

# Chapter 4

## Network Layer

Reti di Elaboratori

Corso di Laurea in Informatica

Università degli Studi di Roma "La Sapienza"

Canale A-L

Prof.ssa Chiara Petrioli

Parte di queste slide sono state prese dal materiale associato al libro  
*Computer Networking: A Top Down Approach*, 5th edition.

All material copyright 1996-2009

J.F Kurose and K.W. Ross, All Rights Reserved

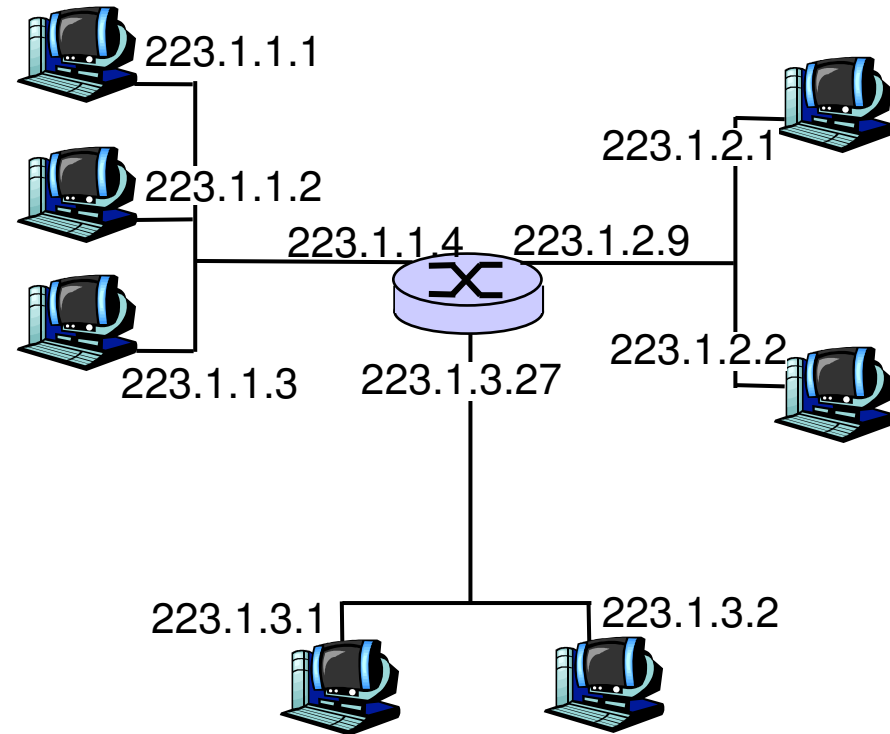
Thanks also to Antonio Capone, Politecnico di Milano, Giuseppe Bianchi and  
Francesco LoPresti, Un. di Roma Tor Vergata

# Chapter 4: Network Layer

- r 4.1 Introduction
- r 4.2 Virtual circuit and datagram networks
- r 4.3 What's inside a router
- r 4.4 IP: Internet Protocol
  - m Datagram format
  - m IPv4 addressing
  - m ICMP
  - m IPv6
- r 4.5 Routing algorithms
  - m Link state
  - m Distance Vector
  - m Hierarchical routing
- r 4.6 Routing in the Internet
  - m RIP
  - m OSPF
  - m BGP
- r 4.7 Broadcast and multicast routing

# IP Addressing: introduction

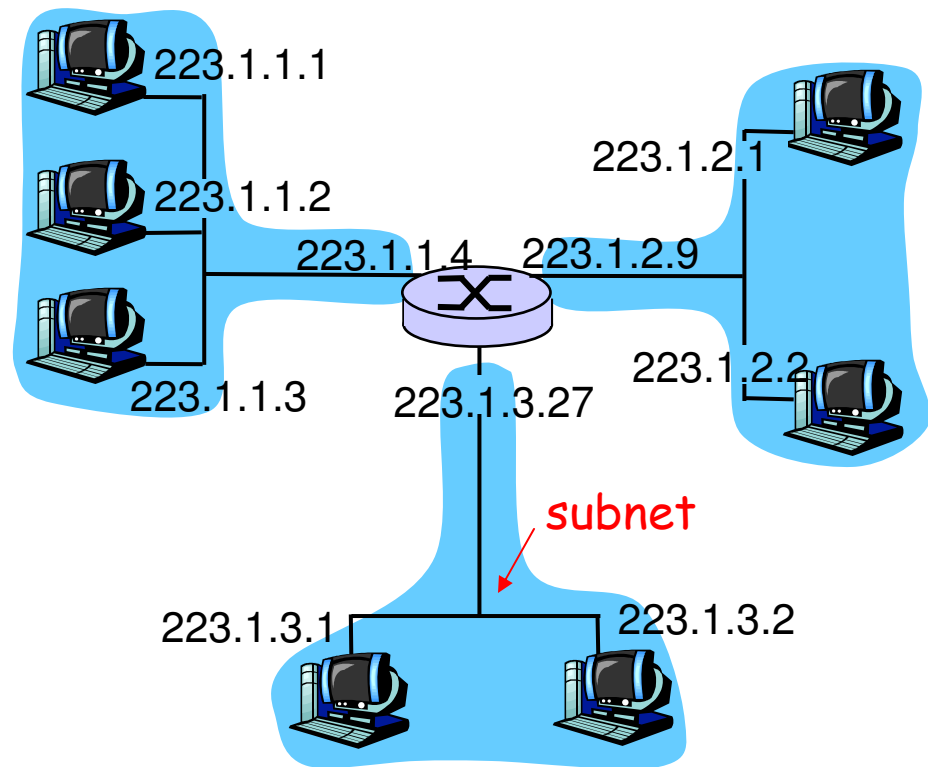
- r IP address: 32-bit identifier for host, router *interface*
- r *interface*: connection between host/router and physical link
  - m router's typically have multiple interfaces
  - m host typically has one interface
  - m IP addresses associated with each interface



$$223.1.1.1 = \underbrace{11011111}_{223} \underbrace{00000001}_1 \underbrace{00000001}_1 \underbrace{00000001}_1$$

# Subnets

- r IP address:
  - m subnet part (high order bits)
  - m host part (low order bits)
- r *What's a subnet?*
  - m device interfaces with same subnet part of IP address
  - m can physically reach each other without intervening router

