# COMP 3331/9331: Computer Networks and Applications

Week 8

Network Layer Part 2

**Reading Guide: Chapter 4: Sections 4.4**

# Network Layer: outline

4.1 introduction

4.2 virtual circuit and datagram networks

4.3 what's inside a router

4.4 IP: Internet Protocol
- datagram format
- IPv4 addressing
- ICMP
- IPv6

4.5 routing algorithms
- link state
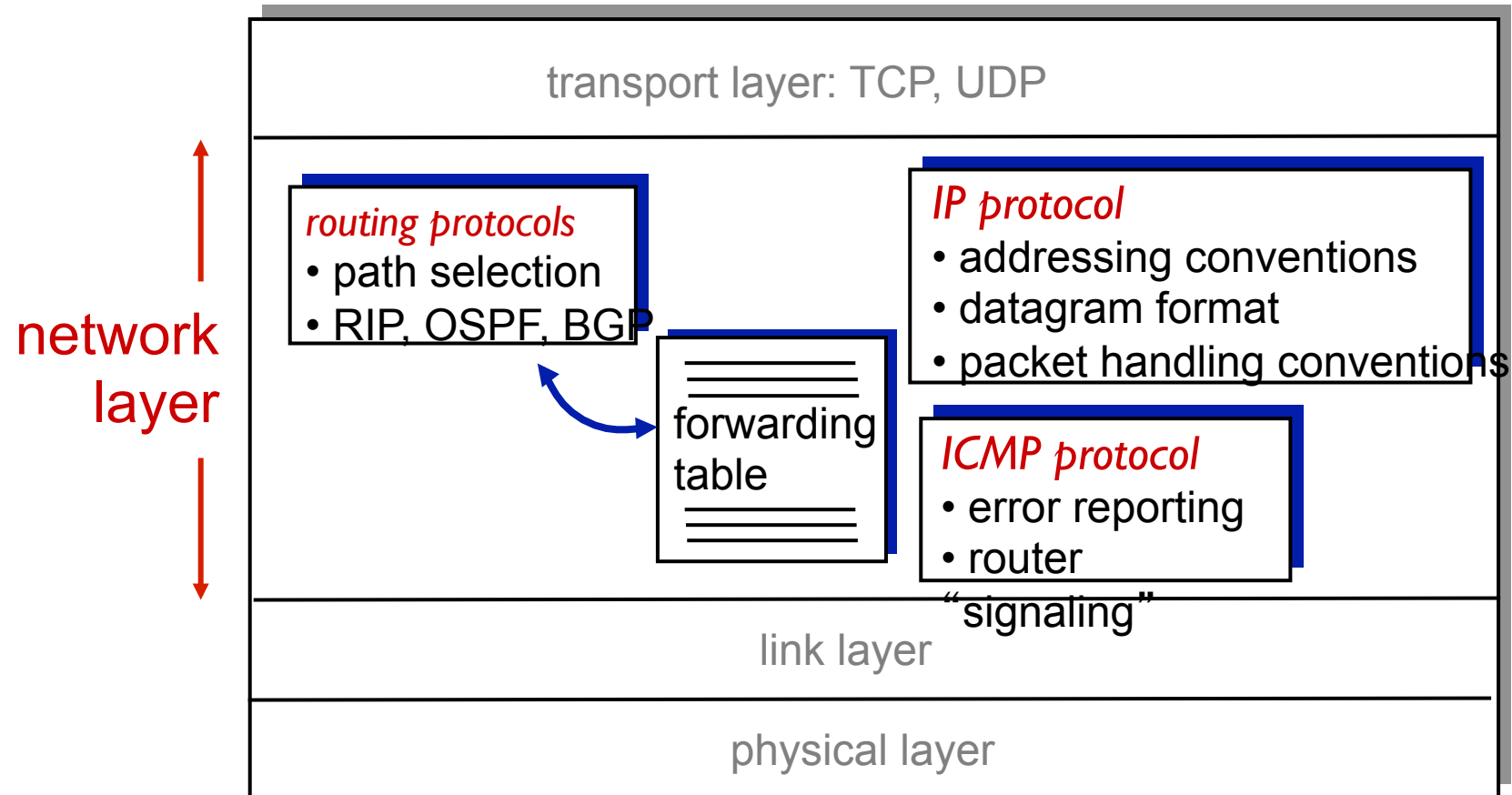- distance vector
- hierarchical routing

4.6 routing in the Internet
- RIP
- OSPF
- BGP
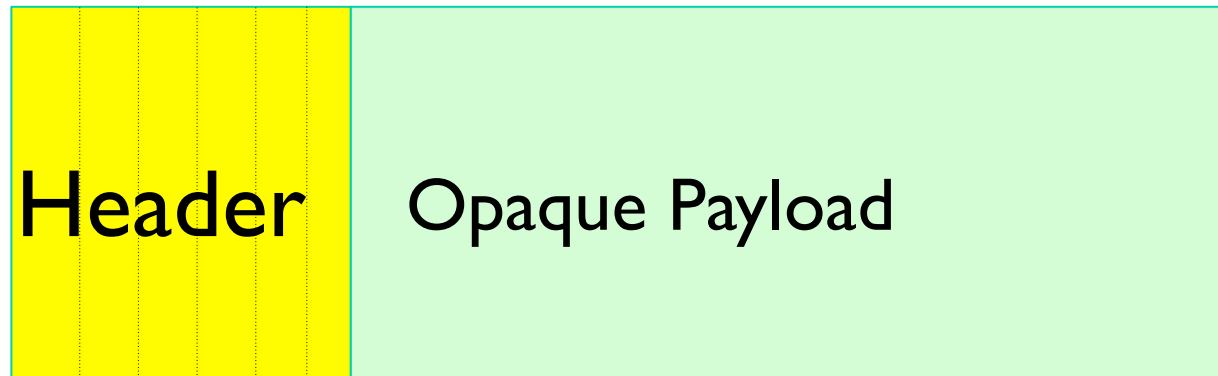
4.7 broadcast and multicast routing

# The Internet network layer

host, router network layer functions:



network layer

| transport layer: TCP, UDP |
| link layer |
| physical layer |

routing protocols
• path selection
• RIP, OSPF, BGP

forwarding table

IP protocol
• addressing conventions
• datagram format
• packet handling conventions

ICMP protocol
• error reporting
• router "signaling"

# What is Designing IP?

❖ Syntax: format of packet
- Nontrivial part: packet "header"
- Rest is opaque payload *(why opaque?)*

| Header | Opaque Payload |
|--------|----------------|

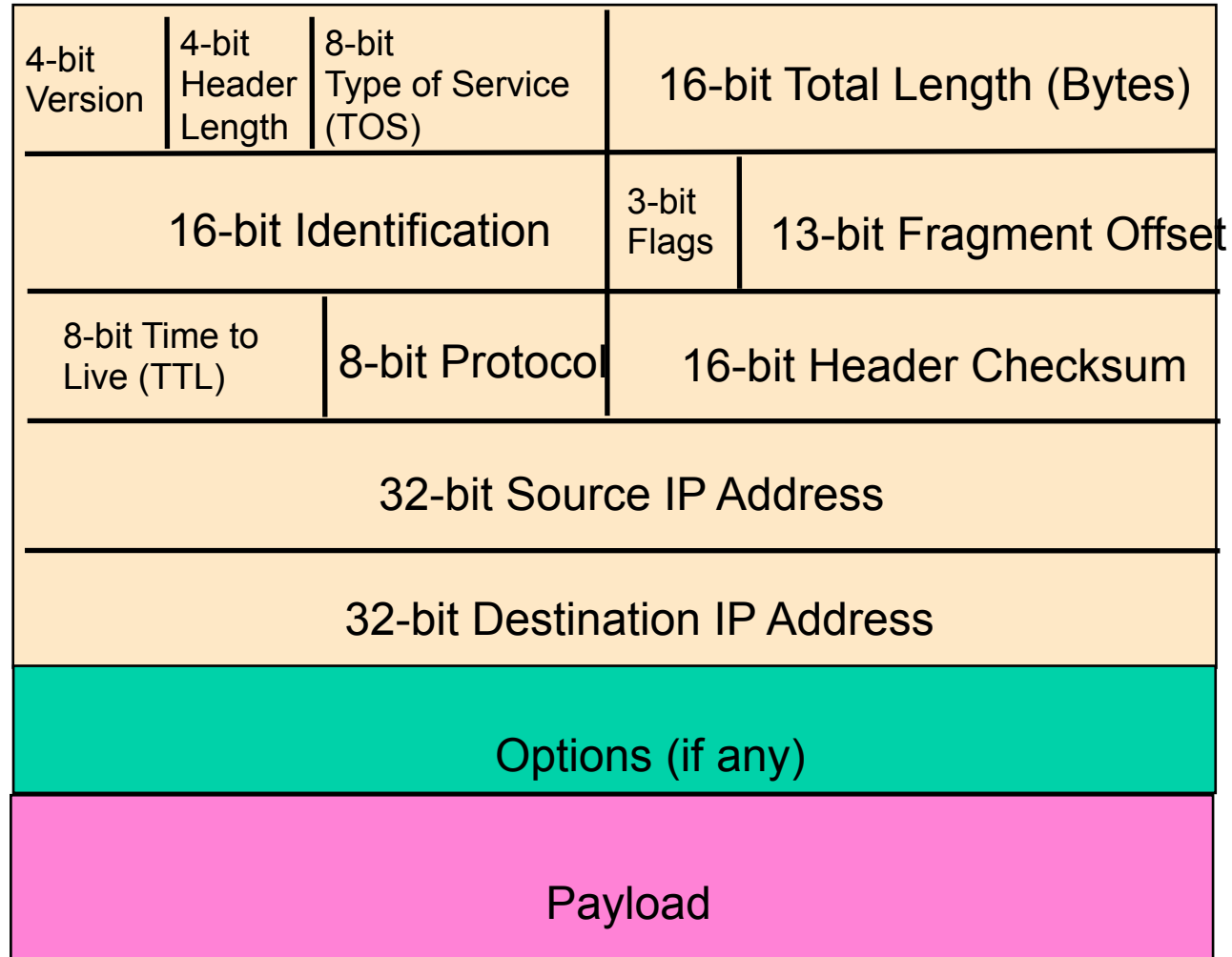❖ Semantics: meaning of header fields
- Required processing

# Packet Headers

❖ Think of packet header as interface
  ▪ Only way of passing information from packet to router

❖ Designing interfaces
  ▪ What task are you trying to perform?
  ▪ What information do you need to accomplish it?
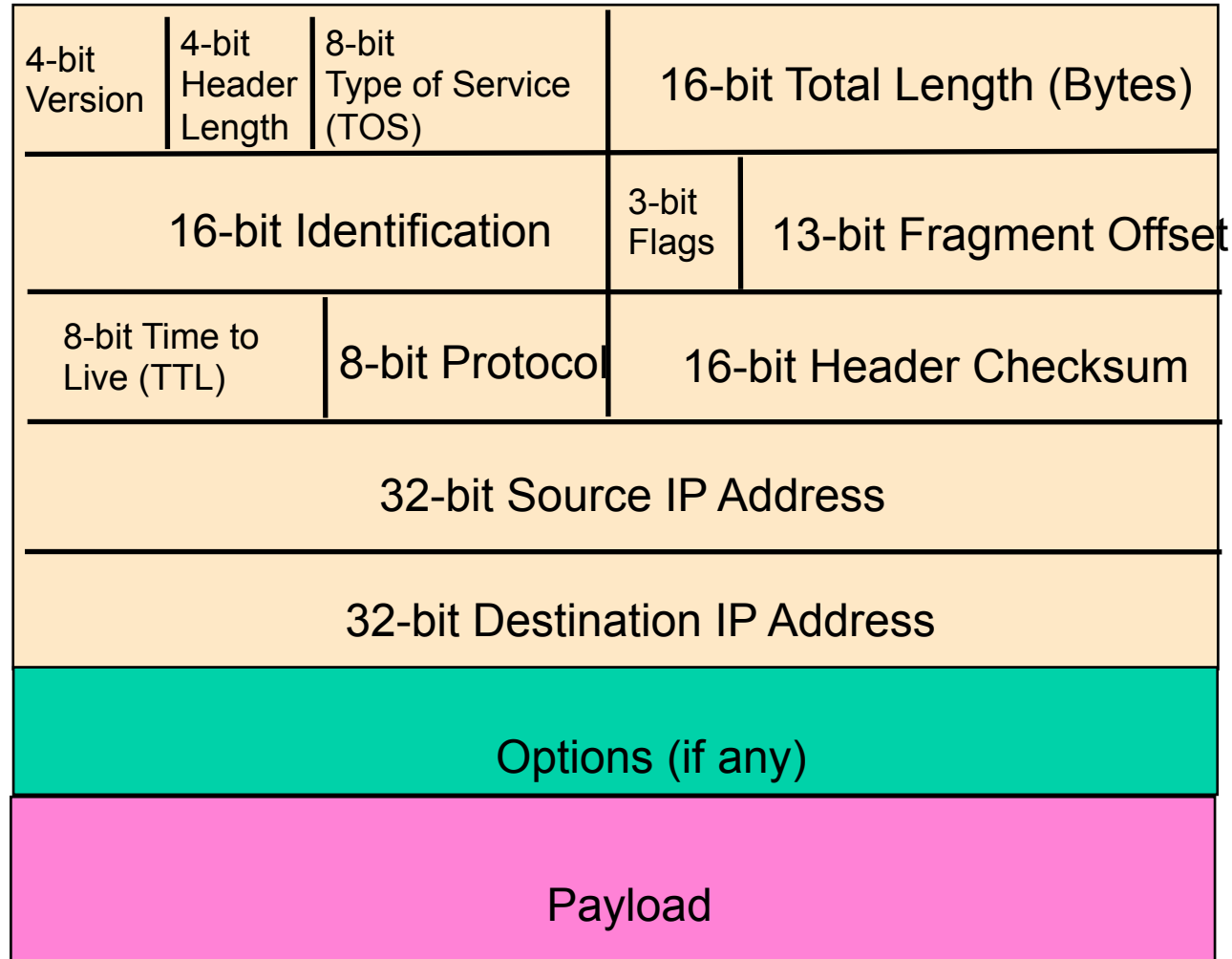
❖ Header reflects information needed for basic tasks

