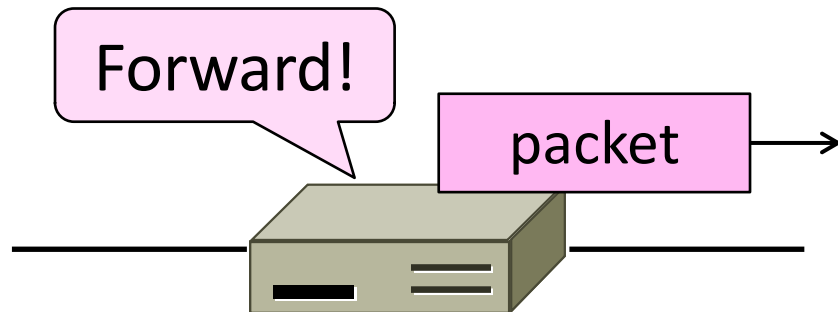


Topic

- How do routers forward packets?
 - We'll look at how IP does it
 - (We'll cover routing later)



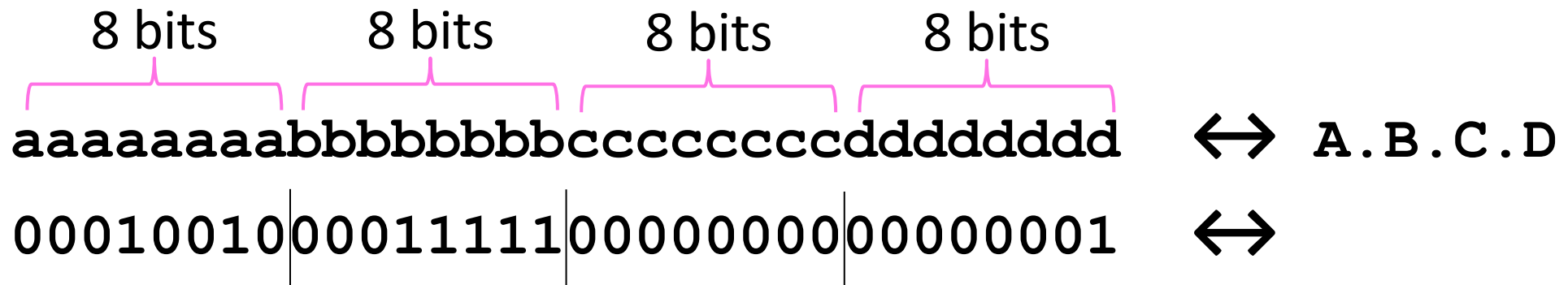
Recap

- We want the network layer to:
 - Scale to large networks
 - Using addresses with hierarchy
 - Support diverse technologies
 - Internetworking with IP
 - Use link bandwidth well
 - Lowest-cost routing
- This lecture
- More later
- Next time



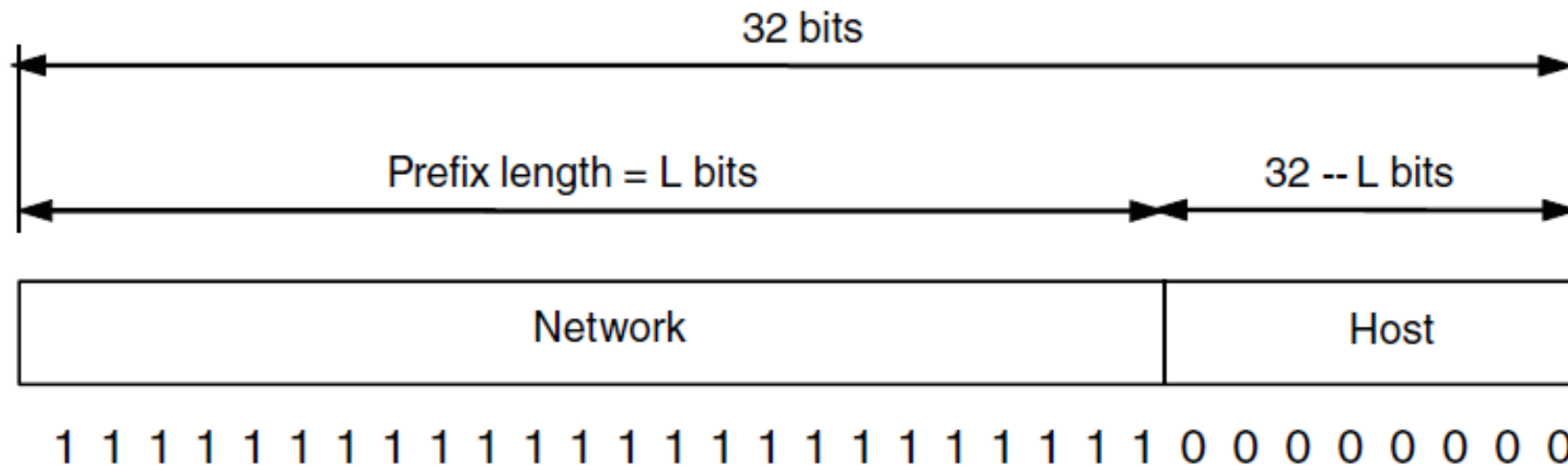
IP Addresses

- IPv4 uses 32-bit addresses
 - Later we'll see IPv6, which uses 128-bit addresses
- Written in “dotted quad” notation
 - Four 8-bit numbers separated by dots



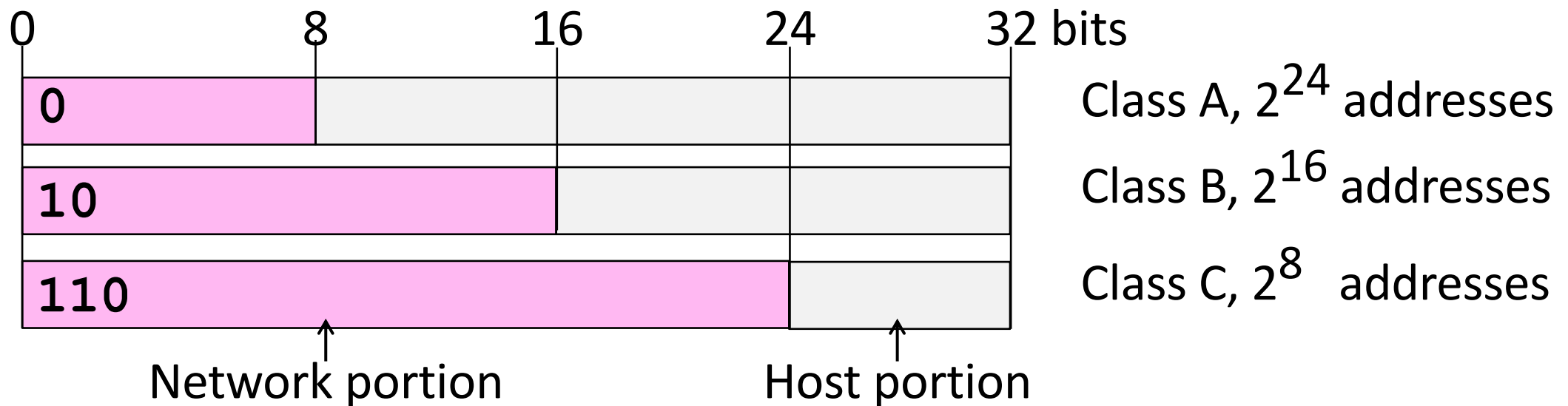
IP Prefixes

- Addresses are allocated in blocks called prefixes
 - Addresses in an L-bit prefix have the same top L bits
 - There are 2^{32-L} addresses aligned on 2^{32-L} boundary



Classful IP Addressing

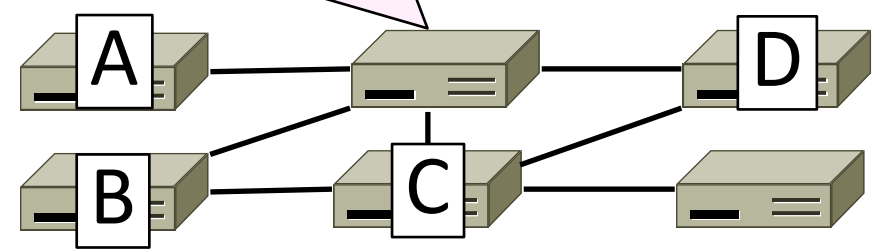
- Originally, IP addresses came in fixed size blocks with the class/size encoded in the high-order bits
 - They still do, but the classes are now ignored



IP Forwarding

- All addresses on one network belong to the same prefix
- Node uses a table that lists the next hop for prefixes

