

# **CSCI-1680**

## **Link Layer Wrap-Up**

Rodrigo Fonseca



# Today: Link Layer (cont.)

- Framing
- Reliability
  - Sliding window
  - Error correction
- **Medium Access Control**
  - (Short)Case study:  
Ethernet
- **Link Layer Switching**



# Medium Access Control

- **Control access to shared physical medium**
  - E.g., who can talk when?
  - If everyone talks at once, no one hears anything
  - Job of the Link Layer
- **Two conflicting goals**
  - Maximize utilization when one node sending
  - Approach  $1/N$  allocation when  $N$  nodes sending



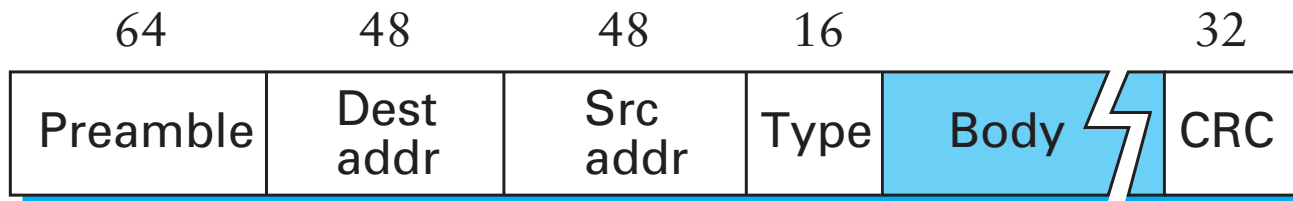
# Different Approaches

- **Partitioned Access**
  - Time Division Multiple Access (TDMA)
  - Frequency Division Multiple Access (FDMA)
  - Code Division Multiple Access (CDMA)
- **Random Access**
  - ALOHA/ Slotted ALOHA
  - Carrier Sense Multiple Access / Collision Detection (CSMA/CD)
  - Carrier Sense Multiple Access / Collision Avoidance (CSMA/CA)
  - RTS/CTS (Request to Send/Clear to Send)
  - Token-based



# Case Study: Ethernet (802.3)

- **Dominant wired LAN technology**
  - 10BASE2, 10BASE5 (Vampire Taps)
  - 10BASET, 100BASE-TX, 1000BASE-T, 10GBASE-T,...
- **Both Physical and Link Layer specification**
- **CSMA/CD**
  - Carrier Sense / Multiple Access / Collision Detection
- **Frame Format (Manchester Encoding):**



# Ethernet Addressing

- **Globally unique, 48-bit unicast address per adapter**
  - Example: **00:1c:43**:00:3d:09 (**Samsung** adapter)
  - 24 msb: organization
  - <http://standards.ieee.org/develop/regauth/oui/oui.txt>
- **Broadcast address: all 1s**
- **Multicast address: first bit 1**
- **Adapter can work in *promiscuous* mode**



# Ethernet MAC: CSMA/CD

- **Problem: shared medium**
  - 10Mbps: up to 2500m, with 4 repeaters at 500m
- **Transmit algorithm**
  - If line is idle, transmit immediately
  - Upper bound message size of 1500 bytes
  - If line is busy: wait until idle and transmit immediately
- **We will just assume you can detect collisions for now**



# When to transmit again?

- Delay and try again: exponential backoff
- $n$ th time:  $k \times 51.2\mu\text{s}$ , for  $k = U\{0..(2^{\min(n,10)}-1)\}$ 
  - 1<sup>st</sup> time: 0 or  $51.2\mu\text{s}$
  - 2<sup>nd</sup> time: 0,  $51.2$ ,  $102.4$ , or  $153.6\mu\text{s}$
- Give up after several times (usually 16)
- Exponential backoff is a useful, general technique





# Capture Effect

- Exponential backoff leads to self-adaptive use of channel
- A and B are trying to transmit, and collide
- Both will back off either 0 or  $51.2\mu\text{s}$
- Say A wins.
- Next time, collide again.
  - A will wait between 0 or 1 slots
  - B will wait between 0, 1, 2, or 3 slots
- ...



# Ethernet Recap

- **Service provided: send frames among stations with specific addresses**
  - Receiver looks at the dest address and decides to receive
- **Addresses are just names, no topology information**
  - Special broadcast and multicast addresses
- **All nodes in the same “broadcast domain”**
- **Earlier versions:**
  - same “collision domain”
  - CSMA/CD with exponential backoff
- **More recent versions:**
  - Switched Ethernet: all nodes have direct connection to switches

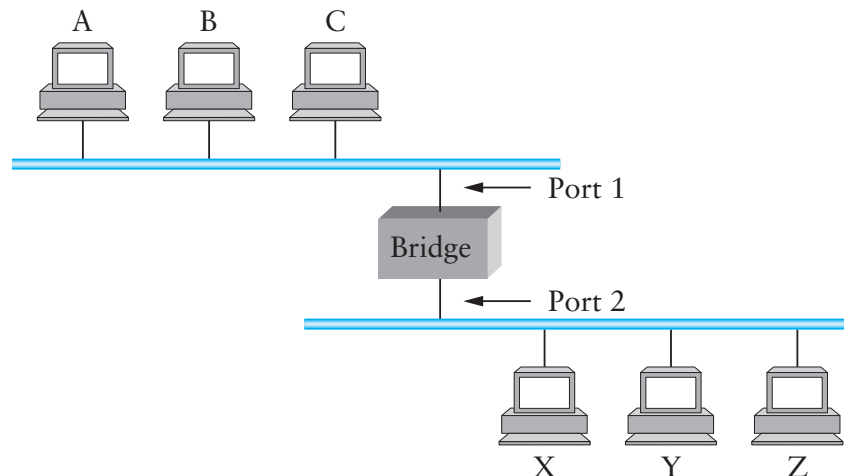


# Bridging



# Bridges and Extended LANs

- **LANs have limitations**
  - E.g. Ethernet < 1024 hosts, < 2500m
- **Connect two or more LANs with a *bridge***
  - Operates on Ethernet addresses
  - Forwards packets from one LAN to the other(s)
- **Ethernet *switch* is just a multi-way bridge**

