# Ch3 Distributed—mem programming with MPI



IN  ex) MPP, COW

msg —passing interface

```
$> mpicc ...
$> mpic++ p1.cpp ⏎
$> ⌠mpi exec  —n 4  ./a.out ⏎
   ⌡mpirun  —np 4  ./a.out ⏎
```



MPI_Init (NULL, NULL); // initialize a MPI session
⋮
○
MPI_Finalize();

MPI_Comm_size (MPI_COMM_WORLD,
&comm_size);

Communicator          int comm_size
Collection of user
defined processes
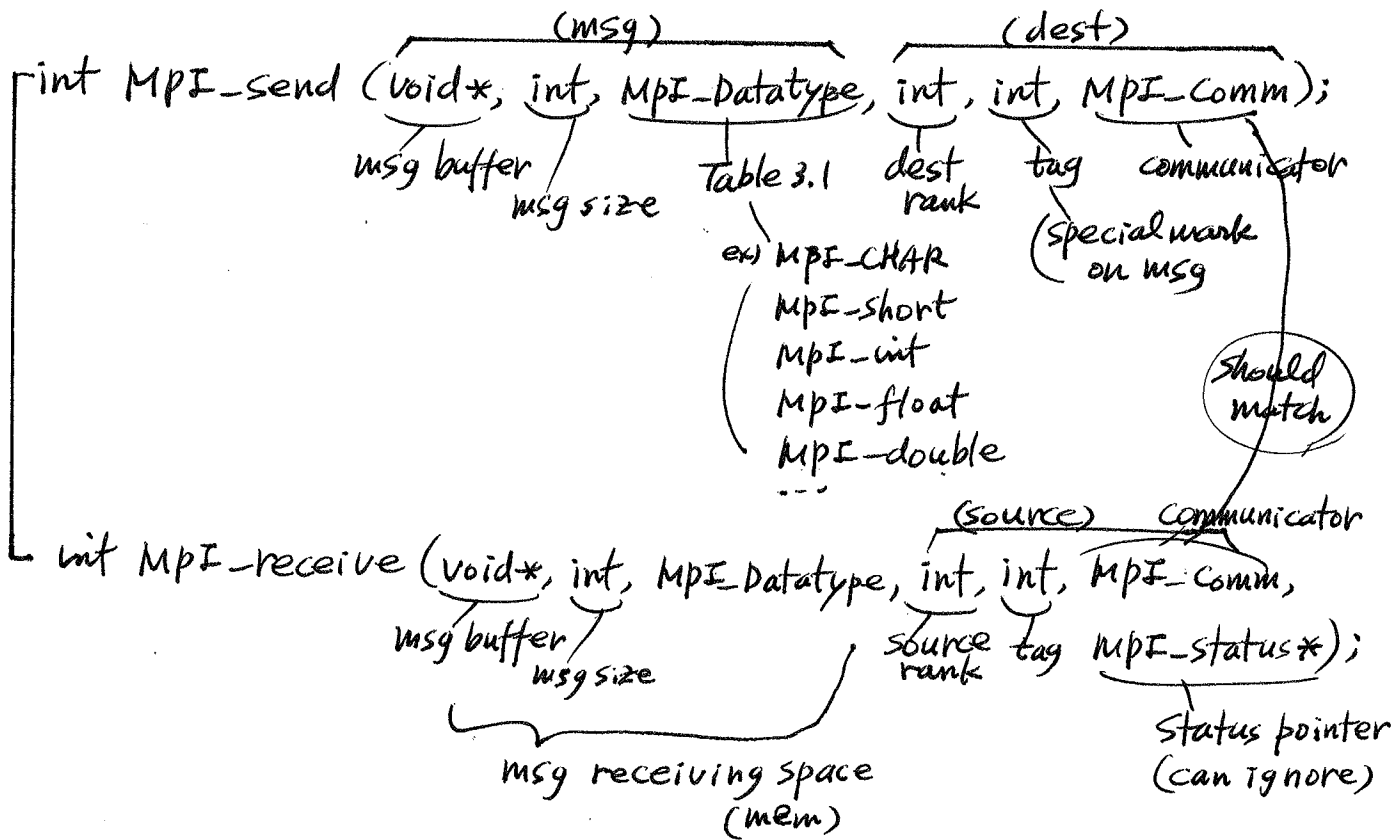gets the number of processes from
—n ④

