

# Networks and Distributed Systems

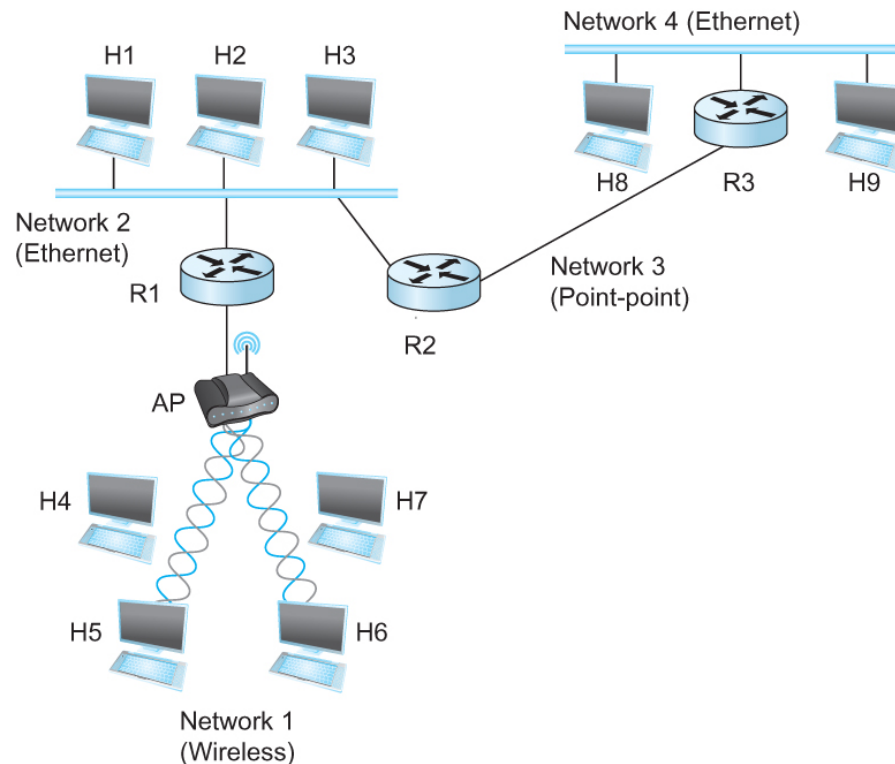
## Lecture 9 – Internetworking (IP)

# Outline

- Switching and Bridging
- Basic Internetworking (IP)
- Routing

# Internetworking

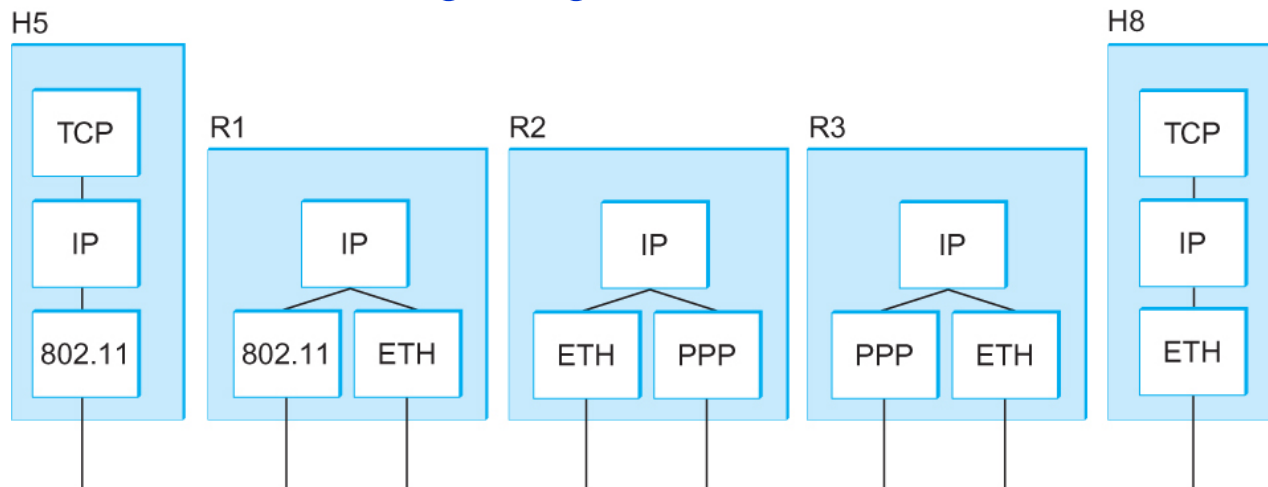
- What is internetwork
  - An arbitrary collection of networks interconnected to provide some sort of host-host to packet delivery service



