

# Network Layer: Data Plane

EECS3214

18-02-25

© All material copyright 1996-2016  
J.F. Kurose and K.W. Ross, All Rights Reserved

4-1

## Chapter 4: outline

### 4.1 Overview of Network layer

- data plane
- control plane

### 4.2 What's inside a router

### 4.3 IP: Internet Protocol

- datagram format
- fragmentation
- IPv4 addressing
- network address translation
- IPv6

### 4.4 Generalized Forward and SDN

- match
- action
- OpenFlow examples of match-plus-action in action

Network Layer: Data Plane 4-2

## Chapter 4: network layer

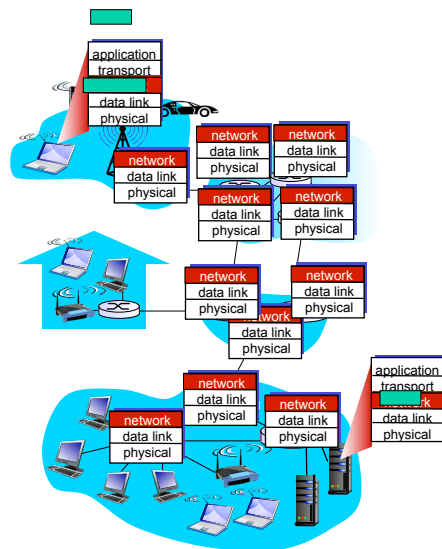
### *chapter goals:*

- understand principles behind network layer services, focusing on data plane:
  - network layer service models
  - forwarding versus routing
  - how a router works
  - generalized forwarding
- instantiation, implementation in the Internet

Network Layer: Data Plane 4-3

## Network layer

- transport segment from sending to receiving host
- on sending side encapsulates segments into datagrams
- on receiving side, delivers segments to transport layer
- network layer protocols in *every* host, router
- router examines header fields in all IP datagrams passing through it



Network Layer: Data Plane 4-4

## Two key network-layer functions

### *network-layer functions:*

- *forwarding*: move packets from router's input to appropriate router output
  - implemented in hardware
- *routing*: determine route taken by packets from source to destination
  - *routing algorithms*
  - implemented in software

### *analogy: taking a trip*

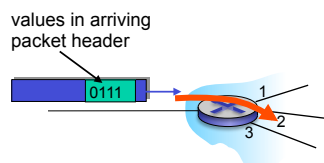
- *forwarding*: process of getting through single interchange
- *routing*: process of planning trip from source to destination

Network Layer: Data Plane 4-5

## Network layer: data plane, control plane

### *Data plane*

- local, per-router function
- determines how datagram arriving on router input port is forwarded to router output port
- forwarding function



- forwarding table (routing table)

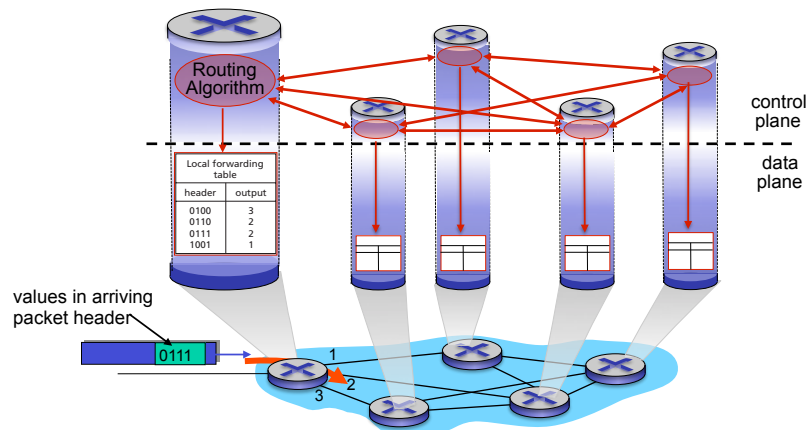
### *Control plane*

- network-wide logic
- determines how datagram is routed among routers along end-end path from source host to destination host
- two control-plane approaches:
  - *traditional routing algorithms*: implemented in routers
  - *software-defined networking (SDN)*: implemented in (remote) servers

