

**ECE 1749H:
Interconnection Networks for
Parallel Computer Architectures:**

Topology

Prof. Natalie Enright Jerger

Announcements

- Tentative Presentation Schedule posted
 - E-mail me if:
 - You have registered/plan to register AND aren't on the list
 - You have a significant conflict with your assigned date
 - E.g. conference travel

Announcements (2)

- Title: Living in interesting times:
Disruptive trends in computer
architecture
- Where: GB405
- When: Wednesday, January 20, 2010,
2-3pm
- <http://www.eecg.utoronto.ca/cider/>
- Speaker: Bob Blainey, IBM Toronto

Last time

- Why on-chip networks?
- Various system architectures
 - Interactions with on-chip network

Topology Overview

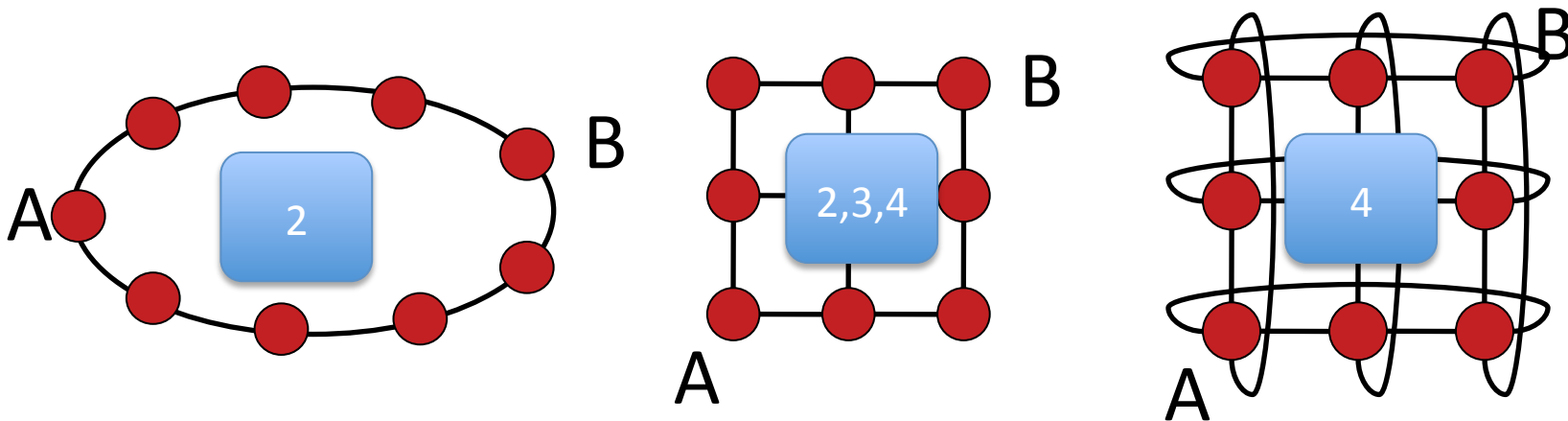
- Definition: determines arrangement of channels and nodes in network
 - Analogous to road map
- Often first step in network design
- Significant impact on network cost-performance
 - Determines number of **hops**
 - Latency
 - Network energy consumption
 - Implementation **complexity**
 - Node degree
 - Ease of layout

Abstract Metrics

- Use metrics to evaluate **performance** and **cost** of topology
- Also influenced by routing/flow control
 - At this stage
 - Assume **ideal** routing (perfect load balancing)
 - Assume **ideal** flow control (no idle cycles on any channel)

Abstract Metrics: Degree

- Switch Degree: number of links at a node
 - Proxy for estimating **cost**
 - Higher degree requires more links and port counts at each router



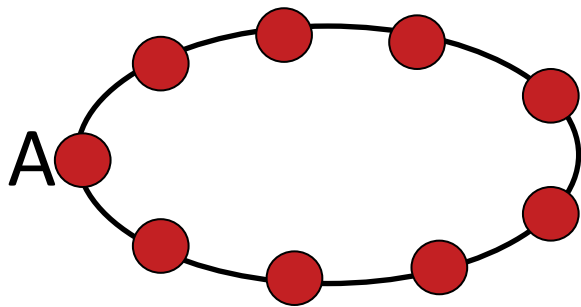
Abstract Metrics: Hop Count

- Path: ordered set of channels between source and destination
- Hop Count: number of hops a message takes from source to destination
 - Simple, useful proxy for network **latency**
 - Every node, link incurs some propagation delay even when no contention
- Minimal hop count: smallest hop count connecting two nodes

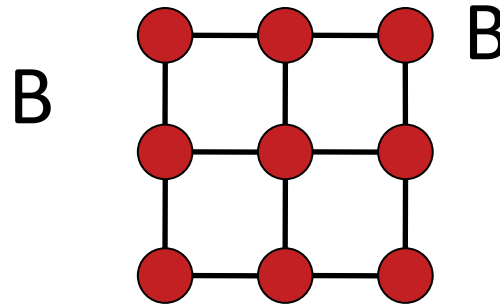
Hop Count

- Network **diameter**: large min hop count in network
- Average minimum hop count: average across all src/dst pairs
 - Implementation may incorporate non-minimal paths
 - Increases average hop count

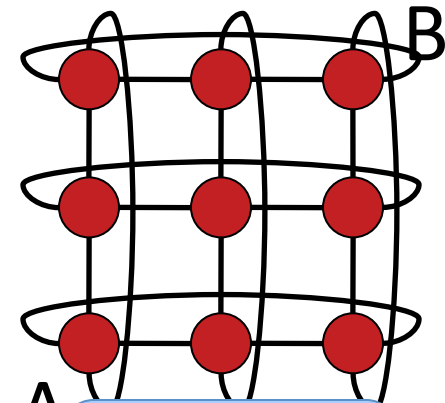
Hop Count



Max = 4
Avg = 2.2



Max = 4
1.77



Max = 2
1.33

- Uniform random traffic
 - Ring > Mesh > Torus
- Derivations later

Latency

- Time for packet to traverse network
 - Start: head arrives at input port
 - End: tail departs output port
- Latency = Head latency + serialization latency
 - Serialization latency: time for packet with Length L to cross channel with bandwidth b (L/b)
- Approximate with hop count
 - Other design choices (routing, flow control) impact latency
 - Unknown at this stage

