

### Research paper

- Measured Capacity of an Ethernet: Myths and Reality
- Theoretical work seems to suggest that Ethernet works saturate at 37%. Realistic networks can offer higher throughputs



Feb-4-03

4/598N: Computer Networks

1

### Lessons learnt

- Don't install long cables: to cover a large area, break up the cable with bridges or gateways (routers), not repeaters.
- Don't put too many hosts on one cable: use gateways to break the network into communities of interest, trading higher delay for inter-community traffic for better intra-community response time and throughput
  - Current ethernets are 10BaseT or 100BaseT and use switches



Feb-4-03

4/598N: Computer Networks

2

### Lessons learnt

- Implement the protocol correctly: proper collision detection and binary exponential backoff in interface or host software is essential to good performance
- Use the largest possible packet size: this keeps the packet count down, reducing the likelihood of collision and not incidentally reducing overheads internal to hosts
  - Especially important for Gigabit Ethernets. They define Jumbo frames (9KB packets). More on this for HWP 2
- Don't mix serious real-time and serious bulk-data applications: it is not possible to simultaneously guarantee the lowest delay and the highest throughput (although for moderate requirements both kinds of applications coexist well)



Feb-4-03

4/598N: Computer Networks

3

### Switching and Forwarding

- Outline
  - Store-and-Forward Switches
  - Bridges and Extended LANs
  - Cell Switching
  - Segmentation and Reassembly



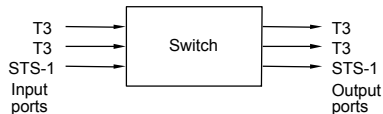
Feb-4-03

4/598N: Computer Networks

4

### Scalable Networks

- Switch
  - forwards packets from input port to output port
  - port selected based on address in packet header



- Advantages
  - cover large geographic area (tolerate latency)
  - support large numbers of hosts (scalable bandwidth)

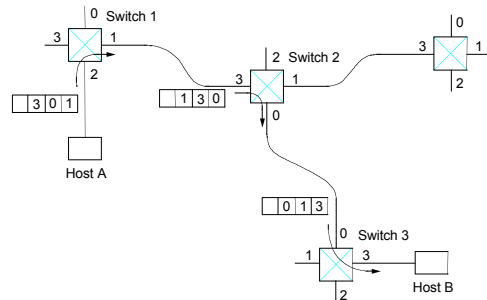


Feb-4-03

4/598N: Computer Networks

5

### Source Routing



Feb-4-03

4/598N: Computer Networks

6

### Virtual Circuit Switching

- Explicit connection setup (and tear-down) phase
- Subsequence packets follow same circuit
- Sometimes called connection-oriented model

- Analogy: phone call
- Each switch maintains a VC table

Feb-4-03
4/598N: Computer Networks
7

### Virtual Circuit Model

- Typically wait full RTT for connection setup before sending first data packet.
- While the connection request contains the full address for destination, each data packet contains only a small identifier, making the per-packet header overhead small.
- If a switch or a link in a connection fails, the connection is broken and a new one needs to be established.
- Connection setup provides an opportunity to reserve resources.

Feb-4-03
4/598N: Computer Networks
8

### Datagram Switching

- No connection setup phase
- Each packet forwarded independently
- Sometimes called connectionless model

- Analogy: postal system
- Each switch maintains a forwarding (routing) table

Feb-4-03
4/598N: Computer Networks
9

### Datagram Model

- There is no round trip time delay waiting for connection setup; a host can send data as soon as it is ready.
- Source host has no way of knowing if the network is capable of delivering a packet or if the destination host is even up.
- Since packets are treated independently, it is possible to route around link and node failures.
- Since every packet must carry the full address of the destination, the overhead per packet is higher than for the connection-oriented model.

Feb-4-03
4/598N: Computer Networks
10

### Bridges and Extended LANs

- LANs have physical limitations (e.g., 2500m)
- Connect two or more LANs with a bridge
  - accept and forward strategy
  - level 2 connection (does not add packet header)

- Ethernet Switch = Bridge on Steroids

Feb-4-03
4/598N: Computer Networks
11

### Learning Bridges

- Do not forward when unnecessary
- Maintain forwarding table

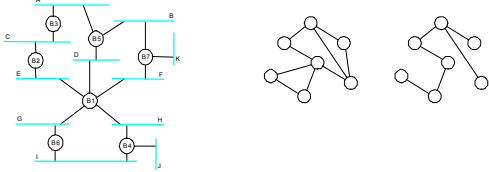
Host	Port
A	1
B	1
C	1
X	2
Y	2
Z	2

- Learn table entries based on source address
- Table is an optimization; need not be complete
- Always forward broadcast frames

Feb-4-03
4/598N: Computer Networks
12

## Spanning Tree Algorithm

- Problem: loops - no mechanism to remove looping frames



- Bridges run a distributed spanning tree algorithm
  - select which bridges actively forward
  - developed by Radia Perlman
  - now IEEE 802.1 specification



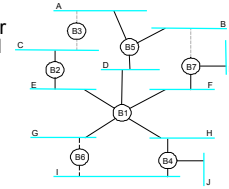
Feb-4-03

4/598N: Computer Networks

13

## Algorithm Overview

- Each bridge has unique id (e.g., B1, B2, B3)
- Select bridge with smallest id as root
- Select bridge on each LAN closest to root as designated bridge (use id to break ties)
  - Each bridge forwards frames over each LAN for which it is the designated bridge



Feb-4-03

4/598N: Computer Networks

14

## Algorithm Details

- Bridges exchange configuration messages
  - id for bridge sending the message
  - id for what the sending bridge believes to be root bridge
  - distance (hops) from sending bridge to root bridge
- Each bridge records current best configuration message for each port
- Initially, each bridge believes it is the root