MPI_comm_rank (MPI_COMM_WORLD,
&my_rank);

—SPMD programming
different task on
different process

gets my_rank
among comm_size

int my_rank

— 2 basic communication functions (p2p: point-to-point)

┌ int MPI_send (void*, int, MPI_Datatype, int, int, MPI_Comm);
│                    (msg)              (dest)

- void* → msg buffer
- int → msg size
- MPI_Datatype → Table 3.1

  ex) MPI_CHAR
      MPI_short
      MPI_int
      MPI_float
      MPI_double
      ...

- int (dest) → dest rank
- int → tag (special mark on msg)
- MPI_Comm → communicator

  (Should match)

└ int MPI_receive (void*, int, MPI_Datatype, int, int, MPI_Comm, MPI_status*);
                        (source)   communicator

- void* → msg buffer
- int → msg size
- int (source) → source rank
- int → tag
- MPI_status* → status pointer (can ignore)

msg receiving space (mem)

— for a successful communication,

$$\begin{cases} Send\_comm = recv\_comm \\ send\_tag = recv\_tag \\ dest = source \end{cases}$$

also, $\begin{cases} send\_type = recv\_type \\ send\_buffer\_size \leq recv\_buffer\_size \end{cases}$

— $\begin{cases} sender — must specify receiver \\ receiver — can use \underline{wild\ card} \end{cases}$

"MPI_ANY_source" — (for source rank)

— problem: receiver can receive a msg with unknown {size, source, tag}

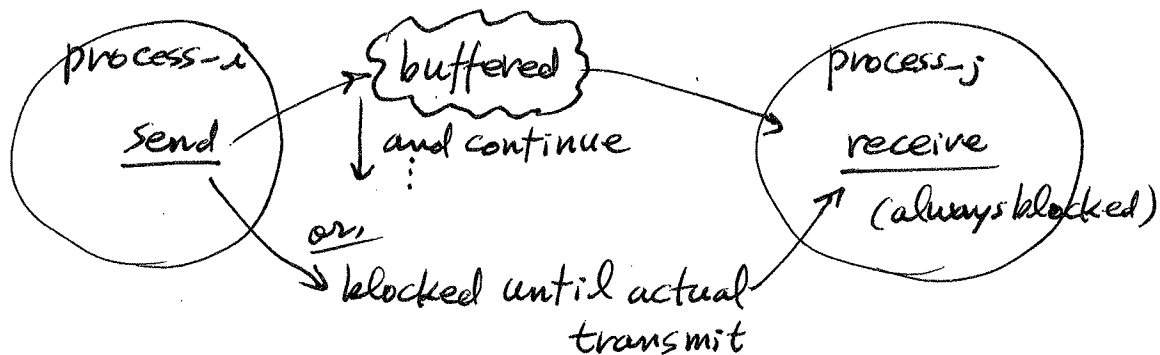└ sol:        MPI_Get_Count(..)    ← using status

↓

ex)
```
          ⋮
    MPI_status  status;
    MPI_Recv (...., &status);
                        ↑
   ⎛ status.MPI_Source ↝ source
   ⎝ status.MPI_Tag ↝ tag
          ⋮
    MPI_Get_count (&status, recv_type, &count);
                       ↓         ↓        ↑
```
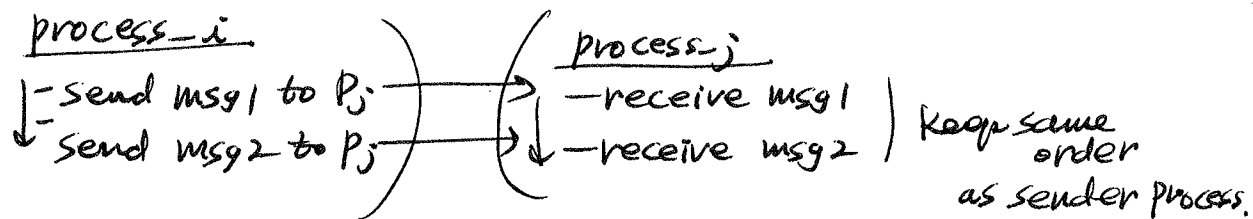└─── size of msg