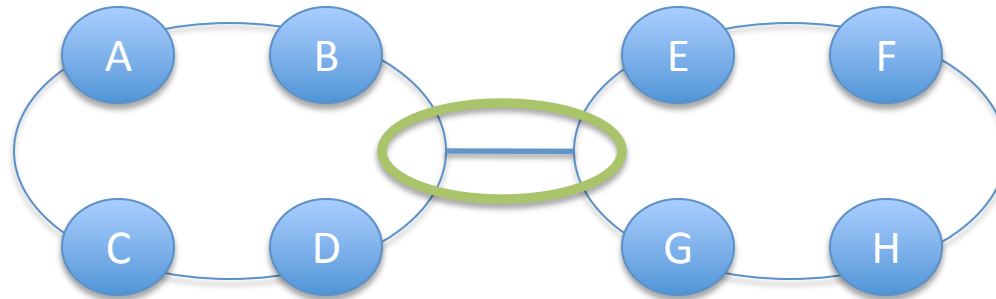
Abstract Metrics: Maximum Channel Load

- Estimate max **bandwidth** the network can support
 - Max bits per second (bps) that can be injected by every node before it saturates
 - **Saturation**: network cannot accept any more traffic
 - Determine most congested link
 - For given traffic pattern
 - Will limit overall network bandwidth
 - Estimate load on this channel

Maximum Channel Load

- Preliminary
 - Don't know specifics of link yet
 - Define relative to injection load
- Channel load of 2
 - Channel is loaded with twice injection bandwidth
 - If each node injects a flit every cycle
 - 2 flits will want to traverse bottleneck channel every cycle
 - If bottleneck channel can only handle 1 flit per cycle
 - Max network bandwidth is $\frac{1}{2}$ link bandwidth
 - A flit can be injected every other cycle

Maximum Channel Load Example

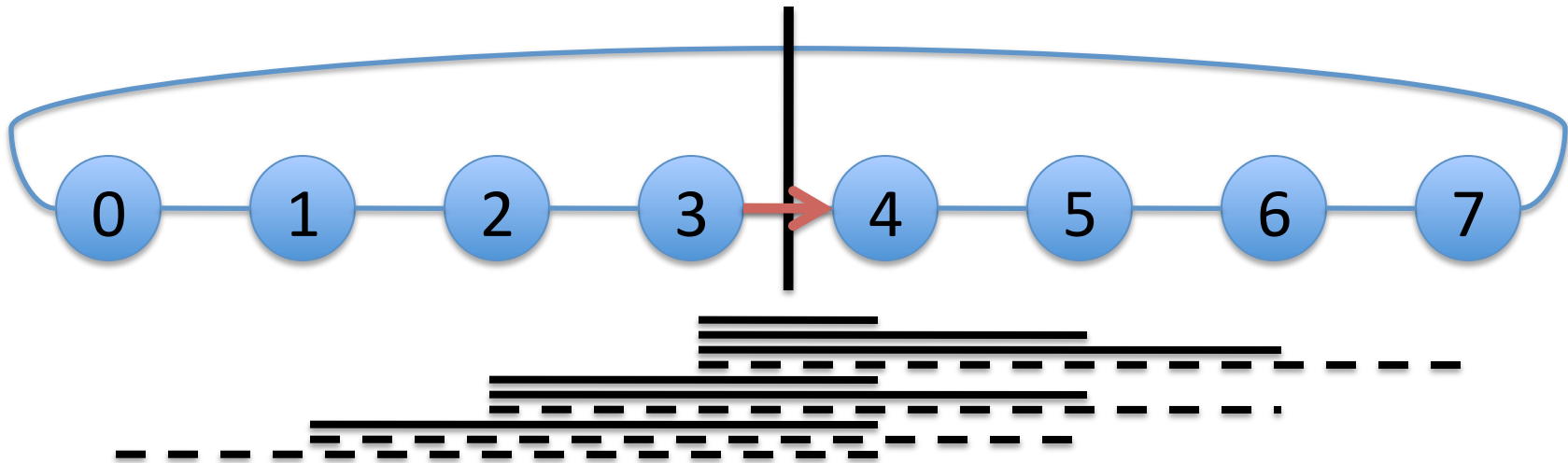


- Uniform random
 - Every node has equal probability of sending to every node
- Identify bottleneck channel
- Half of traffic from every node will cross bottleneck channel
 - $8 \times \frac{1}{2} = 4$
- Network saturates at $\frac{1}{4}$ injection bandwidth

Bisection Bandwidth

- Common off-chip metric
 - Proxy for cost
 - Amount of global wiring that will be necessary
 - Less useful for on-chip
 - Global on-chip wiring considered **abundant**
- Cuts: partition all the nodes into two disjoint sets
 - Bandwidth of a cut
- Bisection
 - A cut which divides all nodes into nearly half
 - Channel bisection \rightarrow min. channel count over all bisections
 - Bisection bandwidth \rightarrow min. bandwidth over all bisections
- With uniform traffic
 - $\frac{1}{2}$ of traffic crosses bisection

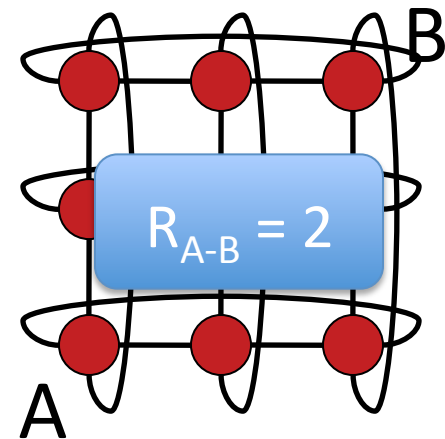
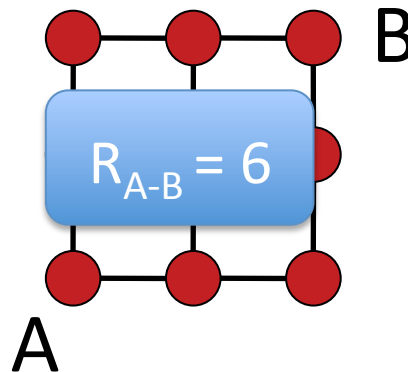
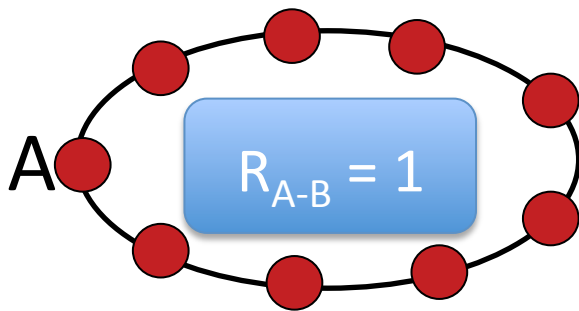
Throughput Example



- Bisection = 4 (2 in each direction)
- With uniform random traffic
 - 3 sends $1/8$ of its traffic to 4,5,6
 - 3 sends $1/16$ of its traffic to 7 (2 possible shortest paths)
 - 2 sends $1/8$ of its traffic to 4,5
 - Etc
- Channel load = 1

