

# Report S1 C

Kengo Nakajima  
Information Technology Center

Technical & Scientific Computing II (4820-1028)  
Seminar on Computer Science II (4810-1205)

# 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

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

# Copying files on Oakleaf-FX

## Copy

```
>$ cd <$O-TOP>  
>$ cp /home/z30088/class_eps/C/s1r-c.tar .  
>$ tar xvf s1r-c.tar
```

## Confirm directory

```
>$ ls  
mpi  
>$ cd mpi/s1-ref
```

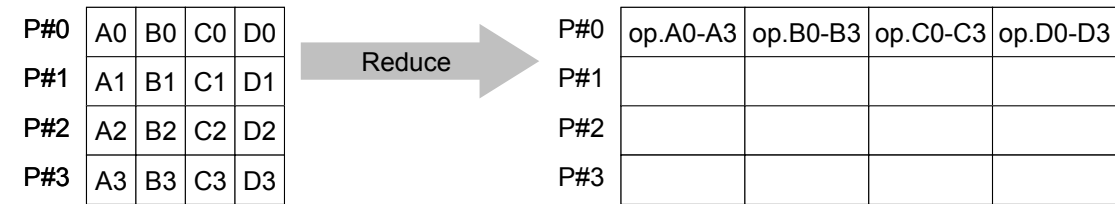
This directory is called as `<$O-s1r>`.

`<$O-s1r> = <$O-TOP>/mpi/s1-ref`

# S1-1 : Reading Local Vector, Calc. Norm

- 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.
- Use MPI\_Allreduce (or MPI\_Reduce)
- Advice
  - Checking each component of variables and arrays !

# 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/recive buffer
    - FORTTRAN 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.

# “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

```
double x0, xsum;
```

```
MPI_Reduce
```

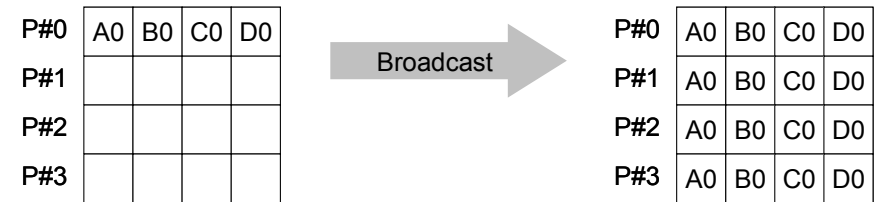
```
(&x0, &xsum, 1, MPI_DOUBLE, MPI_SUM, 0, <comm>)
```

```
double x0[4];
```

```
MPI_Reduce
```

```
(&x0[0], &x0[2], 2, MPI_DOUBLE_PRECISION, MPI_SUM, 0, <comm>)
```

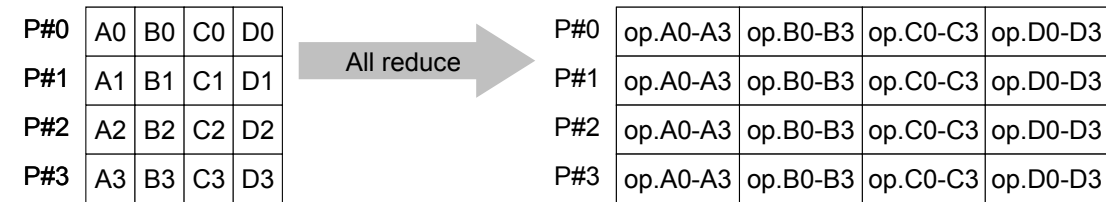
# MPI\_Bcast



- 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/recive buffer  
FORTRAN 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



- MPI\_Reduce + MPI\_Bcast
- Summation (of dot products) and MAX/MIN values are likely to be 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/recv buffer
  
- op            MPI\_Op        I        reduce operation
- comm         MPI\_Comm     I        communicator

