



中国科学院大学
University of Chinese Academy of Sciences



高级计算机系统结构

沈海华

shenhh@ucas.ac.cn

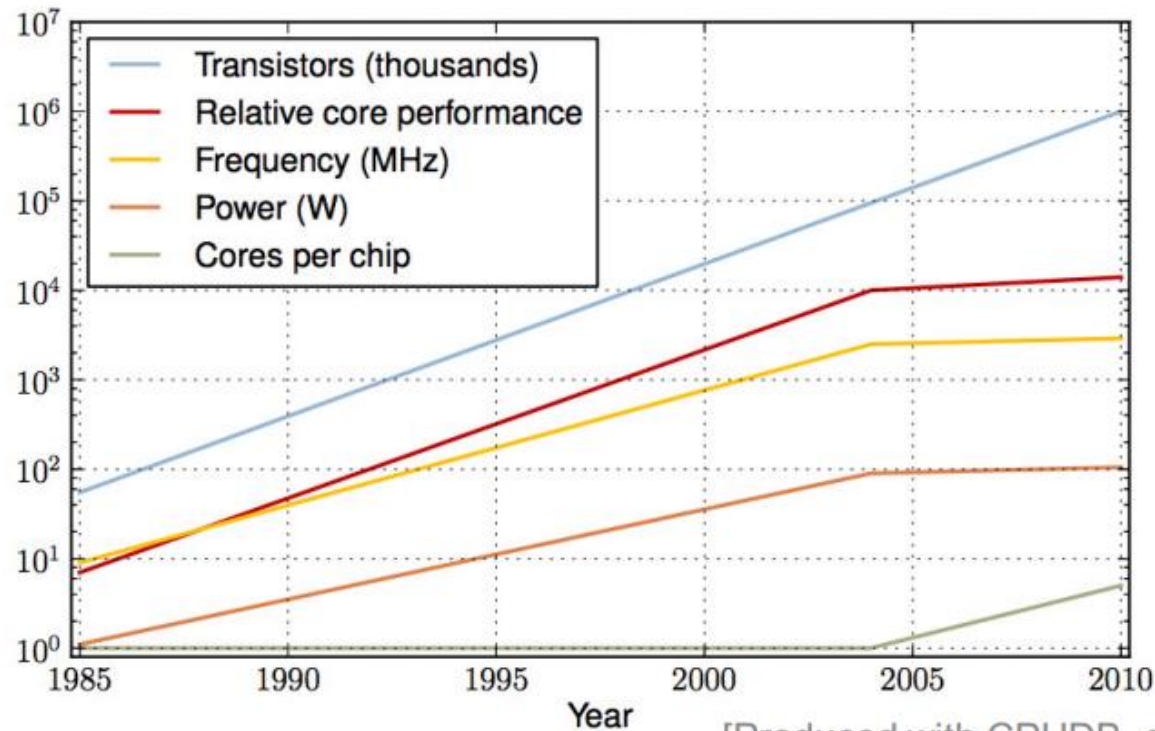
第十一讲 On-Chip Networks I:

Topology/Flow Control

Daniel Sanchez

Computer Science & Artificial Intelligence Lab
M.I.T.

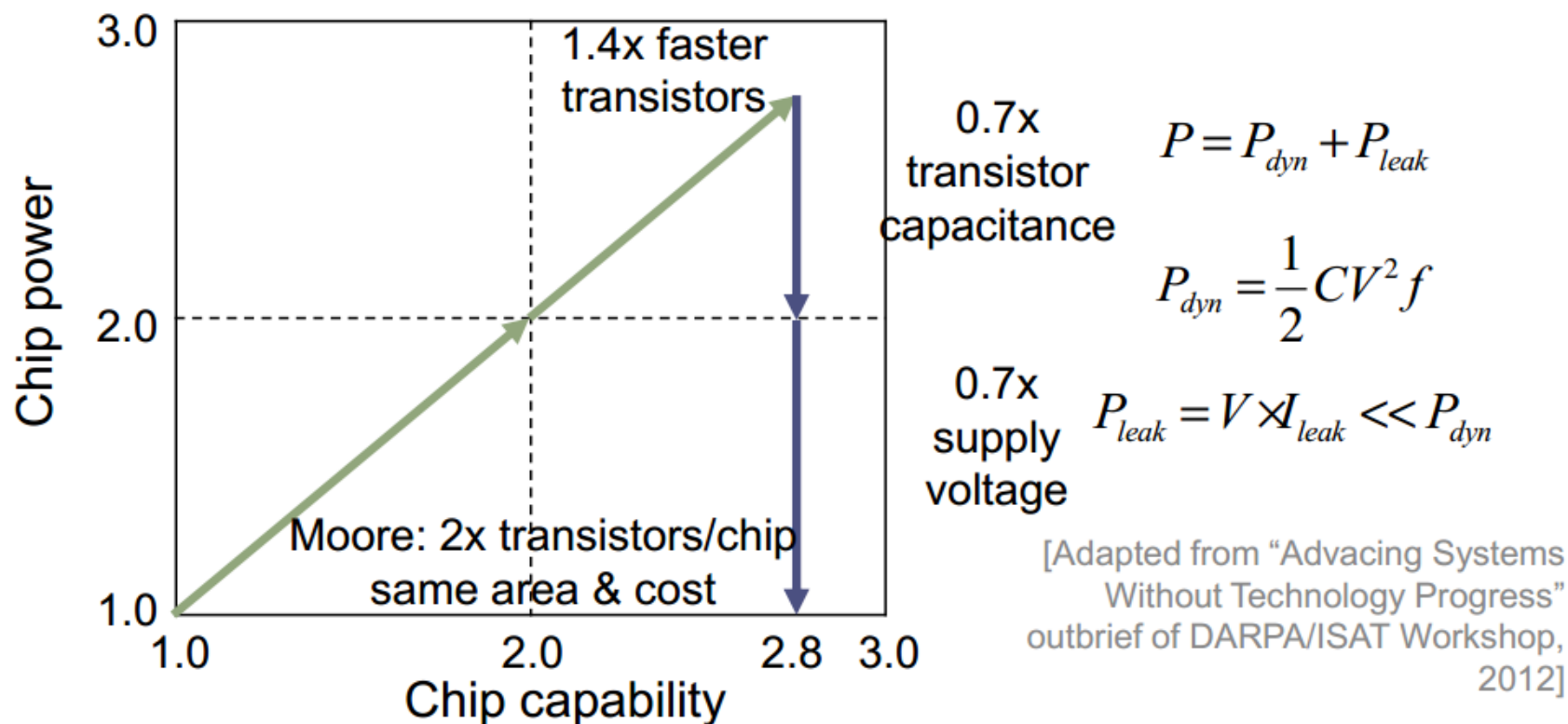
The Shift to Multicore



- Since 2005, improvements in system performance mainly due to increasing cores/chip
- Why? **Limited instruction-level parallelism**
Technology scaling