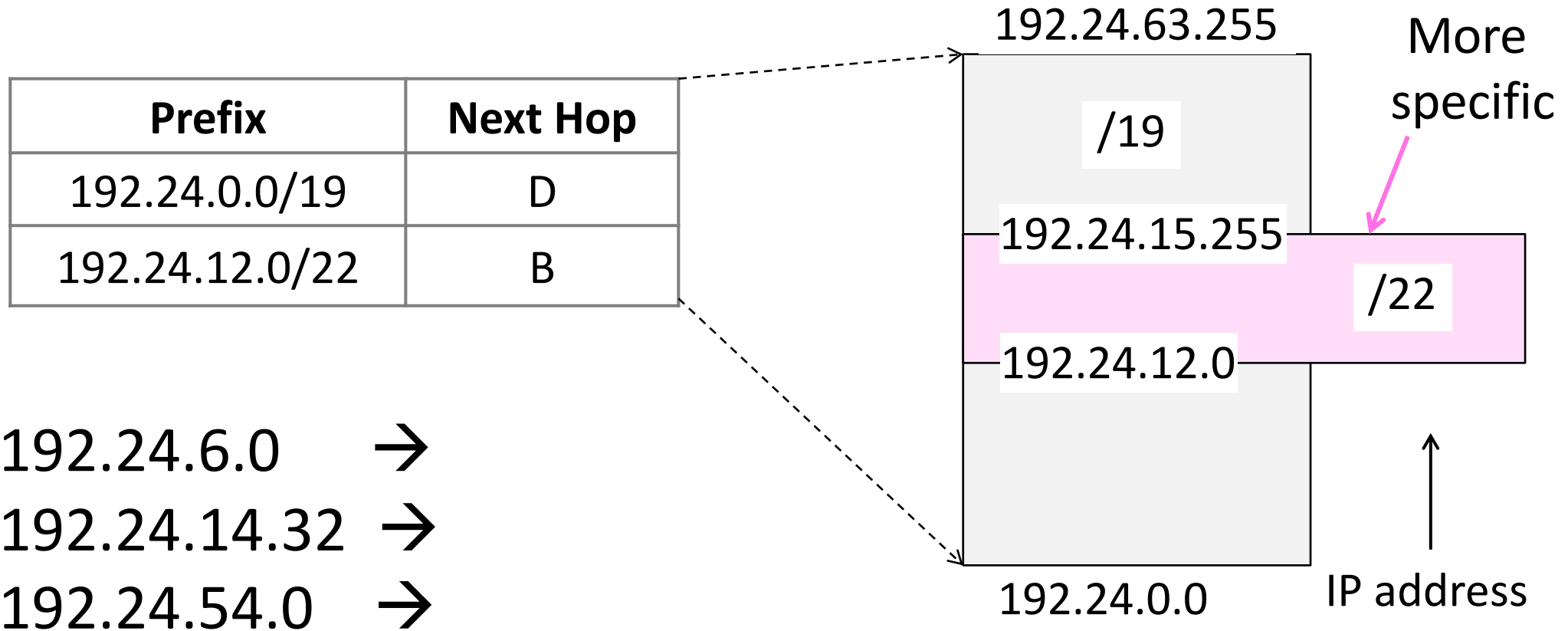
Prefix	Next Hop
192.24.0.0/19	D
192.24.12.0/22	B



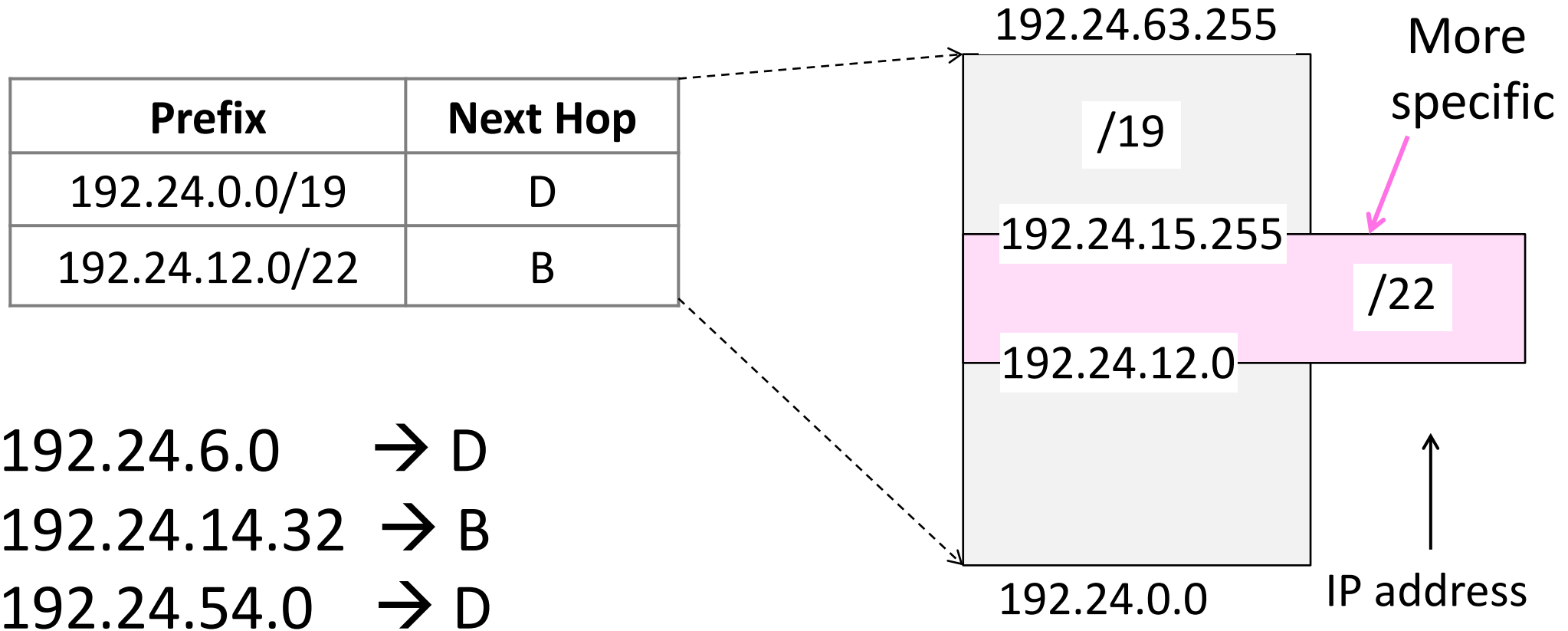
Longest Matching Prefix

- Prefixes in the table might overlap!
 - Combines hierarchy with flexibility
- Longest matching prefix forwarding rule:
 - For each packet, find the longest prefix that contains the destination address, i.e., the most specific entry
 - Forward the packet to the next hop router for that prefix

Longest Matching Prefix (2)

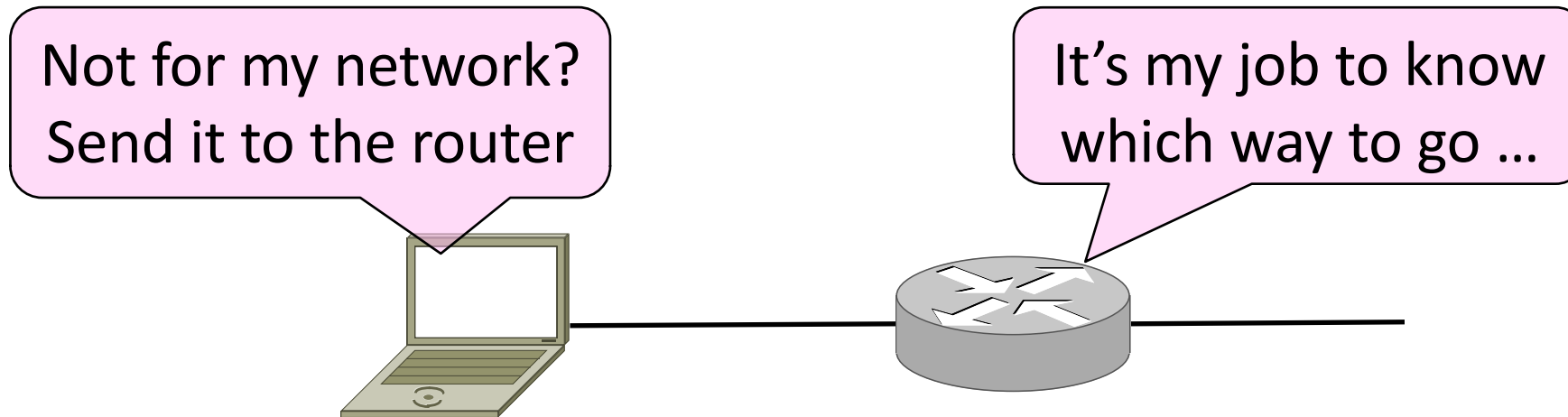


Longest Matching Prefix (2)



Host/Router Distinction

- In the Internet:
 - Routers do the routing, know which way to all destinations
 - Hosts send remote traffic (out of prefix) to nearest router



Host Forwarding Table

- Give using longest matching prefix
 - 0.0.0.0/0 is a default route that catches all IP addresses

Prefix	Next Hop
My network prefix	Send to that IP
0.0.0.0/0	Send to my router



Flexibility of Longest Matching Prefix

- Can provide default behavior, with less specifics
 - To send traffic going outside an organization to a border router
- Can special case behavior, with more specifics
 - For performance, economics, security, ...