Network Layer: Data Plane 4-6

## Per-router control plane

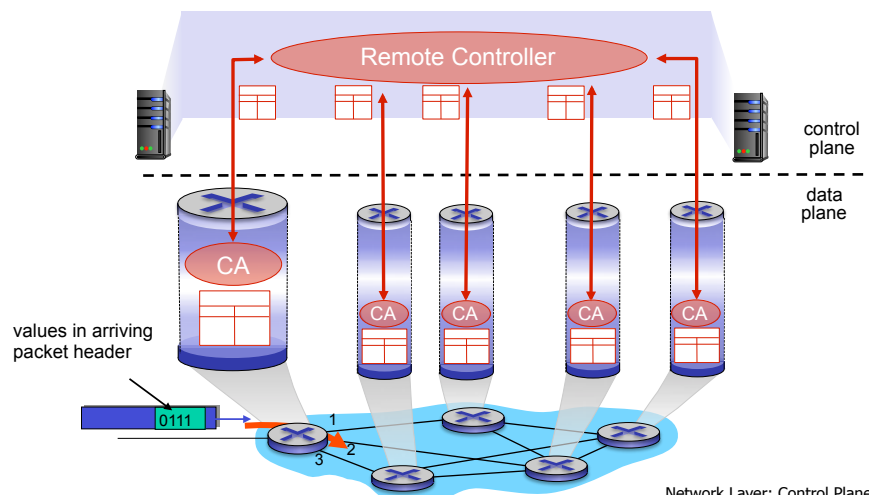
Individual routing algorithm components *in each and every router* interact in the control plane.



Network Layer: Control Plane 5-7

## Logically centralized control plane

A distinct (typically remote) controller interacts with local control agents. Computes and distributes the forwarding tables to local CAs.



Network Layer: Control Plane 5-8

## Network service model

**Q:** What *service model* for “channel” transporting datagrams from sender to receiver?

*example services for individual datagrams:*

- guaranteed delivery
- guaranteed delivery with less than 40 msec delay

*example services for a flow of datagrams:*

- in-order datagram delivery
- guaranteed minimum bandwidth to flow
- restrictions on changes in inter-packet spacing

Network Layer: Data Plane 4-9

## Network layer service models:

Network Architecture	Service Model	Guarantees ?				Congestion feedback
		Bandwidth	Loss	Order	Timing	
Internet	best effort	none	no	no	no	no (inferred via loss)
ATM	CBR	constant rate	yes	yes	yes	no congestion
ATM	VBR	guaranteed rate	yes	yes	yes	no congestion
ATM	ABR	guaranteed minimum	no	yes	no	yes
ATM	UBR	none	no	yes	no	no

Network Layer: Data Plane 4-10

## Chapter 4: outline

### 4.1 Overview of Network layer

- data plane
- control plane

### 4.2 What's inside a router

### 4.3 IP: Internet Protocol

- datagram format
- fragmentation
- IPv4 addressing
- network address translation
- IPv6