Classic CMOS Scaling

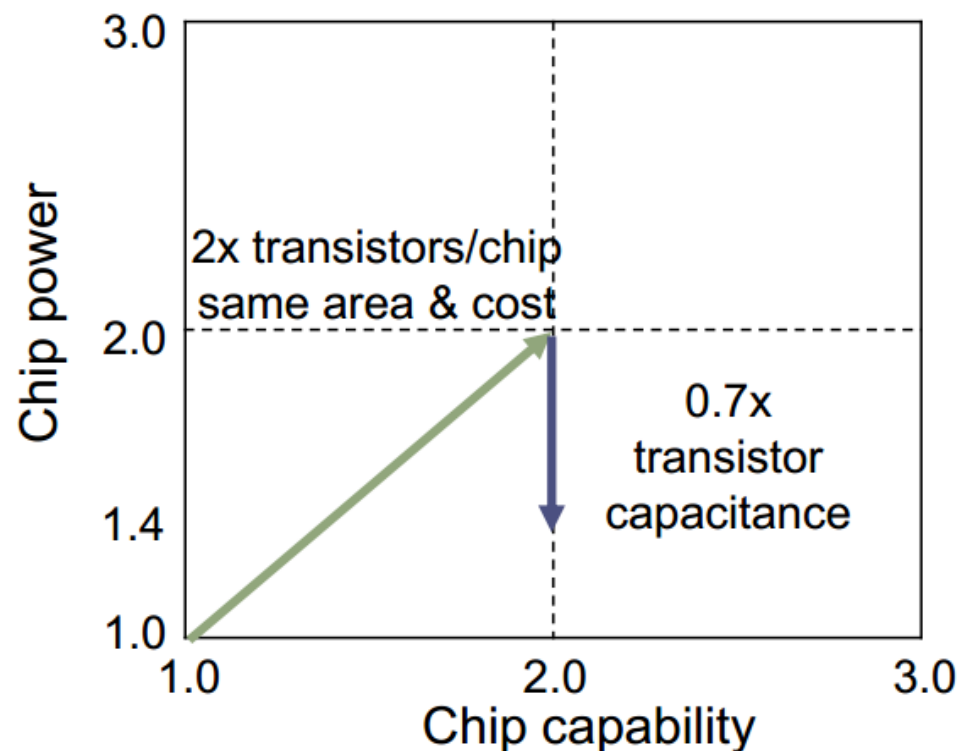
- Moore's law + Dennard scaling: Each generation (e.g., 90→65nm),



2x transistors, 1.4x frequency, same power
→ **area-constrained designs**

Current CMOS Scaling

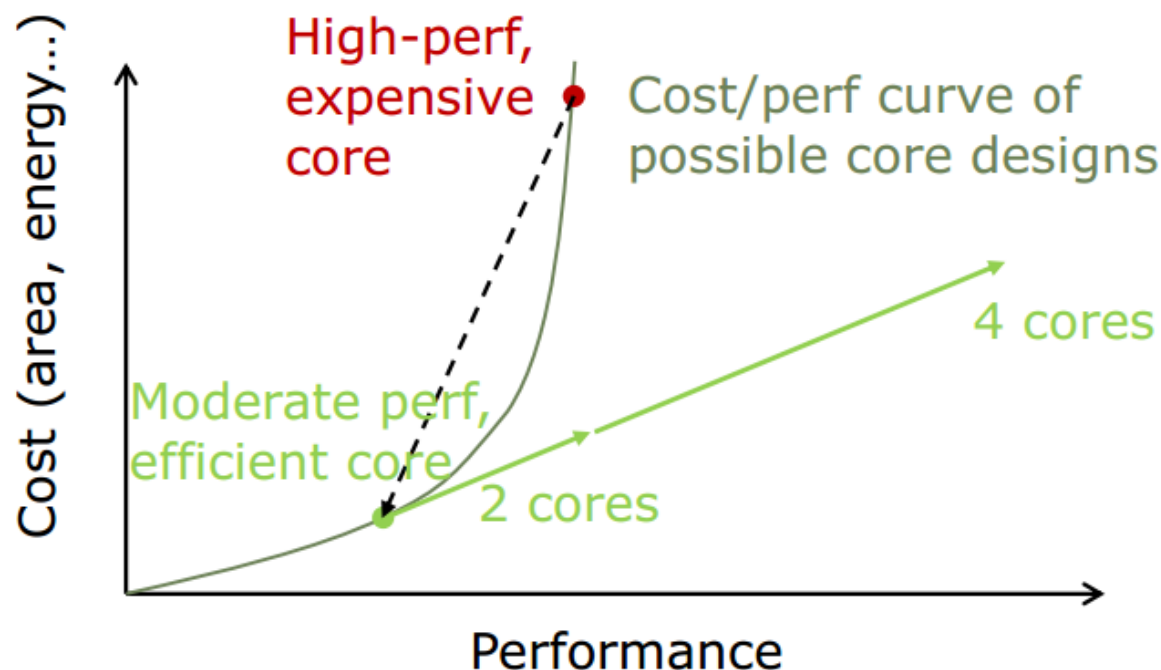
- Frequency and supply voltage scaling are mostly exhausted



[Adapted from "Advancing Systems Without Technology Progress" outbrief of DARPA/ISAT Workshop, 2012]

2x transistors, same frequency, 1.4x power
→ **power-constrained designs**

Multicore Performance



What factors may limit multicore performance?

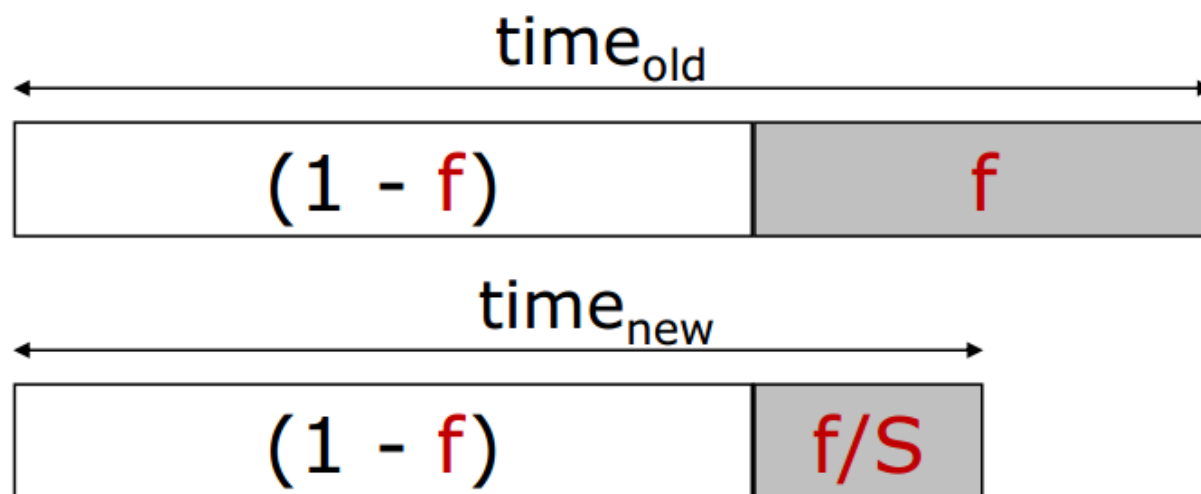
Limited application parallelism
Memory accesses and inter-core communication
Programming complexity