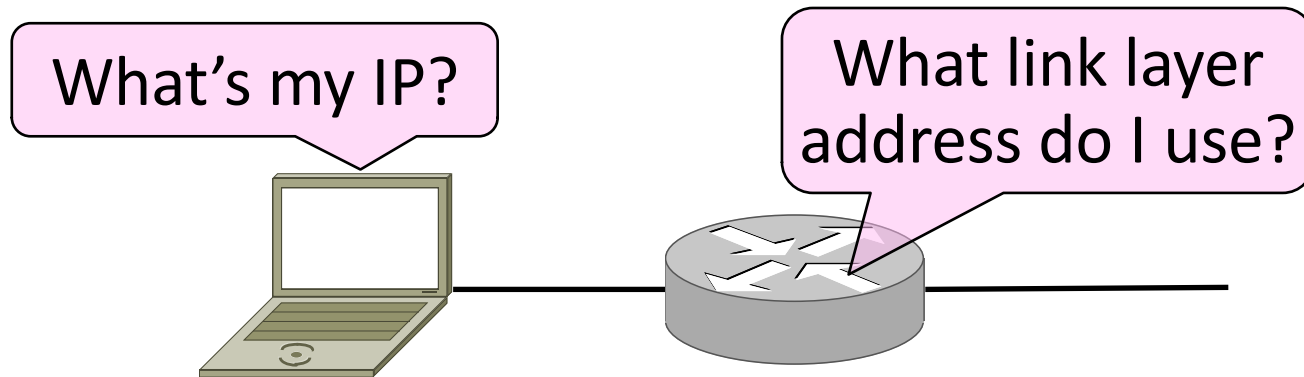
Performance of Longest Matching Prefix

- Uses hierarchy for a compact table
 - Relies on use of large prefixes
- Lookup more complex than table
 - Used to be a concern for fast routers
 - Not an issue in practice these days



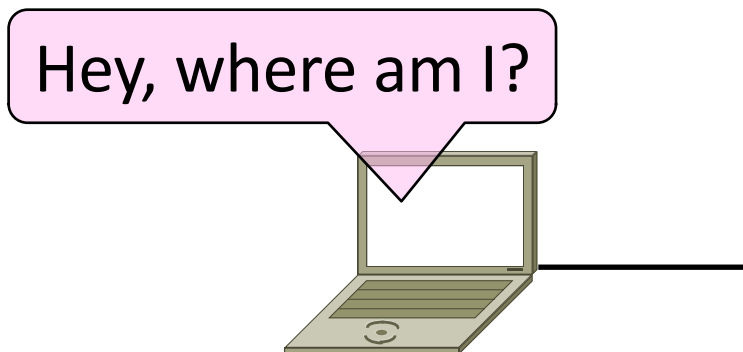
Topic

- Filling in the gaps we need to make for IP forwarding work in practice
 - Getting IP addresses (DHCP) »
 - Mapping IP to link addresses (ARP) »



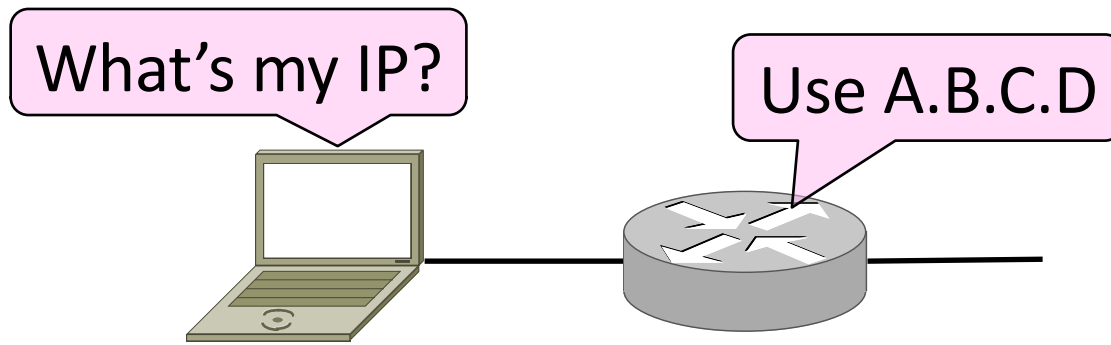
Getting IP Addresses

- Problem:
 - A node wakes up for the first time ...
 - What is its IP address? What's the IP address of its router? Etc.
 - At least Ethernet address is on NIC



Getting IP Addresses (2)

1. Manual configuration (old days)
 - Can't be factory set, depends on use
2. A protocol for automatically configuring addresses (DHCP) »
 - Shifts burden from users to IT folk



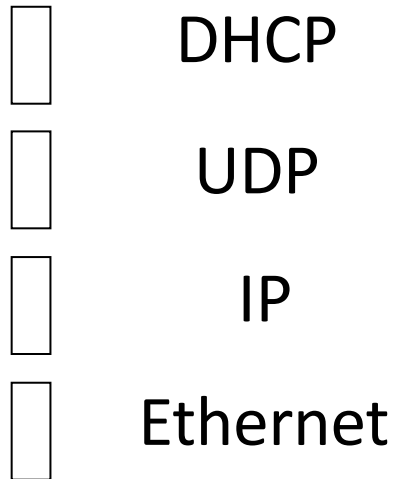
DHCP

- DHCP (Dynamic Host Configuration Protocol), from 1993, widely used
- It leases IP address to nodes
- Provides other parameters too
 - Network prefix
 - Address of local router
 - DNS server, time server, etc.



