## Week 4 – Collective communication, topologies

- Main topics this week
  - Collective communication, topologies
  - Project match making
- Reading
  - Pacheco, Chapters 4 and 5

Research Computing @ CU Boulder        Week 4 - Collective Communication    1    2/8/12

# Lecture 7

Collective Communication and block partitioning

Research Computing @ CU Boulder

## MPI collective communication operations

- MPI_BCAST
- MPI_GATHER
- MPI_SCATTER
- MPI_REDUCE

Arguments similar to those of  MPI_SEND and MPI_RECV
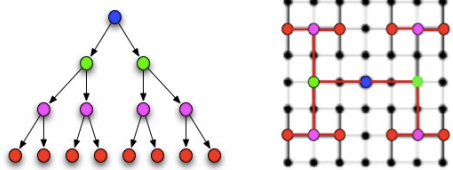
Research Computing @ CU Boulder        Week 4 - Collective Communication    3    2/8/12

## Implement by embedding virtual binary tree into actual topology

broadcast or scatter

## Other operations similarly based on trees

- Broadcast and scatter go from root to leaves
- Gather and reduce go from leaves to roots, reduce operates on data while gather just collects it

- Other ops like MPI_ALLGATHER: all to all, communication goes in one direction then other

- Next week's lab assignment:   butterfly allreduce
  - Topology is butterfly
  - Communication goes in both directions

## All that said

- You probably won't beat MPI_ALLGATHER which uses a better communication algorithm than plain old MPI_BCAST

But you can try!

## Consider an algorithm that requires collective communication

$\alpha = x^T y$

$x = ( x_1, \ x_2, \ \dots, \ x_n )$

$y = ( y_1, \ y_2, \ \dots, \ y_n )$

We have p processors, how to implement?

   first, assume p divides n

      how do we divide the work?

      how do we accumulate the result?

      does replicating data make sense?

Research Computing @ CU Boulder      Week 4 - Collective Communication    7    2/8/12

## Data Decomposition Options

- Interleaved (cyclic)
  - Easy to determine "owner" of each index
- Block
  - Balances loads
  - More complicated to determine owner if $n$ not a multiple of $p$

Research Computing @ CU Boulder      Week 4 - Collective Communication    8    2/8/12

## Block Decomposition Options

- Want to balance workload when $n$ not a multiple of $p$
- Each process gets either $\lceil n/p \rceil$ or $\lfloor n/p \rfloor$ elements
- Seek simple expressions
  - Find low, high indices given an owner
  - Find owner given an index

Research Computing @ CU Boulder      Week 4 - Collective Communication    9    2/8/12

## Method #1

- Let $r = n \bmod p$
- If $r = 0$, all blocks have same size
- Else
  - First $r$ blocks have size $\lceil n/p \rceil$
  - Remaining $p\text{-}r$ blocks have size $\lfloor n/p \rfloor$

## Examples

17 elements divided among 7 processes

17 elements divided among 5 processes

17 elements divided among 3 processes

## Method #1 Calculations

- Indexing starts with 0
- First element controlled by process $i$
  $$i\lfloor n/p \rfloor + \min(i,r)$$
- Last element controlled by process $i$
  $$(i+1)\lfloor n/p \rfloor + \min(i+1,r) - 1$$
- Process controlling element $j$
  $$\min\left(\lfloor j/(\lfloor n/p \rfloor + 1)\rfloor, \lfloor (j-r)/\lfloor n/p \rfloor \rfloor\right)$$

## Method #2

- Scatters larger blocks among processes
- First element controlled by process $i$

$$\lfloor in/p \rfloor$$

- Last element controlled by process $i$

$$\lfloor (i+1)n/p \rfloor - 1$$

- Process controlling element $j$

$$\lfloor p(j+1)-1)/n \rfloor$$

Research Computing @ CU Boulder  Week 4 - Collective Communication  1 3  2/8/12

## Examples

17 elements divided among 7 processes

17 elements divided among 5 processes

17 elements divided among 3 processes

Research Computing @ CU Boulder  Week 4 - Collective Communication  1 4  2/8/12

## Comparing Methods

Our choice

| Operations | Method 1 | Method 2 |
|---|---|---|
| Low index | 4 | 2 |
| High index | 6 | 4 |
| Owner | 7 | 4 |

Assuming no operations for "floor" function

Research Computing @ CU Boulder  Week 4 - Collective Communication  1 5  2/8/12

## Pop Quiz

- Illustrate how block decomposition method #2 would divide 13 elements among 5 processes.

13(0)/ 5 = 0   13(2)/ 5 = 5   13(4)/ 5 = 10

13(1)/5 = 2     13(3)/ 5 = 7

## Block Decomposition Macros

```
#define BLOCK_LOW(id,p,n)   ((i)*(n)/(p))

#define BLOCK_HIGH(id,p,n) \
        (BLOCK_LOW((id)+1,p,n)-1)

#define BLOCK_SIZE(id,p,n) \
        (BLOCK_LOW((id)+1)-BLOCK_LOW(id))

#define BLOCK_OWNER(index,p,n) \
        (((p)*(index)+1)-1)/(n))
```

## Local vs. Global Indices

L 0 1
G 0 1

L 0 1 2
G 2 3 4

L 0 1
G 5 6

L 0 1 2
G 7 8 9

L 0 1 2
G 10 11 12

## Looping over Elements

- Sequential program
```
for (i = 0; i < n; i++) {
    ...
}
```

Index $i$ on this process…

- Parallel program
```
size = BLOCK_SIZE (id,p,n);
for (i = 0; i < size; i++) {
    gi = i + BLOCK_LOW(id,p,n);
}
```

…takes place of sequential program's index $i$

Research Computing @ CU Boulder · Week 4 - Collective Communication · 19 · 2/8/12

## When is a parallel implementation of our algorithm worth it?

- Cost of sending k byte message:   $T = t_s + k\, t_c$

- Cost of a floating-point operation:  $\omega$

- On Frost (very roughly),

- $T = (3 \times 10^{-6}) + k\,(6 \times 10^{-10})$ sec

- $\omega = 4 \times 10^{-14}$ sec
      23 Tflops = $23 \times 10^{12}$ floating-point ops / second

- $T/\omega = 7 \times 10^{7}$

Research Computing @ CU Boulder · Week 4 - Collective Communication · 20 · 2/8/12

## Let's take it up a notch

- Matrix-vector multiply b = A*x  (dense nxn A)

- One algorithm does dot products of rows of A with x
- Can we use what we did with dot products alone?
      Or do we need to reconsider?

- What's the best way to divide up this work?
- Consider costs of
  - Arithmetic
  - Communication

Research Computing @ CU Boulder · Week 4 - Collective Communication · 21 · 2/8/12

## Ways to partition A

- Block row
- Block column
- Checkerboard

- Vectors are partitioned accordingly

## How does the picture change for matrix-matrix multiply?

- C = A*B, both A and B nxn and dense (so C is, too)

- Need to partition all three matrices this time.

## What about matrix transpose?

- It's all about communication—no arithmetic!