— should avoid hanging receive (#sender)
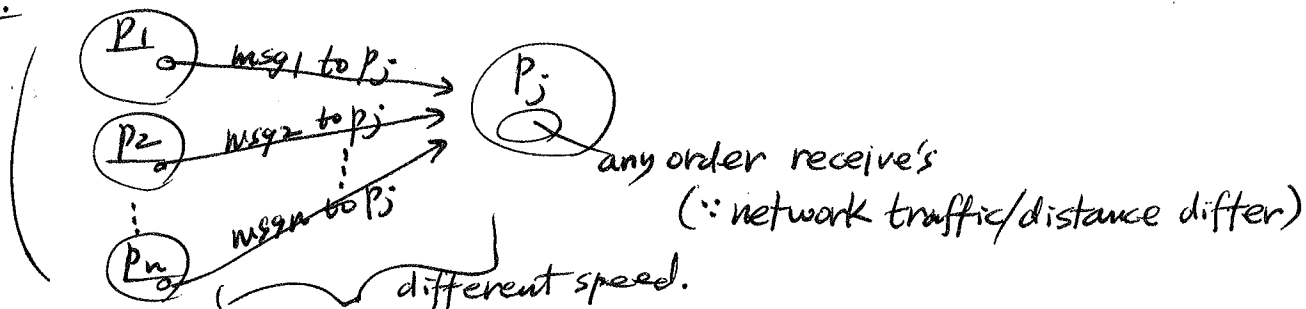
— Semantics of <u>MPI_Send / MPI_Recv</u>

process_i **Send** → buffered → process_j receive (always blocked)

buffered → ↓ and continue

Send → or, → blocked until actual transmit → process_j receive

— hybrid scheme:
   using <u>cut-off</u>, ⎡ msg_size < cut-off
                          ⇒ buffered
                      else
                          ⇒ blocked

— Non_overtaking (used in MPI) — between 2 processes

<u>process_i</u>
  ⎡ send msg1 to Pj ⟶ receive msg1
  ⎣ send msg2 to Pj ⟶ receive msg2   } keep same order as sender process.

<u>process_j</u>
  — receive msg1
  — receive msg2

vs.

P1 → msg1 to Pj ⟶
P2 → msg2 to Pj ⟶   Pj
Pn → msgn to Pj ⟶   any order receive's
                    (∵ network traffic/distance differ)
(different speed.)

— Trapezoid computation in <u>MPI</u>

<u>review</u>



$\left( h = \dfrac{b-a}{n} \right)$

$n$ trapezoids



$a+h \quad a+2h \quad a+3h \cdots \quad a+(n-1)h$

— local integral

$$= \left( \frac{f(x_i) + f(x_i+1)}{2} \right) * h$$

— global integral

$$= \left( \frac{f(x_0) + f(x_1)}{2} \right) * h + \left( \frac{f(x_1) + f(x_2)}{2} \right) * h + \cdots + \left( \frac{f(x_{n-1}) + f(x_n)}{2} \right) * h$$

$$= \frac{h}{2} \left[ (f(x_0) + f(x_1)) + (f(x_1) + f(x_2)) + (f(x_2) + f(x_3)) + \cdots + (f(x_{n-1}) + f(x_n)) \right]$$

$$= h \left[ \frac{f(x_0)}{2} + f(x_1) + f(x_2) + \cdots + f(x_{n-1}) + \frac{f(x_n)}{2} \right]$$

<u>Serial code</u>

```
h = (b-a) / n ;
approx = (f(a) + f(b)) / 2 ;

for(i = 1 ~ n-1)
   ( xi = a + i*h ;
   ( approx += f(xi) ;

approx = approx * h ;
```