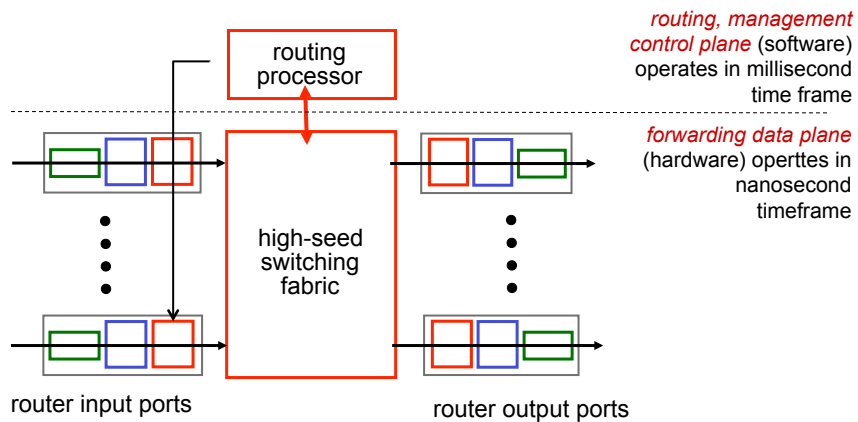
### 4.4 Generalized Forward and SDN

- match
- action
- OpenFlow examples of match-plus-action in action

Network Layer: Data Plane 4-11

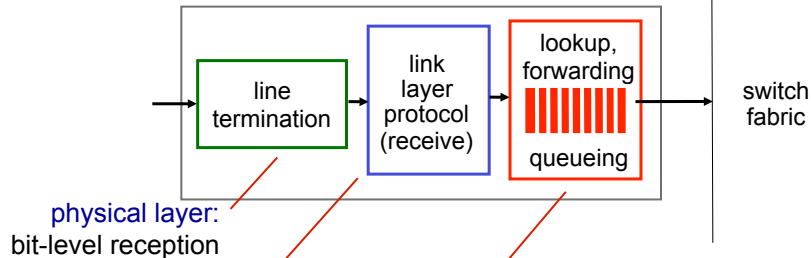
## Router architecture overview

- high-level view of generic router architecture:



Network Layer: Data Plane 4-12

## Input port functions



physical layer:  
bit-level reception

data link layer:  
e.g., Ethernet  
see chapter 5

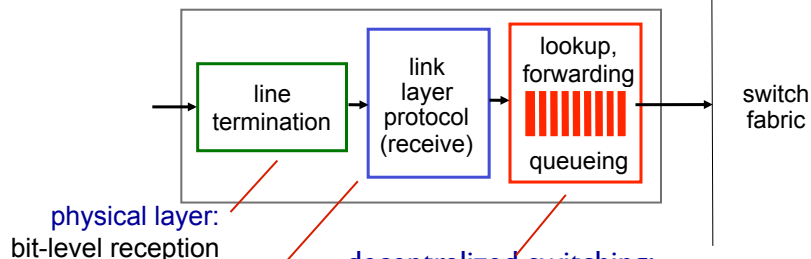
a few to approx. 1,000  
input ports in a router

### decentralized switching:

- using header field values, lookup output port using forwarding table in input port memory (“*match plus action*”)
- goal: complete input port processing at ‘line speed’
- queuing: if datagrams arrive faster than forwarding rate into switch fabric

Network Layer: Data Plane 4-13

## Input port functions



physical layer:  
bit-level reception

data link layer:  
e.g., Ethernet  
see chapter 5

### decentralized switching:

- using header field values, lookup output port using forwarding table in input port memory (“*match plus action*”)
- **destination-based forwarding:** forward based only on destination IP address (traditional)
- **generalized forwarding:** forward based on any set of header field values

Network Layer: Data Plane 4-14

## Destination-based forwarding: an example

<i>forwarding table</i>	
Destination Address Range	Link Interface
11001000 00010111 00010000 00000000 through 11001000 00010111 00010111 11111111	0
11001000 00010111 00011000 00000000 through 11001000 00010111 00011000 11111111	1
11001000 00010111 00011001 00000000 through 11001000 00010111 00011111 11111111	2
otherwise	3

Q: but what happens if ranges don't divide up so nicely?

Network Layer: Data Plane 4-15

## Longest prefix matching: example

### *longest prefix matching*

when looking for forwarding table entry for given destination address, use *longest* address prefix that matches destination address.

Destination Address Range	Link interface
11001000 00010111 00010*** *****	0
11001000 00010111 00011000 *****	1
11001000 00010111 00011*** *****	2
otherwise	3

examples:

DA: 11001000 00010111 00010**110** 10100001

which interface?

DA: 11001000 00010111 0001**1000** 10101010

which interface?