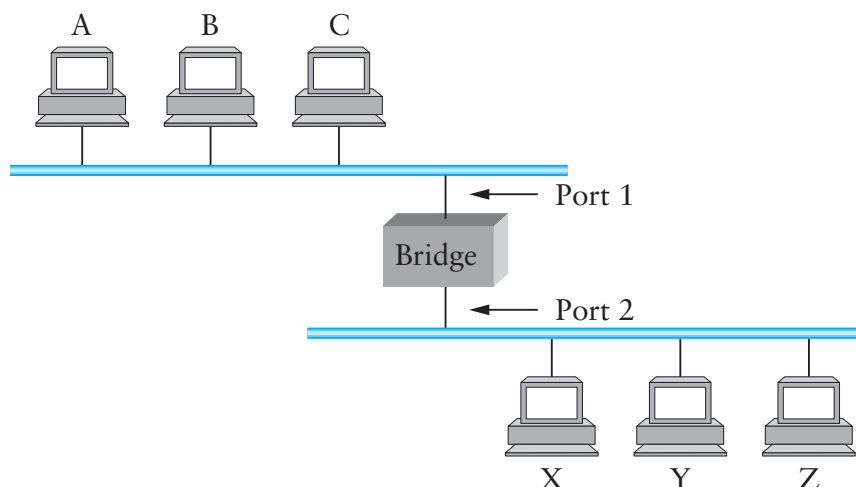


# Bridges

- **Unicast:** forward with filtering
- **Broadcast:** always forward
- **Multicast:** always forward or learn groups
- **Difference between bridges and repeaters?**
  - Bridges: same broadcast domain; copy *frames*
  - Repeaters: same broadcast and *collision domain*; copy *signals*



# Learning Bridges



- **Idea: don't forward a packet where it isn't needed**
  - If you know recipient is not on that port
- **Learn hosts' locations based on source addresses**
  - Build a table as you receive packets
  - Table is a *cache*: if full, evict old entries. Why is this fine?
- **Table says when *not* to forward a packet**
  - Doesn't need to be complete for *correctness*



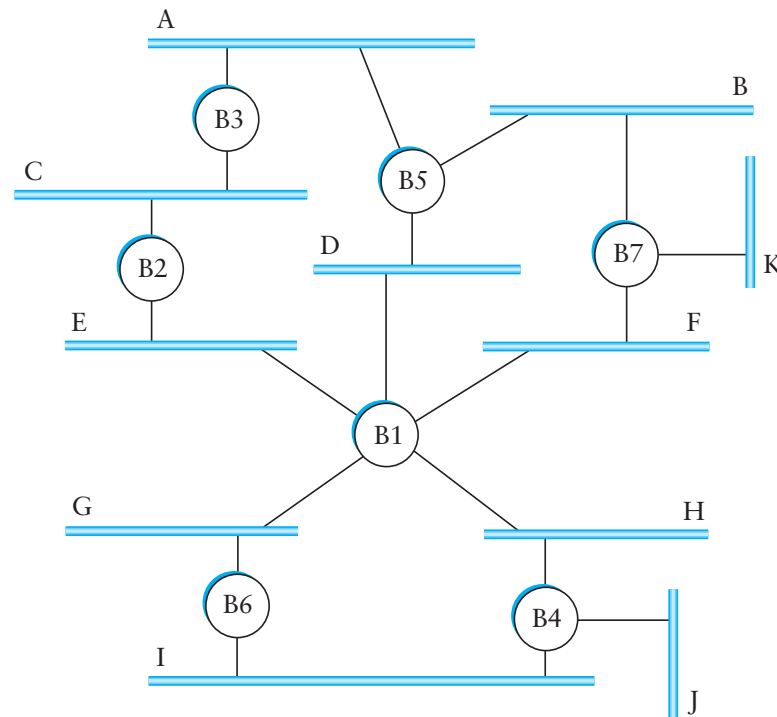
# Attack on a Learning Switch

- **Eve: wants to sniff all packets sent to Bob**
- **Same segment: easy (shared medium)**
- **Different segment on a learning bridge: hard**
  - Once bridge learns Bob's port, stop broadcasting
- **How can Eve force the bridge to keep broadcasting?**
  - Flood the network with frames with spoofed src addr!



# Dealing with Loops

- **Problem: people may create loops in LAN!**
  - Accidentally, or to provide redundancy
  - Don't want to forward packets indefinitely





# Enter Radia Perlman

“...we have designed an algorithm that allows the extended network to consist of an arbitrary topology. (...)

