



15-441: Computer Networking

Lecture 8 – Bridging, Addressing and Forwarding

Copyright ©, 2007-11 Carnegie Mellon University

Scale



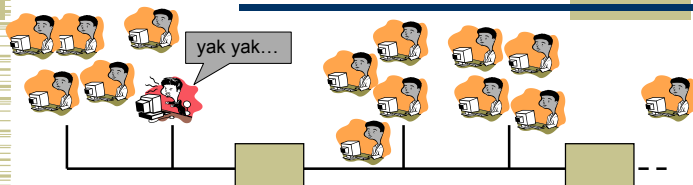
- What breaks when we keep adding people to the same wire?

Fall 2011

Lecture 8: Bridging/Addressing/Forwarding

2

Scale



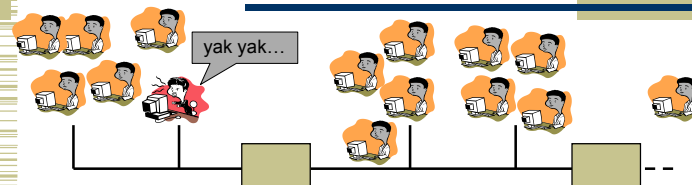
- What breaks when we keep adding people to the same wire?
- Only solution: split up the people onto multiple wires
 - But how can they talk to each other?

Fall 2011

Lecture 8: Bridging/Addressing/Forwarding

3

Problem 1 – Reconnecting LANs



- When should these boxes forward packets between wires?
- How do you specify a destination?
- How does your packet find its way?

Fall 2011

Lecture 8: Bridging/Addressing/Forwarding

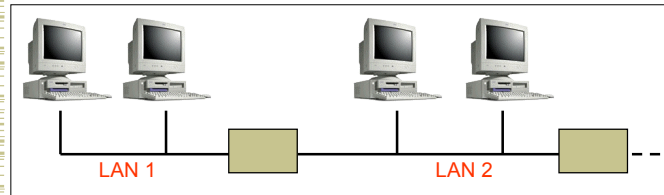
4

Outline

- **Bridging**
- Internetworks
 - Methods for packet forwarding
- Traditional IP addressing

Building Larger LANs: Bridges

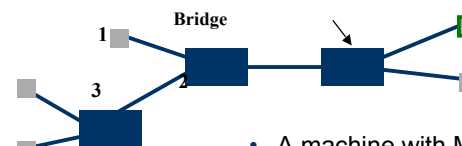
- Extend reach of a single shared medium
- Connect two or more “segments” by copying data frames between them
 - Only copy data when needed → key difference from repeaters/hubs
 - Reduce collision domain compared with single LAN
 - Separate segments can send at once → much greater bandwidth
- Challenge: learning which packets to copy across links



Transparent Bridges

- Design goals:
 - Self-configuring without hardware or software changes
 - Bridge do not impact the operation of the individual LANs
- Three parts to making bridges transparent:
 - Forwarding frames
 - Learning addresses/host locations
 - Spanning tree algorithm

Frame Forwarding



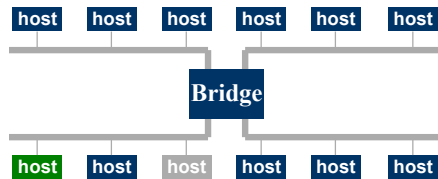
MAC Address	Port	Age
A21032C9A591	1	36
99A323C90842	2	01
8711C98900AA	2	15
301B2369011C	2	16
695519001190	3	11

- A machine with MAC Address lies in the direction of number port of the bridge
- For every packet, the bridge “looks up” the entry for the packets destination MAC address and forwards the packet on that port.
 - Other packets are broadcast – why?
- Timer is used to flush old entries

Learning Bridges



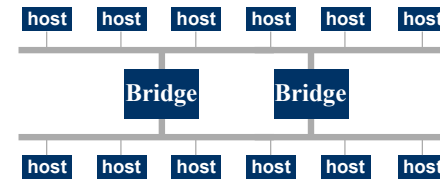
- Manually filling in bridge tables?
 - Time consuming, error-prone
- Keep track of source address of packets arriving on every link, showing what segment hosts are on
 - Fill in the forwarding table based on this information



Spanning Tree Bridges



- More complex topologies can provide redundancy.
 - But can also create loops.
- What is the problem with loops?
- Solution: spanning tree