Network Layer: Data Plane 4-16



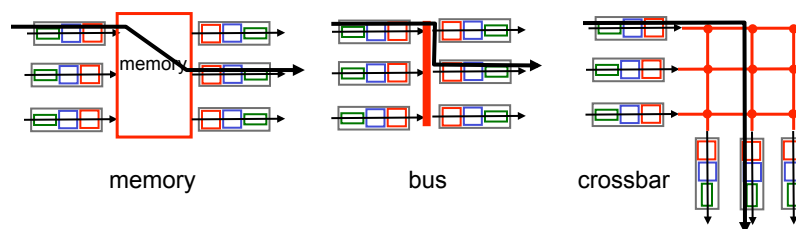
## Longest prefix matching

- we'll see *why* longest prefix matching is used shortly, when we study addressing
- longest prefix matching: often performed using *ternary content addressable memories (TCAMs)*
  - *content addressable*: present address to TCAM: retrieve address in one clock cycle, regardless of table size
  - Cisco Catalyst: can hold up ~1M routing table entries in TCAM

Network Layer: Data Plane 4-17

## Switching fabrics

- transfer packet from input buffer to appropriate output buffer
- switching rate: rate at which packets can be transfer from inputs to outputs
  - often measured as multiple of input/output line rate
  - N inputs: switching rate N times line rate desirable
- three types of switching fabrics

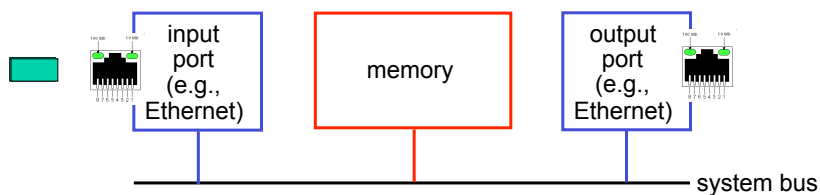


Network Layer: Data Plane 4-18

## Switching via memory

### *first generation routers:*

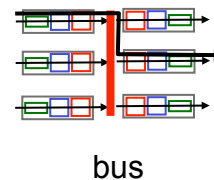
- traditional computers with switching under direct control of CPU
- packet copied to system's memory
- speed limited by memory bandwidth (2 bus crossings per datagram)



Network Layer: Data Plane 4-19

## Switching via a bus

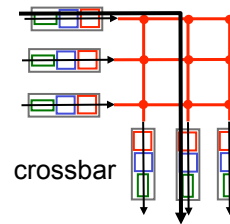
- datagram from input port memory to output port memory via a shared bus
- *bus contention*: switching speed limited by bus bandwidth
- 32 Gbps bus, Cisco 5600: sufficient speed for access and enterprise routers



Network Layer: Data Plane 4-20

## Switching via interconnection network

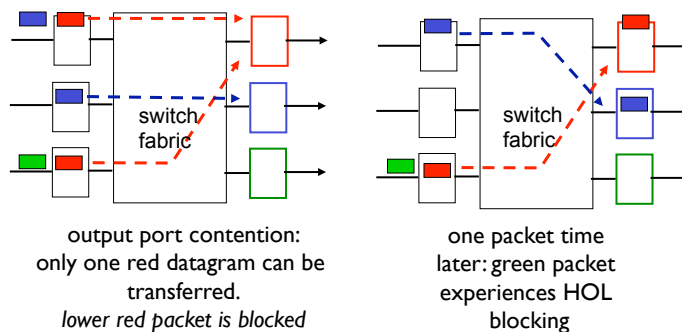
- overcome bus bandwidth limitations
- banyan networks, crossbar, other interconnection nets initially developed to connect processors in multiprocessor
- advanced design: fragmenting datagram into fixed length cells, switch cells through the fabric.
- Cisco 12000: switches 60 Gbps through the interconnection network



Network Layer: Data Plane 4-21

## Input port queuing

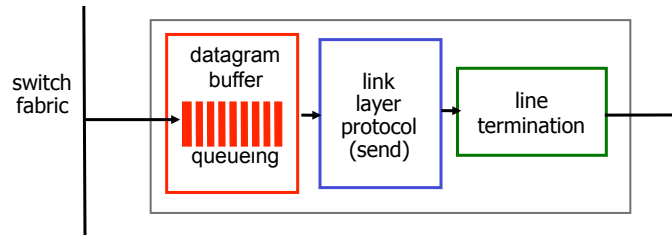
- fabric slower than input ports combined -> queueing may occur at input queues
  - *queueing delay and loss due to input buffer overflow!*
- **Head-of-the-Line (HOL) blocking:** queued datagram at front of queue prevents others in queue from moving forward