A simple internetwork where H represents hosts and R represents routers

# Internetworking

- What is IP
  - IP stands for Internet Protocol
  - Key tool used today to build scalable, heterogeneous internetworks
  - It runs on all the nodes in a collection of networks and defines the infrastructure that allows these nodes and networks to function as a single logical internetwork



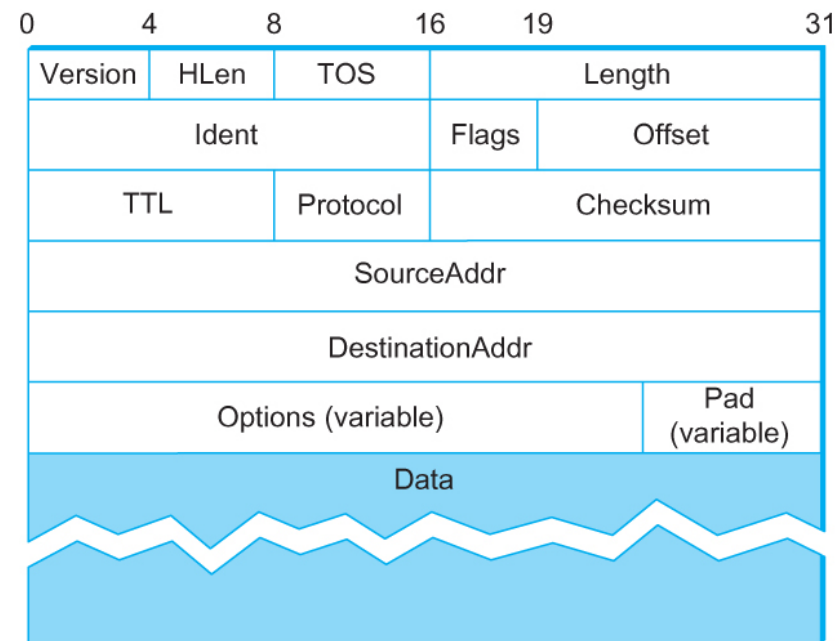
A simple internetwork showing the protocol layers

# IP Service Model

- Packet Delivery Model
  - Connectionless model for data delivery
  - Best-effort delivery (unreliable service)
    - packets are lost
    - packets are delivered out of order
    - duplicate copies of a packet are delivered
    - packets can be delayed for a long time
- Global Addressing Scheme
  - Provides a way to identify all hosts in the network