Amdahl's Law

- $\text{Speedup} = \text{time}_{\text{without enhancement}} / \text{time}_{\text{with enhancement}}$
- Suppose an enhancement speeds up a fraction f of a task by a factor of S

$$\text{time}_{\text{new}} = \text{time}_{\text{old}} \cdot ((1-f) + f/S)$$

$$S_{\text{overall}} = 1 / ((1-f) + f/S)$$



Corollary: Make the common case fast

Amdahl's Law and Parallelism

- Say you write a program that can do 90% of the work in parallel, but the other 10% is sequential
- What is the maximum speedup you can get by running on a multicore machine?

$$S_{\text{overall}} = 1 / ((1-f) + f/S)$$

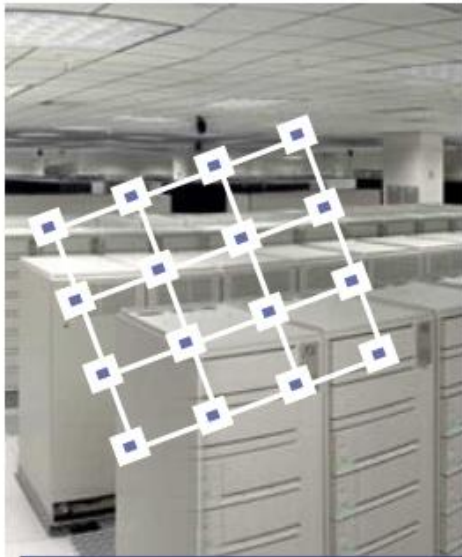
$$f = 0.9, S=\infty \rightarrow S_{\text{overall}} = 10$$

What f do you need to use a 1000-core machine well?

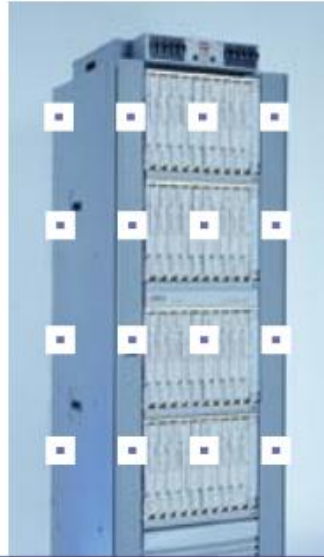
On-Chip Networks

History: From interconnection networks to on-chip networks

Box-to-box networks



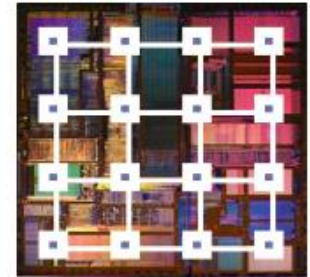
Board-to-board networks



Chip-to-chip networks



On-chip networks



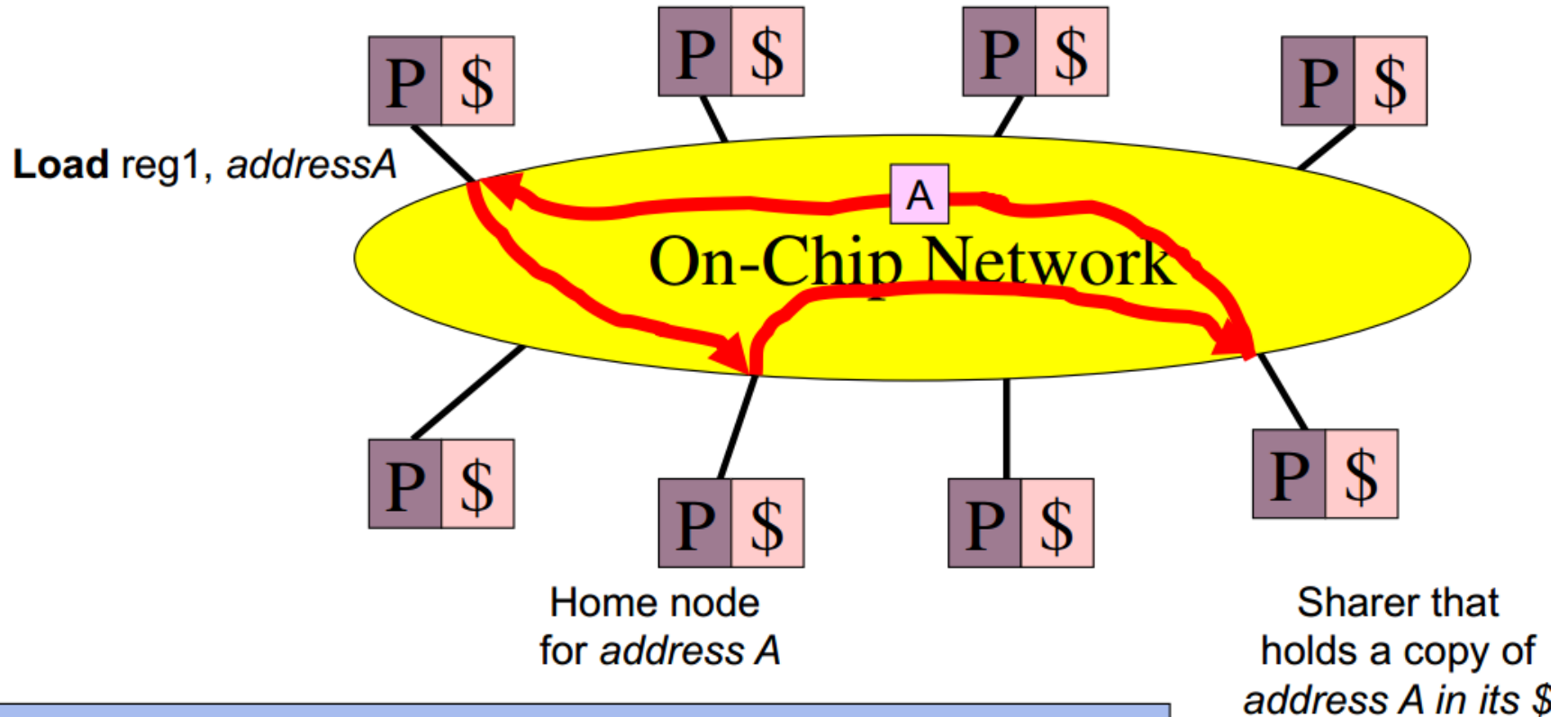
Focus on on-chip networks connecting caches in shared memory processors

