# ECE 4750 Computer Architecture, Fall 2016
# T06 Fundamental Network Concepts

School of Electrical and Computer Engineering
Cornell University
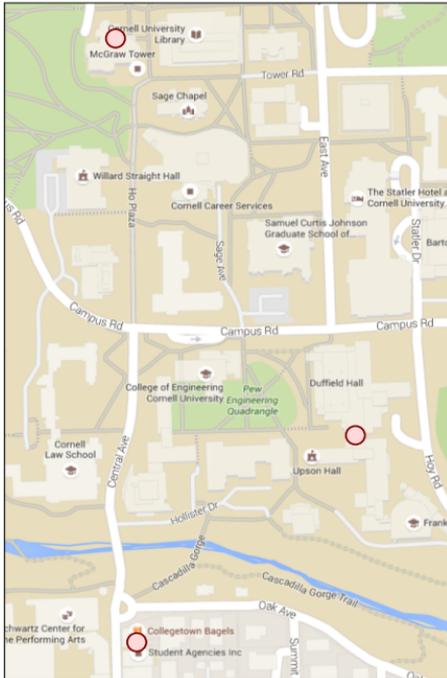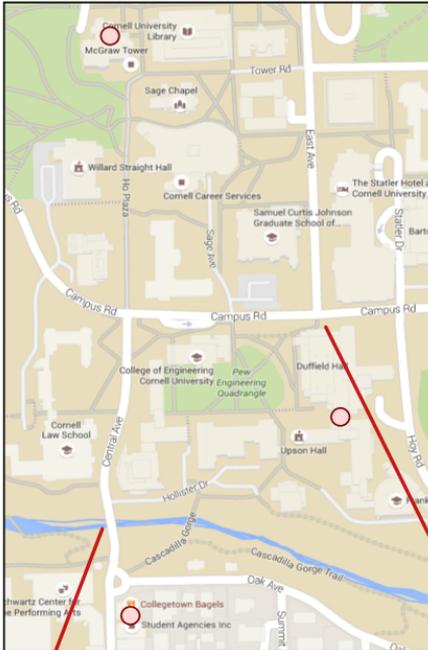
revision: 2016-10-16-22-06

# 1. Network/Roadway Analogy

Our goal is to run some errands around town using our bike. Assume we are studying in the Duffield Atrium and we need to: (1) do some laundry in Collegetown (maybe pickup some coffee?), and (2) pick up some books about computer architecture from the library.
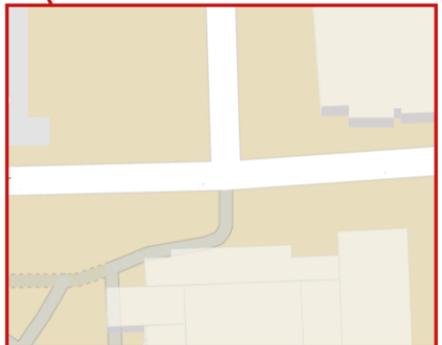


- Duffield Atrium
- Coffee/Laundry
- Library

## 1.1.  Running Errands



- Network Topology
  - Arrangement of roads and intersections to interconnect sources and destinations
  - Wide vs. narrow roads
  - Long vs. short roads
  - Small vs. large intersections

- Network Routing
  - Path from source to destination along roads and intersections
  - Short vs. long paths
  - Common vs. rare paths

- Network Microarchitecture
  - Managing long line of bikes and cars on road
  - Managing many cars and bikes at same intersection
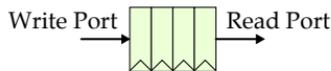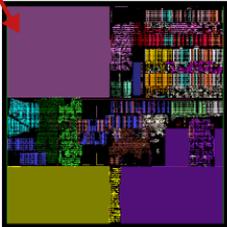
## 1.2. Network Technology

| | AWG24 Twisted Pair | PCB Trace | On-Chip M6 Wire in 0.18μm |
|---|---|---|---|
| |  |  |  |
| **Resistance** | $0.08\,\Omega/m$ | $5\,\Omega/m$ | $40\,k\Omega/m$ |
| **Inductance** | $400\,nH/m$ | $300\,nH/m$ | $4\,\mu H/m$ |
| **Capacitance** | $40\,pF/m$ | $100\,pF/m$ | $300\,pF/m$ |
| **Data Rate** | $\approx Gb/s$ | $\approx Gb/s$ | $\approx Gb/s$ |
| **Critical Length** | 1m | 10cm | <1cm |
| **Pitch** | $\approx mm$ | <mm | <μm |