# S1-1 : Local Vector, Norm Calculation

## Uniform Vectors (a1.\*): s1-1-for\_a1.c

```

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

int main(int argc, char **argv){
    int i, N;
    int PeTot, MyRank;
    MPI_Comm SolverComm;
    double vec[8];
    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, "a1.%d", MyRank);
    fp = fopen(filename, "r");
    assert(fp != NULL);

    N=8;
    for(i=0;i<N;i++){
        fscanf(fp, "%lf", &vec[i]);
    }
    sum0 = 0.0;
    for(i=0;i<N;i++){
        sum0 += vec[i] * vec[i];
    }

    MPI_Allreduce(&sum0, &sum, 1, MPI_DOUBLE, MPI_SUM, MPI_COMM_WORLD);
    sum = sqrt(sum);

    if(!MyRank) printf("%27.20E\n", sum);
    MPI_Finalize();
    return 0;
}

```

# S1-1: Local Vector, Norm Calculation

## Non-uniform Vectors (a2.\*): s1-1-for\_a2.c

```

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

int main(int argc, char **argv){
    int i, PeTot, MyRank, n;
    MPI_Comm SolverComm;
    double *vec, *vec2;
    int * Count, CountIndex;
    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]);}
    sum0 = 0.0;
    for(i=0;i<n;i++){
        sum0 += vec[i] * vec[i];}

    MPI_Allreduce(&sum0, &sum, 1, MPI_DOUBLE, MPI_SUM, MPI_COMM_WORLD);
    sum = sqrt(sum);

    if(!MyRank) printf("%27.20E\n", sum);
    MPI_Finalize();
    return 0;}

```

# S1-1: Running the Codes

```
$ cd <$0-S1r>  
$ mpifccpx -Kfast s1-1-for_a1.c  
$ mpifccpx -Kfast s1-1-for_a2.c  
  
(modify "go4.sh")  
$ pjsub go4.sh
```

# S1-1 : Local Vector, Calc. Norm Results

## Results using one core

```
a1.* 1.62088247569032590000E+03  
a2.* 1.22218492872396360000E+03
```

```
$> frtpx -Kfast dot-a1.f  
$> pjsub gol.sh
```

```
$> frtpx -Kfast dot-a2.f  
$> pjsub gol.sh
```

## Results

```
a1.* 1.62088247569032590000E+03  
a2.* 1.22218492872396360000E+03
```

## gol.sh

```
#!/bin/sh  
#PJM -L "node=1"  
#PJM -L "elapsed=00:10:00"  
#PJM -L "rscgrp=lecture5"  
#PJM -g "gt95"  
#PJM -j  
#PJM -o "test.lst"  
#PJM --mpi "proc=1"  
  
mpiexec ./a.out
```

# S1-1: Local Vector, Calc. Norm

If SENDBUF=RECVBUF, what happens ?

## True

```
MPI_Allreduce(&sum0, &sum, 1, MPI_DOUBLE, MPI_SUM,  
             MPI_COMM_WORLD)
```

## False

```
MPI_Allreduce(&sum0, &sum0, 1, MPI_DOUBLE, MPI_SUM,  
             MPI_COMM_WORLD)
```

# S1-1: Local Vector, Calc. Norm

If SENDBUF=RECVBUF, what happens ?