Multi-Chip: Supercomputers, Data Centers, Internet Routers, Servers

On-Chip: Servers, Laptops, Phones, HDTVs, Access routers

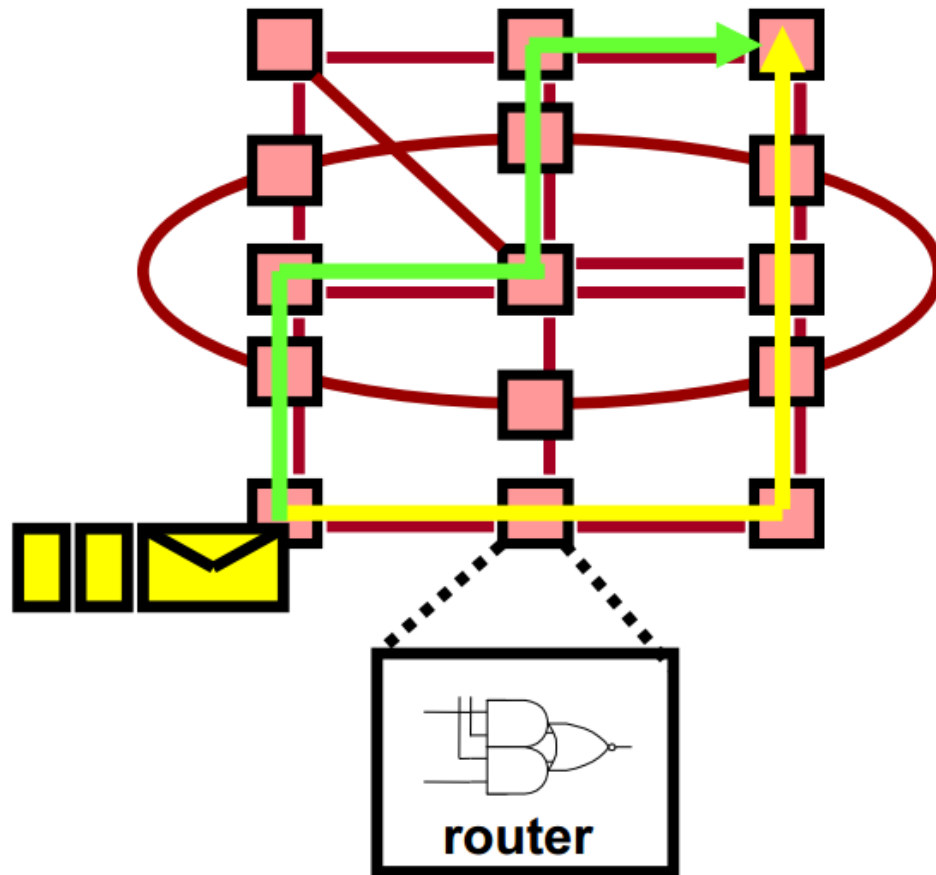
What's an on-chip network?

E.g. Cache-coherent chip multiprocessor



It transports cache coherence messages and cache lines between processor cores

Designing an on-chip network



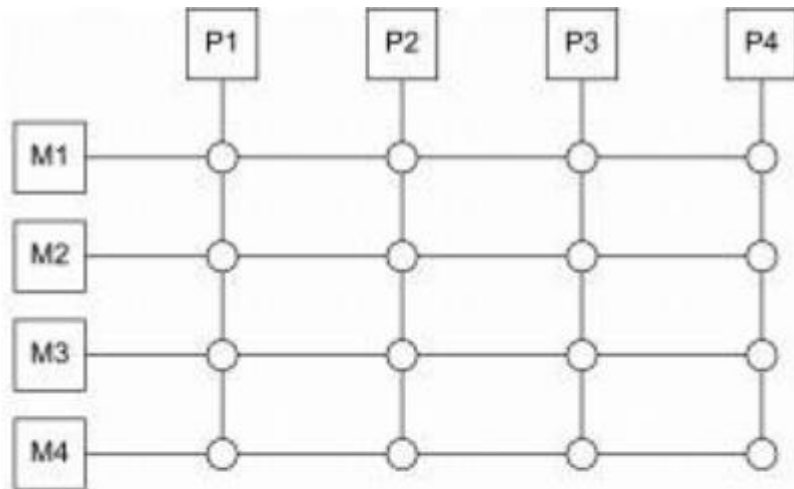
- Topology
- Flow control
- Router microarchitecture
- Routing

Interconnection Network Architecture

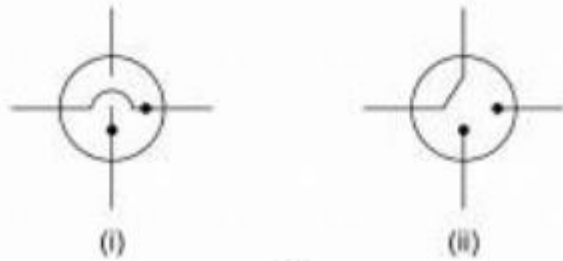
- *Topology*: How to connect the nodes up? (processors, memories, router line cards, ...)
- *Routing*: Which path should a message take?
- *Flow control*: How is the message actually forwarded from source to destination?
- *Router microarchitecture*: How to build the routers?
- *Link microarchitecture*: How to build the links?

Topology

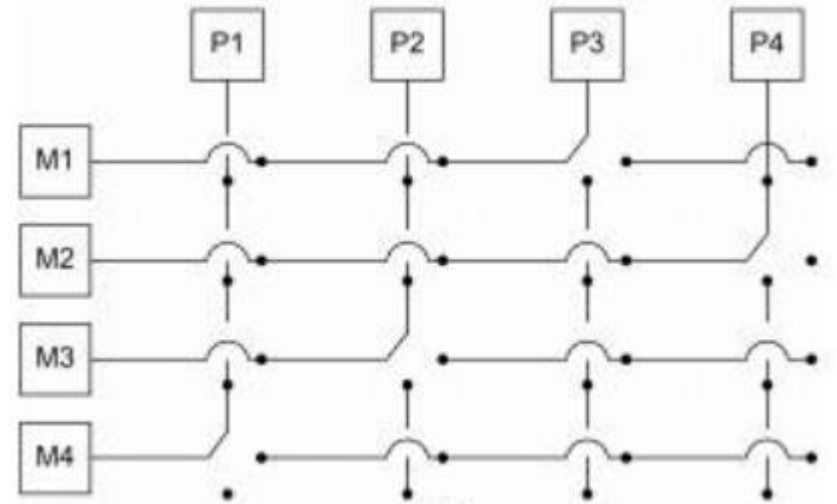
总线（bus）和交叉开关网（crossbar）



a)



b)

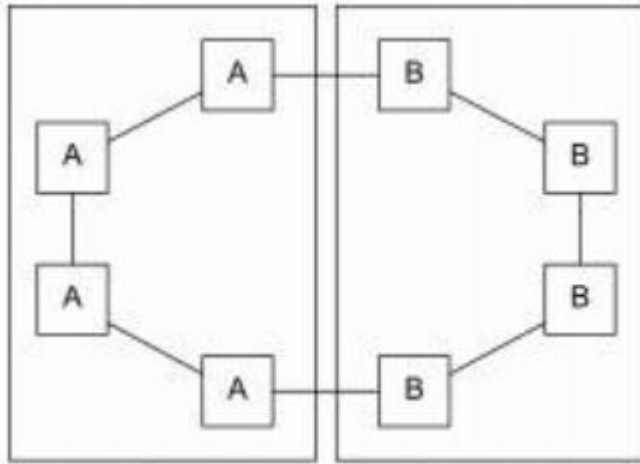


c)

Topological Properties

- *Routing Distance* - number of links on route
- *Diameter* - maximum routing distance
- *Average Distance*
- A network is *partitioned* by a set of links if their removal disconnects the graph
- *Bisection Bandwidth* is the bandwidth crossing a minimal cut that divides the network in half

