

## Week 4 – Collective communication, topologies

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

Research Computing @ CU Boulder

Collective Communication and topologies

1

2/6/12

---

---

---

---

---

---

---

---

## Lecture 6

Collective Communication and topologies  
 Network topology slides adopted from:  
 Networks: Topologies - HPC Advisory Council  
[www.hpcadvisorycouncil.com/events/2011/.../pdf/.../10\\_HPCAC.pdf](http://www.hpcadvisorycouncil.com/events/2011/.../pdf/.../10_HPCAC.pdf)

Research Computing @ CU Boulder

---

---

---

---

---

---

---

---

## The project

Milestone	Due Date
Project Matchmaking	Wednesday, Feb.15
Project Proposal	Wednesday, Feb.24
CSCI 5576 Paper Presentations	March 12-16
Annotated Bibliography	March 23
Project Checkpoint	Monday, Apr. 4
Final Presentations	April 30 – May 4
Final Paper	May 9

Research Computing @ CU Boulder

Collective Communication and topologies

3

2/6/12

---

---

---

---

---

---

---

---

## World leading large scale machines

- National Supercomputing Centre in Shenzhen
  - Fat-tree, 5.2K nodes, 120K cores, NVIDIA GPUs, China (Petaflop)
- Tokyo Institute of Technology
  - Fat-tree, 4K nodes, NVIDIA GPUs, Japan (Petaflop)
- Commissariat à l'Energie Atomique (CEA)
  - Fat-tree, 4K nodes, 140K cores, France (Petaflop)
- Los Alamos National Lab – Roadrunner
  - Fat-tree, 4K nodes, 130K cores, USA (Petaflop)
- NASA
  - Hypercube, 9.2K nodes, 82K cores – NASA, USA
- Jülich JuRoPa
  - Fat-tree, 3K nodes, 30K cores, Germany
- Sandia National Labs – Red Sky
  - 3D-Torus, 5.4K nodes, 43K cores – Sandia "Red Sky", USA



Research Computing @ CU Boulder

Collective Communication and topologies

4

2/6/12

---

---

---

---

---

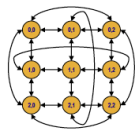
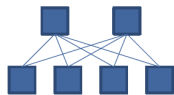
---

---

---

## Network topologies

- Fat-tree (CLOS), Mesh, 3D-Torus topologies
- CLOS (fat-tree)
  - Can be fully non-blocking (1:1) or blocking (x:1)
  - Typically enables best performance
    - Non blocking bandwidth, lowest network latency
- Mesh or 3D Torus
  - Blocking network, cost-effective for systems at scale
  - Great performance solutions for applications with locality
  - Support for dedicate sub-networks
  - Simple expansion for future growth



Research Computing @ CU Boulder

Collective Communication and topologies

5

2/6/12

---

---

---

---

---

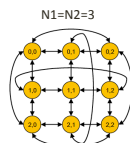
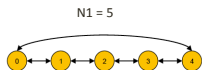
---

---

---

## d-Dimensional Torus Topology

- Formal definition
  - $T=(V,E)$  is said to be d-dimensional torus of size  $N_1 \times N_2 \times \dots \times N_d$  if:
    - $V=\{(v_1, v_2, \dots, v_d) : 0 \leq v_i \leq N_i - 1\}$
    - $E=\{(u \rightarrow v) : \text{exists } j \text{ s.t. } 1) \text{ for each } i \neq j, v_i = u_i \text{ AND } 2) v_j = (u_j \pm 1) \bmod N_j\}$
- Examples



Research Computing @ CU Boulder

Collective Communication and topologies

6

2/6/12

---

---

---

---

---

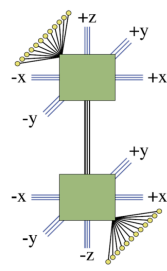
---

---

---

## 3D-Torus System with Infiniband – Key Items

- Multiple server nodes per cube junction
- Smaller 3D cube size the better
  - Lowest latency between remote nodes
- Minimizing throughput contention
- Ability to connect storage
- Support for separate networks
  - Dedicated network (links) for specific applications/usage
  - Example: links dedicated for collectives or specific jobs



