# Operating Systems and Networks
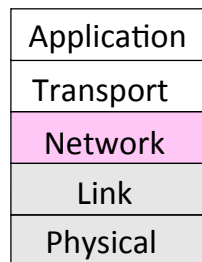
# Network Lecture 5: Network Layer 1

Adrian Perrig
Network Security Group
ETH Zürich

# Pending Issues

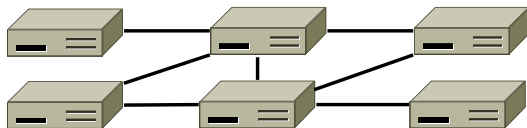- Ethernet performance? See Section 4.3.3 in book. For reasonable parameters, ~85% efficiency.

2

# Where we are in the Course

- Starting the Network Layer!
  - Builds on the link layer. <u>Routers</u> send <u>packets</u> over multiple networks
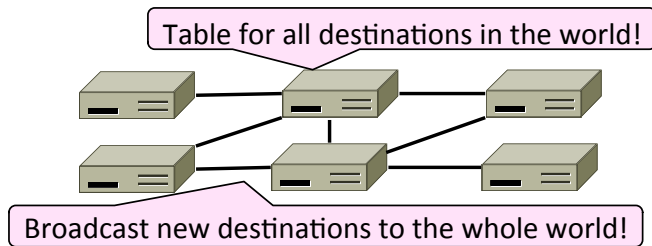
| Application |
| :---: |
| Transport |
| Network |
| Link |
| Physical |

Computer Networks

3

# Why do we need a Network layer?

- We can already build networks  with links and switches and send frames between hosts …

Computer Networks

4

# Shortcomings of Switches

1. Don't scale to large networks
   - Blow up of routing table, broadcast

Table for all destinations in the world!
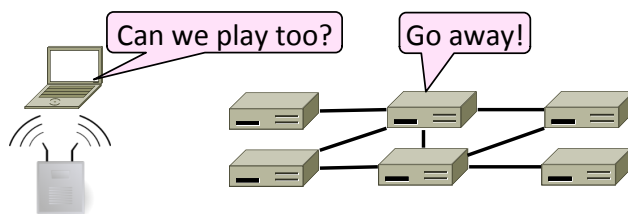
Broadcast new destinations to the whole world!

Computer Networks

5

# Shortcomings of Switches (2)

2. Don't work across more than one link layer technology
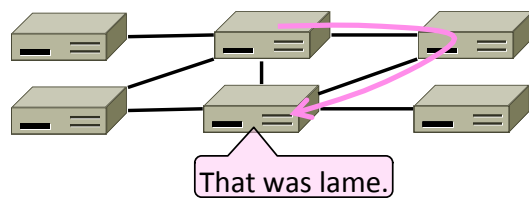   - Hosts on Ethernet + 3G + 802.11  …

Can we play too?

Go away!

Computer Networks

6

# Shortcomings of Switches (3)

3. Don't give much traffic control
   - Want to plan routes / bandwidth



That was lame.

Computer Networks

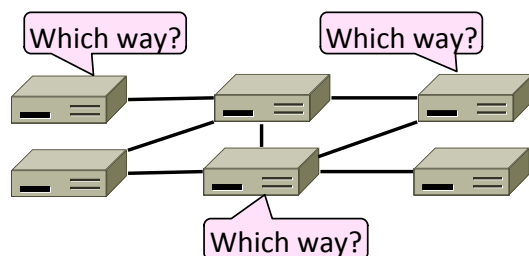# Network Layer Approach

- Scaling:
  - Hierarchy, in the form of prefixes

- Heterogeneity:
  - IP for internetworking

- Bandwidth Control:
  - Lowest-cost routing
  - Later QOS (Quality of Service)

Computer Networks

# Topics

- Network service models
  - Datagrams (packets), virtual circuits
- IP (Internet Protocol)
  - Internetworking
  - Forwarding (Longest Matching Prefix)
  - Helpers: ARP and DHCP
  - Fragmentation and MTU discovery
  - Errors: ICMP (traceroute!)

  This time

- IPv6, the future of IP
- NAT, a "middlebox"
- Routing algorithms

  Next time

Computer Networks

9

---

# Routing vs. Forwarding

- <u>Routing</u> is the process of deciding in which direction to send traffic
  - Network wide (global) and expensive



Which way?   Which way?

Which way?

Computer Networks

10

---

# Routing vs. Forwarding (2)

- <u>Forwarding</u> is the process of sending a packet on its way
  - Node process (local) and fast

# Our Plan

- <u>Forwarding</u> this time
  - What routers do with packets

- <u>Routing</u> next time
  - Logically this comes first
  - But ignore it for now

# Network Services (§5.1)

- What kind of service does the Network layer provide to the Transport layer?
  - How is it implemented at routers?

Service? What's he talking about?

Computer Networks 13

# Two Network Service Models

- Datagrams, or connectionless service
  - Like postal letters
  - (This one is IP)

- Virtual circuits, or connection-oriented service
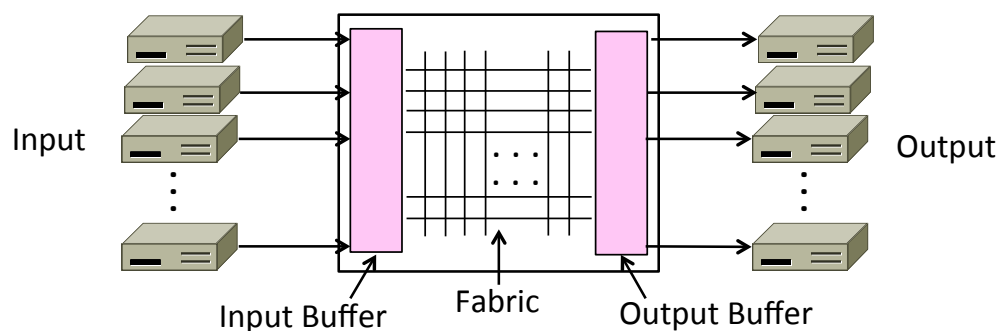  - Like a telephone call

Computer Networks 14

# Store-and-Forward Packet Switching

- Both models are implemented with <u>store-and-forward packet switching</u>
  - Routers receive a complete packet, storing it temporarily if necessary before forwarding it onwards
  - We use statistical multiplexing to share link bandwidth over time

Computer Networks                                                        15

# Store-and-Forward (2)

- Switching element has internal buffering for contention



Input

Output

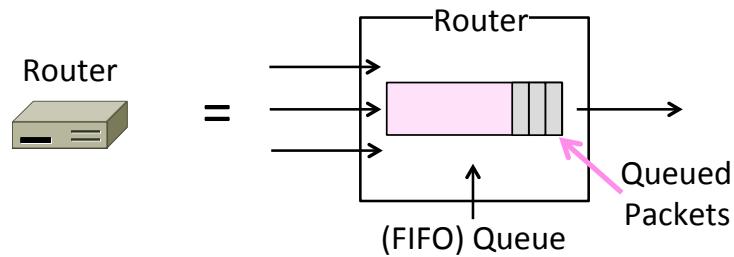Input Buffer          Fabric          Output Buffer

Computer Networks                                                        16

# Store-and-Forward (3)

- Simplified view with per-port output buffering
  - Buffer is typically a FIFO (First In First Out) queue
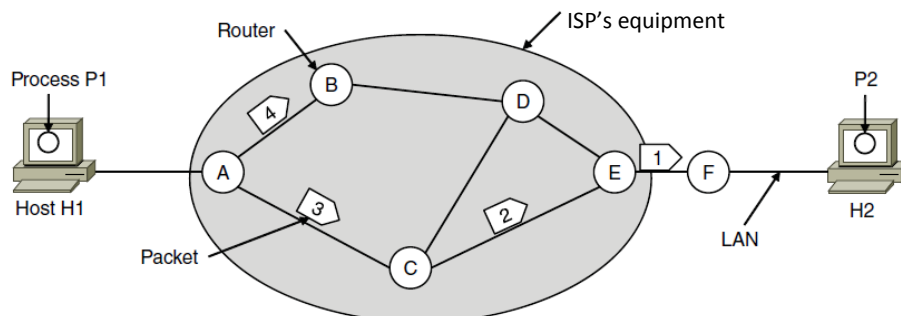  - If full, packets are discarded (congestion, later)



Computer Networks

17

# Datagram Model

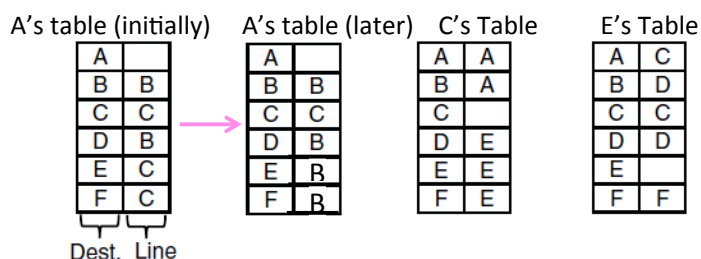- Packets contain a destination address; each router uses it to forward each packet, possibly on different paths
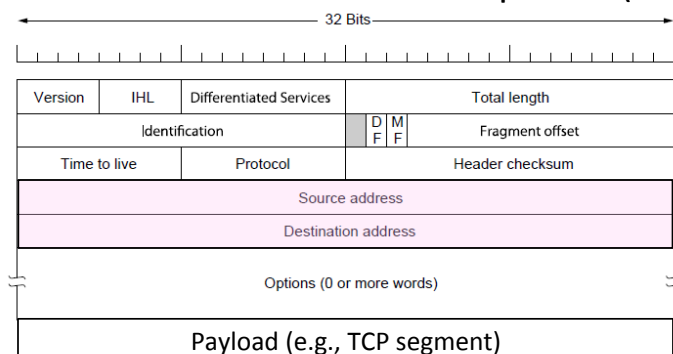