**On-Chip Wires**



**On-Chip Buffers**



Distributed Wire Resistance and Capacitance



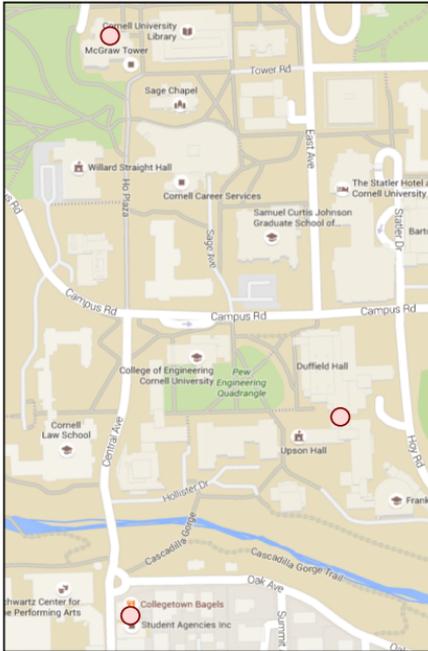Because both wire resistance and wire capacitance increase with length, wire delay grows quadratically with length
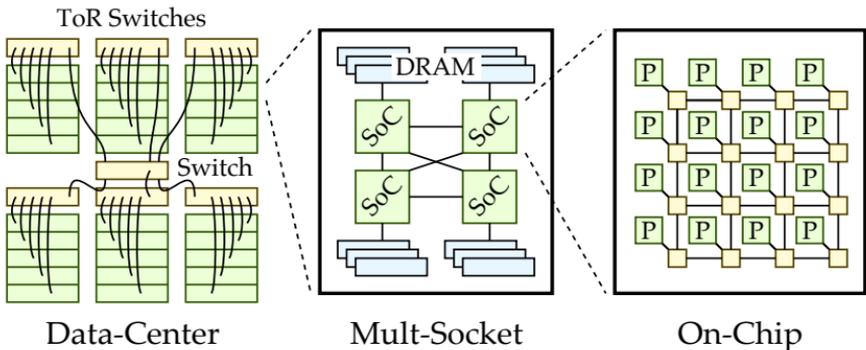


On-chip buffers are 1r1 FIFOs implemented using either a register file or SRAM

On-chip network technology constraints very different from off-chip technology constraints

## 1.3. Networks in Computer Architecture
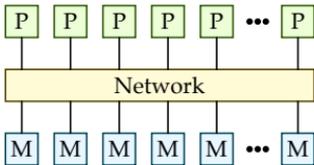


- Network Topology
    - Arrangement of channels and routers to interconnect sources and destinations
    - Roads = channels (implemented with cables, traces, wires)
    - Width of road = channel bw
    - Intersections = routers
    - Size intersection = router radix
- Network Routing
    - Path from source to destination along channels and routers
    - Short vs. long paths = minimal vs. non-minimal paths
    - Common vs. rare paths = channel congestion
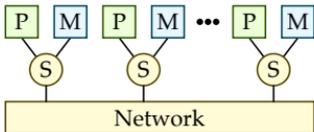- Topology is constrained by packaging (geography)



ToR Switches

Switch

Data-Center       Mult-Socket       On-Chip

**Network Transactions**

We will use processor-memory networks as a driving example throughout the course. Processor-memory networks allow many processors to perform memory read/write transactions on many memories.



Dance-Hall Organization



Integrated-Node Organization

- Message = Logical unit of data transfer provided by network interface

- Packet = Unit of routing within a network

- Flit = Smallest unit of resource allocation in channel/router (flow-control digit)

- Phit = Smallest unit of data processed by a channel/router (physical digit)

For the processor-memory networks we will primarily study:

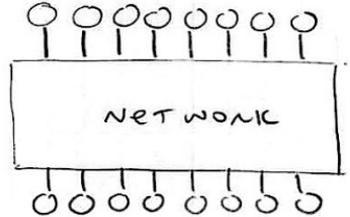$$\text{Message} = \text{Packet} = \text{Flit}$$

The processor will send memory request messages/packets over the network, and the memories will send memory response messages/packets back to the processor. In this context, we consider the network messages/packets to be the "transactions".

| Processors | : | Instructions |
|------------|---|--------------|
| Memories | : | Memory accesses |
| Networks | : | Network packets |

## 2. Network Topology

Network topology is the arrangement of channels and routers to interconnect sources and destinations. We will explore four different topology classes:

- Single-stage bus topologies
- Single-stage crossbar topologies
- Multi-stage butterfly topologies
- Mult-stage torus topologies
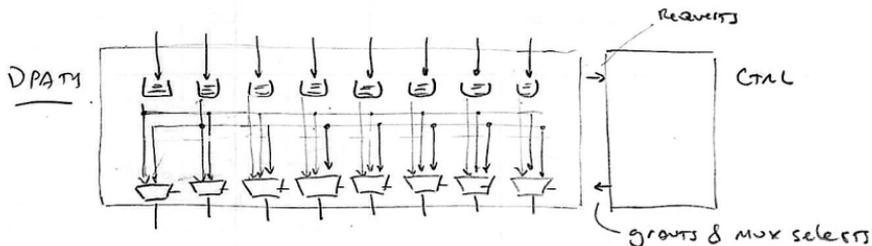


## 2.1. Single-Stage Bus Topology

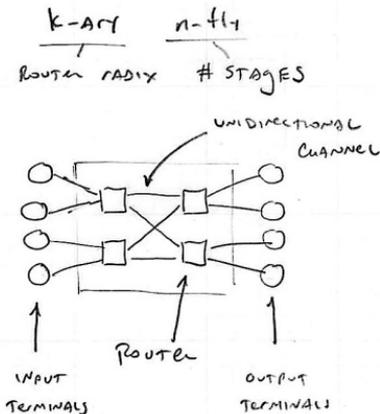**SWMR Bus**

**MWSR Bus**

**MWMR Bus**

**Integrating Multiple Buses**
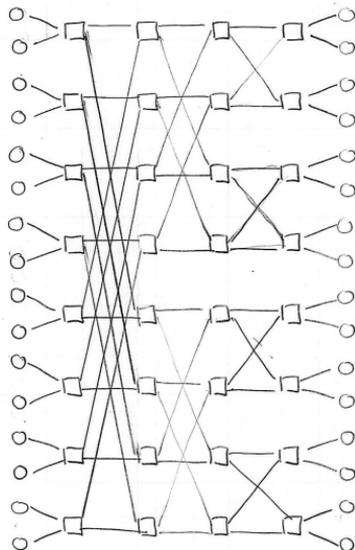
## 2.2. Single-Stage Crossbar Topology



- Single-stage topologies are difficult to scale in terms of cycle time, energy, and area