# What Tasks Do We Need to Do?

❖ Read packet correctly

❖ Get packet to the destination; responses back to the source

❖ Carry data

❖ Tell host what to do with packet once arrived

❖ Specify any special network handling of the packet

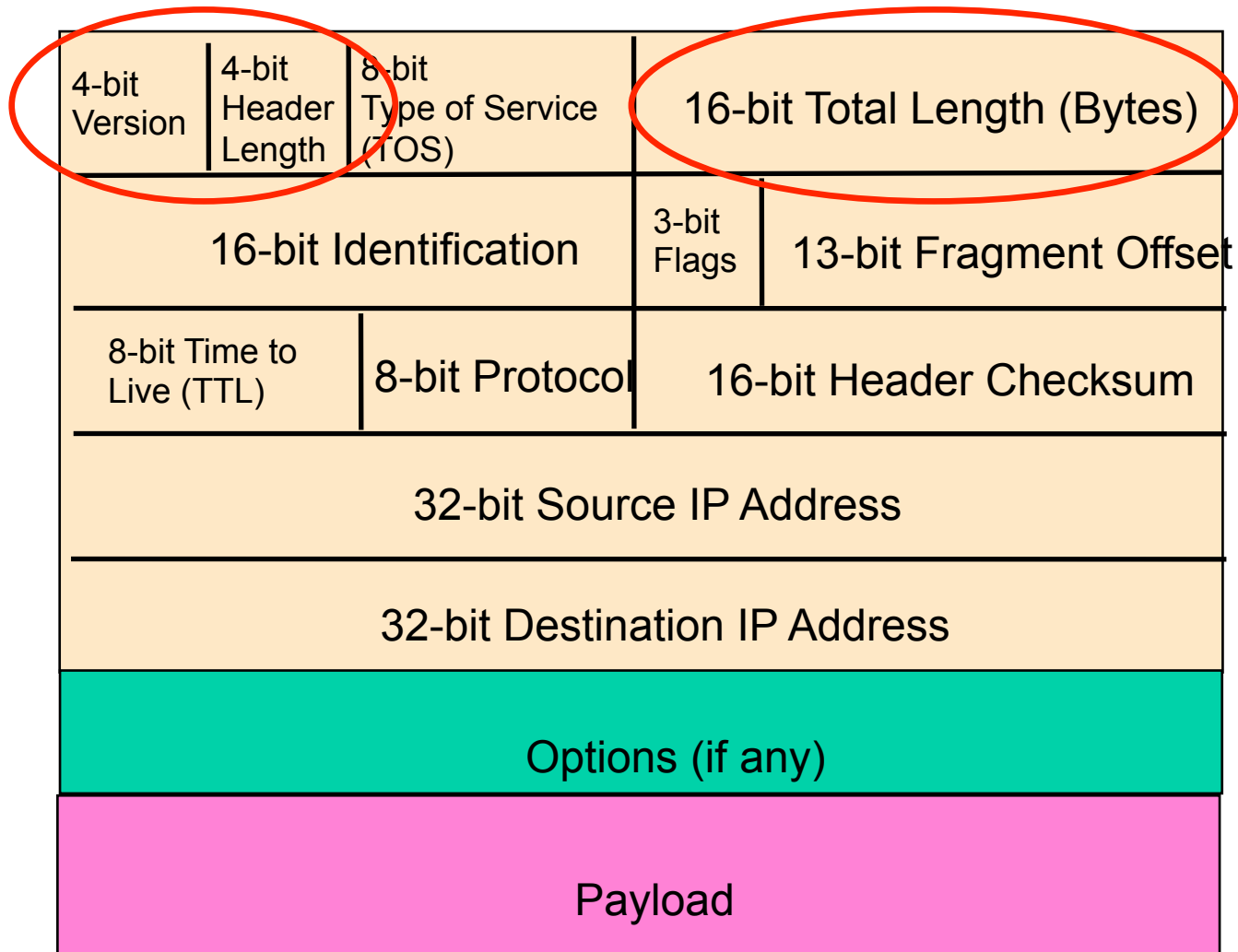❖ Deal with problems that arise along the path

# IP Packet Structure

| 4-bit Version | 4-bit Header Length | 8-bit Type of Service (TOS) | 16-bit Total Length (Bytes) | |
|---|---|---|---|---|
| 16-bit Identification | | | 3-bit Flags | 13-bit Fragment Offset |
| 8-bit Time to Live (TTL) | | 8-bit Protocol | 16-bit Header Checksum | |
| 32-bit Source IP Address | | | | |
| 32-bit Destination IP Address | | | | |
| Options (if any) | | | | |
| Payload | | | | |

# 20 Bytes of Standard Header, then Options

| 4-bit Version | 4-bit Header Length | 8-bit Type of Service (TOS) | 16-bit Total Length (Bytes) | |
|---|---|---|---|---|
| 16-bit Identification | | | 3-bit Flags | 13-bit Fragment Offset |
| 8-bit Time to Live (TTL) | | 8-bit Protocol | 16-bit Header Checksum | |
| 32-bit Source IP Address | | | | |
| 32-bit Destination IP Address | | | | |
| Options (if any) | | | | |
| Payload | | | | |

# Fields for Reading Packet Correctly

| 4-bit Version | 4-bit Header Length | 8-bit Type of Service (TOS) | 16-bit Total Length (Bytes) | |
|---|---|---|---|---|
| 16-bit Identification | | | 3-bit Flags | 13-bit Fragment Offset |
| 8-bit Time to Live (TTL) | | 8-bit Protocol | 16-bit Header Checksum | |
| 32-bit Source IP Address | | | | |
| 32-bit Destination IP Address | | | | |
| Options (if any) | | | | |
| Payload | | | | |

# Reading Packet Correctly

- ❖ **Version number (4 bits)**
  - Indicates the version of the IP protocol
  - Necessary to know what other fields to expect
  - Typically "4" (for IPv4), and sometimes "6" (for IPv6)

- ❖ **Header length (4 bits)**
  - Number of 32-bit words in the header
  - Typically "5" (for a 20-byte IPv4 header)
  - Can be more when IP options are used

- ❖ **Total length (16 bits)**
  - Number of bytes in the packet
  - Maximum size is 65,535 bytes ($2^{16}$ -1)
  - … though underlying links may impose smaller limits
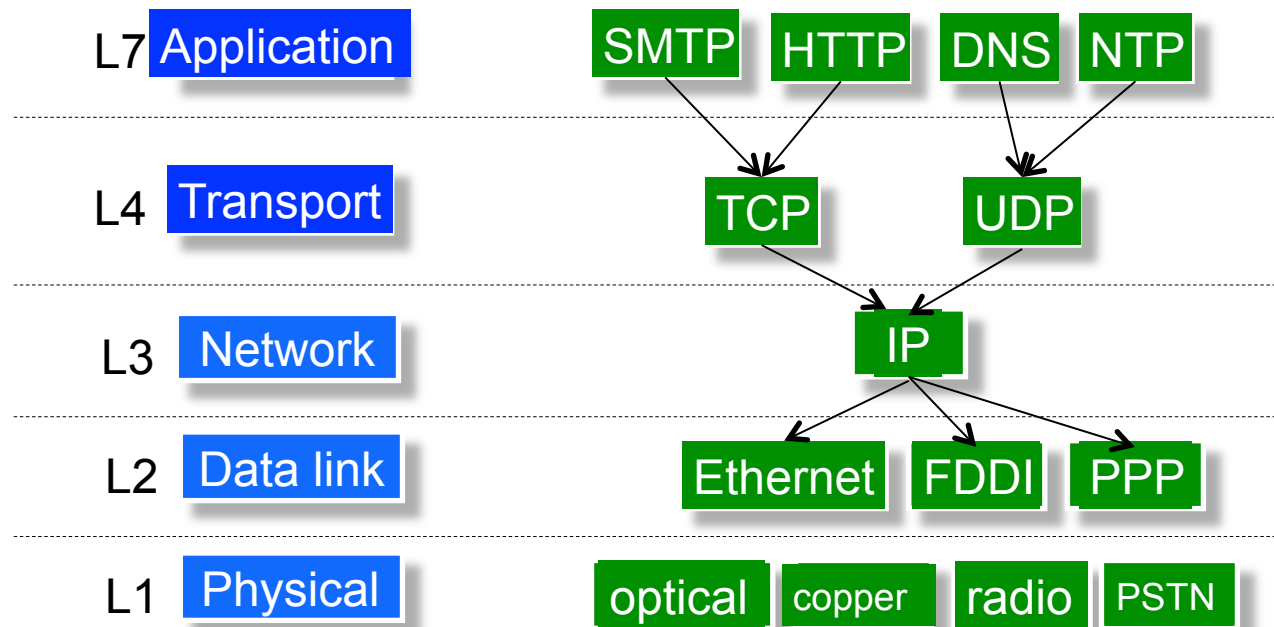
# Fields for Reaching Destination and Back

| 4-bit Version | 4-bit Header Length | 8-bit Type of Service (TOS) | 16-bit Total Length (Bytes) | |
|---|---|---|---|---|
| 16-bit Identification | | | 3-bit Flags | 13-bit Fragment Offset |
| 8-bit Time to Live (TTL) | | 8-bit Protocol | 16-bit Header Checksum | |
| 32-bit Source IP Address | | | | |
| 32-bit Destination IP Address | | | | |
| Options (if any) | | | | |
| Payload | | | | |

# Telling End-Host How to Handle Packet

| 4-bit Version | 4-bit Header Length | 8-bit Type of Service (TOS) | 16-bit Total Length (Bytes) | |
|---|---|---|---|---|
| 16-bit Identification | | | 3-bit Flags | 13-bit Fragment Offset |
| 8-bit Time to Live (TTL) | | 8-bit Protocol | 16-bit Header Checksum | |
| 32-bit Source IP Address | | | | |
| 32-bit Destination IP Address | | | | |
| Options (if any) | | | | |
| Payload | | | | |

# Telling End-Host How to Handle Packet

- ❖ Protocol (8 bits)
  - Identifies the higher-level protocol
  - Important for demultiplexing at receiving host

| | |
|---|---|
| L7 Application | SMTP  HTTP  DNS  NTP |
| L4 Transport | TCP  UDP |
| L3 Network | IP |
| L2 Data link | Ethernet  FDDI  PPP |
| L1 Physical | optical  copper  radio  PSTN |

# Telling End-Host How to Handle Packet

- ❖ Protocol (8 bits)
  - ▪ Identifies the higher-level protocol
  - ▪ Important for demultiplexing at receiving host
- ❖ Most common examples
  - ▪ E.g., "6" for the Transmission Control Protocol (TCP)
  - ▪ E.g., "17" for the User Datagram Protocol (UDP)

protocol=6

| IP header |
| TCP header |
|  |

protocol=17

| IP header |
| UDP header |
|  |

# Potential Problems

❖ Header Corrupted: **Checksum**

❖ Loop: **TTL**

❖ Packet too large: **Fragmentation**

# Checksum, TTL and Fragmentation Fields

| 4-bit Version | 4-bit Header Length | 8-bit Type of Service (TOS) | 16-bit Total Length (Bytes) | | |
|---|---|---|---|---|---|
| 16-bit Identification | | | 3-bit Flags | 13-bit Fragment Offset | |
| 8-bit Time to Live (TTL) | | 8-bit Protocol | 16-bit Header Checksum | | |
| 32-bit Source IP Address | | | | | |
| 32-bit Destination IP Address | | | | | |
| Options (if any) | | | | | |
| Payload | | | | | |

# Header Corruption (Checksum)

❖ Checksum (16 bits)
- Particular form of checksum over packet header

❖ If not correct, router discards packets
- So it doesn't act on bogus information

❖ Checksum recalculated at every router
- **Why?**
- **Why include TTL?**
- **Why only header?**

# Preventing Loops (TTL)

❖ Forwarding loops cause packets to cycle for a looong time
  ▪ As these accumulate, eventually consume **all** capacity



❖ Time-to-Live (TTL) Field  (8 bits)
  ▪ Decremented at each hop, packet discarded if reaches 0
  ▪ …and "time exceeded" message is sent to the source

# IP fragmentation, reassembly

❖ network links have MTU (max.transfer size) - largest possible link-level frame

  ▪ different link types, different MTUs

❖ large IP datagram divided ("fragmented") within net

  ▪ one datagram becomes several datagrams

  ▪ "reassembled" only at final destination

  ▪ IP header bits used to identify, order related fragments

*fragmentation:*
*in:* one large datagram
*out:* 3 smaller datagram

*reassembly*

# IP fragmentation, reassembly

*example:*

- ❖ 4000 byte datagram
- ❖ MTU = 1500 bytes

| | length =4000 | ID =x | fragflag =0 | offset =0 | |

*one large datagram becomes several smaller datagrams*

1480 bytes in data field

offset = 1480/8

| | length =1500 | ID =x | fragflag =1 | offset =0 | |

| | length =1500 | ID =x | fragflag =1 | offset =185 | |

| | length =1040 | ID =x | fragflag =0 | offset =370 | |

Applet:
http://media.pearsoncmg.com/aw/aw_kurose_network_2/applets/ip/ipfragmentation.html

# Quiz: How can we use this for evil?

A: Send fragments that overlap.

B: Send many tiny fragments, none of which have offset 0.

C: Send segments that when assembled, are bigger than the maximum IP datagram.

D: More than one of the above.

E: Nah, networks (and operating systems) are too robust for this to cause problems.