Network Layer: Data Plane 4-22

## Output ports

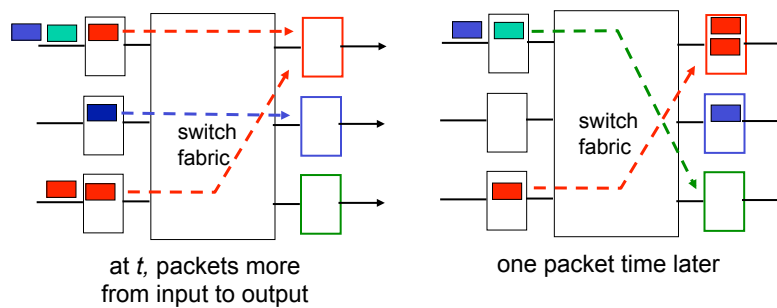
*This slide is HUGE important!*



- **buffering** required from fabric faster rate  
Datagram (packets) can be lost due to congestion, lack of buffers
- **scheduling** datagrams  
Priority scheduling – who gets best performance, network neutrality

Network Layer: Data Plane 4-23

## Output port queueing



- buffering when arrival rate via switch exceeds output line speed
- **queueing (delay) and loss due to output port buffer overflow!**

Network Layer: Data Plane 4-24

## How much buffering?

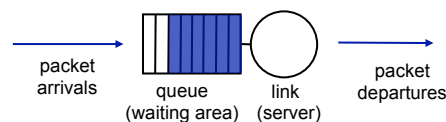
- RFC 3439 rule of thumb: average buffering equal to “typical” RTT (say 250 msec) times link capacity  $C$ 
  - e.g.,  $C = 10$  Gbps link: needs 2.5 Gbit buffer
  - assume a small number of TCP flows
- recent recommendation, assuming a large number of TCP flows: with  $N$  flows, buffering equal to

$$\frac{RTT \cdot C}{\sqrt{N}}$$

Network Layer: Data Plane 4-25

## Scheduling mechanisms

- *scheduling*: choose next packet to send on link
- *FIFO (first in first out) scheduling*: send in order of arrival to queue
  - real-world example?
  - *discard policy*: if packet arrives to full queue: who to discard?
    - *tail drop*: drop arriving packet
    - *priority*: drop/remove on priority basis
    - *random*: drop/remove randomly

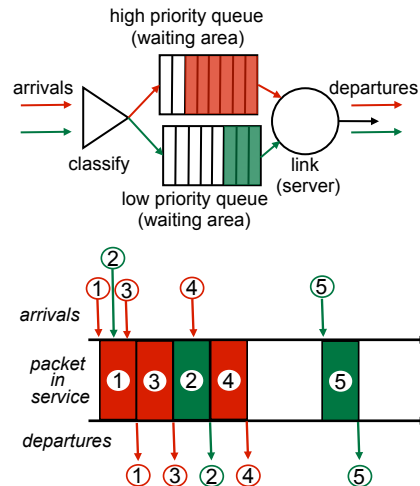


Network Layer: Data Plane 4-26

## Scheduling policies: priority

**priority scheduling:** send highest priority queued packet

- multiple classes, with different priorities
  - class may depend on marking or other header info, e.g. IP source/dest, port numbers, etc.
  - example: VoIP vs. SMTP or IMAP

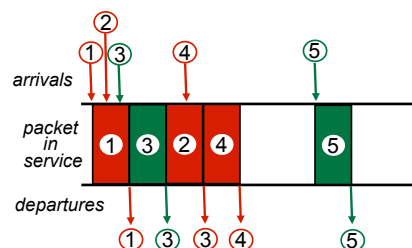


Network Layer: Data Plane 4-27

## Scheduling policies: still more

**Round Robin (RR) scheduling:**

- multiple classes
- cyclically scan class queues, sending one complete packet from each class (if available)