Bisection width



- 衡量“同时通信的链路数目”或者“连接性”的一个标准是等分宽度 (Bisection width)
- 计算等分宽度的方法是：去除最少的链路数从而将节点分成两等份，去除的链路数就是等分宽度。
- 等分带宽 (bisection bandwidth) 通常用来衡量网络的质量。它与等分宽度类似。等分带宽是计算连接两个等分之间的链路总带宽，所以：等分带宽 = 等分宽度 * 链路的带宽。例如，如果在环中，链路的带宽是10亿bit/秒，那么环的等分带宽就是20亿bit/秒。

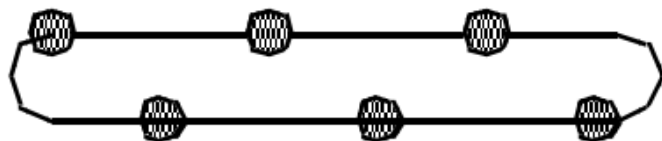
Linear Arrays and Rings



Linear Array



Torus



Torus arranged to use short wires

Route A \rightarrow B given by relative address $R = B - A$

	Linear Array	Ring (1-D Torus)
--	--------------	------------------

Diameter?	$N-1$	$N/2$ (if even N)
-----------	-------	----------------------

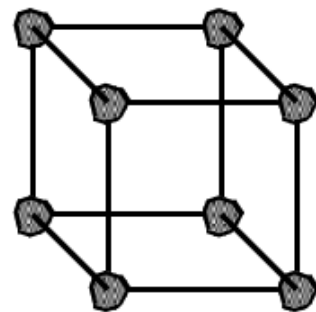
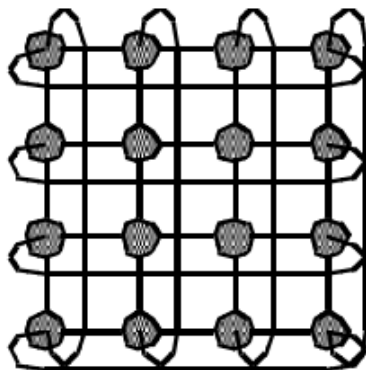
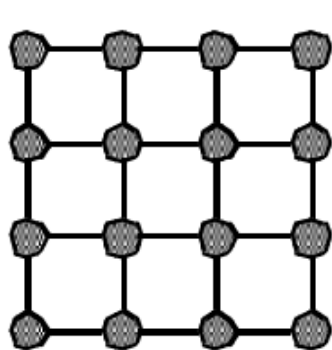
Average distance?	$N/3 - 1/(3N)$	$N/4$ (if even N)
-------------------	----------------	----------------------

Bisection bandwidth?	1	2
----------------------	---	---

- Torus Examples:

- FDDI, SCI, FiberChannel Arbitrated Loop, Intel Xeon

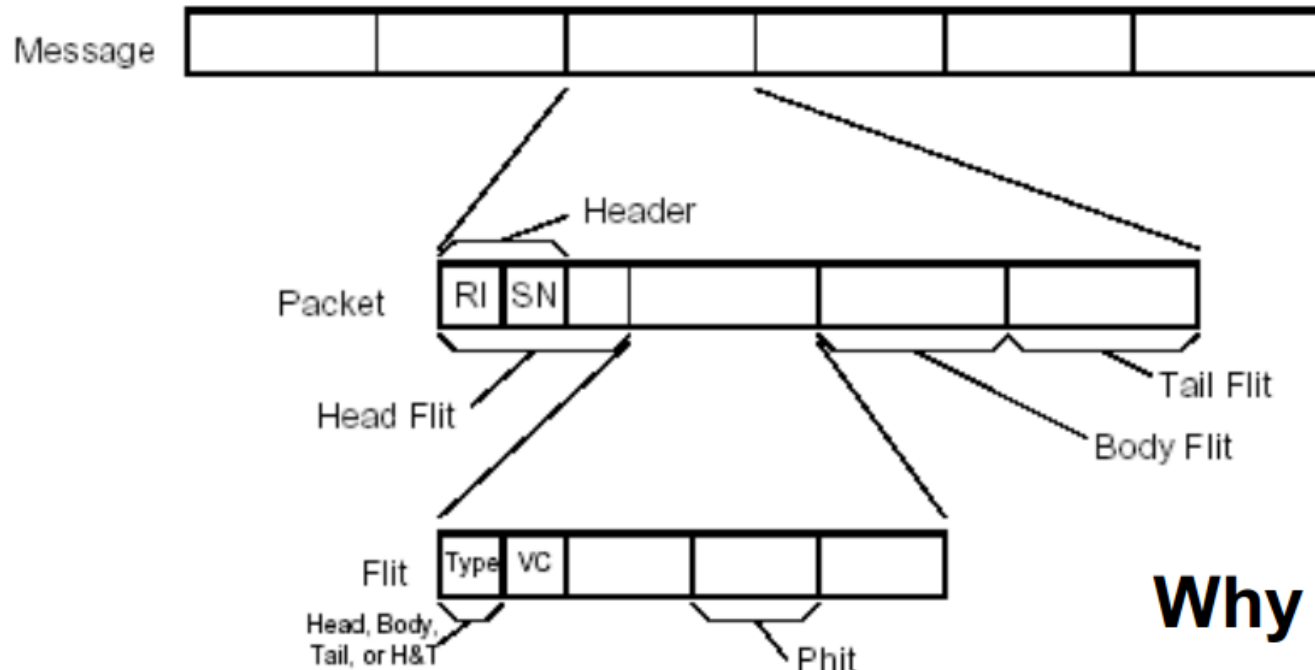
Multidimensional Meshes and Tori



- d -dimensional array
 - $n = k_{d-1} \times \dots \times k_0$ nodes
 - described by d -vector of coordinates (i_{d-1}, \dots, i_0)
- d -dimensional k -ary mesh: $N = k^d$
 - $k = \sqrt[d]{N}$
 - described by d -vector of radix k coordinate
- d -dimensional k -ary torus (or k -ary d -cube)

Routing & Flow Control Overview

Messages, Packets, Flits, Phits



Why flits?

For variable
packet sizes

Packet: Basic unit of routing and sequencing

- Limited size (e.g. 64 bits – 64 KB)

Flit (flow control digit): Basic unit of bandwidth/storage allocation

- All flits in packet follow the same path

Phit (physical transfer digit): data transferred in single clock

Routing vs Flow Control

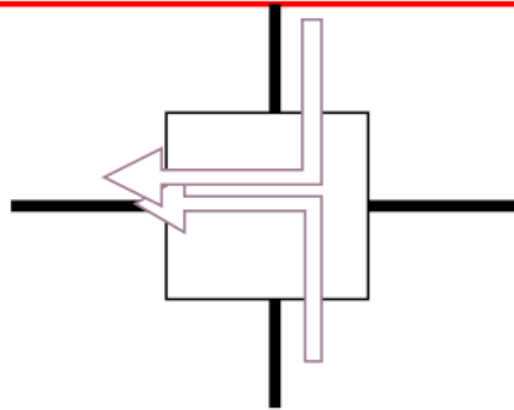
- Routing algorithm chooses path that packets should follow to get from source to destination
- Flow control schemes allocate resources (buffers, links, control state) to packets traversing the network
- Our approach: Bottom-up
 - Today: Flow control, assuming routes are set
 - Next lecture: Routing algorithms