- Multi-stage topologies improve scalability but raise many other interesting challenges

## 2.3.  Multi-Stage Butterfly Topology
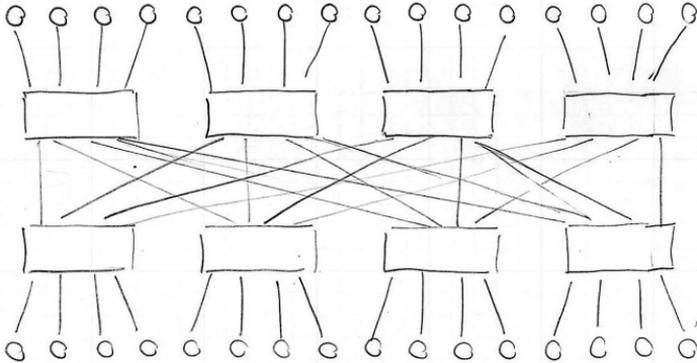


$k$-ARY $n$-FLY

ROUTER RADIX    # STAGES

UNIDIRECTIONAL CHANNEL

2-ARY 2-FLY

EACH ROUTER IS SIMILAR TO A 2×2 CROSSBAR

Router
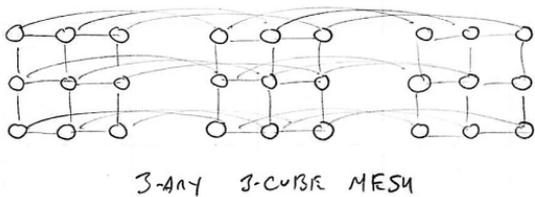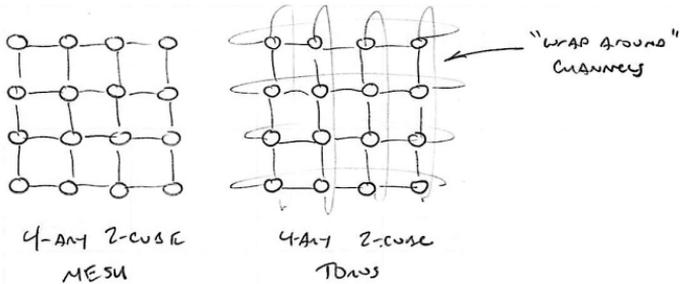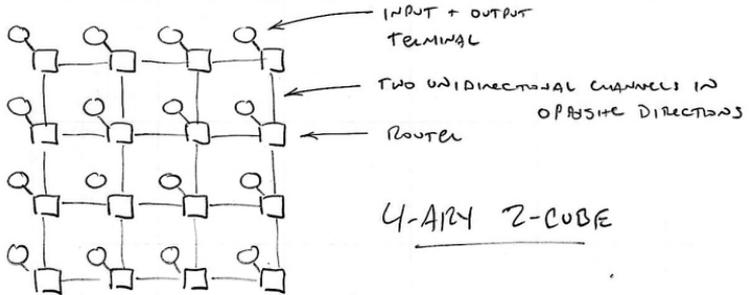
INPUT TERMINALS      OUTPUT TERMINALS

2-ARY 4-FLY

$2^4 = $ NODES

**Example Butterfly Topology: 4-ary 2-fly**



**Example Butterfly Topology: 3-ary 2-fly**

Sketch a 3-ary 2-fly. Use circles for the terminals, squares for the routers, and lines for one uni-directional channel.

## 2.4. Multi-Stage Torus Topology
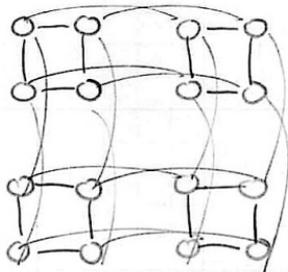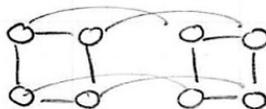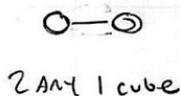
$k$-ARY  $N$-CUBE

NODES IN EACH DIMENSION

$n$ DIM GRID

INPUT + OUTPUT TERMINAL

TWO UNIDIRECTIONAL CHANNELS IN OPPOSITE DIRECTIONS

ROUTER

4-ARY 2-CUBE

"WRAP AROUND" CHANNELS

4-ARY 2-CUBE MESH

4-ARY 2-CUBE TORUS

3-ARY 3-CUBE MESH

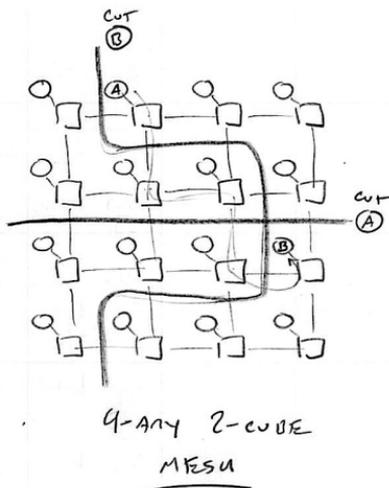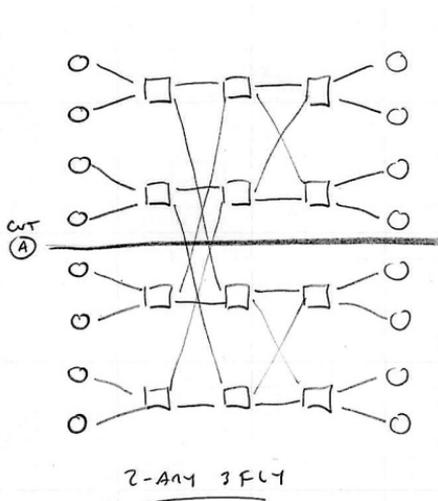**Example Torus Topology: 8-ary 1-cube torus**