## True

```
MPI_Allreduce(&sum0, &sum, 1, MPI_DOUBLE, MPI_SUM,  
             MPI_COMM_WORLD)
```

## False

```
MPI_Allreduce(&sum0, &sum0, 1, MPI_DOUBLE, MPI_SUM,  
             MPI_COMM_WORLD)
```

## True

```
MPI_Allreduce(&sumK[1], &sumK[2], 1, MPI_DOUBLE, MPI_SUM,  
             MPI_COMM_WORLD)
```

SENDBUF .ne. RECVBUF

# S1-2: Integration by Trapezoidal Rule

- 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

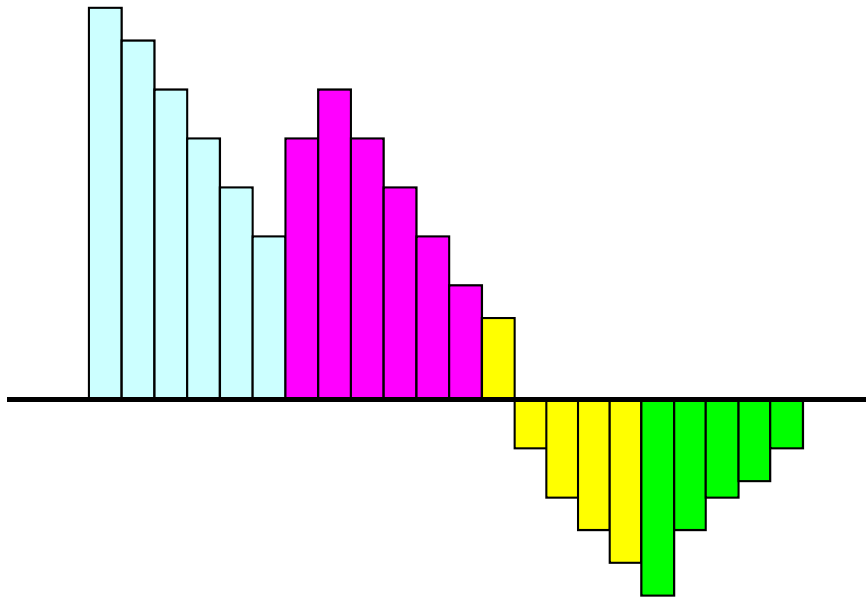
$$\int_0^1 \frac{4}{1+x^2} dx$$



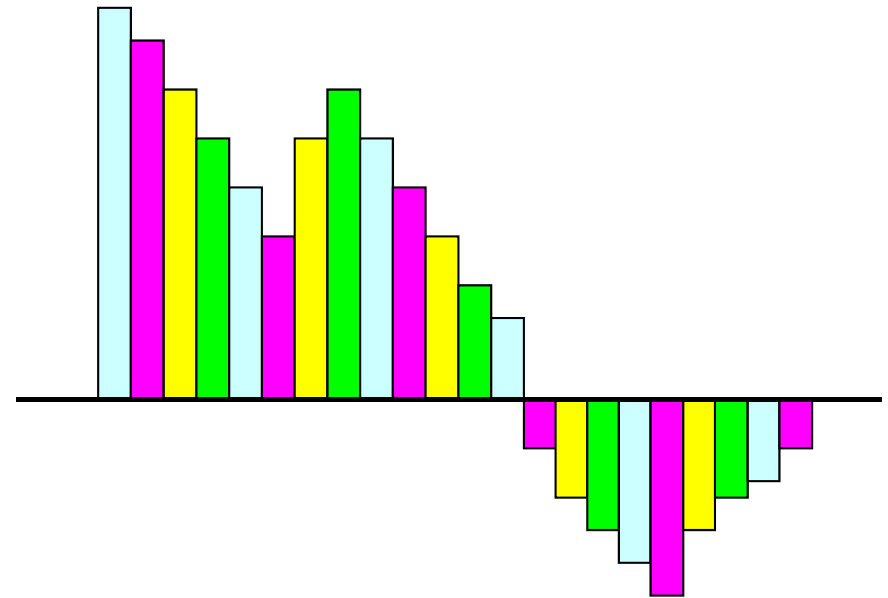
# S1-2: Integration by Trapezoidal Rule

## Two Types of Load Distribution

Type-A



Type-B



$$\frac{1}{2} \Delta x \left( f_1 + f_{N+1} + \sum_{i=2}^N 2f_i \right) \text{ corresponds to "Type-A".}$$

# S1-2: Integration by Trapezoidal Rule

## TYPE-A (1/2): s1-3a.c

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

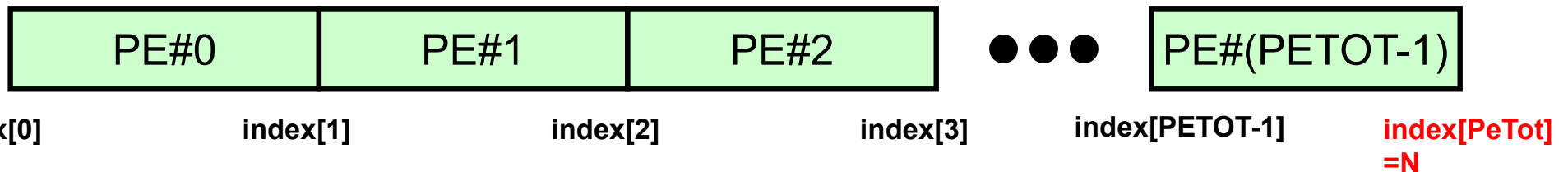
int main(int argc, char **argv){
    int i;
    double TimeStart, TimeEnd, sum0, sum, dx;
    int PeTot, MyRank, n, int *index;
    FILE *fp;