Path Diversity

- Multiple shortest paths between source/destination pair (R)
- **Fault tolerance**
- Better **load balancing** in network
- Routing algorithm should be able to exploit path diversity

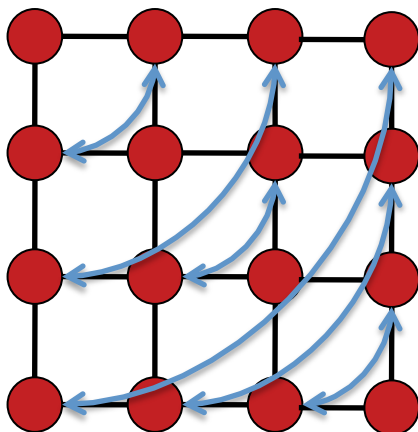


Evaluating Topologies

- Important to consider traffic pattern
- Talked about system architecture impact on traffic
- If actual traffic pattern unknown
 - Synthetic traffic patterns
 - Evaluate common scenarios
 - Stress test network
 - Derive various properties of network

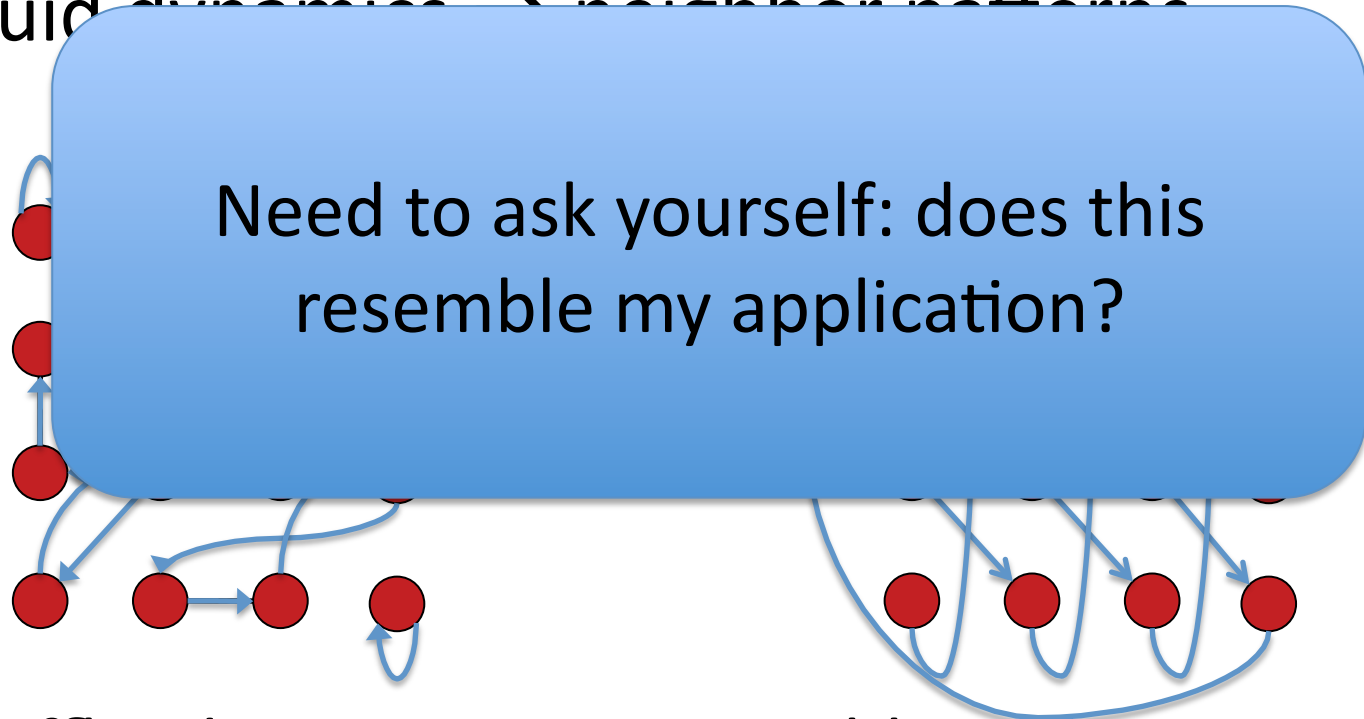
Traffic Patterns

- Historically derived from particular applications of interest
 - Spatial distribution
 - Matrix Transpose \rightarrow Transpose traffic pattern
 - $d_i = s_{i+b/2 \bmod b}$
 - b-bit address, d_i : ith bit of destination



Traffic Patterns (2)

- Fast Fourier Transform (FFT) or sorting application \rightarrow shuffle permutation
- Fluid dynamics \rightarrow neighbor patterns



Shuffle: $d_i = s_{i-1 \bmod b}$

Neighbor: $d_x = s_x + 1 \bmod k$

Traffic Patterns (3)

- Uniform random
 - Each source equally likely to communication with each destination
 - Most commonly used traffic pattern
 - Very **benign**
 - Traffic is uniformly distributed
 - Balances load even if topology/routing algorithm has very poor load balancing
 - Need to be careful
 - But can be good for debugging/verifying implementation
 - Well-understood pattern

Stress-testing Network

- Uniform random can make bad topologies look good
- Permutation traffic will stress-test the network
 - Many types of permutation (ex: shuffle, transpose, neighbor)
 - Each source sends all traffic to single destination
 - Concentration of load on individual pairs
 - Stresses load balancing

Final Thoughts: Traffic Patterns

- For topology/routing discussion
 - Focus on **spatial** distribution
- Traffic patterns also have **temporal** aspects
 - Bursty behavior
 - Important to capture temporal behavior as well

Types of Topologies

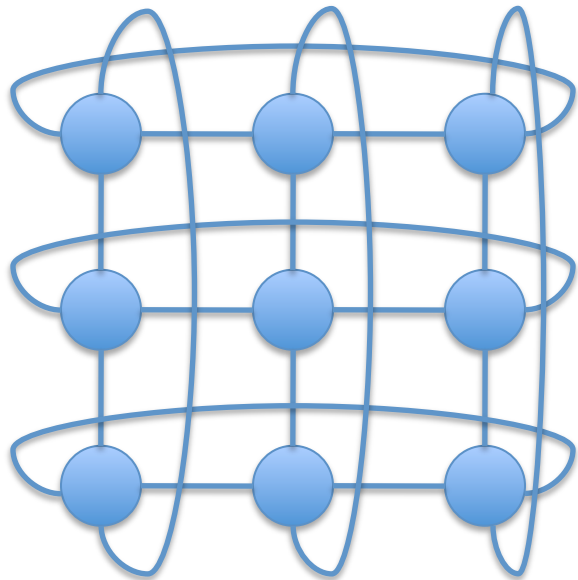
- Focus on switched topologies
 - Alternatives: bus and crossbar
 - Bus
 - Connects a set of components to a single shared channel
 - Effective broadcast medium
 - Crossbar
 - Directly connects n inputs to m outputs without intermediate stages
 - Fully connected, single hop network
 - Component of routers