— Trapezoid — parallel code (MPI) — msg passing

**Concept**
- divide $n$ into $\frac{n}{P}$
- each process computes $\frac{n}{P}$ trapezoids
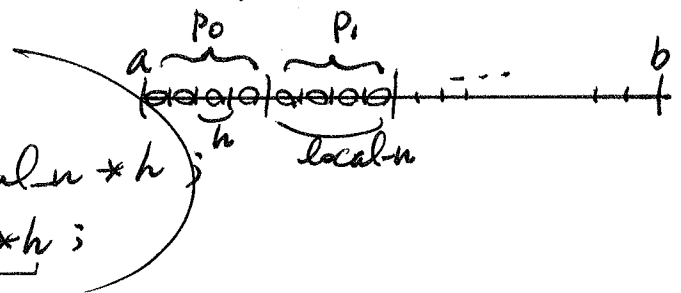- process_$\phi$ combines local integrals

get a, b, n ; { using only 1 process, get inputs and send them to all other processes

$h = (b-a) \times n$ ;

local_n = n/p ;
local_a = a + my_rank * local_n * h ;
local_b = local_a + local_n * h ;



local_integral = Trap(local_a, local_b, local_n, h) ;

if (my_rank $\neq \phi$)
    send local_integral to process_$\phi$ ;
else
    tot_integral $\leftarrow$ local_integral (of $P_0$) ;
    for (i=1 ~ p-1)
        Receive local_integral from Process_i ;
        tot_integral += local_integral ;
    display tot_integral :

— I/o for MPI program (suggestion)
- any process accesses std_out
- only process_$\phi$ accesses std_in

— Collective Comm  ($\neq$ P2P — send/receive)

Communication function involves <u>all the processes</u>
in a communicator.

ex) global sum computation

— <u>MPI_Reduce</u> (& local_integral,  ——— input data  ⎫
&tot_integral,  ——— output data ⎬ must be
⎭ different
1,  ——————— count
MPI_Double,  ——— data type
MPI_Sum,  ——— see Table 3.2
Ø,  ——————— dest rank ⎡ MPI_MAX,
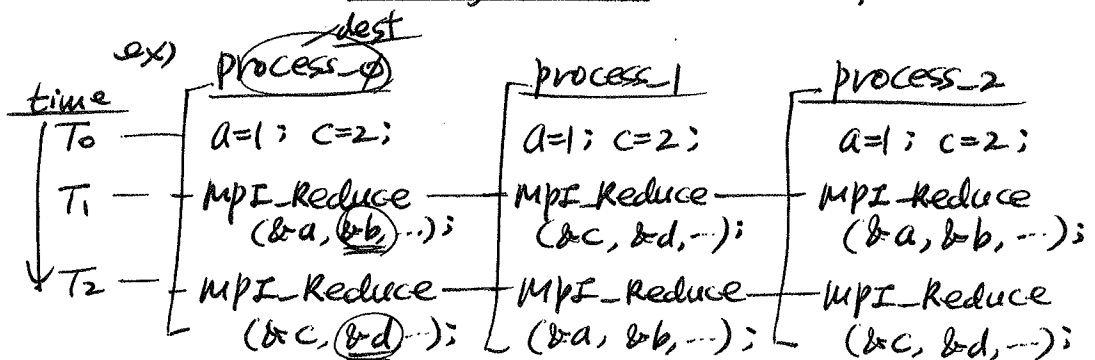MPI_Comm_World);  ⎢ MPI_MIN
⎢ MPI_prod
⎢ MPI_sum
⎢ :

1. all processes must call the same
collective func.
ex) one process calls MPI_Recv(--) ⎫→ crash
others call MPI_Reduce(--) ⎭

2. <u>arguments</u> to collective func. must be compatible
from processes

3. output data arg. is only used in dest-process;
others should use <u>NULL</u>.