DHCP Protocol Stack

- DHCP is a client-server application
 - Uses UDP ports 67, 68

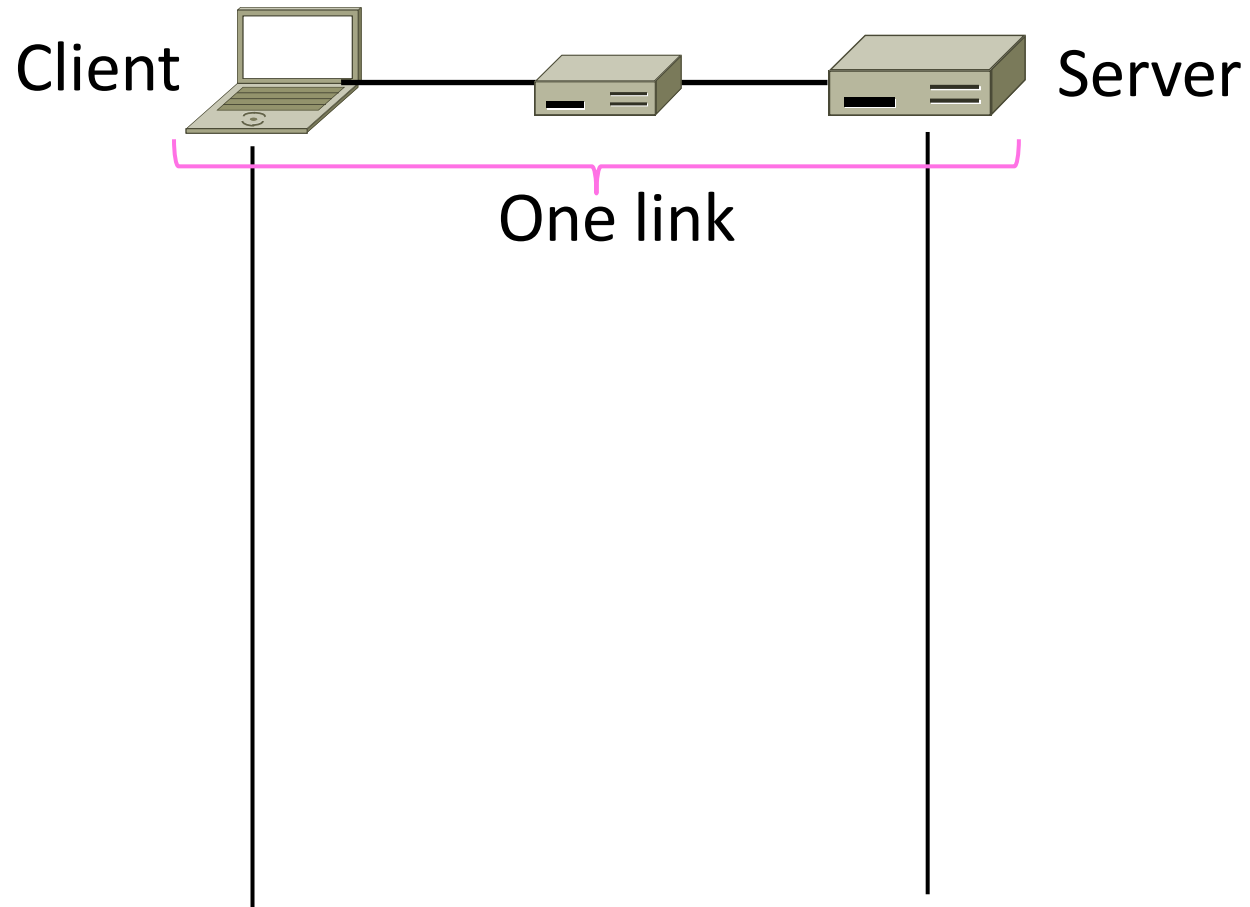


DHCP Addressing

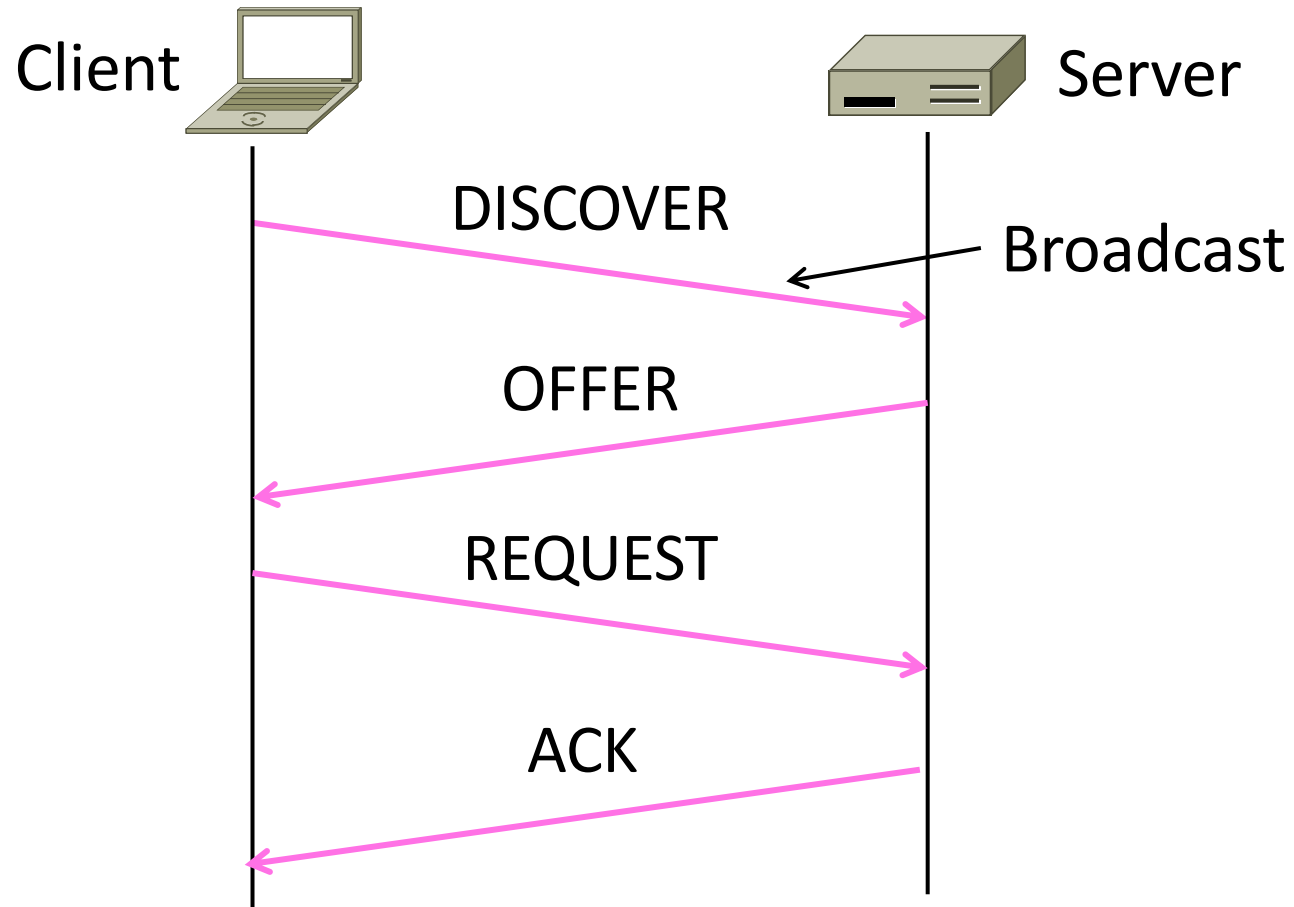
- Bootstrap issue:
 - How does node send a message to DHCP server before it is configured?
- Answer:
 - Node sends broadcast messages that delivered to all nodes on the network
 - Broadcast address is all 1s
 - IP (32 bit): 255.255.255.255
 - Ethernet (48 bit): ff:ff:ff:ff:ff:ff



DHCP Messages

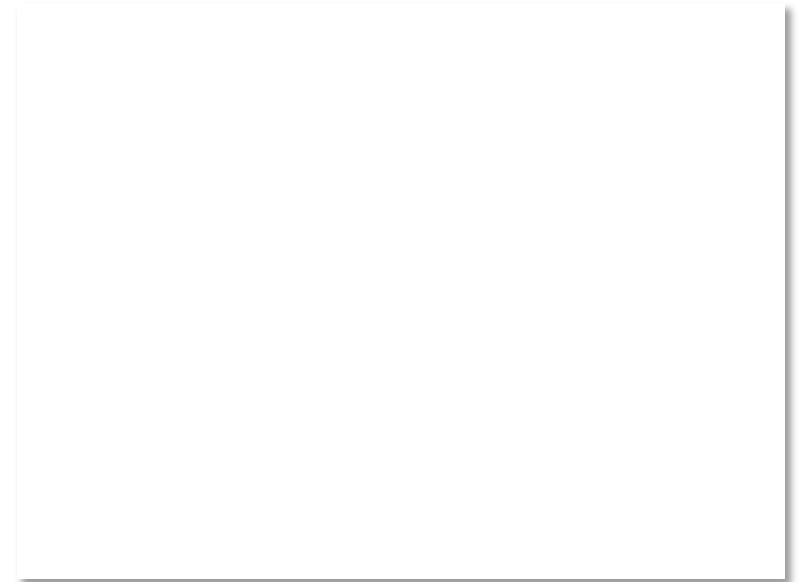


DHCP Messages (2)



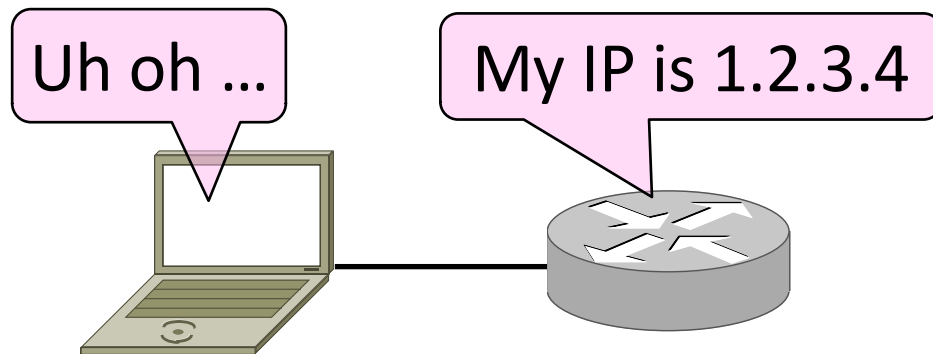
DHCP Messages (3)

- To renew an existing lease, an abbreviated sequence is used:
 - REQUEST, followed by ACK
- Protocol also supports replicated servers for reliability



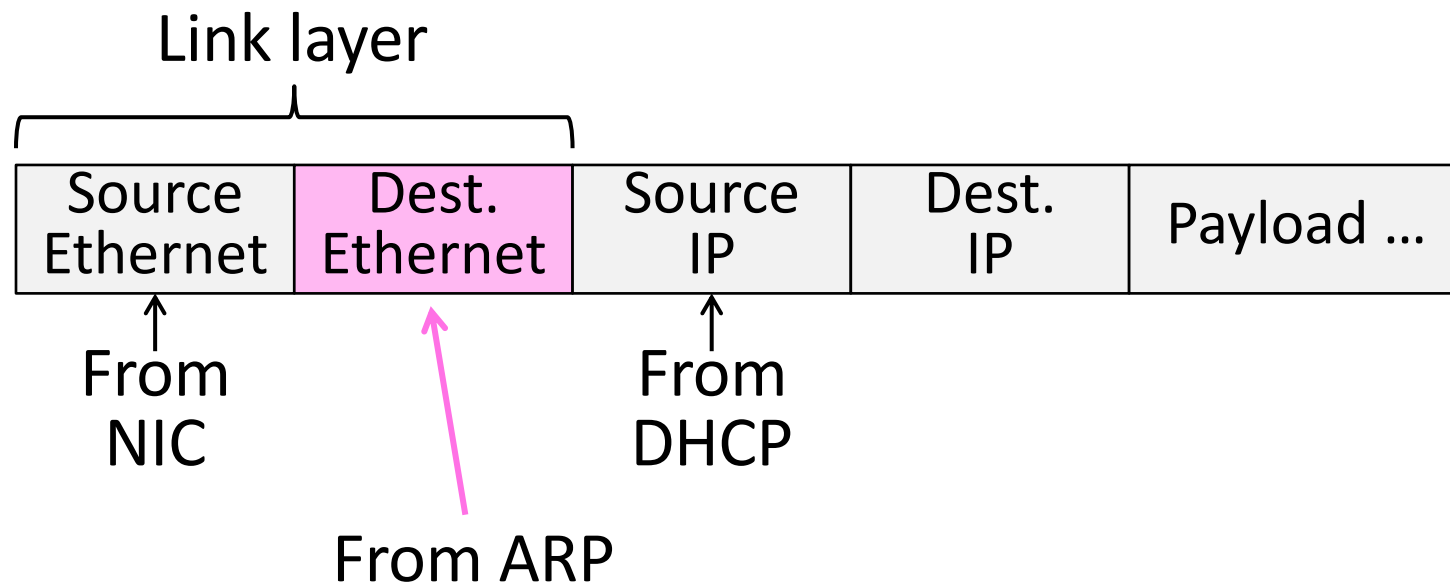
Sending an IP Packet

- Problem:
 - A node needs Link layer addresses to send a frame over the local link
 - How does it get the destination link address from a destination IP address?