Properties of Routing Algorithms

- **Deterministic/Oblivious**
 - Route determined by (source, dest), not intermediate state (i.e. traffic)
- **Adaptive**
 - Route influenced by traffic along the way
- **Minimal**
 - Only selects shortest paths
- **Deadlock-free**
 - No traffic pattern can lead to a situation where no packets move forward

Flow Control

Contention



- Two packets trying to use the same link at the same time
 - Limited or no buffering
- Problem arises because we are sharing resources
 - Sharing bandwidth and buffering

Flow Control Protocols

- **Bufferless**
 - Circuit switching
 - Dropping
 - Misrouting
- **Buffered**
 - Store-and-forward
 - Virtual cut-through
 - Wormhole
 - Virtual-channel

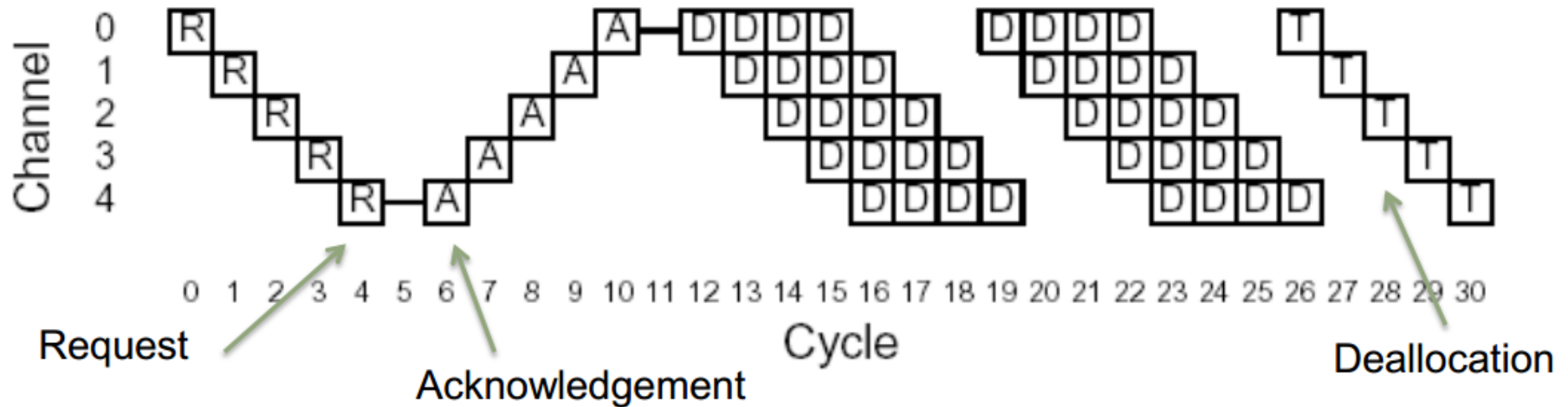


**Complexity
&
Efficiency**

Circuit Switching

- Form a circuit from source to dest
- Probe to set up path through network
- Reserve all links
- Data sent through links
- Bufferless

Time-space View: Circuit Switching



- Why is this good?

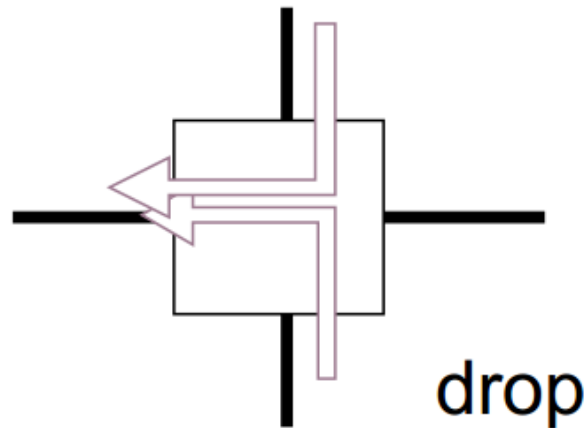
Simple to implement

- Why is it not?

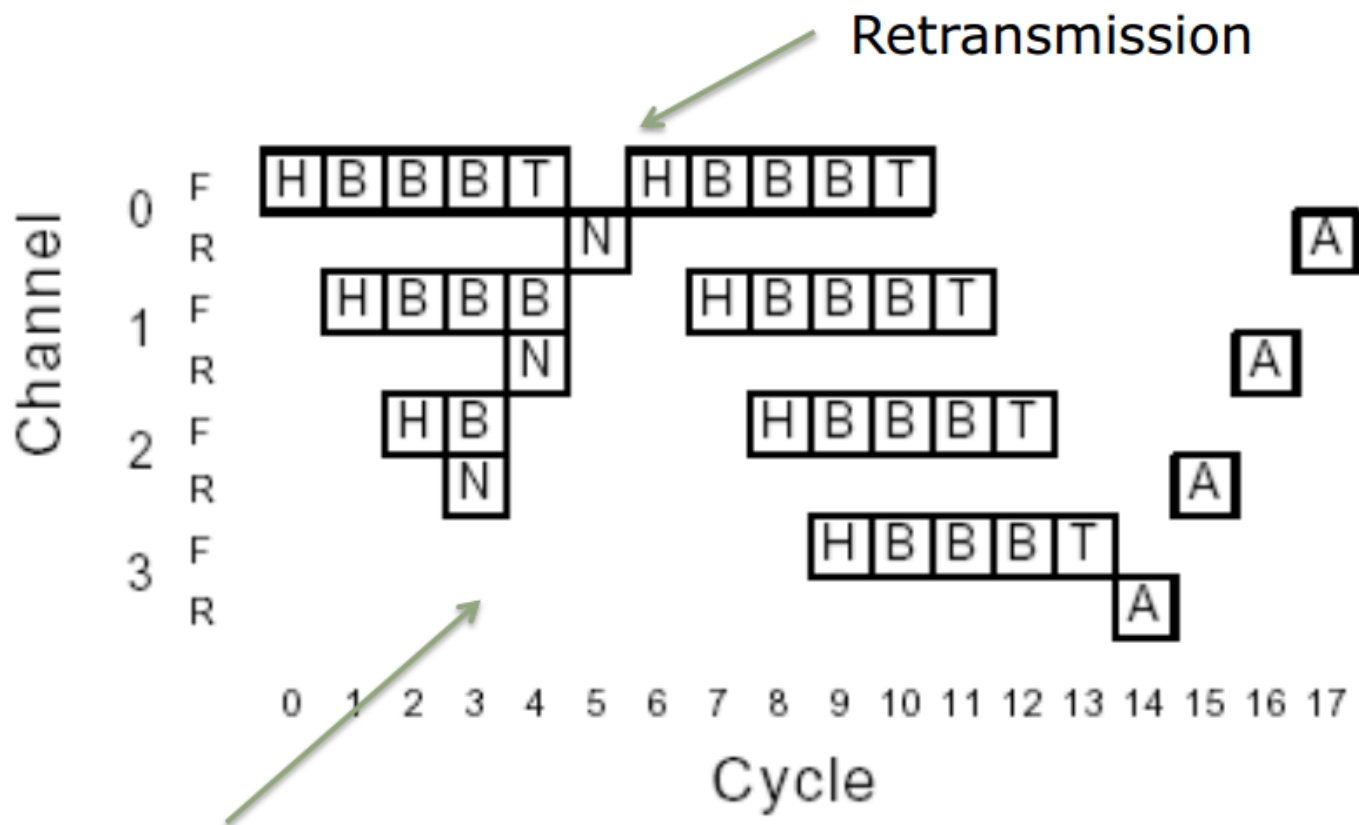
Wasteful, increased 3X latency for short packets