# IP Fragmentation Attacks …

## IP fragmentation exploits  [edit]

### IP fragment overlapped  [edit]

The IP fragment overlapped exploit occurs when two fragments contained within the same IP datagram have offsets that indicate that they overlap each other in positioning within the datagram. This could mean that either fragment A is being completely overwritten by fragment B, or that fragment A is partially being overwritten by fragment B. Some operating systems do not properly handle fragments that overlap in this manner and may throw exceptions or behave in other undesirable ways upon receipt of overlapping fragments. This is the basis for the teardrop Denial of service attacks.

### IP fragmentation buffer full  [edit]

The IP fragmentation buffer full exploit occurs when there is an excessive amount of incomplete fragmented traffic detected on the protected network. This could be due to an excessive number of incomplete fragmented datagrams, a large number of fragments for individual datagrams or a combination of quantity of incomplete datagrams and size/number of fragments in each datagram. This type of traffic is most likely an attempt to bypass security measures or Intrusion Detection Systems by intentional fragmentation of attack activity.

### IP fragment overrun  [edit]

The IP Fragment Overrun exploit is when a reassembled fragmented datagram exceeds the declared IP data length or the maximum datagram length. By definition, no IP datagram should be larger than 65,535 bytes. Systems that try to process these large datagrams can crash, and can be indicative of a denial of service attempt.

### IP fragment overwrite  [edit]

Overlapping fragments may be used in an attempt to bypass Intrusion Detection Systems. In this exploit, part of an attack is sent in fragments along with additional random data; future fragments may overwrite the random data with the remainder of the attack. If the completed datagram is not properly reassembled at the IDS, the attack will go undetected.

### IP fragment too many datagrams  [edit]

The Too Many Datagrams exploit is identified by an excessive number of incomplete fragmented datagrams detected on the network. This is usually either a denial of service attack or an attempt to bypass security measures. An example of "Too Many Datagrams", "Incomplete Datagram" and "Fragment Too Small" is the Rose Attack.[1]
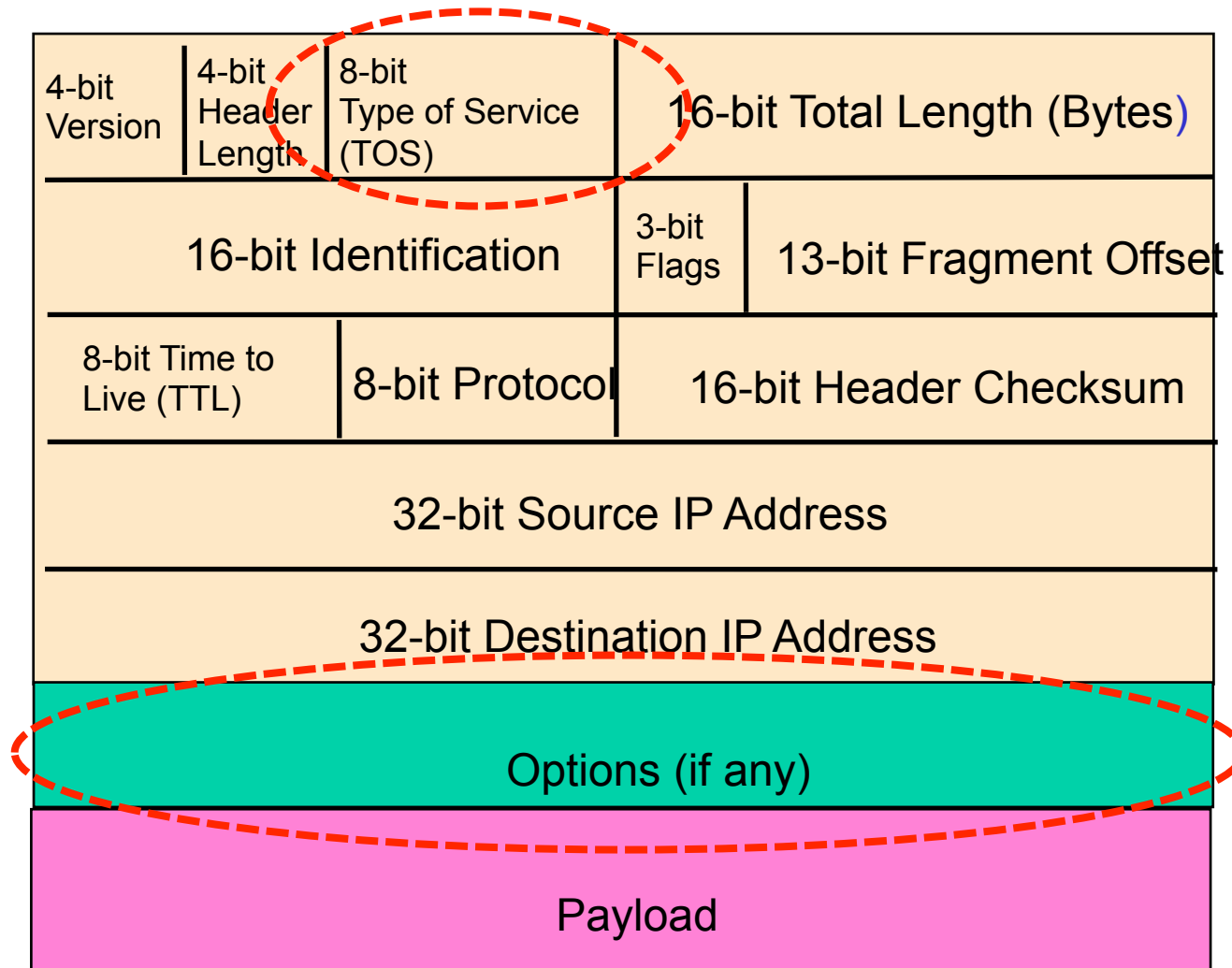
### IP fragment incomplete datagram  [edit]

This exploit occurs when a datagram can not be fully reassembled due to missing data. This can indicate a denial of service attack or an attempt to defeat packet filter security policies.

### IP fragment too small  [edit]

An IP Fragment Too Small exploit is when any fragment other than the final fragment is less than 400 bytes, indicating that the fragment is likely intentionally crafted. Small fragments may be used in denial of service attacks or in an attempt to bypass security measures or detection.
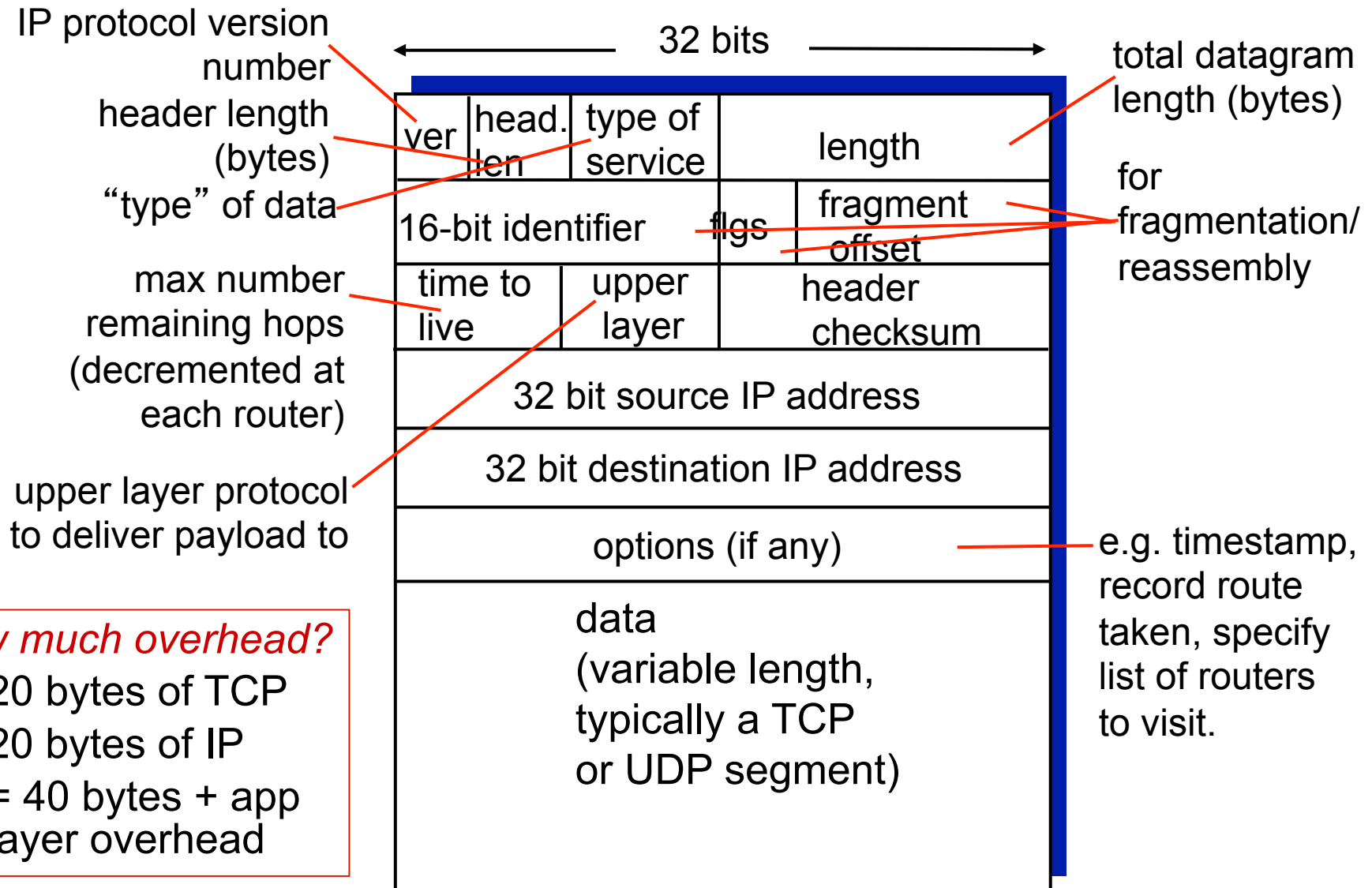
# Fields for Special Handling

| 4-bit Version | 4-bit Header Length | 8-bit Type of Service (TOS) | 16-bit Total Length (Bytes) | |
|---|---|---|---|---|
| 16-bit Identification | | | 3-bit Flags | 13-bit Fragment Offset |
| 8-bit Time to Live (TTL) | | 8-bit Protocol | 16-bit Header Checksum | |
| 32-bit Source IP Address | | | | |
| 32-bit Destination IP Address | | | | |
| Options (if any) | | | | |
| Payload | | | | |

# Special Handling

❖ "Type of Service", or "Differentiated Services Code Point (DSCP)" (8 bits)

- Allow packets to be treated differently based on needs
- E.g., low delay for audio, high bandwidth for bulk transfer
- Has been redefined several times, will cover later in class

❖ Options (not often used)

# Examples of Options

- ❖ Record Route

- ❖ Strict Source Route

- ❖ Loose Source Route

- ❖ Timestamp

- ❖ Traceroute

- ❖ Router Alert

- ❖ …..

# IP datagram format

IP protocol version number

header length (bytes)

"type" of data

max number remaining hops (decremented at each router)

upper layer protocol to deliver payload to

*how much overhead?*
- ❖ 20 bytes of TCP
- ❖ 20 bytes of IP
- ❖ = 40 bytes + app layer overhead

32 bits

| ver | head. len | type of service | length | | |
|---|---|---|---|---|---|
| 16-bit identifier | | | flgs | fragment offset | |
| time to live | | upper layer | header checksum | | |
| 32 bit source IP address | | | | | |
| 32 bit destination IP address | | | | | |
| options (if any) | | | | | |
| data (variable length, typically a TCP or UDP segment) | | | | | |

total datagram length (bytes)

for fragmentation/ reassembly

e.g. timestamp, record route taken, specify list of routers to visit.

# Network Layer: outline

4.1 introduction

4.2 virtual circuit and
datagram networks

4.3 what's inside a router

4.4 IP: Internet Protocol
- datagram format
- IPv4 addressing
- ICMP
- IPv6

4.5 routing algorithms
- link state
- distance vector
- hierarchical routing

4.6 routing in the Internet
- RIP
- OSPF
- BGP

4.7 broadcast and multicast
routing

# IP addressing: introduction

- *IP address:* **32**-bit identifier for host, router *interface*

- *interface:* connection between host/router and physical link
    - router's typically have multiple interfaces
    - host typically has one or two interfaces (e.g., wired Ethernet, wireless 802.11)

- *IP addresses associated with each interface*

223.1.1.1

223.1.2.1

223.1.1.2

223.1.1.4    223.1.2.9

223.1.3.27

223.1.1.3

223.1.2.2

223.1.3.1    223.1.3.2

223.1.1.1 = 11011111 00000001 00000001 000000

     223      1      1      1

# IP addressing: introduction

*Q: how are interfaces actually connected?*

*A: we'll learn about that in the link layer*

*A:* wired Ethernet interfaces connected by Ethernet switches

*For now:* don't need to worry about how one interface is connected to another (with no intervening router)

223.1.1.1

223.1.2.1

223.1.1.2

223.1.1.4    223.1.2.9

223.1.1.3

223.1.3.27

223.1.2.2

223.1.3.1    223.1.3.2

*A:* wireless WiFi interfaces connected by WiFi base station

# Subnets

❖ **IP address:**
  - subnet part - high order bits
  - host part - low order bits

❖ *what's a subnet ?*
  - device interfaces with same subnet part of IP address
  - can physically reach each other *without intervening router*



223.1.1.1

223.1.1.2

223.1.1.4 223.1.2.9

223.1.2.1

223.1.2.2

223.1.1.3 223.1.3.27

subnet

223.1.3.1 223.1.3.2

network consisting of 3 subnets

# Subnets

*recipe*

❖ to determine the subnets, detach each interface from its host or router, creating islands of isolated networks

❖ each isolated network is called a *subnet*

223.1.1.0/24

223.1.2.0/24

223.1.1.1

223.1.1.2

223.1.1.4  223.1.2.9

223.1.2.1

223.1.2.2

223.1.1.3  223.1.3.27

subnet

223.1.3.1

223.1.3.2

223.1.3.0/24

subnet mask: /24

# Subnets

how many?

