# ECE 669

# Parallel Computer Architecture

## Lecture 25

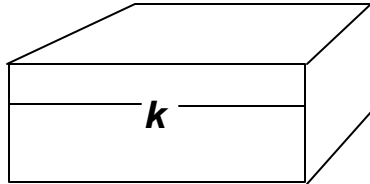## Final Exam Review

# Bandwidth & Latency

**Direct**

**Latency**

- **Average latency $k_d$:**



$$nk \;-\; \text{worst case}$$

$$n \frac{(k-1)}{2} \;-\; torus \;-\; \textbf{1 dir channels}$$

$$\sim \; n \frac{k}{4} \;-\; torus \;-\; 2 \text{ dir channels}$$

$$\sim \; n \frac{k}{3} \;-\; \text{no end conns} \;-\; 2 \text{ dir channels}$$

- **Bandwidth per node - more complex**
  - $\dfrac{1}{B}$ **if all near neighbor messages**
  - **If average travel dist then?**
  - **Let**

$$\text{Avg DIST} \;=\; \frac{nk}{3}$$

$$\text{Msg size} \;=\; B$$

# Analogy

- **If each student takes 8 years to graduate**
- **And if a Prof. can support 10 students max at any time**

**How many new students can Prof. take on in a year?** $\quad 8x = 10$

**Each year take on** $\dfrac{10}{8}$

° **Similarly:**

- **Network has $Nn$ channels**
- **# of flits it can sustain = $Nn$**
- **# of msgs it can concurrently sustain= $\dfrac{Nn}{B}$**
- **Each msg flit uses $\dfrac{nk}{3}$ channels to <u>dest</u>**
- **So**

$$N \times \frac{nk}{3} = \frac{Nn}{B}$$

or

$$BW \text{ per node} = x = \frac{3}{kB}$$

# Another way of getting *BW* is:

- **Max # msgs in net at any time** $= \dfrac{Nn}{B}$

- **These take** $= \dfrac{kn}{3}$ **cycles to get delivered, during which time no new mgs can get in**

- **I.e. we can inject** $\dfrac{Nn}{B}$ msgs every $\dfrac{kn}{3}$ cycles

    **or # injected per node per cycle**

    $$= \dfrac{\frac{Nn}{B}}{3} \cdot \dfrac{3}{kn} \cdot \dfrac{1}{N}$$

    $$= \dfrac{3}{Bk}$$

° **Note: We have not considered contention thus far.**

° **In practice, latency shoots up much before we achieve the theoretical due to contention.**

# Deriving *U* is not easy

° **Notice**

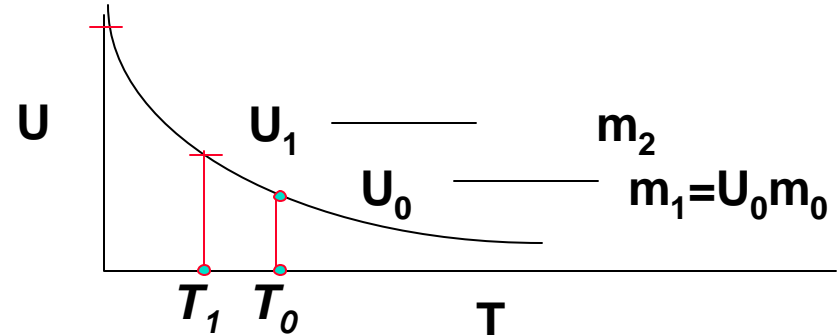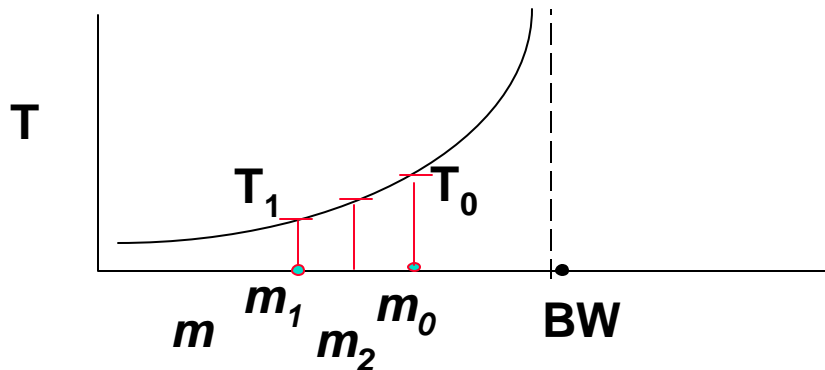*m* = **Probability of a message on a** useful **processor cycle**

$\mathbf{m_{eff}} = \boldsymbol{m} \cdot \boldsymbol{U}$ = **probability of msg on any cycle**

**T** = $\boldsymbol{f}(\boldsymbol{m_{eff}})$ = **network delay as a function of** *m*

$$U = \frac{1}{1 + mT}$$

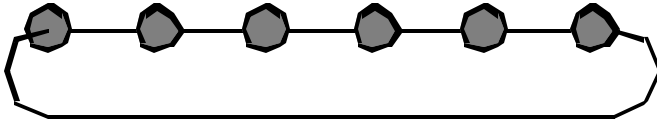**or # of useful processor cycles depends on** *T* **and** $m_{eff}$
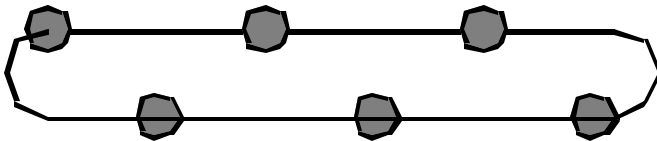
° **Cyclic dependence!**

# Linear Arrays and Rings
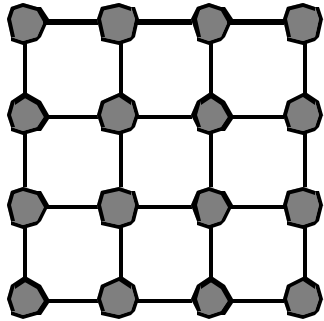
Linear Array

Torus

Torus arranged to use short wires

° **Linear Array**

- **Diameter?**
- **Average Distance?**
- **Bisection bandwidth?**
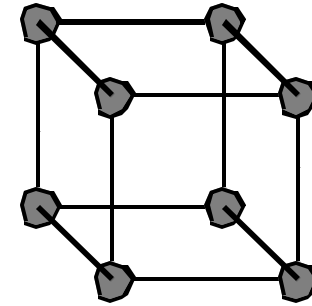- **Route A -> B given by relative address R = B-A**
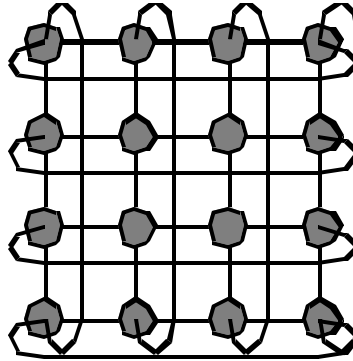
° **Torus?**

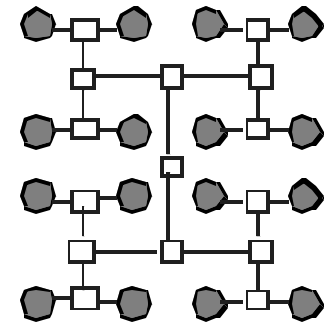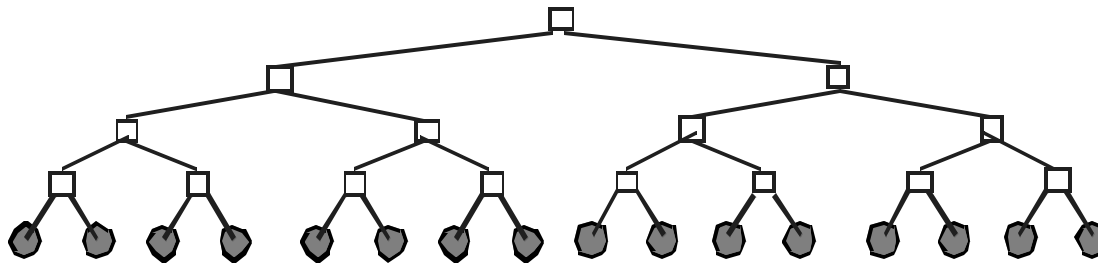° **Examples: FDDI, SCI, KSR1**