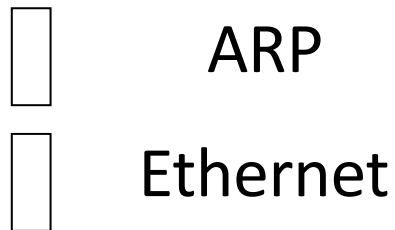
ARP (Address Resolution Protocol)

- Node uses to map a local IP address to its Link layer addresses

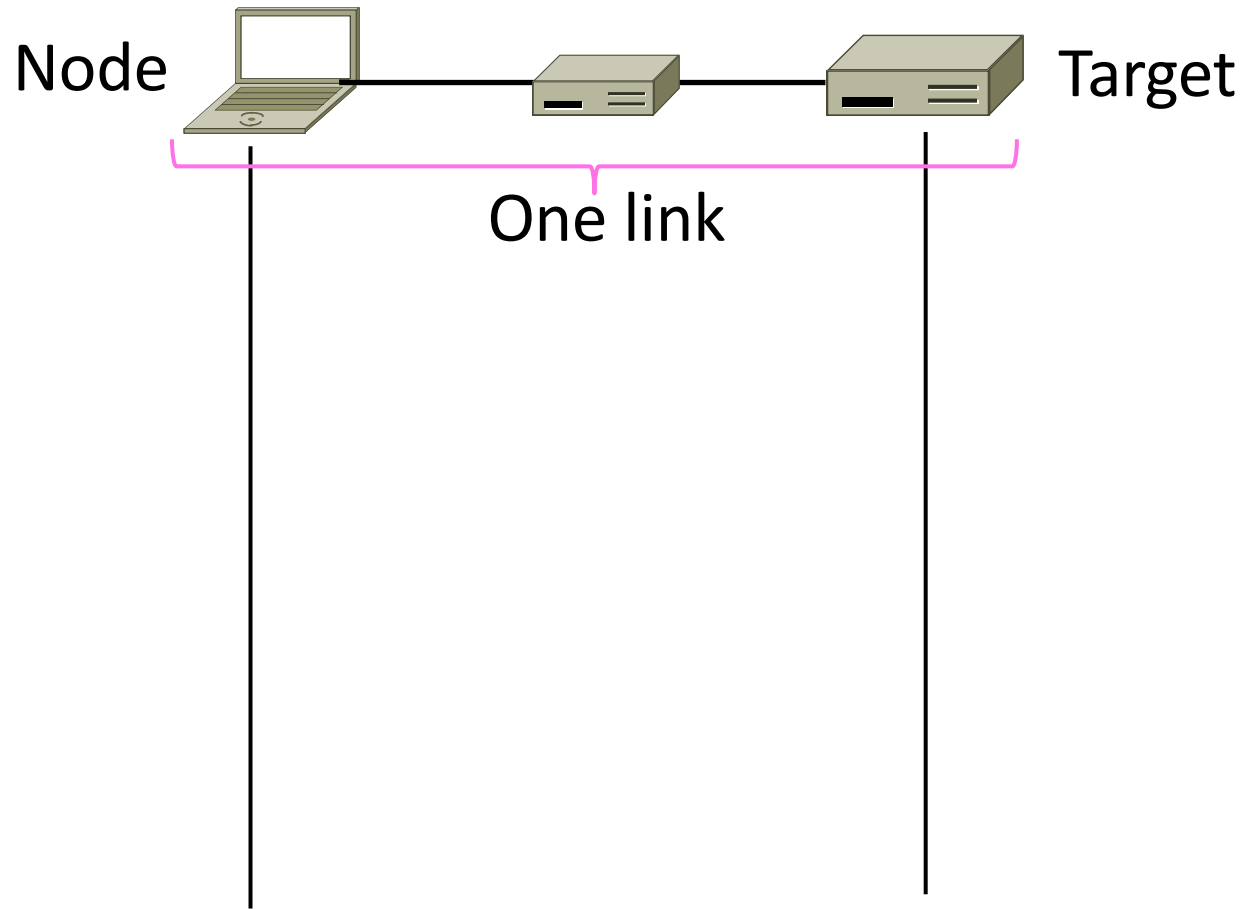


ARP Protocol Stack

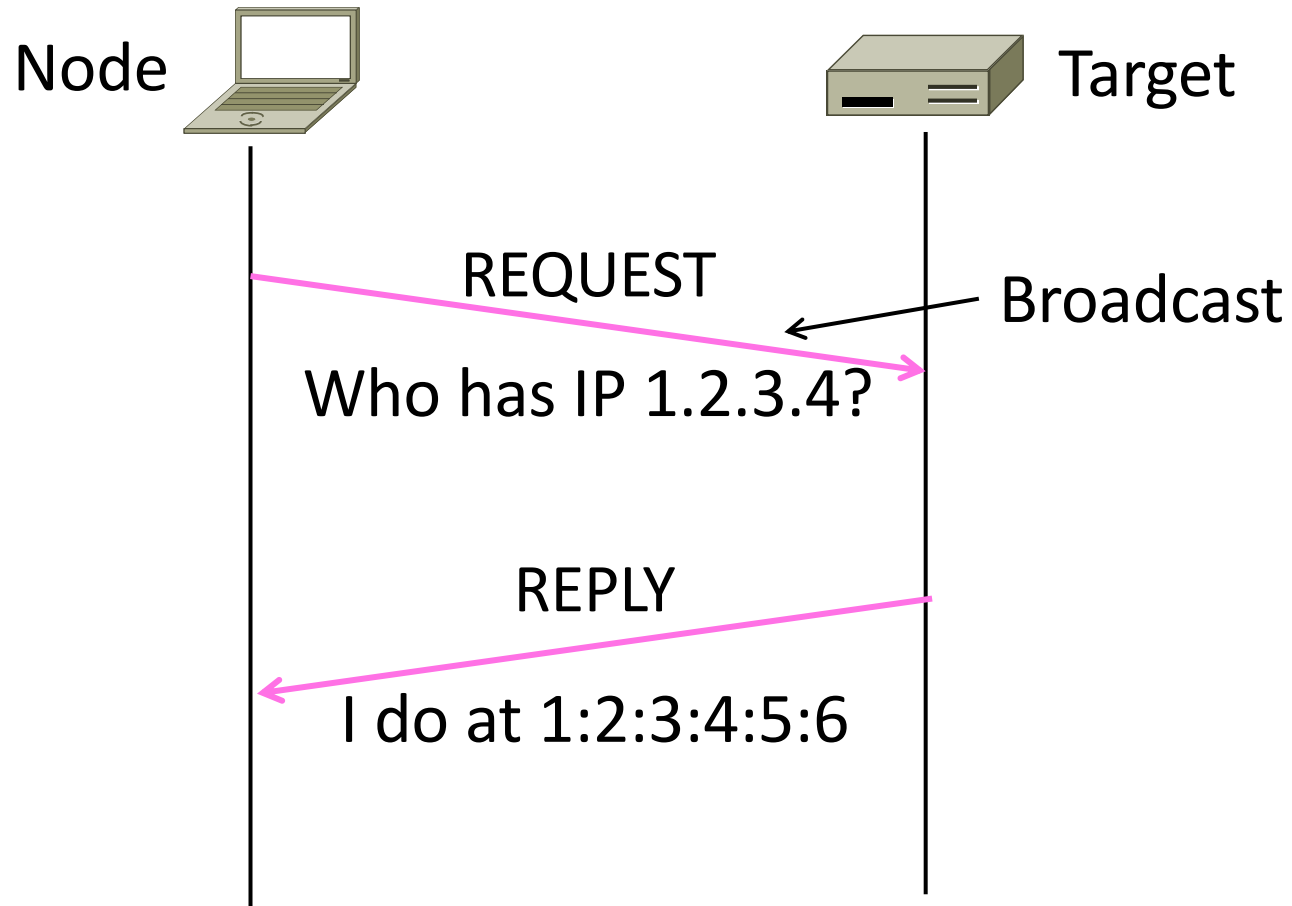
- ARP sits right on top of link layer
 - No servers, just asks node with target IP to identify itself
 - Uses broadcast to reach all nodes