# Packet Format

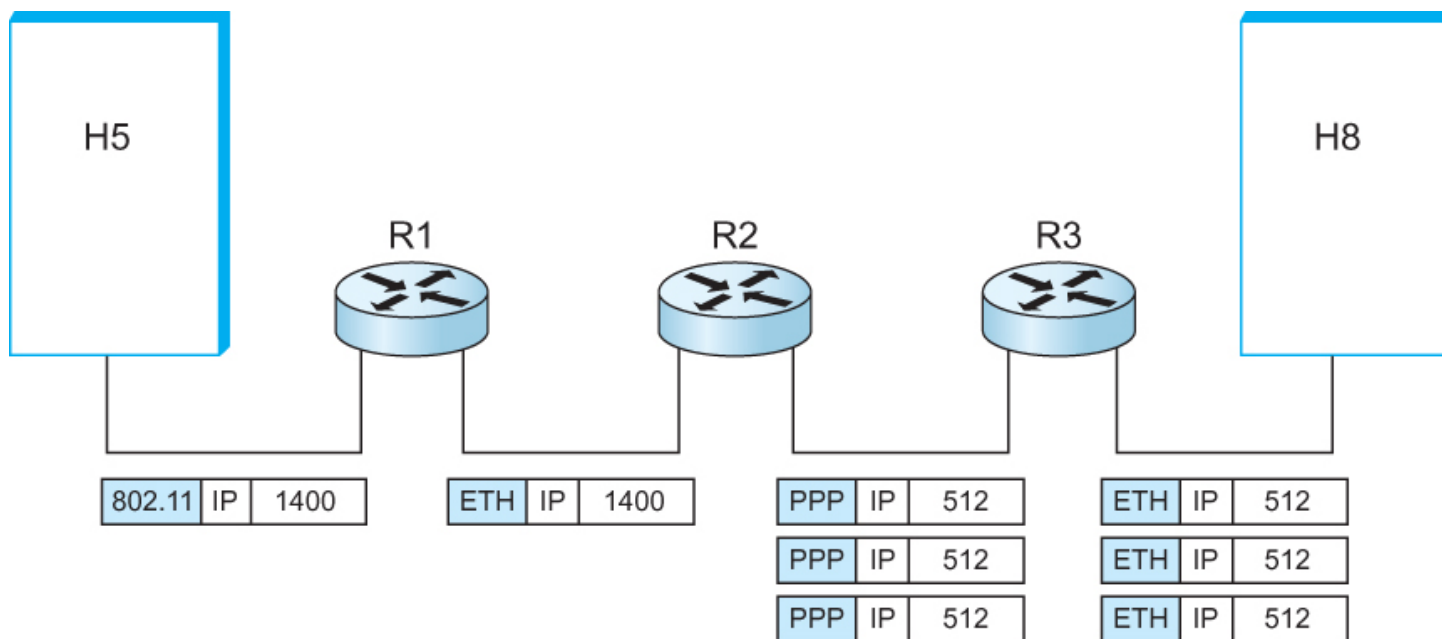
- Version (4): currently 4
- Hlen (4): number of 32-bit words in header
- TOS (8): type of service (not widely used)
- Length (16): number of bytes in this datagram
- Ident (16): used by fragmentation
- Flags/Offset (16): used by fragmentation
- TTL (8): number of hops this datagram has traveled
- Protocol (8): demux key (TCP=6, UDP=17)
- Checksum (16): of the header only
- DestAddr & SrcAddr (32)



# IP Fragmentation and Reassembly

- Each network has some MTU (Maximum Transmission Unit)
  - Ethernet (1500 bytes), FDDI (4500 bytes)
- Strategy
  - Fragmentation occurs in a router when it receives a datagram that it wants to forward over a network which has (MTU < datagram)
  - Reassembly is done at the receiving host
  - All the fragments carry the same identifier in the *Ident* field
  - Fragments are self-contained datagrams
  - IP does not recover from missing fragments

# IP Fragmentation and Reassembly



IP datagrams traversing the sequence of physical networks



# IP Fragmentation and Reassembly

(a)

|                 |  |  |   |            |
|-----------------|--|--|---|------------|
| Start of header |  |  |   |            |
| Ident = x       |  |  | 0 | Offset = 0 |
| Rest of header  |  |  |   |            |
| 1400 data bytes |  |  |   |            |

(b)

|                 |  |  |   |            |
|-----------------|--|--|---|------------|
| Start of header |  |  |   |            |
| Ident = x       |  |  | 1 | Offset = 0 |
| Rest of header  |  |  |   |            |
| 512 data bytes  |  |  |   |            |

|                 |  |  |   |             |
|-----------------|--|--|---|-------------|
| Start of header |  |  |   |             |
| Ident = x       |  |  | 1 | Offset = 64 |
| Rest of header  |  |  |   |             |
| 512 data bytes  |  |  |   |             |

|                 |  |  |   |              |
|-----------------|--|--|---|--------------|
| Start of header |  |  |   |              |
| Ident = x       |  |  | 0 | Offset = 128 |
| Rest of header  |  |  |   |              |
| 376 data bytes  |  |  |   |              |

Header fields used in IP fragmentation. (a) Unfragmented packet; (b) fragmented packets.