Feb-4-03

4/598N: Computer Networks

15

## Algorithm Detail (cont)

- When learn not root, stop generating config messages
  - in steady state, only root generates configuration messages
- When learn not designated bridge, stop forwarding config messages
  - in steady state, only designated bridges forward config messages
- Root continues to periodically send config messages
- If any bridge does not receive config message after a period of time, it starts generating config messages claiming to be the root



Feb-4-03

4/598N: Computer Networks

16

## Broadcast and Multicast

- Forward all broadcast/multicast frames
  - current practice
- Learn when no group members downstream
- Accomplished by having each member of group G send a frame to bridge multicast address with G in source field



Feb-4-03

4/598N: Computer Networks

17

## Limitations of Bridges

- Do not scale
  - spanning tree algorithm does not scale - traffic gets bridged through the root bridge
    - Spanning tree is designed to avoid loops, not traffic balancing: redundant routes are ignored
  - broadcast does not scale
- Do not accommodate heterogeneity
- Caution: beware of transparency



Feb-4-03

4/598N: Computer Networks

18

### SmartBridges

- [http://www.researchchannel.com/programs/uw/Asx/cse\\_smbr\\_1300k.asx](http://www.researchchannel.com/programs/uw/Asx/cse_smbr_1300k.asx)
- <http://www.uwv.org/programs/displayevent.asp?rid=782>
- Hybrid between IP routing and bridging



Feb-4-03

4/598N: Computer Networks

19

### Cell Switching (ATM)

- Connection-oriented packet-switched network
- Used in both WAN and LAN settings
- Signaling (connection setup) Protocol: Q.2931
- Specified by ATM forum
- Packets are called cells
  - 5-byte header + 48-byte payload
- Commonly transmitted over SONET
  - other physical layers possible



Feb-4-03

4/598N: Computer Networks

20

### Variable vs Fixed-Length Packets

- No Optimal Length
  - if small: high header-to-data overhead
  - if large: low utilization for small messages
- Fixed-Length Easier to Switch in Hardware
  - simpler
  - enables parallelism



Feb-4-03

4/598N: Computer Networks

21

### Big vs Small Packets

- Small Improves Queue behavior
  - finer-grained pre-emption point for scheduling link
    - maximum packet = 4KB
    - link speed = 100Mbps
    - transmission time =  $4096 \times 8/100 = 327.68\mu s$
    - high priority packet may sit in the queue  $327.68\mu s$
    - in contrast,  $53 \times 8/100 = 4.24\mu s$  for ATM
  - near cut-through behavior
    - two 4KB packets arrive at same time
    - link idle for  $327.68\mu s$  while both arrive
    - at end of  $327.68\mu s$ , still have 8KB to transmit
    - in contrast, can transmit first cell after  $4.24\mu s$
    - at end of  $327.68\mu s$ , just over 4KB left in queue



Feb-4-03

4/598N: Computer Networks

22

### Big vs Small (cont)

- Small Improves Latency (for voice)
  - voice digitally encoded at 64KBps (8-bit samples at 8KHz)
  - need full cell's worth of samples before sending cell
  - example: 1000-byte cells implies 125ms per cell (too long)
  - smaller latency implies no need for echo cancellors
- ATM Compromise: 48 bytes =  $(32+64)/2$



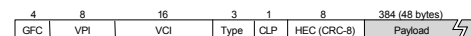
Feb-4-03

4/598N: Computer Networks

23

### Cell Format

- User-Network Interface (UNI)



- host-to-switch format
  - GFC: Generic Flow Control (still being defined)
  - VCI: Virtual Circuit Identifier
  - VPI: Virtual Path Identifier
  - Type: management, congestion control, AAL5 (later)
  - CLPL Cell Loss Priority
  - HEC: Header Error Check (CRC-8)
- Network-Network Interface (NNI)
    - switch-to-switch format
    - GFC becomes part of VPI field



Feb-4-03

4/598N: Computer Networks

24

### Segmentation and Reassembly

- ATM Adaptation Layer (AAL)
  - AAL 1 and 2 designed for applications that need guaranteed rate (e.g., voice, video)
  - AAL 3/4 designed for packet data
  - AAL 5 is an alternative standard for packet data

Feb-4-03
4/598N: Computer Networks
25

### AAL 3/4

- Convergence Sublayer Protocol Data Unit (CS-PDU)
 

8	8	16	< 64 KB	0-24	8	8	16
CPI	Btag	BAsize	User data	Pad	0	Etag	Len

  - CPI: commerce part indicator (version field)
  - Btag/Etag: beginning and ending tag
  - BAsize: hint on amount of buffer space to allocate
  - Length: size of whole PDU

Feb-4-03
4/598N: Computer Networks
26

### Cell Format

- Type
  - BOM: beginning of message
  - COM: continuation of message
  - EOM end of message
- SEQ: sequence of number
- MID: message id
- Length: number of bytes of PDU in this cell

40	2	4	10	352 (44 bytes)	6	10
ATM header	Type	SEQ	MID	Payload	Length	CRC-10

Feb-4-03
4/598N: Computer Networks
27

### AAL5

- CS-PDU Format
 

< 64 KB	0-47 bytes	16	16	32
Data	Pad	Reserved	Len	CRC-32

  - pad so trailer always falls at end of ATM cell
  - Length: size of PDU (data only)
  - CRC-32 (detects missing or misordered cells)
- Cell Format
  - end-of-PDU bit in Type field of ATM header

Feb-4-03
4/598N: Computer Networks
28