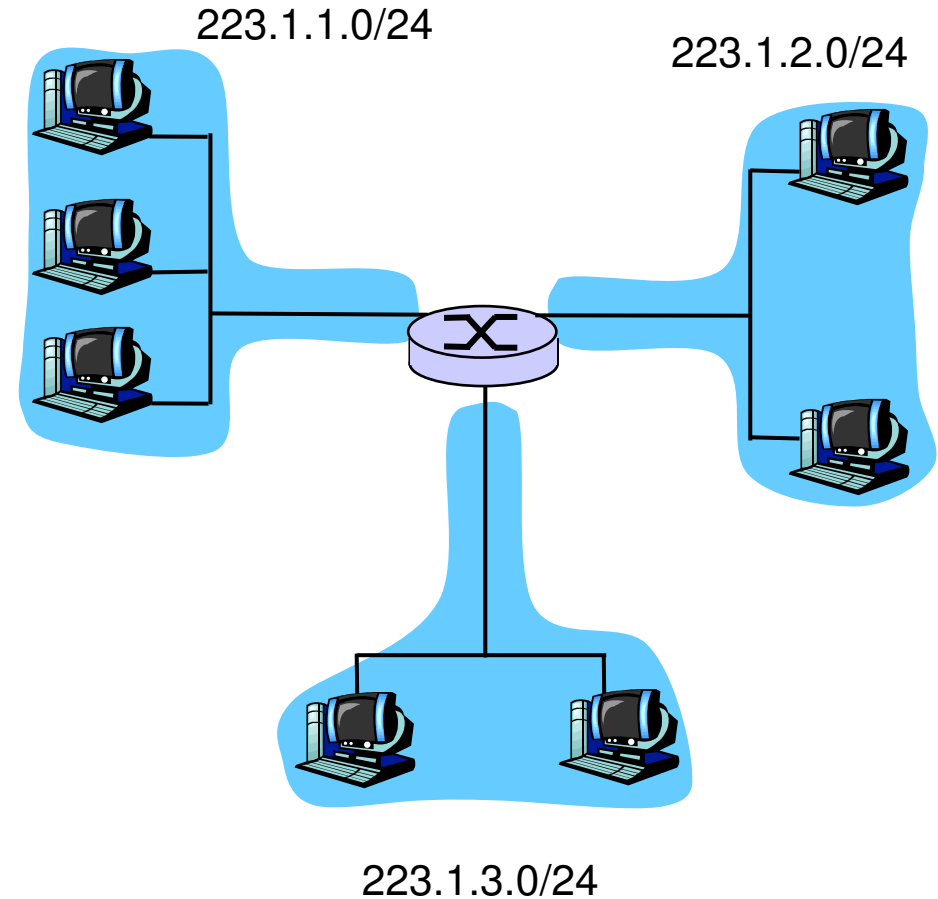


network consisting of 3 subnets

# Subnets

## Recipe

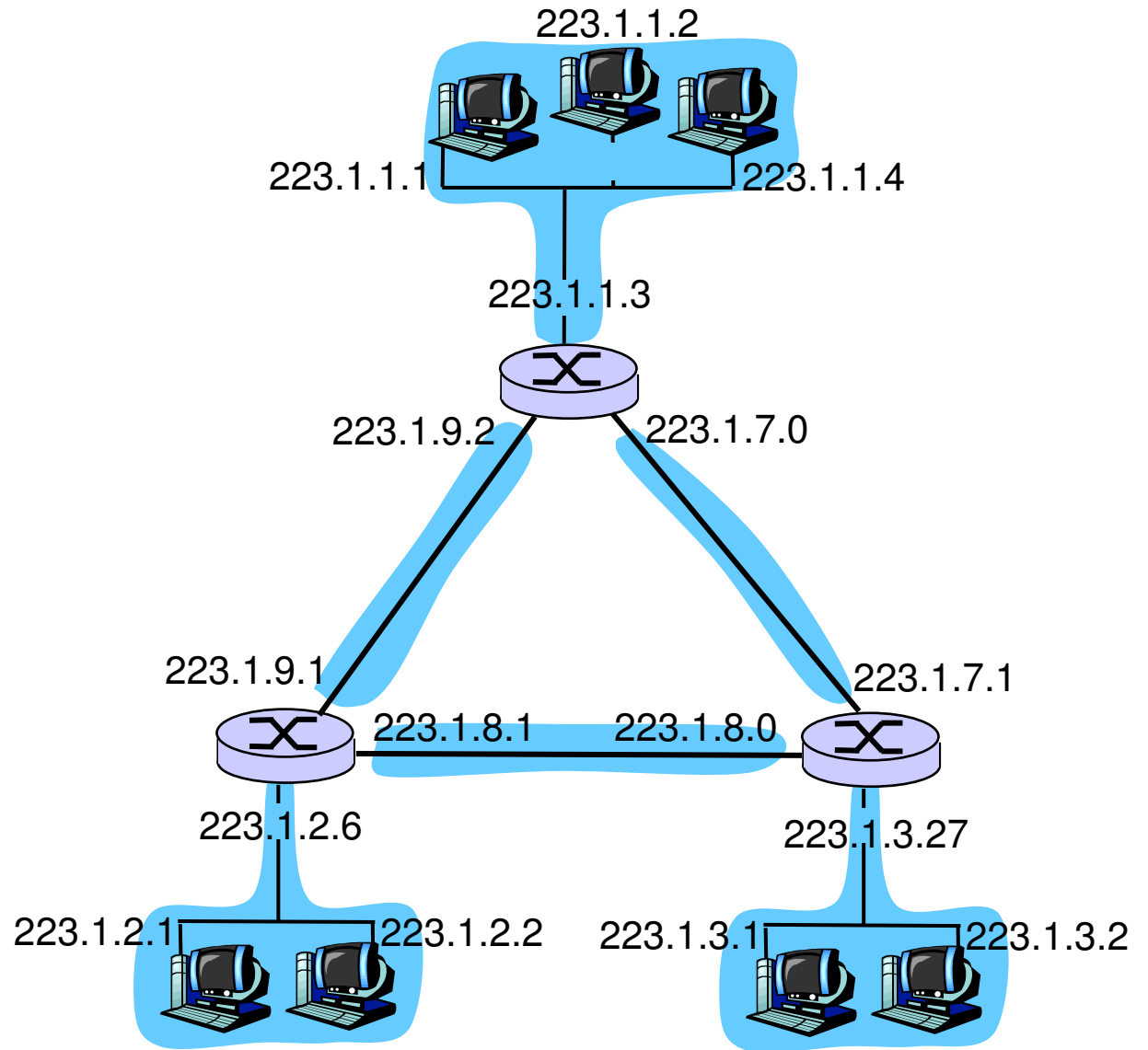
- r To determine the subnets, detach each interface from its host or router, creating islands of isolated networks. Each isolated network is called a **subnet**.



Subnet mask: /24

# Subnets

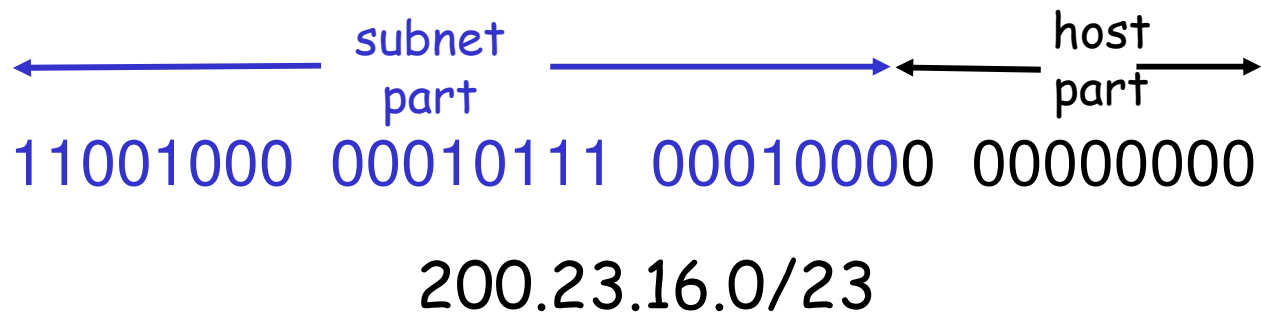
How many?



# IP addressing: CIDR

## CIDR: Classless InterDomain Routing

- m subnet portion of address of arbitrary length
- m address format: **a.b.c.d/x**, where x is # bits in subnet portion of address



# IP addresses: how to get one?

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

- r hard-coded by system admin in a file
  - m Windows: control-panel->network->configuration->tcp/ip->properties
  - m UNIX: /etc/rc.config
- r **DHCP: Dynamic Host Configuration Protocol:** dynamically get address from as server
  - m "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 and "on")

Support for mobile users who want to join network (more shortly)

DHCP overview:

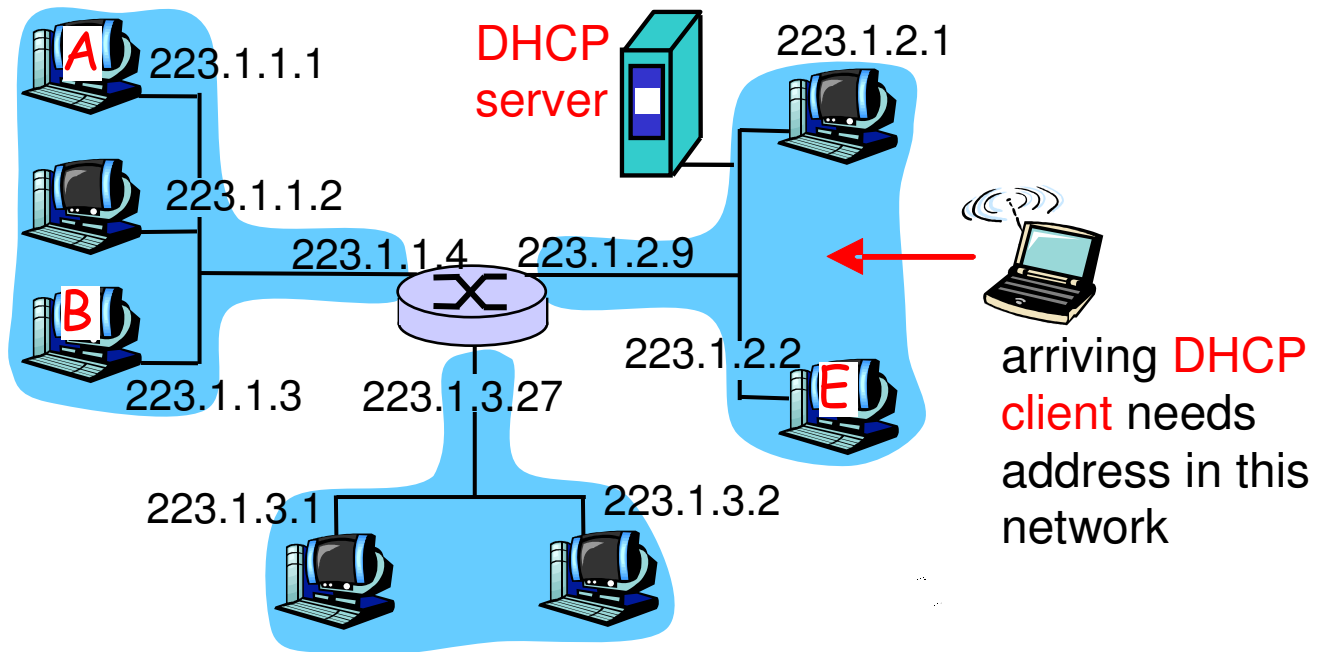
m host broadcasts "DHCP discover" msg [optional]

