

# CS 356: Computer Network Architectures

## Lecture 9: The Internet Protocol (IP) Ch 3.2

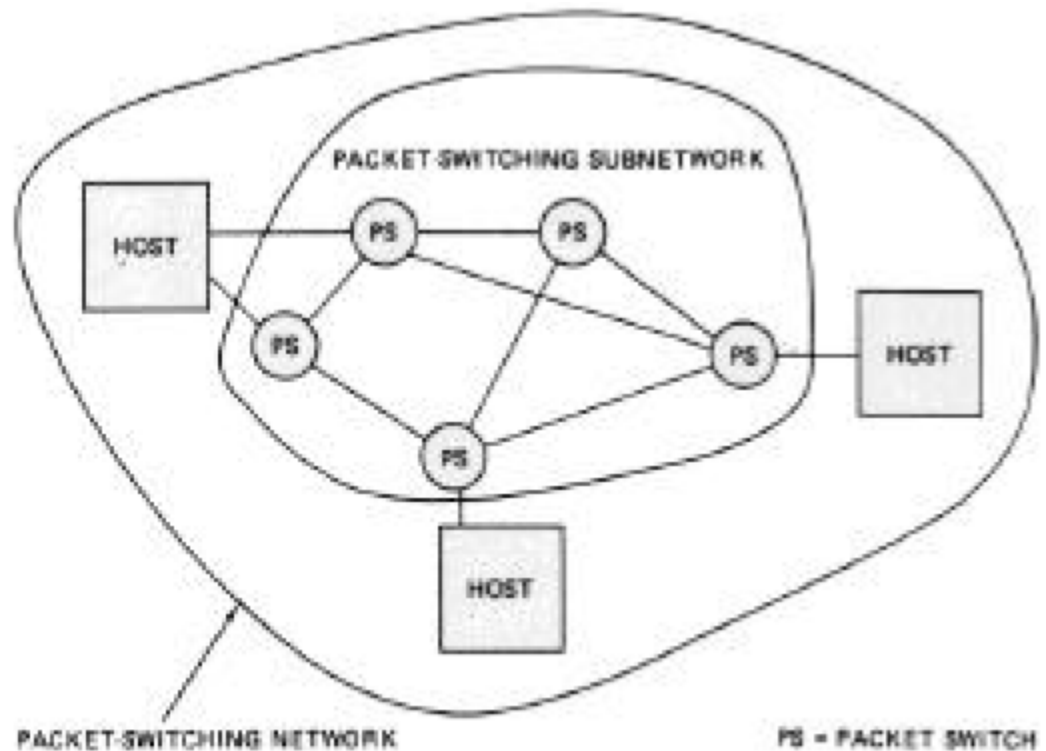
Xiaowei Yang  
xwy@cs.duke.edu

# Overview

- History of IP
- IP header format
- IP addressing
- IP forwarding
  - Forwarding algorithm
  - Fragmentation

# History of the Internet

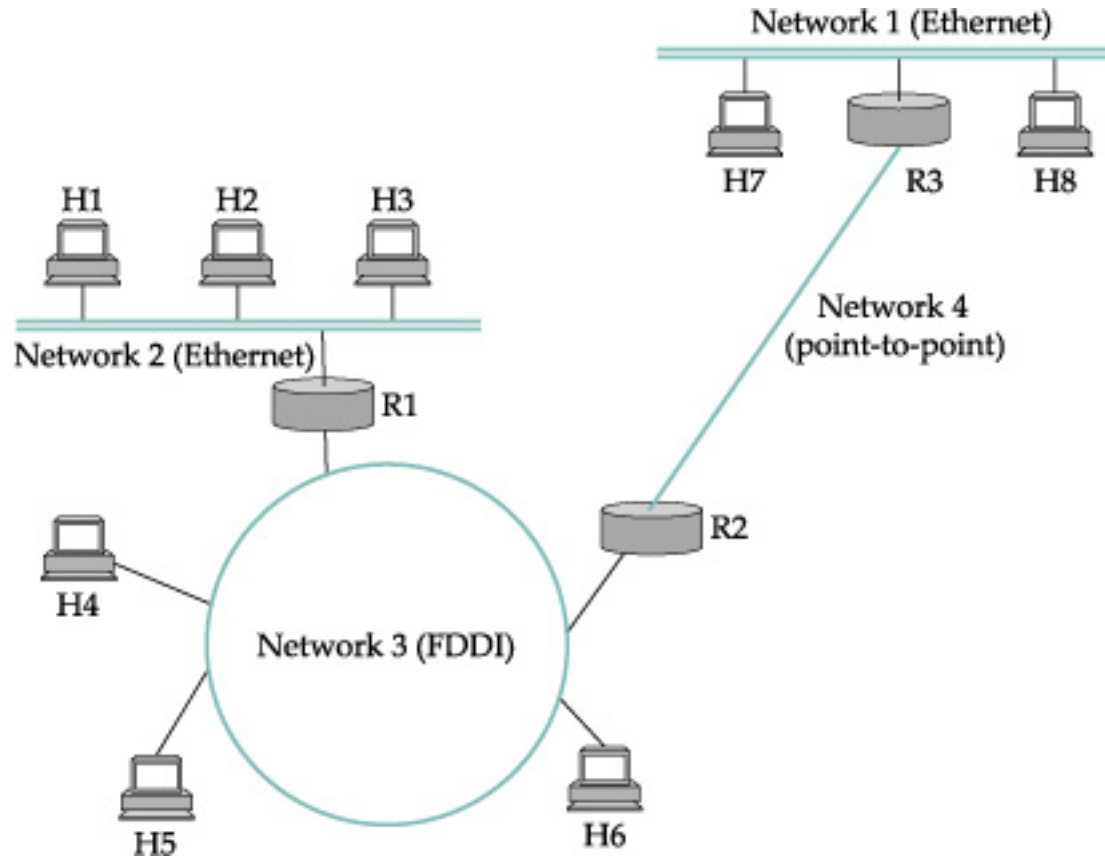
- Original design goal:  
Interconnecting different networks
- Many different types of packet switch networks
  - ARPANET, packet satellite networks, ground-based packet radio networks, and other networks.
- Each has
  - Hosts, packet switches, processes
  - A protocol for communication
- Q: what would you do differently given such a design task?



# Challenges

1. Different addressing schemes and host communication protocols
  - Ethernet, FDDI, ATM
2. Different Maximum Transmission Units (MTUs)
3. Different success or failure indicators
4. End-to-end reliability: failures may occur at each network
5. Different control protocols
  - Status information, routing, fault detection/isolation

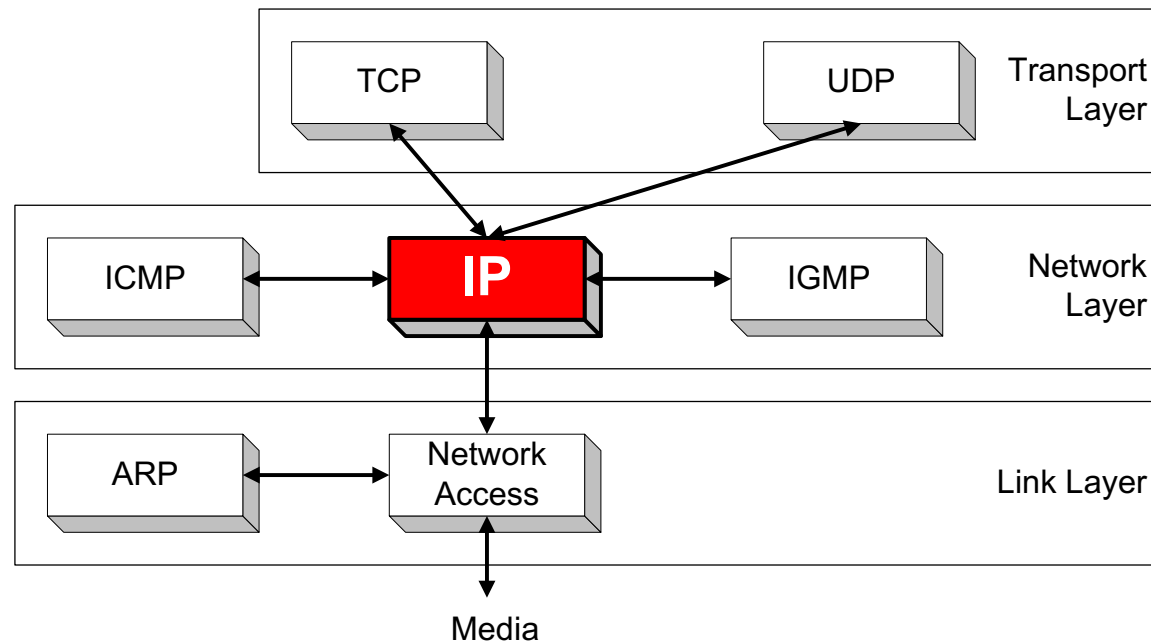
# Inter-networking



- One level of indirection
  - Routers interface different networks
- Uniform addressing (IP)
- Routers send packets to their destination IP addresses