Sketch a 8-ary 1-cube torus. Use circles for a terminal+router and lines for the two uni-directional channels in opposite directions.

**Constructing k=ary n-cube from k k-ary (n-1) cubes**



2 Ary 1 cube

2 Ary 2cube

2 Ary 3 cube

2 Ary 4 cube

HYPER CUDES

BINARY n-cube

## 2.5. Terminology



2-ANY 3FLY

4-ARY 2-CUBE MESH

**Nodes and Channels**

- Uni-directional channels in butterfly
- Bi-directional channels in mesh
- Indirect Network: node is either a terminal or router (butterfly)
- Direct Network: node combines a terminal and router (mesh)
- Channel parameters

  - $w_c$ = channel width (number of pins/wires)
  - $f_c$ = channel frequency
  - $t_c$ = channel latency
  - $l_c$ = channel length
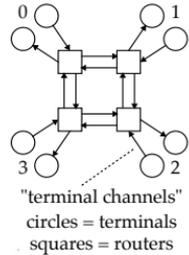  - $b_c$ = channel bandwidth ($w_c \times f_c$)

**Bisection Cuts**

- Cut: set of channels that partitions terminals into two sets

- Bisection Cut: cut that partitions terminals in two equal halves

- Min Bisection Cut ($B_c$): min channel count over all bisection cuts

- Bisection Bandwidth ($B_B$): min bandwidth over all bisection cuts

- For networks with uniform channel bandwidth: $B_B = B_c \times b_c$

- Bisection bandwidth is a good way to estimate global wiring resources (i.e., technology constraints)

- Example from previous butterfly/mesh topologies
    - Cut A on butterfly is a minimum bisection cut
    - Both cut A and B on mesh are bisection cuts
    - Cut A on mesh is a minimum bisection cut

**Paths**

- Channel Hop Count ($H_c$): number of channels on a path

- Channels from terminal to first router may or may not be included

- Router Hop Count ($H_r$): number of routers on a path

- Minimal Path: smallest hop count between two terminals

- Diameter ($H_{max}$): largest minimal path between all terminal pairs

- Average min hop count ($H_{min}$): average over all terminal pairs

- Example from previous butterfly/mesh topologies
    - Path from A to B on mesh: $H_r = 5$, $H_c = 4$
    - $H_{max,bfly} = 4$, $H_{max,mesh} = 8$
    - $H_{min,bfly} = 4$

Calculating $H_{min}$ usually requires enumerating the minimal hop count for every source/destination pair. The book states that $H_{min}$ for a torus with even $k$ is $nk/4$ and for a mesh with even $k$ is $nk/3$. Where does this come from?



"terminal channels"
circles = terminals
squares = routers

CALCULATING $H_{min, R}$

src

|     | 0 | 1 | 2 | 3 |
|-----|---|---|---|---|
| 0   | 1 | 2 | 3 | 2 |
| 1   | 2 | 1 | 2 | 3 |
| 2   | 3 | 2 | 1 | 2 |
| 3   | 2 | 3 | 2 | 1 |

Dest

$1 \times 4 = 4$
$2 \times 8 = 16$
$3 \times 4 = 12$

including $i \rightarrow i$     $32/16 = 2$
excluding $i \rightarrow i$     $28/12 = 2.3$

CALCULATING $H_{min, C}$ WITH TERMINAL CHANNELS

src

|     | 0 | 1 | 2 | 3 |
|-----|---|---|---|---|
| 0   | 2 | 3 | 4 | 3 |
| 1   | 3 | 2 | 3 | 4 |
| 2   | 4 | 3 | 2 | 3 |
| 3   | 3 | 4 | 3 | 2 |

Dest

$2 \times 4 = 8$
$3 \times 8 = 24$
$4 \times 4 = 16$

including $i \rightarrow i$     $40/16 = 2.5$
excluding $i \rightarrow i$     $40/12 = 3.3$

CALCULATING $H_{min, C}$ WITHOUT TERMINAL CHANNELS

src

|     | 0 | 1 | 2 | 3 |
|-----|---|---|---|---|
| 0   | 0 | 1 | 2 | 1 |
| 1   | 1 | 0 | 1 | 2 |
| 2   | 2 | 1 | 0 | 1 |
| 3   | 1 | 2 | 1 | 0 |

Dest

$1 \times 8 = 8$
$2 \times 4 = 8$

including $i \rightarrow i$     $16/16 = 1$
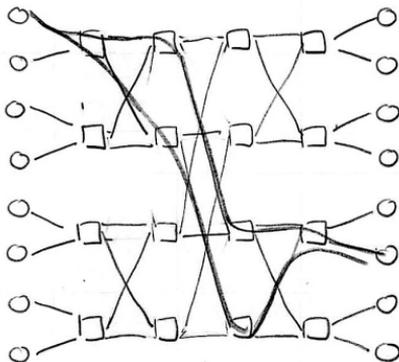excluding $i \rightarrow i$     $16/12 = 1.3$