ARP Messages



ARP Messages (2)



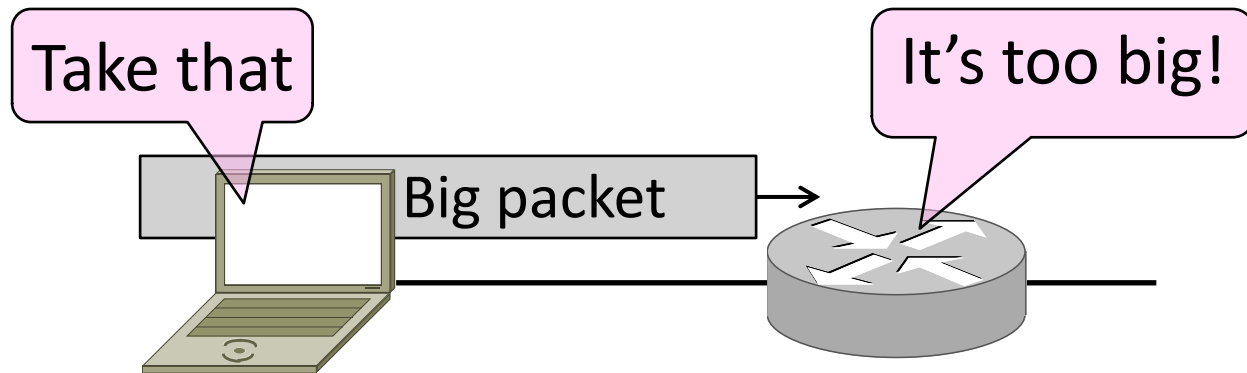
Discovery Protocols

- Help nodes find each other
 - There are more of them!
 - E.g., zeroconf, Bonjour
- Often involve broadcast
 - Since nodes aren't introduced
 - Very handy glue



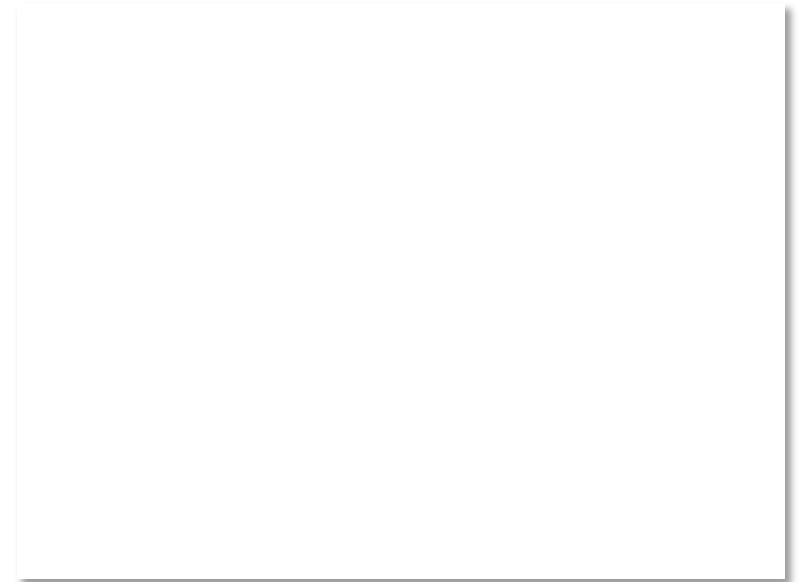
Topic

- How do we connect networks with different maximum packet sizes?
 - Need to split up packets, or discover the largest size to use



Packet Size Problem

- Different networks have different maximum packet sizes
 - Or MTU (Maximum Transmission Unit)
 - E.g., Ethernet 1.5K, WiFi 2.3K
- Prefer large packets for efficiency
 - But what size is too large?
 - Difficult because node does not know complete network path



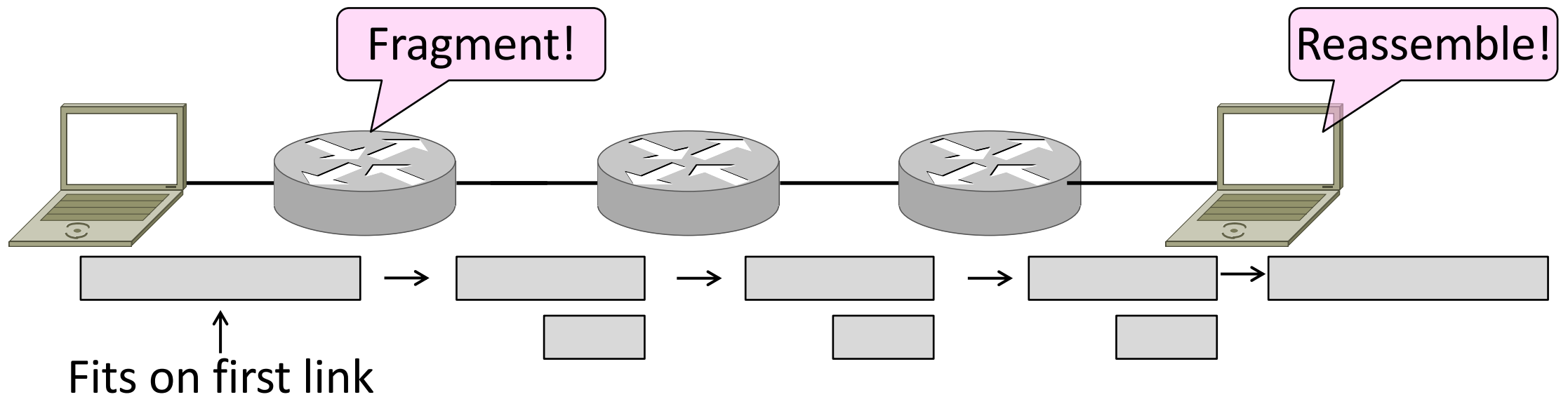
Packet Size Solutions

- Fragmentation (now)
 - Split up large packets in the network if they are too big to send
 - Classic method, dated
- Discovery (next)
 - Find the largest packet that fits on the network path and use it
 - IP uses today instead of fragmentation



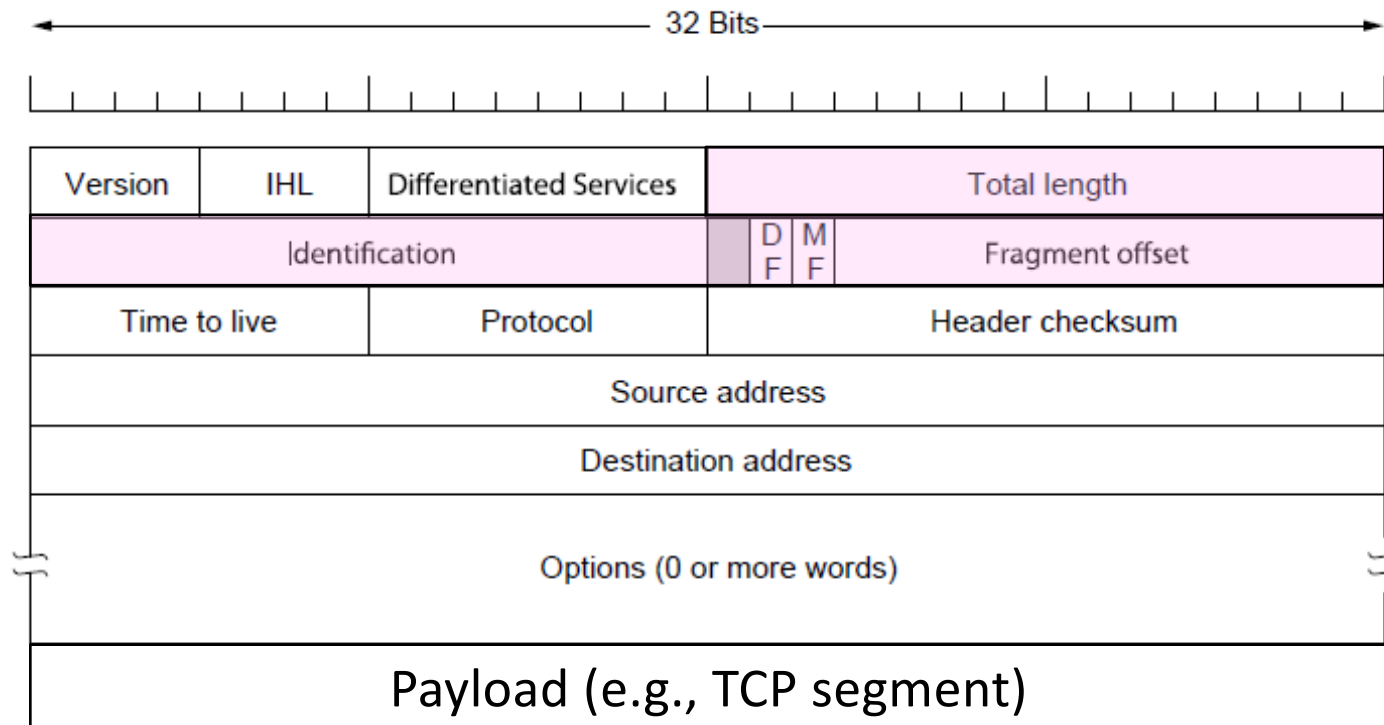
IPv4 Fragmentation

- Routers fragment packets that are too large to forward
- Receiving host reassembles to reduce load on routers