Spanning Tree Protocol Overview



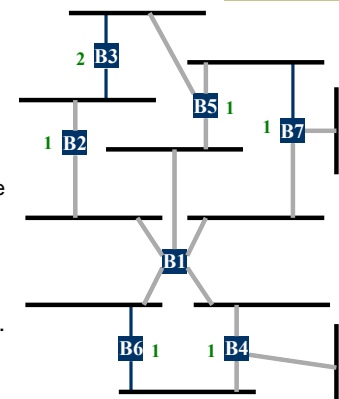
Embed a tree that provides a single unique path to each destination:

- Elect a single bridge as a root bridge
- Each bridge calculates the distance of the shortest path to the root bridge
- Each LAN identifies a *designated bridge*, the bridge closest to the root. It will forward packets to the root.
- Each bridge determines a *root port*, which will be used to send packets to the root
- Identify the ports that form the spanning tree

Spanning Tree Algorithm Steps

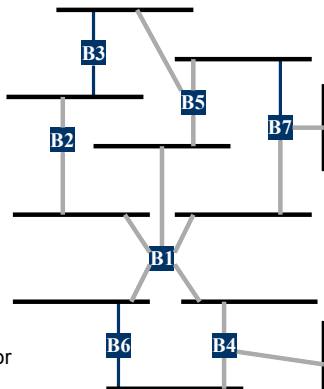


- Root of the spanning tree is the bridge with the lowest identifier.
 - All ports are part of tree
- Each bridge finds shortest path to the root.
 - Remembers port that is on the shortest path
 - Used to forward packets
- Select for each LAN the designated bridge that has the shortest path to the root.
 - Identifier as tie-breaker
 - Responsible for that LAN



Spanning Tree Algorithm

- Each node sends configuration message to all neighbors.
 - Identifier of the sender
 - Id of the presumed root
 - Distance to the presumed root
 - E.g. B5 sends (B5, B5, 0)
- When B receive a message, it decide whether the solution is better than their local solution.
 - A root with a lower identifier?
 - Same root but lower distance?
 - Same root, distance but sender has lower identifier?
- After convergence, each bridge knows the root, distance to root, root port, and designated bridge for each LAN.



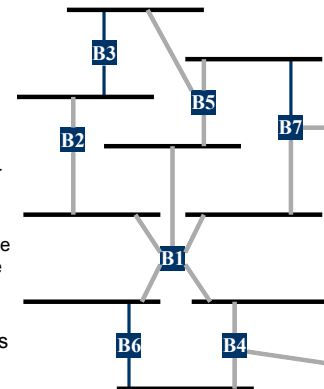
Fall 2011

Lecture 8: Bridging/Addressing/Forwarding

13

Spanning Tree Algorithm (part 2)

- Each bridge B can now select which of its ports make up the spanning tree:
 - B's root port
 - All ports for which B is the designated bridge on the LAN
- Bridges can now configure their ports.
 - Forwarding state* or *blocked state*, depending on whether the port is part of the spanning tree
- Root periodically sends configuration messages and bridges forward them over LANs they are responsible for.



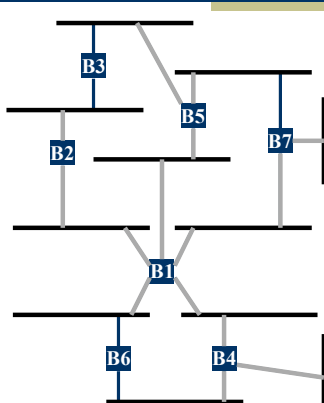
Fall 2011

Lecture 8: Bridging/Addressing/Forwarding

14

Spanning Tree Algorithm Example

- Node B2:
 - Sends (B2, B2, 0)
 - Receives (B1, B1, 0) from B1
 - Sends (B2, B1, 1) "up"
 - Continues the forwarding forever
- Node B1:
 - Will send notifications forever
- Node B7:
 - Sends (B7, B7, 0)
 - Receives (B1, B1, 0) from B1
 - Sends (B7, B1, 1) "up" and "right"
 - Receives (B5, B5, 0) - ignored
 - Receives (B5, B1, 1) - better
 - Continues forwarding the B1 messages forever to the "right"

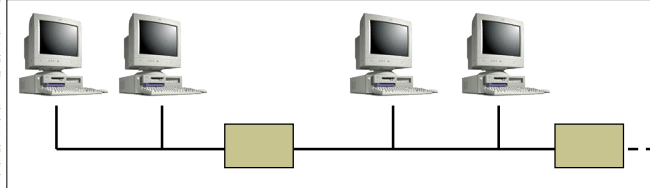


Fall 2011

Lecture 8: Bridging/Addressing/Forwarding

15

Problem 2 – Bridging Weaknesses



- Doesn't handle incompatible LAN technologies
- How well does it scale?