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

    index = calloc(PeTot+1, sizeof(int));
    fp = fopen("input.dat", "r");
    fscanf(fp, "%d", &n);
    fclose(fp);
    if(MyRank==0) printf("%s%8d¥n", "N=", n);
    dx = 1.0/n;

    for(i=0;i<=PeTot;i++){
        index[i] = ((long long)i * n)/PeTot;}
}
```

“N (number of segments) “ is specified in “input.dat”



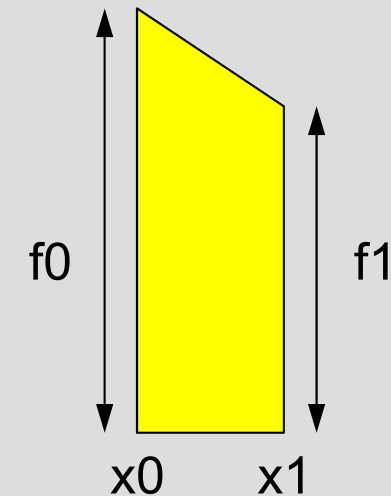
# S1-2: Integration by Trapezoidal Rule

TYPE-A (2/2): s1-3a.c

```

TimeS = MPI_Wtime();
sum0 = 0.0;
for(i=index[MyRank]; i<index[MyRank+1]; i++)
{
    double x0, x1, f0, f1;
    x0 = (double)i * dx;
    x1 = (double)(i+1) * dx;
    f0 = 4.0/(1.0+x0*x0);
    f1 = 4.0/(1.0+x1*x1);
    sum0 += 0.5 * (f0 + f1) * dx;
}

```



```

MPI_Reduce(&sum0, &sum, 1, MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD);
TimeE = MPI_Wtime();

```

```

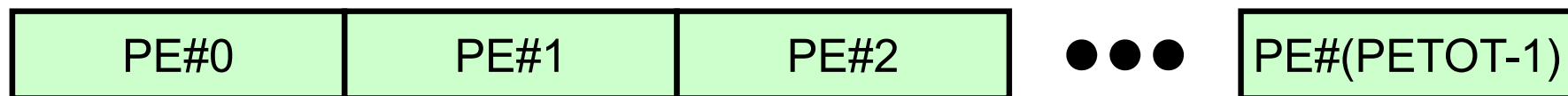
if(!MyRank) printf("%24.16f%24.16f%24.16f¥n", sum, 4.0*atan(1.0), TimeE - TimeS);

```

```

MPI_Finalize();
return 0;
}

```



index[0]

index[1]

index[2]

index[3]

index[PETOT-1]

index[PeTot]  
=N

# S1-2: Integration by Trapezoidal Rule

TYPE-B: s1-3b.c

```

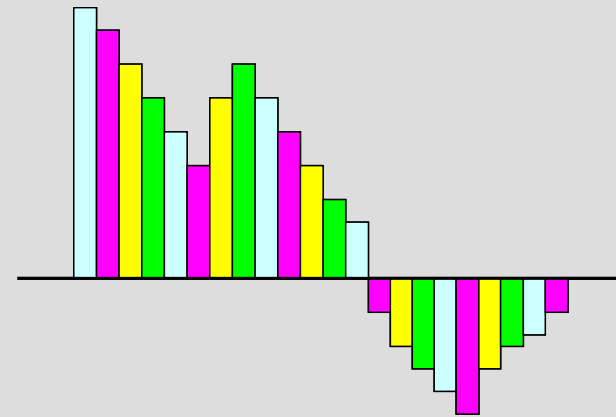
TimeS = MPI_Wtime();
sum0 = 0.0;
for(i=MyRank; i<n; i+=PeTot)
{
    double x0, x1, f0, f1;
    x0 = (double)i * dx;
    x1 = (double)(i+1) * dx;
    f0 = 4.0/(1.0+x0*x0);
    f1 = 4.0/(1.0+x1*x1);
    sum0 += 0.5 * (f0 + f1) * dx;
}