# Multidimensional Meshes and Tori



**2D Grid**

**3D Cube**

° *n*-dimensional *k*-ary mesh: $N = k^n$

  • $k = \sqrt[n]{N}$

  • described by *n*-vector of radix k coordinate

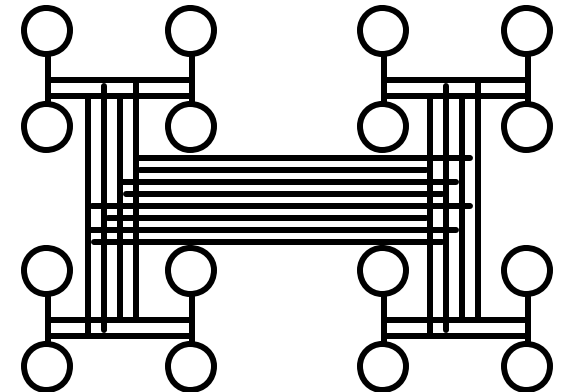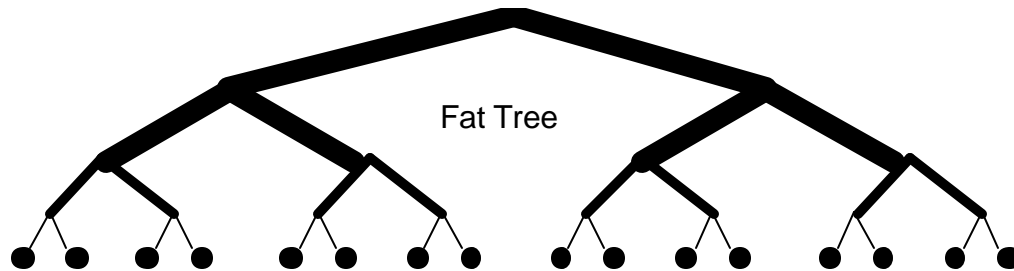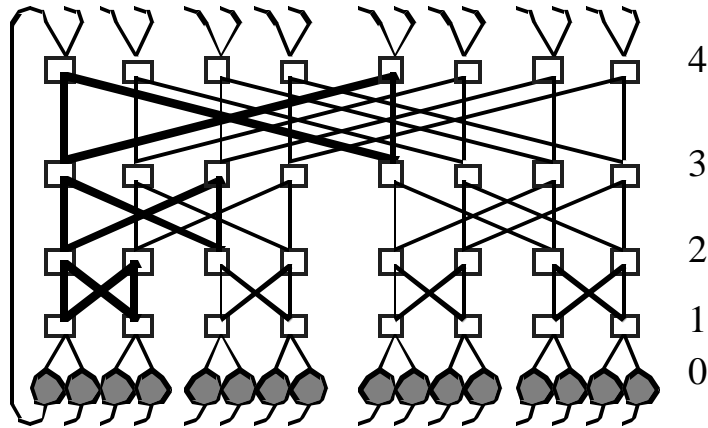° *n*-dimensional *k*-ary torus (or *k*-ary *n*-cube)?

# Trees



- ° **Diameter and ave distance logarithmic**
  - • **k-ary tree, height $d = \log_k N$**
  - • **address specified d-vector of radix k coordinates describing path down from root**

- ° **Fixed degree**

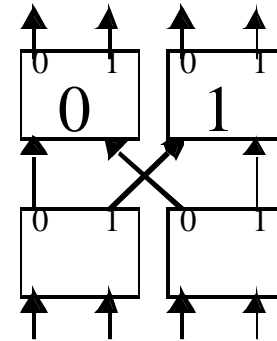- ° **H-tree space is O(N) with O($\ddot{I}$N) long wires**

- ° **Bisection BW?**

# Fat-Trees

° **Fatter links (really more of them) as you go up, so bisection BW scales with N**

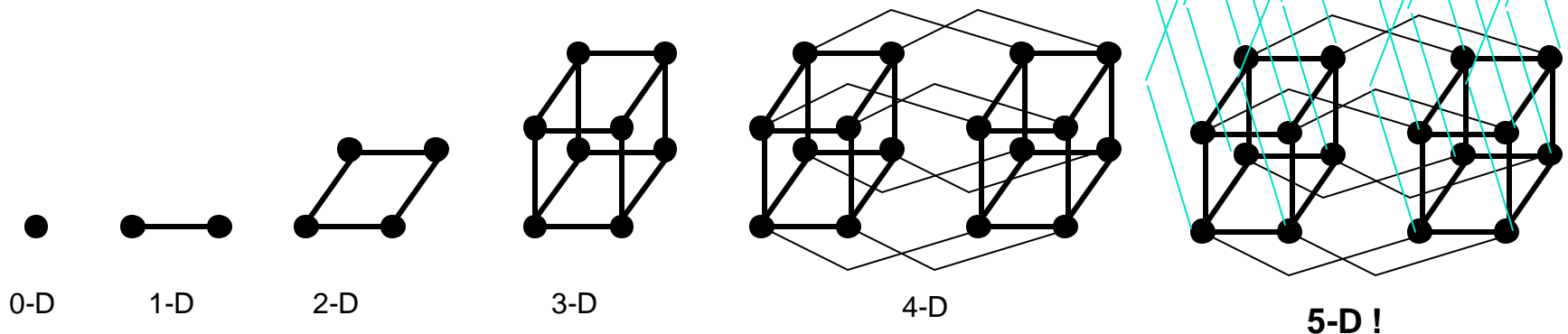Fat Tree

# Butterflies



**16 node butterfly**



**building block**

° **Tree with lots of roots!**

° **N log N   (actually N/2 x logN)**

° **Exactly one route from any source to any dest**

° **Bisection N/2**

# Hypercubes

- ° **Also called binary n-cubes.   # of nodes = N = $2^n$.**

- ° **O(logN) Hops**

- ° **Good bisection BW**

- ° **Complexity**

  - • **Out degree is n = logN**

  **correct dimensions in order**

  - • **with random comm. 2 ports per processor**

0-D        1-D        2-D        3-D        4-D        **5-D !**

# Toplology Summary

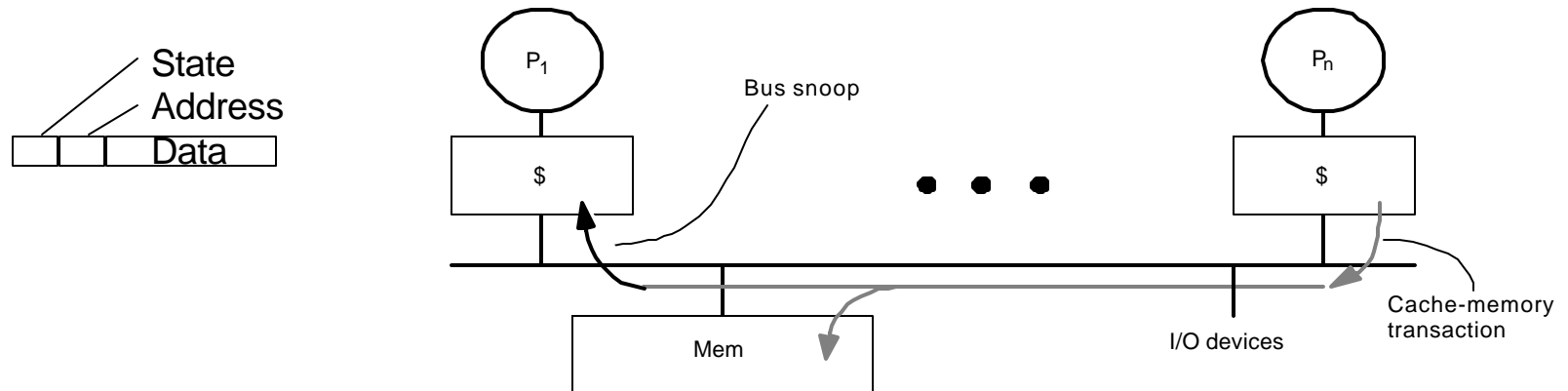| Topology | Degree | Diameter | Ave Dist | Bisection | D (D ave) @ P=1024 |
|---|---|---|---|---|---|
| 1D Array | 2 | N-1 | N / 3 | 1 | huge |
| 1D Ring | 2 | N/2 | N/4 | 2 | |
| 2D Mesh | 4 | $2 (N^{1/2} - 1)$ | $2/3 N^{1/2}$ | $N^{1/2}$ | 63 (21) |
| 2D Torus | 4 | $N^{1/2}$ | $1/2 N^{1/2}$ | $2N^{1/2}$ | 32 (16) |
| k-ary n-cube | 2n | nk/2 | nk/4 | nk/4 | 15 (7.5)  @n=3 |
| Hypercube | n =log N | | n | n/2 | N/2      10 (5) |

° **All have some "bad permutations"**

  - many popular permutations are very bad for meshs (transpose)