Fall 2011

Lecture 8: Bridging/Addressing/Forwarding

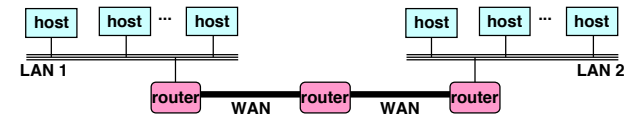
16

Outline

- Bridging
- **Internetworks**
 - Methods for packet forwarding
- Traditional IP addressing

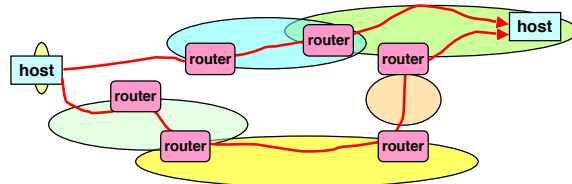
What is an Internetwork?

- Multiple incompatible LANs can be physically connected by specialized computers called **routers**
- The connected networks are called an **internetwork**
 - The "**Internet**" is one (very big & successful) example of an internetwork



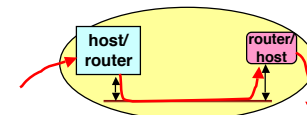
LAN 1 and LAN 2 might be completely different, totally incompatible LANs (e.g., Ethernet and ATM)

Logical Structure of Internet



- Ad hoc interconnection of networks
 - No particular topology
 - Vastly different router & link capacities
- Send packets from source to destination by hopping through networks
 - Router connect one network to another
 - Different paths to destination may exist

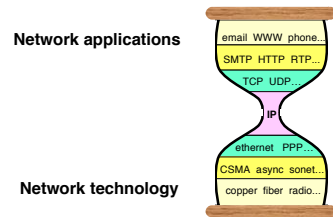
Routing Through Single Network



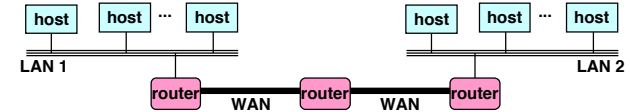
- Path Consists of Series of Hops
 - Source – Router
 - Router – Router (typically high-speed, point-to-point link)
 - Router – Destination
- Each Hop Uses Link-Layer Protocol
 - Determine hop destination
 - Based on destination
 - Send over local network
 - Put on header giving MAC address of intermediate router (or final destination)

Internet Protocol (IP)

- Hour Glass Model
 - Create abstraction layer that hides underlying technology from network application software
 - Make as minimal as possible
 - Allows range of current & future technologies
 - Can support many different types of applications



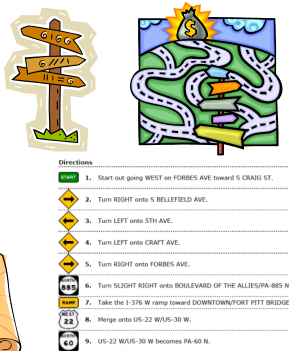
Problem 3: Internetwork Design



- How do I designate a distant host?
 - Addressing / naming
- How do I send information to a distant host?
 - What gets sent?
 - What route should it take?
- Must support:
 - Heterogeneity LAN technologies
 - Scalability → ensure ability to grow to worldwide scale

Getting to a Destination

- How do you get driving directions?
- Intersections → routers
- Roads → links/networks
- Roads change slowly



Forwarding Packets: Choices

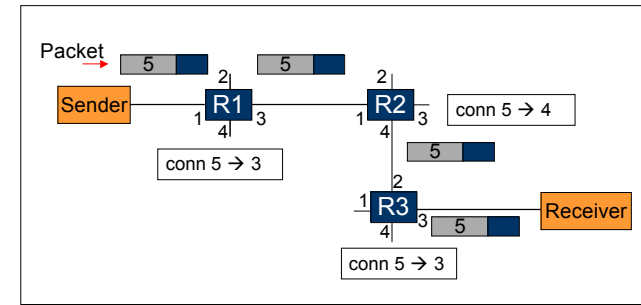
- Table of virtual circuits
 - Connection routed through network to set up state
 - Packets forwarded using connection state
- Source routing
 - Packet carries path
- Table of global addresses (IP)
 - Routers keep next hop for destination
 - Packets carry destination address