Computer Networks

18

9

# Datagram Model (2)

- Each router has a forwarding table keyed by address
  - Gives next hop for each destination address; may change

A's table (initially)  A's table (later)  C's Table  E's Table

| A |   |
|---|---|
| B | B |
| C | C |
| D | B |
| E | C |
| F | C |

| A |   |
|---|---|
| B | B |
| C | C |
| D | B |
| E | B |
| F | B |

| A | A |
|---|---|
| B | A |
| C |   |
| D | E |
| E | E |
| F | E |

| A | C |
|---|---|
| B | D |
| C | C |
| D | D |
| E |   |
| F | F |

Dest. Line

Computer Networks                                                          19

---

# IP (Internet Protocol)

- Network layer of the Internet, uses datagrams (next)
  - IPv4 carries 32 bit addresses on each packet (often 1.5 KB)

32 Bits

| Version | IHL | Differentiated Services | | | Total length |
|---------|-----|-------------------------|---|---|--------------|
| Identification | | | D F | M F | Fragment offset |
| Time to live | | Protocol | | | Header checksum |
| Source address | | | | | |
| Destination address | | | | | |
| Options (0 or more words) | | | | | |

Payload (e.g., TCP segment)

Computer Networks                                                          20
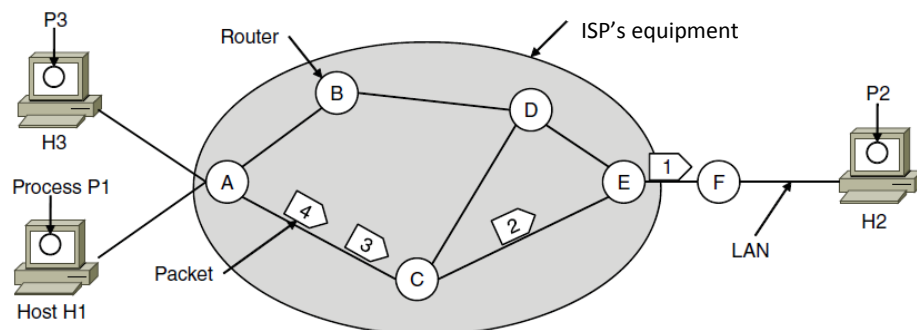
10

# Virtual Circuit Model

- Three phases:
    1. Connection establishment, circuit is set up
        - Path is chosen, circuit information stored in routers
    2. Data transfer, circuit is used
        - Packets are forwarded along the path
    3. Connection teardown, circuit is deleted
        - Circuit information is removed from routers

- Just like a telephone circuit, but virtual in the sense that no bandwidth need be reserved; statistical sharing of links
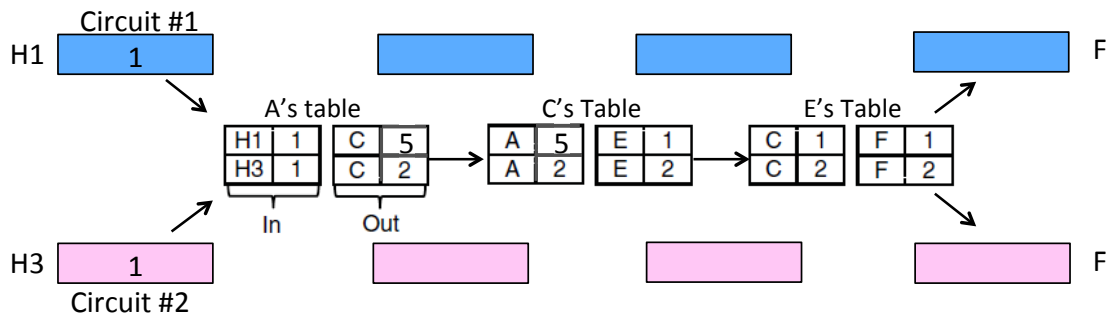
Computer Networks

21

# Virtual Circuits (2)

- Packets only contain a short label to identify the circuit
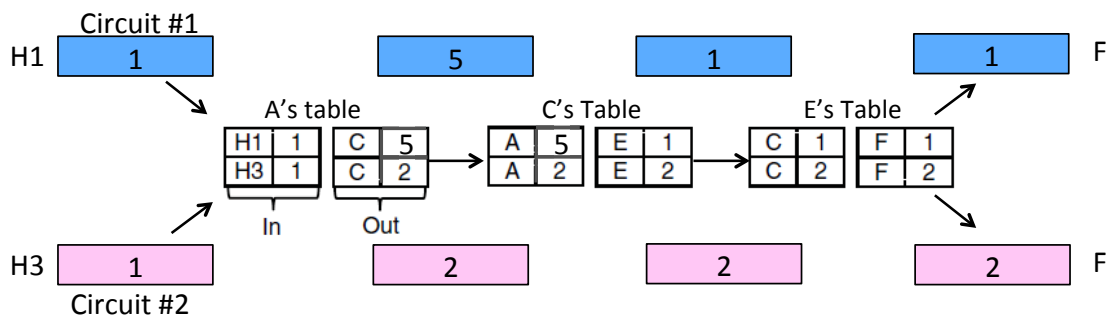    - Labels don't have any global meaning, only unique for a link



Computer Networks

22

# Virtual Circuits (3)

- Each router has a forwarding table keyed by circuit
  - Gives output line and next label to place on packet

Circuit #1

H1 [ 1 ]        [        ]        [        ]        [        ] F

A's table
| H1 | 1 |   | C | 5 |
| H3 | 1 |   | C | 2 |
In  Out

C's Table
| A | 5 |   | E | 1 |
| A | 2 |   | E | 2 |

E's Table
| C | 1 |   | F | 1 |
| C | 2 |   | F | 2 |

H3 [ 1 ]        [        ]        [        ]        [        ] F
Circuit #2

Computer Networks                                    23
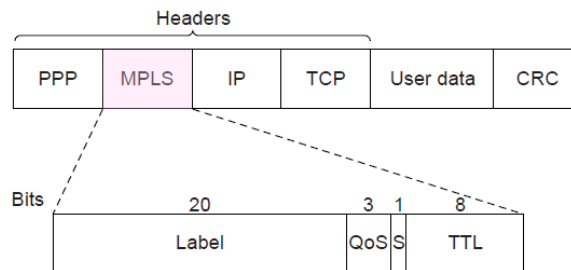
# Virtual Circuits (4)

- Each router has a forwarding table keyed by circuit
  - Gives output line and next label to place on packet

Circuit #1

H1 [ 1 ]        [ 5 ]        [ 1 ]        [ 1 ] F

A's table
| H1 | 1 |   | C | 5 |
| H3 | 1 |   | C | 2 |
In  Out

C's Table
| A | 5 |   | E | 1 |
| A | 2 |   | E | 2 |

E's Table
| C | 1 |   | F | 1 |
| C | 2 |   | F | 2 |

H3 [ 1 ]        [ 2 ]        [ 2 ]        [ 2 ] F
Circuit #2

Computer Networks                                    24

12

# MPLS (Multi-Protocol Label Switching, §5.6.5)

- A virtual-circuit like technology widely used by ISPs
  - ISP sets up circuits inside their backbone ahead of time
  - ISP adds MPLS label to IP packet at ingress, undoes at egress
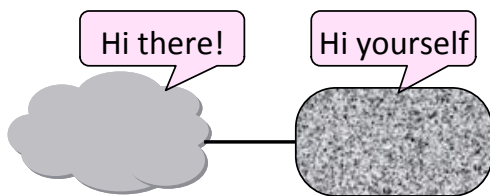


Computer Networks

# Datagrams vs Virtual Circuits

- Complementary strengths

| Issue | Datagrams | Virtual Circuits |
|---|---|---|
| Setup phase | Not needed | Required |
| Router state | Per destination | Per connection |
| Addresses | Packet carries full address | Packet carries short label |
| Routing | Per packet | Per circuit |
| Failures | Easier to mask | Difficult to mask |
| Quality of service | Difficult to add | Easier to add |

Computer Networks

# Internetworking (§5.5, 5.6.1)

- How do we connect different networks together?
  - This is called <u>internetworking</u>
  - We'll look at how IP does it



Computer Networks
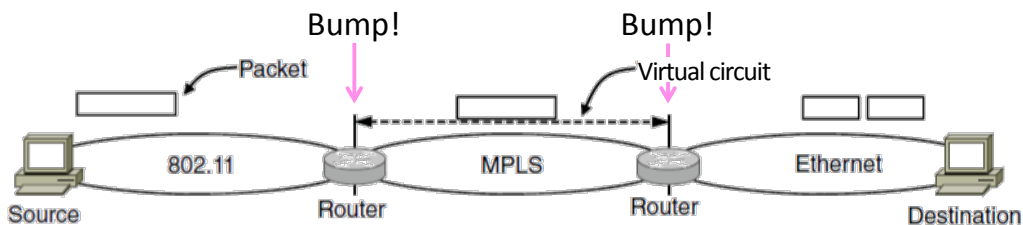
27

# How Networks May Differ

- Basically, in a lot of ways:
  - Service model (datagrams, VCs)
  - Addressing (what kind)
  - QOS (priorities, no priorities)
  - Packet sizes
  - Security (whether encrypted)

- Internetworking hides the differences with a common protocol. (Uh oh.)