m DHCP server responds with "DHCP offer" msg [optional]

m host requests IP address: "DHCP request" msg

m DHCP server sends address: "DHCP ack" msg

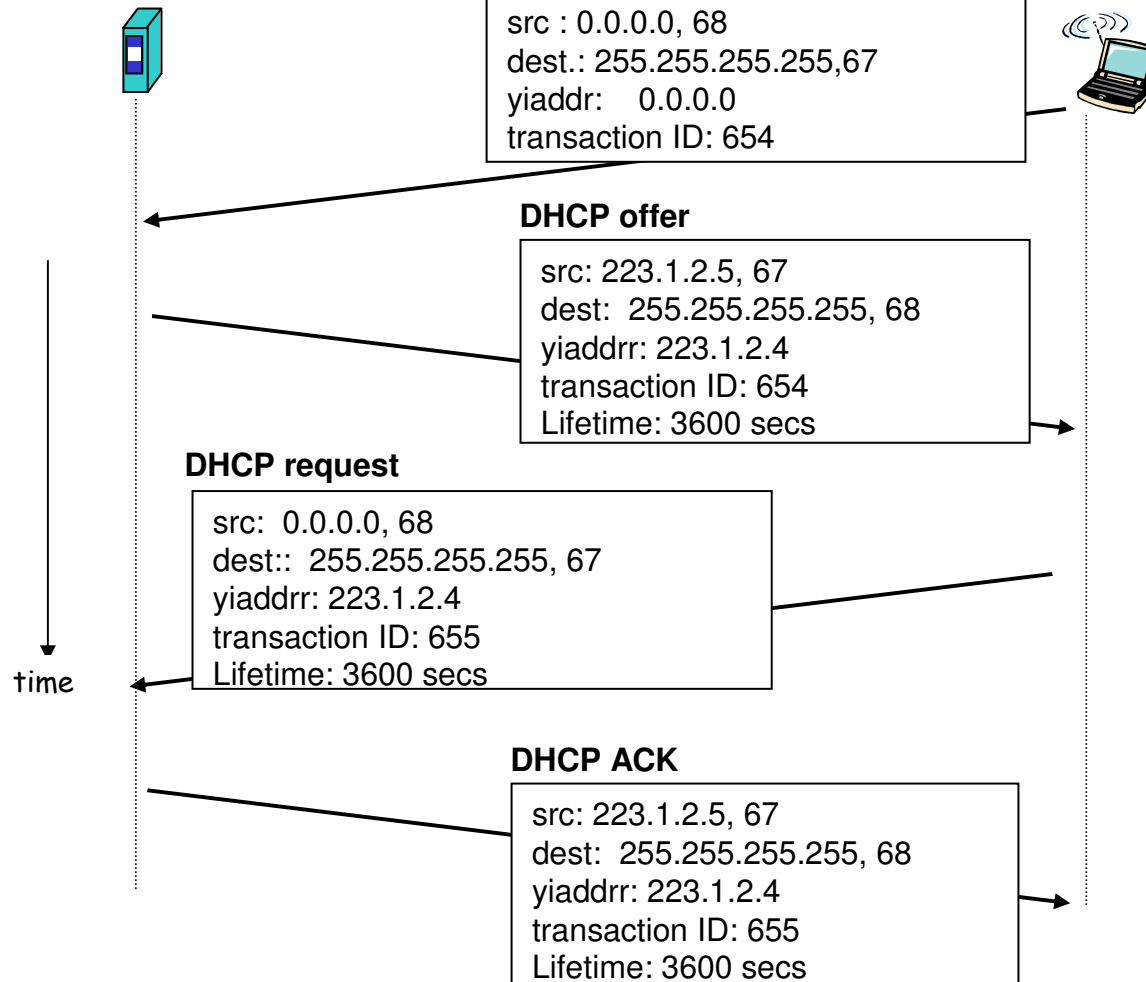
# DHCP client-server scenario



# DHCP client-server scenario

DHCP server: 223.1.2.5

arriving client



# DHCP

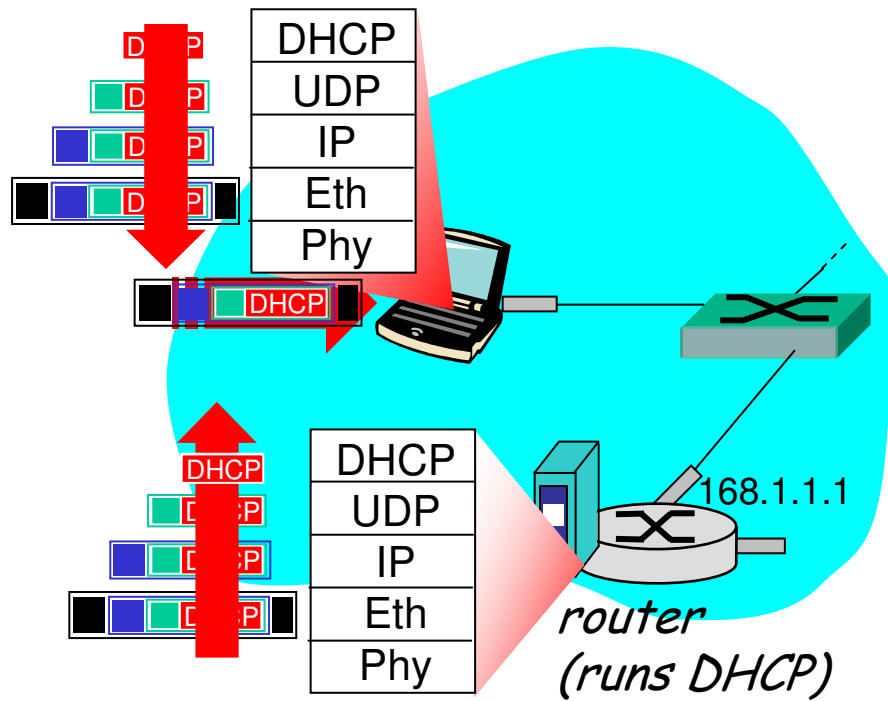
1. The client broadcasts a DHCPDISCOVER message on its local physical subnet. The DHCPDISCOVER message **MAY** include options that suggest values for the network address and lease duration. BOOTP relay agents may pass the message on to DHCP servers not on the same physical subnet.
2. Each server may respond with a DHCPOFFER message that includes an available network address in the 'yiaddr' field (and other configuration parameters in DHCP options). Servers need not reserve the offered network address, although the protocol will work more efficiently if the server avoids allocating the offered network address to another client. When allocating a new address, servers **SHOULD** check that the offered network address is not already in use; e.g., the server may probe the offered address with an ICMP Echo Request.

# DHCP: more than IP address

DHCP can return more than just allocated IP address on subnet:

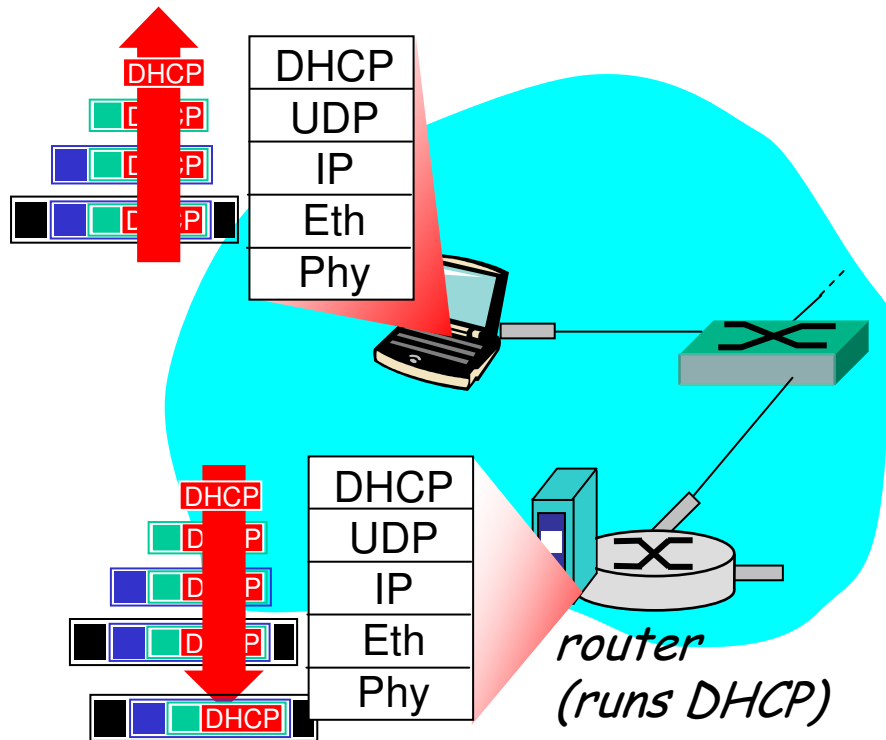
- m address of first-hop router for client
- m name and IP address of DNS sever
- m network mask (indicating network versus host portion of address)