Simplified Virtual Circuits



- Connection setup phase
 - Use other means to route setup request
 - Each router allocates flow ID on local link
- Each packet carries connection ID
 - Sent from source with 1st hop connection ID
- Router processing
 - Lookup flow ID – simple table lookup
 - Replace flow ID with outgoing flow ID
 - Forward to output port

Simplified Virtual Circuits Example



Virtual Circuits



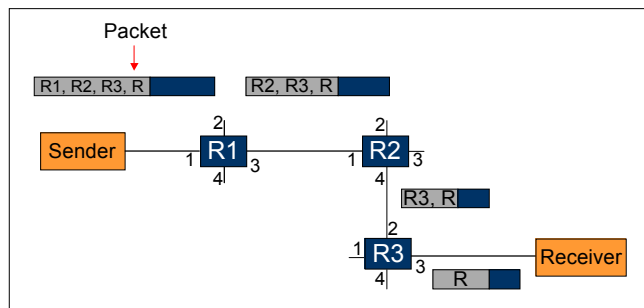
- Advantages
 - Efficient lookup (simple table lookup)
 - Can reserve bandwidth at connection setup
 - Easier for hardware implementations
- Disadvantages
 - Still need to route connection setup request
 - More complex failure recovery – must recreate connection state
- Typical use → fast router implementations
 - ATM – combined with fix sized cells
 - MPLS – tag switching for IP networks

Source Routing



- List entire path in packet
 - Driving directions (north 3 hops, east, etc..)
- Router processing
 - Strip first step from packet
 - Examine next step in directions
 - Forward to next step

Source Routing Example



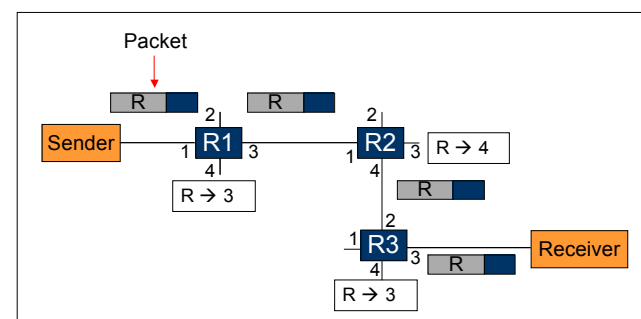
Source Routing

- Advantages
 - Switches can be very simple and fast
- Disadvantages
 - Variable (unbounded) header size
 - Sources must know or discover topology (e.g., failures)
- Typical uses
 - Ad-hoc networks (DSR)
 - Machine room networks (Myrinet)

Global Addresses (IP)

- Each packet has destination address
- Each router has forwarding table of destination → next hop
 - At v and x: destination → east
 - At w and y: destination → south
 - At z: destination → north
- Distributed routing algorithm for calculating forwarding tables

Global Address Example



Global Addresses



- Advantages
 - Stateless – simple error recovery
- Disadvantages
 - Every switch knows about every destination
 - Potentially large tables
 - All packets to destination take same route
 - Need routing protocol to fill table

Comparison



	Source Routing	Global Addresses	Virtual Circuits
Header Size	Worst	OK – Large address	Best
Router Table Size	None	Number of hosts (prefixes)	Number of circuits
Forward Overhead	Best	Prefix matching (Worst)	Pretty Good
Setup Overhead	None	None	Connection Setup
Error Recovery	Tell all hosts	Tell all routers	Tell all routers and Tear down circuit and re-route

Problem 4: Router Table Size



- Global addressing networks (e.g., Internet, Ethernet bridging) require switches/routers to know next hop for all destinations
- How do we avoid large tables?