Computer Networks

28

# Connecting Datagram and VC networks

- An example to show that it's not so easy
  - Need to map destination address to a VC and vice-versa
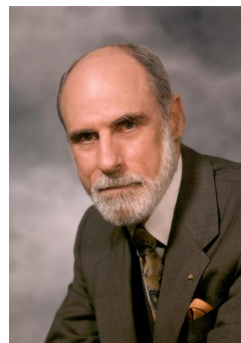  - A bit of a "road bump", e.g., might have to set up a VC

# Internetworking – Cerf and Kahn

- Pioneered by Cerf and Kahn, the "fathers of the Internet"
  - In 1974, later led to TCP/IP

- Tackled the problems of interconnecting networks
  - Instead of mandating a single network technology

Vint Cerf          Bob Kahn
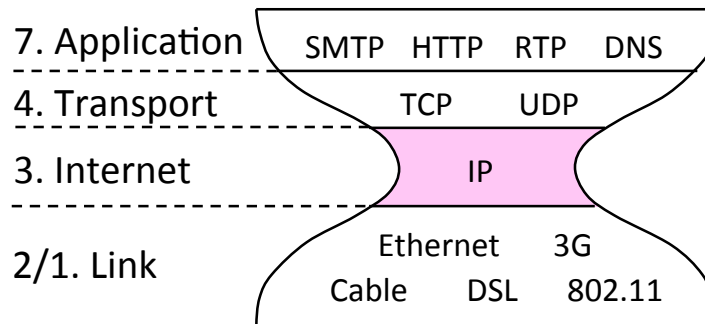


© 2009 IEEE          © 2009 IEEE

# Internet Reference Model

- IP is the "narrow waist" of the Internet
  - Supports many different links below and apps above

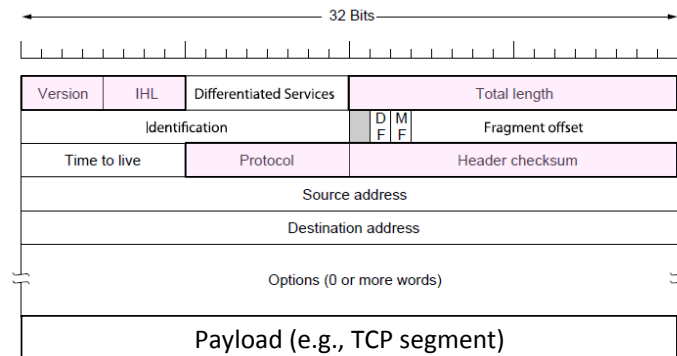| 7. Application | SMTP   HTTP   RTP   DNS |
| 4. Transport | TCP       UDP |
| 3. Internet | IP |
| 2/1. Link | Ethernet      3G<br>Cable      DSL      802.11 |

Computer Networks

31

# IP as a Lowest Common Denominator

- Suppose only some networks support QOS or security etc.
  - Difficult for internetwork to support

- Pushes IP to be a "lowest common denominator" protocol
  - Asks little of lower-layer networks
  - Gives little as a higher layer service

Computer Networks

32

16

# IPv4 (Internet Protocol)

- Various fields to meet straightforward needs
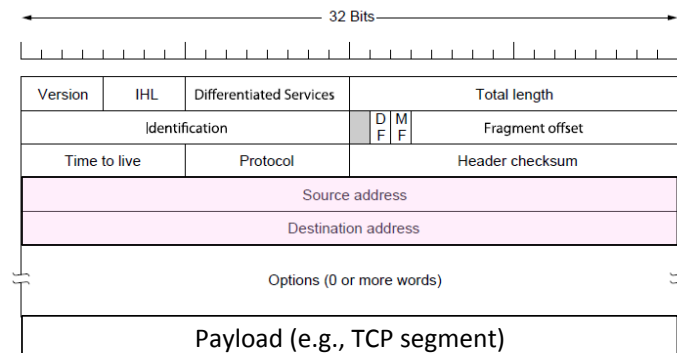  - Version, Header (IHL) and Total length, Protocol, and Header Checksum



Computer Networks

33

# IPv4 (2)

- Network layer of the Internet, uses datagrams
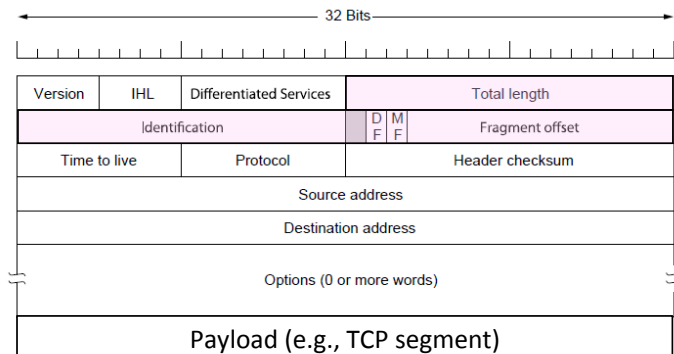  - Provides a layer of addressing above link addresses (next)



Computer Networks

34

# IPv4 (3)

- Some fields to handle packet size differences (later)
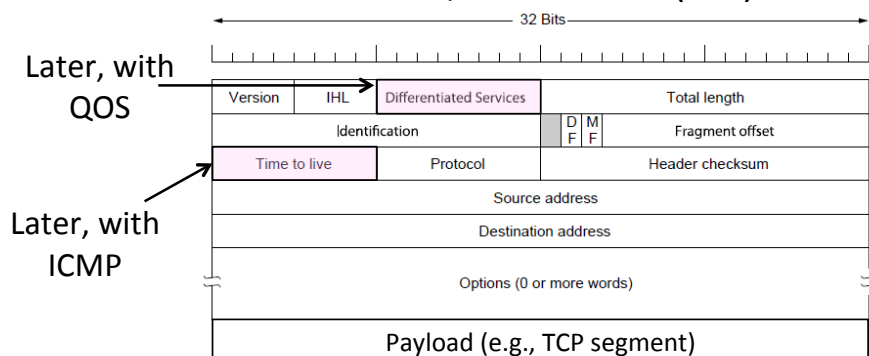  - Identification, Fragment offset, Fragment control bits

| 32 Bits | | | | |
|---|---|---|---|---|
| Version | IHL | Differentiated Services | | Total length |
| Identification | | | D F M F | Fragment offset |
| Time to live | | Protocol | | Header checksum |
| Source address | | | | |
| Destination address | | | | |
| Options (0 or more words) | | | | |
| Payload (e.g., TCP segment) | | | | |

Computer Networks

35

# IPv4 (4)

- Other fields to meet other needs (later, later)
  - Differentiated Services, Time to live (TTL)

Later, with QOS

Later, with ICMP

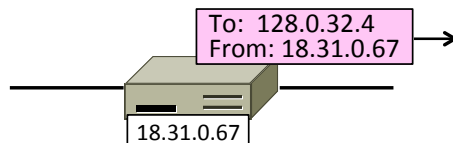| 32 Bits | | | | |
|---|---|---|---|---|
| Version | IHL | Differentiated Services | | Total length |
| Identification | | | D F M F | Fragment offset |
| Time to live | | Protocol | | Header checksum |
| Source address | | | | |
| Destination address | | | | |
| Options (0 or more words) | | | | |
| Payload (e.g., TCP segment) | | | | |

Computer Networks

36

18

# IP Prefixes (§5.6.1-5.6.2)

- What do IP addresses look like?
  - And IP prefixes, or blocks of addresses
  - (This is IPv4; we'll cover IPv6 later.)

To:  128.0.32.4
From: 18.31.0.67

18.31.0.67

Computer Networks

37

# IP Addresses

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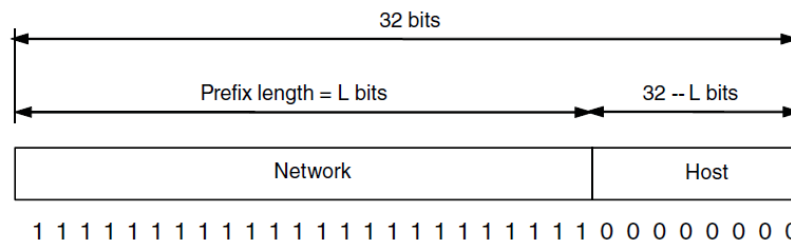  - Four 8-bit numbers separated by dots

| 8 bits | 8 bits | 8 bits | 8 bits | | |
|---|---|---|---|---|---|
| aaaaaaaa | bbbbbbbb | cccccccc | dddddddd | ↔ | A.B.C.D |
| 00010010 | 00011111 | 00000000 | 00000001 | ↔ | |

Computer Networks

38

# IP Prefixes – Modern

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

32 bits

Prefix length = L bits | 32 -- L bits

| Network | Host |

1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0 0 0 0 0

Computer Networks
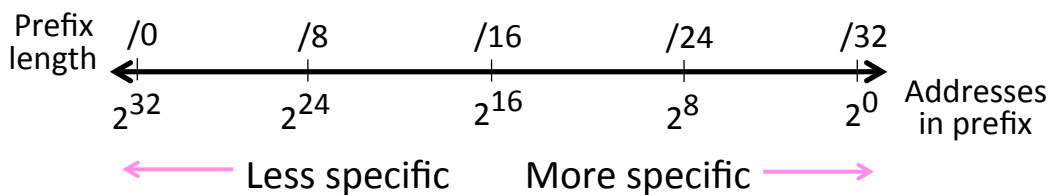
39

# IP Prefixes (2)

- Written in "IP address/length" notation
  - Address is lowest address in the prefix, length is prefix bits
  - E.g., 128.13.0.0/16 is 128.13.0.0 to 128.13.255.255
  - So a /24 ("slash 24") is 256 addresses, and a /32 is one address

```
00010010 00011111 00000000 xxxxxxxx  ↔
```