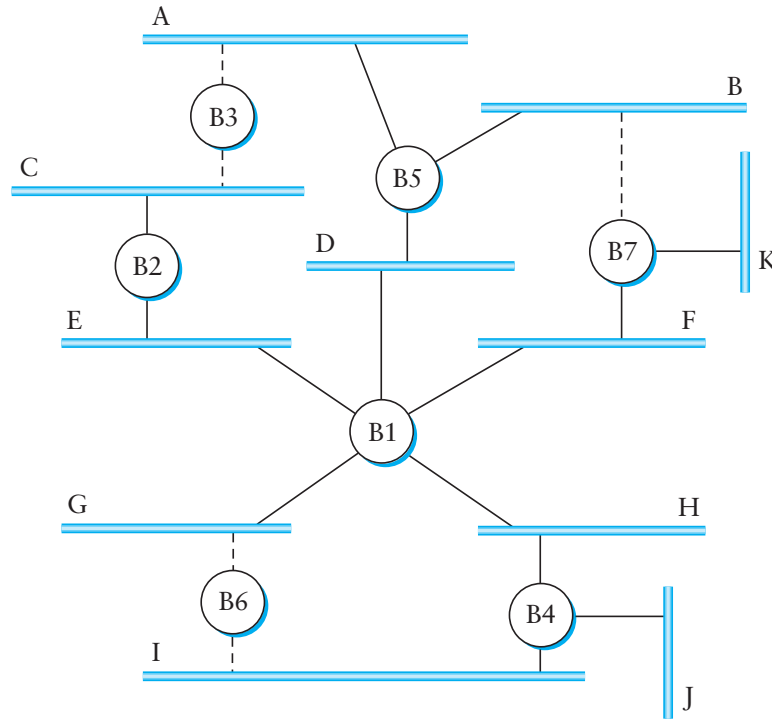
The algorithm (...) computes a subset of the topology that connects all LANs yet is loop-free (a spanning tree).”



Perlman, Radia (1985). "An Algorithm for Distributed Computation of a Spanning Tree in an Extended LAN". *ACM SIGCOMM Computer Communication Review*. **15** (4): 44–53. [doi:10.1145/318951.319004](https://doi.org/10.1145/318951.319004)



# Spanning Tree



- **Need to disable ports, so that no loops in network**
- **Like creating a spanning tree in a graph**
  - View switches and networks as nodes, ports as edges



# Distributed Spanning Tree Algorithm

- **Every bridge has a unique ID (Ethernet address)**
- **Goal:**
  - Bridge with the smallest ID is the root
  - Each segment has one designated bridge, responsible for forwarding its packets towards the root
    - Bridge closest to root is designated bridge
    - If there is a tie, bridge with lowest ID wins



# Spanning Tree Protocol

- **Send message when you think you are the root**
- **Otherwise, forward messages from best known root**
  - Add one to distance before forwarding
  - Don't forward over discarding ports (see next slide)
- **Spanning Tree messages contain:**
  - ID of bridge sending the message
  - ID sender believes to be the root
  - Distance (in hops) from sender to root
- **Bridges remember best config msg on each port**
- **In the end, only root is generating messages**