223.1.1.2

223.1.1.1

223.1.1.4

223.1.1.3

223.1.9.2        223.1.7.0

223.1.9.1                                223.1.7.1

223.1.8.1        223.1.8.0

223.1.2.6                                223.1.3.27

223.1.2.1        223.1.2.2    223.1.3.1        223.1.3.2

# Original Internet Addresses

❖ First eight bits: network address (/8)
❖ Last 24 bits: host address

*Assumed 256 networks were more than enough!*

# Next Design: "Classful" Addressing

❖ Three main classes

■ Class A

| 0 | network | host |
|---|---------|------|

0 ... 8

{ 126 nets
~16M hosts

■ Class B

| 1 | 0 | network | host |
|---|---|---------|------|

0 ... 16

{ ~16K nets
~65K hosts

■ Class C

| 1 | 1 | 0 | network | host |
|---|---|---|---------|------|

0 ... 24

{ ~2M nets
254 hosts

**Problem: Networks only come in three sizes!**

# Today's addressing: CIDR

CIDR: Classless InterDomain Routing
- subnet portion of address of arbitrary length
- address format: a.b.c.d/x, where x is # bits in subnet portion of address

subnet part ← → host part

11001000 00010111 00010000 00000000

200.23.16.0/23

200.23.17.255

/23

200.23.16.0

IP address

# Subnet Address

❖ **Subnet Mask**
  - Used in conjunction to with the network address to indicate how many higher order bits are used for the network part of the address (i.e. network prefix)
    - Bit-wise AND
  - 223.1.1.0/24 is equivalent to 223.1.1.0 with subnet mask 255.255.255.0

❖ **Broadcast Address**
  - host part is all 111's
  - E.g. 223.1.3.255

223.1.1.0/24

223.1.2.0/24

B: 223.1.1.2

223.1.3.0/24

| Host B | Dot-decimal address | Binary |
|---|---|---|
| IP address | 223.1.1.2 | 11111101.00000001.00000001.00000010 |
| Subnet Mask | 255.255.255.0 | 11111111.11111111.11111111.00000000 |
| Network Part | 223.1.1.0 | 11111101.00000001.00000001.00000000 |
| Host Part | 0.0.0.2 | 00000000.00000000.00000000.00000010 |

# IP addresses: how to get one?

Q: How does a *host* get IP address?

❖ hard-coded by system admin in a file
  ▪ Windows: control-panel->network->configuration->tcp/ip->properties
  ▪ UNIX: /etc/rc.config

❖ DHCP: Dynamic Host Configuration Protocol: dynamically get address from as server
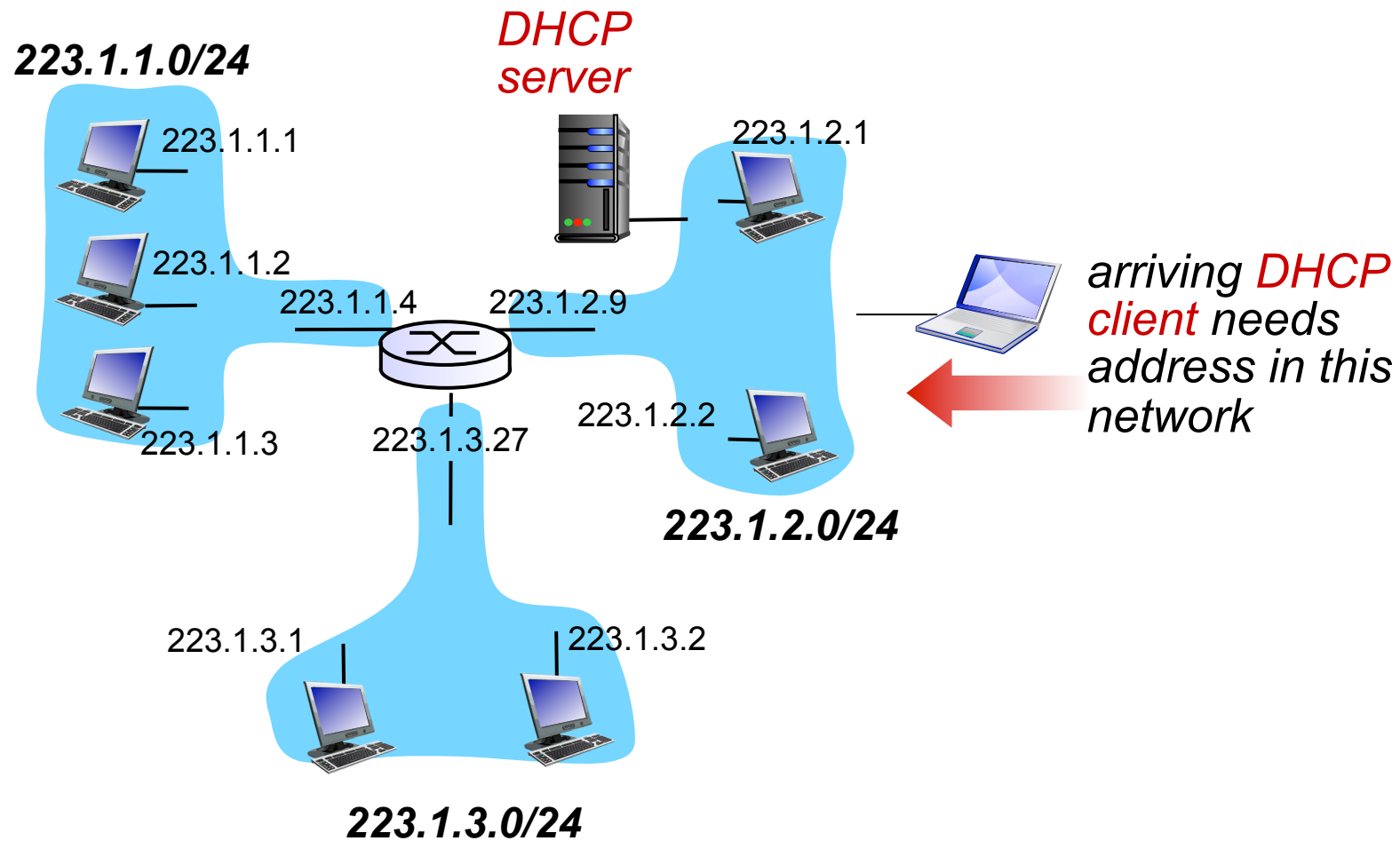  ▪ "plug-and-play"

# DHCP: Dynamic Host Configuration Protocol

*goal:* allow host to *dynamically* obtain its IP address from network server when it joins network

- can renew its lease on address in use
- allows reuse of addresses (only hold address while connected/"on")
- support for mobile users who want to join network (more shortly)

*DHCP overview:*

- host broadcasts "DHCP discover" msg [optional]
- DHCP server responds with "DHCP offer" msg [optional]
- host requests IP address: "DHCP request" msg
- DHCP server sends address: "DHCP ack" msg

# DHCP client-server scenario

223.1.1.0/24

DHCP server

223.1.1.1

223.1.2.1

223.1.1.2

223.1.1.4    223.1.2.9

arriving DHCP client needs address in this network

223.1.1.3    223.1.3.27    223.1.2.2

223.1.2.0/24

223.1.3.1    223.1.3.2

223.1.3.0/24

# DHCP client-server scenario

DHCP server: 223.1.2.5

**DHCP discover**

arriving client

```
src : 0.0.0.0, 68
dest.: 255.255.255.255,67
yiaddr:    0.0.0.0
transaction ID: 654
```

**DHCP offer**

```
src: 223.1.2.5, 67
dest:  255.255.255.255, 68
yiaddrr: 223.1.2.4
transaction ID: 654
lifetime: 3600 secs
```

**DHCP request**

```
src:  0.0.0.0, 68
dest:: 255.255.255.255, 67
yiaddrr: 223.1.2.4
transaction ID: 655
lifetime: 3600 secs
```

**DHCP ACK**

```
src: 223.1.2.5, 67
dest:  255.255.255.255, 68
yiaddrr: 223.1.2.4
transaction ID: 655
lifetime: 3600 secs
```

# DHCP: more than IP addresses

DHCP can return more than just allocated IP address on subnet:
- address of first-hop router for client
- name and IP address of DNS sever
- network mask (indicating network versus host portion of address)

# DHCP: example



router with DHCP
server built into
router

- ❖ **connecting laptop needs its IP address, addr of first-hop router, addr of DNS server: use DHCP**

- ❖ **DHCP request encapsulated in UDP, encapsulated in IP, encapsulated in 802.1 Ethernet**

- ❖ **Ethernet frame broadcast (dest: FFFFFFFFFFFF) on LAN, received at router running DHCP server**

- ❖ **Ethernet demuxed to IP demuxed, UDP demuxed to DHCP**

# DHCP: example



router with DHCP server built into router

- ❖ DCP server formulates DHCP ACK containing client's IP address, IP address of first-hop router for client, name & IP address of DNS server

- ❖ encapsulation of DHCP server, frame forwarded to client, demuxing up to DHCP at client

- ❖ client now knows its IP address, name and IP address of DSN server, IP address of its first-hop router

# DHCP: Wireshark output (home LAN)

## request

Message type: **<u>Boot Request (1)</u>**
Hardware type: Ethernet
Hardware address length: 6
Hops: 0
**Transaction ID: 0x6b3a11b7**
Seconds elapsed: 0
Bootp flags: 0x0000 (Unicast)
Client IP address: 0.0.0.0 (0.0.0.0)
Your (client) IP address: 0.0.0.0 (0.0.0.0)
Next server IP address: 0.0.0.0 (0.0.0.0)
Relay agent IP address: 0.0.0.0 (0.0.0.0)
**Client MAC address: Wistron_23:68:8a (00:16:d3:23:68:8a)**
Server host name not given
Boot file name not given
Magic cookie: (OK)
Option: (t=53,l=1) **DHCP Message Type = DHCP Request**
Option: (61) Client identifier
    Length: 7; Value: 010016D323688A;
    Hardware type: Ethernet
    Client MAC address: Wistron_23:68:8a (00:16:d3:23:68:8a)
Option: (t=50,l=4) Requested IP Address = 192.168.1.101
Option: (t=12,l=5) Host Name = "nomad"
**Option: (55) Parameter Request List**
    Length: 11; Value: 010F03062C2E2F1F21F92B
    **1 = Subnet Mask; 15 = Domain Name**
    **3 = Router; 6 = Domain Name Server**
    44 = NetBIOS over TCP/IP Name Server
    ……

Message type: **Boot Reply (2)**
Hardware type: Ethernet
Hardware address length: 6
Hops: 0
**Transaction ID: 0x6b3a11b7**
Seconds elapsed: 0
Bootp flags: 0x0000 (Unicast)
**Client IP address: 192.168.1.101 (192.168.1.101)**
Your (client) IP address: 0.0.0.0 (0.0.0.0)
**Next server IP address: 192.168.1.1 (192.168.1.1)**
Relay agent IP address: 0.0.0.0 (0.0.0.0)
Client MAC address: Wistron_23:68:8a (00:16:d3:23:68:8a)
Server host name not given
Boot file name not given
Magic cookie: (OK)
**Option: (t=53,l=1) DHCP Message Type = DHCP ACK**
**Option: (t=54,l=4) Server Identifier = 192.168.1.1**
**Option: (t=1,l=4) Subnet Mask = 255.255.255.0**
**Option: (t=3,l=4) Router = 192.168.1.1**
**Option: (6) Domain Name Server**
    **Length: 12; Value: 445747E2445749F244574092;**
    **IP Address: 68.87.71.226;**
    **IP Address: 68.87.73.242;**
    **IP Address: 68.87.64.146**
**Option: (t=15,l=20) Domain Name = "hsd1.ma.comcast.net."**

# DHCP: further details

❖ DHCP uses UDP and port numbers 67 (server side) and 68 (client side)

❖ Usually the MAC address is used to identify clients
  ▪ DHCP server can be configured with a "registered list" of acceptable MAC addresses

❖ DHCP offer message includes ip address, length of lease, subnet mask, DNS servers, default gateway

❖ DHCP security holes
  ▪ DoS attack by exhausting pool of IP addresses
  ▪ Masquerading as a DHCP server
  ▪ Authentication for DHCP - RFC 3118

# IP addresses: how to get one?

*Q:* how does *network* get subnet part of IP addr?