# Global Addresses

## ■ Properties

- globally unique
- hierarchical: network + host
- 4 Billion IP address, half are A type,  $\frac{1}{4}$  is B type, and  $\frac{1}{8}$  is C type

## ■ Format



## ■ Dot notation

- 10.3.2.4
- 128.96.33.81
- 192.12.69.77

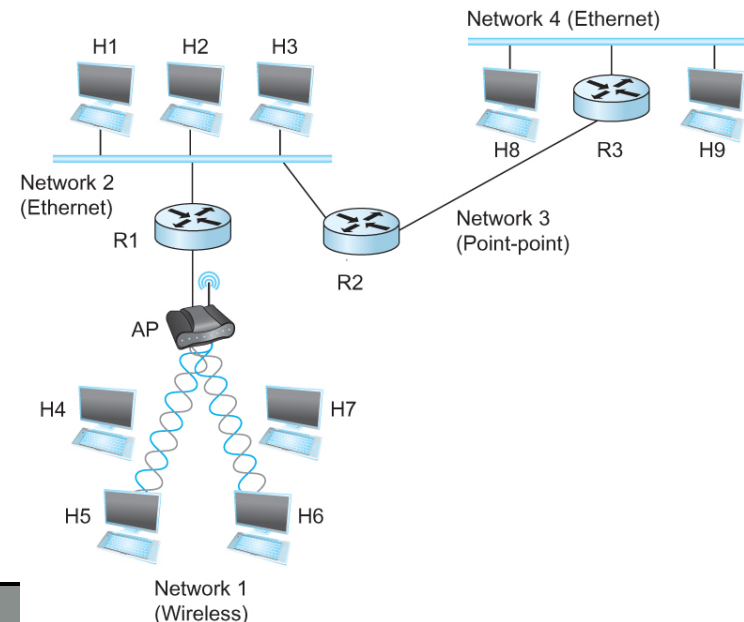
# IP Datagram Forwarding

## ■ Strategy

- every datagram contains destination's address
- if directly connected to destination network, then forward to host
- if not directly connected to destination network, then forward to some router
- forwarding table maps network number into next hop
- each host has a default router
- each router maintains a forwarding table

## ■ Example (router R2)

| NetworkNum | NextHop     |
|------------|-------------|
| 1          | R1          |
| 2          | Interface 1 |
| 3          | Interface 0 |
| 4          | R3          |



# IP Datagram Forwarding

## ■ Algorithm

```
if (NetworkNum of destination = NetworkNum of one of my  
    interfaces) then  
    deliver packet to destination over that interface  
else  
    if (NetworkNum of destination is in my forwarding table)  
    then  
        deliver packet to NextHop router  
    else  
        deliver packet to default router
```

For a host with only one interface and only a default router in its forwarding table, this simplifies to

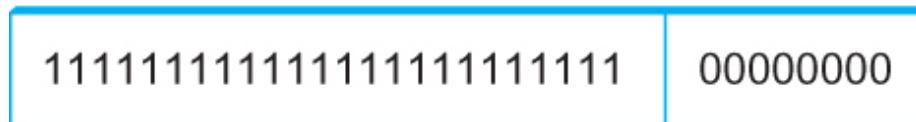
```
if (NetworkNum of destination = my NetworkNum) then  
    deliver packet to destination directly  
else  
    deliver packet to default router
```

# Subnetting

- Add another level to address/routing hierarchy: *subnet*
- *Subnet masks* define variable partition of host part of class A and B addresses
- Subnets visible only within site