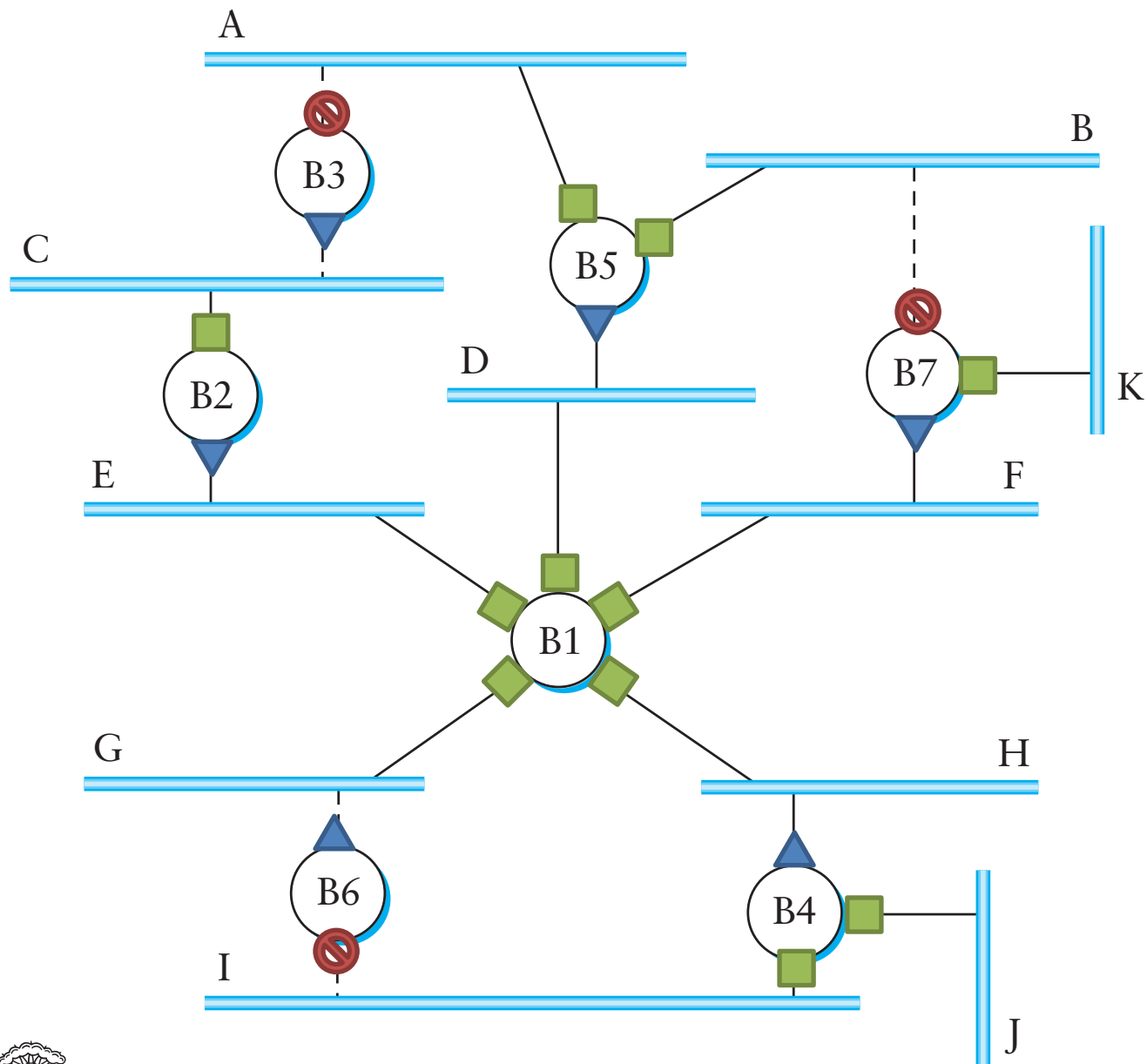


# Spanning Tree Protocol (cont.)

- **Forwarding and Broadcasting**
- **Port states\***:
  - **Root port**: a port the bridge uses to reach the root
  - **Designated port**: the lowest-cost port attached to a single segment
  - If a port is not a root port or a designated port, it is a **discarding port**.



\* In a later protocol RSTP, there can be ports configured as backups and alternates.



# Algorhyme

**I think that I shall never see  
a graph more lovely than a tree.  
A tree whose crucial property  
is loop-free connectivity.  
A tree that must be sure to span  
so packet can reach every LAN.  
First the root must be selected.  
By ID, it is elected.  
Least cost paths from root are traced.  
In the tree, these paths are placed.  
A mesh is made by folks like me,  
then bridges find a spanning tree.**

**Radia Perlman**



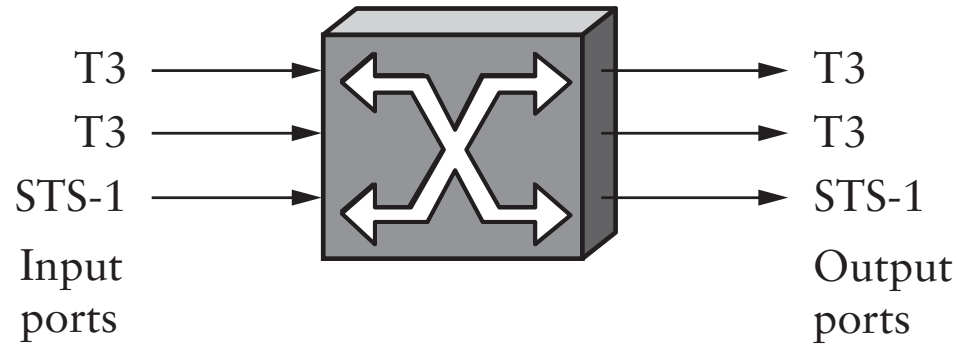
# Limitations of Bridges

- **Scaling**
  - Spanning tree algorithm doesn't scale
  - Broadcast does not scale
  - No way to route around congested links, *even if path exists*
- **May violate assumptions**
  - Could confuse some applications that assume single segment
    - Much more likely to drop packets
    - Makes latency between nodes non-uniform
  - Beware of transparency





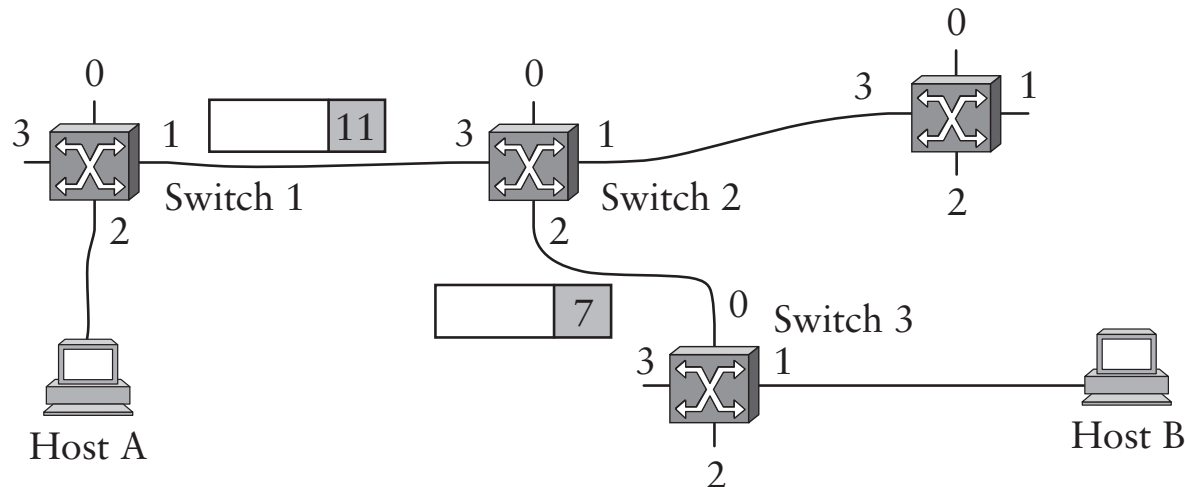
# Switching



- **Switches must be able to, given a packet, determine the outgoing port**
- **3 ways to do this:**
  - Virtual Circuit Switching
  - Datagram Switching
  - Source Routing