*A:* gets allocated portion of its provider ISP's address space

| ISP's block | 11001000  00010111  00010000 | 00000000 | 200.23.16.0/20 |
|---|---|---|---|
| Organization 0 | 11001000  00010111  00010000 | 00000000 | 200.23.16.0/23 |
| Organization 1 | 11001000  00010111  00010010 | 00000000 | 200.23.18.0/23 |
| Organization 2 | 11001000  00010111  00010100 | 00000000 | 200.23.20.0/23 |
| ... | ….. | …. | …. |
| Organization 7 | 11001000  00010111  00011110 | 00000000 | 200.23.30.0/23 |

# CIDR: Addresses allocated in contiguous prefix chunks

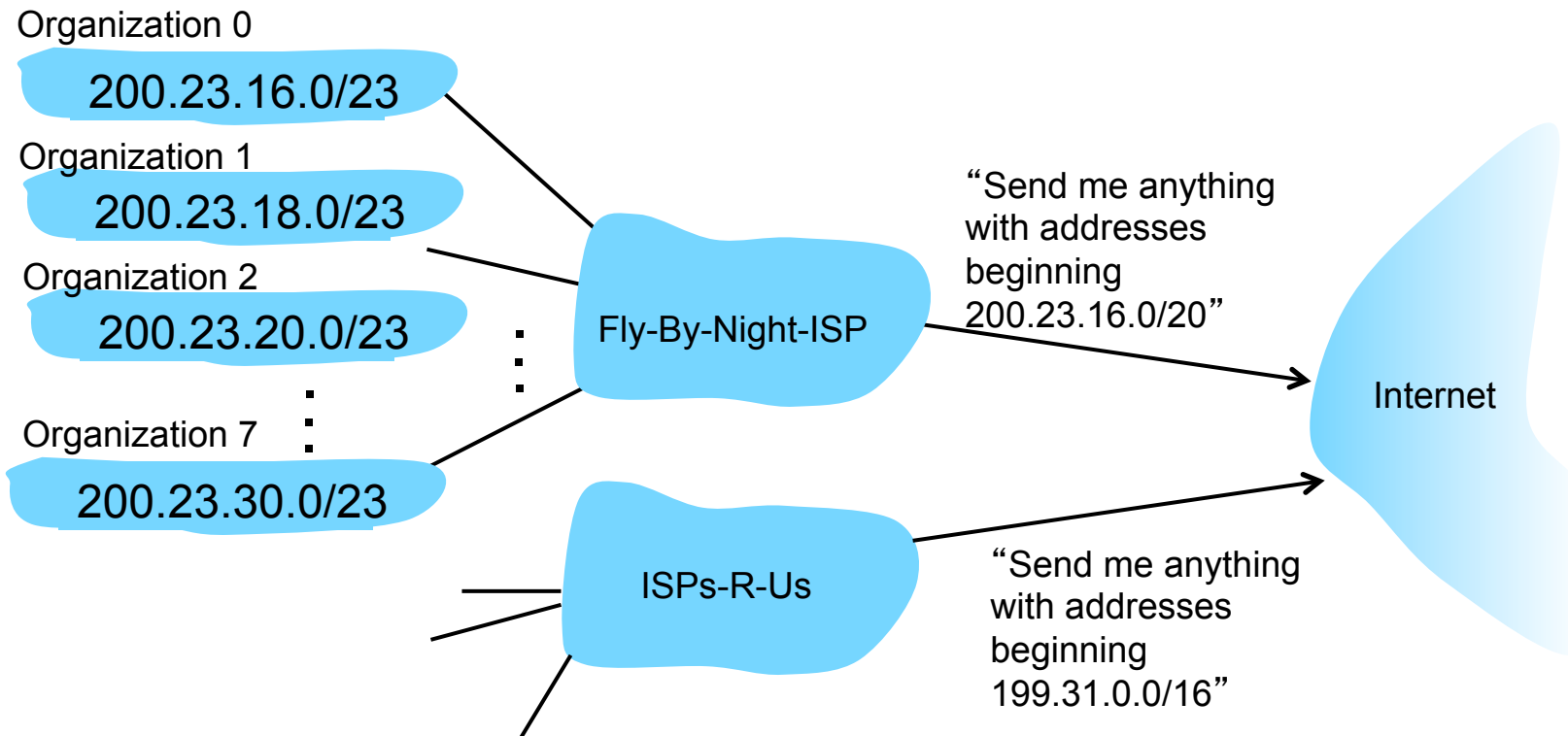Recursively break down chunks as get closer to host

12.0.0.0/8

12.0.0.0/15
12.2.0.0/16
12.3.0.0/16

⋮
⋮
⋮
⋮

12.253.0.0/16
⋮

12.3.0.0/22
12.3.4.0/24
⋮
⋮
12.3.254.0/23

⋮
⋮
⋮

12.253.0.0/19
12.253.32.0/19
12.253.64.0/19
12.253.64.108/30
12.253.96.0/18
12.253.128.0/17

# Hierarchical addressing: route aggregation

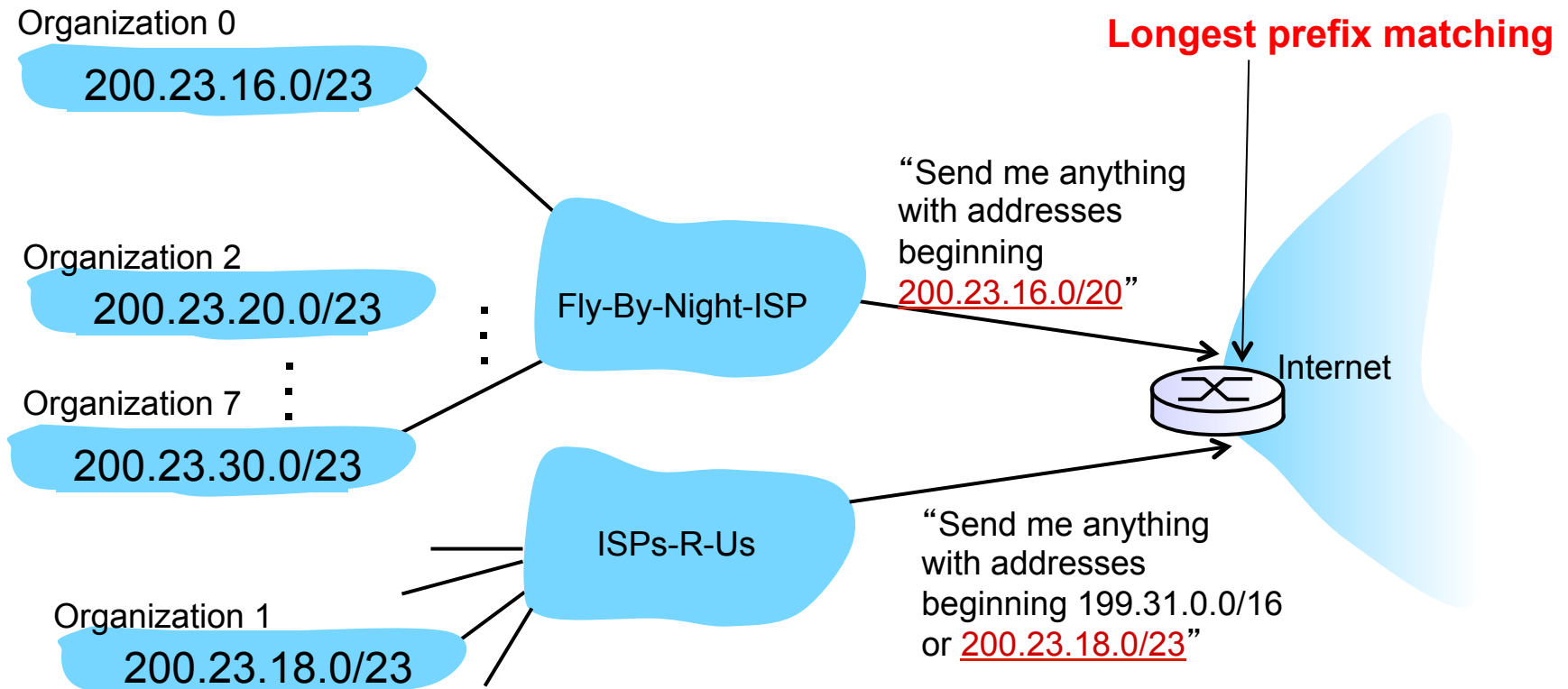hierarchical addressing allows efficient advertisement of routing information:

Organization 0

200.23.16.0/23

Organization 1

200.23.18.0/23

Organization 2

200.23.20.0/23

Organization 7

200.23.30.0/23

Fly-By-Night-ISP

ISPs-R-Us

"Send me anything with addresses beginning 200.23.16.0/20"

"Send me anything with addresses beginning 199.31.0.0/16"

Internet

# Quiz: What should we do if organization 1 decides to switch to ISPs-R-Us

Organization 0
200.23.16.0/23

Organization 1
200.23.18.0/23

Organization 2
200.23.20.0/23

Organization 7
200.23.30.0/23

Fly-By-Night-ISP

ISPs-R-Us

"Send me anything with addresses beginning 200.23.16.0/20"

"Send me anything with addresses beginning 199.31.0.0/16"

Internet

A: Move 200.23.18.0/23 to ISPs-R-Us (and break up Fly-By-Night's/20 block).
B: Give new addresses to Organization 1 (and force them to change all their addresses)
C: Some other solution

# Hierarchical addressing: more specific routes

ISPs-R-Us has a more specific route to Organization 1

Organization 0
200.23.16.0/23

**Longest prefix matching**

"Send me anything with addresses beginning 200.23.16.0/20"

Organization 2
200.23.20.0/23

Fly-By-Night-ISP

Organization 7
200.23.30.0/23

Internet

ISPs-R-Us

"Send me anything with addresses beginning 199.31.0.0/16 or 200.23.18.0/23"

Organization 1
200.23.18.0/23

# Example: continued

❖ But how will this work?

❖ Routers in the Internet will have two entries in their tables
   ▪ 200.23.16.0/20 (Fly-by-Night-ISP)
   ▪ 200.23.18.0/23 (ISPs-R-Us)

❖ Longest prefix match

200.23.31.255

/20

200.23.19.255

/23

200.23.18.0

200.23.16.0

IP address

White Paper on IP addresses linked to page - Very informative

# More on IP addresses

❖ IP addresses are allocated as blocks and have geographical significance

❖ It is possible to determine the geographical location of an IP address

http://www.geobytes.com/IpLocator.htm



MAP OF THE INTERNET
THE IPv4 SPACE, 2006

THIS CHART SHOWS THE IP ADDRESS SPACE ON A PLANE USING A FRACTAL MAPPING WHICH PRESERVES GROUPING -- ANY CONSECUTIVE STRING OF IPs WILL TRANSLATE TO A SINGLE COMPACT, CONTIGUOUS REGION ON THE MAP. EACH OF THE 256 NUMBERED BLOCKS REPRESENTS ONE /8 SUBNET (CONTAINING ALL IPs THAT START WITH THAT NUMBER). THE UPPER LEFT SECTION SHOWS THE BLOCKS SOLD DIRECTLY TO CORPORATIONS AND GOVERNMENTS IN THE 1990's BEFORE THE RIRs TOOK OVER ALLOCATION.

= UNALLOCATED BLOCK

Network Layer  52

# IP addressing: the last word...

*Q:* how does an ISP get block of addresses?

*A:* ICANN: Internet Corporation for Assigned

Names and Numbers http://www.icann.org/

- allocates addresses
- manages DNS
- assigns domain names, resolves disputes

❖ Regional Internet Registries (RIR) act as intermediaries

- RIPE NCC (Riseaux IP Europiens Network Coordination Center) for Europe, Middle East, Africa
- APNIC (Asia Pacific Network Information Center) for Asia and Pacific
- ARIN (American Registry for Interent Numbers) for the Americas, Caribbean, sub-Saharan Africa
- LACNIC (Latin America and Caribbean)

# Made-up Example in More Detail

- ❖ ICANN gives APNIC several /8s
- ❖ APNIC gives Telstra one /8, **129.0/8**
  - ▪ Network Prefix: **10000001**
- ❖ Telstra gives UNSW a /16, **129.94/16**
  - ▪ Network Prefix: **1000000101011110**
- ❖ UNSW gives CSE a /24, **12.197.242/24**
  - ▪ Network Prefix: **100000010101111011110010**
- ❖ CSE gives me a specific address **129.94.242.51**
  - ▪ Address: **10000001010111101111001000110011**

# Quiz: Header Fields

❖ Which of the following fields is not part of either a TCP or UDP header?

A. Source port

B. Source IP address

C. Receive window

D. Length

E. Checksum

# Quiz: DHCP

❖ What transport protocol does DHCP use?
A. UDP
B. TCP
C. IP
D. HTTP

# Quiz: IP Addressing

❖ How many IP addresses belong to the subnet 128.119.254.0/25 ? What are the IP addresses at the two end-points of this range ?

# Quiz: Subnets

❖ How many subnets are there in this network?

# Quiz: Subnets

❖ The two subnets 128.119.245.129/25 and 128.119.245.4/26 have overlapping IP addresses.

A. True
B. False

# Recall IPv4 addresses scarcity

## ARIN Finally Runs Out of IPv4 Addresses

IPv4 Address Cupboards are Bare in North America.

Network World | Sep 22, 2015 7:25 AM PT

It is often said, "the Internet is running out of phone numbers," as a way to express that the Internet is running out of IPv4 addresses, to those who are unfamiliar with Internet technologies. IPv4 addresses, like phone numbers are assigned hierarchically, and thus, have inherent inefficiency. The world's Internet population has been growing and the number of Internet-connected devices continues to rise, with no end in sight. In the next week, the American Registry for Internet Numbers (ARIN) will have exhausted their supply of IPv4 addresses. The metaphorical IPv4 cupboards are bare. This long-predicted Internet historical event marks opening a new chapter of the Internet's evolution. However, it is somehow anti-climactic now that this date has
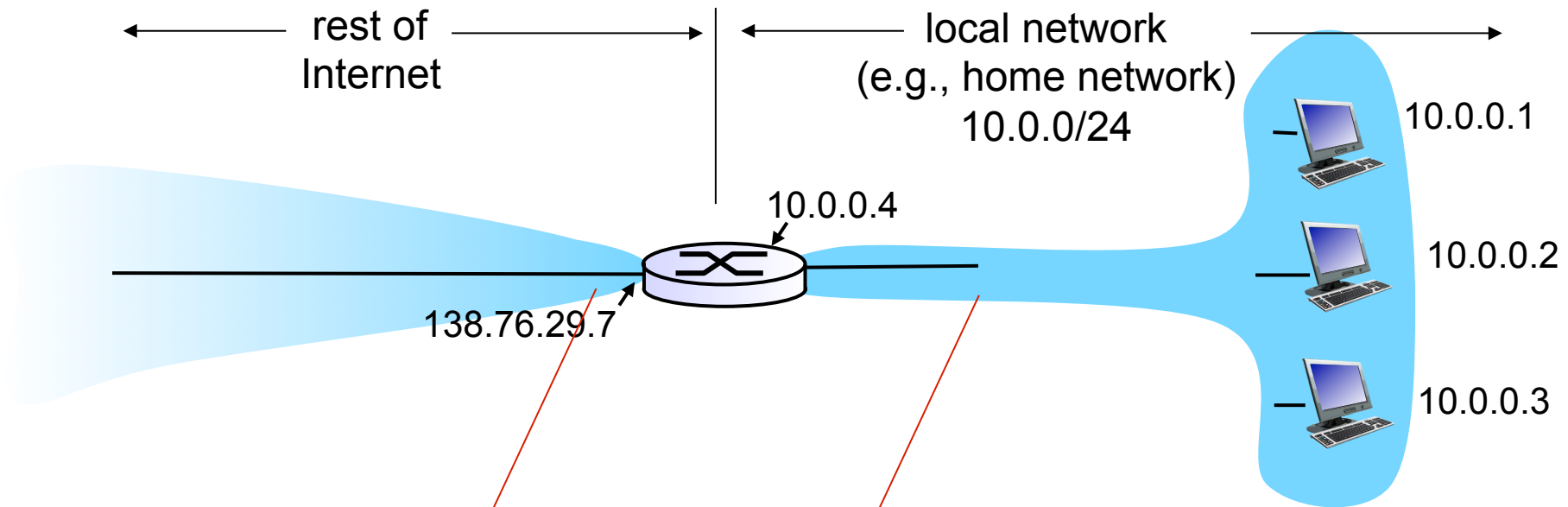
# Private Addresses