MPI_Reduce(&sum0, &sum, 1, MPI_DOUBLE, MPI_SUM, 0, MPI_COMM_WORLD);
TimeE = MPI_Wtime();

if(!MyRank) printf("%24.16f%24.16f%24.16f\n", sum, 4.0*atan(1.0), TimeE-TimeS);

MPI_Finalize();
return 0;
}

```

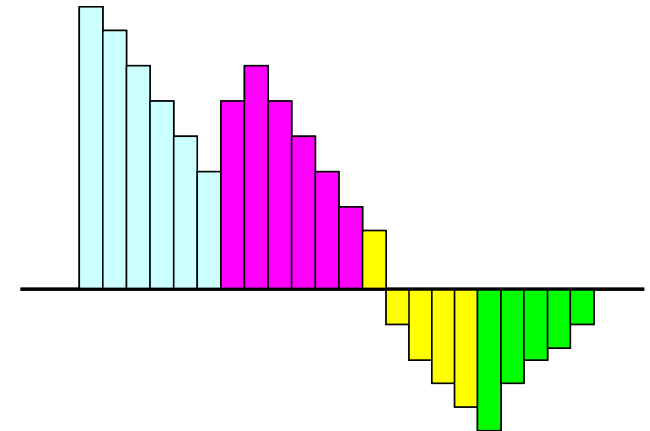


# S1-2: Running the Codes

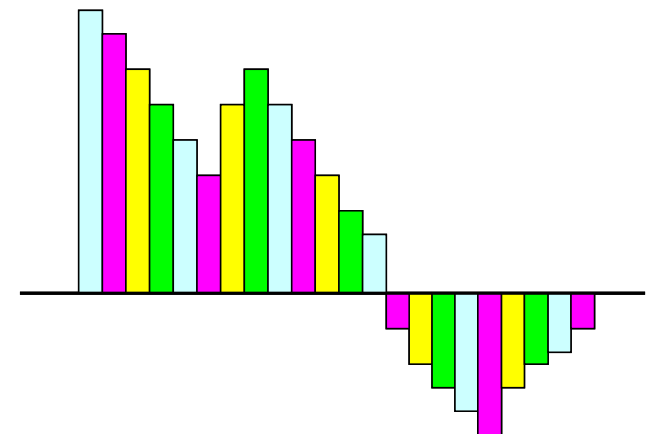
```
$ mpifccpx -Kfast s1-3a.c  
$ mpifccpx -Kfast s1-3b.c
```

```
(modify "go.sh")  
$ pjsub go.sh
```

Type-A



Type-B



# go.sh

```

#!/bin/sh
#PJM -L "node=1"           Node # (.1e.12)
#PJM -L "elapse=00:10:00"  Comp.Time (.1e.15min)
#PJM -L "rscgrp=lecture5"  "Queue" (or lecture4)
#PJM -g "gt95"             "Wallet"
#PJM -
#PJM -o "test.lst"         Standard Output
#PJM --mpi "proc=8"        MPI Process # (.1e.192)