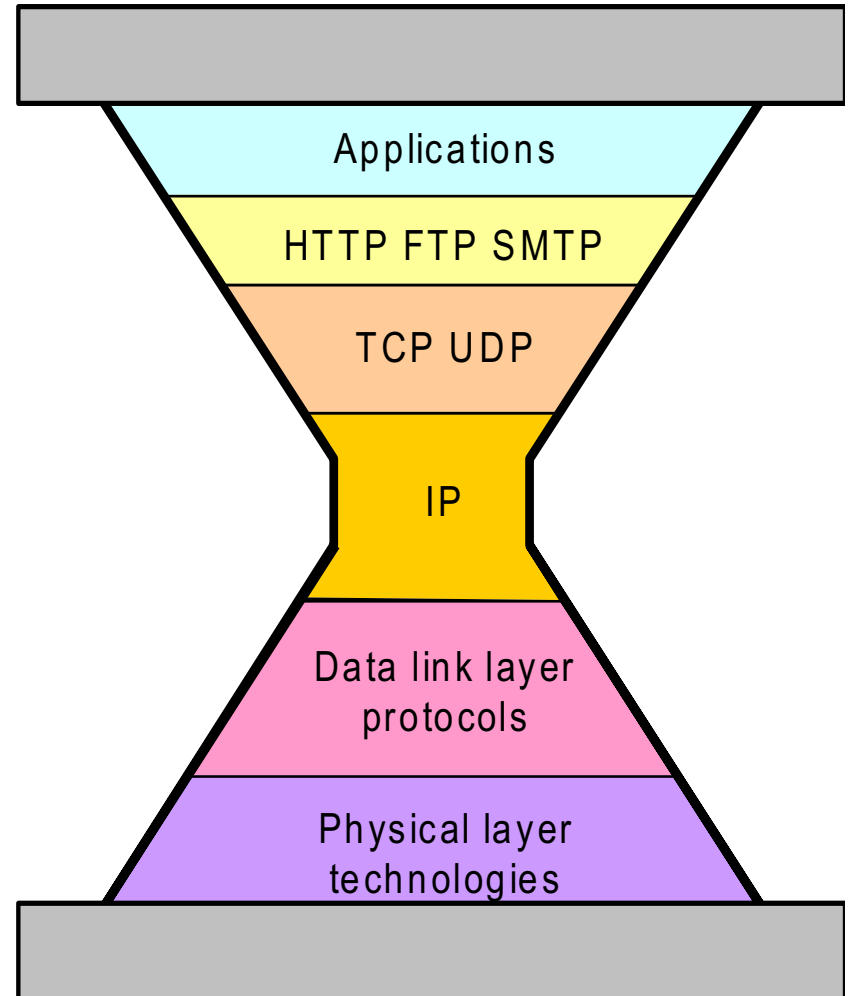
# Internet Protocol

- IP (Internet Protocol) is a Network Layer Protocol
- IP's current version is Version 4 (IPv4). It is specified in RFC 791.
- IPv6 is also deployed



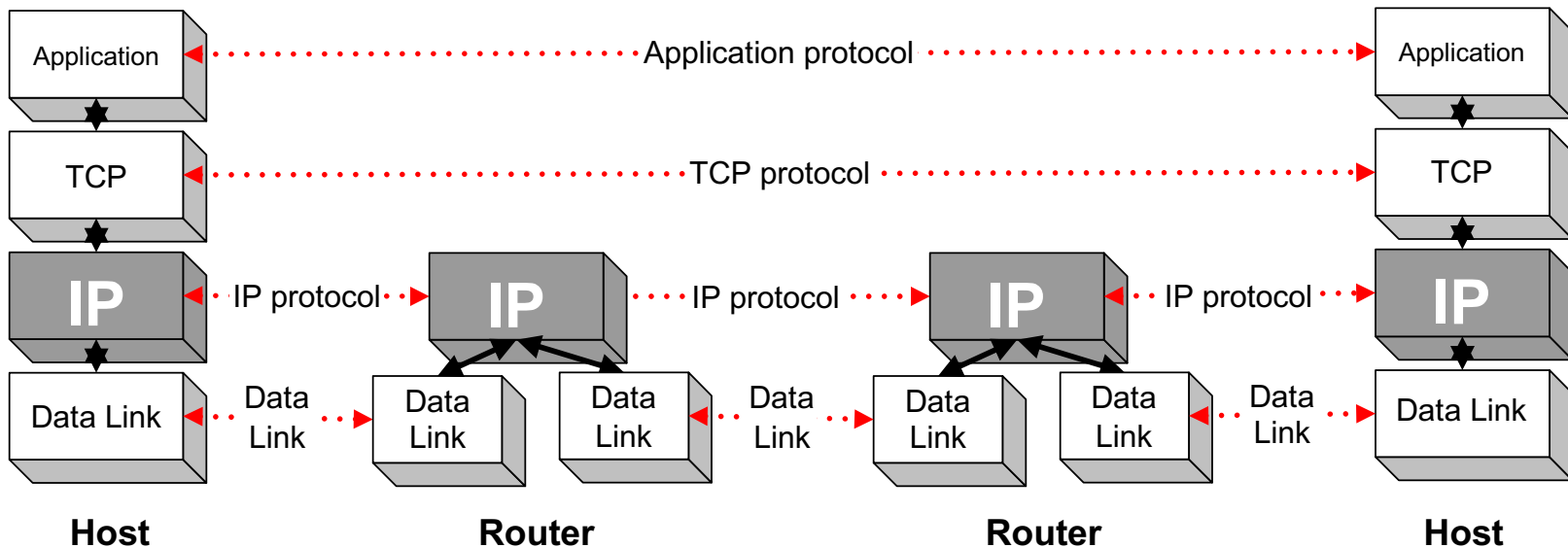
# IP: the thin waist of the hourglass

- **IP is the waist of the hourglass of the Internet protocol architecture**
- Multiple higher-layer protocols
- Multiple lower-layer protocols
- Only one protocol at the network layer.
- What is the advantage of this architecture?
  - To avoid the  $N * M$  problem



