

Lec 14

- MPI Blocking Collectives

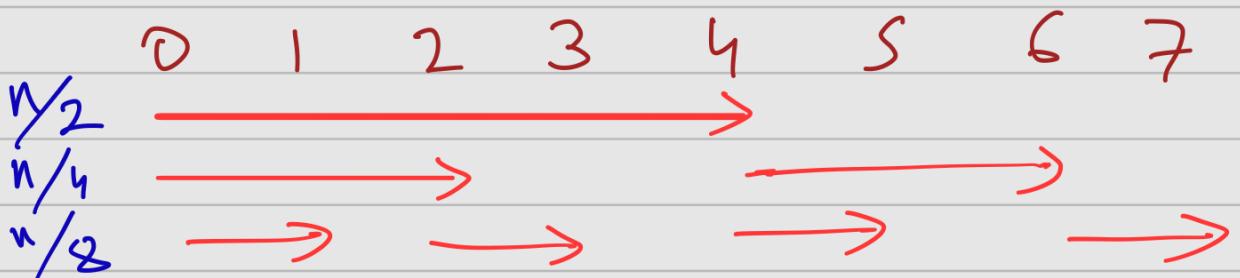
Barrier, Reduce , Gather, Scatter,
 Broadcast, Allgather, Allreduce, Alltoall

→ Scatter

Recursive Halving
 Vector Halving
 Distance Halving

(finally
 each
 process
 needs
 n/p)

$$\log p \times L + (p-1) \times \left(\frac{n}{p}\right) \times \left(\frac{1}{B}\right)$$



→ All Gather

① Ring

- every comm b/w direct neighbour
- first cycle: left nei \rightarrow right nei apne data
 second cycle: left nei \rightarrow right nei ko
 ... also left i ka data

n Bytes P processes \Rightarrow

$\frac{n}{P}$ bytes

$$(P-1) \left(L + \frac{n}{P} \times \frac{1}{B} \right)$$

send/recv every time

(2) Recursive Doubling

$(2^{k-1}) \times \frac{n}{P}$ Bytes at k^{th} step.

(Data size, dist double every step)

0 ↔ 1	2 ↔ 3	4 ↔ 5	6 ↔ 7
0 ↔ 2	1 ↔ 3	4 ↔ 6	5 ↔ 7
0 ↔ 4	1 ↔ 5	2 ↔ 6	3 ↔ 7

$$\log P \times L + (P-1) \left(\frac{n}{P} \right) \times \frac{1}{B}$$

\rightarrow Allgather Short msg = 80 KB
Long msg = 512 KB

$S < 80\text{KB}$

$M < 512\text{KB}$

	Power of 2	Not Power of 2
S	Recursive Doub	Bruck
M	Recursive Doub	Ring
L	Ring	Ring

→ Broadcast — Van de Geijn

$$\text{Binomial} : \log p \times \left(L + \frac{n}{B} \right)$$

Scatter + Allgather

$$\hookrightarrow \log p L + \frac{(p-1)}{p} \frac{n}{B}$$

$$\text{Allg (ring)} : (p-1) L + (p-1) \frac{n}{pB}$$

$$\text{Allg (rec'd)} \quad \log p L + (p-1) \frac{n}{pB}$$

Bcast - Short - Msg - Size (x)
Bcast - Long - Msg - Size (y)

S < X = binomial

M = scatter + allg (rec'd)

L > Y = scatter + allg (ring)

→ Reduce

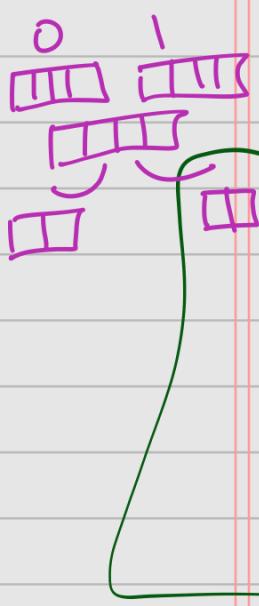
$$\text{Rec'Doub} \Rightarrow \log p \left(L + \frac{n}{B} + n \times C \right)$$

- each process has 'n' Bytes data

compute
cost
per
Byte

(2) Raben Seifner

Recursive Vector Halving & Dist Doubling



then Gather

$$\log P \times L + \frac{(P-1)}{P} \frac{n}{B} + \frac{(P-1)}{P} \times n \times c \quad (\text{reduce + scatter})$$
$$+ \log P \times L + \frac{(P-1)}{P} \times \frac{n}{B} \quad (\text{gather using binomial})$$

→ at each step processes P, Q communicate
reduce first half of data at P
and second half of data at Q.

