# COLLECTIVE COMMUNICATION ①

issues with the trap. rule program: (on 8 procs)

(a) procs 1-7 are idle, while proc 0 performs I/O

(b) after proc 0 has collected data (INPUT), higher rank procs ~~×~~ continue to wait, until proc 0 sends input data to lower rank procs.

(c) proc 0 collects all partial answers and performs addition

main point of parallel computing: multiple procs collaborate on solving a problem

## TREE-STRUCTURED comm/tion

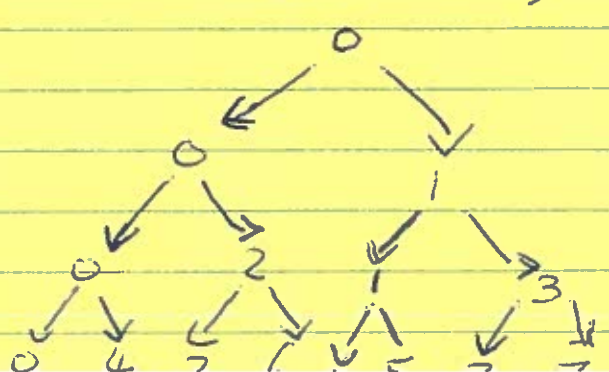focus on the distribution of the input data

divide the work more evenly among procs

imagine that we have a tree of procs with 0 at the root.

| | | |
|---|---|---|
| 1st stage | 0 sends the data to 1 | |
| 2nd " | 0 ~> 2 & 1 ~> 3 | |
| 3rd " | 0 ~> 4 & 1 ~> 5, 2 ~> 6, 3 ~> 7 | |

we reduced the input distrib. loop, from 7 stages ~~×~~ to 3 stages

if we have p proc's, the input data is distributed in $\lceil \log_2 p \rceil$ stages (rather than p-1)

case study: $\log_2 8 = 3$ (previous example)

$\log_2 1024 = 10$ (reduction by a factor of 100)

modify Get_data fct to use a tree-structured distribution scheme:

```
for (stage = first ; stage <= last ; stage++)
    if (I_receive(stage, my-rank, &source))
        Receive (data, source);
    else if (I_send (stage, my-rank, p, &dest))
        Send (data, dest);
```

I_receive fct returns $\begin{cases} 1, \text{ calling process receives data during current stage} \\ 0, \text{ otherwise} \end{cases}$

if the calling proc. receives data, the parameter "source" is used to return the rank of the sender.

I_send fct returns $\begin{cases} 1, \text{ proc send's during current stage} \\ 0, \text{ otherwise} \end{cases}$

implementation: we need to calculate
(1) whether a proc. receives, and the source
(2)    "        "    "  sends,    "    "  destination

( several possible tree-schemes are
possible. <u>no canonical choice of order</u>ing
deciding the best scheme requires
some knowledge of the <u>topology</u>
of our system.

<span style="color:red">stages numbered starting at 0,...</span>

<u>general tree scheme:</u>

if $2^{stage} < $ my-rank $< 2^{stage+1}$
then I-receive from proc my-rank$-2^{stage}$

if my-rank $< 2^{stage}$
then I-send to proc my-rank$+2^{stage}$

int I-receive ( $\quad$ ( 2nd version of <u>Get-data fct</u> )

$\qquad$ int stage
$\qquad$ int my-rank $\qquad$ ] IN PARAMS
OUT PARAM [ int* source-ptr ) { $\qquad$ $2^{stage}$
$\overbrace{\qquad}$
$1 <<$ stage

int power_2_stage; $\qquad$ power_2_stage =
if ((power_2_stage $<=$ my-rank)
$\&\&$ (my-rank $<$ 2*power_2_stage))
{ *source-ptr = my-rank - power_2_stage;
return 1;
}
else return 0;
}

int <u>I-send</u> (int stage, int my-rank, int p, $\overbrace{\qquad}^{IN\ PARAMS}$
... $\qquad$ int* dest-ptr) PARAM OUT

```
void Send (
                    float a,        ⌉ 4
                    float b',       ⌉ ⟶ IN
                    int n',         ⌉    PARAMS
                    int dest ) {
```

```
MPI_Send (&a, 1 MPI_FLOAT, dest, 0 MPI_C_W)
MPI_Send (&b', 1'    ,,      ,,  , 1',  ,,  );
MPI_Send (&n', 1' MPI_INT, ,, , 2', ,, );
}
}
```

```
    void Receive (
                    float* a_ptr,    ⌉ OUT
                    float* b_ptr,    ⌉ PARAMS
                    int*   n_ptr,    ⌋
                    int    source  ⟶ IN PARAM
```

```
                              OUT PARAMS
                      ⌈ float*   float*  int* ⌉
void Get_data1 ( a_ptr, b_ptr, n_ptr,
                int my_rank, int p ) {
                      IN PARAMS
```

```
    int source, dest, stage;

    if (my_rank == 0) { scanf (a_ptr, b_ptr,
                                n_ptr);
    }

    for (stage = 0; stage < log2p; stage++)
        if ( I_receive (stage, my_rank, &source))
            Receive (a_ptr, b_ptr, n_ptr, source)
    ? else if I_send (...)
        Send (a_ptr, *b_ptr, *n_int, dest
```