# Virtual Circuit Switching



- **Explicit set-up and tear down phases**
  - Establishes Virtual Circuit Identifier on each link
  - Each switch stores VC table
- **Subsequent packets follow same path**
  - Switches map [in-port, in-VCI] : [out-port, out-VCI]
- **Also called *connection-oriented* model**

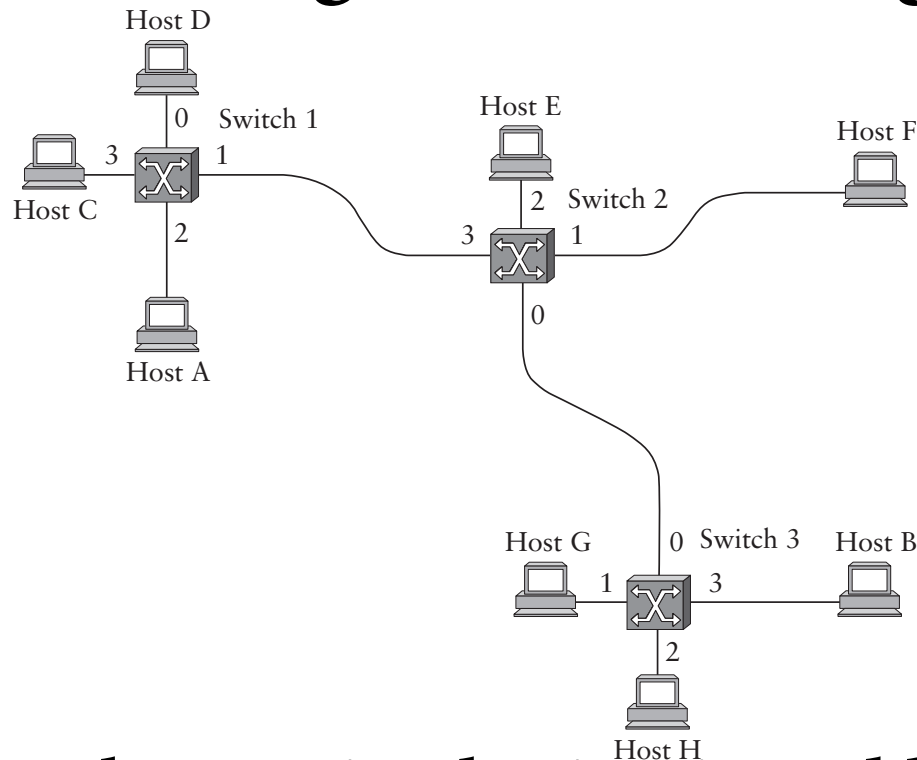


# Virtual Circuit Model

- **Requires one RTT before sending first packet**
- **Connection request contain full destination address, subsequent packets only small VCI**
- **Setup phase allows reservation of resources, such as bandwidth or buffer-space**
  - Any problems here?
- **If a link or switch fails, must re-establish whole circuit**
- **Example: ATM, MPLS**



# Datagram Switching



Switch 2

Addr	Port
A	3
B	0
C	3
D	3
E	2
F	1
G	0
H	0

- Each packet carries destination address
- Switches maintain address-based tables
  - Maps [destination address]:[out-port]
- Also called *connectionless* model

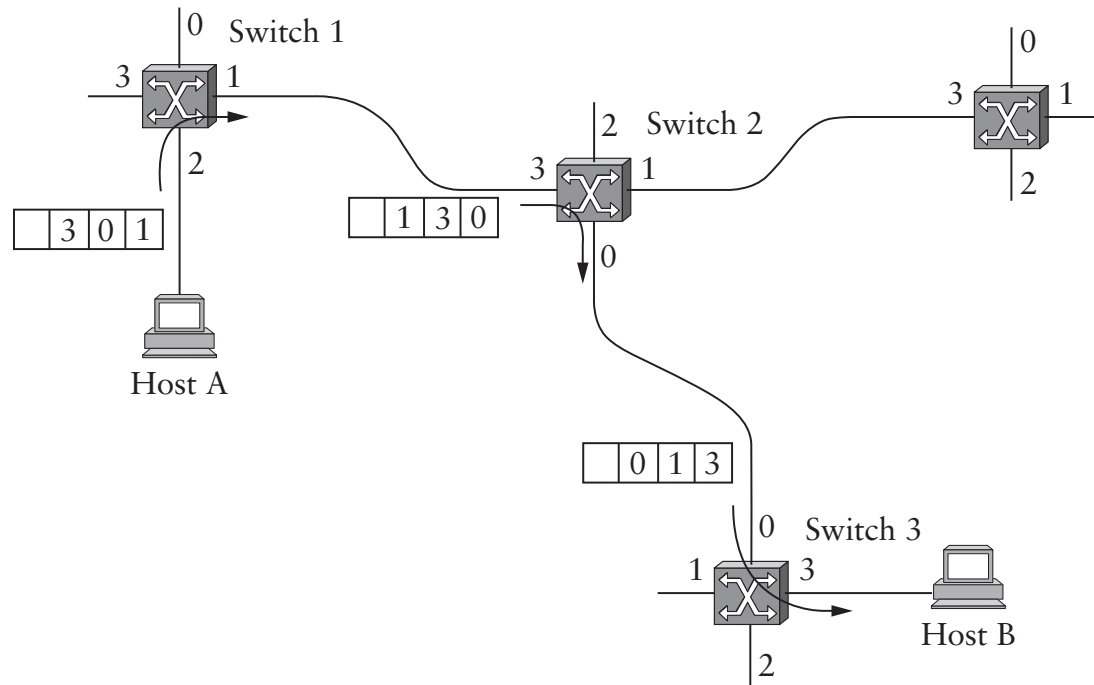


# Datagram Switching

- **No delay for connection setup**
- **Source can't know if network can deliver a packet**
- **Possible to route around failures**
- **Higher overhead per-packet**
- **Potentially larger tables at switches**