EQUATION FROM Book IS $\frac{nk}{3}$ for even $n$

$\frac{nk}{3} = \frac{2 \cdot 2}{3} = 1.33$     THIS IS for $H_{min,C}$
WITHOUT TERMINAL CHANNELS
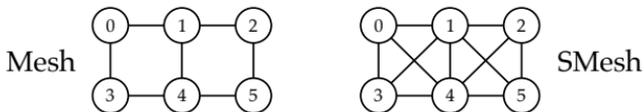+ ignoring $i$ sending to itself

**Path Diversity**

- Path diversity is number of paths between any two terminals
- Mesh has high path diversity, butterfly has no path diversity



ADDING EXTRA
BFLY STAGES
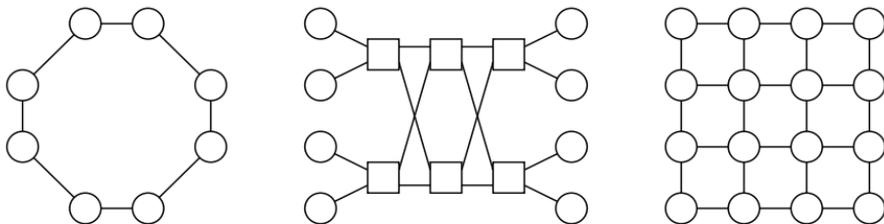↑ PAM DIVERSITY

**Example Topologies**



Mesh          SMesh

Assume $b_c$ for mesh is 32 bits/cycle, while $b_c$ for smesh is 16 bits/cycle. smesh has reduced channel bandwidth due to increased number of channels, assumes specific technology constraint. Calculate the following parameters for each topology.

|       | $B_C$ | $B_B$ | $H_{max}$ | $H_{min}$ |
|-------|-------|-------|-----------|-----------|
| mesh  |       |       |           |           |
| smesh |       |       |           |           |

# 3. Network Routing

Network routing is the choice of path from source to destination along channels and routers. Routing algorithms can be oblivious (not factor in network state) or adaptive (factor in network state), and can also be deterministic or non-deterministic.
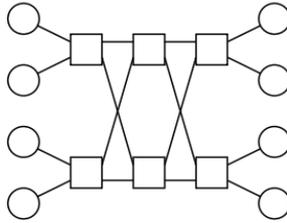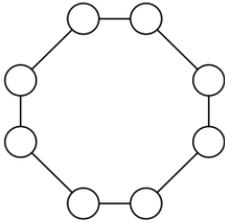
## 3.1. Oblivious Deterministic Routing



- Ring: always use minimal path, equidistant choose CW
- Butterfly: use destination to choose middle router
- Mesh: route in X first, then route in Y (Dimension-Ordered-Routing)

## 3.2. Oblivious Non-Deterministic Routing



- Ring: randomly choose CW vs CCW
- Butterfly: randomly choose middle router
- Mesh: randomly choose between XY-DOR and YX-DOR (O1-TURN)
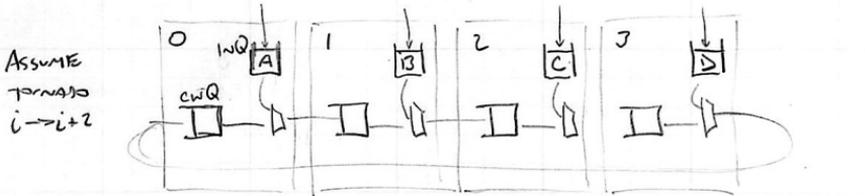
## 3.3. Adaptive Routing



- Ring: look at queues, choose direction with least congestion
- Butterfly: look at queues, choose middle router with least congestion
- Mesh: look at queues during each hop to choose X vs Y