Outline



- Bridging
- Internetworks
 - Methods for packet forwarding
- Traditional IP addressing

Addressing in IP



- IP addresses are names of interfaces
 - E.g., 128.2.1.1
- Domain Name System (DNS) names are names of hosts
 - E.g., www.cmu.edu
- DNS binds host names to interfaces
- Routing binds interface names to paths

Router Table Size



- One entry for every host on the Internet
 - 630M (1/09) entries, doubling every 2.5 years
- One entry for every LAN
 - Every host on LAN shares prefix
 - Still too many and growing quickly
- One entry for every organization
 - Every host in organization shares prefix
 - Requires careful address allocation

Addressing Considerations



- Hierarchical vs. flat
 - Pennsylvania / Pittsburgh / Oakland / CMU / Seshan
vs.
Srinivasan Seshan: 123-45-6789
vs.
Srinivasan Seshan: (412) 268-0000
- What information would routers need to route to Ethernet addresses?
 - Need hierarchical structure for designing scalable binding from interface name to route!
- What type of Hierarchy?
 - How many levels?
 - Same hierarchy depth for everyone?
 - Same segment size for similar partition?

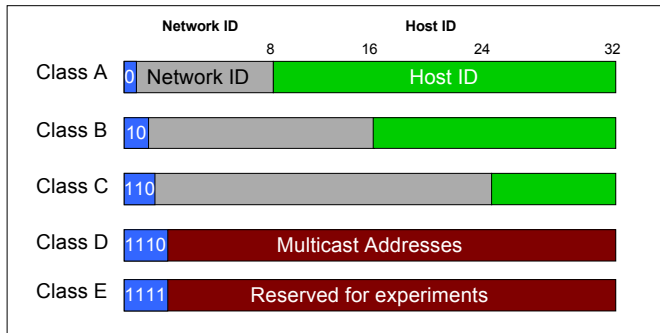
IP Addresses



- Fixed length: 32 bits
- Initial classful structure (1981) (not relevant now!!!)
- Total IP address size: 4 billion
 - Class A: 128 networks, 16M hosts
 - Class B: 16K networks, 64K hosts
 - Class C: 2M networks, 256 hosts

High Order Bits	Format	Class
0	7 bits of net, 24 bits of host	A
10	14 bits of net, 16 bits of host	B
110	21 bits of net, 8 bits of host	C

IP Address Classes (Some are Obsolete)



Original IP Route Lookup



- Address would specify prefix for forwarding table
 - Simple lookup
- www.cmu.edu address 128.2.11.43
 - Class B address – class + network is 128.2
 - Lookup 128.2 in forwarding table
 - Prefix – part of address that really matters for routing
- Forwarding table contains
 - List of class+network entries
 - A few fixed prefix lengths (8/16/24)
- Large tables
 - 2 Million class C networks

Subnet Addressing RFC917 (1984)

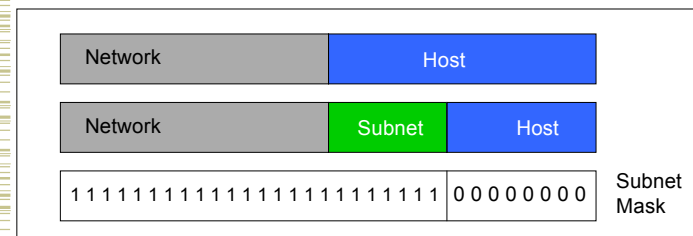


- Class A & B networks too big
 - Very few LANs have close to 64K hosts
 - For electrical/LAN limitations, performance or administrative reasons
- Need simple way to get multiple “networks”
 - Use bridging, multiple IP networks or split up single network address ranges (subnet)
- CMU case study in RFC
 - Chose not to adopt – concern that it would not be widely supported ☺

Subnetting



- Add another layer to hierarchy
- Variable length subnet masks
 - Could subnet a class B into several chunks



Subnetting Example

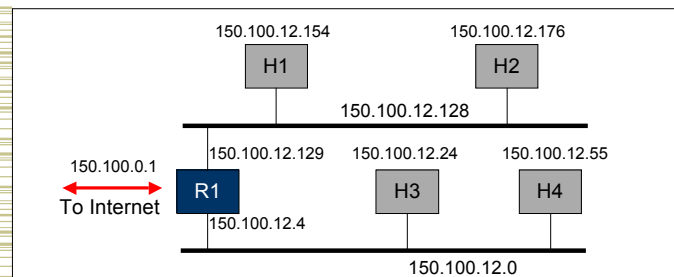


- Assume an organization was assigned address 150.100
- Assume < 100 hosts per subnet
- How many host bits do we need?
 - Seven
- What is the network mask?
 - 11111111 11111111 11111111 10000000
 - 255.255.255.128