4. # tags and <u>calling order</u> is used for matching

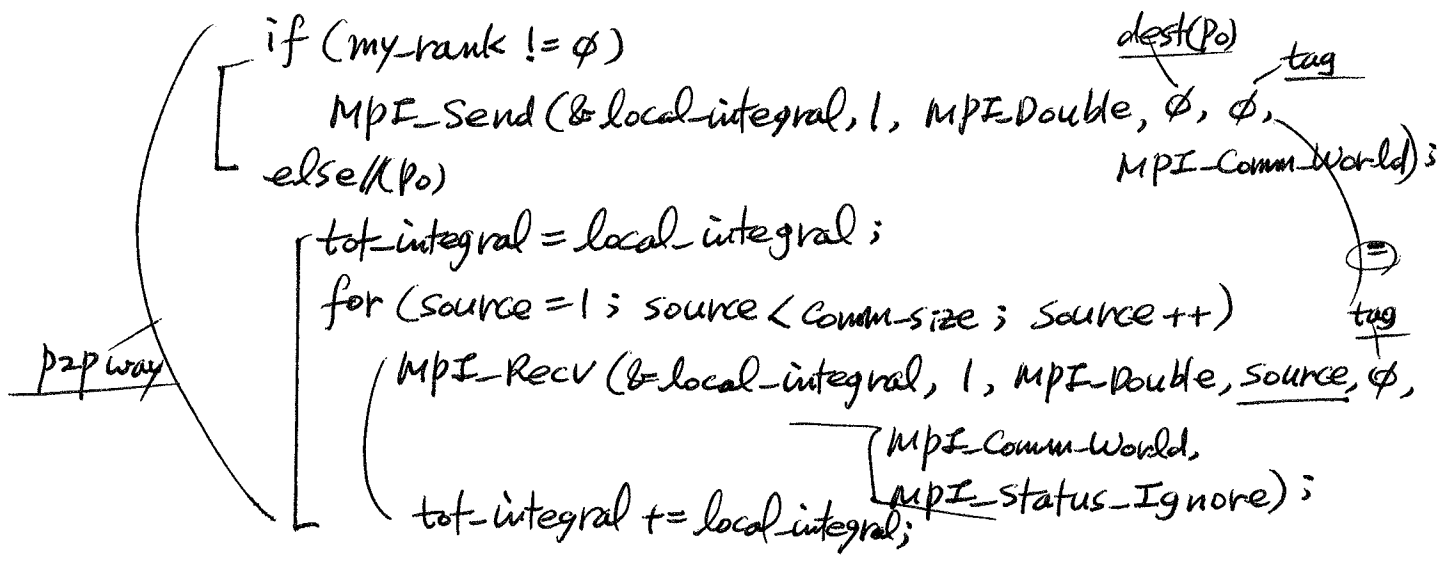ex)  process_0 (dest)      process_1         process_2

| time | process_0 | process_1 | process_2 |
|---|---|---|---|
| T0 | a=1; c=2; | a=1; c=2; | a=1; c=2; |
| T1 | MPI_Reduce (&a, &b, ..); | MPI_Reduce (&c, &d, --); | MPI_Reduce (&a, &b, --); |
| T2 | MPI_Reduce (&c, &d ..); | MPI_Reduce (&a, &b, --); | MPI_Reduce (&c, &d, --); |

Assume, MPI_SUM

⎡ expected output: b=3, d=6
⎣ actual output: b=1+2+1=④, d=2+1+2=⑤

# ⊕ Collective Communication

## — MPI_Reduce

ex) trapezoid computing

p2p way
```
if (my_rank != ∅)
    MPI_Send(&local_integral, 1, MPI_Double, ∅, ∅,
                                        dest(P0)        tag
                                        MPI_Comm_World);
else(P0)
    tot_integral = local_integral;
    for (source = 1; source < comm_size; source++)
        MPI_Recv(&local_integral, 1, MPI_Double, source, ∅,
                                        MPI_Comm_World,        tag
                    tot_integral += local_integral;  MPI_Status_Ignore);
```
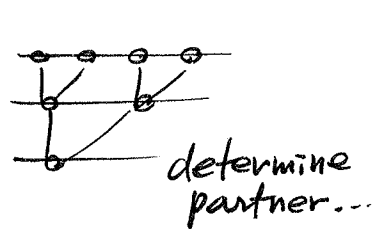
→ it looks like:



but, P0 performs serial
    tot_int ⊖ local_int



⟹ better way is using __full-tree__ reduction.