Speculative Flow Control: Dropping

- If two things arrive and I don't have resources, drop one of them
- Flow control protocol on the Internet



Time-space Diagram: Dropping



Unable to allocate channel 3

Disadvantages:

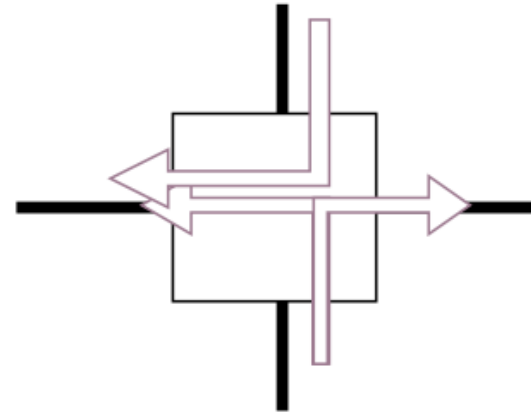
Poor tradeoff of traffic
and buffering

Less Simple Flow Control: Misrouting

- If only one message can enter the network at each node, and one message can exit the network at each node, the network can never be congested. Right?

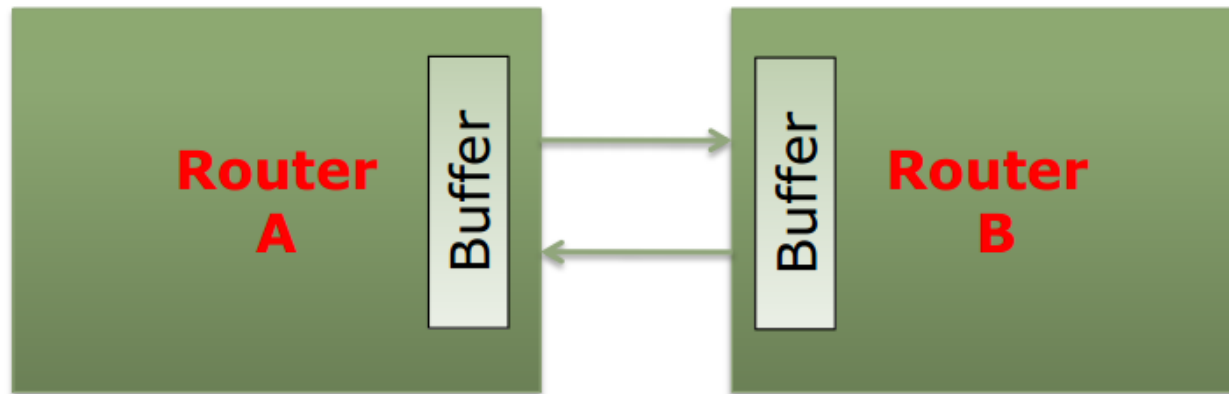
Wrong! Multiple hops cause congestion

- Philosophy behind misrouting: intentionally route away from congestion
- No need for buffering
- Problems?



Livelock: need to guarantee that progress is made

Buffered Routing



- Link-level flow control:
 - Given that you can't drop packets, how to manage the buffers?
When can you send stuff forward, when not?
- Metrics of interest:
 - Throughput/Latency
 - Buffer utilization (turnaround time)

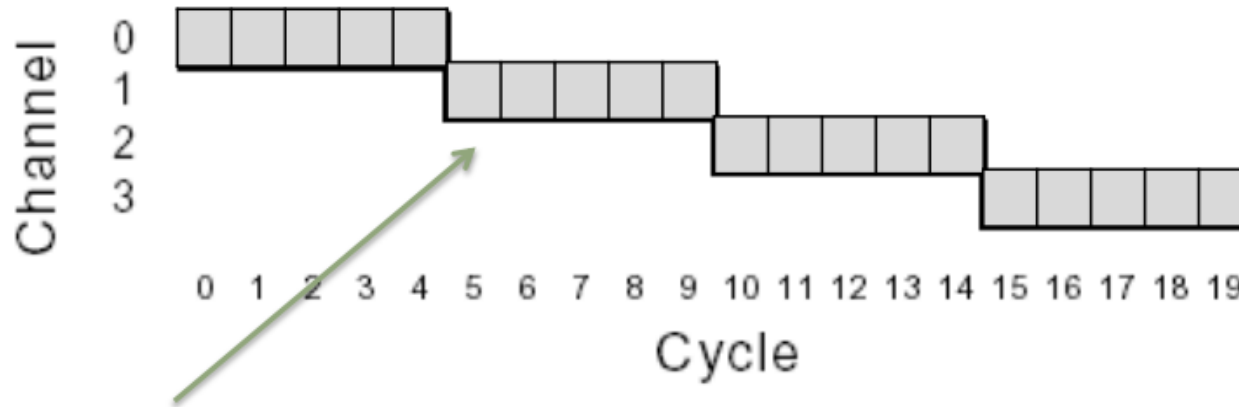
Techniques for link backpressure

- Naïve stall-based (on/off):
 - Can source send or not?
- Sophisticated stall-based (credit-based):
 - How many flits can be sent to the next node?
- Speculative (ack/nack):
 - Guess can always send, but keep copy
 - Resolve if send was successful (ack/nack)
 - On ack – drop copy
 - On nack - resend

Store-and-Forward (packet-based, no flits)

- **Strategy:**
 - Make intermediate stops and wait until the entire packet has arrived before you move on
- **Advantage:**
 - Other packets can use intermediate links