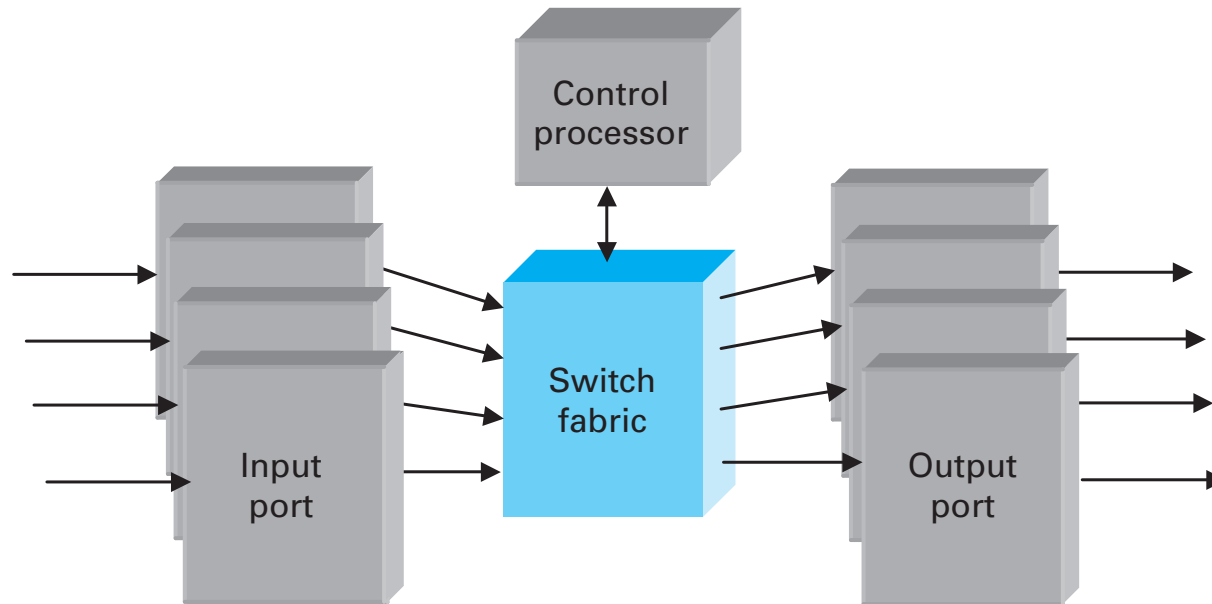
# Source Routing



- **Packets carry entire route: ports**
- **Switches need no tables!**
  - But end hosts must obtain the path information
- **Variable packet header**



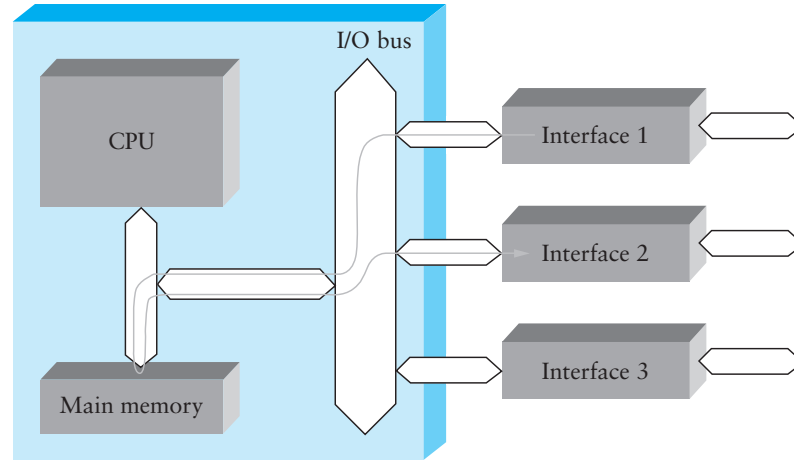
# Generic Switch Architecture



- **Goal: deliver packets from input to output ports**
- **Three potential performance concerns:**
  - Throughput in bytes/second
  - Throughput in packets/second
  - Latency