❖ Defined in RFC 1918:
- 10.0.0./8 (16,777,216 hosts)
- 172.16.0.0/12 (1,048,576 hosts)
- 192.168.0.0/16 (65536 hosts)

❖ These addresses cannot be routed
- Anyone can use them
- Often used for NAT

❖

# NAT: network address translation



rest of Internet ←→ local network (e.g., home network) 10.0.0/24

10.0.0.1
10.0.0.2
10.0.0.3

10.0.0.4

138.76.29.7

*all* datagrams *leaving* local network have *same* single source NAT IP address: 138.76.29.7,different source port numbers
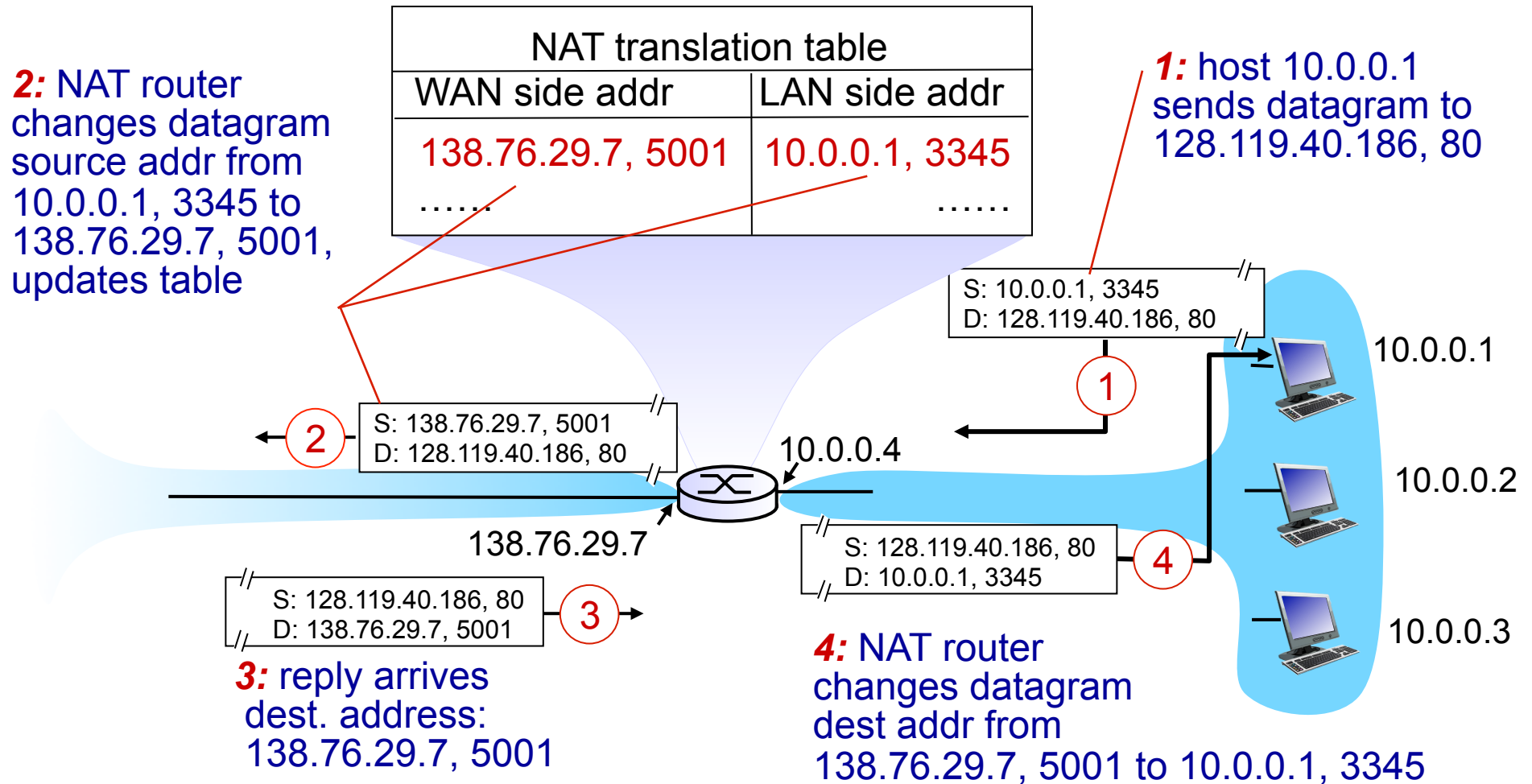
datagrams with source or destination in this network have 10.0.0/24 address for source, destination (as usual)

# NAT: network address translation

*implementation*: NAT router must:

- *outgoing datagrams: replace* (source IP address, port #) of every outgoing datagram to (NAT IP address, new port #)

  . . . remote clients/servers will respond using (NAT IP address, new port #) as destination addr

- *remember (in NAT translation table)* every (source IP address, port #)  to (NAT IP address, new port #) translation pair

- *incoming datagrams: replace* (NAT IP address, new port #) in dest fields of every incoming datagram with corresponding (source IP address, port #) stored in NAT table

# NAT: network address translation

**2:** NAT router changes datagram source addr from 10.0.0.1, 3345 to 138.76.29.7, 5001, updates table

| NAT translation table | |
|---|---|
| WAN side addr | LAN side addr |
| 138.76.29.7, 5001 | 10.0.0.1, 3345 |
| …… | …… |

**1:** host 10.0.0.1 sends datagram to 128.119.40.186, 80

S: 10.0.0.1, 3345
D: 128.119.40.186, 80

①

10.0.0.1

S: 138.76.29.7, 5001
D: 128.119.40.186, 80

②

10.0.0.4

10.0.0.2

138.76.29.7

S: 128.119.40.186, 80
D: 10.0.0.1, 3345

④

S: 128.119.40.186, 80
D: 138.76.29.7, 5001

③

10.0.0.3

**3:** reply arrives dest. address: 138.76.29.7, 5001

**4:** NAT router changes datagram dest addr from 138.76.29.7, 5001 to 10.0.0.1, 3345

# NAT Advantages

Local network uses just one IP address as far as outside world is concerned:

- range of addresses not needed from ISP: just one IP address for all devices

- can change addresses of devices in local network without notifying outside world

- can change ISP without changing addresses of devices in local network

# Quiz: NAT

❖ Devices inside the local network are not explicitly addressable or visible by outside world.
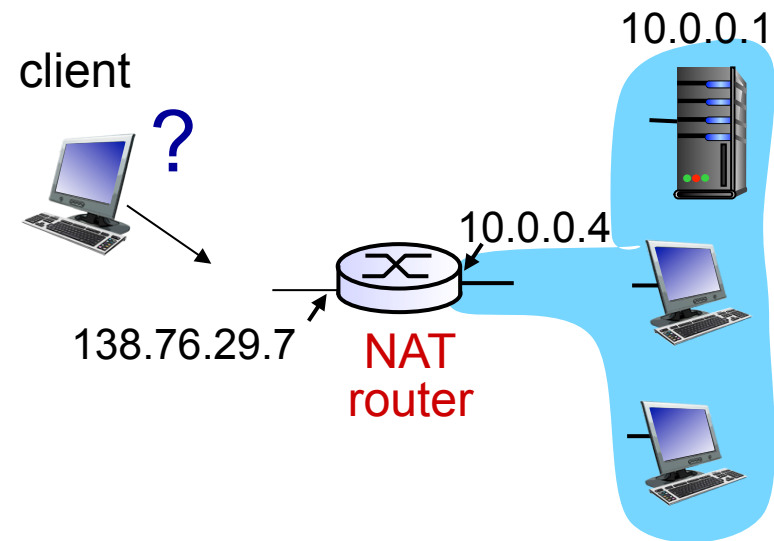
A: This is an advantage

B: This is a disadvantage

# NAT: network address translation

❖ 16-bit port-number field:

- 60,000 simultaneous connections with a single LAN-side address!

❖ NAT is controversial:

- routers should only process up to layer 3

- violates end-to-end argument

    • NAT possibility must be taken into account by app designers, e.g., P2P applications

- address shortage should instead be solved by IPv6

# NAT: Practical Issues

❖ NAT modifies port # and IP address
  ▪ *Requires recalculation of TCP and IP checksum*

❖ Some applications embed IP address or port numbers in their message payloads
  ▪ DNS, FTP (PORT command), SIP, H.323
  ▪ For legacy protocols, NAT must look into these packets and translate the embedded IP addresses/port numbers
  ▪ Duh, What if these fields are encrypted ?? (SSL/TLS, IPSEC, etc)
  ▪ **Q: In some cases why may NAT need to change TCP sequence number??**

❖ If applications change port numbers periodically, the NAT must be aware of this

❖ NAT Traversal Problems
  ▪ E.g: How to setup a server behind a NAT router?
  ▪ How to talk to a Skype user behind a NAT router?
  ▪ Possible workarounds in next few slides
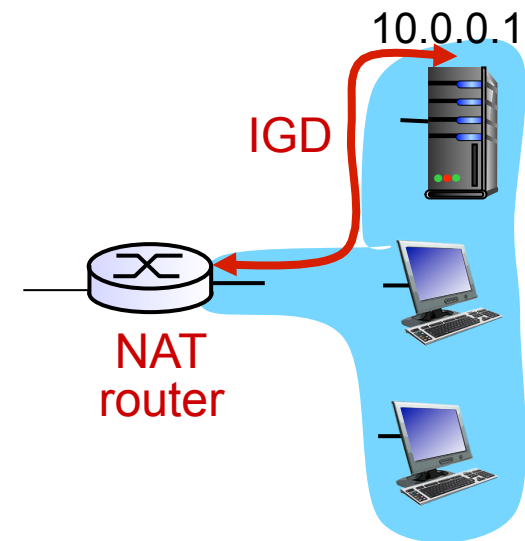
# NAT traversal problem

- ❖ **client wants to connect to server with address 10.0.0.1**
  - ▪ server address 10.0.0.1 local to LAN (client can't use it as destination addr)
  - ▪ only one externally visible NATed address: 138.76.29.7
- ❖ *solution1:* statically configure NAT to forward incoming connection requests at given port to server
  - ▪ e.g., (123.76.29.7, port 2500) always forwarded to 10.0.0.1 port 25000

client

?

138.76.29.7

NAT
router

10.0.0.1

10.0.0.4

# NAT traversal problem
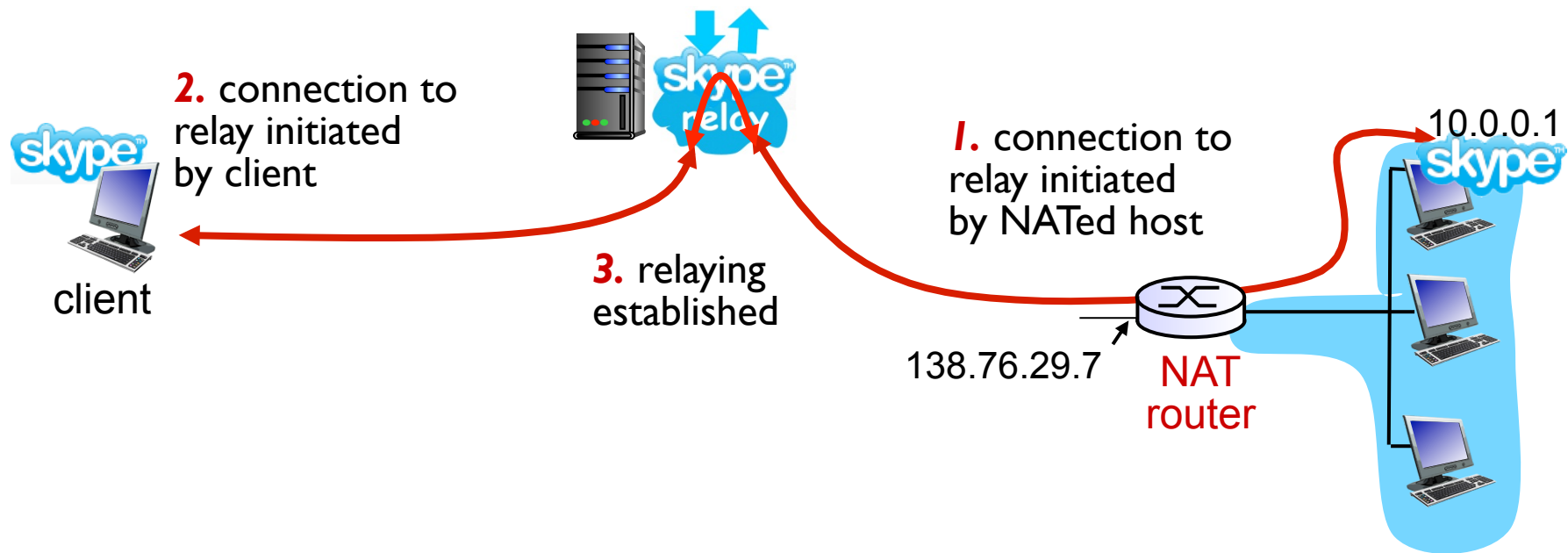
❖ *solution 2:* Universal Plug and Play (UPnP) Internet Gateway Device (IGD) Protocol.  Allows NATed host to:
  ❖ learn public IP address (138.76.29.7)
  ❖ add/remove port mappings (with lease times)

  i.e., automate static NAT port map configuration

10.0.0.1

IGD

NAT router

# NAT traversal problem

❖ *solution 3:* relaying (used in Skype)
   ■ NATed client establishes connection to relay
   ■ external client connects to relay
   ■ relay bridges packets between to connections

**2.** connection to relay initiated by client

**1.** connection to relay initiated by NATed host

10.0.0.1

**3.** relaying established

client

138.76.29.7

NAT router

# NAT: Devil in the details

❖ Despite the problems, NAT has been widely deployed

❖ Most protocols can be successfully passed through a NAT, including VPN

❖ Modern hardware can easily perform NAT functions at > 100 Mbps

❖ IPv6 is still not widely deployed commercially, so the need for NAT is real

❖ After years of refusing to work on NAT, the IETF has been developing "NAT control protocols" for hosts

❖ Lot of practical variations

▪ Full-cone NAT, Restricted Cone NAT, Port Restricted Cone NAT, Symmetric NAT, …..