  - ramdomness in wiring or routing makes it hard to find a bad one!
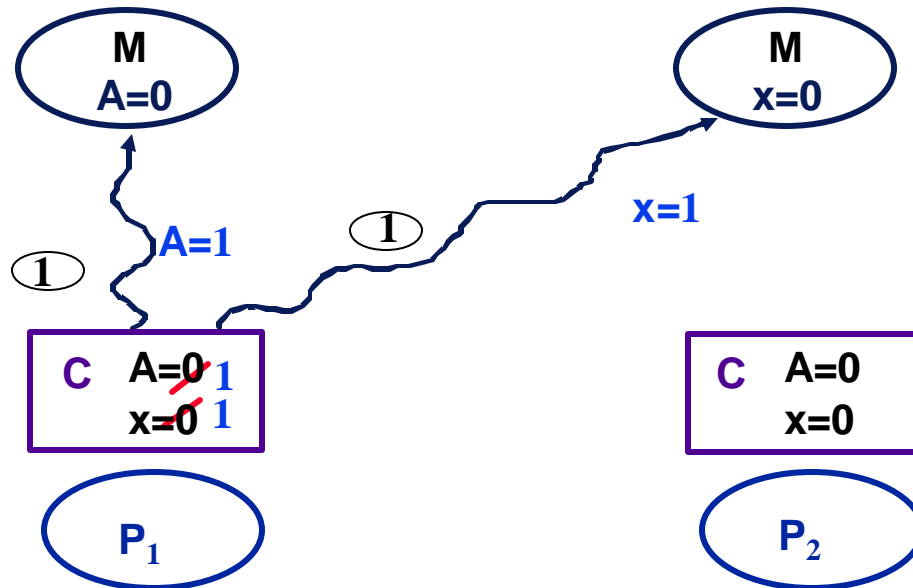
# Example Cache Coherence Problem



- **Processors see different values for u after event 3**
- **With write back caches, value written back to memory depends on happenstance of which cache flushes or writes back value when**
    - **Processes accessing main memory may see very stale value**
- **Unacceptable to programs, and frequent!**

# Snoopy Cache-Coherence Protocols



- **Bus is a broadcast medium & Caches know what they have**

- **Cache Controller "snoops" all transactions on the shared bus**

  - **relevant transaction if for a block it contains**
  - **take action to ensure coherence**
    - **invalidate, update, or supply value**
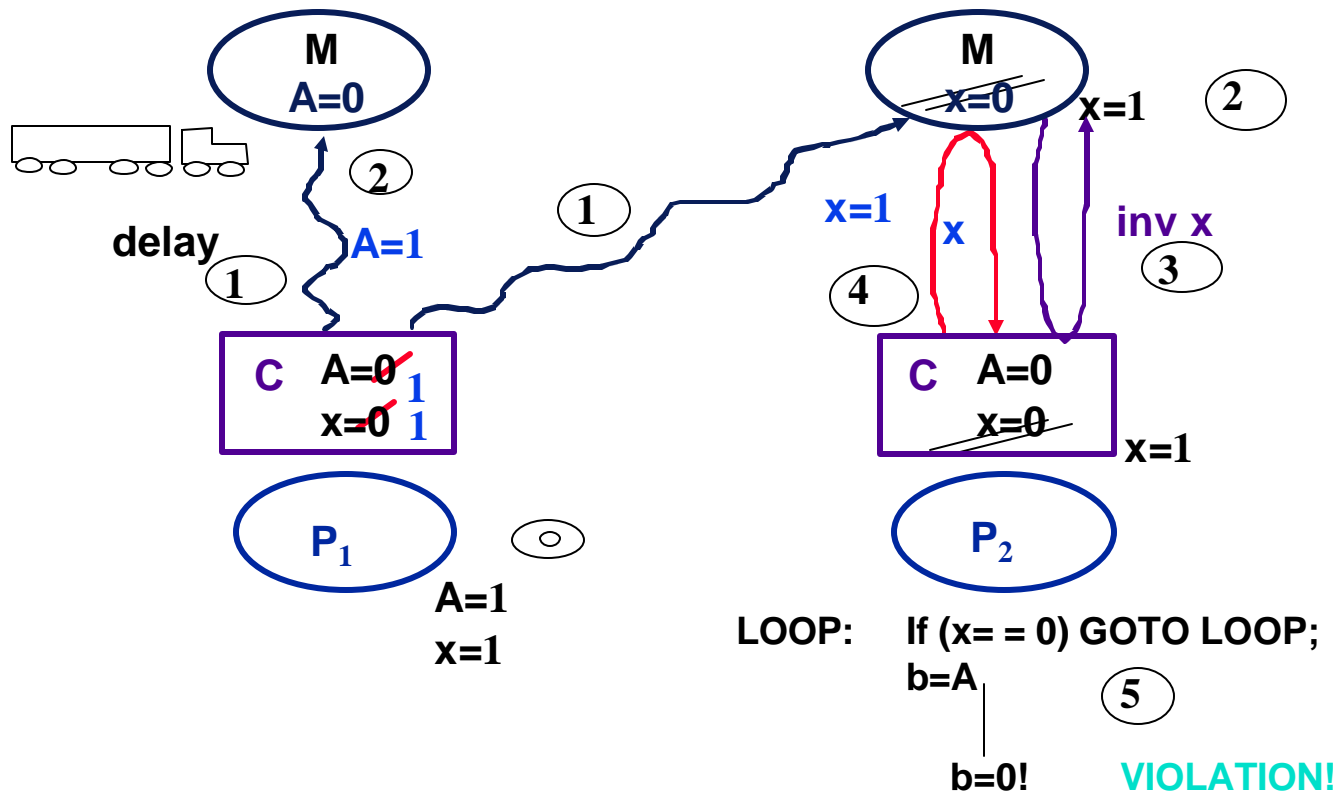  - **depends on state of the block and the protocol**