↔ `128.13.0.0/16`

Computer Networks

40

# IP Prefixes (3)

- <u>More specific</u> prefix
  - Has longer prefix, hence a smaller number of IP addresses
- <u>Less specific</u> prefix
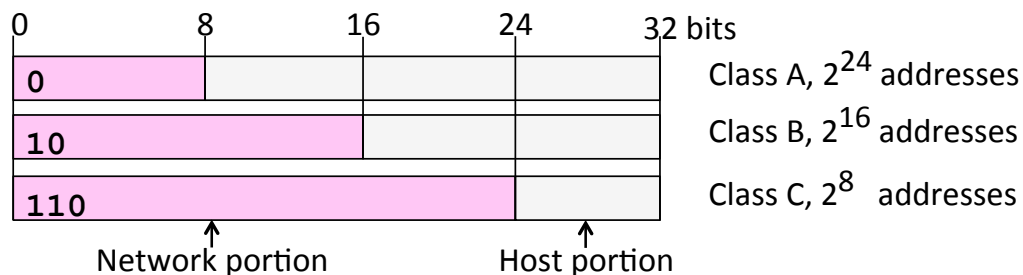  - Has shorter prefix, hence a larger number of IP addresses

Prefix length  /0      /8      /16      /24      /32

$2^{32}$      $2^{24}$      $2^{16}$      $2^8$      $2^0$   Addresses in prefix

←—— Less specific      More specific ——→

Computer Networks

41

# IP Address Classes – Historical

- 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

0        8        16        24        32 bits

| 0 | | | |
| 10 | | | |
| 110 | | | |

Class A, $2^{24}$ addresses

Class B, $2^{16}$ addresses

Class C, $2^8$ addresses

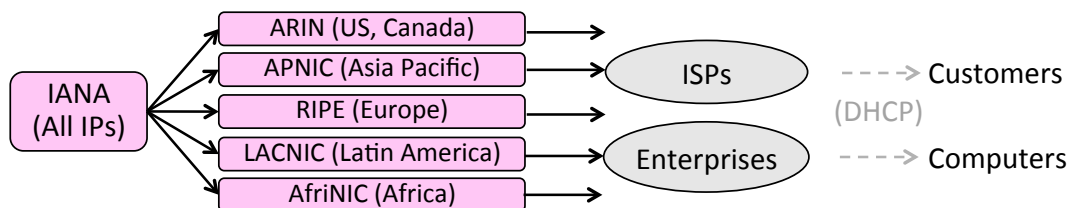Network portion      Host portion

Computer Networks

42

# Public / Private IP Addresses

- Public IP addresses, e.g., 18.31.0.1
  - Valid destination on the global Internet
  - Must be allocated to you before use
  - Now exhausted … time for IPv6!

- Private IP addresses
  - Can be used freely within private networks (home, small company)
  - 10.0.0.0/8, 172.16.0.0/12, 192.168.0.0/16
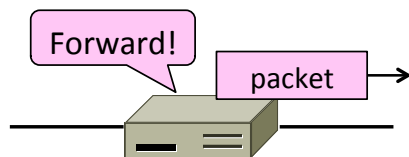  - Need public IP address(es) and NAT to connect to global Internet

Computer Networks

43

---

# Allocating Public IP Addresses

- Follows a hierarchical process
  - IANA delegates to regional bodies (RIRs)
  - RIRs delegate to companies in their region
  - Companies assign to their customers/computers (later, DHCP)



| IANA (All IPs) | ARIN (US, Canada) | | |
| | APNIC (Asia Pacific) | ISPs | - - -> Customers (DHCP) |
| | RIPE (Europe) | | |
| | LACNIC (Latin America) | Enterprises | - - -> Computers |
| | AfriNIC (Africa) | | |

Computer Networks

44

# IP Forwarding (§5.6.1-5.6.2)

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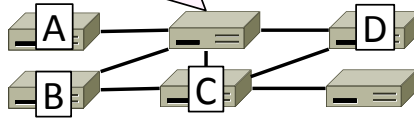  - We'll look at how IP does it
  - (We'll cover routing later)



Computer Networks

45

# Recap

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

Computer Networks

46

# IP Forwarding

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

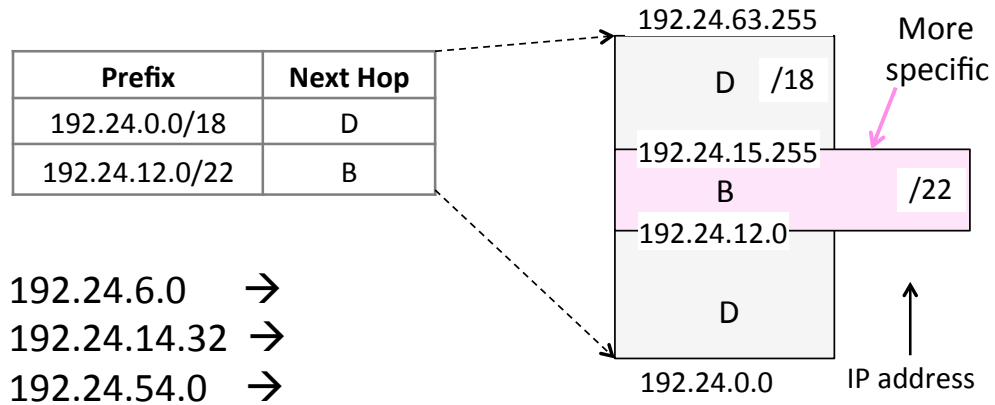| Prefix | Next Hop |
|---|---|
| 192.24.0.0/18 | D |
| 192.24.12.0/22 | B |



Computer Networks

47

# Longest Matching Prefix

- Prefixes in the table might overlap!
  - Combines hierarchy with flexibility

- <u>Longest matching prefix</u> forwarding rule:
  - For each packet, find the longest prefix that contains the destination address, i.e., the most specific entry
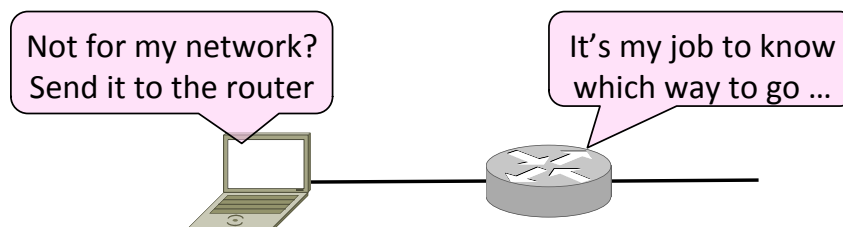  - Forward the packet to the next hop router for that prefix

Computer Networks

48

# Longest Matching Prefix (2)

| Prefix | Next Hop |
|--------|----------|
| 192.24.0.0/18 | D |
| 192.24.12.0/22 | B |

192.24.6.0    →
192.24.14.32  →
192.24.54.0   →

192.24.63.255

D    /18

More specific

192.24.15.255

B    /22

192.24.12.0

D

192.24.0.0    IP address

Computer Networks                    49

# 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

Not for my network?
Send it to the router

It's my job to know
which way to go …

Computer Networks                    50

25

# 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 direct to that IP |
| 0.0.0.0/0 | Send to my router |

Computer Networks

51
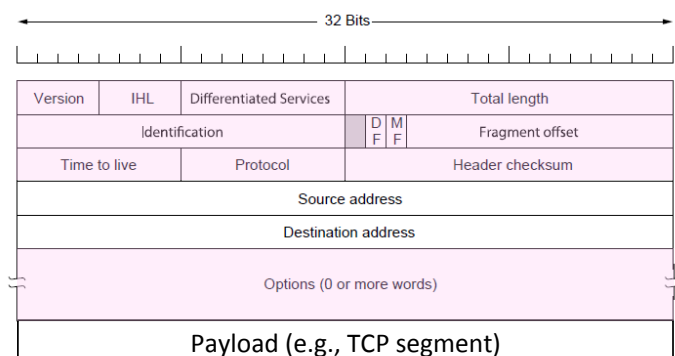
# Flexibility of Longest Matching Prefix

- Can provide default behavior, with less specific prefixes
  - To send traffic going outside an organization to a border router

- Can special case behavior, with more specific prefixes
  - For performance, economics, security, …

Computer Networks

52

# Performance of Longest Matching Prefix

- Uses hierarchy for a compact table
    - Benefits from less specific prefixes

- Lookup more complex than table
    - Was a concern for fast routers, but not an issue in practice these days

Computer Networks

53

---

# Other Aspects of Forwarding

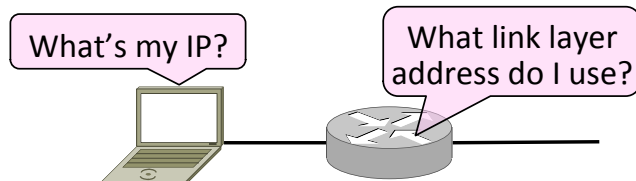- It's not all about addresses …



Computer Networks

54

# Other Aspects (2)

- Decrement TTL value
  - Protects against loops
- Checks header checksum
  - To add reliability
- Fragment large packets
  - Split to fit it on next link
- Send congestion signals
  - Warns hosts of congestion
- Generates error messages
  - To help mange network
- Handle various options

Coming later

Computer Networks

55

# Helping IP with ARP, DHCP (§5.6.4)

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

What's my IP?

What link layer address do I use?

Computer Networks

56

# 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.
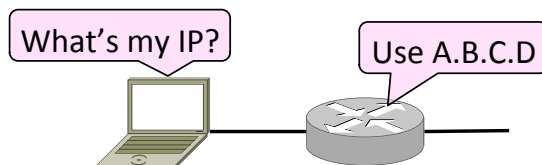  - At least Ethernet address is on NIC

Hey, where am I?