# BROADCAST

A communication pattern that involves all the proc^s in a comm/tor is a <u>collective communication</u>.

A <u>broadcast</u> is a collective comm/tion in which a single process sends the same data to every process in the comm/tor.

int MPI_Bcast (

[ tree-structured com/tion is much more efficient than a sequential from the root node

IN on root proc
OUT or other procs ⟶ ( IN OUT PARAM ) ~> void * message,

IN PARAMS [ int count,
MPI_Datatype datatype,
int root,
MPI_Comm comm )

(hand-coded is effici? )

send a copy of the data in "message" on the proc with rank "root" to each process in the comm/tor "comm".

should be called by ALL proc^s in the comm/tor, with the same arg^s for root/comm.

count/datatype: specify how much memory is needed for the message.

⌐⟶ these two parameters should be the same on all the proc^s of the comm/tor.

( no tags in collective comm/tions mechanism )

the reason is that in some cases a single proc receives data from many

```
void Get_data2( a_ptr, b_ptr, n_ptr,
                               my_rank)
    {
        if (my_rank ==0) {  scanf...}
```
*zero also executes this part*
```
        MPI_Bcast (a_ptr, ! MPI_FLOAT 0 MCW);
           "      "   (b_ptr, ", "      ,'0,' ", );
           "      "   (n_ptr, ",MPI_INT, 0,' ", );
    }
```

faster & more easily comprehensible
version.

TAGS, SAFETY, BUFFERING, SYNCHRONIZATION

MPI_Bcast  does not use tags,  WHY?

MPI_Send /MPI_Recv
use tags : proc A sends several msgs
                to proc B, and B handles
                  them according to their
example1                                      tags.
consider this seq. of events

| TIME | Proc A | Proc B |
|------|--------|--------|
| 1 | MPI_Send to B | local work |
| 2 | tag=0 MPI_Send to B | local work |
| 3 | tag=1 local work | MPI_Recv from A tag=1 |
| 4 | local work | MPI_Recv from A tag=0 |

this sequence requires <u>buffering</u>

(set aside memory for storing msg's,
   before a "receive" has been executed)

msg. env. ~> {rank of sender/receiver}
             {tag, comm/tor              }

until B calls MPI_Recv, the system
does not know where the msg that A
is sending should be stored.

when B calls MPI_Recv, the system
looks for any buffered msg's that
has an env. that matches the recv. param's.
if there is no such message,
   then it will wait until one arrives.

if no buffering is available A cannot
   send data to B, until it knows that
   B is ready to receive.

send cannot complete until receiver is ready to receive ‖ <u>send uses synchronous mode</u>

if a program that assumes buffering is
available, is run on a system that
does not provide buffering ~> DEADLOCK

<u>UNSAFE MPI Program</u>

A hangs while it waits for B to receive
               1st send

B hangs while it waits for A to execute
               2nd send

# example 2 (uses Bcasts) ( ~ ~ 1→ LOCAL WORK )

| TIME | Proc A | Proc B | Proc C |
|------|--------|--------|--------|
| 1 | MPI_Bcast &x | ~~~ | ~~~ |
| 2 | MPI_Bcast &y | ~~~ | ~~~ |
| 3 | ~~~ | MPI_Bcast &y | MPI_Bcast &x |
| 4 | ~~~ | MPI_Bcast &x | MPI_Bcast &y |

suppose A broadcasts two floats x,y
in (that x=5, y=10 ) to B,C
( on proc A )

when bcasts are completed on all
3 proc$^s$, x=5, y=10 on proc$^s$ A, C

but on proc B, x=10 // reversed
            y=5  // values WHY?

⚠️ first parameter of Bcast is IN/OUT

broadcasts assumed synchorization:
on a given process, the bcast would
not return until every process had
received the bcast data.

this restriction is relaxed, when buffering
is available, A can complte its bcasts,
before B, C 'begin their bcast calls

BUT the EFFECT in terms of data communicated, must be
the same as if there was sychronization.
so 1st bcast on B matches 1st BCAST on A and stores 5 in y