# Does caching violate this model?



**M**
**A=0**

**M**
**x=0**

①

**A=1**

①

**x=1**

**C** **A=0** ~~**1**~~
**x=0** **1**

**C** **A=0**
**x=0**

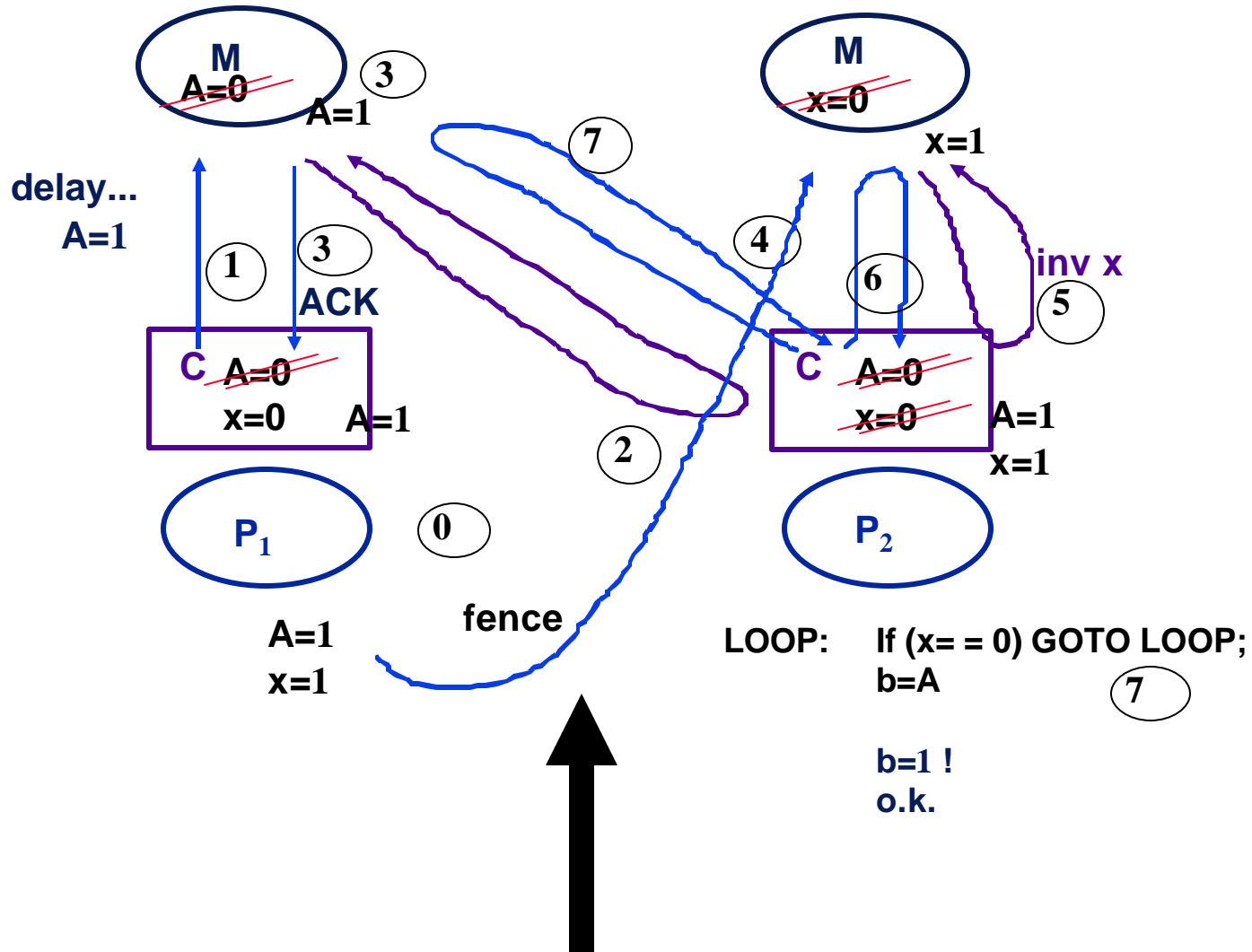**P₁**

**P₂**

LOOP:    If (x= = 0) GOTO LOOP;
b=A

**If b = = 0 at the end, sequential consistency is violated**
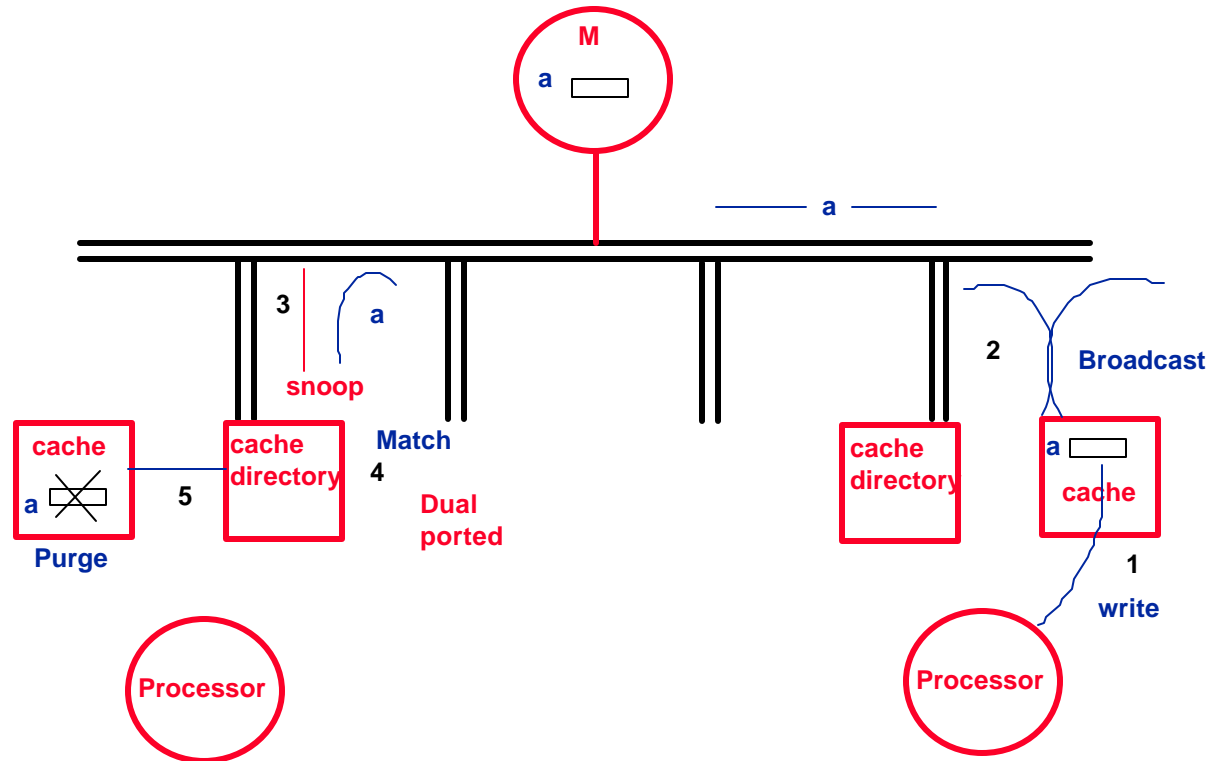
# Does caching violate this model?



If b = = 0 at the end, sequential consistency is violated
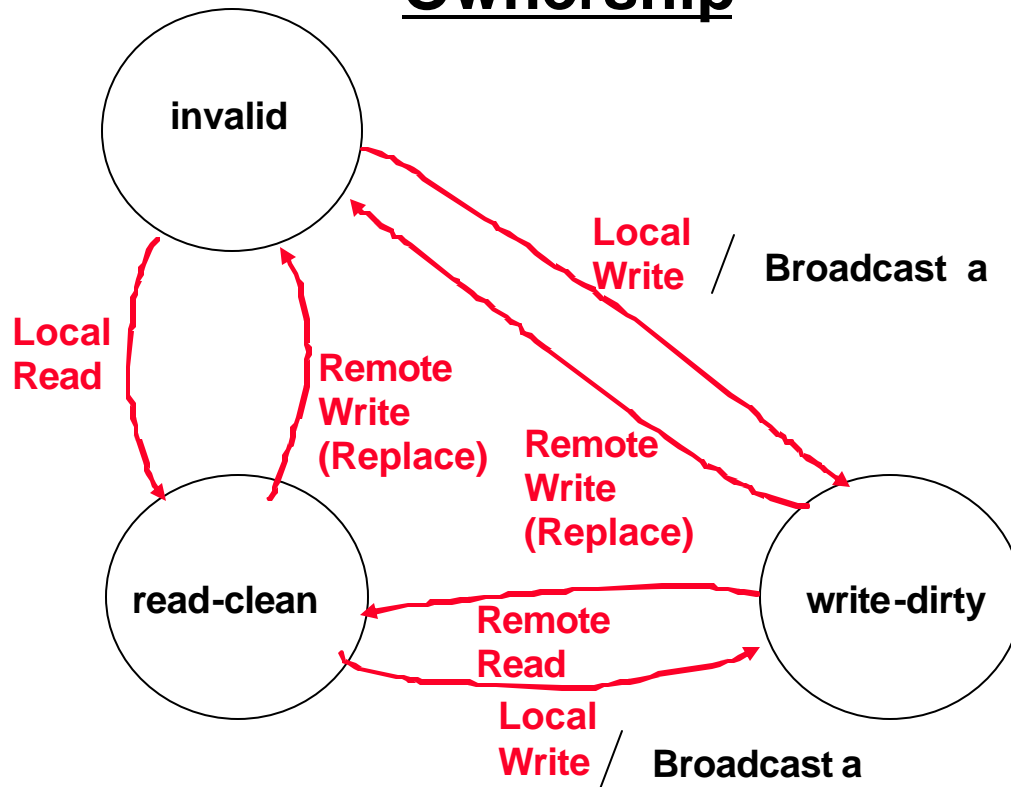
# Does caching violate this model?

# Coherence in small machines: Snooping Caches



- **Broadcast address on shared write**

- **Everyone listens (snoops) on bus to see if any of their own addresses match**

- **How do you know when to broadcast, invalidate...**
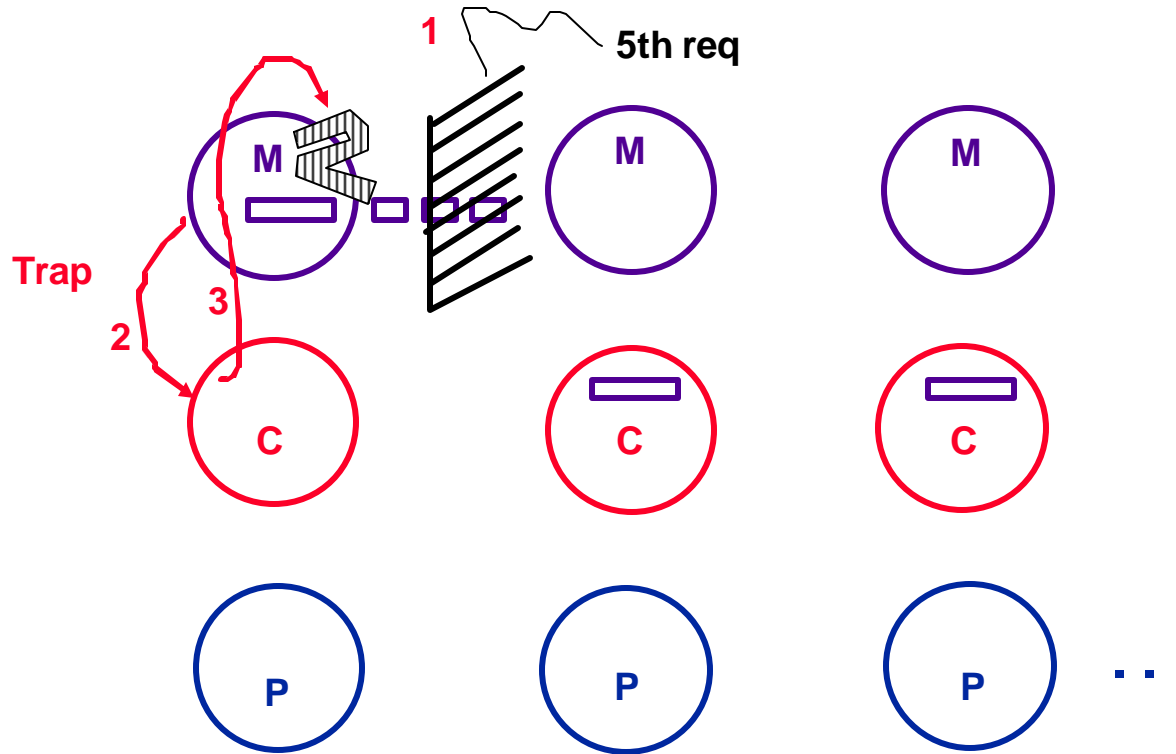    - **State associated with each cache line**