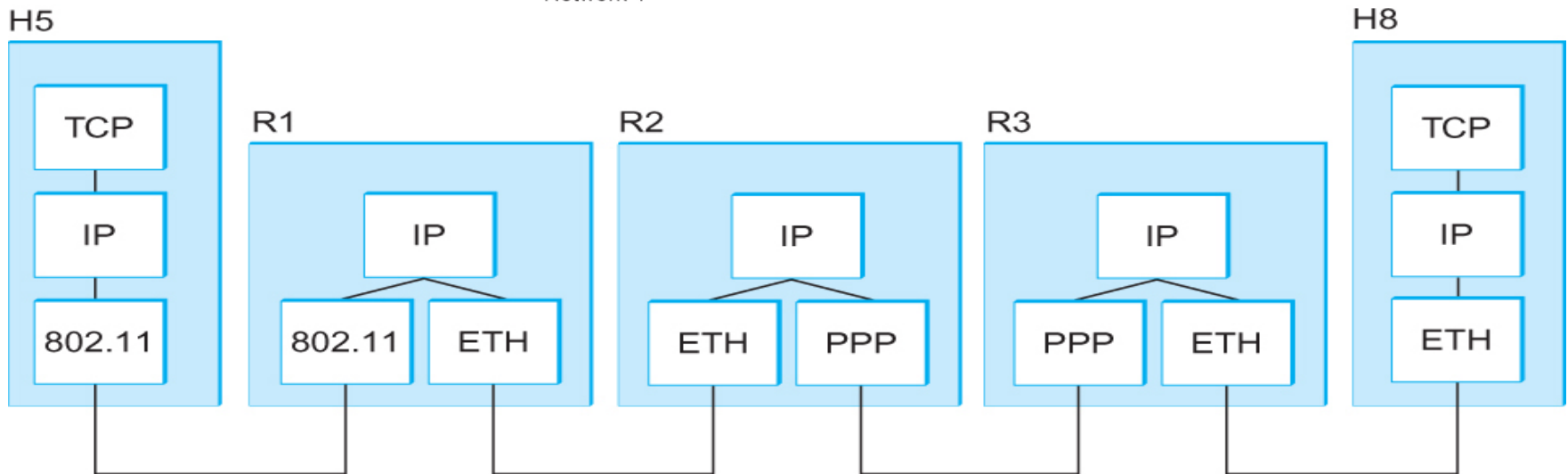
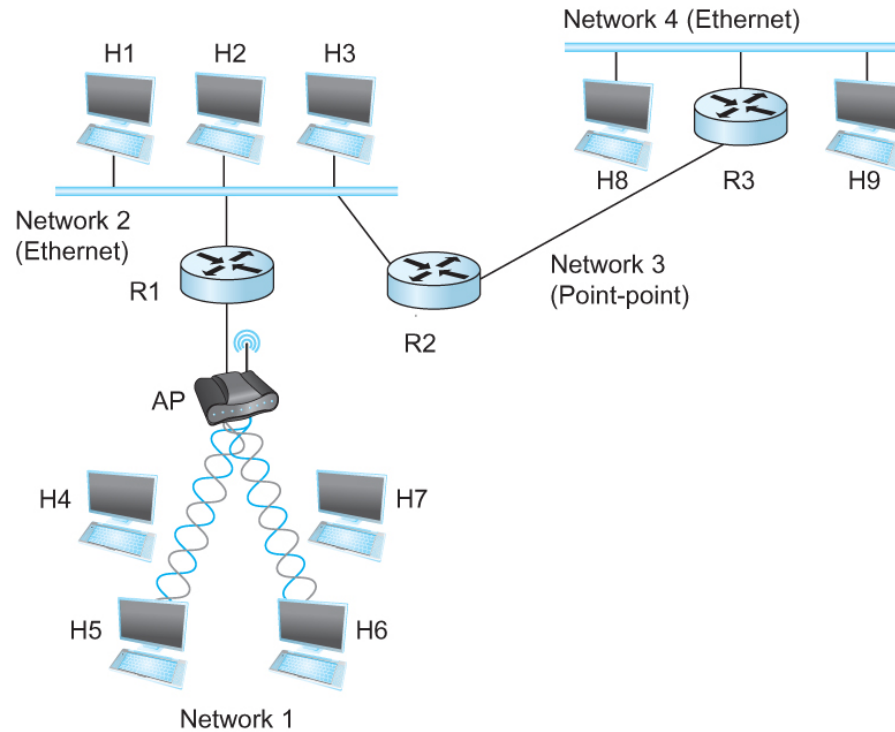
# Application protocol

- Routers look at a packet's IP header and link layer header





# A simple network



# IP Service Model

- Delivery service of IP is minimal
- IP provides an **unreliable connectionless** best effort service (also called: “datagram service”).
  - **Unreliable**
  - **Connectionless**
  - **Best effort**
- Consequences
  - Loss, out of order, and duplicate must be handled at the upper layer

# Basic IP router functions

- Things you need to understand to do lab2
  - Internet protocol
    - IP header
    - IP addressing
    - IP forwarding
  - Address resolution protocol
  - Error reporting and control
    - Internet Control Message Protocol

# IP header format

| IPv4 Header Format |       |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
|--------------------|-------|------------------------|---|---|---|-----|---|---|---|----------|---|----|----|----|----|-----|----|-----------------|----|----|-----------------|----|----|----|----|----|----|----|----|----|----|----|----|
| Offsets            | Octet | 0                      |   |   |   |     |   |   |   | 1        |   |    |    |    |    |     |    | 2               |    |    |                 |    |    |    |    | 3  |    |    |    |    |    |    |    |
| Octet              | Bit   | 0                      | 1 | 2 | 3 | 4   | 5 | 6 | 7 | 8        | 9 | 10 | 11 | 12 | 13 | 14  | 15 | 16              | 17 | 18 | 19              | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
| 0                  | 0     | Version                |   |   |   | IHL |   |   |   | DSCP     |   |    |    |    |    | ECN |    | Total Length    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 4                  | 32    | Identification         |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    | Flags           |    |    | Fragment Offset |    |    |    |    |    |    |    |    |    |    |    |    |
| 8                  | 64    | Time To Live           |   |   |   |     |   |   |   | Protocol |   |    |    |    |    |     |    | Header Checksum |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 12                 | 96    | Source IP Address      |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 16                 | 128   | Destination IP Address |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 20                 | 160   | Options (if IHL > 5)   |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 24                 | 192   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 28                 | 224   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 32                 | 256   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |

- 20 bytes fixed length header + variable length options
- **Internet Header Length (IHL 4 bits):** the length of header in 32-bit words
  - Maximum header length?

# IP header format

| IPv4 Header Format |       |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
|--------------------|-------|------------------------|---|---|---|-----|---|---|---|----------|---|----|----|----|----|-----|----|-----------------|----|----|-----------------|----|----|----|----|----|----|----|----|----|----|----|----|
| Offsets            | Octet | 0                      |   |   |   |     |   |   |   | 1        |   |    |    |    |    |     |    | 2               |    |    |                 |    |    |    |    | 3  |    |    |    |    |    |    |    |
| Octet              | Bit   | 0                      | 1 | 2 | 3 | 4   | 5 | 6 | 7 | 8        | 9 | 10 | 11 | 12 | 13 | 14  | 15 | 16              | 17 | 18 | 19              | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
| 0                  | 0     | Version                |   |   |   | IHL |   |   |   | DSCP     |   |    |    |    |    | ECN |    | Total Length    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 4                  | 32    | Identification         |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    | Flags           |    |    | Fragment Offset |    |    |    |    |    |    |    |    |    |    |    |    |
| 8                  | 64    | Time To Live           |   |   |   |     |   |   |   | Protocol |   |    |    |    |    |     |    | Header Checksum |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 12                 | 96    | Source IP Address      |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 16                 | 128   | Destination IP Address |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 20                 | 160   | Options (if IHL > 5)   |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 24                 | 192   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 28                 | 224   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 32                 | 256   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |

- **DSCP (Differentiated Services Code Point 6 bits):**  
old Type of Service
  - Real-time, VoIP
- **Explicit Congestion Notification (ECN)**
  - Early Congestion notice

# IP header format

| IPv4 Header Format |       |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |
|--------------------|-------|------------------------|---|---|---|-----|---|---|---|----------|---|----|----|----|----|-----|----|-----------------|----|----|----|-----------------|----|----|----|----|----|----|----|----|----|----|----|
| Offsets            | Octet | 0                      |   |   |   |     |   |   |   | 1        |   |    |    |    |    |     |    | 2               |    |    |    |                 |    |    |    | 3  |    |    |    |    |    |    |    |
| Octet              | Bit   | 0                      | 1 | 2 | 3 | 4   | 5 | 6 | 7 | 8        | 9 | 10 | 11 | 12 | 13 | 14  | 15 | 16              | 17 | 18 | 19 | 20              | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
| 0                  | 0     | Version                |   |   |   | IHL |   |   |   | DSCP     |   |    |    |    |    | ECN |    | Total Length    |    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |
| 4                  | 32    | Identification         |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    | Flags           |    |    |    | Fragment Offset |    |    |    |    |    |    |    |    |    |    |    |
| 8                  | 64    | Time To Live           |   |   |   |     |   |   |   | Protocol |   |    |    |    |    |     |    | Header Checksum |    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |
| 12                 | 96    | Source IP Address      |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |
| 16                 | 128   | Destination IP Address |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |
| 20                 | 160   | Options (if IHL > 5)   |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |
| 24                 | 192   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |
| 28                 | 224   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |
| 32                 | 256   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |

- **Total length (16 bits):** packet length in bytes, including the header
  - 65535 bytes
  - Fragmentation and reassembly

# IP header format

| IPv4 Header Format |       |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |
|--------------------|-------|------------------------|---|---|---|-----|---|---|---|----------|---|----|----|----|----|-----|----|-----------------|----|----|----|-----------------|----|----|----|----|----|----|----|----|----|----|----|
| Offsets            | Octet | 0                      |   |   |   |     |   |   |   | 1        |   |    |    |    |    |     |    | 2               |    |    |    |                 |    |    |    | 3  |    |    |    |    |    |    |    |
| Octet              | Bit   | 0                      | 1 | 2 | 3 | 4   | 5 | 6 | 7 | 8        | 9 | 10 | 11 | 12 | 13 | 14  | 15 | 16              | 17 | 18 | 19 | 20              | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
| 0                  | 0     | Version                |   |   |   | IHL |   |   |   | DSCP     |   |    |    |    |    | ECN |    | Total Length    |    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |
| 4                  | 32    | Identification         |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    | Flags           |    |    |    | Fragment Offset |    |    |    |    |    |    |    |    |    |    |    |
| 8                  | 64    | Time To Live           |   |   |   |     |   |   |   | Protocol |   |    |    |    |    |     |    | Header Checksum |    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |
| 12                 | 96    | Source IP Address      |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |
| 16                 | 128   | Destination IP Address |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |
| 20                 | 160   | Options (if IHL > 5)   |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |
| 24                 | 192   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |
| 28                 | 224   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |
| 32                 | 256   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |

- **Identification (16 bits):** Unique datagram identifier from a host
  - Incremented whenever a datagram is transmitted (in some OS)
  - Used by many researchers for various purposes

# IP header format

| IPv4 Header Format |       |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
|--------------------|-------|------------------------|---|---|---|-----|---|---|---|----------|---|----|----|----|----|-----|----|-----------------|----|----|-----------------|----|----|----|----|----|----|----|----|----|----|----|----|
| Offsets            | Octet | 0                      |   |   |   |     |   |   |   | 1        |   |    |    |    |    |     |    | 2               |    |    |                 |    |    |    |    | 3  |    |    |    |    |    |    |    |
| Octet              | Bit   | 0                      | 1 | 2 | 3 | 4   | 5 | 6 | 7 | 8        | 9 | 10 | 11 | 12 | 13 | 14  | 15 | 16              | 17 | 18 | 19              | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
| 0                  | 0     | Version                |   |   |   | IHL |   |   |   | DSCP     |   |    |    |    |    | ECN |    | Total Length    |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 4                  | 32    | Identification         |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    | Flags           |    |    | Fragment Offset |    |    |    |    |    |    |    |    |    |    |    |    |
| 8                  | 64    | Time To Live           |   |   |   |     |   |   |   | Protocol |   |    |    |    |    |     |    | Header Checksum |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 12                 | 96    | Source IP Address      |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 16                 | 128   | Destination IP Address |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 20                 | 160   | Options (if IHL > 5)   |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 24                 | 192   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 28                 | 224   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |
| 32                 | 256   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |    |                 |    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |

- **Flags (3 bits):**
  - First bit always set to 0
  - DF bit (Do not fragment)
  - MF bit (More fragments)
- **Fragment offset (13 bits)**
- **Identification, Flags, Fragment offset**
  - fragmentation and assembly



# IP header format

| IPv4 Header Format |       |                        |   |   |   |     |   |   |          |      |   |    |    |    |    |                 |       |              |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |  |
|--------------------|-------|------------------------|---|---|---|-----|---|---|----------|------|---|----|----|----|----|-----------------|-------|--------------|----|-----------------|----|----|----|----|----|----|----|----|----|----|----|----|----|--|
| Offsets            | Octet | 0                      |   |   |   |     |   |   |          | 1    |   |    |    |    |    |                 |       | 2            |    |                 |    |    |    |    |    | 3  |    |    |    |    |    |    |    |  |
| Octet              | Bit   | 0                      | 1 | 2 | 3 | 4   | 5 | 6 | 7        | 8    | 9 | 10 | 11 | 12 | 13 | 14              | 15    | 16           | 17 | 18              | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |  |
| 0                  | 0     | Version                |   |   |   | IHL |   |   |          | DSCP |   |    |    |    |    | ECN             |       | Total Length |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |  |
| 4                  | 32    | Identification         |   |   |   |     |   |   |          |      |   |    |    |    |    |                 | Flags |              |    | Fragment Offset |    |    |    |    |    |    |    |    |    |    |    |    |    |  |
| 8                  | 64    | Time To Live           |   |   |   |     |   |   | Protocol |      |   |    |    |    |    | Header Checksum |       |              |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |  |
| 12                 | 96    | Source IP Address      |   |   |   |     |   |   |          |      |   |    |    |    |    |                 |       |              |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |  |
| 16                 | 128   | Destination IP Address |   |   |   |     |   |   |          |      |   |    |    |    |    |                 |       |              |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |  |
| 20                 | 160   | Options (if IHL > 5)   |   |   |   |     |   |   |          |      |   |    |    |    |    |                 |       |              |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |  |
| 24                 | 192   |                        |   |   |   |     |   |   |          |      |   |    |    |    |    |                 |       |              |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |  |
| 28                 | 224   |                        |   |   |   |     |   |   |          |      |   |    |    |    |    |                 |       |              |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |  |
| 32                 | 256   |                        |   |   |   |     |   |   |          |      |   |    |    |    |    |                 |       |              |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |  |

- **Time To Live (TTL) (1byte):**
  - Specifies the longest path before a datagram is dropped
  - Role of TTL field: Ensure that a packet is eventually dropped when a routing loop occurs

Used as follows:

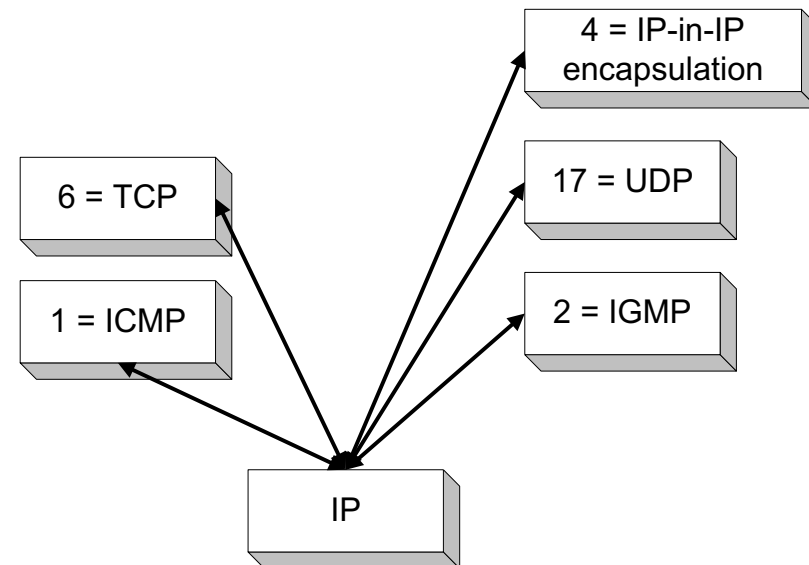
- Sender sets the value (e.g., 64)
- Each router decrements the value by 1
- When the value reaches 0, the datagram is dropped

# IP header format

IPv4 Header Format

| Offsets | Octet | 0                      |   |   |   |     |   |   |   | 1        |   |    |    |    |    |     |       | 2               |    |                 |    |    |    |    |    | 3  |    |    |    |    |    |    |    |
|---------|-------|------------------------|---|---|---|-----|---|---|---|----------|---|----|----|----|----|-----|-------|-----------------|----|-----------------|----|----|----|----|----|----|----|----|----|----|----|----|----|
| Octet   | Bit   | 0                      | 1 | 2 | 3 | 4   | 5 | 6 | 7 | 8        | 9 | 10 | 11 | 12 | 13 | 14  | 15    | 16              | 17 | 18              | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
| 0       | 0     | Version                |   |   |   | IHL |   |   |   | DSCP     |   |    |    |    |    | ECN |       | Total Length    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 4       | 32    | Identification         |   |   |   |     |   |   |   |          |   |    |    |    |    |     | Flags |                 |    | Fragment Offset |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 8       | 64    | Time To Live           |   |   |   |     |   |   |   | Protocol |   |    |    |    |    |     |       | Header Checksum |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 12      | 96    | Source IP Address      |   |   |   |     |   |   |   |          |   |    |    |    |    |     |       |                 |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 16      | 128   | Destination IP Address |   |   |   |     |   |   |   |          |   |    |    |    |    |     |       |                 |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 20      | 160   | Options (if IHL > 5)   |   |   |   |     |   |   |   |          |   |    |    |    |    |     |       |                 |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 24      | 192   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |       |                 |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 28      | 224   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |       |                 |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 32      | 256   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |       |                 |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |

- **Protocol (1 byte):**
  - Specifies the higher-layer protocol.
  - De-multiplexing to higher layers.