IPv4 Fragmentation Fields

- Header fields used to handle packet size differences
 - Identification, Fragment offset, MF/DF control bits



IPv4 Fragmentation Procedure

- Routers split a packet that is too large:
 - Typically break into large pieces
 - Copy IP header to pieces
 - Adjust length on pieces
 - Set offset to indicate position
 - Set MF (More Fragments) on all pieces except last
- Receiving hosts reassembles the pieces:
 - Identification field links pieces together, MF tells receiver when it has all pieces



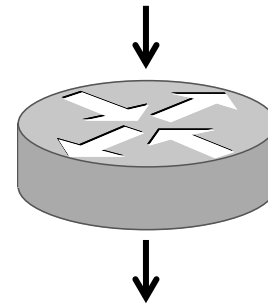
IPv4 Fragmentation (2)

Before
MTU = 2300

ID = 0x12ef
Data Len = 2300
Offset = 0
MF = 0



(Ignore length
of headers)



After
MTU = 1500

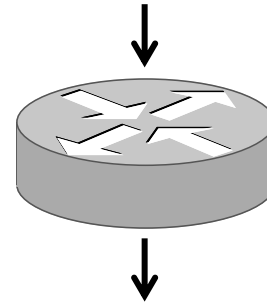
ID =
Data Len =
Offset =
MF =

ID =
Data Len =
Offset =
MF =

IPv4 Fragmentation (3)

Before
MTU = 2300

ID = 0x12ef
Data Len = 2300
Offset = 0
MF = 0



After
MTU = 1500

ID = 0x12ef
Data Len = 1500
Offset = 0
MF = 1



ID = 0x12ef
Data Len = 800
Offset = 1500
MF = 0



IPv4 Fragmentation (4)

- It works!
 - Allows repeated fragmentation
- But fragmentation is undesirable
 - More work for routers, hosts
 - Tends to magnify loss rate
 - Security vulnerabilities too

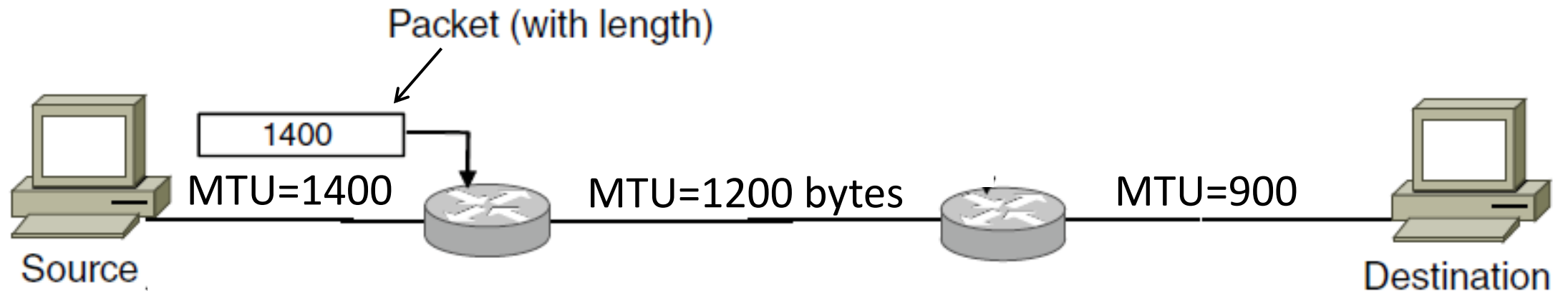


Path MTU Discovery

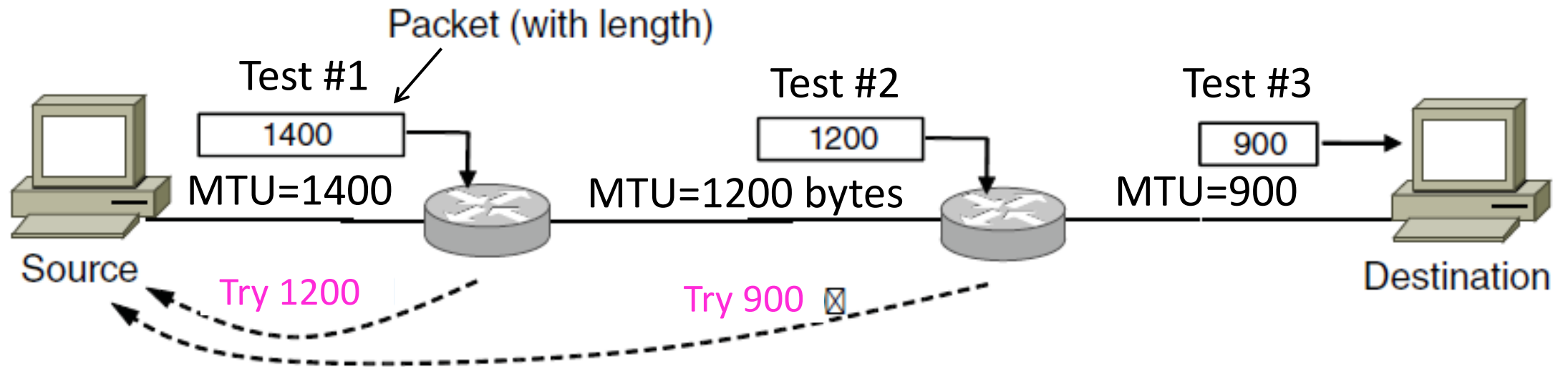
- Discover the MTU that will fit
 - So we can avoid fragmentation
 - The method in use today
- Host tests path with large packet
 - Routers provide feedback if too large; they tell host what size would have fit



Path MTU Discovery (2)



Path MTU Discovery (3)



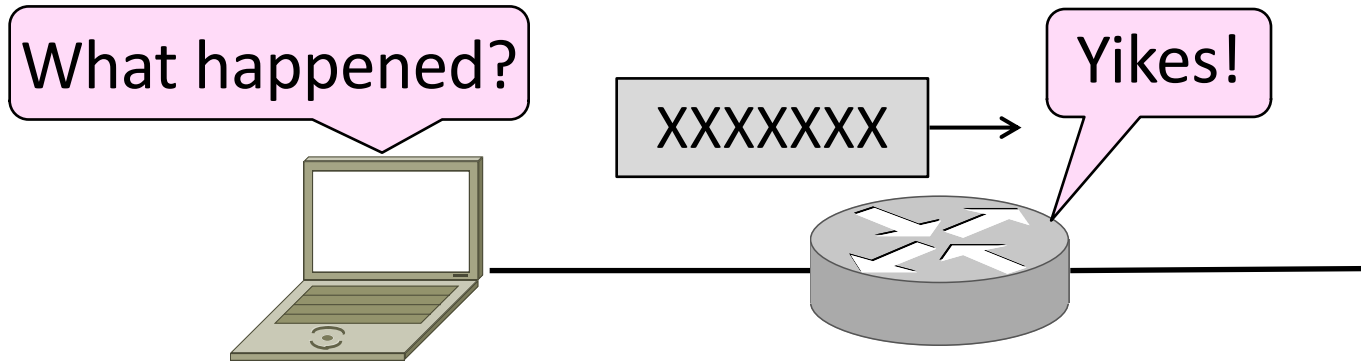
Path MTU Discovery (4)

- Process may seem involved
 - But usually quick to find right size
- Path MTU depends on the path and so can change over time
 - Search is ongoing
- Implemented with ICMP (next)
 - Set DF (Don't Fragment) bit in IP header to get feedback messages