Computer Networks

57

# 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 folks

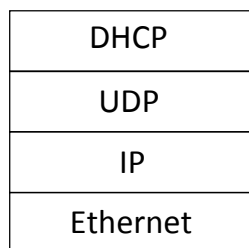What's my IP?   Use A.B.C.D

Computer Networks

58

29

2014-03-27

# 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.

Computer Networks                                                                59
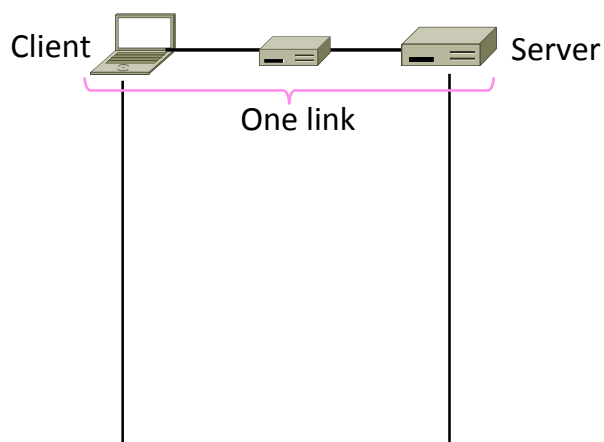
---

# DHCP Protocol Stack

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

| DHCP |
| :---: |
| UDP |
| IP |
| Ethernet |

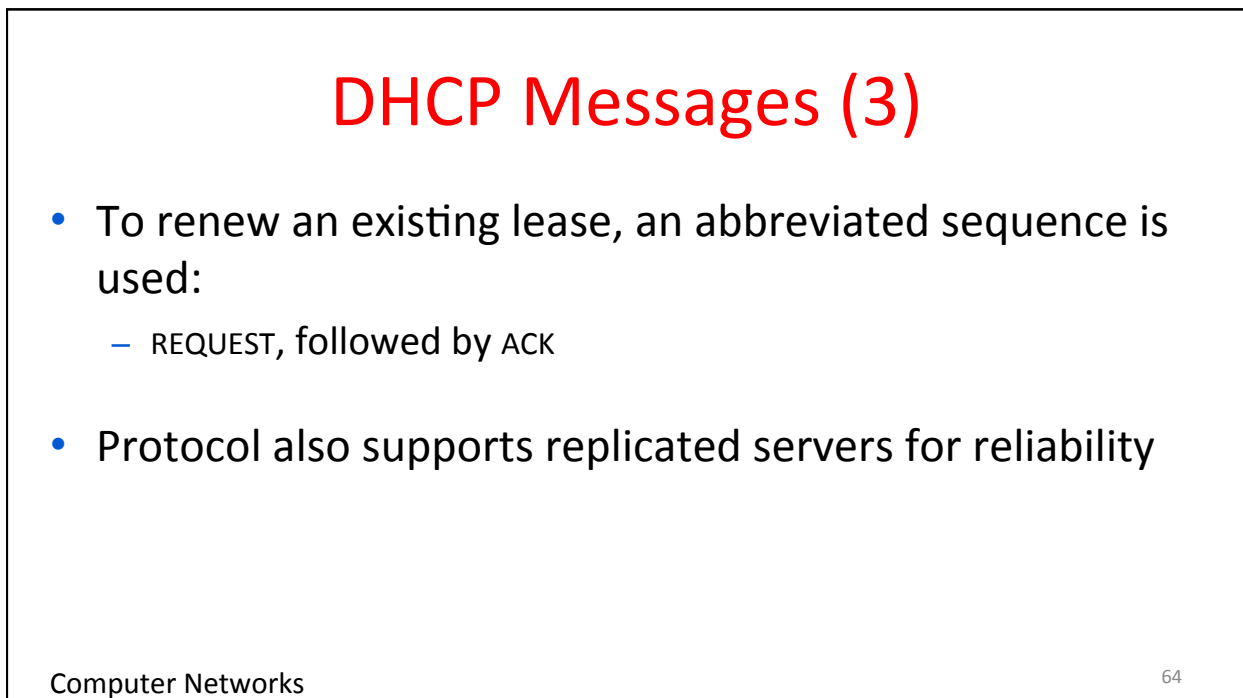Computer Networks                                                                60

# DHCP Addressing

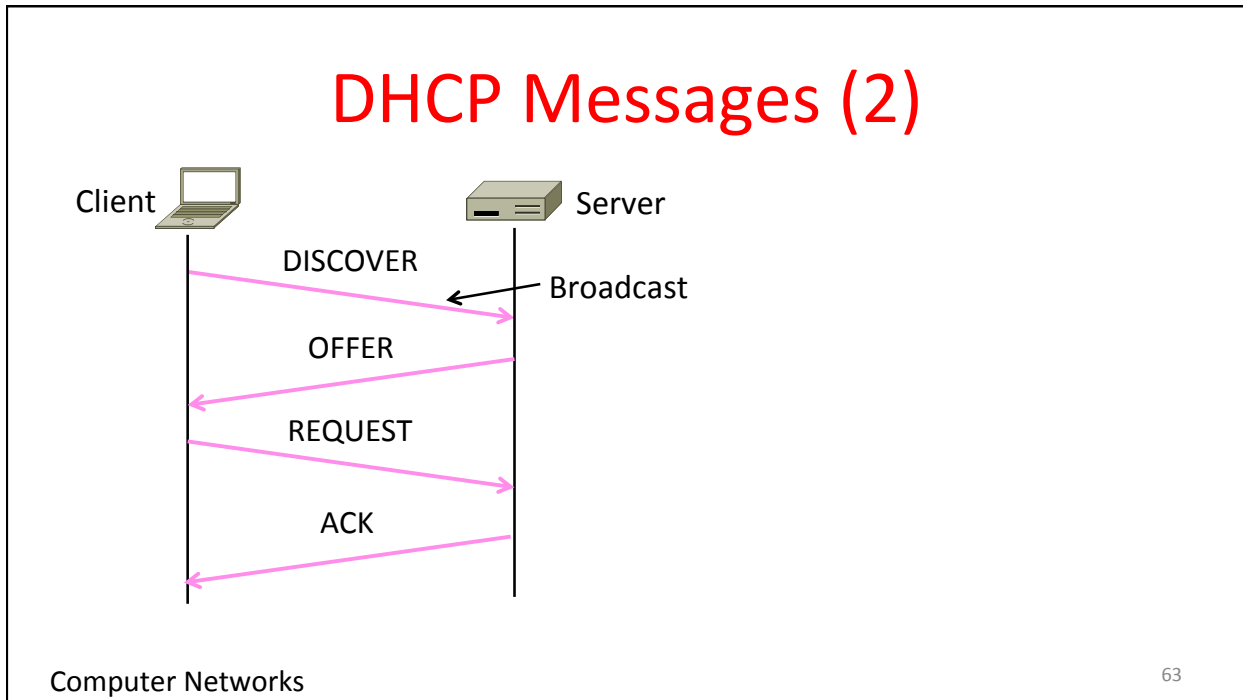- Bootstrap issue:
  - How does node send a message to DHCP server before it is configured?

- Answer:
  - Node sends <u>broadcast</u> messages that delivered to all nodes on the network
  - <u>Broadcast address</u> is all 1s
  - IP (32 bit): 255.255.255.255
  - Ethernet (48 bit): ff:ff:ff:ff:ff:ff

Computer Networks

61

# DHCP Messages

Client               Server

One link

Computer Networks

62

# DHCP Messages (2)

Client                     Server

DISCOVER

Broadcast

OFFER

REQUEST

ACK

Computer Networks

63

---

# DHCP Messages (3)

- To renew an existing lease, an abbreviated sequence is used:
  - REQUEST, followed by ACK

- Protocol also supports replicated servers for reliability
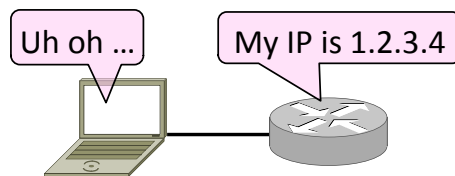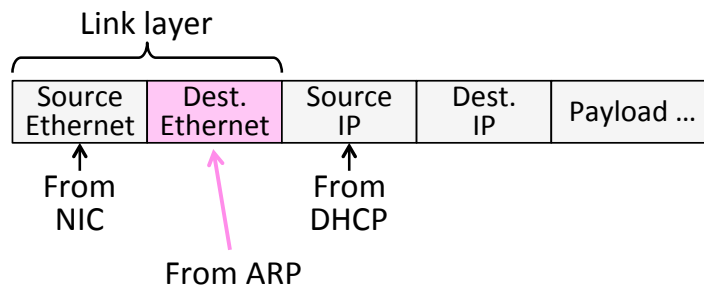
Computer Networks

64

32

# 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?

Uh oh …

My IP is 1.2.3.4

Computer Networks

65

# ARP (Address Resolution Protocol)

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

Link layer

| Source Ethernet | Dest. Ethernet | Source IP | Dest. IP | Payload … |
|---|---|---|---|---|