# IP header format

IPv4 Header Format

| Offsets | Octet | 0                      |   |   |   |     |   |   |   | 1        |   |    |    |    |    |     |       | 2               |    |                 |    |    |    |    |    | 3  |    |    |    |    |    |    |    |
|---------|-------|------------------------|---|---|---|-----|---|---|---|----------|---|----|----|----|----|-----|-------|-----------------|----|-----------------|----|----|----|----|----|----|----|----|----|----|----|----|----|
| Octet   | Bit   | 0                      | 1 | 2 | 3 | 4   | 5 | 6 | 7 | 8        | 9 | 10 | 11 | 12 | 13 | 14  | 15    | 16              | 17 | 18              | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
| 0       | 0     | Version                |   |   |   | IHL |   |   |   | DSCP     |   |    |    |    |    | ECN |       | Total Length    |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 4       | 32    | Identification         |   |   |   |     |   |   |   |          |   |    |    |    |    |     | Flags |                 |    | Fragment Offset |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 8       | 64    | Time To Live           |   |   |   |     |   |   |   | Protocol |   |    |    |    |    |     |       | Header Checksum |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 12      | 96    | Source IP Address      |   |   |   |     |   |   |   |          |   |    |    |    |    |     |       |                 |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 16      | 128   | Destination IP Address |   |   |   |     |   |   |   |          |   |    |    |    |    |     |       |                 |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 20      | 160   | Options (if IHL > 5)   |   |   |   |     |   |   |   |          |   |    |    |    |    |     |       |                 |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 24      | 192   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |       |                 |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 28      | 224   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |       |                 |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 32      | 256   |                        |   |   |   |     |   |   |   |          |   |    |    |    |    |     |       |                 |    |                 |    |    |    |    |    |    |    |    |    |    |    |    |    |

- **Header checksum (16 bits):** header checksum
  - Header only
  - Must be computed at every hop!

# Fields of the IP Header

- **Options:**

- **Record Route:** each router that processes the packet adds its IP address to the header.
- **Timestamp:** each router that processes the packet adds its IP address and time to the header.
- **(loose) Source Routing:** specifies a list of routers that must be traversed.
- **(strict) Source Routing:** specifies a list of the only routers that can be traversed.
- IP options increase routers processing overhead

- **Padding:** Padding bytes are added to ensure that header ends on a 4-byte boundary

# Global IP addresses

# What is an IP Address?

- An IP address is a unique global identifier for a network interface
  - An IP address uniquely identifies a network location
- Routers forwards a packet based on the destination address of the packet
- Uniqueness ensures global reachability

# IP versions

- IPv4 (32-bit)
  - Classful IP addresses (obsolete)
  - Classless inter-domain routing (CIDR) (RFC 854, current standard)
- IP Version 6 addresses (128-bit)

# Dotted Decimal Notation

- Each byte is identified by a decimal number in the range  $[0...255]$ :

|          |          |          |          |
|----------|----------|----------|----------|
| 10000000 | 10001111 | 10001001 | 10010000 |
|----------|----------|----------|----------|

1<sup>st</sup> Byte

2<sup>nd</sup> Byte

3<sup>rd</sup> Byte

4<sup>th</sup> Byte

**= 128**

**= 143**

**= 137**

**= 144**

**128.143.137.144**



# Structure of an IP address



- An IP address has a structure
  - Network prefix identifies a network
  - Host number identifies a specific host interface
- Improves the scalability of routing
  - Scales better than flat addresses

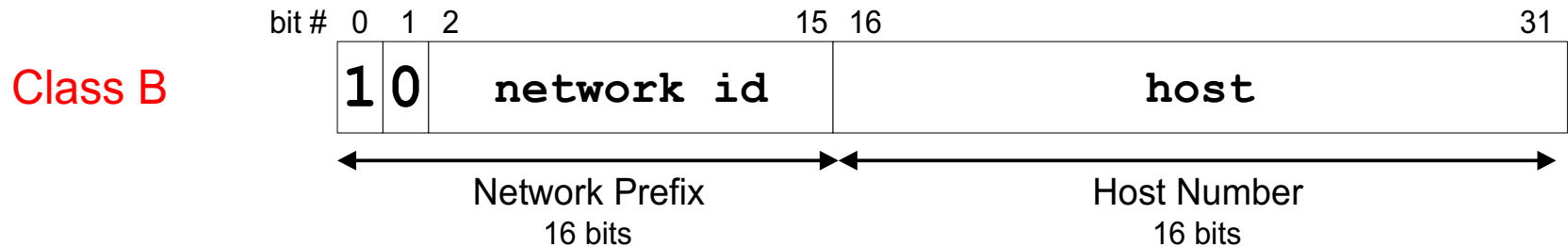
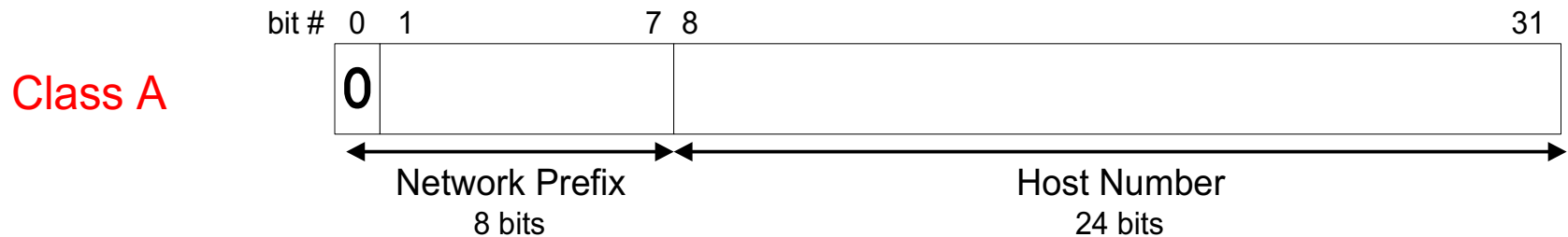
# How long is a network prefix?

- **Before 1993:** The network prefix is implicitly defined (**class-based addressing**)
- **After 1993:** The network prefix is indicated by a **netmask**

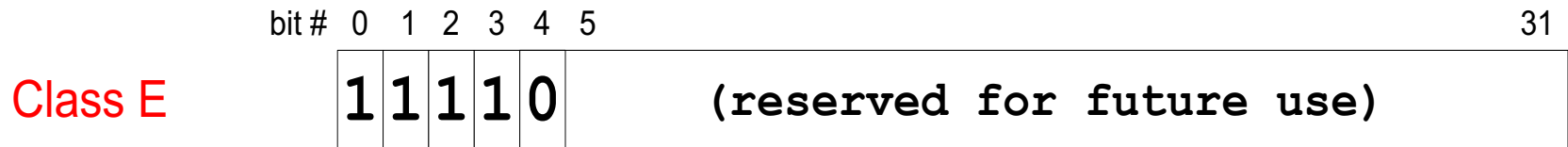
# Before 1993: Class-based addressing

- The Internet address space was divided up into classes:
  - **Class A:** Network prefix is 8 bits long
  - **Class B:** Network prefix is 16 bits long
  - **Class C:** Network prefix is 24 bits long
  - Class D is multicast address
  - Class E is reserved

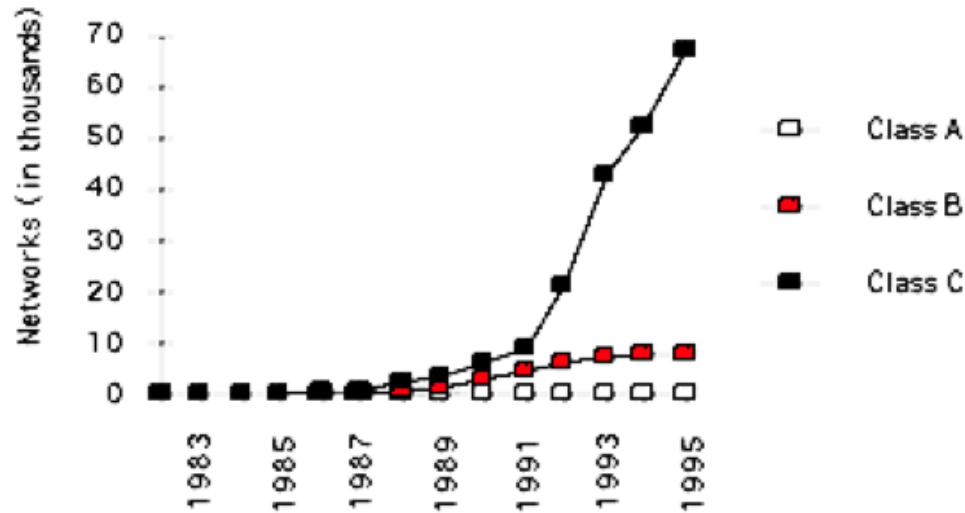
# Classful IP Addresses (before 1993)