# State diagram for ownership protocols

## <u>Ownership</u>



- **In ownership protocol: writer owns exclusive copy**
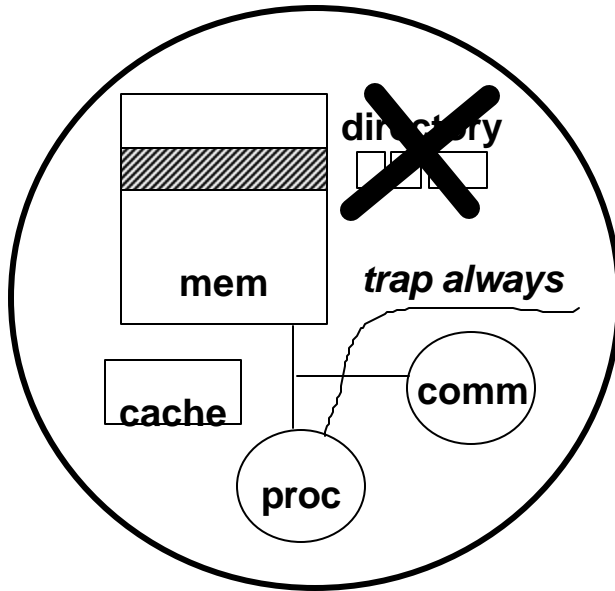- For each shared data cache block
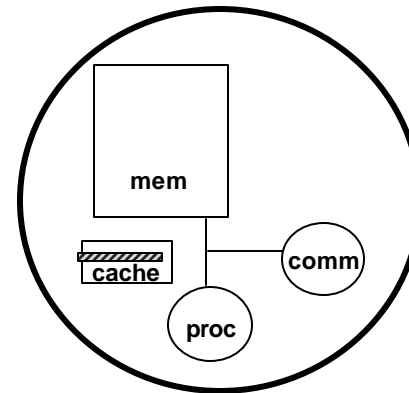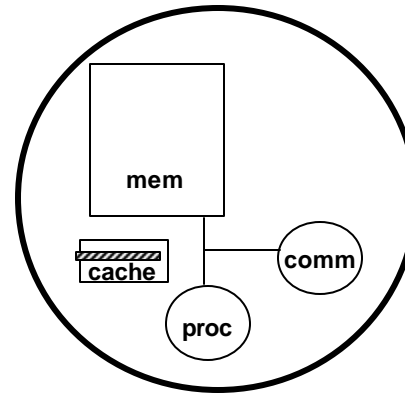
# Network



## ° **LimitLESS directories:**

**Limited directories Locally Extended through Software Support**

- **Trap processor when 5th request comes**
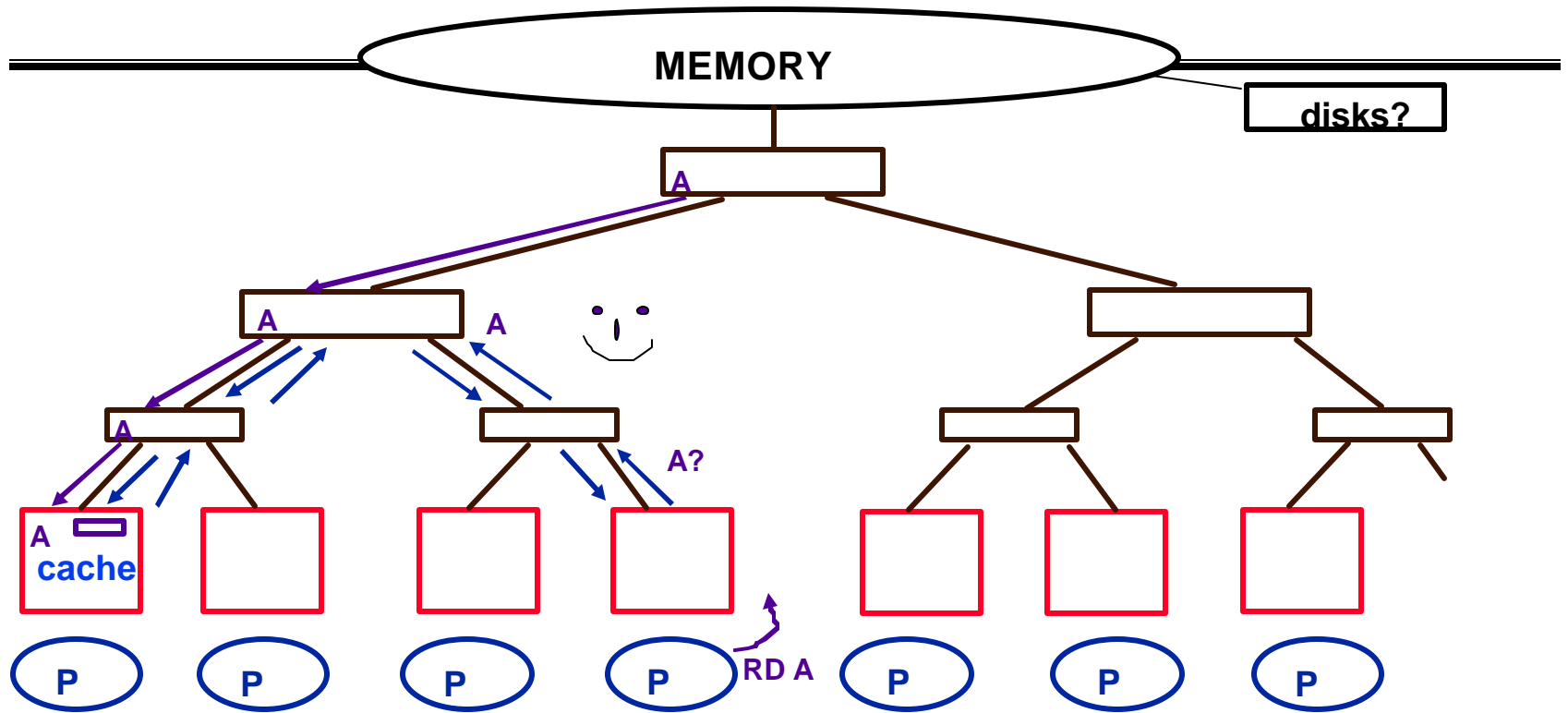- **Processor extends directory into local memory**
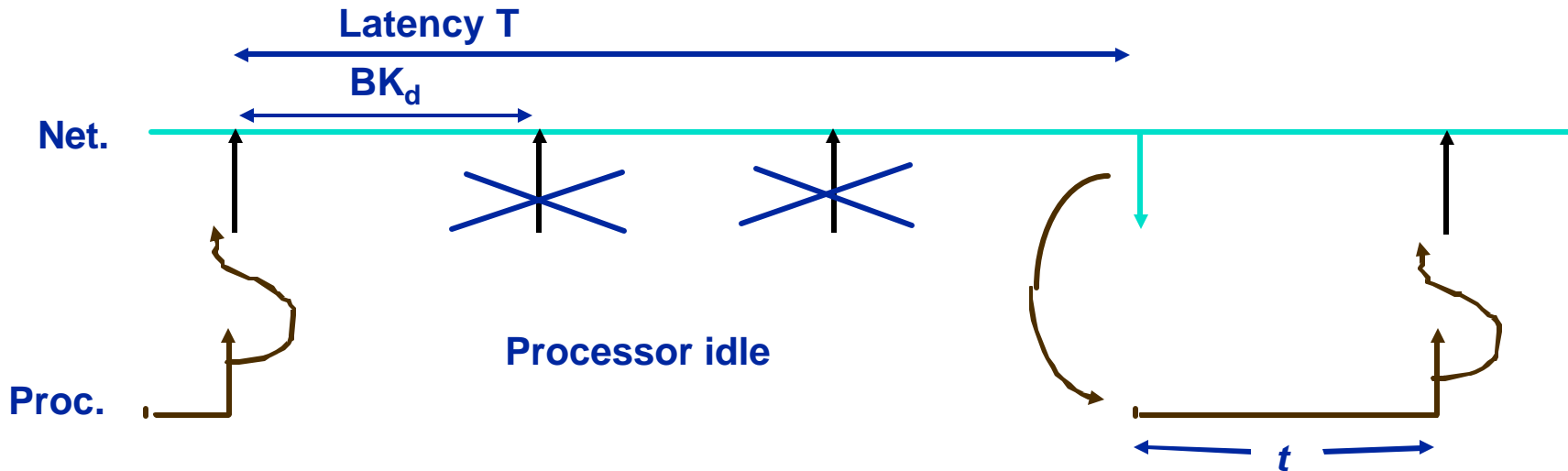
# Zero pointer LimitLESS: All software coherence