Forwarding Example



- Assume a packet arrives with address 150.100.12.176
- Step 1: AND address with class + subnet mask



Aside: Interaction with Link Layer



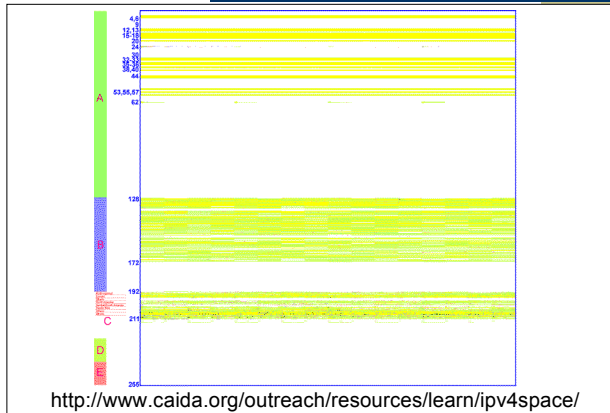
- How does one find the Ethernet address of a IP host?
- ARP - Address Resolution Protocol
 - Broadcast search for IP address
 - E.g., “who-has 128.2.184.45 tell 128.2.206.138” sent to Ethernet broadcast (all FF address)
 - Destination responds (only to requester using unicast) with appropriate 48-bit Ethernet address
 - E.g, “reply 128.2.184.45 is-at 0:d0:bc:f2:18:58” sent to 0:c0:4f:d:ed:c6

IP Address Problem (1991)



- Address space depletion
 - In danger of running out of classes A and B
 - Why?
 - Class C too small for most domains
 - Very few class A – very careful about giving them out
 - Class B – greatest problem
- Class B sparsely populated
 - But people refuse to give it back
- Large forwarding tables
 - 2 Million possible class C groups

IP Address Utilization ('97)



<http://www.caida.org/outreach/resources/learn/ipv4space/>

Important Concepts



- Hierarchical addressing critical for scalable system
 - Don't require everyone to know everyone else
 - Reduces number of updates when something changes

EXTRA SLIDES

How is IP Design Standardized?



- IETF
 - Voluntary organization
 - Meeting every 4 months
 - Working groups and email discussions
- “We reject kings, presidents, and voting; we believe in rough consensus and running code” (Dave Clark 1992)
 - Need 2 independent, interoperable implementations for standard

Addressing Considerations



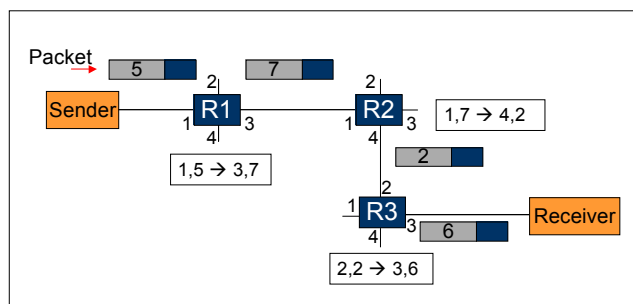
- Fixed length or variable length?
- Issues:
 - Flexibility
 - Processing costs
 - Header size
- Engineering choice: IP uses fixed length addresses

Virtual Circuits/Tag Switching



- Connection setup phase
 - Use other means to route setup request
 - Each router allocates flow ID on local link
 - Creates mapping of inbound flow ID/port to outbound flow ID/port
- Each packet carries connection ID
 - Sent from source with 1st hop connection ID
- Router processing
 - Lookup flow ID – simple table lookup
 - Replace flow ID with outgoing flow ID
 - Forward to output port

Virtual Circuits Examples



Virtual Circuits

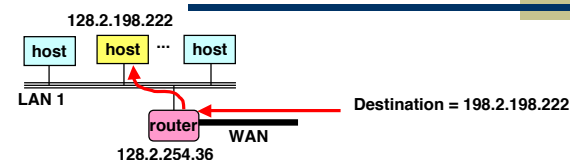


- Advantages
 - More efficient lookup (simple table lookup)
 - More flexible (different path for each flow)
 - Can reserve bandwidth at connection setup
 - Easier for hardware implementations
- Disadvantages
 - Still need to route connection setup request
 - More complex failure recovery – must recreate connection state
- Typical uses
 - ATM – combined with fix sized cells
 - MPLS – tag switching for IP networks