# Classful IP Addresses (before 1993)



# Problems with Classful IP Addresses

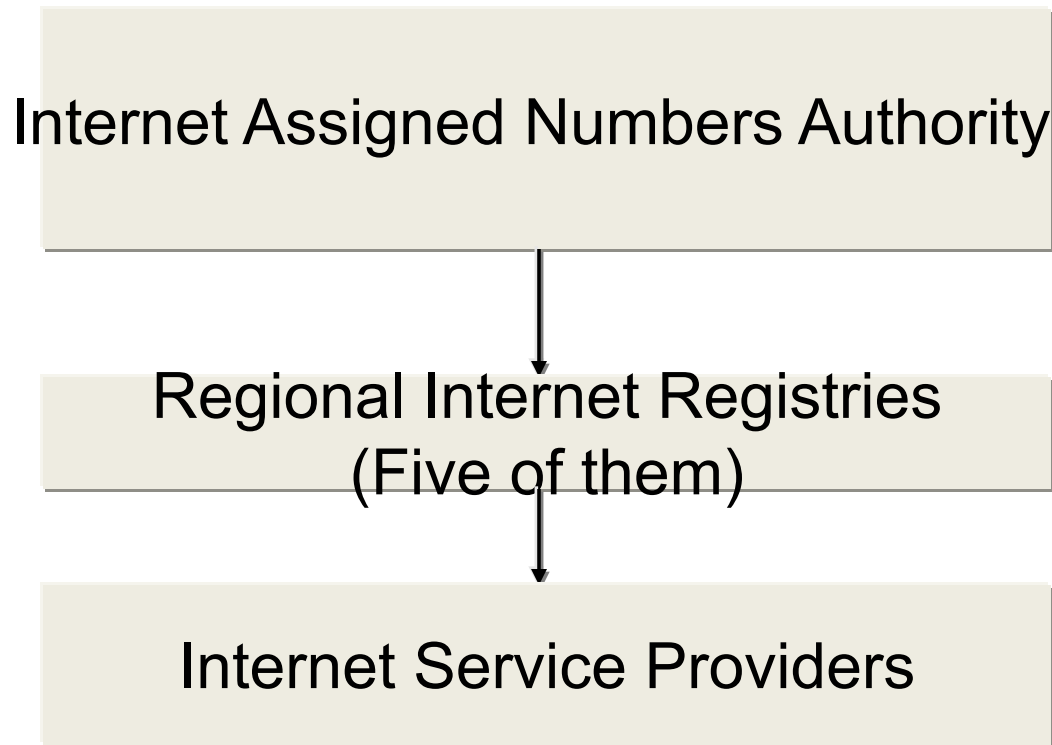


- Fast growing routing table size
  - Each router must have an entry for every network prefix
  - $\sim 2^{21} = 2,097,152$  class C networks
  - In 1993, the size of routing tables started to outgrow the capacity of routers
- Local admins must request another network number before installing a new network at their site

## Solution: Classless Inter-domain routing (CIDR)

- Network prefix is of variable length
  - No rigid class boundary
- Addresses are allocated hierarchically
- Routers can aggregate multiple address prefixes into one routing entry
- Hierarchy is the key

# Hierarchical IP Address Allocation



- American Registry for Internet Numbers (ARIN)
- RIPE, APNIC, LACNIC, AfriNIC



## CIDR network prefix has variable length

|      |          |          |          |          |
|------|----------|----------|----------|----------|
|      | 128      | 143      | 137      | 144      |
| Addr | 10000000 | 10001111 | 10001001 | 10010000 |
|      | 255      | 255      | 255      | 0        |
| Mask | 11111111 | 11111111 | 11111111 | 00000000 |

- A network mask specifies the number of bits used to identify a network in an IP address.

# CIDR notation

- CIDR notation of an IP address:
  - 128.143.137.144/24
  - /24 is the prefix length. It states that the first 24 bits are the network prefix of the address (and the remaining 8 bits are available for specific host addresses)
- CIDR notation can nicely express blocks of addresses
  - An address block  
[128.195.0.0, 128.195.255.255]  
can be represented by an address prefix  
128.195.0.0/16
  - How many IP addresses are there in a /x address block?
    - $2^{(32-x)}$

# Output of ifconfig

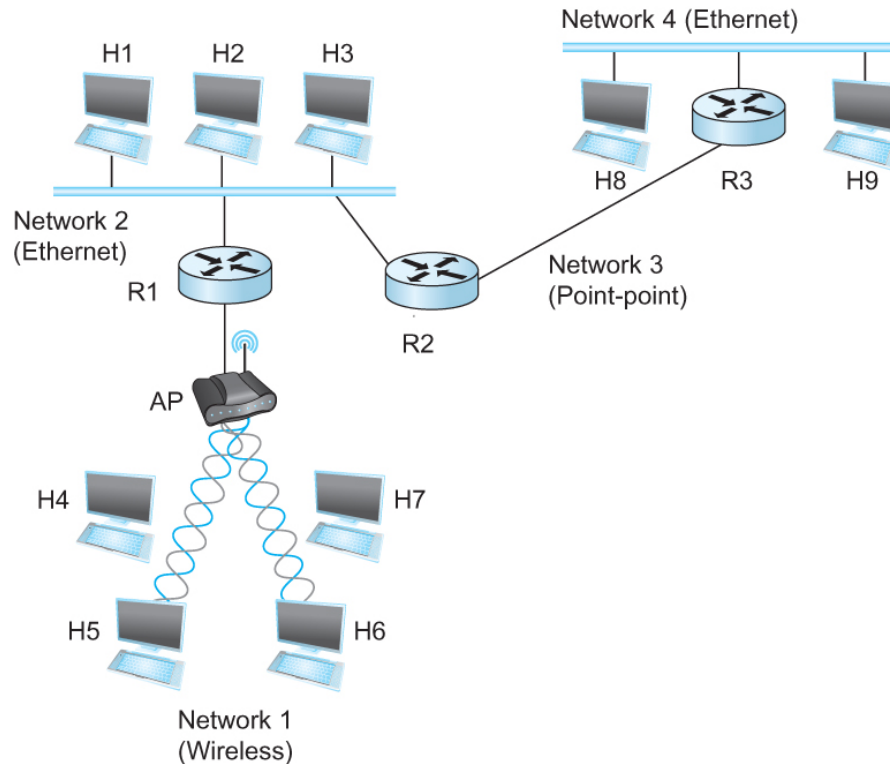
```
en0: flags=8863<UP,BROADCAST,SMART,RUNNING,SIMPLEX,MULTICAST> mtu 1500
    ether a8:66:7f:16:02:08
    inet6 fe80::10cf:731b:1d54:e775%en0 prefixlen 64 secured scopeid 0x5
    inet 10.194.131.251 netmask 0xffffe000 broadcast 10.194.159.255
    nd6 options=201<PERFORMNUD,DAD>
    media: autoselect
    status: active
```