remote

local

**Hierarchical - E.g. KSR (actually has rings...)**

# Include network latency

° **Each request suffers T cycles of latency**

**Latency T**

**BK$_d$**

**Net.**

**Processor idle**

**Proc.**

$t$

$$\frac{t}{t + T} = \frac{1}{1 + mT}$$

° **Processor utilization =**

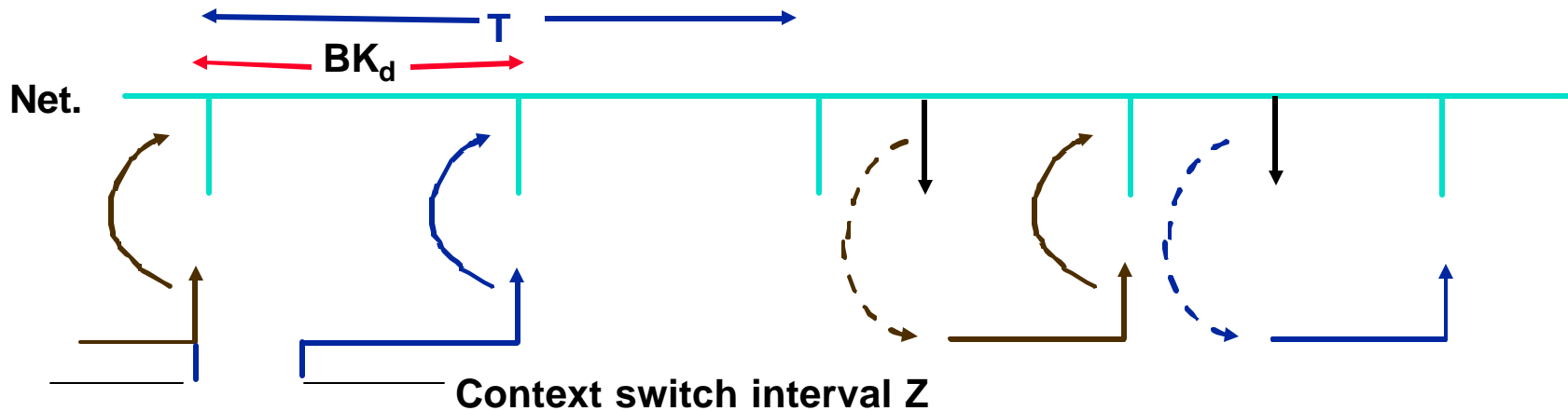$$\left( t = \frac{1}{m} \right)$$

**Processor utilization**

**Network bandwidth also wasted because of lost issue opportunities!**

FIX?

# One solution

## Overlap communication with computation.

$T$

$BK_d$

**Net.**
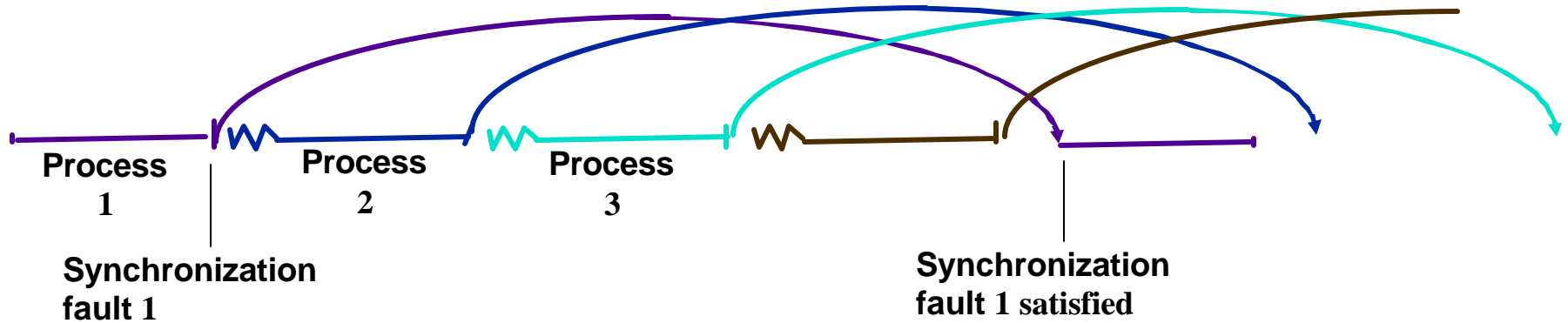
**Context switch interval Z**

- Multithread" the processor

**Processor utilization** $= \dfrac{pt}{t + T}$ if $pt < (t + T)$

**or** $= \dfrac{t}{t + Z}$ otherwise

- And/or allow multiple outstanding requests -- non-blocking memory
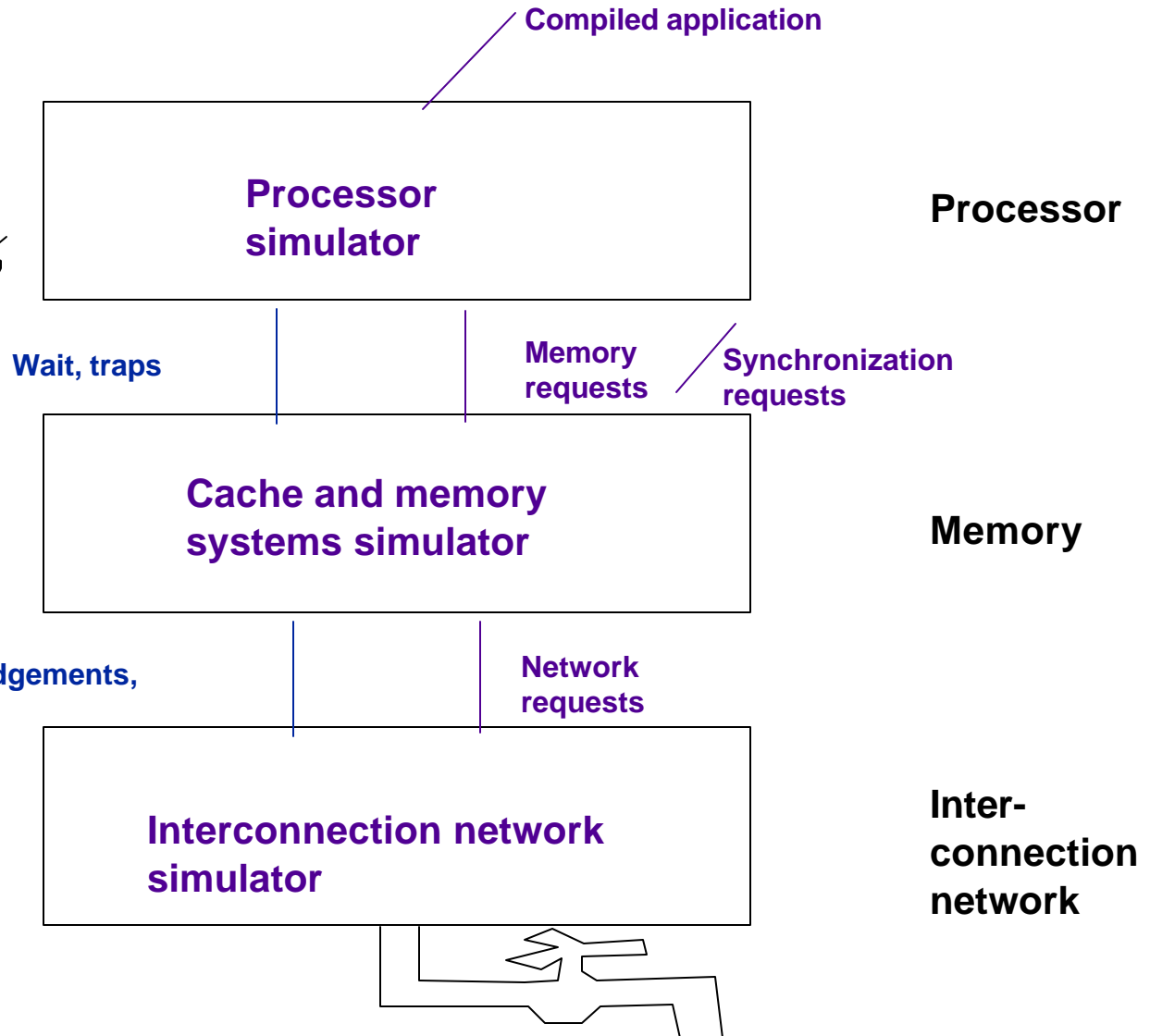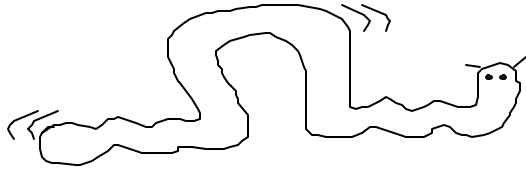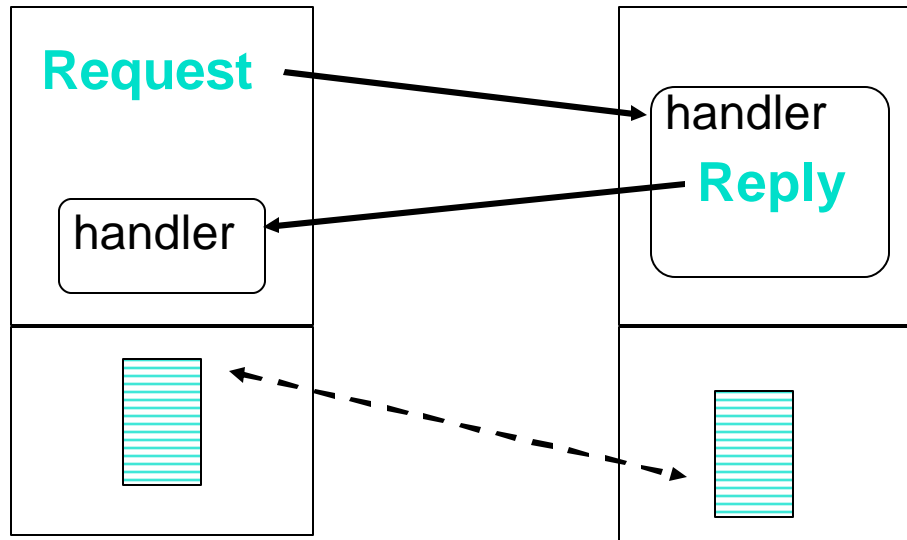
# Synchronization delays



