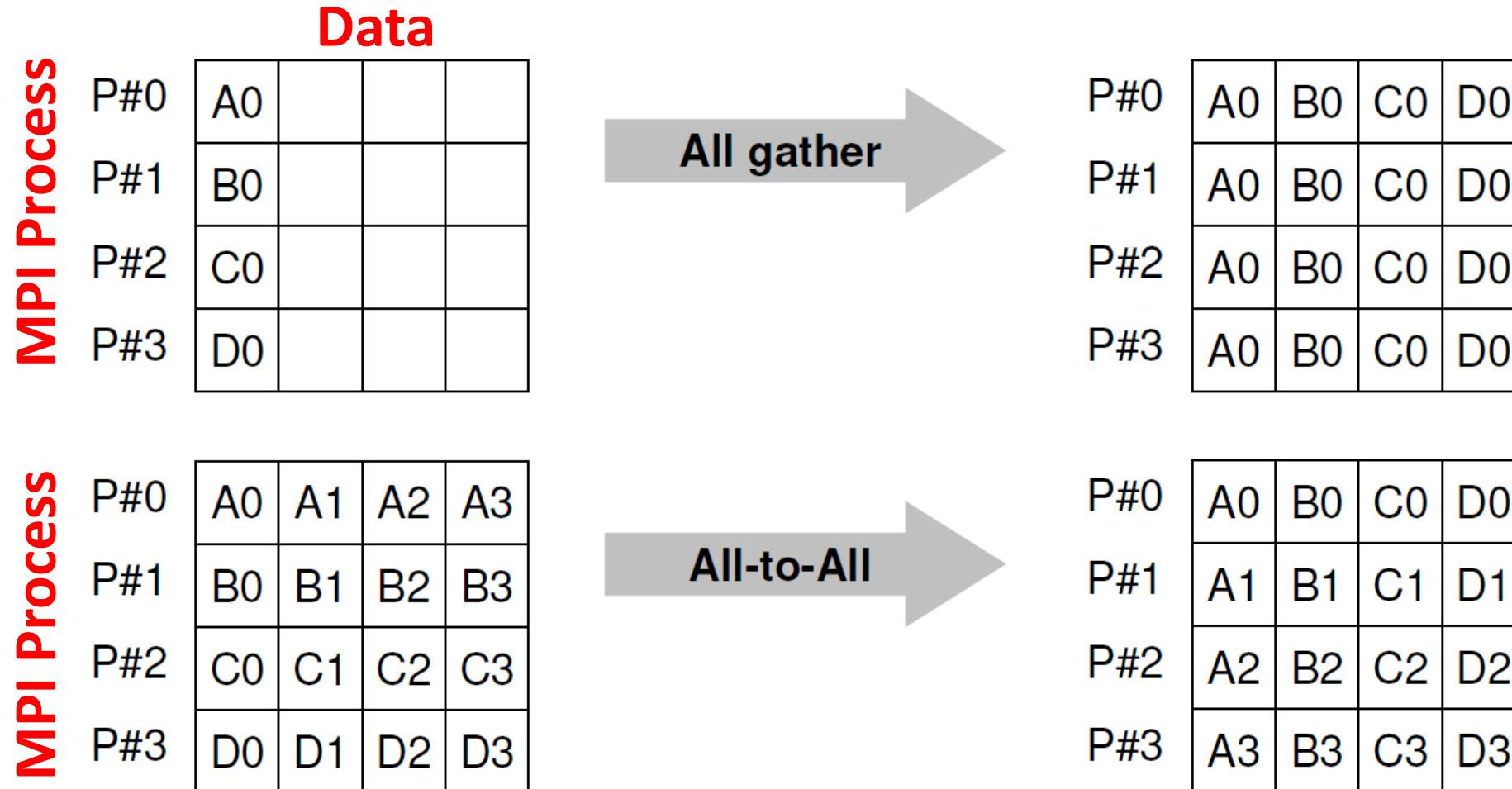
```
$ pjstat
```

```
$ cat output.xxxxxxxxx/x/x/stdout.x.x
```

Collective Communication



Collective Communication



Collective Communication

MPI Process

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				

op: Operation

MPI Process

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

Collective Communication

op: Operation

MPI Process

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			

Some predefined reduce operations

MPI_MAX	MPI_LAND	MPI_LXOR
MPI_MIN	MPI_BAND	MPI_BXOR
MPI_SUM	MPI_LOR	MPI_MAXLOC
MPI_PROD	MPI_BOR	MPI_MINLOC
...

Example Codes (Collective Communication)

MPI_Bcast

```
$ cat mpi_bcast_random.c
$ mpifccpx -Nclang -Kfast -o mpi_bcast_random mpi_bcast_random.c
$ pbsub job_mpi_bcast_random.sh
```

MPI_Allreduce

```
$ cat mpi_allreduce_random.c
$ mpifccpx -Nclang -Kfast -o mpi_allreduce_random mpi_allreduce_random.c
$ pbsub job_mpi_allreduce_random.sh
```

Send the output of the MPI_AllReduce job

```
/vol0601/share/ra020021  
ComputerScience/20221205_Nonaka  
job_output
```

```
$ cp -r output.xxxxxx /vol0601/share/ra020021/  
ComputerScience/20221205_Nonaka/job_output
```