Some Special IP Addresses

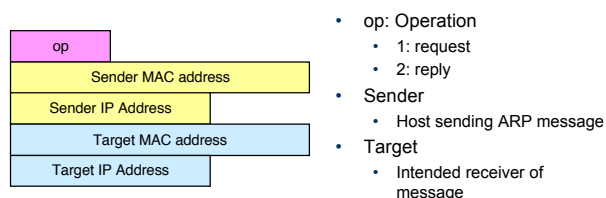
- 127.0.0.1: local host (a.k.a. the loopback address)
- Host bits all set to 0: network address
- Host bits all set to 1: broadcast address

Finding a Local Machine



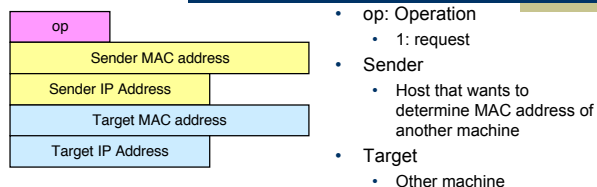
- Routing Gets Packet to Correct Local Network
 - Based on IP address
 - Router sees that destination address is of local machine
- Still Need to Get Packet to Host
 - Using link-layer protocol
 - Need to know hardware address
- Same Issue for Any Local Communication
 - Find local machine, given its IP address

Address Resolution Protocol (ARP)



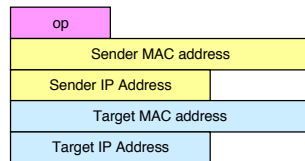
- op: Operation
 - 1: request
 - 2: reply
- Sender
 - Host sending ARP message
- Target
 - Intended receiver of message
- Diagrammed for Ethernet (6-byte MAC addresses)
- Low-Level Protocol
 - Operates only within local network
 - Determines mapping from IP address to hardware (MAC) address
 - Mapping determined dynamically
 - No need to statically configure tables
 - Only requirement is that each host know its own IP address

ARP Request



- op: Operation
 - 1: request
- Sender
 - Host that wants to determine MAC address of another machine
- Target
 - Other machine
- Requestor
 - Fills in own IP and MAC address as "sender"
 - Why include its MAC address?
- Mapping
 - Fills desired host IP address in target IP address
- Sending
 - Send to MAC address `ff:ff:ff:ff:ff:ff`
 - Ethernet broadcast

ARP Reply



- op: Operation
 - 2: reply
- Sender
 - Host with desired IP address
- Target
 - Original requestor
- Responder becomes “sender”
 - Fill in own IP and MAC address
 - Set requestor as target
 - Send to requestor’s MAC address

ARP Example



```

Time           Source MAC   Dest MAC
09:37:53.729185 0:2:b3:8a:35:bf ff:ff:ff:ff:ff:ff 0806 60:
arp who-has 128.2.222.198 tell 128.2.194.66
09:37:53.729202 0:3:47:b8:e5:f3 0:2:b3:8a:35:bf 0806 42:
arp reply 128.2.222.198 is-at 0:3:47:b8:e5:f3
    
```

- Exchange Captured with windump
 - Windows version of tcpdump
- Requestor:
 - blackhole-ad.scs.cs.cmu.edu (128.2.194.66)
 - MAC address 0:2:b3:8a:35:bf
- Desired host:
 - bryant-tp2.vlsi.cs.cmu.edu (128.2.222.198)
 - MAC address 0:3:47:b8:e5:f3

Caching ARP Entries



- Efficiency Concern
 - Would be very inefficient to use ARP request/reply every time need to send IP message to machine
- Each Host Maintains Cache of ARP Entries
 - Add entry to cache whenever get ARP response
 - Set timeout of ~20 minutes

ARP Cache Example



```

• Show using command “arp -a”
Interface: 128.2.222.198 on Interface 0x1000003
Internet Address      Physical Address      Type
128.2.20.218          00-b0-8e-83-df-50     dynamic
128.2.102.129         00-b0-8e-83-df-50     dynamic
128.2.194.66          00-02-b3-8a-35-bf     dynamic
128.2.198.34          00-06-5b-f3-5f-42     dynamic
128.2.203.3           00-90-27-3c-41-11     dynamic
128.2.203.61          08-00-20-a6-ba-2b     dynamic
128.2.205.192         00-60-08-1e-9b-fd     dynamic
128.2.206.125         00-d0-b7-c5-b3-f3     dynamic
128.2.206.139         00-a0-c9-98-2c-46     dynamic
128.2.222.180         08-00-20-a6-ba-c3     dynamic
128.2.242.182         08-00-20-a7-19-73     dynamic
128.2.254.36          00-b0-8e-83-df-50     dynamic
    
```