## 3.4.  Deadlock

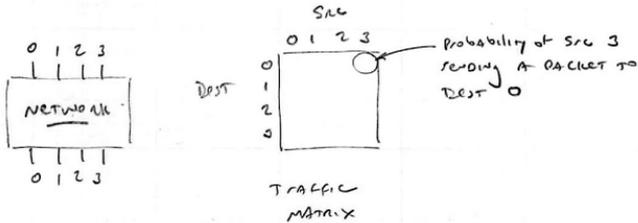CYCLIC DEPENDENCY IN "WAITS FOR" AND "HOLDS" RELATION



ASSUME
TORNADO
$i \rightarrow i+2$

Cyclic Dependency

ACTORS

Resources

- DEADLOCK AVOIDANCE vs DEADLOCK DETECTION / RECOVERY
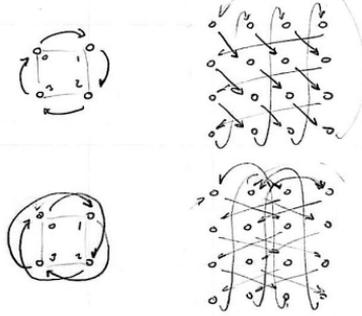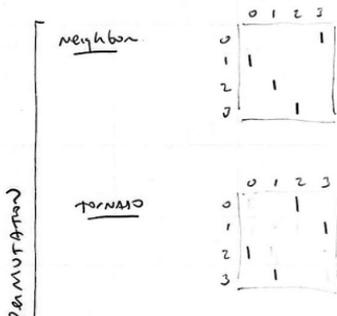
CONTRAST TO LIVELOCK
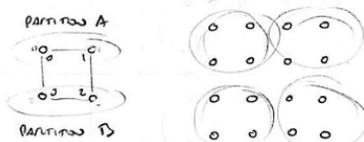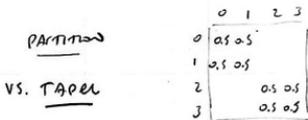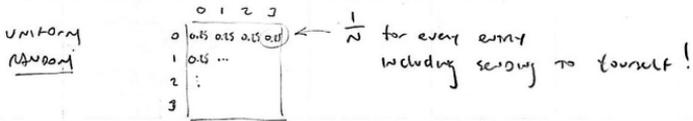  CHAOS ROUTING ALGORITHM

# 4. Analyzing Network Performance

Similar to analyzing processors and memories, we will use simple first-order equations to help build intuition about throughput and latency trade-offs in networks.
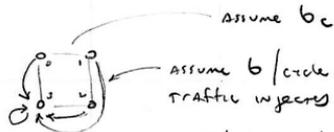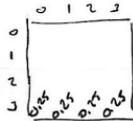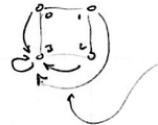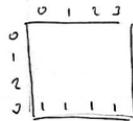
## 4.1. Traffic Patterns

HOT SPOT

assume $b_c$

assume $b$ /circle Traffic injects

$0.256 + 0.256 = 0.56$
going over $2 \to 3$
channel

INADMISSIBLE TRAFFIC PATTERNS

oversubscribed
Hotspot

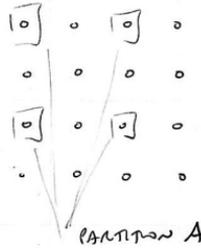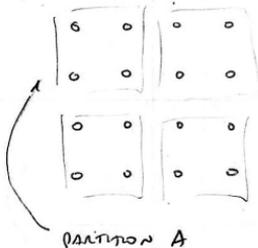$1b + 1b = 2b$
going over this $2 \to 3$
channel

CHANNEL $2 \to 3$ oversubscribed
but so is ejection channel
to output terminal at node 3

logical to physical mapping

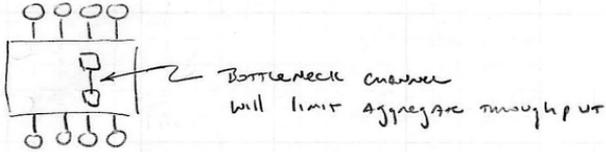ASSUME PARTITION TRAFFIC PATTERN
How are logical src/dest IDs in Traffic Pattern
mapped to physical terminal IDs?

PARTITION A               PARTITION A

mapping can turn any logical permutation pattern into
any other pattern.

gives permutation pattern usually assumes a specific mapping

22

## 4.2.  Ideal Throughput



Bottleneck channel will limit aggregate throughput

Channel Load $(\gamma_c)$ is amount of traffic that crosses channel $c$ if each input injects one unit of traffic according to given traffic pattern

Ex Traffic Pattern
Src $i \rightarrow$ Dest $i+1$



Channel load ranges from 0-2

Max channel load $(\gamma_{MAX})$ is 2. These are the bottleneck channels that will limit throughput

Alternative way to think about $\gamma$.

Channel load is ratio of BW demanded from channel to input BW injected by one terminal
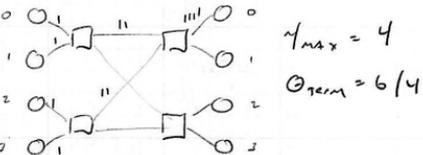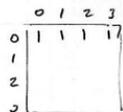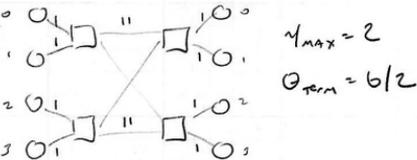
IDEAL THROUGHPUT $\quad \Theta_{Term} = \dfrac{b_c}{\gamma_{MAX}} \longleftarrow$ BW of bottleneck channel

$\longrightarrow$ MAX CHANNEL LOAD

IN PREVIOUS EXAMPLE

$$\Theta_{Term} = \frac{b}{\gamma_{MAX}} = \frac{b}{2} \qquad \Theta_{TOT} = N \cdot \frac{b}{\gamma_{MAX}} = 8 \cdot \frac{b}{2} = 4b$$

SMALL BFLY EXAMPLES



each TICK = 1/4 UNIT

$\gamma_{MAX} = 1$

$\Theta_{Term} = b$

$\gamma_{MAX} = 1$

$\Theta_{Term} = b$

$\gamma_{MAX} = 2$

$\Theta_{Term} = b/2$

$\gamma_{MAX} = 4$

$\Theta_{Term} = b/4$

More generally for uniform random traffic

- on average with uniform random traffic, half the traffic crosses the bisection

- N total units of traffic, N/2 cross bisection

- Ideal routing will evenly balance load across bisection channels (this was an issue in ring example)