Types of Topologies

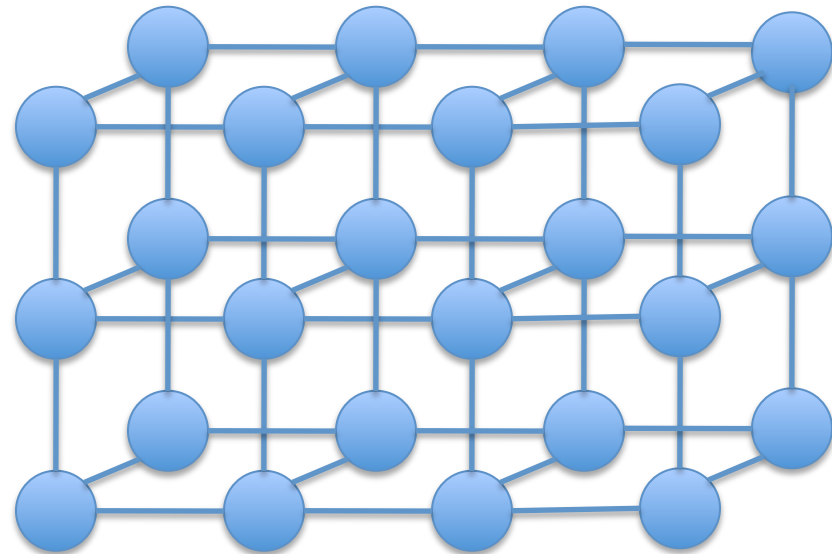
- **Direct**
 - Each router is associated with a terminal node
 - All routers are sources and destinations of traffic
- **Indirect**
 - Routers are distinct from terminal nodes
 - Terminal nodes can source/sink traffic
 - Intermediate nodes switch traffic between terminal nodes
- Most on-chip network use direct topologies

Torus (1)

- K-ary n-cube: k^n network nodes
- N-Dimensional grid with k nodes in each dimension



3-ary 2-mesh



2,3,4-ary 3-mesh

Torus (2)

- Map well to planar substrate for on-chip
- Topologies in Torus Family
 - Ex: Ring -- k-ary 1-cube
- Edge Symmetric
 - Good for load balancing
 - Removing wrap-around links for mesh loses edge symmetry
 - More traffic concentrated on center channels
- Good path diversity
- Exploit locality for near-neighbor traffic

Torus (3)

- Hop Count:
$$H_{\min} = \begin{cases} \frac{nk}{4} & k \text{ even} \\ n\left(\frac{k}{4} - \frac{1}{4k}\right) & k \text{ odd} \end{cases}$$

– For uniform random traffic

- Packet travels $k/4$ hops in each of n dimensions

- For Mesh
$$H_{\min} = \begin{cases} \frac{nk}{3} & k \text{ even} \\ n\left(\frac{k}{3} - \frac{1}{3k}\right) & k \text{ odd} \end{cases}$$

Torus (4)

- Degree = $2n$, 2 channels per dimension
 - All nodes have same degree
- Total channels = $2nN$

Channel Load for Torus

- Even number of k-ary (n-1)-cubes in outer dimension
- Dividing these k-ary (n-1)-cubes gives a 2 sets of k^{n-1} bidirectional channels or $4k^{n-1}$
- $\frac{1}{2}$ Traffic from each node cross bisection

$$\text{channel load} = \frac{N}{2} \times \frac{k}{4N} = \frac{k}{8}$$

- Mesh has $\frac{1}{2}$ the bisection bandwidth of torus

Torus Path Diversity

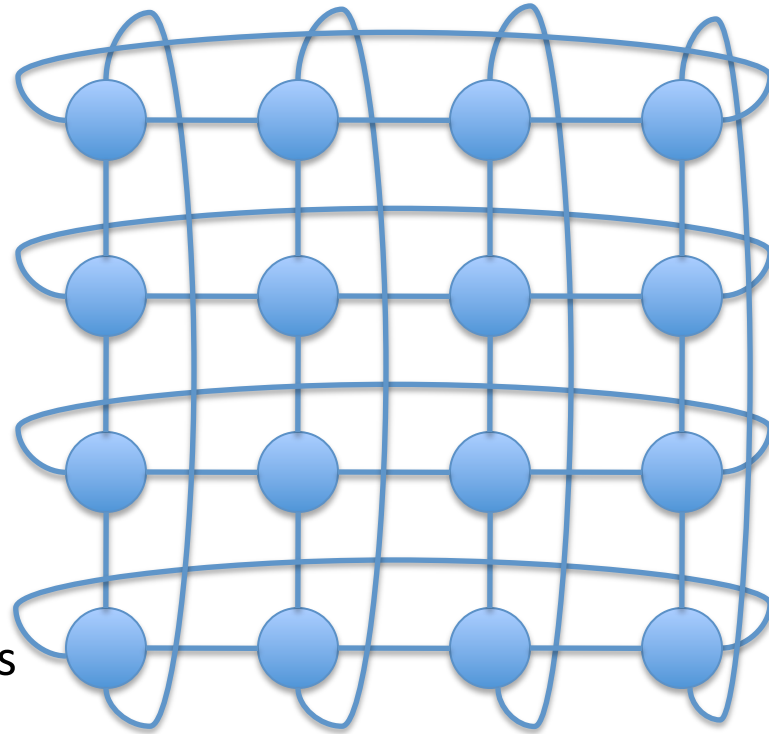
$$|R_{xy}| = \binom{\Delta x + \Delta y}{\Delta x}$$