Time-space View: Store-and-Forward



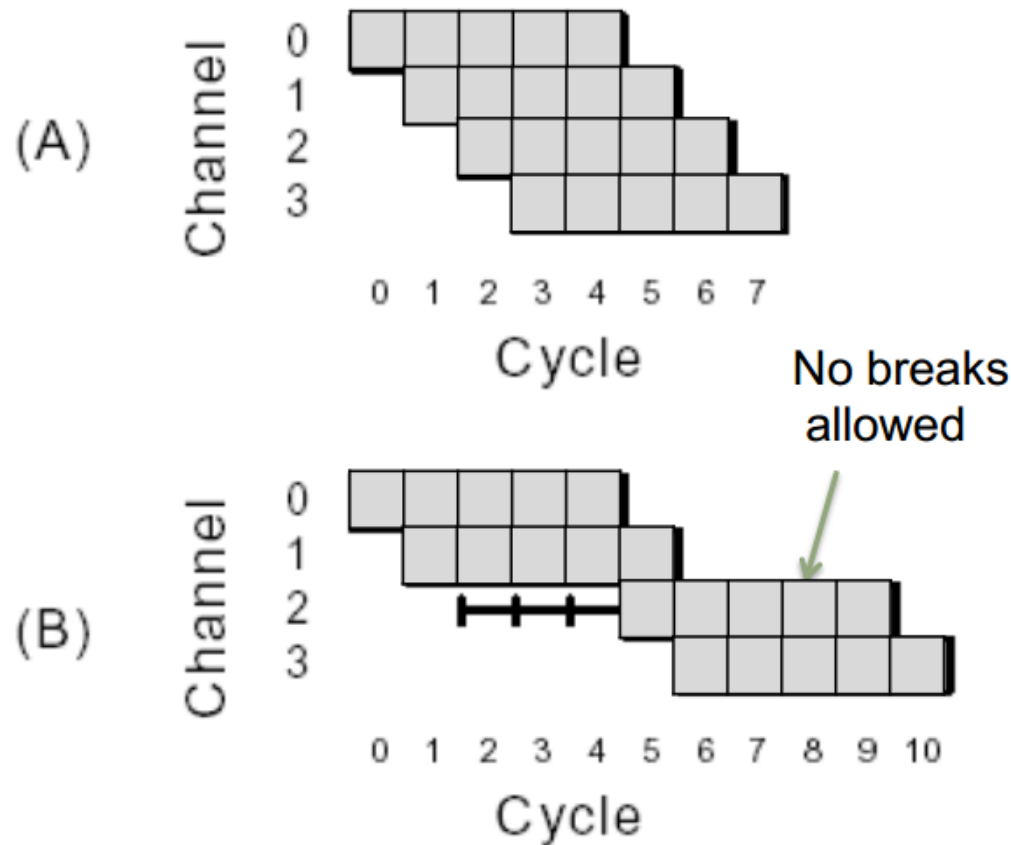
Could be allocated at a much later time without packet dropping

- Buffering allows packet to wait for channel
- Drawback? Serialization latency experienced at each hop/channel

Virtual Cut-through (packet-based)

- Why wait till entire message has arrived at each intermediate stop?
- The head flit of the packet can dash off first
- When the head gets blocked, whole packet gets blocked at one intermediate node
- Used in Alpha 21364

Time-space View: Virtual Cut-through



- Advantages?

Lower latency

- Disadvantages?

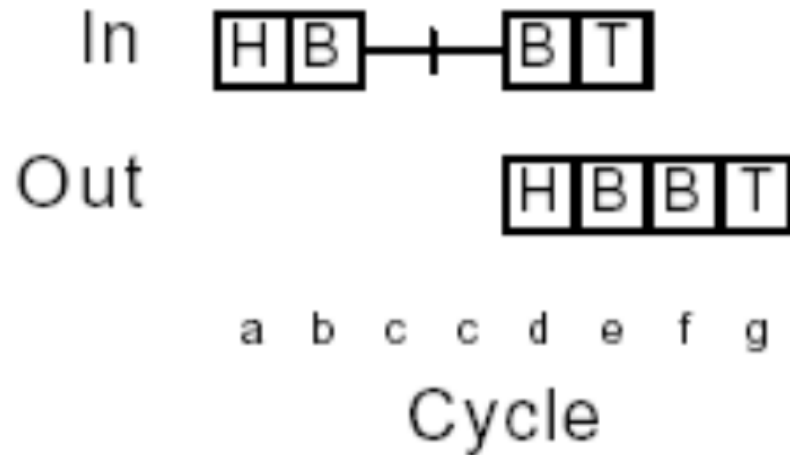
Buffers allocated in packets
→ large buffers & low utilization

Channels allocated in packets
→ unfairness & low utilization

Flit-Buffer Flow Control: Wormhole

- When a packet blocks, just block wherever the pieces (flits) of the message are at that time.
- Operates like cut-through but with channel and buffers allocated to flits rather than packets
 - Channel state (virtual channel) allocated to packet so body flits can follow head flit

Time-space View: Wormhole



- Advantages?
- Disadvantages?

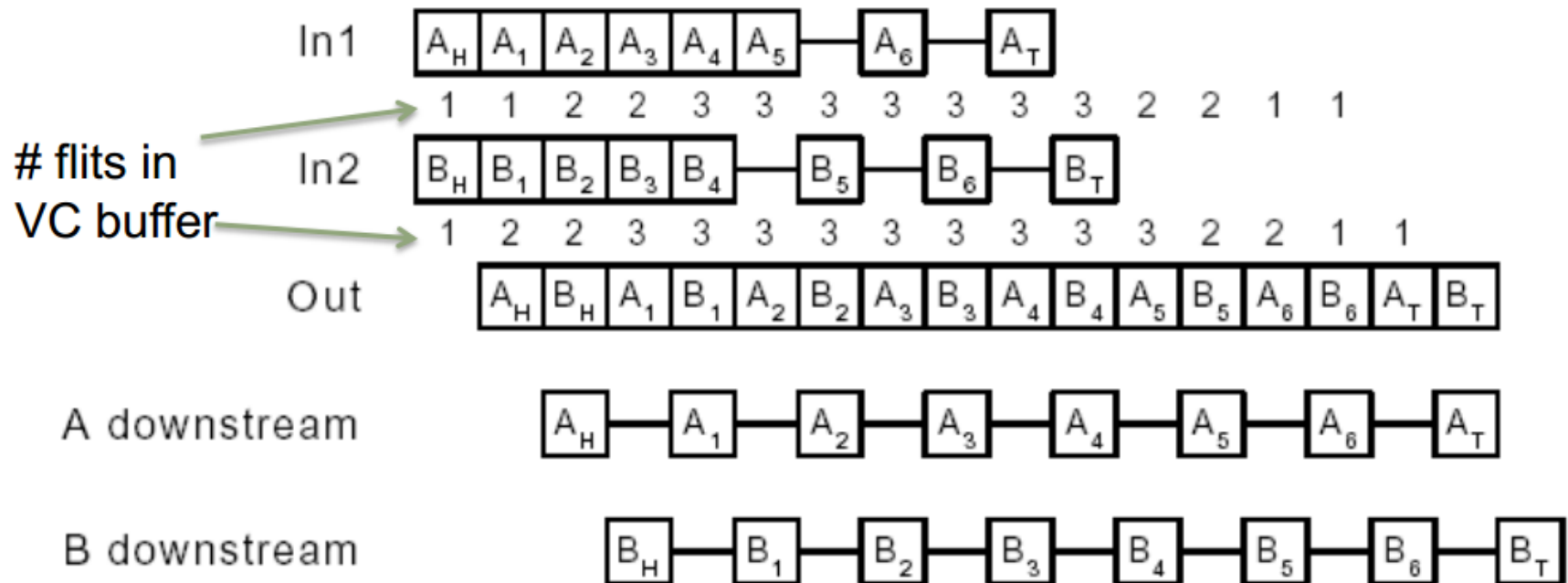
Smaller amount of buffer space required

May block a channel mid-packet, another packet cannot use bandwidth

Virtual-Channel (VC) Flow Control

- When a message blocks, instead of holding on to links so others can't use them, hold on to **virtual** links
- Multiple queues in buffer storage
 - Lanes on the highway
- Virtual channel can be thought of as channel state and flit buffers

Time-space View: Virtual-Channel



- Advantages?
- Disadvantages?

Significantly reduces blocking

More complex router,
fair VC allocation required

Next Time:

Router (Switch) Microarchitecture
Routing Algorithms