at K^{th} step, each process has $\frac{n}{2^K}$ data
after $\log P$ steps, gather at root.

$$\begin{array}{llll} 0 \leftrightarrow 1 & 2 \leftrightarrow 3 & 4 \leftrightarrow 5 & 6 \leftrightarrow 7 \\ 0 \leftrightarrow 2 & 1 \leftrightarrow 3 & 4 \leftrightarrow 6 & 5 \leftrightarrow 7 \\ 0 \leftrightarrow 4 & 1 \leftrightarrow 5 & 2 \leftrightarrow 6 & 3 \leftrightarrow 7 \end{array}$$

(like all gather)

→ All reduce (Raben Seifner)

reduce - scatter using recursive vector halving and dist doubling

then

all gather using vector doubling, dist

$$\log p \times L + \frac{p-1}{P} \times \frac{n}{B} + \frac{p-1}{P} \times n \times c$$

$$+ \log p \times L + \frac{p-1}{P} \times \frac{n}{B}$$

- Short / Medium :

Reduce (RecDoub) + Bcast (Binomial)

$$\log p \times L + \log p \times \frac{n}{B} + \log p \times n \times c +$$

$$\log p \times L + \log p \frac{n}{B}$$

Long

Reduce Scatter + allgather (see doub.)

$$\log p \times L + \frac{p-1}{P} \times \frac{n}{B} + \frac{p-1}{P} \times n \times c$$

$$+ \log p \times L + (p-1) \frac{n}{P} \times \frac{1}{B}$$

Lee 15

• col-wise chunk \Rightarrow col wise cutters ||

- row-wise decomp \Rightarrow rows are cutted =
decomp along x axis \Rightarrow col wise decomp
" " y axis \Rightarrow row wise "

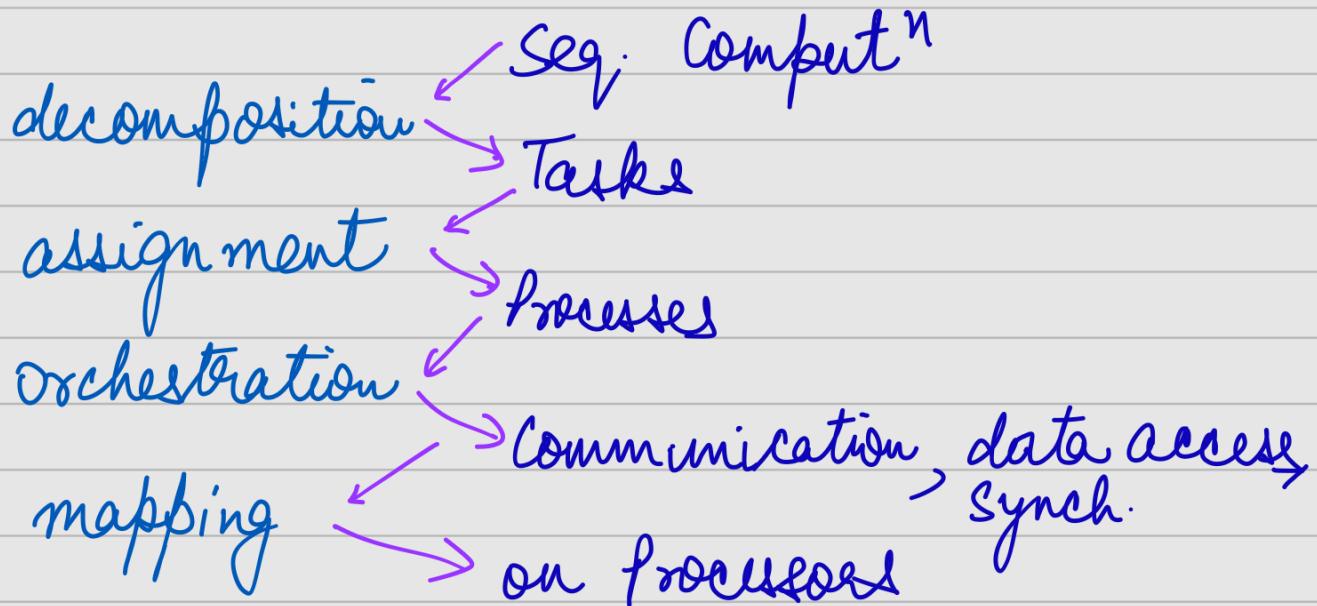
Gatherv - recvcounts [i], displ [i]

Scatterv - sendcount [i], displ [i]

Allgatherv - recvct [i], displ [i]

Allto allv - send ct [i], displs [i]
recvct [i]; displs [i]