2 dimensions*

$$\Delta x = 2, \Delta y = 2$$

$$|R_{xy}| = 6$$

$$|R_{xy}| = 24 \quad \text{NW, NE, SW, SE combos}$$



2 edge and node disjoint minimum paths

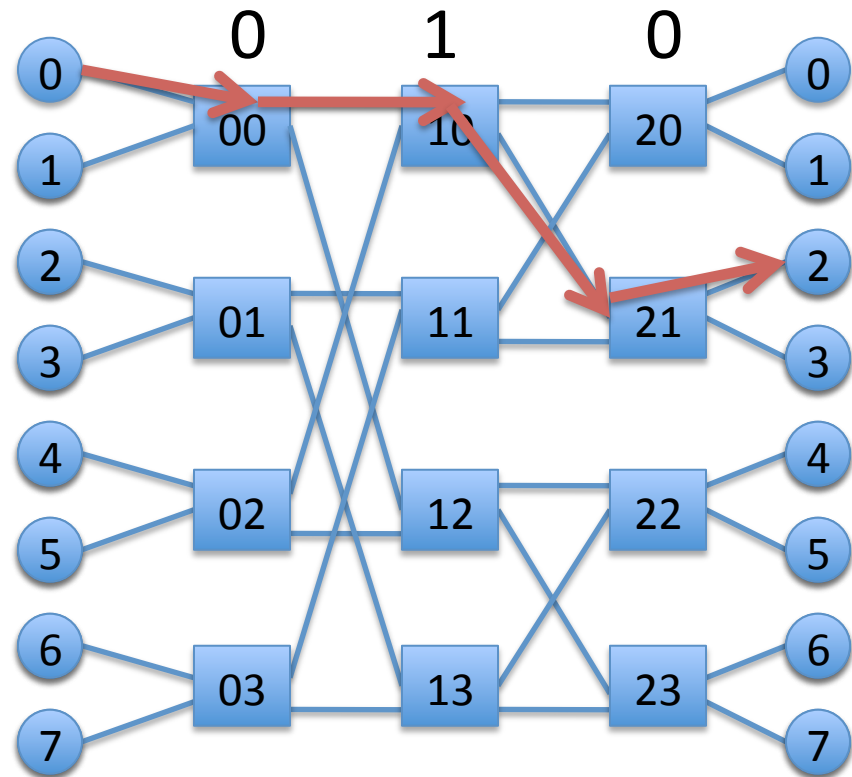
*assume single direction for x and y

Mesh

- A torus with end-around connection removed
- Same node degree
- Bisection channels halved
 - Max channel load = $k/4$
- Higher demand for central channels
 - Load imbalance

Butterfly

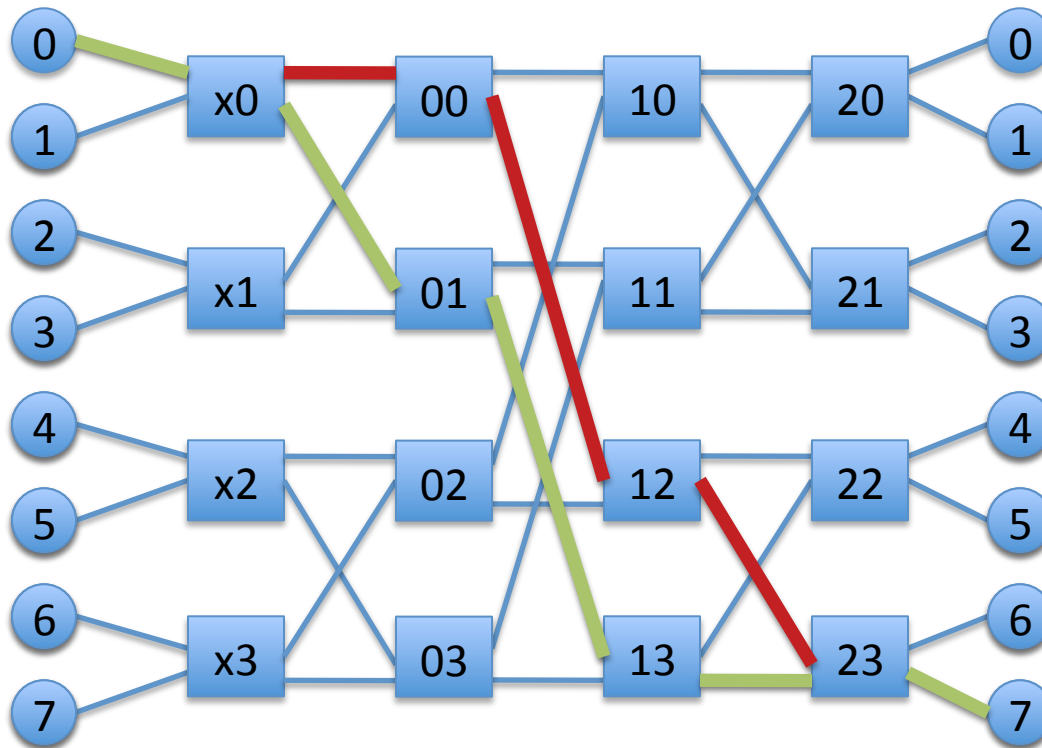
- Indirect network
- K-ary n-fly: k^n network nodes
- Routing from 000 to 010
 - Dest address used to directly route packet
 - Bit n used to select output port at stage n



2-ary 3-fly
2 input switch, 3 stages

Butterfly (2)

- No path diversity $|R_{xy}| = 1$
 - Can add extra stages for diversity
 - Increase network diameter



Butterfly (3)

- Hop Count
 - $\log_k N + 1$
 - Does **not** exploit **locality**
 - Hop count same regardless of location
- Switch Degree = $2k$
- Requires long wires to implement

Butterfly: Channel Load

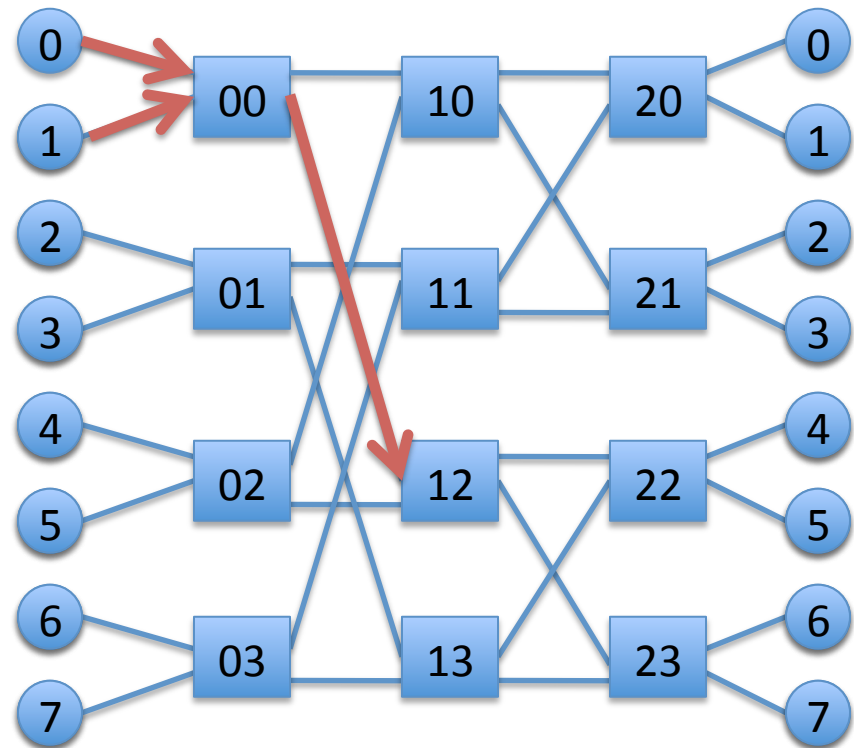
- $H_{\min} \times N$: Channel demand
 - Number of channel traversals required to deliver one round of packets
- Channel Load \rightarrow uniform traffic
 - Equally loads channels