# DHCP: example



- r connecting laptop needs its IP address, addr of first-hop router, addr of DNS server: use DHCP
- r DHCP request encapsulated in UDP, encapsulated in IP, encapsulated in 802.1 Ethernet
- r Ethernet frame broadcast (dest: FFFFFFFFFFFFFFFF) on LAN, received at router running DHCP server
- r Ethernet demux'ed to IP demux'ed, UDP demux'ed to DHCP

# DHCP: example



- r DCP server formulates DHCP ACK containing client's IP address, IP address of first-hop router for client, name & IP address of DNS server
- r encapsulation of DHCP server, frame forwarded to client, demux'ing up to DHCP at client
- r client now knows its IP address, name and IP address of DSN server, IP address of its first-hop router

# DHCP

**DHCPDISCOVER** - Client broadcast to locate available servers.

**DHCPOFFER** - Server to client in response to DHCPDISCOVER with offer of configuration parameters.

**DHCPREQUEST** - Client message to servers either (a) requesting offered parameters from one server and implicitly declining offers from all others, (b) confirming correctness of previously allocated address after, e.g., system reboot, or (c) extending the lease on a particular network address.

**DHCPACK** - Server to client with configuration parameters, including committed network address.

**DHCPNAK** - Server to client indicating client's notion of network address is incorrect (e.g., client has moved to new subnet) or client's lease as expired

**DHCPDECLINE** - Client to server indicating network address is already in use.

**DHCPRELEASE** - Client to server relinquishing network address and cancelling remaining lease.

**DHCPINFORM** - Client to server, asking only for local configuration parameters; client already has externally configured network address.



# DHCP: wireshark output (home LAN)

Message type: **Boot Request (1)**  
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  
    .....

request

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."**

reply

# IP addresses: how to get one?

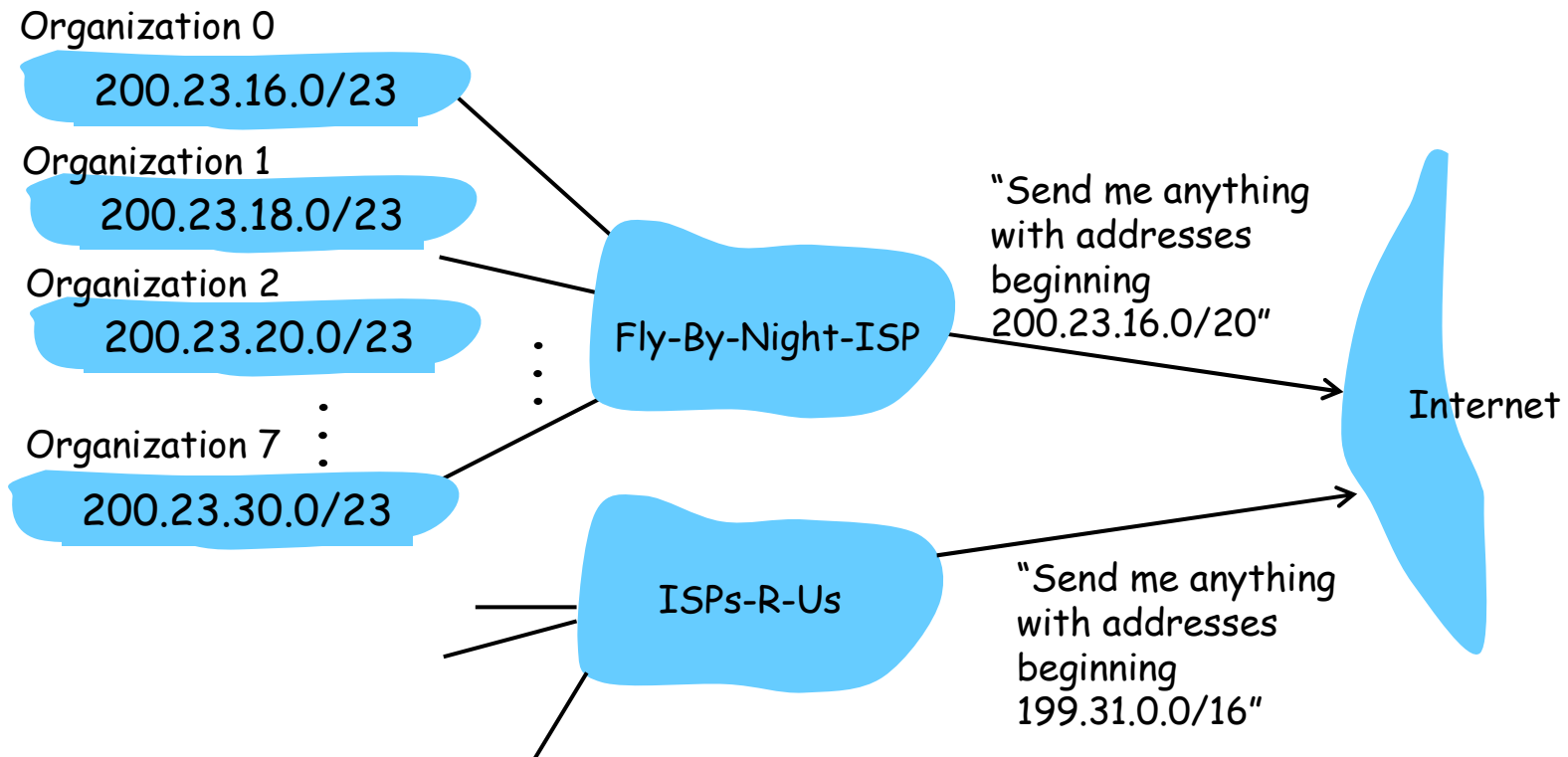
Q: How does *network* get subnet part of IP addr?

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

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

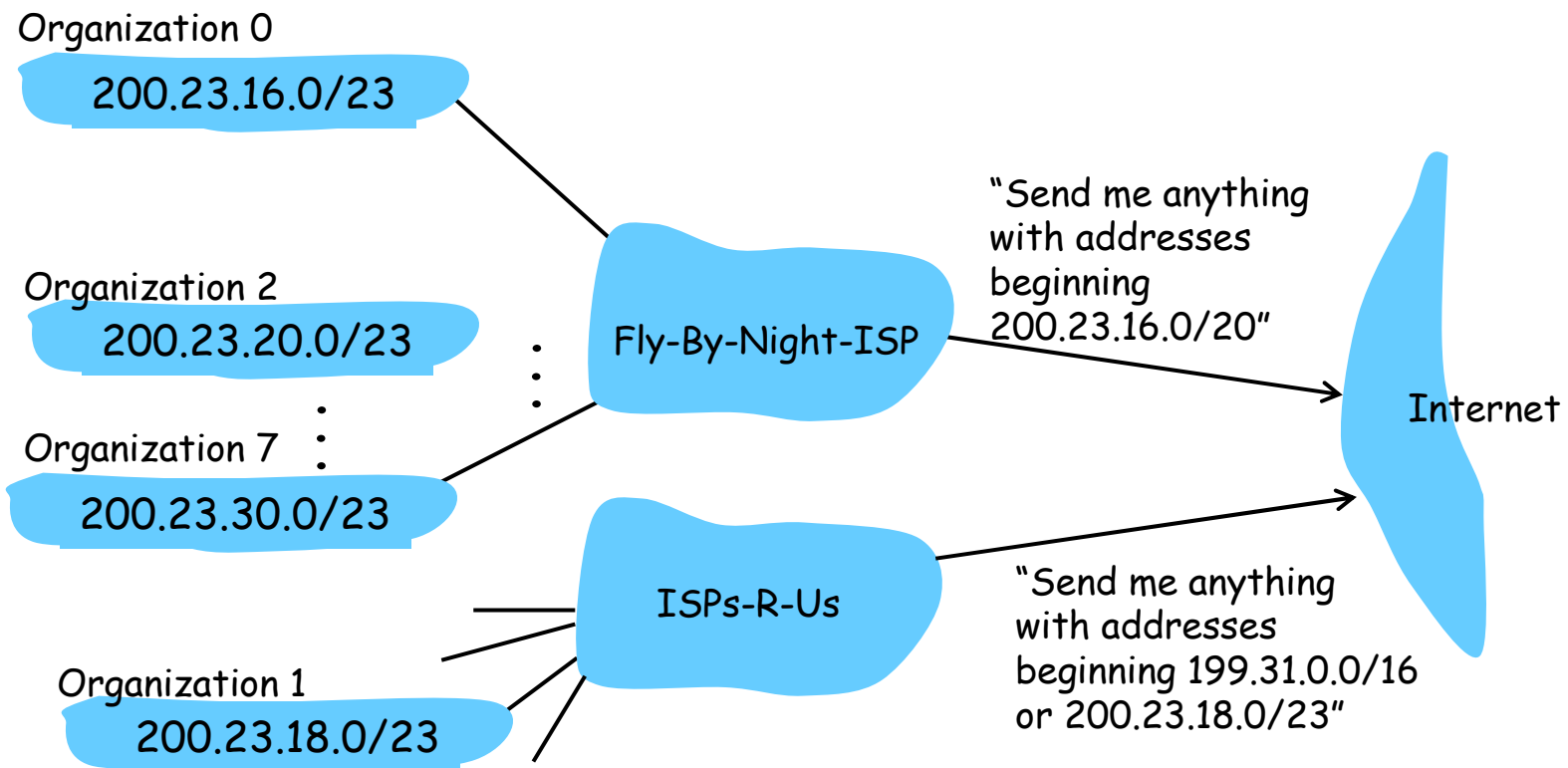
# Hierarchical addressing: route aggregation