Class B address

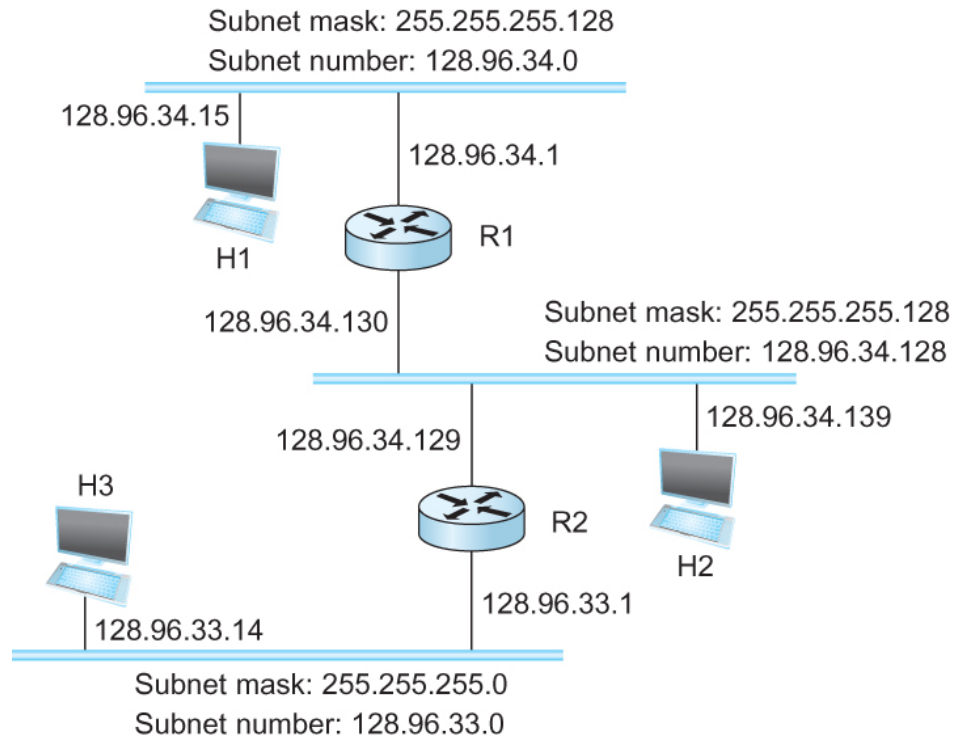


Subnet mask (255.255.255.0)



Subnetted address

# Subnetting



## ■ Forwarding Table at Router R1

| SubnetNumber  | SubnetMask      | NextHop     |
|---------------|-----------------|-------------|
| 128.96.34.0   | 255.255.255.128 | Interface 0 |
| 128.96.34.128 | 255.255.255.128 | Interface 1 |
| 128.96.33.0   | 255.255.255.0   | R2          |

# Subnetting

## Forwarding Algorithm

```
D = destination IP address
for each entry < SubnetNum, SubnetMask, NextHop>
    D1 = SubnetMask & D
    if D1 = SubnetNum
        if NextHop is an interface
            deliver datagram directly to destination
        else
            deliver datagram to NextHop (a router)
```

# Subnetting

## Notes

- Would use a default router if nothing matches
- Not necessary for all ones in subnet mask to be contiguous
- Can put multiple subnets on one physical network
- Subnets not visible from the rest of the Internet



# Classless Addressing

- Classless Inter-Domain Routing
  - A technique that addresses two scaling concerns in the Internet
    - The growth of backbone routing table as more and more network numbers need to be stored in them
    - Potential exhaustion of the 32-bit address space
  - Address assignment efficiency
    - Arises because of the IP address structure with class A, B, and C addresses
    - Forces us to hand out network address space in fixed-size chunks of three very different sizes
      - A network with two hosts needs a class C address
        - Address assignment efficiency =  $2/255 = 0.78$
      - A network with 256 hosts needs a class B address
        - Address assignment efficiency =  $256/65535 = 0.39$

# Classless Addressing

- Exhaustion of IP address space centers on exhaustion of the class B network numbers
- Solution
  - Say “NO” to any Autonomous System (AS) that requests a class B address unless they can show a need for something close to 64K addresses
  - Instead give them an appropriate number of class C addresses
  - For any AS with at least 256 hosts, we can guarantee an address space utilization of at least 50%
- What is the problem with this solution?

# Classless Addressing

- Problem with this solution
  - Excessive storage requirement at the routers.
- If a single AS has, say 16 class C network numbers assigned to it,
  - Every Internet backbone router needs 16 entries in its routing tables for that AS
  - This is true, even if the path to every one of these networks is the same
- If we had assigned a class B address to the AS
  - The same routing information can be stored in one entry
  - Efficiency =  $16 \times 255 / 65,536 = 6.2\%$

# Classless Addressing

- CIDR tries to balance the desire to minimize the number of routes that a router needs to know against the need to hand out addresses efficiently.
- CIDR uses aggregate routes
  - Uses a single entry in the forwarding table to tell the router how to reach a lot of different networks
  - Breaks the rigid boundaries between address classes

# Classless Addressing

- Consider an AS with 16 class C network numbers.
- Instead of handing out 16 addresses at random, hand out a block of contiguous class C addresses
- Suppose we assign the class C network numbers from 192.4.16 through 192.4.31
- Observe that top 20 bits of all the addresses in this range are the same (1 1000000 00000100 0001)
  - We have created a 20-bit network number (which is in between class B network number and class C number)
- Requires to hand out blocks of class C addresses that share a common prefix