# IP Forwarding

# Forwarding of IP datagrams

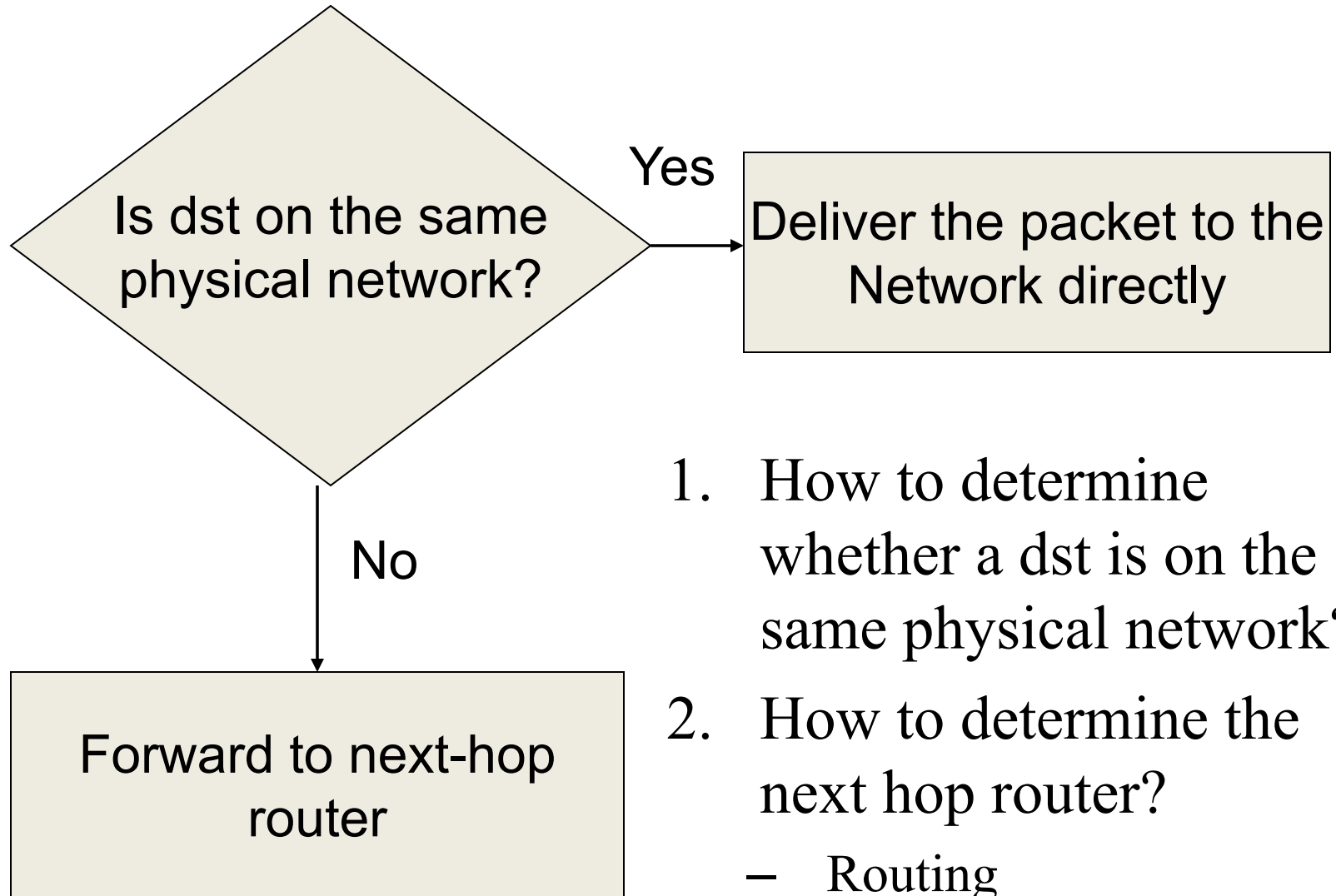
- There are two distinct processes to delivering IP datagrams:
  1. **Forwarding (data plane):** How to pass a packet from an input interface to the output interface?
  2. **Routing (control plane):** How to find and setup the forwarding tables?

# Forwarding basics



- Routers forward according to network prefixes
- All interfaces on the same network have the same network prefixes

# Forwarding algorithm



# Detailed forwarding algorithm

- If (networkNum == networkNum of one of my interfaces) then
  - Deliver packet over the interface
- Else
  - if (NetworkNum is in my forwarding table) then
    - Deliver to the NextHop router
  - Else
    - Deliver packet to the default router



# How does a host/router determine the network number of a destination address?

- Destination address & network mask = NetworkNumOfDestination
- If (NetworkNumOfDestination == my network Number) then
  - Send through my direct interfaces

# Forwarding table lookup

- **Forwarding table lookup:** Use the IP destination address as a key to search the routing table
- Result of the lookup is the IP address of a next hop router, and/or the name of a network interface

| Destination<br>address  | Next hop/<br>interface   |
|---|--|
| network prefix<br><i>or</i><br>host IP address<br><i>or</i><br>loopback address<br><i>or</i><br>default route | IP address of<br>next hop router<br><br><i>or</i><br><br>Name of a<br>network<br>interface |

# Type of forwarding table entries

- **Network route**
  - Destination addresses is a network address (e.g., 10.0.2.0/24)
  - Most entries are network routes
- **Host route**
  - Destination address is an interface address (e.g., 10.0.1.2/32)
  - Used to specify a separate route for certain hosts
- **Default route**
  - Used when no network or host route matches
- **Loopback address**
  - Routing table for the loopback address (127.0.0.1)
  - The next hop lists the loopback (lo0) interface as outgoing interface

# Unified forwarding algorithm

- Observation:
  - A directly physical network can be an entry in the forwarding table
  - A default route can be an entry
- 1. Look up destination address in the forwarding table using longest prefix match
- 2. Forward the packet to the next hop indicated by the matched entry

# The longest prefix matching algorithm

1. Search for a match on all 32 bits
2. Search for a match for 31 bits
- ....
32. Search for a match on 0 bits

Host route, loopback entry

→ 32-bit prefix match

Default route is represented as 0.0.0.0/0

→ 0-bit prefix match

# Why longest prefix match?

- Longest → smallest network
- Network prefixes may be aggregated

# Example

**128.143.71.21**



| Destination address | Next hop |
|---------------------|----------|
| 10.0.0.0/8          | eth0     |
| 128.143.0.0/16      | R2       |
| 128.143.64.0/20     | R3       |
| 128.143.192.0/20    | R3       |
| 128.143.71.0/24     | R4       |
| 128.143.71.55/32    | R3       |
| 0.0.0.0/0 (default) | R5       |



**The longest prefix match for  
128.143.71.21 is for 24 bits  
with entry 128.143.71.0/24**

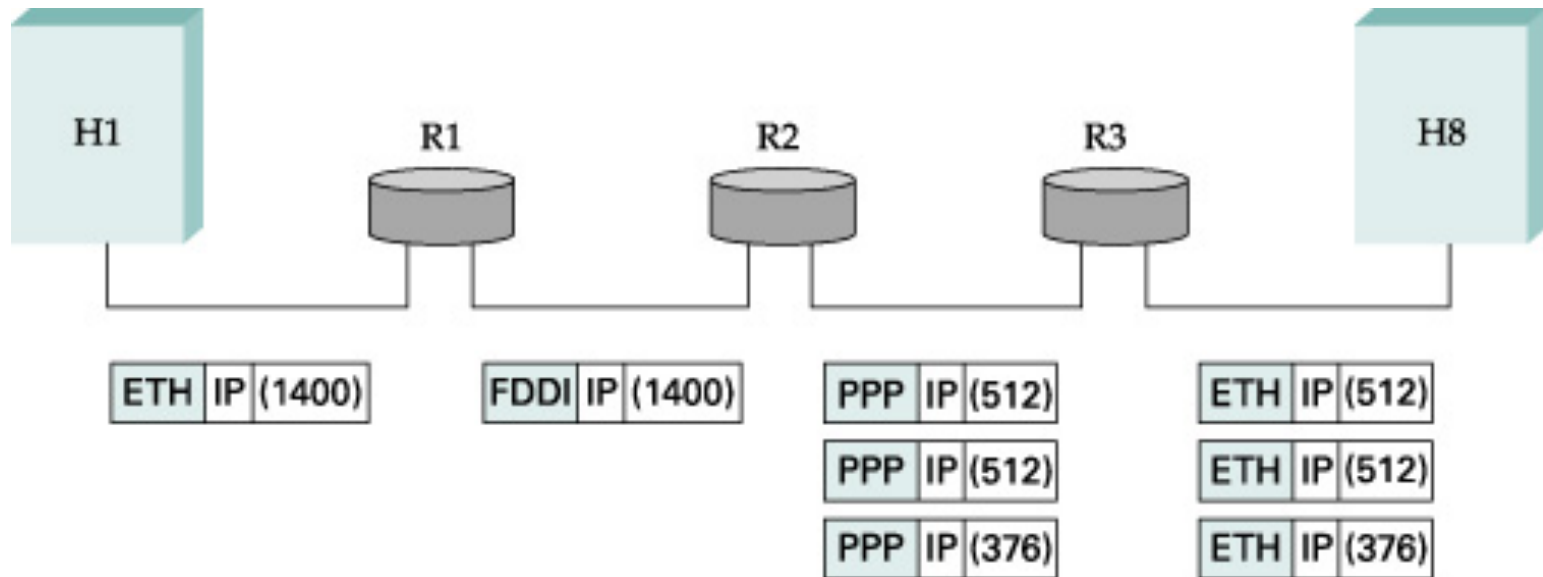
**Datagram will be sent to R4**

# Fragmentation and Reassembly

(not required for Lab 2)