**Process 1**

**Process 2**

**Process 3**

**Synchronization fault 1**

**Synchronization fault 1 satisfied**

° **If no multithreading**



**Wasted processor cycles**

**Fault**

**Satisfied**

# Full system simulation (coupled)

**Compiled application**

**4000x slowdown per node**

**Very accurate**

**Processor simulator**

**Processor**

**Wait, traps**

**Memory requests**

**Synchronization requests**

**Cache and memory systems simulator**

**Memory**

**Acknowledgements, responses**

**Network requests**

**Interconnection network simulator**

**Inter-connection network**

**Measures:**
   **Speedup, runtime**
   **proc. util.**
   **net latency**
   **cache miss rate**
      **.**
      **.**

# Active Messages



- **User-level analog of network transaction**
  - **transfer data packet and invoke handler to extract it from the network and integrate with on-going computation**

- **Request/Reply**

- **Event notification: interrupts, polling, events?**

- **May also perform memory-to-memory transfer**
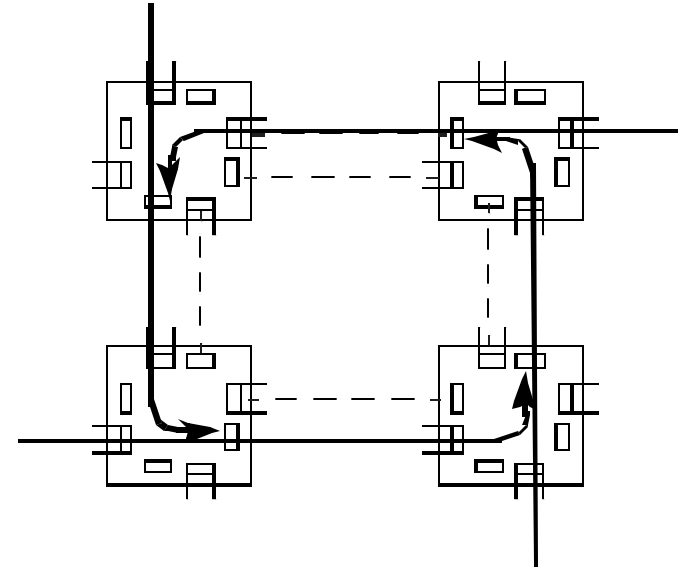
# Deadlock Freedom

° **How can it arise?**

  - **necessary conditions:**
    - **shared resource**
    - **incrementally allocated**
    - **non-preemptible**
  - **think of a channel as a shared resource that is acquired incrementally**
    - **source buffer then dest. buffer**
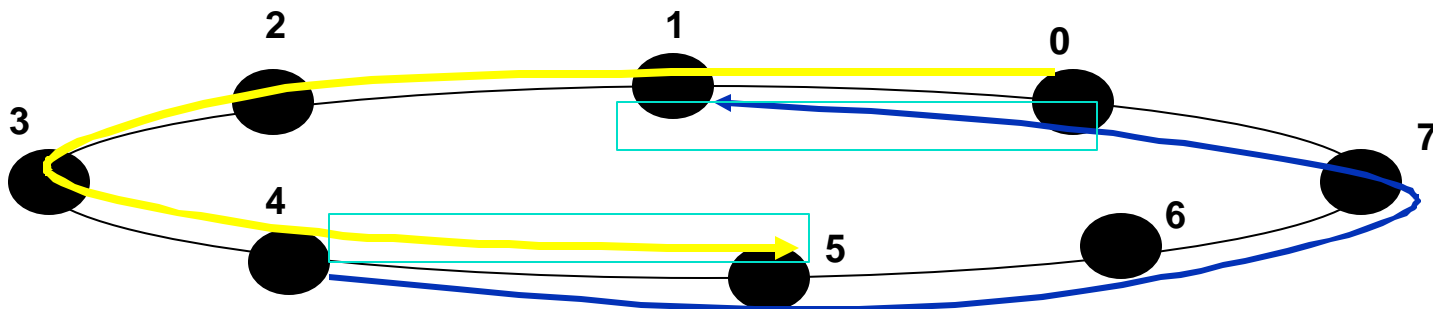    - **channels along a route**

° **How do you avoid it?**

  - **constrain how channel resources are allocated**
  - **ex: dimension order**

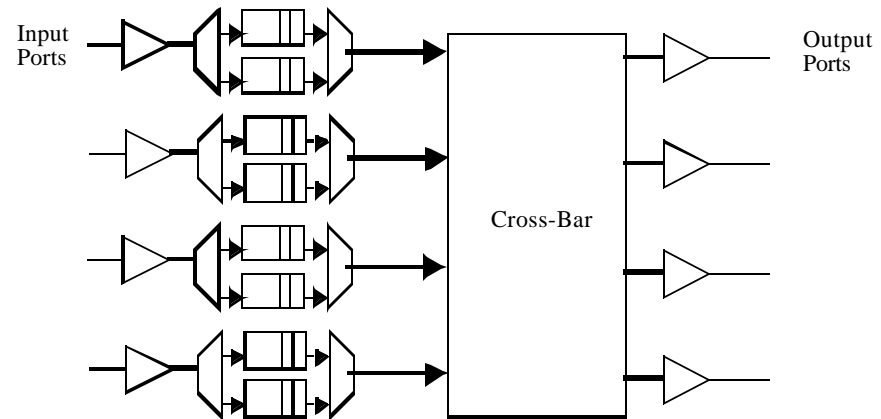° **How do you prove that a routing algorithm is deadlock free**

# More examples:

○ **Consider other topologies**
  - **butterfly?**
  - **tree?**
  - **fat tree?**

○ **Any assumptions about routing mechanism? amount of buffering?**

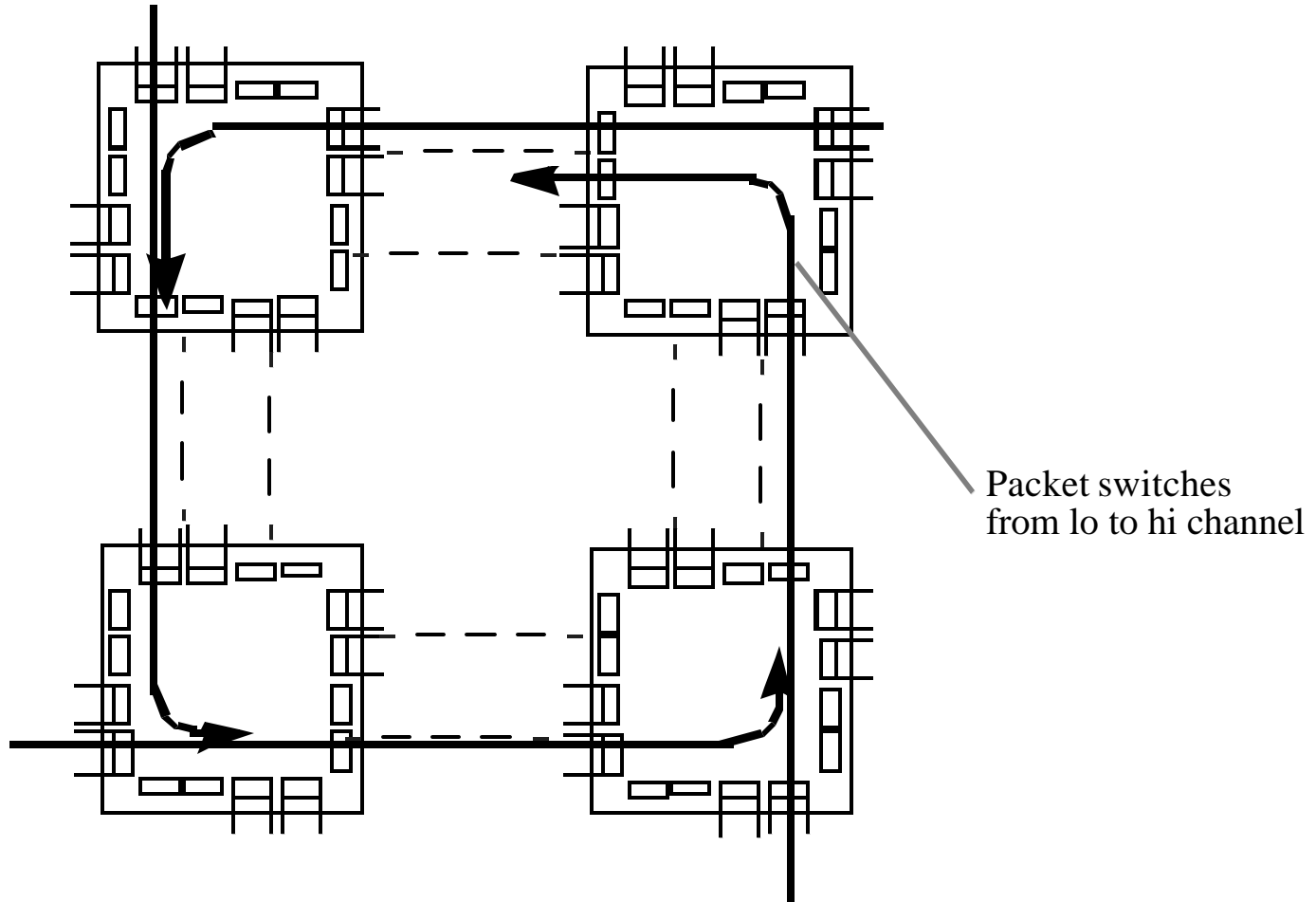○ **What about wormhole routing on a ring?**