Research Computing @ CU Boulder

Collective Communication and topologies

7

2/6/12

---

---

---

---

---

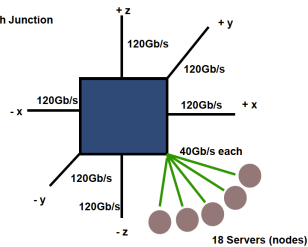
---

---

---

## Example 3D Torus

3D Torus Switch Junction



3D Torus size: 8x8x8 (512 36-port switches)  
Total number of servers: 9216

Research Computing @ CU Boulder

Collective Communication and topologies

8

2/6/12

---

---

---

---

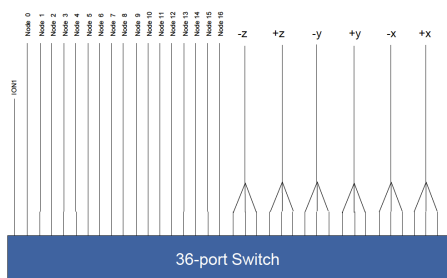
---

---

---

---

## 3D Torus Connections Example



Research Computing @ CU Boulder

Collective Communication and topologies

9

2/6/12

---

---

---

---

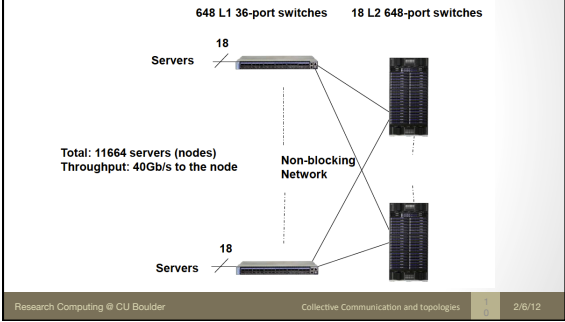
---

---

---

---

### Example: Non-blocking, fat tree, Infiniband



---

---

---

---

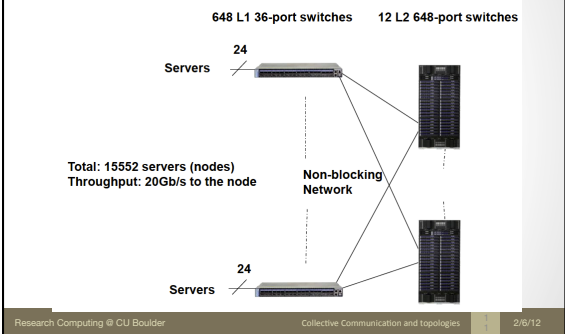
---

---

---

---

### Example: 2:1 Oversubscription



---

---

---

---

---

---

---

---

### Determining the network topology

Interconnection Network (API)	Topology Tool(s)	Query
Myrinet (MX)	fm_db2wirelist	
InfiniBand (OFED)	ibdiagnet & ibnet-discover	
SeaStar (Cray XT)	xtprocadmin & xtdb2proc	
BlueGene/P (DCMF)	DCMF API	

---

---

---

---

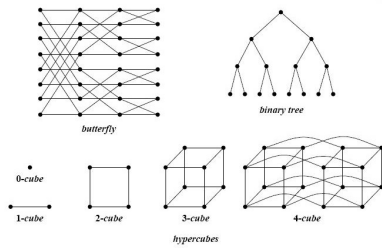
---

---

---

---

## Topologies from algorithm design



Research Computing @ CU Boulder

Collective Communication and topologies

2/6/12

---

---

---

---

---

---

---

---

## Collective communication

- Broadcast
  - one to all: data from root to all others
- Gather
  - all to one: all send data to root (result is vector)
- Scatter
  - one to all: reverse of gather
- Reduce
  - all to one: combine results on all at root via specified operation

Research Computing @ CU Boulder

Collective Communication and topologies

2/6/12

---

---

---

---

---

---

---

---

## MPI operations

- MPI\_BCAST
- MPI\_GATHER
- MPI\_SCATTER
- MPI\_REDUCE

arguments similar to those of MPI\_SEND and MPI\_RECV

Research Computing @ CU Boulder

Collective Communication and topologies

2/6/12

---

---

---

---

---

---

---

---

## How to do gather without collective routine

- The simple way of where all send to 0 (root) at once
- More efficiently, take advantage of network topology
  - Can reduce number of steps
  - Does reduce contention for wires

Research Computing @ CU Boulder

Collective Communication and topologies

2/6/12

---

---

---

---

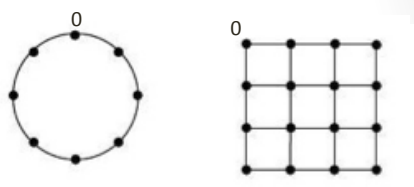
---

---

---

---

## Contention



Two simple and important topologies show  
topologies are graphs: nodes and edges

Research Computing @ CU Boulder

Collective Communication and topologies

2/6/12

---

---

---

---

---

---

---

---

## Consider a ring of $p$ processors

- Gather requires  $p/2$  communication steps
- Cost per step =  $T = t_s + k t_c$ 
  - $t_s$  = startup or latency
  - $k$  = bytes in message (constant size here)
  - $t_c$  = communication cost
  - $1/t_c$  = bandwidth
- Total cost =  $(p/2) * T$

Research Computing @ CU Boulder

Collective Communication and topologies

2/6/12

---

---

---

---

---

---

---

---

## What's a more efficient topology?

Do we immediately see how to do a gather on a mesh?