Hierarchical addressing allows efficient advertisement of routing information:



# Hierarchical addressing: more specific routes

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



## IP addressing: the last word...

Q: How does an ISP get block of addresses?

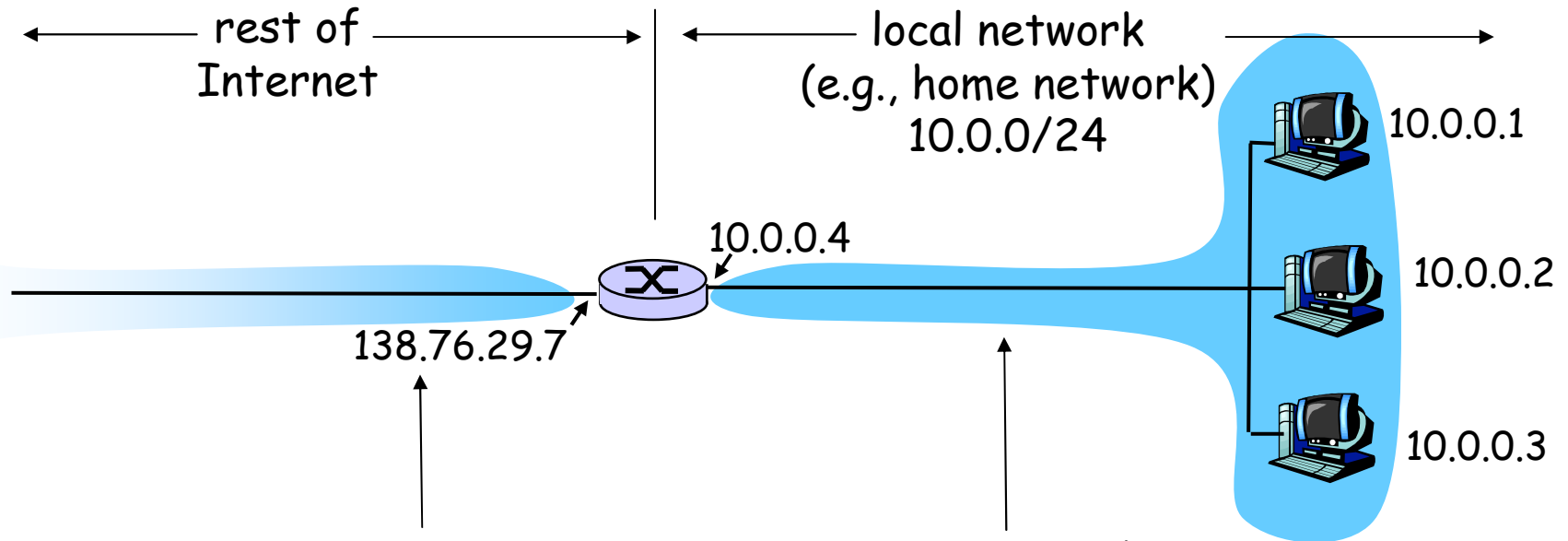
A: **ICANN**: Internet **C**orporation for **A**ssigned  
**N**ames and **N**umbers

m allocates addresses

m manages DNS

m assigns domain names, resolves disputes

# NAT: Network Address Translation



*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

- r **Motivation:** local network uses just one IP address as far as outside world is concerned:
  - m range of addresses not needed from ISP: just one IP address for all devices
  - m can change addresses of devices in local network without notifying outside world
  - m can change ISP without changing addresses of devices in local network
  - m devices inside local net not explicitly addressable, visible by outside world (a security plus).

# NAT: Network Address Translation

**Implementation:** NAT router must:

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

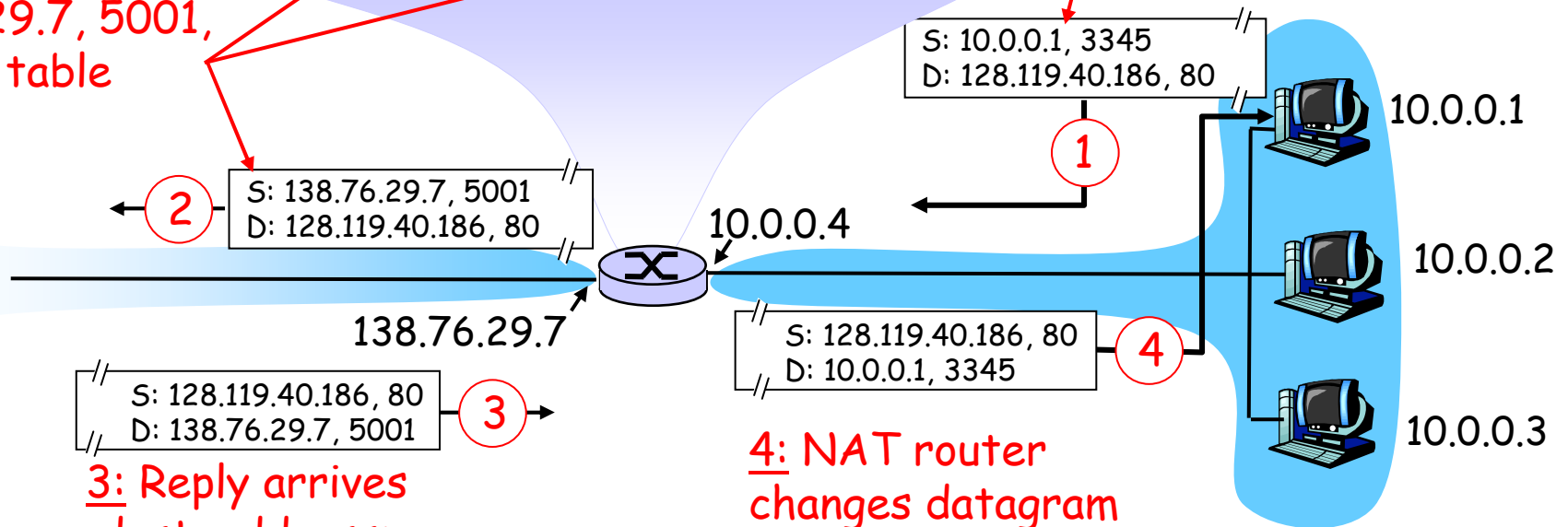


# NAT: Network Address Translation

**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



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

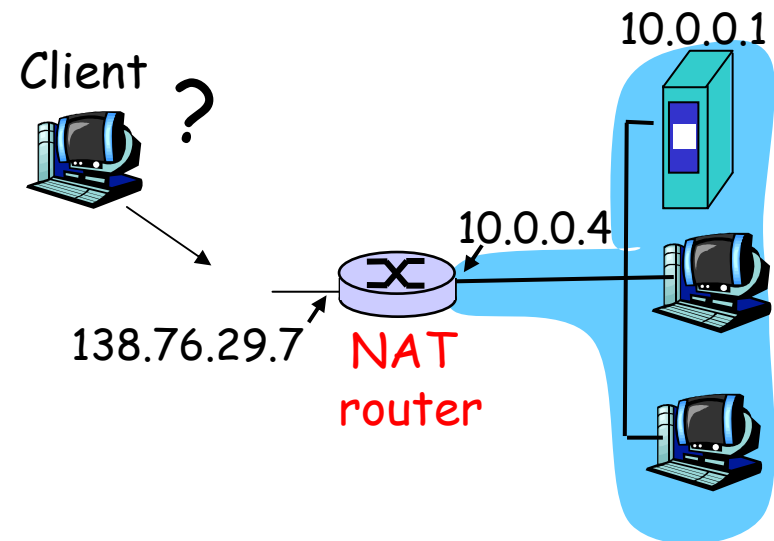
**4:** NAT router changes datagram dest addr from 138.76.29.7, 5001 to 10.0.0.1, 3345