$$\frac{NH_{\min}}{C} = \frac{k^n (n + 1)}{k^n (n + 1)} = 1$$

- Increases for adversarial traffic

Butterfly: Channel Load

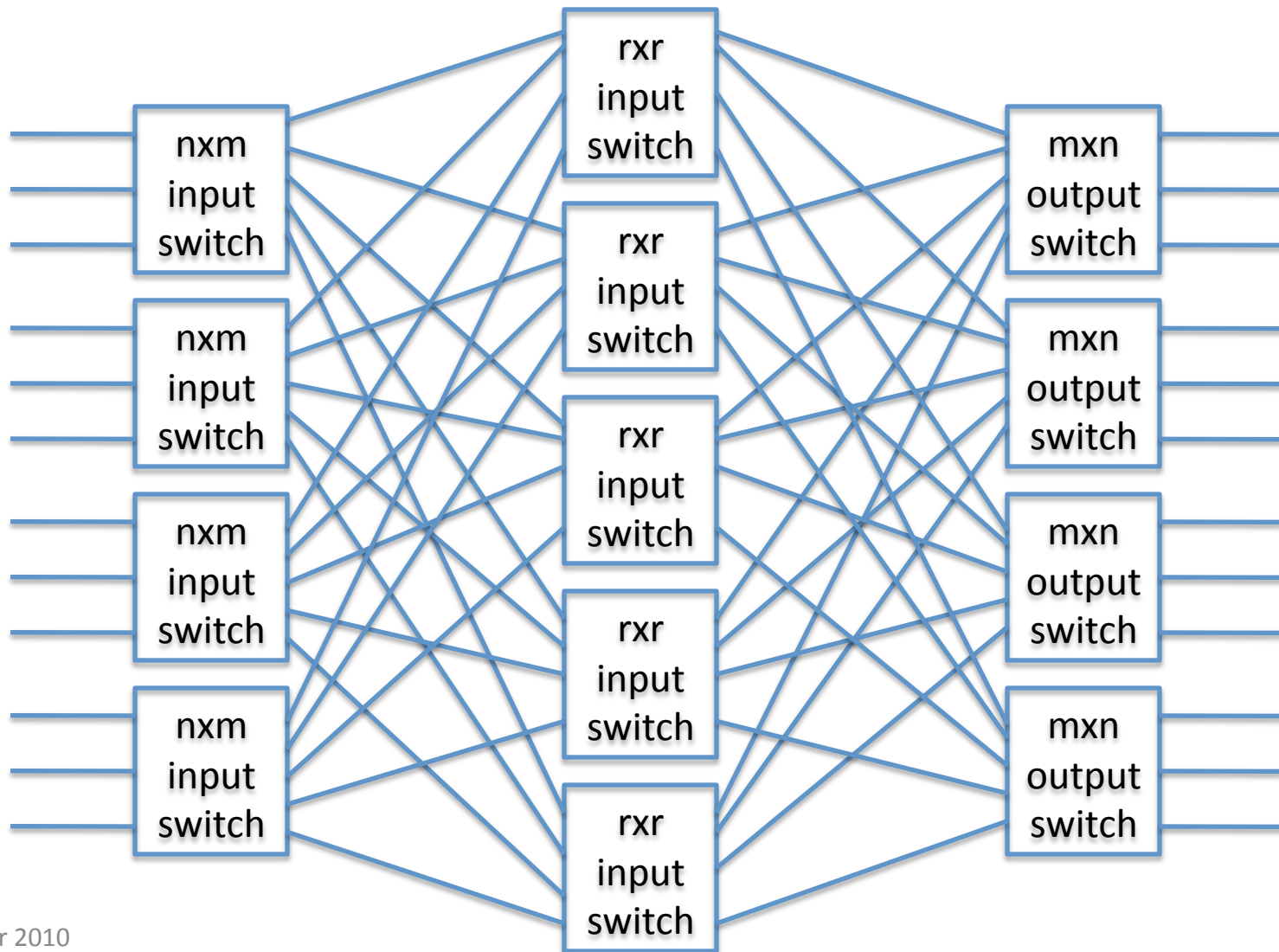
- Adversarial traffic
 - All traffic from top half sent to bottom half
 - E.g. 0 sends to 4, 1 sends to 5



Clos Network

- 3-stage indirect network
 - Larger number of stages: built recursively by replacing middle stage with 3-stage Clos
- Characterized by triple (m, n, r)
 - M : # of middle stage switches
 - N : # of input/output ports on input/output switches
 - R : # of input/output switches
- Hop Count = 4