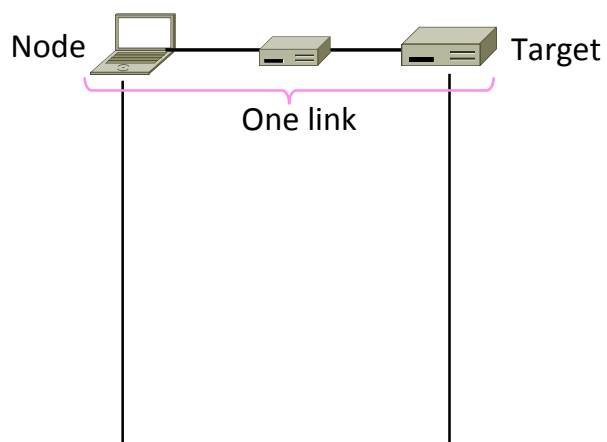From NIC

From ARP

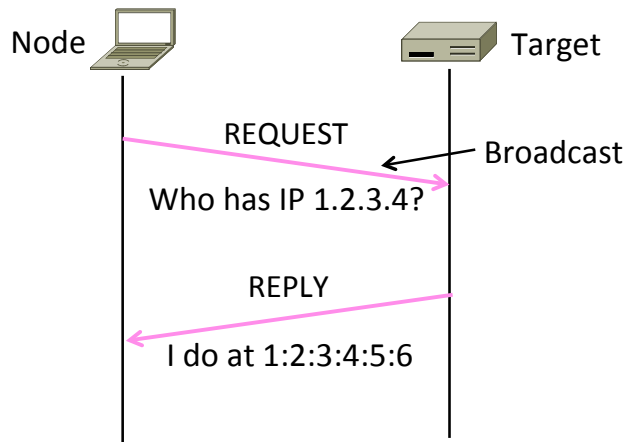From DHCP

Computer Networks

66

33

# 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 |
|:---:|
| Ethernet |

Computer Networks

67

---

# ARP Messages

Node ———— Target

One link

Computer Networks

68

# ARP Messages (2)

Node            Target

```
        REQUEST        ← Broadcast
     Who has IP 1.2.3.4?

         REPLY
     I do at 1:2:3:4:5:6
```

---

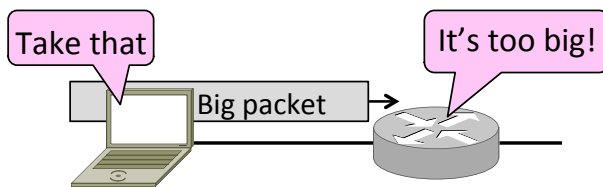# Discovery Protocols

- 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

# Packet Fragmentation (§5.5.5)

- 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 MTUs
  - 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, but dated

- Discovery (next)
  - Find the largest packet that fits on the network path and use it
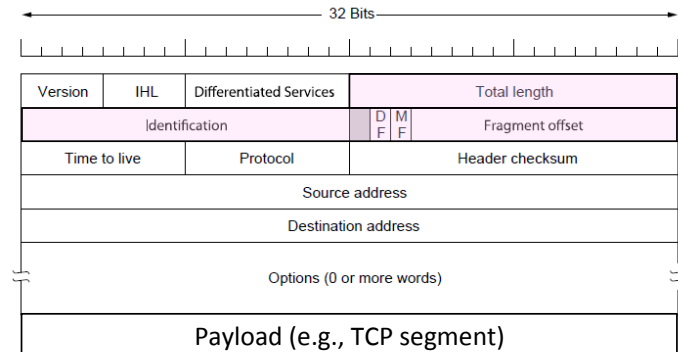  - IP uses today instead of fragmentation

Computer Networks

73

# IPv4 Fragmentation

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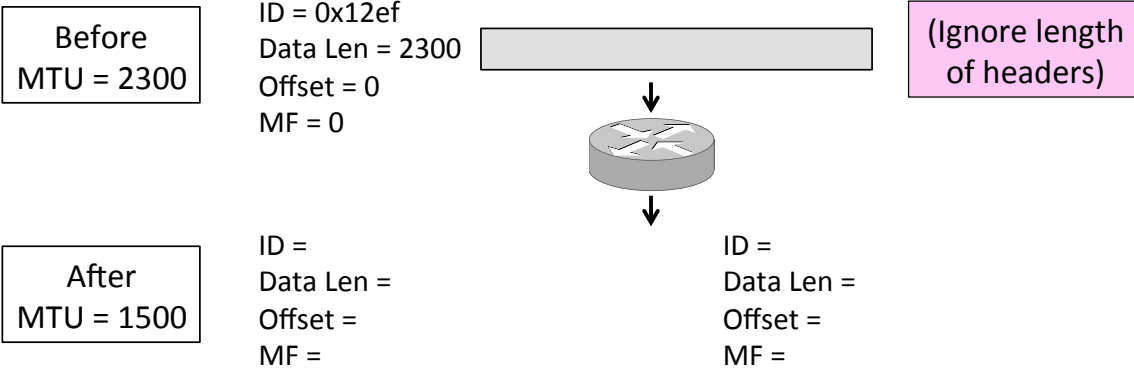- Receiving host reassembles to reduce load on routers



Computer Networks

74

# IPv4 Fragmentation Fields

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

| 32 Bits | | | | | | |
|---|---|---|---|---|---|---|
| Version | IHL | Differentiated Services | | Total length | | |
| Identification | | | DF MF | Fragment offset | | |
| Time to live | | Protocol | | Header checksum | | |
| Source address | | | | | | |
| Destination address | | | | | | |
| Options (0 or more words) | | | | | | |
| Payload (e.g., TCP segment) | | | | | | |

Computer Networks

75

---

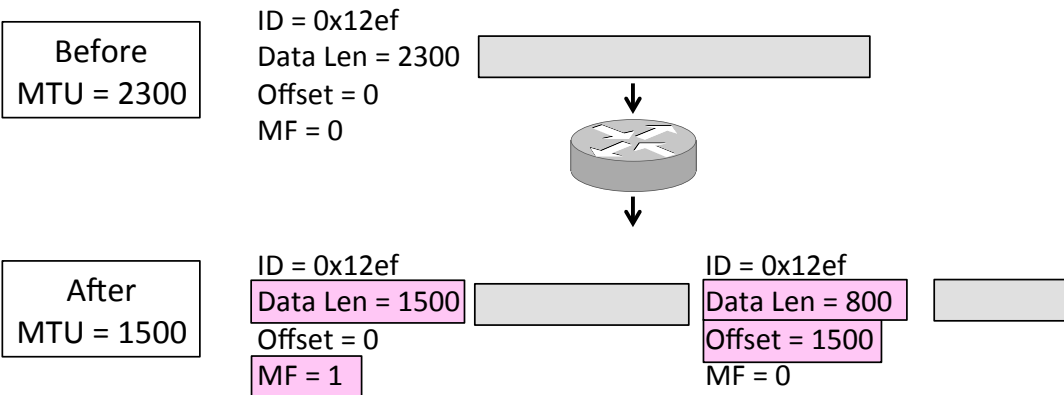# IPv4 Fragmentation Procedure

- 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 pieces:
  - Identification field links pieces together, MF tells receiver when it has all pieces

Computer Networks

76

# IPv4 Fragmentation (2)

| Before MTU = 2300 | ID = 0x12ef<br>Data Len = 2300<br>Offset = 0<br>MF = 0 | | (Ignore length of headers) |

| After MTU = 1500 | ID =<br>Data Len =<br>Offset =<br>MF = | ID =<br>Data Len =<br>Offset =<br>MF = |

Computer Networks

77

# IPv4 Fragmentation (3)

| Before MTU = 2300 | ID = 0x12ef<br>Data Len = 2300<br>Offset = 0<br>MF = 0 |

| After MTU = 1500 | ID = 0x12ef<br>Data Len = 1500<br>Offset = 0<br>MF = 1 | ID = 0x12ef<br>Data Len = 800<br>Offset = 1500<br>MF = 0 |

Computer Networks
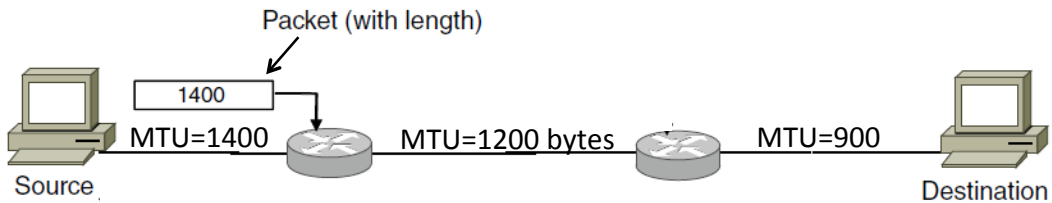
78

# 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

Computer Networks

79

# 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

Computer Networks

80

# Path MTU Discovery (2)

Packet (with length)

1400

MTU=1400    MTU=1200 bytes    MTU=900

Source    Destination

Computer Networks    81

# Path MTU Discovery (3)

Packet (with length)

Test #1    Test #2    Test #3

1400    1200    900

MTU=1400    MTU=1200 bytes    MTU=900

Source    Destination
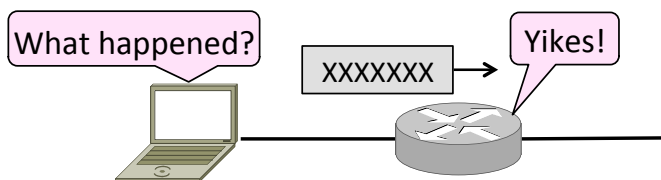
Try 1200    Try 900

Computer Networks    82

# Path MTU Discovery (4)

- Process may seem involved
  - But usually quick to find right size

- Path MTU depends on the path, 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

Computer Networks

83

# Error Handling with ICMP (§5.6.4)

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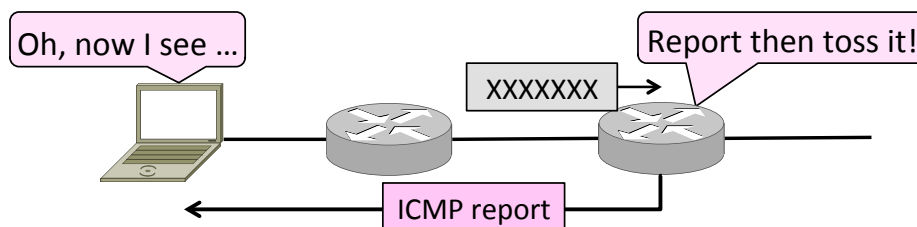


Computer Networks

84

# 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

Computer Networks

85

# ICMP Errors

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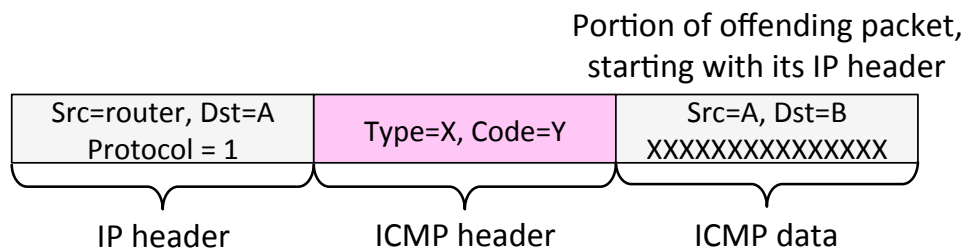  - It discards the problematic packet; host needs to rectify

Oh, now I see …

Report then toss it!