# NAT: Network Address Translation

- r 16-bit port-number field:
  - m 60,000 simultaneous connections with a single LAN-side address!
- r NAT is controversial:
  - m routers should only process up to layer 3
  - m violates end-to-end argument
    - NAT possibility must be taken into account by app designers, eg, P2P applications
  - m address shortage should instead be solved by IPv6

# NAT traversal problem

- r client wants to connect to server with address 10.0.0.1
  - m server address 10.0.0.1 local to LAN (client can't use it as destination addr)
  - m only one externally visible NATted address: 138.76.29.7
- r solution 1: statically configure NAT to forward incoming connection requests at given port to server
  - m e.g., (138.76.29.7, port 2500) always forwarded to 10.0.0.1 port 25000

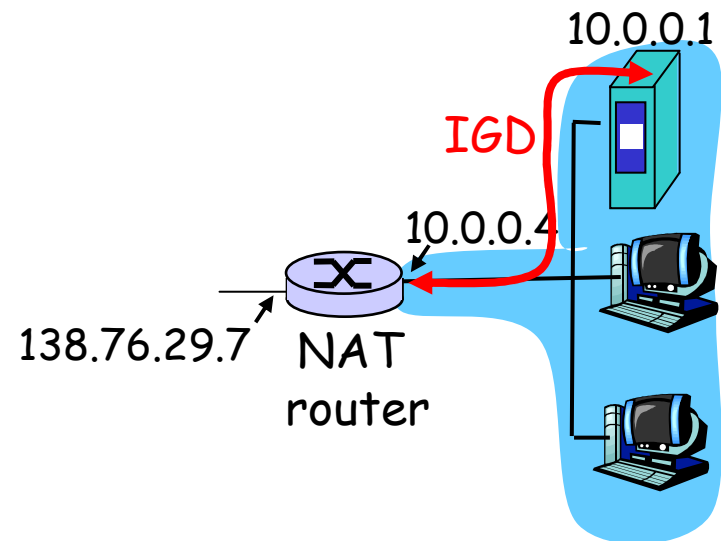


# NAT traversal problem

r solution 2: Universal Plug and Play (UPnP) Internet Gateway Device (IGD) Protocol. Allows NATted host to:

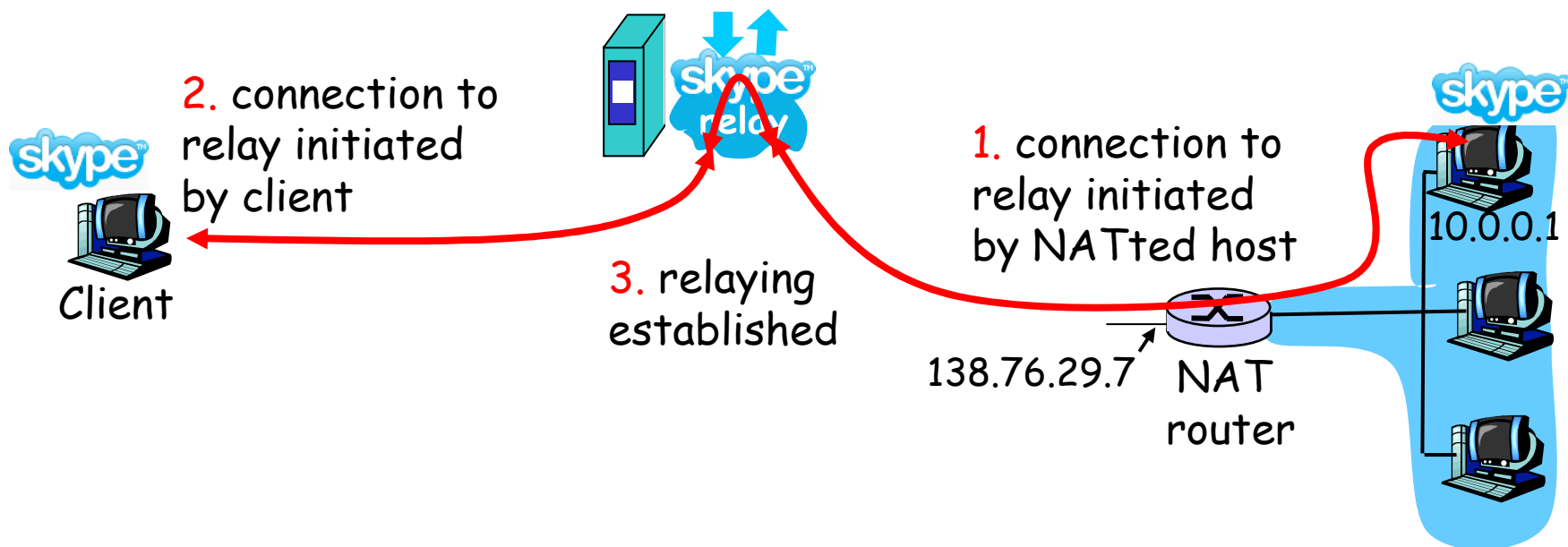
- ❖ learn public IP address (138.76.29.7)
- ❖ add/remove port mappings (with lease times)

i.e., automate static NAT port map configuration



# NAT traversal problem

- r solution 3: relaying (used in Skype)
  - m NATed client establishes connection to relay
  - m External client connects to relay
  - m relay bridges packets between to connections



# Chapter 4: Network Layer

- r 4.1 Introduction
- r 4.2 Virtual circuit and datagram networks
- r 4.3 What's inside a router
- r 4.4 IP: Internet Protocol
  - m Datagram format
  - m IPv4 addressing
  - m ICMP
  - m IPv6
- r 4.5 Routing algorithms
  - m Link state
  - m Distance Vector
  - m Hierarchical routing
- r 4.6 Routing in the Internet
  - m RIP
  - m OSPF
  - m BGP
- r 4.7 Broadcast and multicast routing

# ICMP: Internet Control Message Protocol

r used by hosts & routers to communicate network-level information

m error reporting:  
unreachable host, network,  
port, protocol

m echo request/reply (used  
by ping)

r network-layer "above" IP:

m ICMP msgs carried in IP  
datagrams

r **ICMP message:** type, code plus  
first 8 bytes of IP datagram  
causing error

<u>Type</u>	<u>Code</u>	<u>description</u>
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

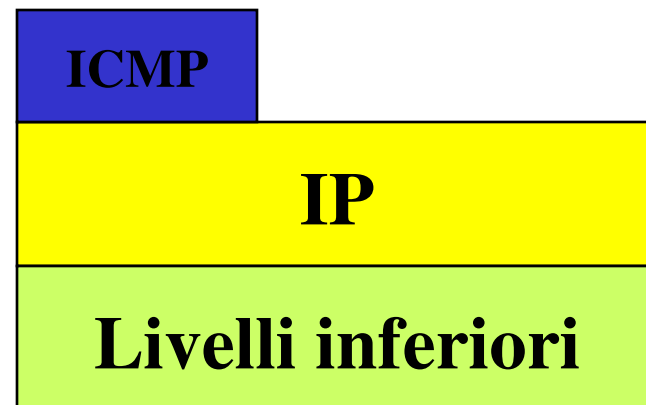
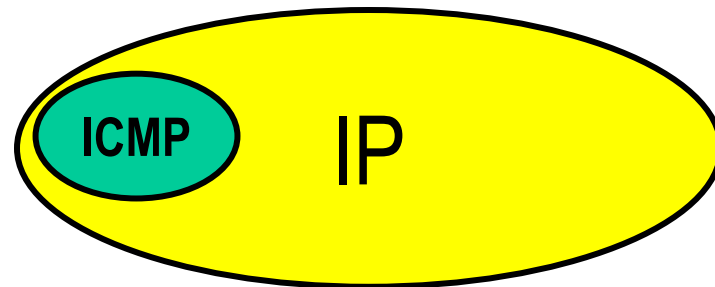
# Traceroute and ICMP

- r Source sends series of UDP segments to dest
    - m First has TTL =1
    - m Second has TTL=2, etc.
    - m Unlikely port number
  - r When nth datagram arrives to nth router:
    - m Router discards datagram
    - m And sends to source an ICMP message (type 11, code 0)
    - m Message includes name of router & IP address
  - r When ICMP message arrives, source calculates RTT
  - r Traceroute does this 3 times
- Stopping criterion
- r UDP segment eventually arrives at destination host
  - r Destination returns ICMP "host unreachable" packet (type 3, code 3)
  - r When source gets this ICMP, stops.