• The devil is in the detail

❖ External link under lecture notes for further reading (not examinable)

# Network Layer: outline

4.1 introduction

4.2 virtual circuit and
    datagram networks

4.3 what's inside a router

4.4 IP: Internet Protocol
- datagram format
- IPv4 addressing
- ICMP
- IPv6

4.5 routing algorithms
- link state
- distance vector
- hierarchical routing

4.6 routing in the Internet
- RIP
- OSPF
- BGP

4.7 broadcast and multicast
    routing

# ICMP: internet control message protocol

- ❖ **used by hosts & routers to communicate network-level information**
  - error reporting: unreachable host, network, port, protocol
  - echo request/reply (used by ping)
- ❖ **network-layer "above" IP:**
  - ICMP msgs carried in IP datagrams
- ❖ **ICMP message:** type, code plus first 8 bytes of IP datagram causing error

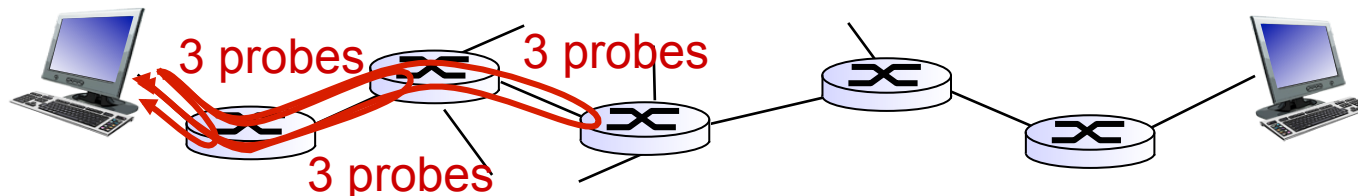| Type | Code | description |
|------|------|-------------|
| 0 | 0 | echo reply (ping) |
| 3 | 0 | dest. network unreachable |
| 3 | 1 | dest host unreachable |
| 3 | 2 | dest protocol unreachable |
| 3 | 3 | dest port unreachable |
| 3 | 6 | dest network unknown |
| 3 | 7 | dest host unknown |
| 4 | 0 | source quench (congestion control - not used) |
| 8 | 0 | echo request (ping) |
| 9 | 0 | route advertisement |
| 10 | 0 | router discovery |
| 11 | 0 | TTL expired |
| 12 | 0 | bad IP header |

# Traceroute and ICMP

- source sends series of UDP segments to dest
  - first set has TTL =1
  - second set has TTL=2, etc.
  - unlikely port number
- when *n*th set of datagrams arrives to nth router:
  - router discards datagrams
  - and sends source ICMP messages (type 11, code 0)
  - ICMP messages includes name of router & IP address

- when ICMP messages arrives, source records RTTs

*stopping criteria:*

- UDP segment eventually arrives at destination host
- destination returns ICMP "port unreachable" message (type 3, code 3)
- source stops

3 probes    3 probes

3 probes

# IPv6: motivation

❖ *initial motivation:* **32-bit address space soon to be completely allocated.**

❖ additional motivation:
  ▪ header format helps speed processing/forwarding
  ▪ header changes to facilitate QoS

*IPv6 datagram format:*
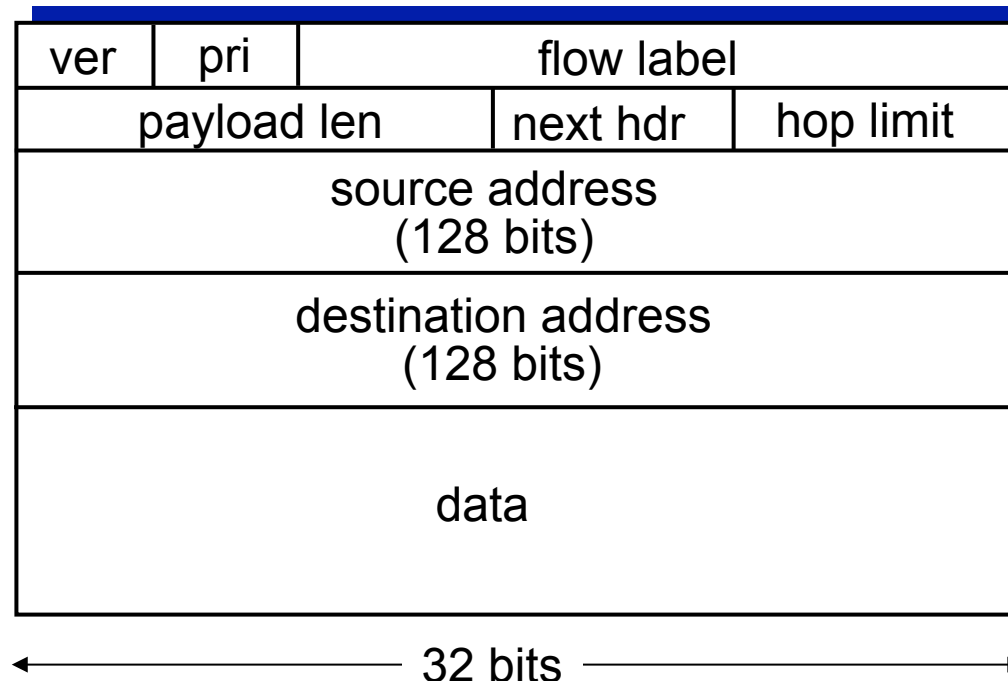  ▪ fixed-length 40 byte header
  ▪ no fragmentation allowed

https://www.google.com/intl/en/ipv6/statistics.html

# IPv6 datagram format

*priority:* identify priority among datagrams in flow (traffic class)
*flow Label:* identify datagrams in same "flow."
          (concept of "flow" not well defined).
*next header:* identify upper layer protocol for data

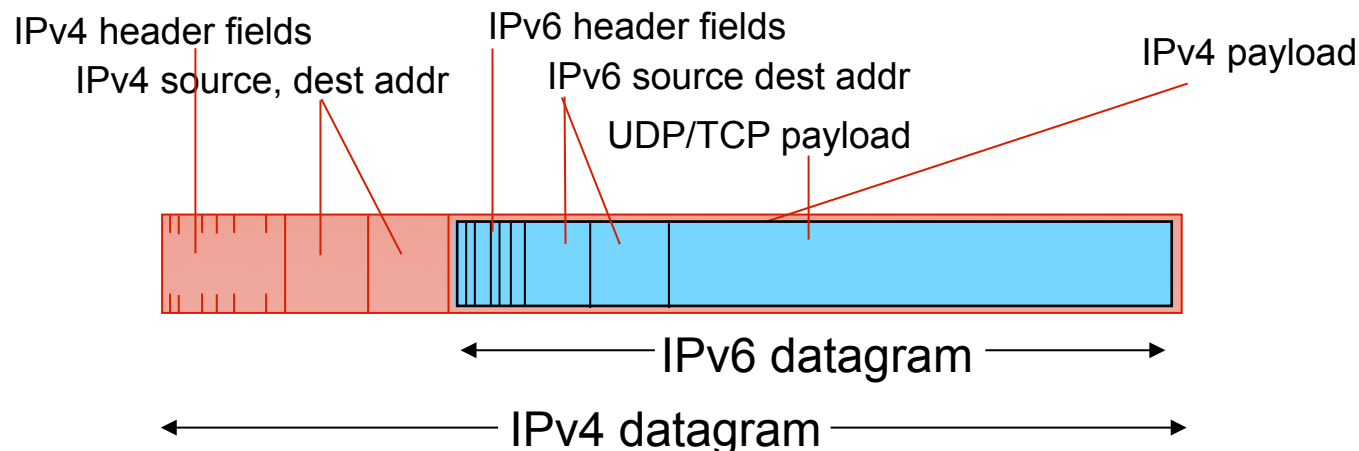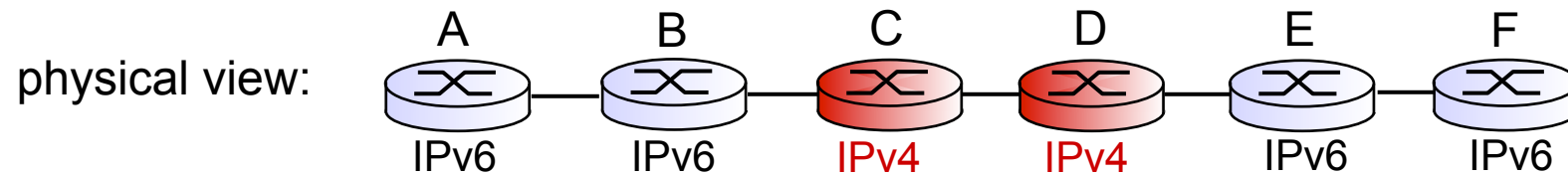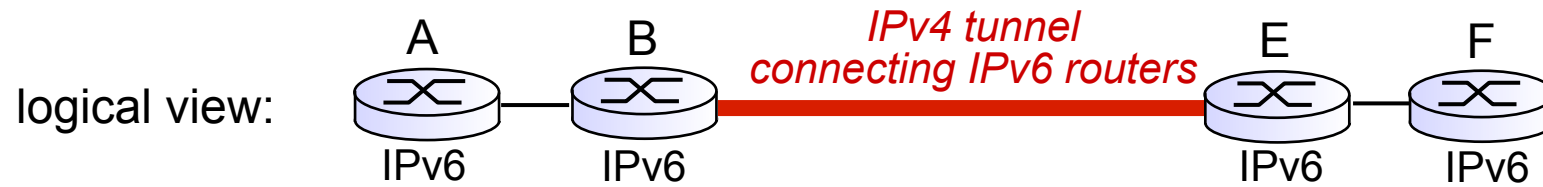| ver | pri | flow label | | |
|---|---|---|---|---|
| payload len | | next hdr | | hop limit |
| source address (128 bits) | | | | |
| destination address (128 bits) | | | | |
| data | | | | |

← 32 bits →

# Other changes from IPv4

❖ *checksum:* removed entirely to reduce processing time at each hop

❖ *options:* allowed, but outside of header, indicated by "Next Header" field

❖ *ICMPv6:* new version of ICMP

  ▪ additional message types, e.g. "Packet Too Big"