# Shared Memory Switch

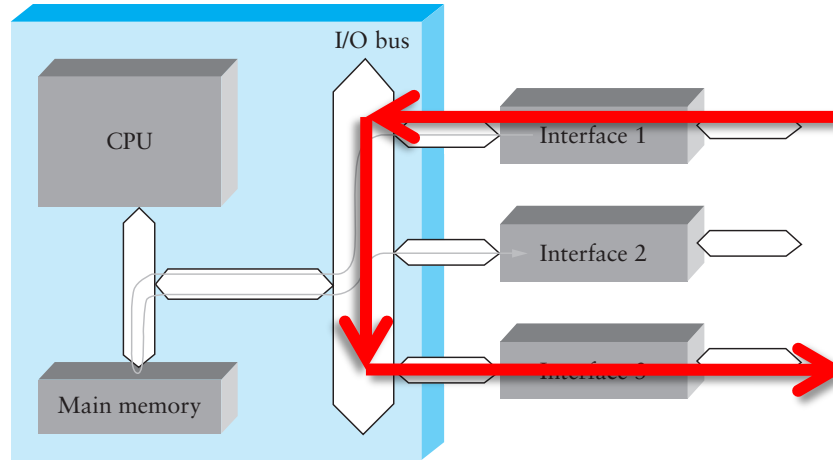


- **1<sup>st</sup> Generation – like a regular PC**
  - NIC DMAs packet to memory over I/O bus
  - CPU examines header, sends to destination NIC
  - I/O bus is serious bottleneck
  - For small packets, CPU may be limited too
  - Typically < 0.5 Gbps





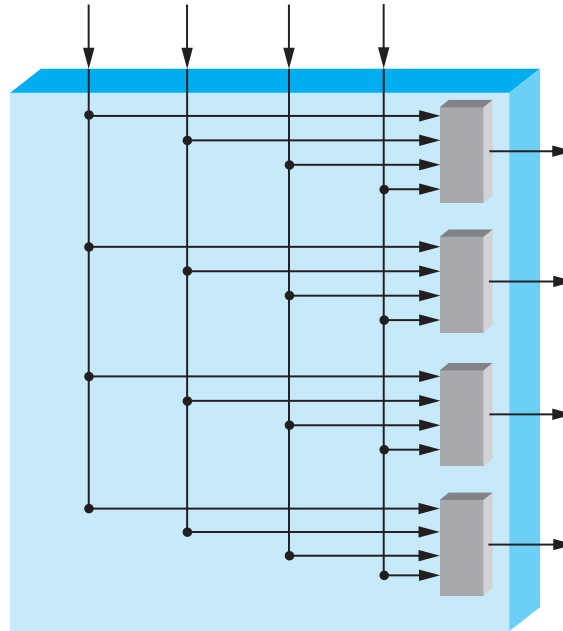
# Shared Bus Switch



- **2<sup>st</sup> Generation**
  - NIC has own processor, cache of forwarding table
  - Shared bus, doesn't have to go to main memory
  - Typically limited to bus bandwidth
    - (Cisco 5600 has a 32Gbps bus)



# Point to Point Switch



- **3<sup>rd</sup> Generation: overcomes single-bus bottleneck**
- **Example: Cross-bar switch**
  - Any input-output permutation
  - Multiple inputs to same output requires trickery
  - Cisco 12000 series: 60Gbps



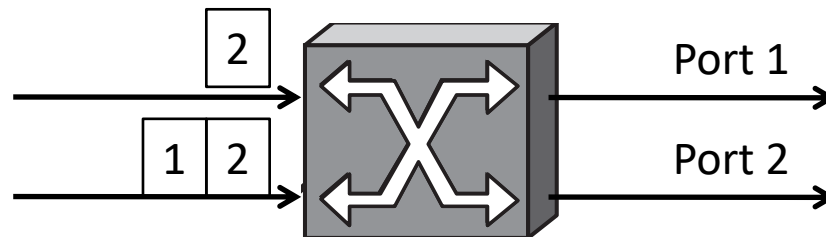
# Cut through vs. Store and Forward

- **Two approaches to forwarding a packet**
  - Receive a full packet, then send to output port
  - Start retransmitting as soon as you know output port, before full packet
- **Cut-through routing can greatly decrease latency**
- **Disadvantage**
  - Can waste transmission (classic *optimistic* approach)
    - CRC may be bad
    - If Ethernet collision, may have to send runt packet on output link



# Buffering

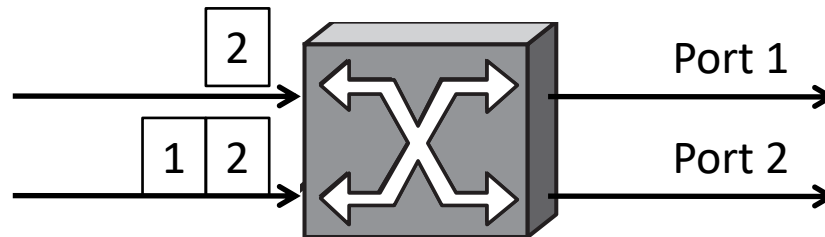
- Buffering of packets can happen at input ports, fabric, and/or output ports
- Queuing discipline is very important
- Consider FIFO + input port buffering
  - Only one packet per output port at any time
  - If multiple packets arrive for port 2, they may block packets to other ports that are free
  - *Head-of-line blocking*: can limit throughput to  $\sim 58\%$  under some reasonable conditions\*



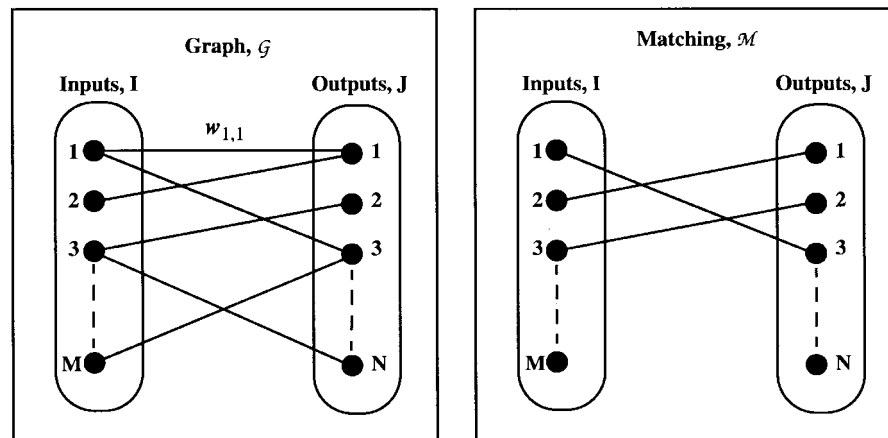
\* For independent, uniform traffic, with same-size frames