$$\gamma_{max} = \frac{N/2}{B_c} = \frac{N}{2B_c}$$

$$\Theta_{Term} = \frac{6}{N/2B_c} = \frac{2bB_c}{N} = \frac{2B_B}{N}$$

$$\Theta_{Tot} = N \frac{2B_B}{N} = 2B_B$$

## 4.3. Zero-Load Latency



3L b/cycle = packet length L
8 b/cycle = $b_c$  Phit

Serialize packet into four phits

Pkt 0     S S S S
Pkt 1           S S S S

Serializer
Deserializer
Router Pipeline : R0 R1 R2
Link Traversal : L0 L1

PKT 1   HEAD PKIT → S R0 R1 R2 L0 L1 R0 R1 R2 D
        BODY PKIT     S R0 R1 R2 L0 L1 R0 R1 R2 D
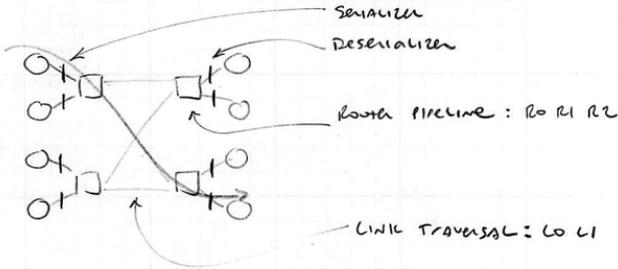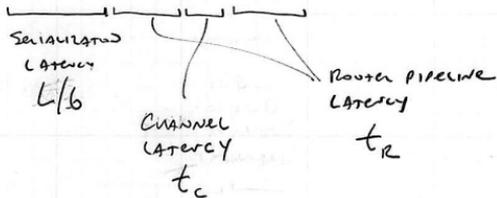        BODY PKIT     S R0 R1 R2 L0 L1 R0 R1 R2 D
        TAIL PKIT     S R0 R1 R2 L0 L1 R0 R1 R2 D →

Serialization Latency $\frac{L}{b}$

Channel Latency $t_c$

Router Pipeline Latency $t_R$

$$T = T_{HEAD} + \frac{L}{b}$$

→ serialization Latency

↳ head pkit Latency
  includes $t_c$, $t_\sim$, hop count, + contention

$$T_\phi = H_R t_R + H_c t_c + \frac{L}{b}$$

Latency Due to Router hops

Latency Due to channel hops

Serialization Latency

Zero Load Latency (no contention)

FOUR WAYS TO IMPROVE LATENCY

$$T_\phi = H_R t_R + H_c t_c + \frac{L}{b}$$

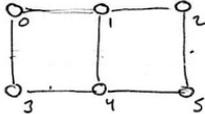Wider channels or shorter msgs

Shorter Routes        Faster routes        Faster Channels

AVG LATENCY VS OFFERED BW



Avg Latency

new IDEAL ROUTING + Flow control

$T_\phi$

offered BW
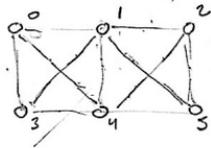
$\Theta_{TOT}$

## 4.4.  Comparing Topologies



MESH                                    SMESH

ASSUME :

$L = 128\,b$

$b_c$ for mesh is 32 b/cycle

$b_c$ for smesh is 16 b/cycle  $\leftarrow$ reduces channel bandwidth due to increases # channels assuming constant resources

$t_r = 3$

$t_c = 1$.

WHICH TOPOLOGY CAN ACHIEVE HIGHER IDEAL THROUGHPUT UNDER UNIFORM RANDOM TRAFFIC?
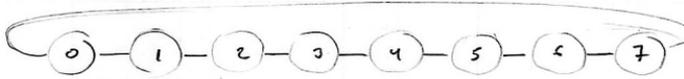
(a)   $\Theta_{term, mesh} > \Theta_{term, smesh}$
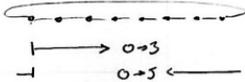
(b)   $\Theta_{term, mesh} < \Theta_{term, smesh}$

(c)   $\Theta_{term, mesh} == \Theta_{term, smesh}$

CALCULATE ZERO LOAD LATENCY UNDER UNIFORM RANDOM TRAFFIC

## 4.5. Comparing Routing Algorithms
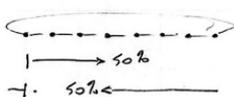


- GREEDY : Always use minimal path, equidistant choose randomly



→ 0→3

→|     0→5 ←——

- UNIFORM RANDOM : RANDOMLY PICK DIRECTION



|——→ 50%

→|.   50%←——

for $0 \to 3$
50% TAKE 4 hops
50% TAKE 6 hops

- Weighted RANDOM : RANDOMLY PICK DIRECTION BUT
WEIGHT PROBABILITY BY DISTANCE



|——→ 63%

→|   37%←——

for $0 \to 3$
5/8 TIME TAKE 3 hop path
3/8 TIME TAKE 5 hop path

PROBABILITY OF TAKING SHORT PATH IS

$$\frac{N - H_{MW,R} + 1}{N}$$

- ADAPTIVE : LOOK AT QUEUES IN EITHER DIRECTION
SEND IN DIRECTION OF QUEUE THAT HAS MOST FREE
ENTRIES.