__MPI_Reduce__(&local_integral, &tot_integral, 1, MPI_Double,
                                        MPI_Sum, ∅, MPI_Comm_World);
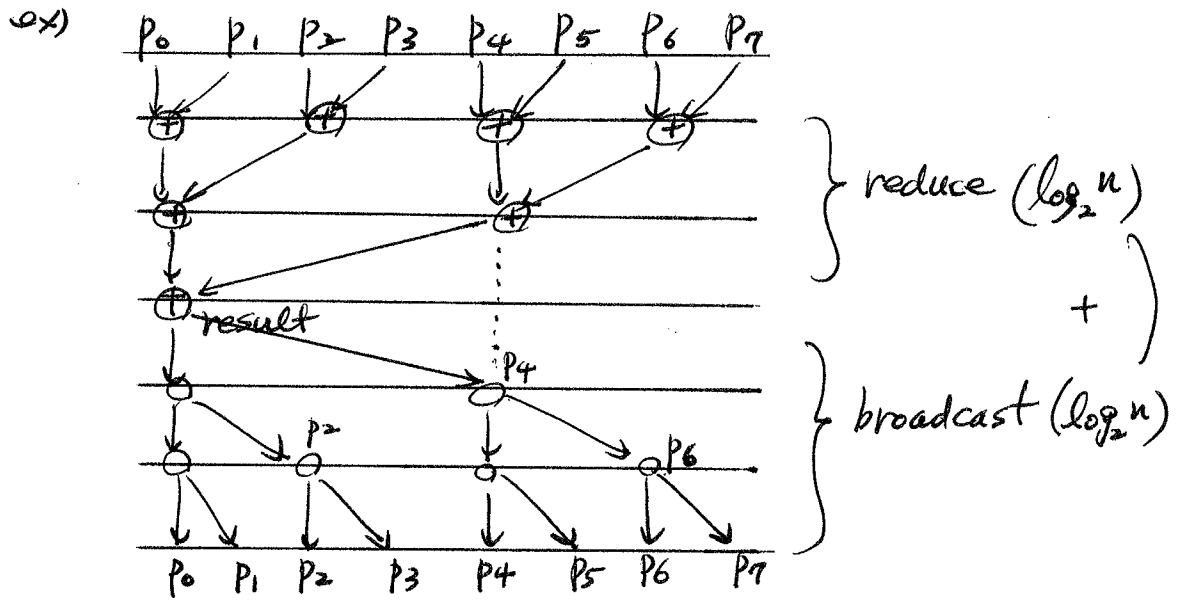                                        reduction    dest
                                        operation    P0
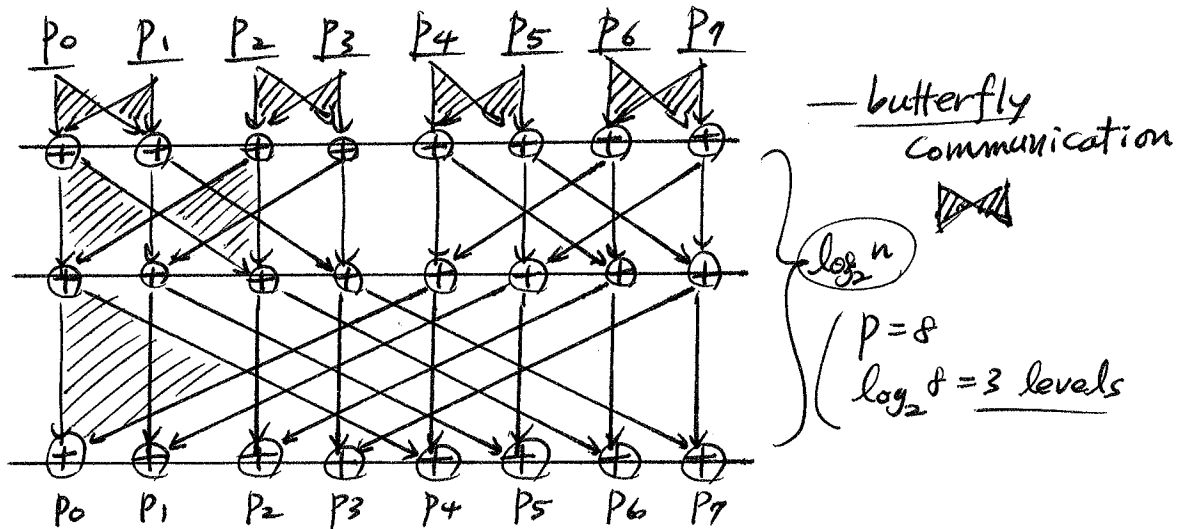


determine
partner...