mpiexec ./a.out

```

N=8

"node=1"

"proc=8"

N=16

"node=1"

"proc=16"

N=32

"node=2"

"proc=32"

N=64

"node=4"

"proc=64"

N=192

"node=12"

"proc=192"

# S1-2: Performance on Oakleaf-FX

- ◆ :  $N=10^6$ , ● :  $10^8$ , ▲ :  $10^9$ , — : Ideal
- Based on results (sec.) using a single core

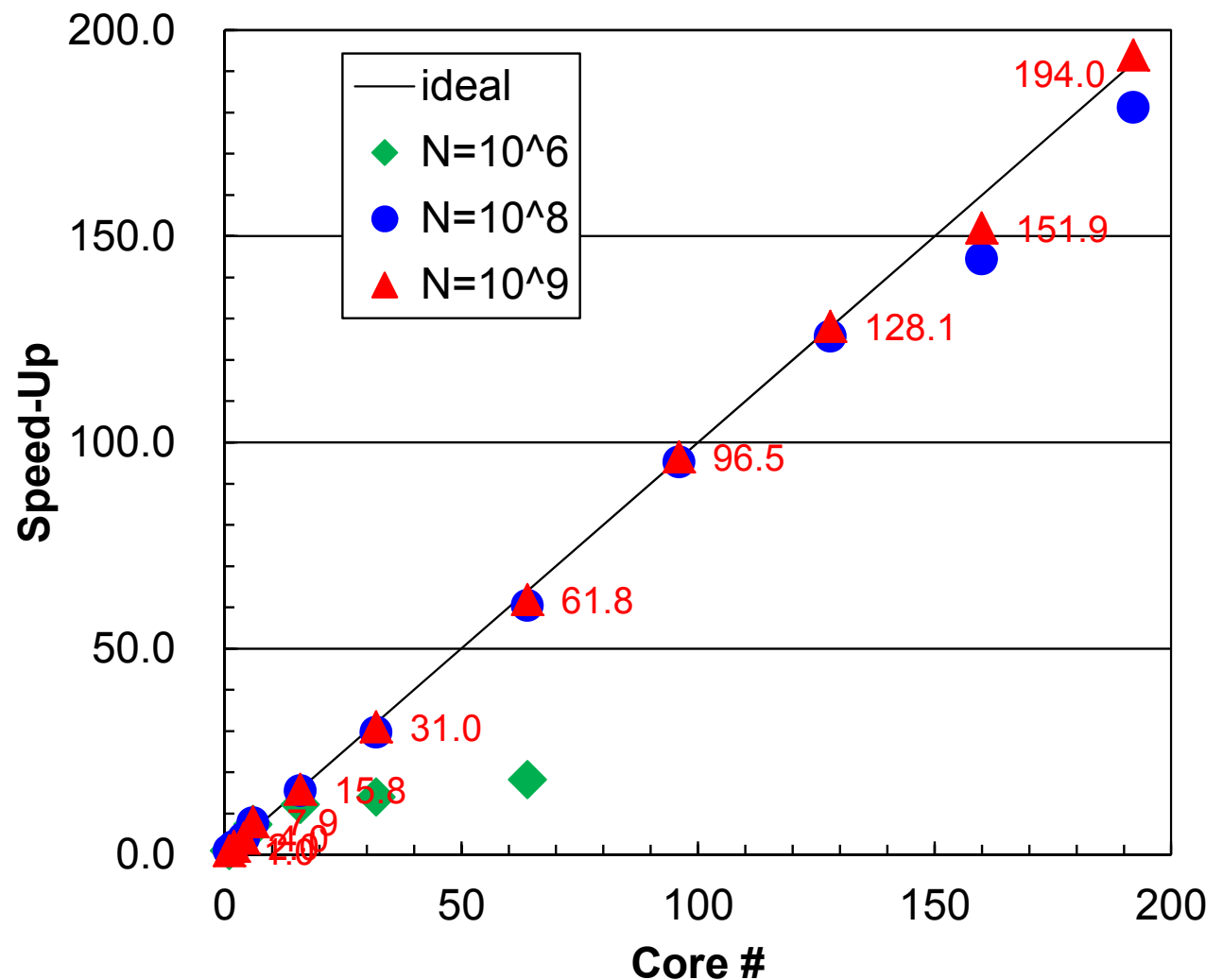
- Strong Scaling**

- Entire problem size fixed
- $1/N$  comp. time using  $N$ -x cores

- Weak Scaling**

- Problem size/core is fixed
- Comp. time is kept constant for  $N$ -x scale problems

**S1-3** using  $N$ -x cores



# Performance is lower than ideal one

- Time for MPI communication
  - Time for sending data
  - Communication bandwidth between nodes
  - Time is proportional to size of sending/receiving buffers
- Time for starting MPI
  - latency
  - does not depend on size of buffers
    - depends on number of calling, increases according to process #
  - $O(10^0)$ - $O(10^1)$   $\mu$ sec.
- Synchronization of MPI
  - Increases according to number of processes



# Performance is lower than ideal one (cont.)

- If computation time is relatively small ( $N$  is small in S1-3), these effects are not negligible.
  - If the size of messages is small, effect of “latency” is significant.