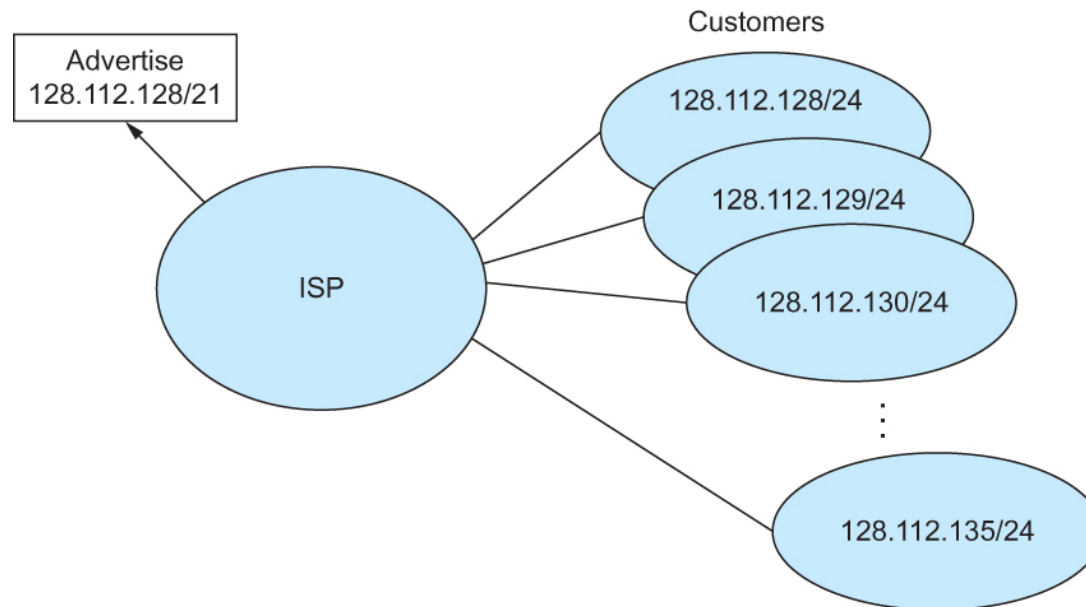
# Classless Addressing

- Requires to hand out blocks of class C addresses that share a common prefix
- The convention is to place a /X after the prefix where X is the prefix length in bits
- For example, the 20-bit prefix for all the networks 192.4.16 through 192.4.31 is represented as 192.4.16/20
- By contrast, if we wanted to represent a single class C network number, which is 24 bits long, we would write it 192.4.16/24

# Classless Addressing

- How do the routing protocols handle this classless addresses
  - It must understand that the network number may be of any length
- Represent network number with a single pair  
`<length, value>`
- All routers must understand CIDR addressing

# Classless Addressing



Route aggregation with CIDR



# IP Forwarding Revisited

- IP forwarding mechanism assumes that it can find the network number in a packet and then look up that number in the forwarding table
- We need to change this assumption in case of CIDR
- CIDR means that prefixes may be of any length, from 2 to 32 bits

# IP Forwarding Revisited

- It is also possible to have prefixes in the forwarding tables that overlap
  - Some addresses may match more than one prefix
- For example, we might find both 171.69 (a 16 bit prefix) and 171.69.10 (a 24 bit prefix) in the forwarding table of a single router
- A packet destined to 171.69.10.5 clearly matches both prefixes.
  - The rule is based on the principle of “longest match”
    - 171.69.10 in this case
- A packet destined to 171.69.20.5 would match 171.69 and not 171.69.10

# Address Translation Protocol (ARP)

- Map IP addresses into physical addresses
  - destination host
  - next hop router
- Techniques
  - encode physical address in host part of IP address
  - table-based
- ARP (Address Resolution Protocol)
  - table of IP to physical address bindings
  - broadcast request if IP address not in table
  - target machine responds with its physical address
  - table entries are discarded if not refreshed