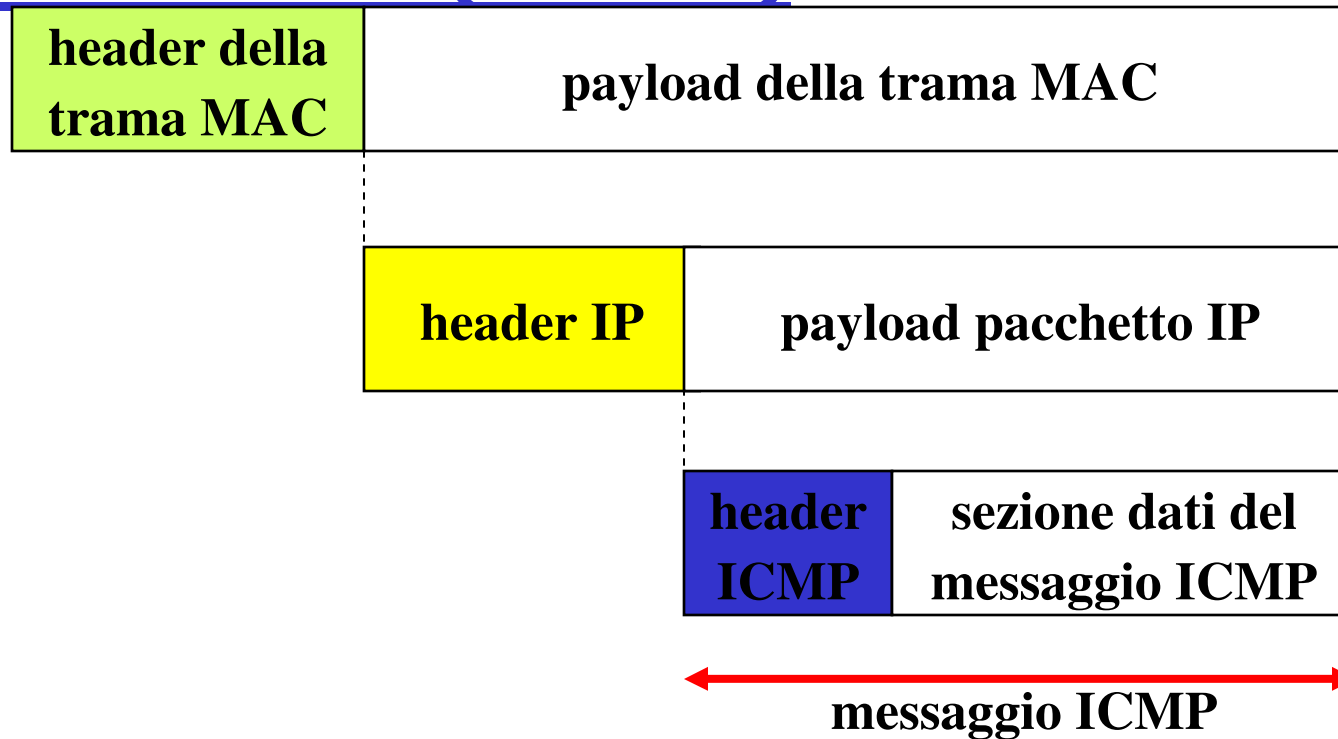


# Internet Control Message Protocol (ICMP)

- r E' un protocollo per messaggi di servizio fra host e router per informazioni su errori e fasi di attraversamento della rete
- r da questo punto di vista può essere considerato come parte di IP
- r i messaggi ICMP sono incapsulati e trasportati da IP, e quindi da questo punto di vista può essere considerato un utente di IP



# Internet Control Message Protocol (ICMP)



- r Nel pacchetto IP il campo protocol indica il codice dell'ICMP
- r il messaggio ICMP viaggia all'interno del pacchetto IP

# Messaggi ICMP

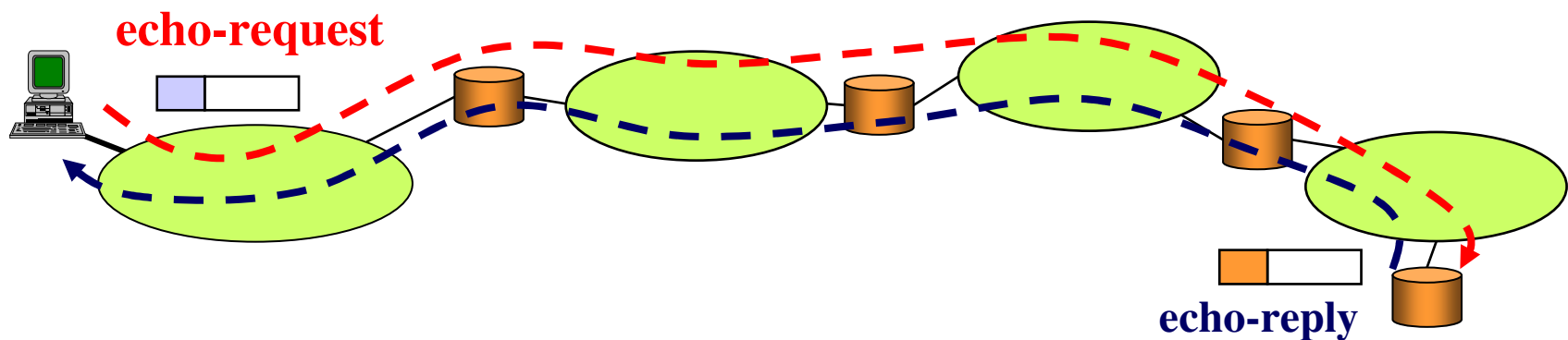
Da informazioni sul codice  
Dell'errore (info piu' dettagliate  
Sull'errore verificatosi)



<b>Type</b>		<b>Type</b>	
<b>0</b>	<b>Echo reply</b>	<b>11</b>	<b>Parameter problem</b>
<b>3</b>	<b>Destination unreachable</b>	<b>13</b>	<b>Timestamp request</b>
<b>4</b>	<b>Source Quench</b>	<b>14</b>	<b>Timestamp reply</b>
<b>5</b>	<b>Redirect (change a route)</b>	<b>17</b>	<b>Address mask request</b>
<b>8</b>	<b>Echo request</b>	<b>18</b>	<b>Address mask reply</b>
<b>11</b>	<b>Time exceeded</b>		

# Echo

- r I messaggi di Echo-request e Echo-reply sono usati per verificare la raggiungibilità e lo stato di un host o un router
- r quando un nodo IP riceve un messaggio di Echo-request risponda immediatamente con un messaggio di Echo reply



# Echo

<b>type</b> (8 request, 0 reply)	<b>code</b> (0)	<b>checksum</b>
<b>identifier</b>		<b>sequence number</b>
<b>optional data</b>		

- r Il campo identifier viene scelto dal mittente della richiesta
- r nella risposta viene ripetuto lo stesso identifier della richiesta
- r più richieste consecutive possono avere lo stesso identifier e differire per il sequence number
- r una sequenza arbitraria può essere aggiunta dal mittente nel campo optional data e deve essere riportata uguale nella risposta

# Uso dei messaggi di echo: PING

```
Prompt di MS-DOS
Auto
C:\>ping 131.175.123.96
Esecuzione di Ping 131.175.123.96 con 32 byte di dati:
Risposta da 131.175.123.96: byte=32 durata<10ms TTL=128
Risposta da 131.175.123.96: byte=32 durata<10ms TTL=128
Risposta da 131.175.123.96: byte=32 durata<10ms TTL=128
Risposta da 131.175.123.96: byte=32 durata<10ms TTL=128
Statistiche Ping per 131.175.123.96:
  Pacchetti: Trasmessi = 4, Ricevuti = 4, Persi = 0 (0% persi),
Tempo approssimativo percorsi andata/ritorno in millisecondi:
  Minimo = 0ms, Massimo = 0ms, Medio = 0ms
C:\>
C:\>
C:\>
C:\>
C:\>
C:\>
C:\>
C:\>
C:\>
C:\>
```

# Destination unreachable

type (3)	code (0-12)	checksum
non usato (0)		
header + primi 64 bit del pacchetto IP che ha causato il problema		

- r Quando un router scarta un pacchetto per qualche motivo normalmente genera un messaggio di errore che invia alla sorgente del pacchetto
- r nel campo code è codificato il motivo che ha causato l'errore
- r ovviamente la generazione del messaggio avviene solo nei casi in cui il router può accorgersi del problema
- r il motivo più comune è il fatto che la destinazione non è presente nella tabella di routing. Un messaggio destination unreachable può anche essere inviato se è necessario frammentare un pacchetto ed il flag DF è a 1. Network Layer 4-39

# Destination unreachable

type (3)	code (0-12)	checksum
non usato (0)		
header + primi 64 bit del pacchetto IP che ha causato il problema		

Alcuni Code:

- 0 network unreachable
- 1 host unreachable
- 2 protocol unreachable
- 3 port unreachable
- 4 fragmentation needed and DF set
- 5 source route failed
- ...