# Head-of-Line Blocking

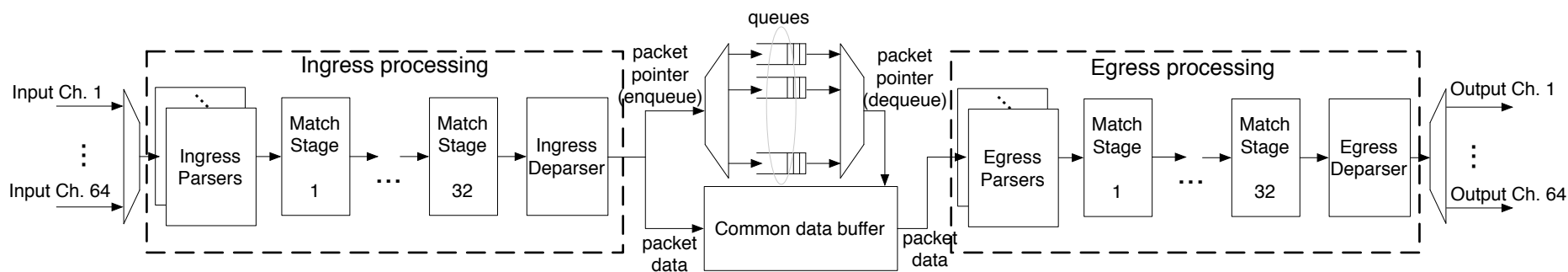


- **Solution: Virtual Output Queueing**
  - Each input port has  $n$  FIFO queues, one for each output
  - Switch using matching in a bipartite graph
  - Shown to achieve 100% throughput\*



# Current Developments

- **Switches are becoming programmable**
  - Match-action paradigm
  - Custom protocols, encapsulation, metering, monitoring



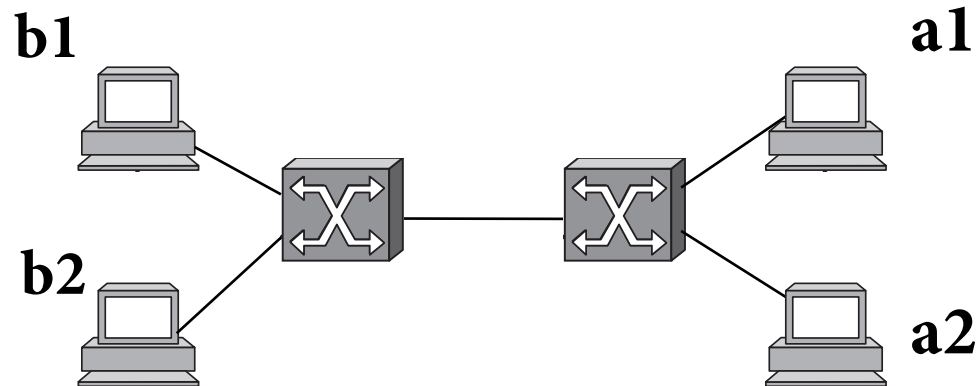
- **Current speeds reach 12.8Tbps (32x400Gbps or 256x50Gbps) on a single programmable switching chip**



**We did not cover these...**



# VLANs

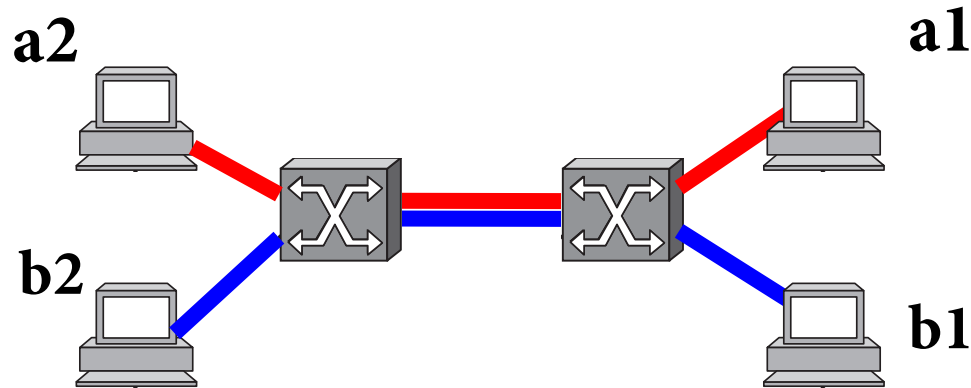


- **Company network, A and B departments**
  - Broadcast traffic does not scale
  - May not *want* traffic between the two departments
  - Topology has to mirror physical locations
  - What if employees move between offices?





# VLANs



- **Solution: Virtual LANs**
  - Assign switch ports to a VLAN ID (color)
  - Isolate traffic: only same color
  - Trunk links may belong to multiple VLANs
  - Encapsulate packets: add 12-bit VLAN ID
- **Easy to change, no need to rewire**



# Coming Up

- **Connecting multiple networks: IP and the Network Layer**