XXXXXXX

ICMP report

Computer Networks

86

# 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

Computer Networks

87

# 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
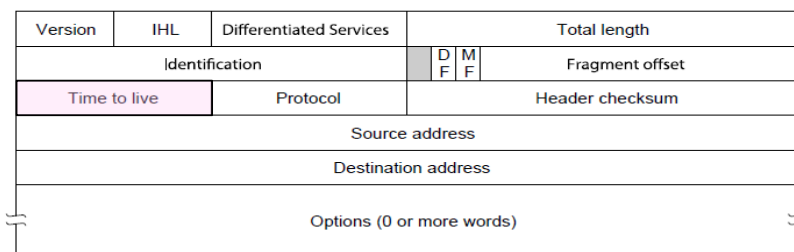
Portion of offending packet,
starting with its IP header

| Src=router, Dst=A<br>Protocol = 1 | Type=X, Code=Y | Src=A, Dst=B<br>XXXXXXXXXXXXXXX |
|---|---|---|
| IP header | ICMP header | ICMP data |

Computer Networks

88

44

# 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 | Traceroute |
| Echo Request or Reply | 8 or 0 / 0 | Ping |

Testing, not a forwarding error: Host sends Echo
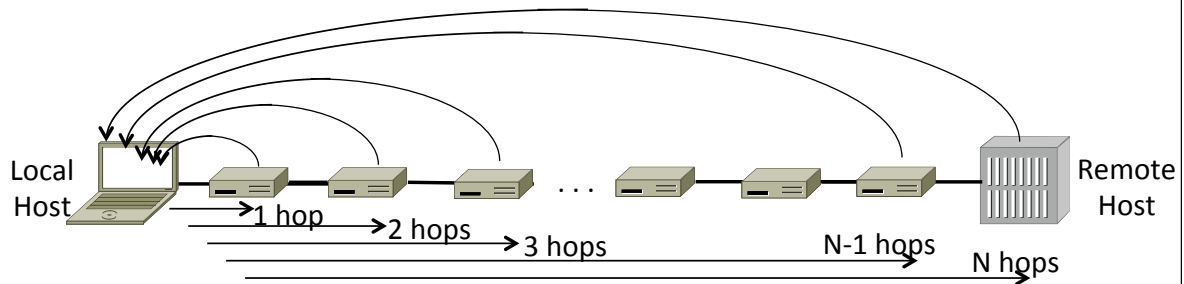Request, and destination responds with an Echo Reply

Computer Networks

89

# Traceroute

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

| Version | IHL | Differentiated Services | | | Total length | |
|---------|-----|-------------------------|---|---|--------------|---|
| Identification | | | | DF MF | Fragment offset | |
| Time to live | | Protocol | | | Header checksum | |
| Source address | | | | | | |
| Destination address | | | | | | |
| Options (0 or more words) | | | | | | |

Computer Networks

90

# Traceroute (2)

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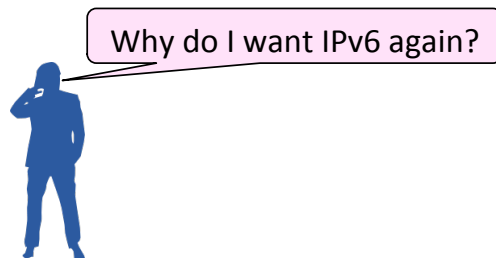
Local Host

Remote Host

1 hop
2 hops
3 hops
N-1 hops
N hops

. . .

Computer Networks

91

# IP Version 6 (§5.6.3)

- IP version 6, the future of IPv4 that is now (still) being deployed

Why do I want IPv6 again?
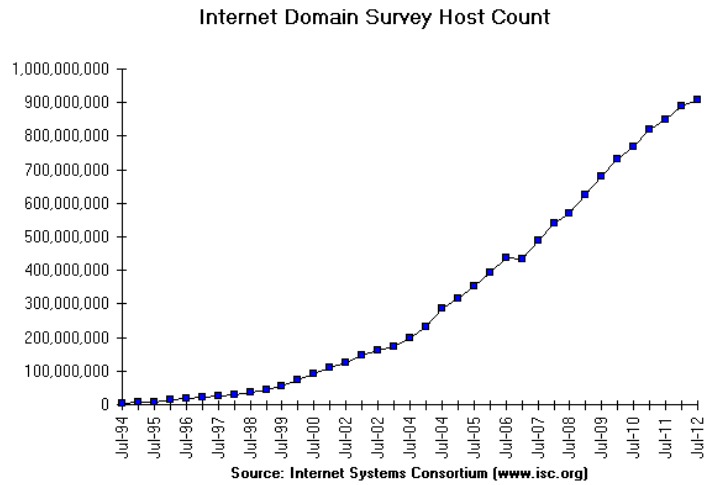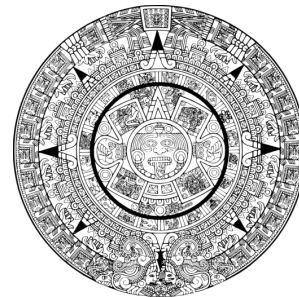
Computer Networks

92

# Internet Growth

Internet Domain Survey Host Count

- At least a billion Internet hosts and growing …

- And we're using 32-bit addresses!

Source: Internet Systems Consortium (www.isc.org)

Computer Networks

93

# The End of New IPv4 Addresses

- Now running on leftover blocks held by the regional registries; much tighter allocation policies

Exhausted on 4/11 and 9/12!

ARIN (US, Canada)

APNIC (Asia Pacific)

IANA (All IPs)

RIPE (Europe)

ISPs

Companies

LACNIC (Latin Amer.)

AfriNIC (Africa)

Exhausted on 2/11!

End of the world ? 12/21/12?

Computer Networks

94

47

# IP Version 6 to the Rescue

- Effort started by the IETF in 1994
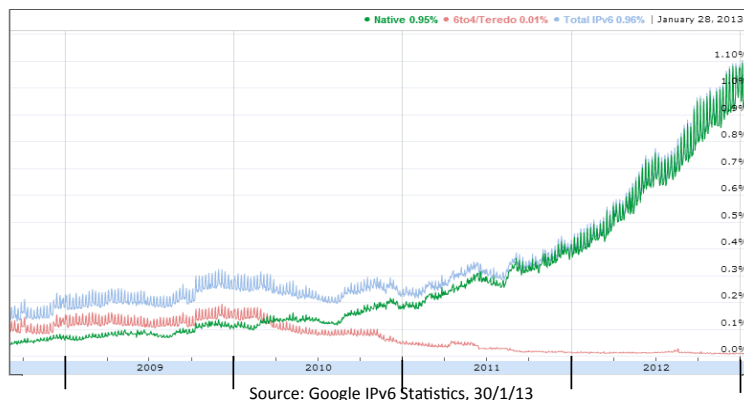  - Much larger addresses (128 bits)
  - Many sundry improvements

- Became an IETF standard in 1998
  - Nothing much happened for a decade
  - Hampered by deployment issues, and a lack of adoption incentives
  - Big push ~2011 as exhaustion looms

Computer Networks                                                      95

# IPv6 Deployment

Percentage of users accessing Google via IPv6

Time for growth!



Source: Google IPv6 Statistics, 30/1/13

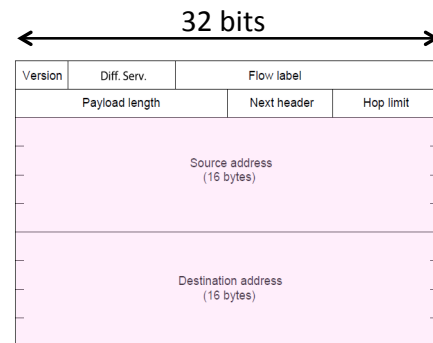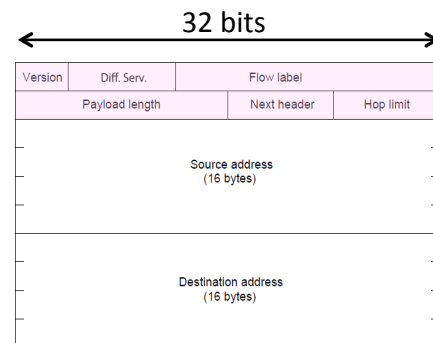Computer Networks                                                      96

# IPv6

- Features large addresses
  - 128 bits, most of header
- New notation
  - 8 groups of 4 hex digits (16 bits)
  - Omit leading zeros, groups of zeros

  Ex:   2001:0db8:0000:0000:0000:ff00:0042:8329
  →

| | 32 bits | | |
|---|---|---|---|
| Version | Diff. Serv. | Flow label | |
| Payload length | | Next header | Hop limit |
| Source address (16 bytes) | | | |
| Destination address (16 bytes) | | | |