Lec 16



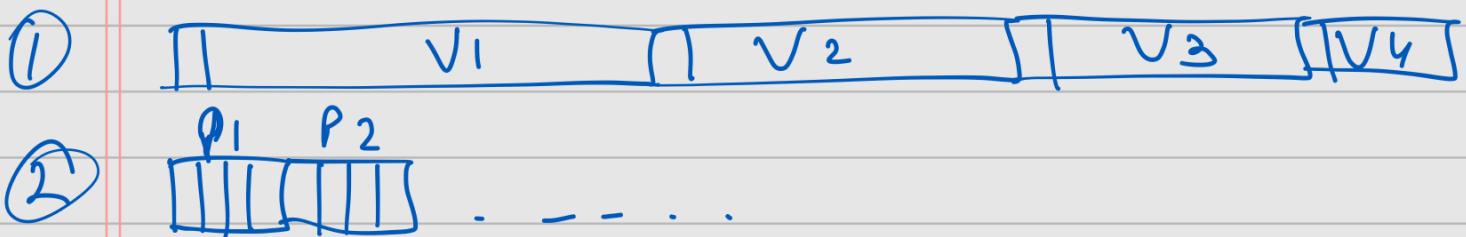
\rightarrow Performance Goals

i) Expose concurrency

- ii) Reduce interprocess comm.
 - iii) Load - balance
 - iv) Reduce synchronisation
 - v) Reduce idling
 - vi) Reduce management overhead
 - vii) Preserve data locality
 - viii) Exploit network topology
-

Lec 17

$M \times N \rightarrow 4$ vars each, $P \times Q$ processes



MPI_File_open (CommWorld, filename, create+rdwr,
Null, fh)

for ($i = 0$ to 4) {
 MPI_Offset

$$\text{offset}[i] = \left\lceil \left(\frac{\text{rank}}{Q} \right) \times Q + (\text{rank} \% Q) \right\rceil \times$$

$i \times M \times N \times$
 sizeof(float)

$\frac{M \times N}{2} \times \text{sizeof(float)}$

$P \times Q$

MPI-File-Seek (fh, offset[i], \rightarrow)
 MPI-File-write (fh, buf[i], \rightarrow) $\Bigg\}$

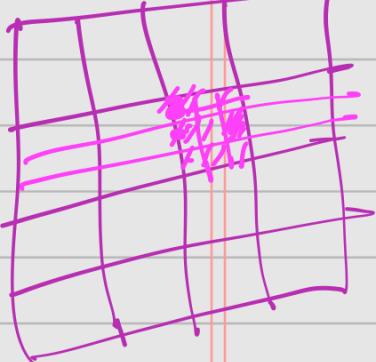
$(0 \ 1 \ 2 \ 3 \ 4) \ Q$
 $\times M/P$

for ($i = 0$ to 4) { for ($l:idx = 0$ to $\frac{M \times N}{P \times Q}$) {

$off[i][idx] = \left[\begin{array}{l} \frac{(rank)}{Q} \times \left(\frac{M \times N}{P} \right) \times \text{sizeof (float)} \\ + \frac{(idx)}{N/Q} \times N \times \text{sizeof (float)} \\ + ((rank) \% Q) \times \left(\frac{N}{Q} \right) \times (\text{float}) \\ + ((idx) \% Q) \times (\text{float}) \\ + i \times \left(\frac{M \times N}{P} \right) \times (\text{float}) \end{array} \right]$

$idx \% \left(\frac{N}{Q} \right)$

Fwrite ($off[i][idx]$, 1 , float);



0	1	2	3	4
5	6	7	8	9
10	11	12	13	14
10	11	12	13	14

0	1	2	3	4	0	1	2	3	4
5	6	7	8	9	5	6	7	8	9
10	11	12	13	14	10	11	12	13	14
10	11	12	13	14	10	11	12	13	14

$off[i][M/P]$

$$off[0][0] = \left(\frac{(rank)}{Q} \times \frac{M \times N}{P} + \frac{(rank \% Q)}{N/Q} \times N \right)$$

$+ N \times f$

$\rightarrow + NX M \times f$

Part 2 :

for ($id = 0$ to $\frac{M \times N}{P \times Q}$) {

 offset =
$$\left(\frac{\text{rank}}{Q} \times \frac{M \times N}{P} + (\text{rank} \% Q) \times \frac{N}{Q} \right)$$

$$+ \left(\left(id / \left(N / Q \right) \right) \times N + \left(id \% \left(\frac{N}{Q} \right) \right) \right)$$

 for ($i = 0$ to 4)

 fwrite_at (fh, offset + $i \times \text{sizeof}(\text{int})$, var [id][i],)