  ▪ multicast group management functions

# Transition from IPv4 to IPv6

- ❖ not all routers can be upgraded simultaneously
  - ▪ no "flag days"
  - ▪ how will network operate with mixed IPv4 and IPv6 routers?
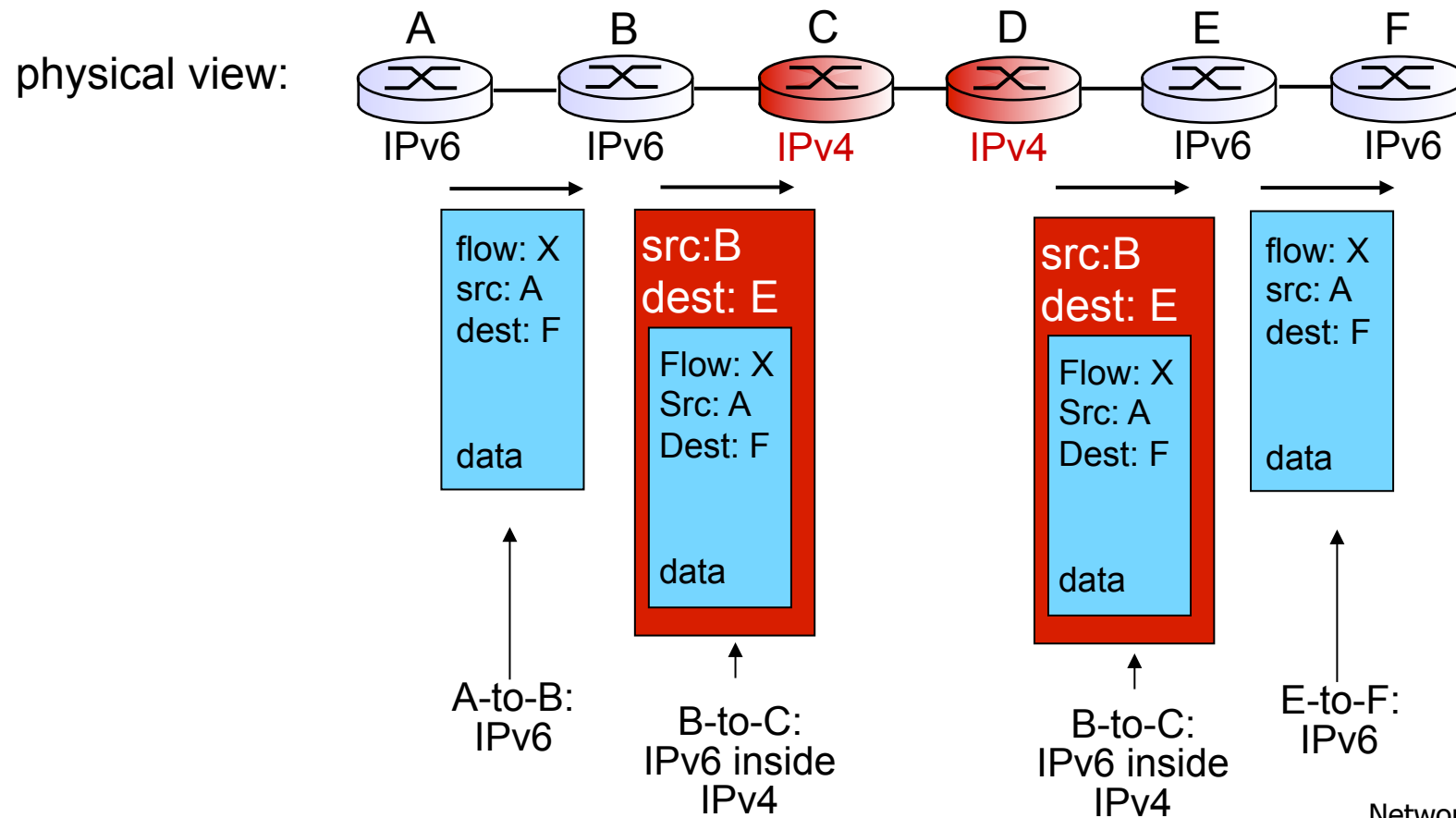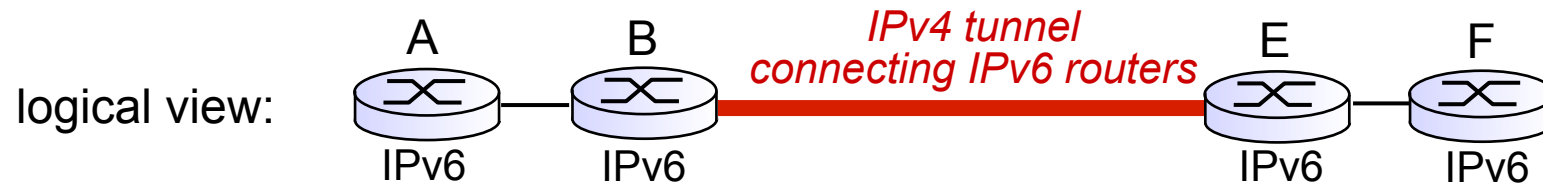- ❖ *tunneling:* IPv6 datagram carried as *payload* in IPv4 datagram among IPv4 routers

IPv4 header fields
IPv4 source, dest addr

IPv6 header fields
IPv6 source dest addr
UDP/TCP payload

IPv4 payload

IPv6 datagram

IPv4 datagram

# Tunneling

logical view:

A     B       *IPv4 tunnel connecting IPv6 routers*       E     F

IPv6    IPv6                            IPv6    IPv6

physical view:

A     B     C     D     E     F

IPv6    IPv6    IPv4    IPv4    IPv6    IPv6

# Tunneling

logical view:

A       B        *IPv4 tunnel connecting IPv6 routers*       E       F

IPv6       IPv6                               IPv6       IPv6

physical view:

A       B       C       D       E       F

IPv6       IPv6       IPv4       IPv4       IPv6       IPv6

flow: X
src: A
dest: F

data

src:B
dest: E

Flow: X
Src: A
Dest: F

data

src:B
dest: E

Flow: X
Src: A
Dest: F

data

flow: X
src: A
dest: F

data

A-to-B:
IPv6

B-to-C:
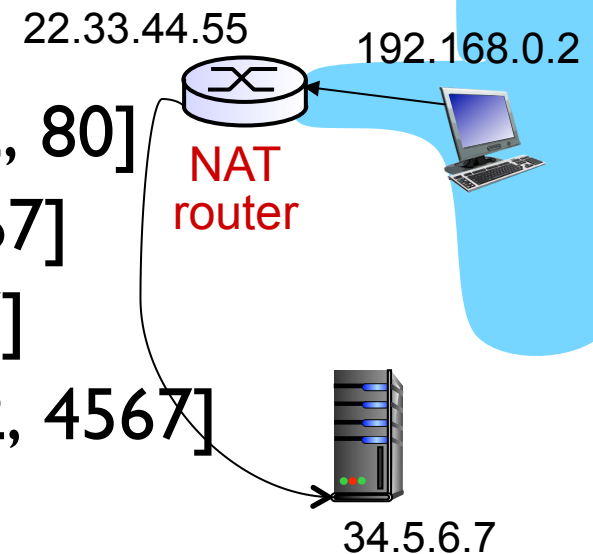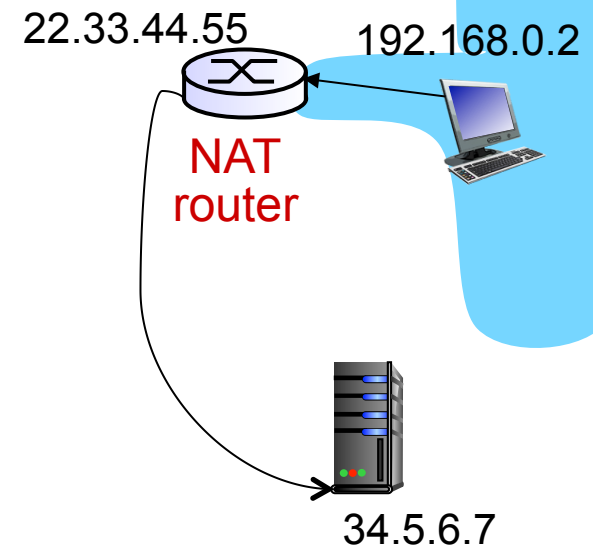IPv6 inside
IPv4

B-to-C:
IPv6 inside
IPv4

E-to-F:
IPv6

# Quiz: NAT

❖ A host with a private IP address 192.168.0.2 opens a TCP socket on its local port 4567 and connects to a web server at 34.5.6.7. The NAT's public IP address is 22.33.44.55. Which of the following mapping entries *could* the NAT create as a result?

A. [22.33.44.55, 3333]→[192.168.0.2, 80]

B. [34.5.6.7, 80] → [22.33.44.55, 4567]

C. [192.168.0.2, 80]→[34.5.6.7, 4567]

D. [22.33.44.55, 3967]→[192.168.0.2, 4567]

22.33.44.55

192.168.0.2

NAT router

34.5.6.7

# Quiz: NAT

❖ A host with a private IP address 192.168.0.2 opens a TCP socket on its local port 4567 and connects to a web server at 34.5.6.7. The NAT's public IP address is 22.33.44.55. Suppose the NAT created the mapping [22.33.44.55, 3967]→[192.168.0.2, 4567] as a result. What are the source and destination port numbers in the SYNACK response from the server?

A. 80, 3967

B. 4567, 80

C. 3967, 80

D. 3967, 4567
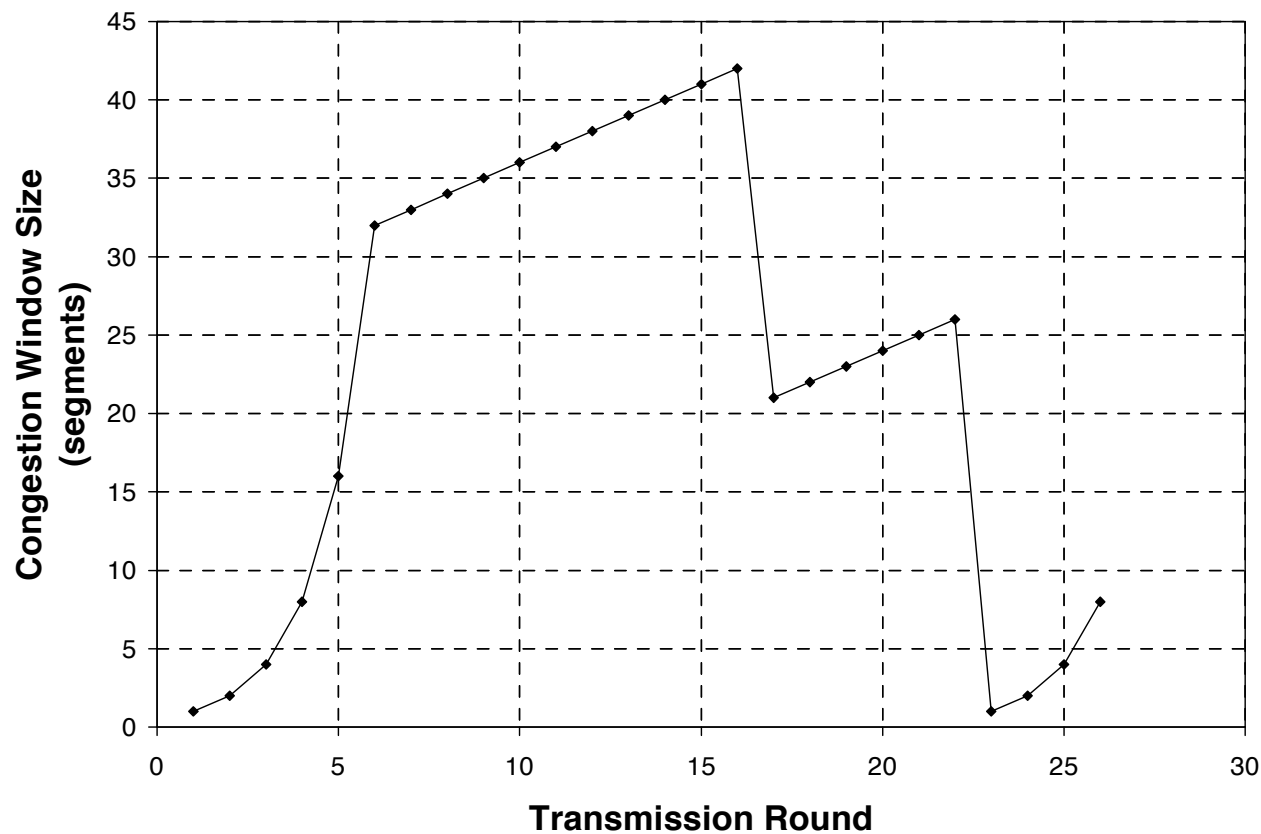
22.33.44.55          192.168.0.2

NAT
router

34.5.6.7

# Quiz: IPv6

❖ Which of the following is not true?

A. IPv6 increases the size of the IP address space from 2^32 to 2^128.

B. IPv6 removes checksums and fragmentation compared to IPv4.

C. IPv6 has fixed length headers.

D. IPv6 adds reliability at the network layer.

# Practice Problem

Consider the following plot of TCP window size as a function of time. Assume that the version is Reno.

# Practice Problem

- ❖ Identify the intervals of time when TCP slow start is operating.

- ❖ Identify the intervals of time when TCP congestion avoidance is operating.

- ❖ After the 16$^{th}$ transmission round, is segment loss detected by a triple duplicate ACK or a timeout?

- ❖ What is the initial value of ssthreshold at the first transmission round?

# Practice Problem

❖ What is the value of ssthreshold at the 18th transmission round?

❖ What is the value of ssthreshold at the 24th transmission round?

❖ During which transmission round is the 70th segment sent?

❖ Assuming a packet loss is detected after the 26th round by the receipt of triple duplicate ACK, what will be value of the congestion window size and of ssthreshold?
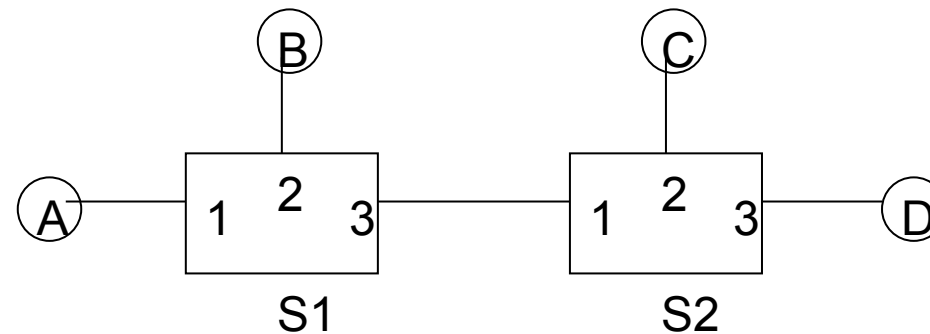
- Consider hosts A, B, C and D connected together via virtual circuit switches S1 and S2 as shown in the figure below (on next page). The tables list, for each switch, the forwarding entries indicating which (port, VC #) pairs are connected to which other (port, VC #) pairs. Assume that the connections are bi-directional. List all host-to-host connections that have been established by these entries. Note that, you will lose marks if you include host-to-host connections that have not been established so do not list arbitrary combinations. You DO NOT have to provide any explanation, simply write down the list of connections that you think have been established.

# Past Exam Question



**S1 Table**

| Incoming Port | VC # | Outgoing Port | VC # |
|---------------|------|---------------|------|
| 1 | 2 | 3 | 1 |
| 1 | 1 | 2 | 3 |
| 2 | 1 | 3 | 2 |

**S2 Table**

| Incoming Port | VC # | Outgoing Port | VC # |
|---------------|------|---------------|------|
| 1 | 1 | 3 | 3 |
| 1 | 2 | 3 | 2 |
| | | | |

Solution: Following connections have been established:

# Forwarding

❖ A forwarding table for a router is as follows:

| Address prefix | Next hop |
|----------------|----------|
| 196.94.2.0/24 | A |
| 196.94.2.128/25 | B |
| 196.94.0.0/16 | C |
| 196.94.64.0/18 | D |
| 196.76.0.0/14 | E |
| 140.0.0.0/8 | F |
| 128.0.0.0/2 | G |
| 0.0.0.0/1 | H |

❖ State the next hop for the following destinations:
  - 139.1.1.1
  - 196.94.2.100
  - 196.94.2.200
  - 196.94.3.100