# Deadlock free wormhole networks?

° **Basic dimension order routing techniques don't work for k-ary n-cubes**

  • **only for k-ary n-arrays (bi-directional)**

° **Idea: add channels!**

  • **provide multiple "virtual channels" to break the dependence cycle**

  • **good for BW too!**

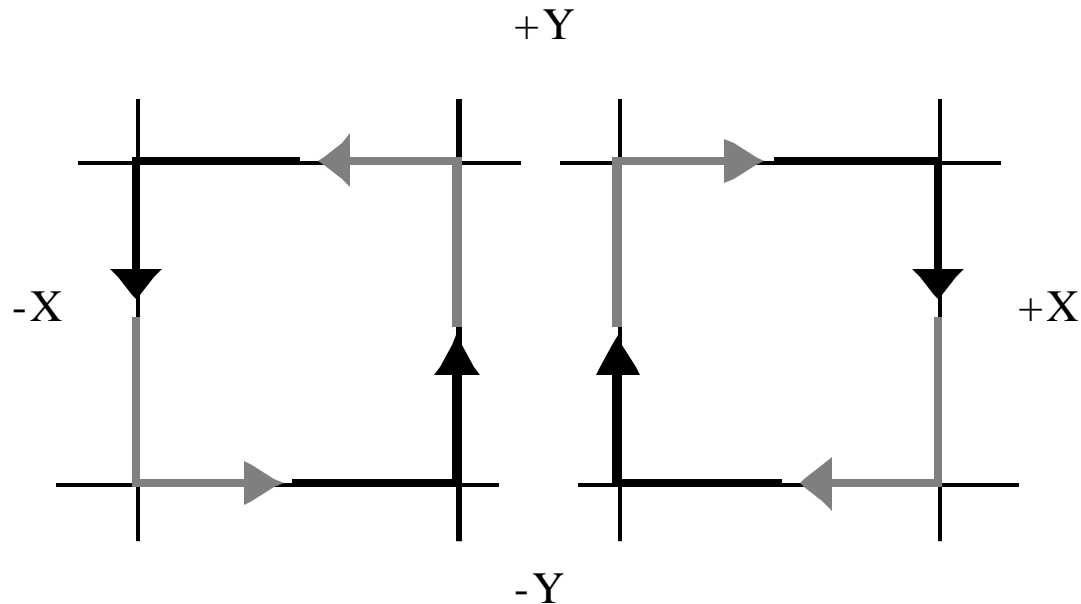Input Ports

Output Ports

Cross-Bar

  • **Do not need to add links, or xbar, only buffer resources**

# Breaking deadlock with virtual channels



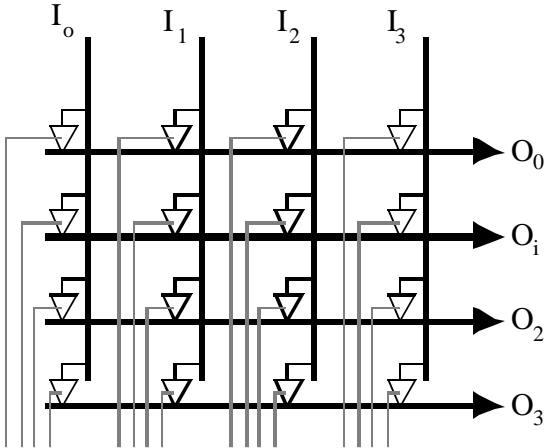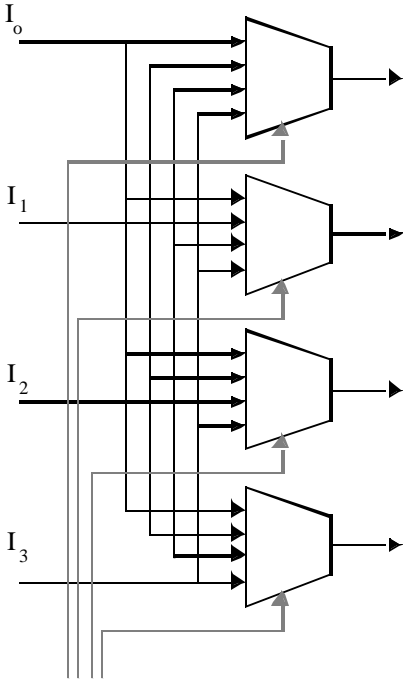Packet switches
from lo to hi channel
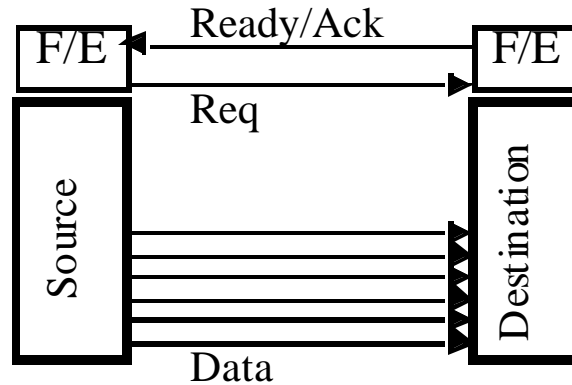
# Turn Restrictions in X,Y

+Y

-X

+X

-Y

- ° **XY routing forbids 4 of 8 turns and leaves no room for adaptive routing**

- ° **Can you allow more turns and still be deadlock free**
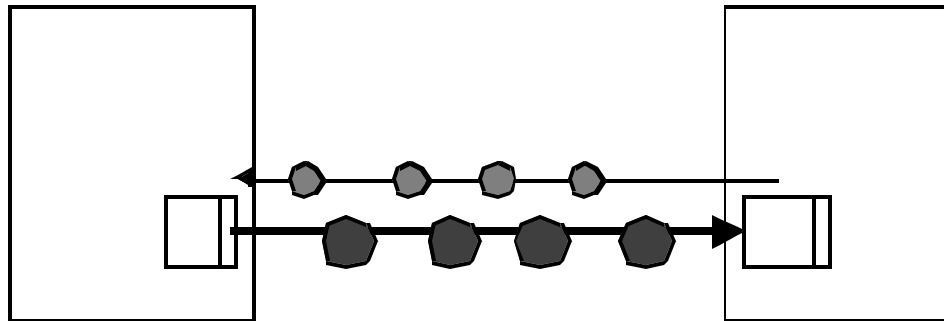
# How do you build a crossbar

# Examples

° **Short Links**



° **long links**

  • **several flits on the wire**

# Modulo Unrolling – Smart Memory

- Loop unrolling relies on dependencies

- Allow maximum parallelism

- Minimize communication



(a) Data
A[]
B[]

Code
...
for(i=0;i<100;i++)
  A[i]=A[i]*B[i+1]
...

(b)
A[]   B[]

for(i=0;i<100;i++)
  A[i]=A[i]*B[i+1]

(c)
$A_0[]$  $B_0[]$   $A_1[]$  $B_1[]$
$A_2[]$  $B_2[]$   $A_3[]$  $B_3[]$

for(i=0;i<100;i+=4) {
  A[i]=A[i]*B[i+1]
  A[i+1]=A[i+1]*B[i+2]
  A[i+2]=A[i+2]*B[i+3]
  A[i+3]=A[i+3]*B[i+4]
}