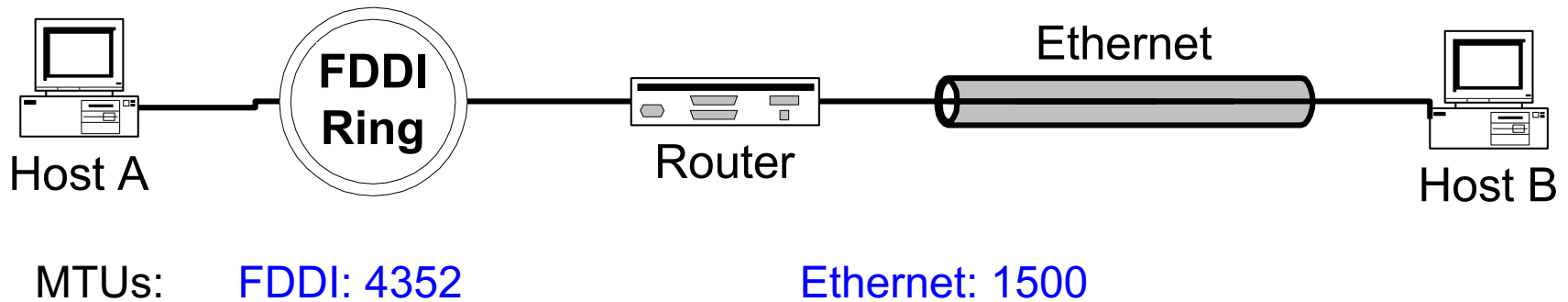


# Different networks have different Maximum Transmission Units (MTUs)



# IP Fragmentation and Reassembly

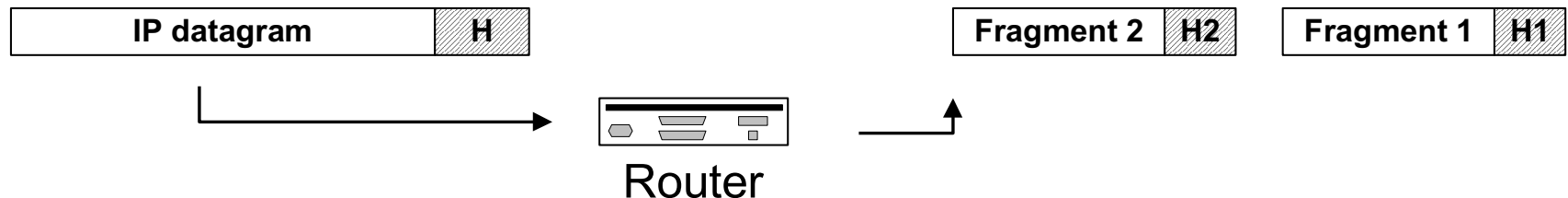
- What if the size of an IP datagram exceeds the MTU?  
IP datagram is fragmented into smaller units.
- What if the route contains networks with different MTUs?



- **Fragmentation:**
  - IP router splits the datagram into several datagrams

# Design question: Where is Fragmentation/reassembly done?

- Fragmentation can be done at the sender or at intermediate routers
- The same datagram can be fragmented several times.
- Reassembly of original datagram is only done at destination hosts !! (why?)



# What's involved in Fragmentation?

- The following fields in the IP header are involved:

|                    |               |          |     |                         |        |                 |
|--------------------|---------------|----------|-----|-------------------------|--------|-----------------|
| version            | header length | DS       | ECN | total length (in bytes) |        |                 |
| Identification     |               |          | 0   | D<br>F                  | M<br>F | Fragment offset |
| time-to-live (TTL) |               | protocol |     | header checksum         |        |                 |

- Identification
  - When a datagram is fragmented, the identification is the same in all fragments
  - Used to reassemble the original packet
- Flags
  - DF bit is set: datagram cannot be fragmented and must be discarded if MTU is too small
    - ICMP sent
  - MF bit:
    - 1: this is not the last fragment
    - 0: last fragment

# What's involved in Fragmentation?

- The following fields in the IP header are involved:

|                          |               |    |                 |                         |        |        |
|--------------------------|---------------|----|-----------------|-------------------------|--------|--------|
| version                  | header length | DS | ECN             | total length (in bytes) |        |        |
| Identification           |               |    |                 | 0                       | D<br>F | M<br>F |
| Fragment offset (13-bit) |               |    |                 |                         |        |        |
| time-to-live (TTL)       | protocol      |    | header checksum |                         |        |        |

- Fragment offset*

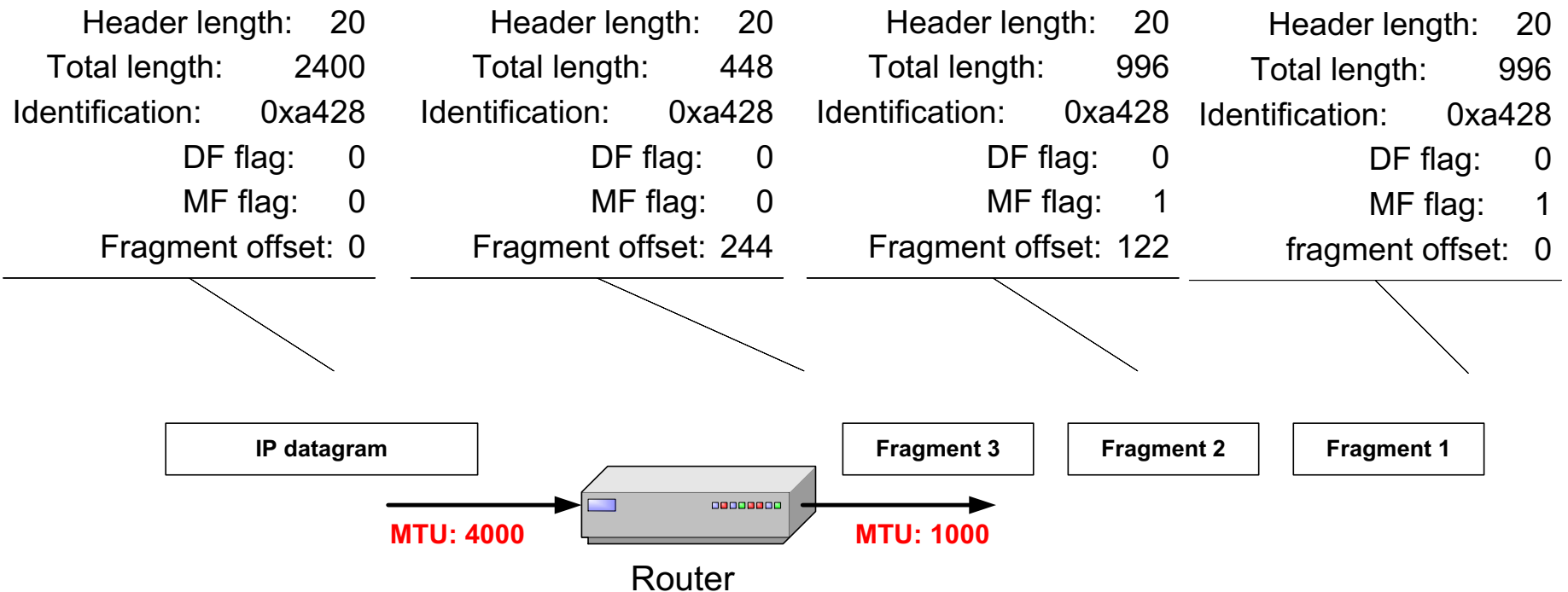
- Offset of the payload of the current fragment in the original datagram in units of 8 bytes
  - Why?
  - Because the field is only 13 bits long, while the total length is 16 bits.

- Total length*

- Total length of the current fragment

# Example of Fragmentation

- A datagram with size 2400 bytes must be fragmented according to an MTU limit of 1000 bytes



# Determining the length of fragments

- Maximum payload length =  $1000 - 20 = 980$  bytes
- Offset specifies the bytes in multiple of 8 bytes. So the payload must be a multiple of 8 bytes.
- $980 - 980 \% 8 = 976$  (the largest number that is less than 980 and divisible by 8)
- The payload for the first fragment is 976 and has bytes 0 ~ 975 of the original IP datagram. The offset is 0.
- The payload for the second fragment is 976 and has bytes 976 ~ 1951 of the original IP datagram. The offset is  $976 / 8 = 122$ .
- The payload of the last fragment is  $2400 - 976 * 2 = 448$  bytes and has bytes 1952 ~ 2400 of the original IP datagram. The offset is 244.
- Total length of three fragments:  $996 + 996 + 448 = 2440 > 2400$ 
  - Why?
  - Two additional IP headers.

# Path MTU discovery

- Fragmentation slows down the router
- → should be done by end hosts
- How does a sender know the MTU of a path?
  - A host only knows the MTU of its links
- Solution
  - send large packets with DF set
  - If receive ICMP Fragmentation needed messages, reduce maximum segment size



# Summary

- History of IP
- IP header format
- IP addressing
- IP forwarding
  - Forwarding algorithm
  - Fragmentation