Computer Networks

97

# IPv6 (2)

- Lots of other, smaller changes
  - Streamlined header processing
  - Flow label to group of packets
  - Better fit with "advanced" features (mobility, multicasting, security)

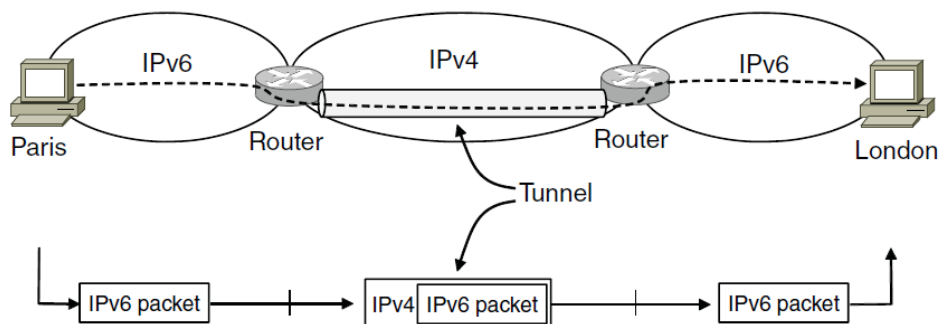| | 32 bits | | |
|---|---|---|---|
| Version | Diff. Serv. | Flow label | |
| Payload length | | Next header | Hop limit |
| Source address (16 bytes) | | | |
| Destination address (16 bytes) | | | |

Computer Networks

98

# IPv6 Transition

- The Big Problem:
  - How to deploy IPv6?
  - Fundamentally incompatible with IPv4

- Dozens of approaches proposed
  - Dual stack (speak IPv4 and IPv6)
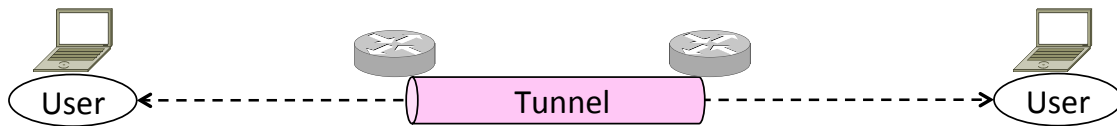  - Translators (convert packets)
  - Tunnels (carry IPv6 over IPv4)

Computer Networks

99

# Tunneling

- Native IPv6 islands connected via IPv4
  - Tunnel carries IPv6 packets across IPv4 network
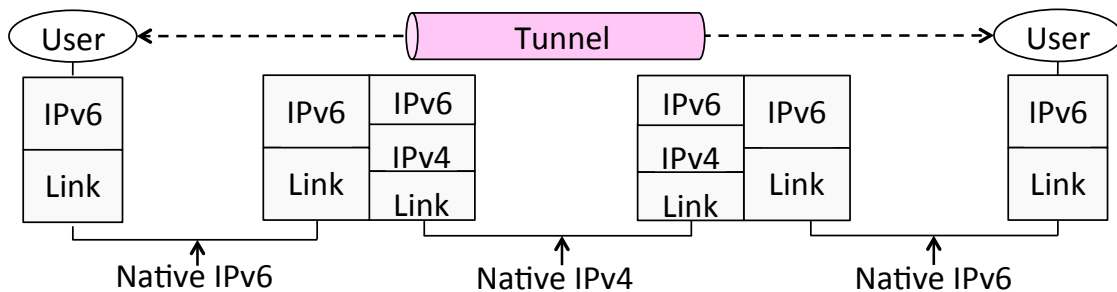


Computer Networks

100

50

# Tunneling (2)

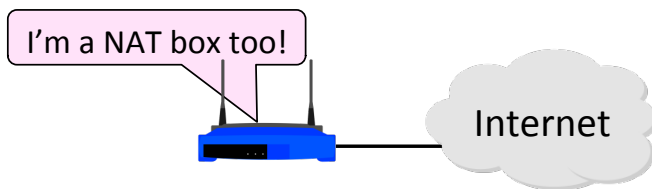- Tunnel acts as a single link across IPv4 network



Computer Networks

101

# Tunneling (3)

- Tunnel acts as a single link across IPv4 network
  - Difficulty is to set up tunnel endpoints and routing



Computer Networks

102

# Network Address Translation (§5.6.2)

- What is NAT (Network Address Translation)? How does it work?
  - NAT is widely used at the edges of the network, e.g., homes



I'm a NAT box too!

Internet

Computer Networks

103

# Layering Review

- Remember how layering is meant to work?
  - "Routers don't look beyond the IP header." Well …



App

Router

App

| TCP | | | | TCP |
| IP | IP | IP | IP | IP | IP |
| 802.11 | 802.11 | Ethernet | Ethernet | 802.11 | 802.11 |

Computer Networks

104

52

# Middleboxes

- Sit "inside the network" but perform "more than IP" processing on packets to add new functionality
  - NAT box, Firewall / Intrusion Detection System

App                          Middlebox                          App

| TCP | | App / TCP | | | | | TCP |
| IP | | IP | IP | | IP | IP | | IP |
| 802.11 | | 802.11 | Ethernet | | Ethernet | 802.11 | | 802.11 |

Computer Networks                                                                 105

---

# Middleboxes (2)

- Advantages
  - A possible rapid deployment path when there is no other option
  - Control over many hosts (IT)

- Disadvantages
  - Breaking layering interferes with connectivity; strange side effects
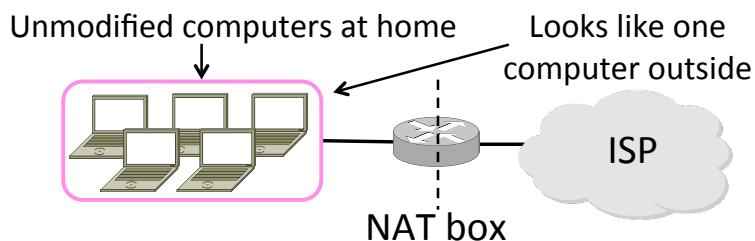  - Poor vantage point for many tasks

Computer Networks                                                                 106

# NAT (Network Address Translation) Box

- NAT box connects an internal network to an external network
  - Many internal hosts are connected using few external addresses
  - Middlebox that "translates addresses"

- Motivated by IP address scarcity
  - Controversial at first, now accepted

Computer Networks

107

# NAT (2)

- Common scenario:
  - Home computers use "private" IP addresses
  - NAT (in AP/firewall) connects home to ISP using a single external IP address

Unmodified computers at home    Looks like one computer outside

ISP

NAT box

Computer Networks

108

# How NAT Works

- Keeps an internal/external table
    - Typically uses IP address + TCP port
    - This is address and port translation

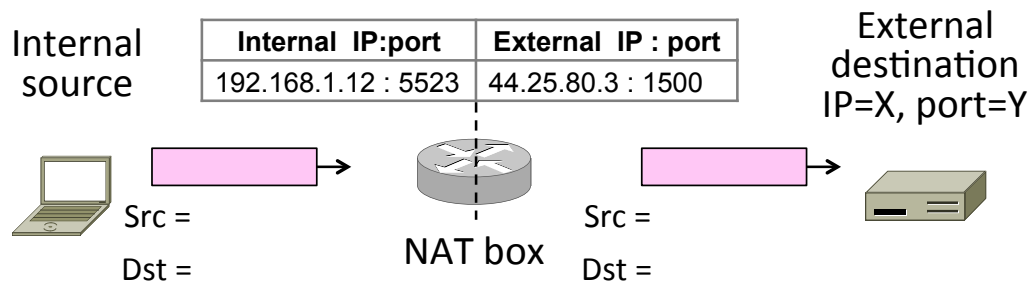<span style="color:magenta">What host thinks</span>     <span style="color:magenta">What ISP thinks</span>

| Internal  IP:port | External  IP : port |
|---|---|
| 192.168.1.12 : 5523 | 44.25.80.3 : 1500 |
| 192.168.1.13 : 1234 | 44.25.80.3 : 1501 |
| 192.168.2.20 : 1234 | 44.25.80.3 : 1502 |

- Need ports to make mapping 1-1 since there are fewer external IPs

Computer Networks

109

---

# How NAT Works (2)

- Internal → External:
    - Look up and rewrite Source IP/port

Internal source

| Internal  IP:port | External  IP : port |
|---|---|
| 192.168.1.12 : 5523 | 44.25.80.3 : 1500 |

External destination
IP=X, port=Y

Src =

Dst =

NAT box

Src =

Dst =

Computer Networks

110

55

# How NAT Works (3)

- External → Internal
  - Look up and rewrite Destination IP/port

Internal destination

| Internal IP:port | External IP : port |
|---|---|
| 192.168.1.12 : 5523 | 44.25.80.3 : 1500 |

External source
IP=X, port=Y

Src =

Dst =

NAT box

Src =

Dst =

Computer Networks

111

# How NAT Works (4)

- Need to enter translations in the table for it to work
  - Create external name when host makes a TCP connection

Internal source

| Internal IP:port | External IP : port |
|---|---|
| 192.168.1.12 : 5523 | |

External destination
IP=X, port=Y

Src =

Dst =

NAT box

Src =

Dst =

Computer Networks

112

# NAT Downsides

- Connectivity has been broken!
  - Can only send incoming packets after an outgoing connection is set up
  - Difficult to run servers or peer-to-peer apps (Skype) at home

- Doesn't work so well when there are no connections (UDP apps)

- Breaks apps that unwisely expose their IP addresses (FTP)

Computer Networks                                                    113

# NAT Upsides

- Relieves much IP address pressure
  - Many home hosts behind NATs
- Easy to deploy
  - Rapidly, and by you alone
- Useful functionality
  - Firewall, helps with privacy

- Kinks will get worked out eventually
  - "NAT Traversal" for incoming traffic

Computer Networks                                                    114