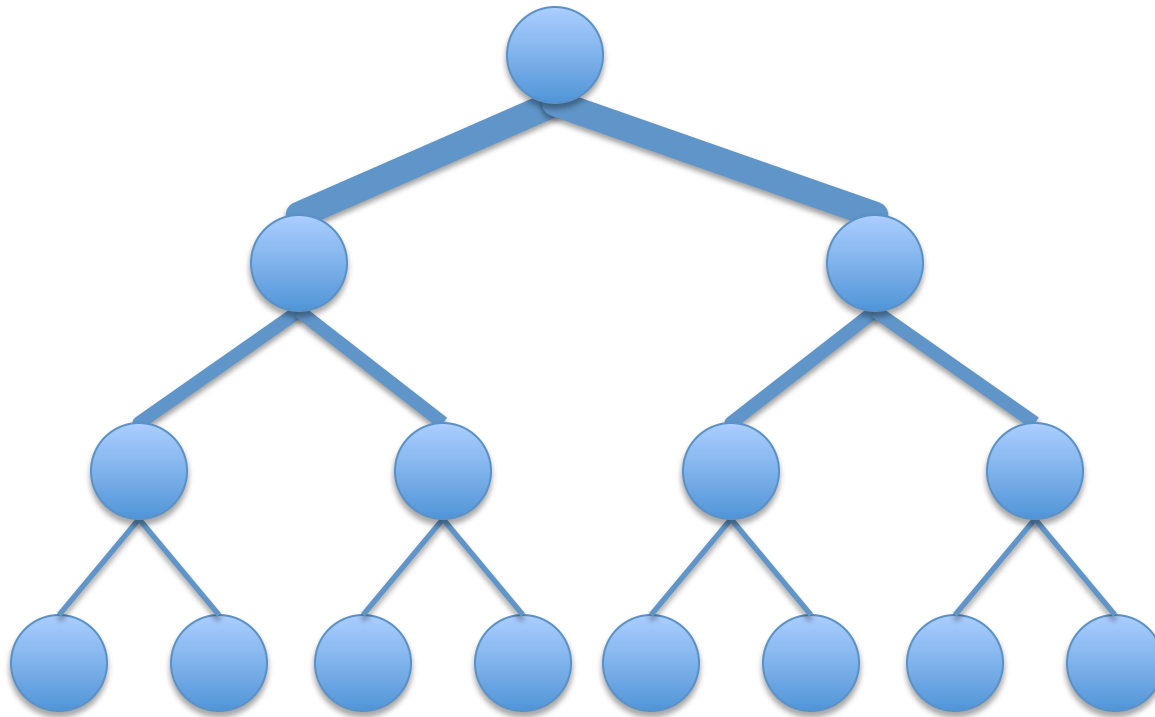
Clos Network



Clos Network

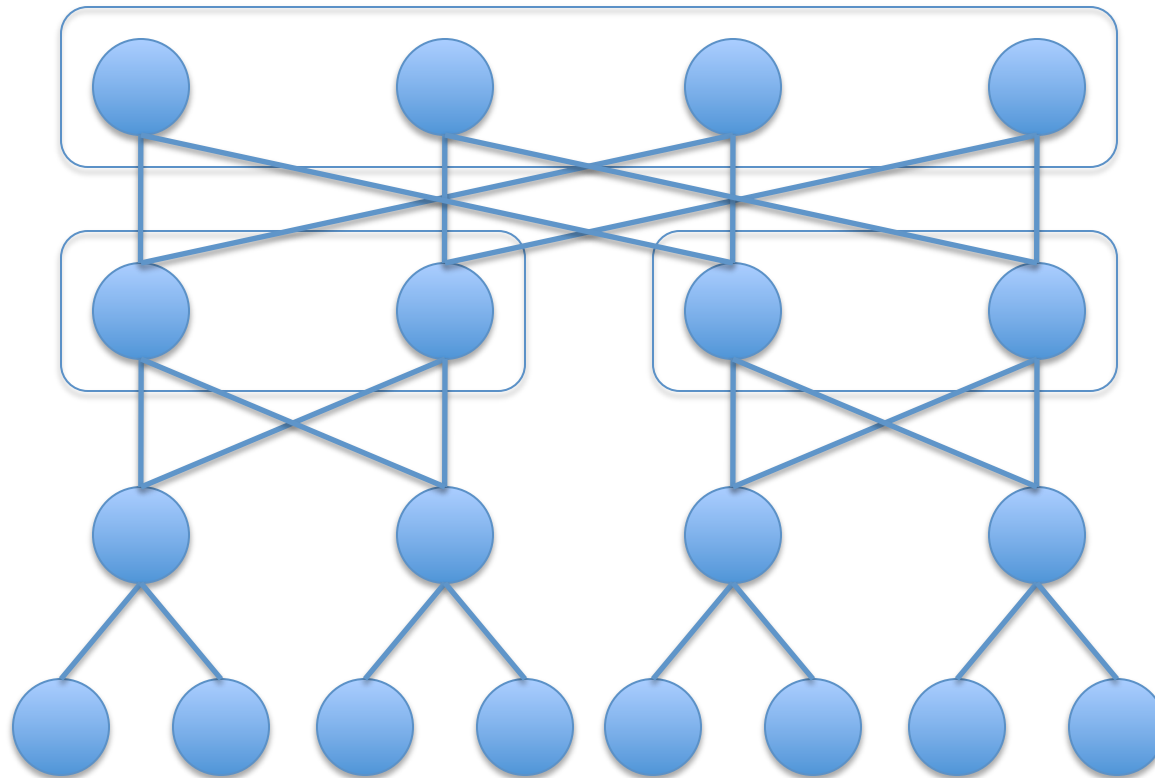
- Non-blocking when $m > 2n-1$
 - Any input can connect to any unique output port
- $r \times n$ nodes
- Degree
 - First and last stages: $n + m$, middle stage: $2r$
- Path diversity: m
- Can be folded along middle switches
 - Input and output switches are shared

Folded Clos (Fat Tree)



- Bandwidth remains constant at each level
- Regular Tree: Bandwidth decreases closer to root

Fat Tree (2)