# Time exceeded

<b>type</b> (11)	<b>code</b> (0-1)	<b>checksum</b>
<b>non usato</b> (0)		
<b>header + primi 64 bit del pacchetto IP che ha causato il problema</b>		

## r Code 0

- m Il messaggio di *time exceeded* viene usato quando il router decrementando il TTL lo pone a 0
- m il messaggio di *time exceeded* viene inviato alla sorgente del pacchetto

## r Code 1

- m viene usato dalla destinazione quando non tutti i frammenti di un pacchetto arrivano entro un tempo massimo

# Parameter problem

<b>type</b> (12)	<b>code</b> (0-1)	<b>checksum</b>
<b>pointer</b>	<b>non usato</b> (0)	
<b>header + primi 64 bit del pacchetto IP che ha causato il problema</b>		

Punta al  
Byte del  
datagramma  
IP originale  
dove si e'  
verificato  
il problema

r

## Code 0

m se l'header di un pacchetto IP ha una incongruenza in qualcuno dei suoi campi viene inviato il messaggio di parameter problem; il campo pointer punta al byte del pacchetto che ha causato il problema

## Code 1

m viene usato quando un'opzione non è implementata e non può essere soddisfatta

# Timestamp request e reply

<b>type</b> (13 request, 14 reply)	<b>code</b> (0)	<b>checksum</b>
<b>identifier</b>		<b>sequence number</b>
<b>originate timestamp</b>		
<b>receive timestamp</b>		
<b>transmit timestamp</b>		

- r Questo messaggio viene usato per scambiarsi informazioni sul clock di sorgente e destinazione
- r *originate timestamp*: viene riempito dalla sorgente
- r *receive timestamp*: viene riempito dalla destinazione appena ricevuto il pacchetto
- r *transmit timestamp*: viene riempito dalla destinazione immediatamente prima di inviare il pacchetto di risposta

# Address mask request e reply

<b>type</b> (17 request, 18 reply)	<b>code</b> (0)	<b>checksum</b>
<b>identifier</b>		<b>sequence number</b>
<b>address mask</b>		

- r Questo messaggio viene usato per conoscere la netmask di un host/router
- r Il campo address mask viene riempito da invia la risposta

# Chapter 4: Network Layer

- r 4.1 Introduction
- r 4.2 Virtual circuit and datagram networks
- r 4.3 What's inside a router
- r 4.4 IP: Internet Protocol
  - m Datagram format
  - m IPv4 addressing
  - m ICMP
  - m IPv6
- r 4.5 Routing algorithms
  - m Link state
  - m Distance Vector
  - m Hierarchical routing
- r 4.6 Routing in the Internet
  - m RIP
  - m OSPF
  - m BGP
- r 4.7 Broadcast and multicast routing

# IPv6

- r **Initial motivation:** 32-bit address space soon to be completely allocated.
  - r **Additional motivation:**
    - m header format helps speed processing/forwarding
    - m header changes to facilitate QoS
- IPv6 datagram format:**
- m fixed-length 40 byte header
  - m no fragmentation allowed

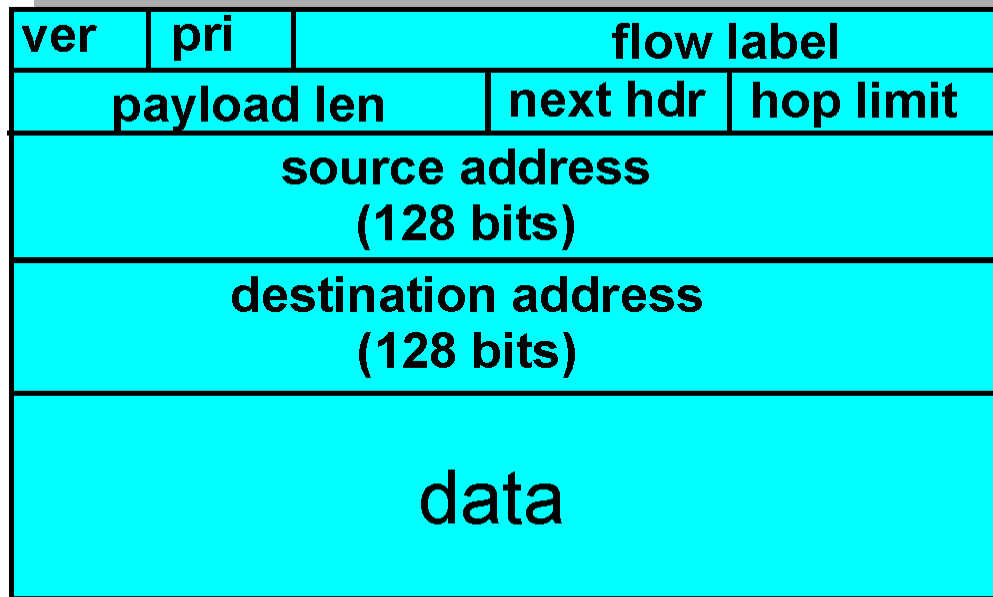
# IPv6 Header (Cont)

*Priority:* identify priority among datagrams in flow

*Flow Label:* identify datagrams in same "flow."

(concept of "flow" not well defined).

*Next header:* identify upper layer protocol for data



← 32 bits →

# Other Changes from IPv4

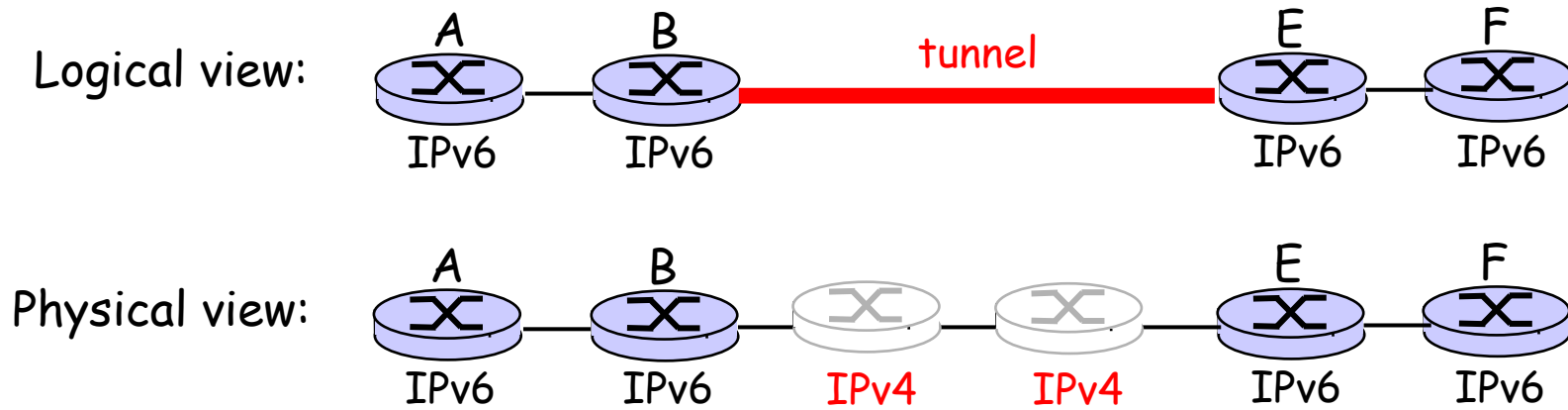
- r *Checksum*: removed entirely to reduce processing time at each hop
- r *Options*: allowed, but outside of header, indicated by "Next Header" field
- r *ICMPv6*: new version of ICMP
  - m additional message types, e.g. "Packet Too Big"
  - m multicast group management functions



# Transition From IPv4 To IPv6

- r Not all routers can be upgraded simultaneous
  - m no "flag days"
  - m How will the network operate with mixed IPv4 and IPv6 routers?
- r *Tunneling*: IPv6 carried as payload in IPv4 datagram among IPv4 routers

# Tunneling



# Tunneling