— <u>MPI_Allreduce</u>  { useful for the situation that
all processes need the reduced result.
(for further operation)

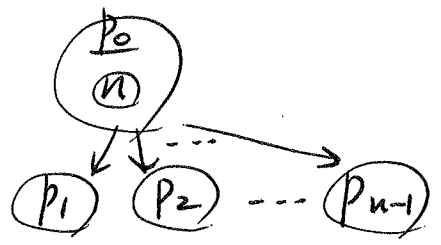— naive method: <u>reduce</u> to $P_0$ and <u>broadcast</u> to all

ex)



} reduce $(\log_2 n)$

$+$

} broadcast $(\log_2 n)$

ex) trapezoid computing

<u>MPI_Allreduce</u> (&local_integral, &tot_integral, 1, MPI_Double,
MPI_SUM, MPI_Comm_World);

# dest_process (all processes are
destination)



— <u>butterfly</u>
communication

} $\log_2 n$

{ $p = 8$
$\log_2 8 = 3$ <u>levels</u>

# — MPI_Bcast

ex) $P_0$ accesses input $n$ and distributes it to all others



— in fact, this is serial in p-2-p comm.
with for-loop way, i.e.,

send to p1;
send to p2;
⋮
send to $P_{n-1}$;

### P2p way

```
if (my_rank ==∅) //Po
    n= atoi (argv[1]);
    for(dest=1 ; dest<comm_sz; dest++)
        MPI_Send (&n, 1, MPI_Int, dest, ∅, MPI_Comm_World);
                                        tag
else //other processes
    MPI_Recv (&n, 1, MPI_Int, ∅, ∅, MPI_Comm_World);
                              Source  tag
                              Po
```

vs.

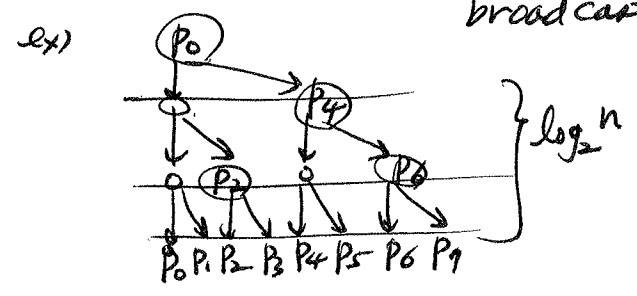### broadcast way

```
if (my_rank ==∅)
    n= atoi (argv[1]);
MPI_Bcast (&n, 1, MPI_Int, ∅, MPI_Comm_World);
```

in para — for source process
out para — for other processes

Source Po

※ Should Call this out of Po
∵ ( Po — sends
    others — receive

— MPI_Bcast performs **full tree** broadcasting.

ex)



$\log_2 n$

P0 P1 P2 P3 P4 P5 P6 P7

— <u>Data distribution</u>

      ex) $p=3$

    — block partition —   $\underbrace{0,1,2,3}_{p_0}, \underbrace{4,5,6,7}_{p_1}, \underbrace{8,9,10,11}_{p_2}$

    — cyclic partition :
$$\begin{cases} p_0 : 0,3,6,9 \\ p_1 : 1,4,7,10 \\ p_2 : 2,5,8,11 \end{cases}$$



    — block-cyclic partition
$$\begin{cases} p_0 : \underline{0,1}, \underline{6,7} \\ p_1 : \underline{2,3}, \underline{8,9} \\ p_2 : \underline{4,5}, \underline{10,11} \end{cases}$$



— <u>MPI_Scatter</u>

    Assume: $n$ is evenly divisible by $p$, and <u>block partition</u>.