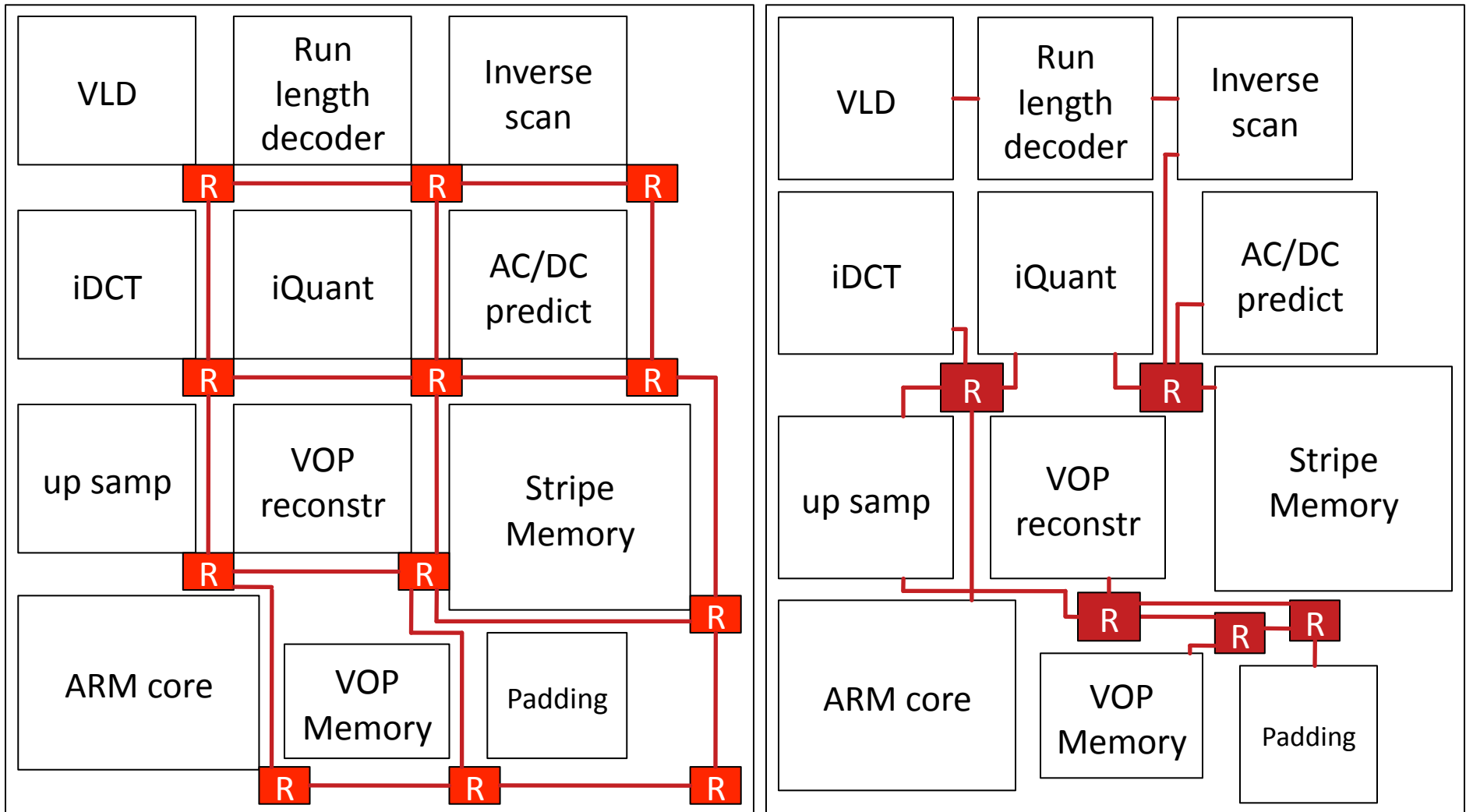


- Provides path diversity

Irregular Topologies

- MPSoC design leverages wide variety of IP blocks
 - Regular topologies may not be appropriate given heterogeneity
 - Customized topology
 - Often more power efficient and deliver better performance
- Customize based on traffic characterization

Irregular Topology Example

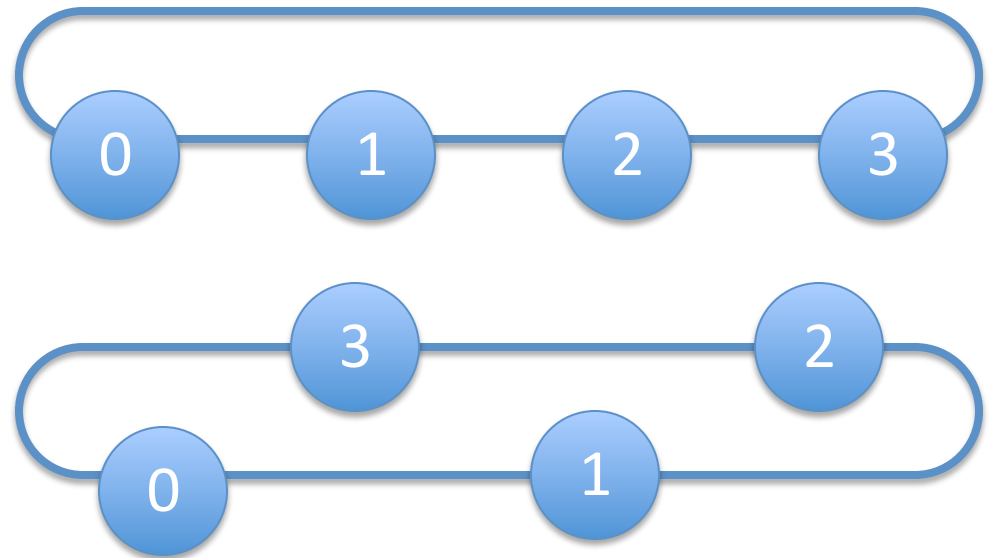


Topology Customization

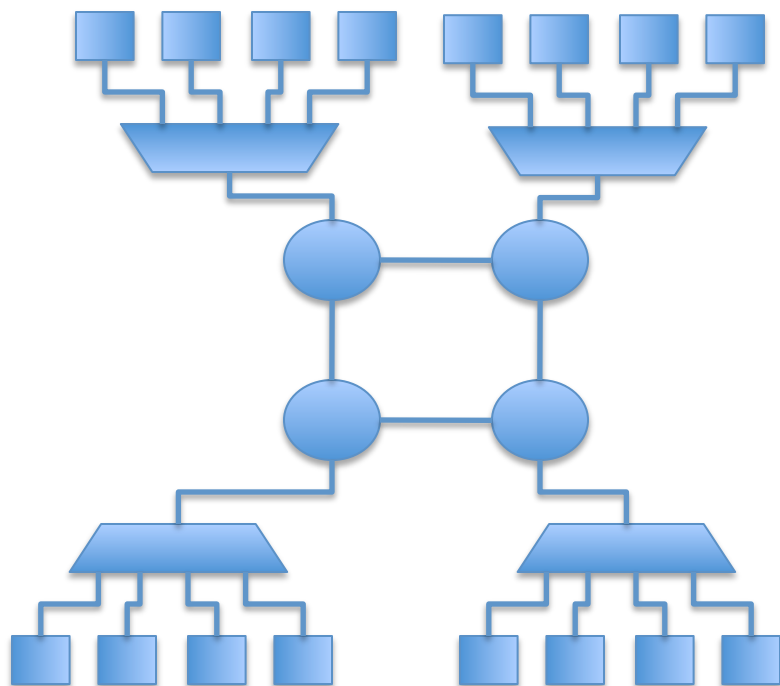
- Merging
 - Start with large number of switches
 - Merge to adjacent routers reduce area and power
- Splitting
 - Large crossbar connecting all nodes
 - Iteratively split into multiple small switches
 - Accommodate design constraints

Implementation

- Folding
 - Equalize path lengths
 - Reduces max link length
 - Increases length of other links



Concentration



- Don't need 1:1 ratio of routers to cores
 - Ex: 4 cores concentrated to 1 router
- Can save area and power
- Increases network complexity
 - Concentrator must implement policy for sharing injection bandwidth
 - During bursty communication
 - Can bottleneck

Implication of Abstract Metrics on Implementation

- Degree: useful proxy for router complexity
 - Increasing ports requires additional buffer queues, requestors to allocators, ports to crossbar
 - All contribute to critical path delay, area and power
 - Link complexity does not correlate with degree
 - Link complexity depends on link width
 - Fixed number of wires, link complexity for 2-port vs 3-port is same

Implications (2)

- Hop Count: useful proxy for overall latency and power

– Doe

- D
p

Hop Count says A is better than B
But A has 18 cycle latency vs 6 cycle
latency for B

– Exa

- Network A with 2 hops, 9 stage pipeline, 4 cycle link traversal vs.
- Network B with 3 hops, 1 stage pipeline, 1 cycle link traversal

Implications (3)

- Topologies typically trade-off hop count and node degree
- Max channel load useful proxy for network saturation and max power
 - Higher max channel load → greater network congestion
 - Traffic pattern impacts max channel load
 - Representative traffic patterns important
 - Max power: dynamic power is highest with peak switching activity and utilization in network

Topology Summary

- First network design decision
- Critical impact on network latency and throughput
 - Hop count provides first order approximation of message latency
 - Bottleneck channels determine saturation throughput