# ARP Packet Format

|                                |           |                                |    |
|--------------------------------|-----------|--------------------------------|----|
| 0                              | 8         | 16                             | 31 |
| Hardware type = 1              |           | ProtocolType = 0x0800          |    |
| HLen = 48                      | PLen = 32 | Operation                      |    |
| SourceHardwareAddr (bytes 0–3) |           |                                |    |
| SourceHardwareAddr (bytes 4–5) |           | SourceProtocolAddr (bytes 0–1) |    |
| SourceProtocolAddr (bytes 2–3) |           | TargetHardwareAddr (bytes 0–1) |    |
| TargetHardwareAddr (bytes 2–5) |           |                                |    |
| TargetProtocolAddr (bytes 0–3) |           |                                |    |

- HardwareType: type of physical network (e.g., Ethernet)
- ProtocolType: type of higher layer protocol (e.g., IP)
- HLEN & PLEN: length of physical and protocol addresses
- Operation: request or response
- Source/Target Physical/Protocol addresses

# Host Configurations

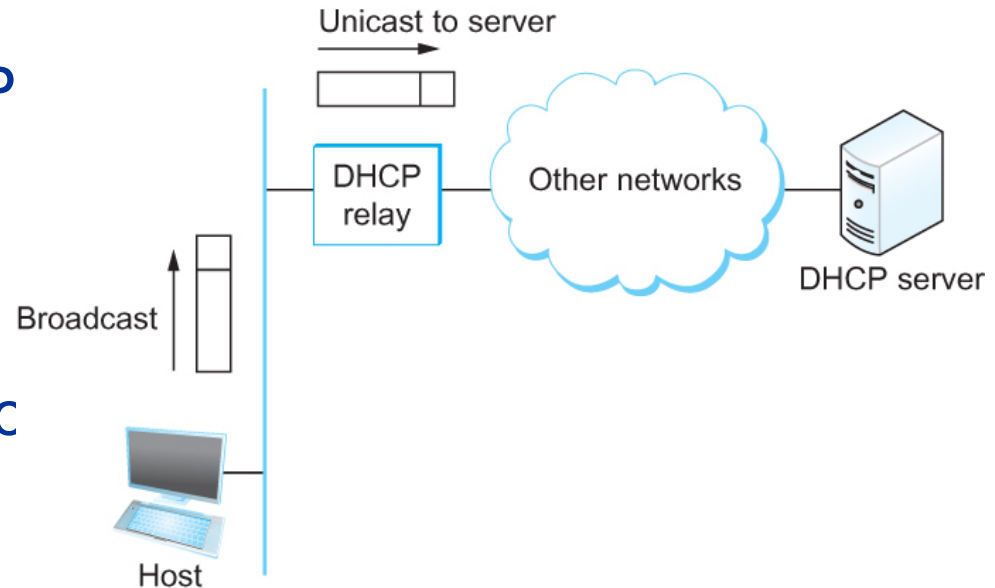
- Notes
  - Ethernet addresses are configured into network by manufacturer and they are unique
  - IP addresses must be unique on a given internetwork but also must reflect the structure of the internetwork
  - Most host Operating Systems provide a way to manually configure the IP information for the host
  - Drawbacks of manual configuration
    - A lot of work to configure all the hosts in a large network
    - Configuration process is error-prone
  - Automated Configuration Process is required

# Dynamic Host Configuration Protocol (DHCP)

- DHCP server is responsible for providing configuration information to hosts
- There is at least one DHCP server for an administrative domain
- DHCP server maintains a pool of available addresses

# DHCP

- Newly booted or attached host sends DHCPDISCOVER message to a special IP address (255.255.255.255)
- DHCP relay agent unicasts the message to DHCP server and waits for the response



# Internet Control Message Protocol (ICMP)

- Defines a collection of error messages that are sent back to the source host whenever a router or host is unable to process an IP datagram successfully
  - Destination host unreachable due to link /node failure
  - Reassembly process failed
  - TTL had reached 0 (so datagrams don't cycle forever)
  - IP header checksum failed
- ICMP-Redirect
  - From router to a source host
  - With a better route information



# Internet Control Message Protocol (ICMP)

- Defines a collection of error messages that are sent back to the source host whenever a router or host is unable to process an IP datagram successfully
  - Destination host unreachable due to link /node failure
  - Reassembly process failed
  - TTL had reached 0 (so datagrams don't cycle forever)
  - IP header checksum failed
- ICMP-Redirect
  - From router to a source host
  - With a better route information