    — process_0 receives (accesses) size $n$ vector, and

        send the $\begin{cases} \text{1st } \frac{n}{p} \text{ elements to } p_0 \\ \text{2nd } n/p \text{ elements to } p_1 \\ \text{3rd } n/p \text{ elements to } p_2, \text{ and so on} \end{cases}$

             ...

  ex) if (my_rank == 0) // $p_0$      (ref: not evenly divisible case,
                               use MPI_Scatterv )

    { int* a = new int[n];
      get input values for a[n] array (vector); (receiving buffer size

      <u>MPI_Scatter</u> (a, local_n, MPI_Int, local_a, local_n,

(p_0)      (sending     (sending      type     (receiving
            buffer*      buffer size   (sender)    buffer*
                    —local size—            —int local_a[]
                    to each process

      delete[] a;             MPI_Int, 0, comm);
    }
                        type      source $p_0$
  else                    (receiver)

<u>others</u>   MPI_Scatter (a, local_n, MPI_Int, local_a, local_n,
           sending buffer info    (MPI_Int, 0, comm);
                             receiving buffer

— Performance/time checking

— time checking for only computation part (not for elapsed time)

```
#include "timer.h"  ──→ macro def
    ⋮
double start, finish;
    ⋮
┌ Get_Time (start);
│    ⋮  ////
└ Get_Time (finish);
```

↳ display (finish – start); — in $\mu$ seconds

macro def:
```
#define Get_Time (now) \
{ struct timeval t; \
  gettimeofday (&t, NULL); \
  now = t.tv_sec + \
        t.tv_usec/1000000.0; }
```

( posix lib. func )

needs <sys/time.h>

— MPI supports MPI_Wtime ()

```
double start, finish;
┌ start = MPI_Wtime ();
│    ////
└ finish = MPI_Wtime ();
```
} each process reports exec time.

↳ display my_rank, (finish – start); — in $\mu$ seconds

— parallel time checking — slowest process's exec. time

```
┌ MPI_Barrier (MPI_Comm_World);
│ ┌ local_start = MPI_Wtime ();
│ │    ////
│ └ local_finish = MPI_Wtime ();
│ local_elapsed = local_finish – local_start;
│ MPI_Reduce (&local_elapsed, &elapsed, 1, MPI_Double,
│             MPI_Max, 0, MPI_Comm_World);
│
│ If (my_rank == 0)
└      ↳ display elapsed;
```

P_0  P_1  P_2 ... barrier

parallel elapsed time (longest)

( reduction operation )   dest = P_0

— Running MPI program on a hybrid system, each node is multicored.
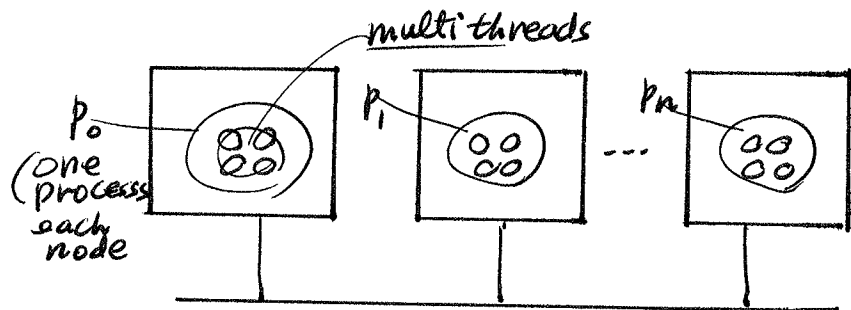
— by default, only 1 process works on each node.

$\Rightarrow$ MPI + OpenMP model

in each <u>node</u> (multicored)
1 process

ex) pragma omp parallel for ---

— How to compile/run?

$> mpic++ -fopenmp p1.cpp ↵
$> mpirun -n 4 ./a.out ↵

multi threads

$P_0$
(one process each node)

$P_1$

$P_n$

...

ex) #include <mpi.h>
    #include <omp.h>
    ---

    MPI_Init (NULL, NULL); //MPI starts
    MPI_Comm_size (MPI_Comm_World, &comm_sz); //get #process
    MPI_Comm_rank (MPI_Comm_World, &my_rank);
    //get my rank
    mpi operations

    #pragma omp parallel for num_threads(4) ---
        omp operations

    mpi operations.