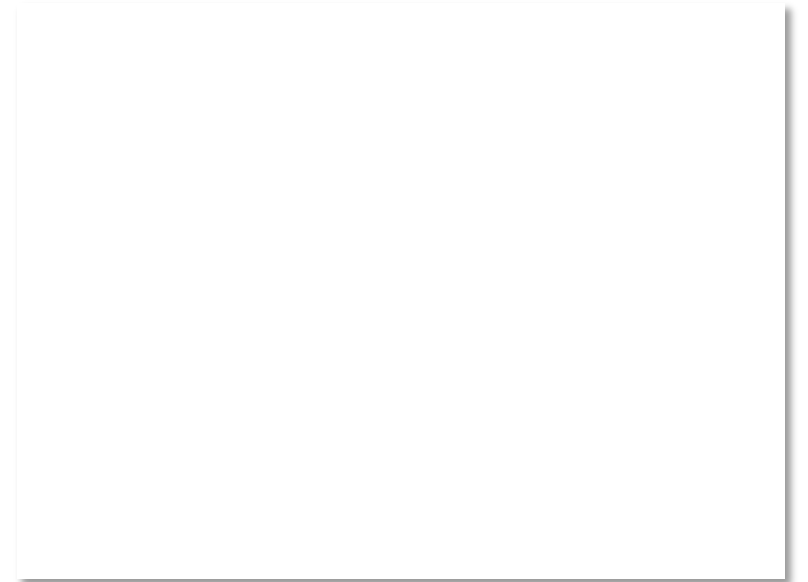
Topic

- What happens when something goes wrong during forwarding?
 - Need to be able to find the problem



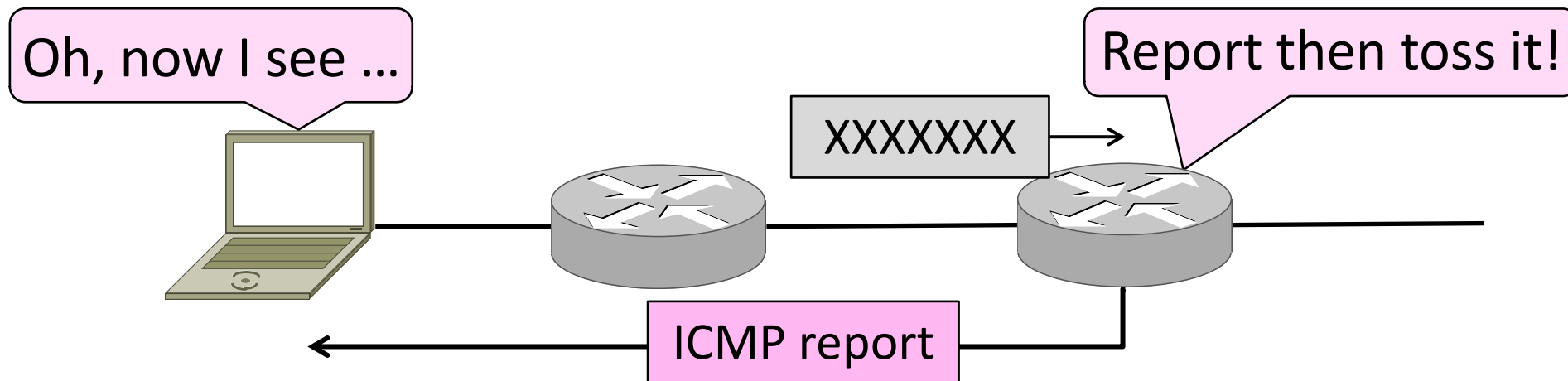
Internet Control Message Protocol

- ICMP is a companion protocol to IP
 - They are implemented together
 - Sits on top of IP (IP Protocol=1)
- Provides error report and testing
 - Error is at router while forwarding
 - Also testing that hosts can use



ICMP Errors

- When router encounters an error while forwarding:
 - It sends an ICMP error report back to the IP source address
 - It discards the problematic packet; host needs to rectify

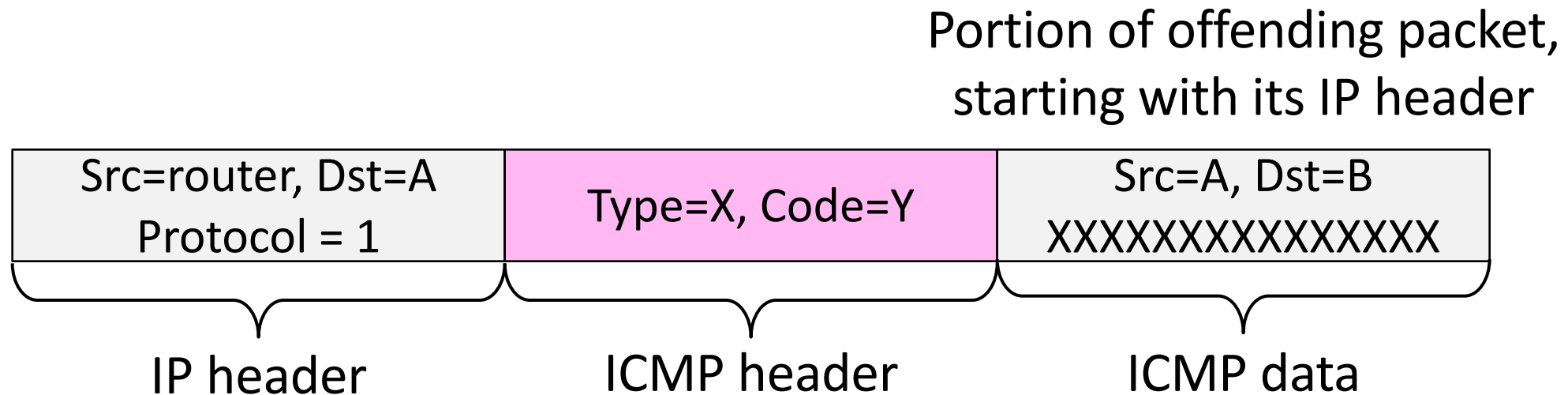


ICMP Message Format

- Each ICMP message has a Type, Code, and Checksum
- Often carry the start of the offending packet as payload
- Each message is carried in an IP packet

ICMP Message Format (2)

- Each ICMP message has a Type, Code, and Checksum
- Often carry the start of the offending packet as payload
- Each message is carried in an IP packet

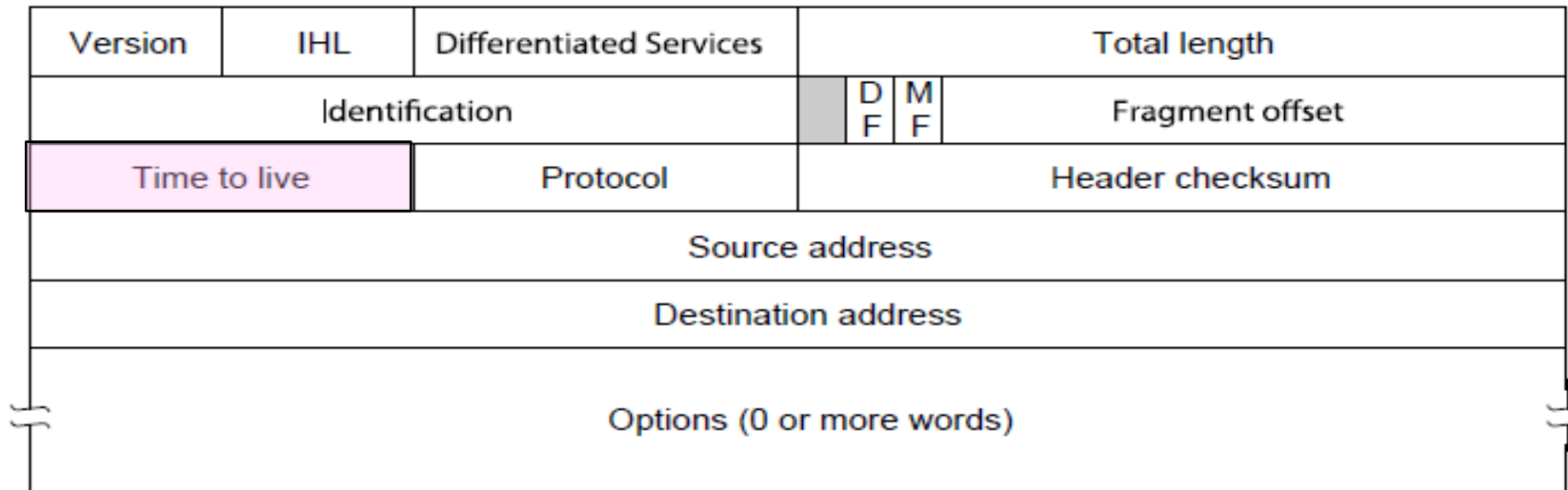


Example ICMP Messages

Name	Type / Code	Usage
Dest. Unreachable (Net or Host)	3 / 0 or 1	Lack of connectivity
Dest. Unreachable (Fragment)	3 / 4	Path MTU Discovery
Time Exceeded (Transit)	11 / 0	Testing, not a forwarding error: Host sends Echo Request, and destination responds with an Echo Reply Traceroute
Echo Request or Reply	8 or 0 / 0	Ping

Traceroute

- IP header contains TTL (Time to live) field
 - Decrement every router hop, with ICMP error if it hits zero
 - Protects against forwarding loops



Traceroute (2)

- Traceroute repurposes TTL and ICMP functionality
 - Sends probe packets increasing TTL starting from 1
 - ICMP errors identify routers on the path