Research Computing @ CU Boulder

Collective Communication and topologies

1

2/6/12

---

---

---

---

---

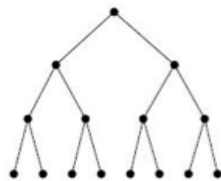
---

---

---

## Tree

- Gather requires \_\_\_\_ steps
- Cost per step =
- Total cost =



Research Computing @ CU Boulder

Collective Communication and topologies

2

2/6/12

---

---

---

---

---

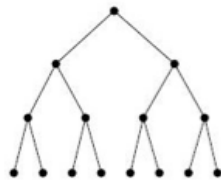
---

---

---

## Tree

- Gather requires  $\log_2 p$  steps
- Cost per step =  $T = t_s + k t_c$
- Total cost =  $\log_2 p \cdot T$



Research Computing @ CU Boulder

Collective Communication and topologies

3

2/6/12

---

---

---

---

---

---

---

---

## But we don't have a tree!

- Graph embeddings allow us to find one graph in another  
(one topology in another)  
= virtual circuit
- Designing efficient parallel programs requires us to understand graph embeddings
  - What are efficient communication schemes?
  - What schemes avoid contention?
  - What schemes permit nearest neighbor communication?

Research Computing @ CU Boulder

Collective Communication and topologies

2/6/12

---

---

---

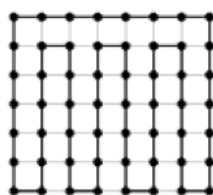
---

---

---

---

## One example



ring in 2-D mesh

Research Computing @ CU Boulder

Collective Communication and topologies

2/6/12

---

---

---

---

---

---

---

## Graph embedding

- Graph embedding gives us flexibility in how we design algorithm: use convenient communication graph  $V_s$  and embed it into network graph  $V_t$
- Concerns
  - **load** : maximum number of nodes in  $V_s$  mapped to same node in  $V_t$  (CPU use)
  - **congestion** : maximum number of edges in  $E_s$  mapped to paths containing same edge in  $E_t$  (bandwidth, contention)
  - **dilation** : maximum distance between any two nodes  $(\Omega(u), \Omega(v)) \in V_t$  such that  $(u, v) \in E_s$  (nearest neighbors)

Research Computing @ CU Boulder

Collective Communication and topologies

2/6/12

---

---

---

---

---

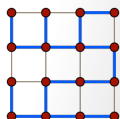
---

---



## How to embed

- Finding optimal embedding NP-complete, in general, so heuristics used to determine good embedding
- Special graphs
  - Tree : connected graph containing no cycles
  - Spanning tree : subgraph that includes all nodes of given graph and is also a tree
- Spanning tree algorithm prevents redundant communication



Research Computing @ CU Boulder

Collective Communication and topologies

2/6/12

---

---

---

---

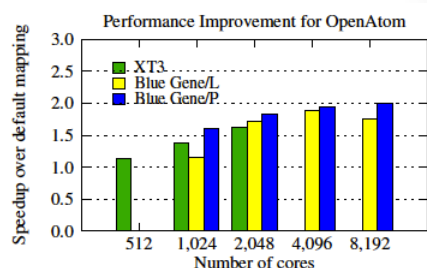
---

---

---

---

## Example of speedup by mapping



Research Computing @ CU Boulder

Collective Communication and topologies

2/6/12

---

---

---

---

---

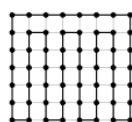
---

---

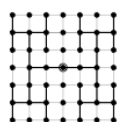
---

## Examples: graph embeddings

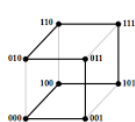
- For some important cases, good or optimal embeddings are known:



ring in 2-D mesh  
dilation 1



binary tree in 2-D mesh  
dilation  $\lfloor (k-1)/2 \rfloor$



ring in hypercube  
dilation 1

Research Computing @ CU Boulder

Collective Communication and topologies

2/6/12

---

---

---

---

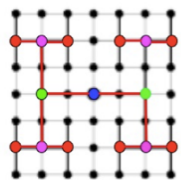
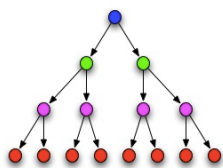
---

---

---

---

So focus on embedding tree into mesh (torus)



binary tree (not spanning tree) into 2-D mesh  
dilation =  $\text{ceil}((k-1)/2)$

Research Computing @ CU Boulder

Collective Communication and topologies

2/6/12

---

---

---

---

---

---

---

---

Other operations similarly based on trees

- Reduce has the same communication scheme as gather but with operations on values
- Broadcast and scatter go from root to leaves
- Other ops like MPI\_ALLGATHER: all to all, communication goes in one direction then other
- Next weeks' lab assignment: butterfly allreduce
  - Topology is butterfly
  - Communication goes in both directions

Research Computing @ CU Boulder

Collective Communication and topologies

2/6/12

---

---

---

---

---

---

---

---