Do NOT CHANGE DIRECTION AFTER INITIAL CHOICE

- Evaluate tornado traffic pattern on the 8-node ring
- Recall that in tornado, node $i$ sends to $i + ((N-1)/2) \bmod N$
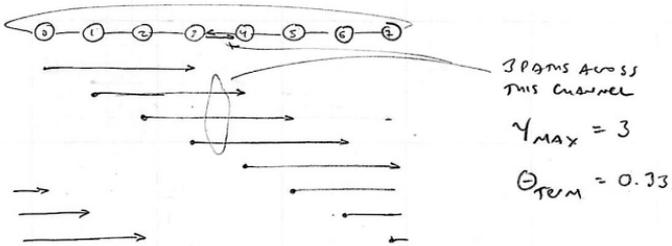
greedy Algo



3 PATHS ACROSS THIS CHANNEL

$\gamma_{MAX} = 3$

$\Theta_{TorN} = 0.33$

CLOCKWISE CHANNELS NOT USED AT ALL!
POOR LOAD BALANCING

$T_0 = 4$

URANDOM ALGO



STILL 3 PATHS CROSS THIS CHANNEL BUT EACH PATH ONLY CARRIES 0.5 UNITS OF TRAFFIC

$\gamma_{MAX, ccw} = 1.5$

WHAT ABOUT CLOCKWISE CHANNELS THOUGH?
5 PATHS CROSS EACH CW CHANNEL
EACH IS 0.5 UNIT OF TRAFFIC

$\gamma_{MAX, cw} = 2.5 > \gamma_{MAX, ccw} = 1.5$      $\Theta_{Torn} = 0.4$

NOW THE CLOCKWISE CHANNELS ARE THE BOTTLENECK!

$T_\phi = 0.5 \times 4 + 0.5 \times 6 = 5$

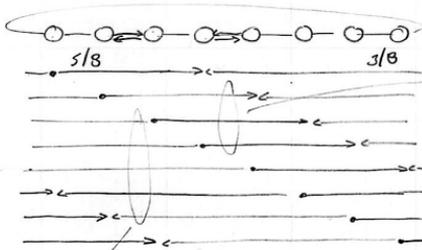## Weighted Random Algo

SAME NUMBER OF PATHS CROSS CW AND CCW CHANNELS AS
BEFORE EXCEPT NOW WITH DIFFERENT AMOUNTS OF
TRAFFIC PR PATH.



3 PATHS EACH WITH
5/8 TRAFFIC

$$\gamma_{max, ccw} = 3 \cdot \frac{5}{8} = 1.875$$

5 PATHS EACH WITH 3/8 TRAFFIC

$$\gamma_{max, cw} = 5 \cdot \frac{3}{8} = 1.875$$

$$\gamma_{max, cw} = \gamma_{max, ccw} \qquad \text{BALANCED !}$$

$$\Theta_{TEM} = \frac{1}{1.875} = 0.53$$

$$T_0 = \frac{5}{8} \times 4 + \frac{3}{8} \, 6 = 2.5 + 2.25 = 4.75$$
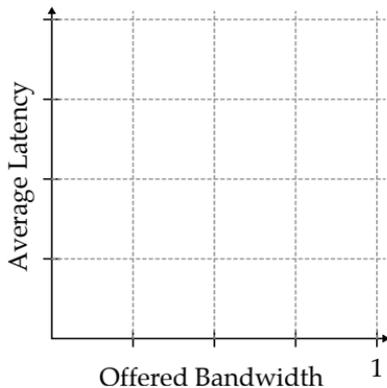
## ADAPTIVE

MINIMAL LATENCY AT LIGHT LOAD
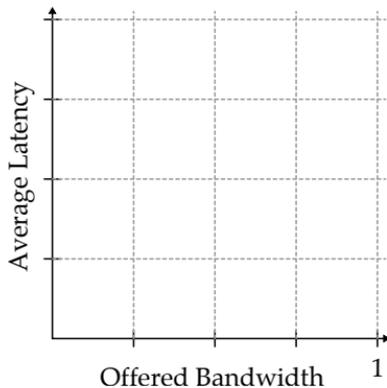MAX THROUGHPUT AT HIGH LOAD

Let's fill in the following table to summarize the performance of the four routing algorithms on the 8-node given the tornado and uniform random traffic patterns.

| Routing Algo | Tornado | | Uniform Random | |
|---|---|---|---|---|
| | $\Theta_{term}$ | $T_0$ | $\Theta_{term}$ | $T_0$ |
| Greedy | | | 1.00 | 3.0 |
| Uniform Random | | | 0.57 | 4.5 |
| Weighted Random | | | 0.76 | 3.1 |
| Adaptive | | | 1.00 | 3.0 |

**Tornado Traffic Pattern**
Bandwidth vs. Latency



**Uniform Random Traffic Pattern**
Bandwidth vs. Latency

**Activity: Routing on butterfly with extra stage**

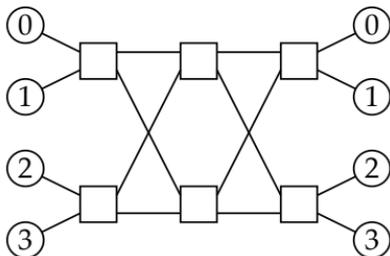Consider two routing algorithms for a 2-ary 2-fly with an extra stage:

- Destination Tag Routing: use destination to choose middle router
- Random Middle Router: randomly choose middle router

Compare these two routing algorithms in terms of ideal terminal throughput ($\Theta_{term}$) and zero-load latency ($T_0$) for the following permutation traffic pattern.

- src $0 \rightarrow$ dest 1
- src $1 \rightarrow$ dest 0
- src $2 \rightarrow$ dest 3
- src $3 \rightarrow$ dest 2

*Hint: Use the tick-mark method to calculate the max channel load for each routing algorithm.*

**Destination-Tag Routing**          **Random Middle Routing**