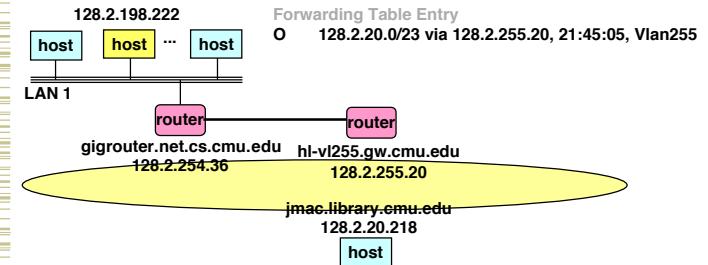

ARP Cache Surprise

- How come 3 machines have the same MAC address?

Interface: 128.2.222.198 on Interface 0x1000003

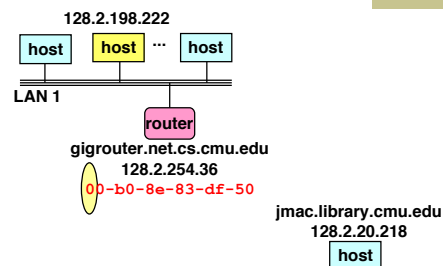
Internet Address	Physical Address	Type
128.2.20.218	00-b0-8e-83-df-50	dynamic
128.2.102.129	00-b0-8e-83-df-50	dynamic
128.2.194.66	00-02-b3-8a-35-bf	dynamic
128.2.198.34	00-06-5b-f3-5f-42	dynamic
128.2.203.3	00-90-27-3c-41-11	dynamic
128.2.203.61	08-00-20-a6-ba-2b	dynamic
128.2.205.192	00-60-08-1e-9b-fd	dynamic
128.2.206.125	00-d0-b7-c5-b3-f3	dynamic
128.2.206.139	00-a0-c9-98-2c-46	dynamic
128.2.222.180	08-00-20-a6-ba-c3	dynamic
128.2.242.182	08-00-20-a7-19-73	dynamic
128.2.254.36	00-b0-8e-83-df-50	dynamic

CMU's Internal Network Structure



- CMU Uses Routing Internally
 - Maintains forwarding tables using OSPF
 - Most CMU hosts cannot be reached at link layer

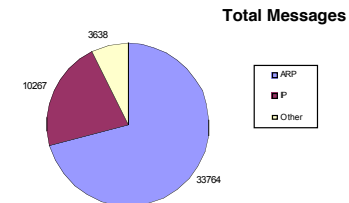
Proxy ARP



- Provides Link-Layer Connectivity Using IP Routing
 - Local router (gigrouter) sees ARP request
 - Uses IP addressing to locate host
 - Becomes "Proxy" for remote host
 - Using own MAC address
 - Requestor thinks that it is communicating directly with remote host

Monitoring Packet Traffic

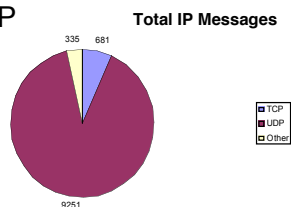
- Experiment
 - Ran windump for 15 minutes connected to CMU network
 - No applications running
 - But many background processes use network
 - Lots of ARP traffic (71% of total)
 - Average 37 ARP requests / second (why all from CS hosts?)
 - Only see responses from own machine (why?)



Monitoring Packet Traffic



- Other Traffic
 - Mostly UDP
 - Encode low-level protocols such as bootp
 - Nothing very exciting (why?)
- Answers for UDP and ARP
 - On a switched network you only see broadcast traffic or traffic sent to/from you
 - TCP is never sent broadcast



Some People Have Too Much Time...



- Everything I needed to know about networks I learned from ~~TV~~ Google video
 - [Ethernet collision animation](#)

AND.....

- Just to make sure...
 1. Packets really can't catch fire. That is not why we have insulation on wires
 2. Don't answer "what happens after a collision" on the exam/HW with "the packets catch on fire!"