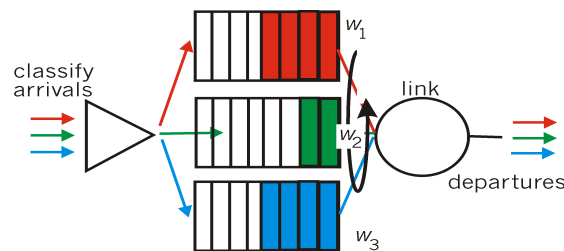


Network Layer: Data Plane 4-28

## Scheduling policies: still more

### *Weighted Fair Queuing (WFQ):*

- generalized Round Robin
- each class gets weighted amount of service in each cycle:  $w_i / (w_1 + w_2 + \dots + w_K)$



Network Layer: Data Plane 4-29

## Chapter 4: outline

### 4.1 Overview of Network layer

- data plane
- control plane

### 4.2 What's inside a router

### 4.3 IP: Internet Protocol

- datagram format
- fragmentation
- IPv4 addressing
- network address translation
- IPv6

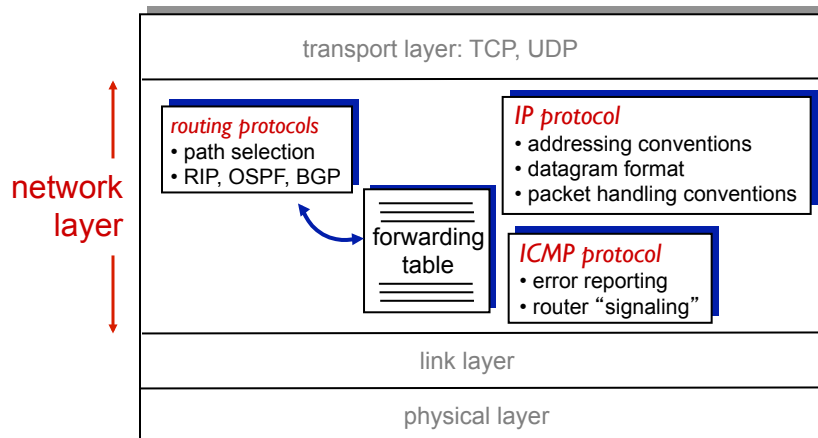
### 4.4 Generalized Forward and SDN

- match
- action
- OpenFlow examples of match-plus-action in action

Network Layer: Data Plane 4-30

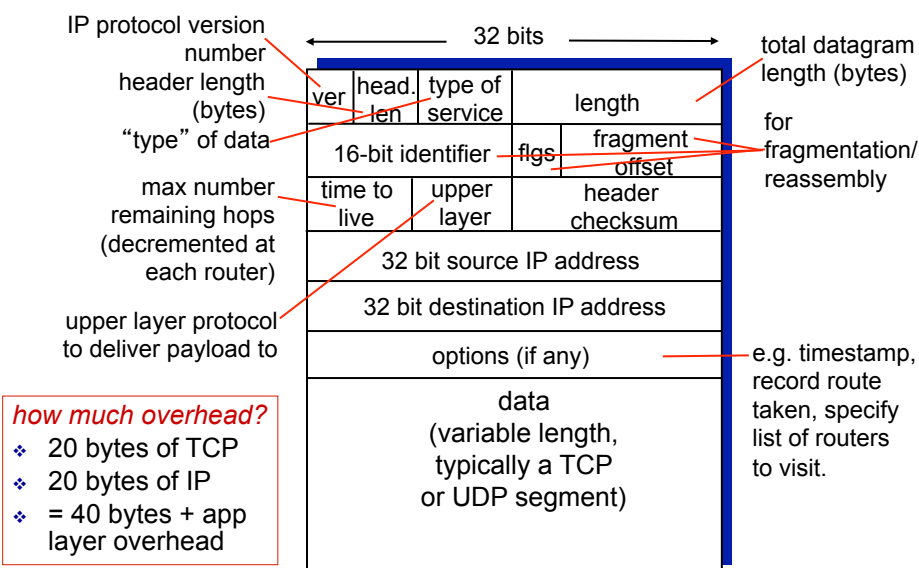
## The Internet network layer

host, router network layer functions:



Network Layer: Data Plane 4-31

## IP datagram format

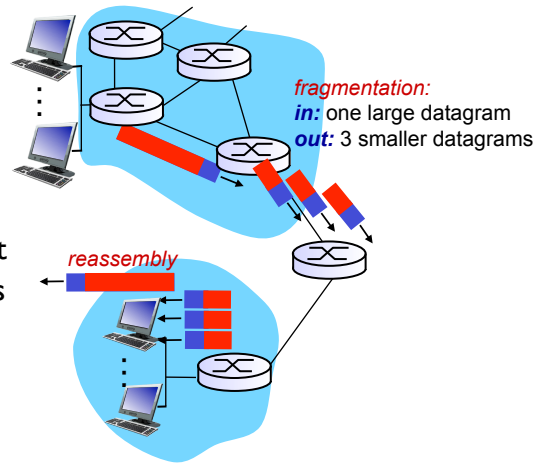


Network Layer: Data Plane 4-32



## 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



Network Layer: Data Plane 4-33

## IP fragmentation, reassembly

### example:

- ❖ 4000 byte datagram
- ❖ MTU = 1500 bytes

1480 bytes in  
data field

offset =  
 $1480/8$

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

one large datagram becomes  
several smaller datagrams

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

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

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

Network Layer: Data Plane 4-34

## Chapter 4: outline

### 4.1 Overview of Network layer

- data plane
- control plane

### 4.2 What's inside a router

### 4.3 IP: Internet Protocol

- datagram format
- fragmentation
- **IPv4 addressing**
- network address translation
- IPv6

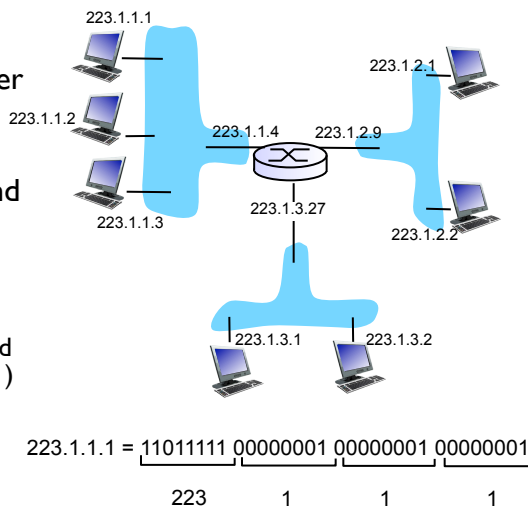
### 4.4 Generalized Forward and SDN

- match
- action
- OpenFlow examples of match-plus-action in action

Network Layer: Data Plane 4-35

## 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**



Network Layer: Data Plane 4-36

## IP addressing: introduction (2)

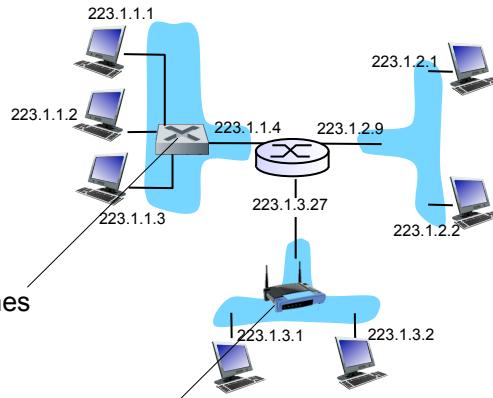
**Q:** how are interfaces actually connected?

**A:** we'll learn about that in chapter 5, 6.

**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)

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



Network Layer: Data Plane 4-37

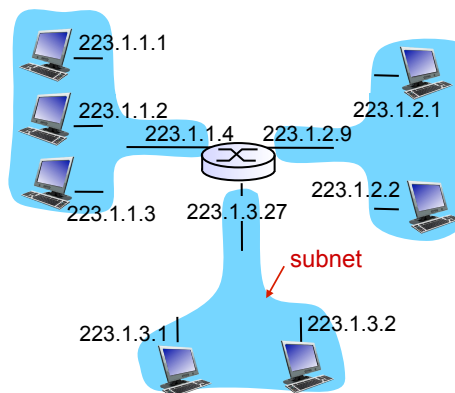
## Subnets

### ■ 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*



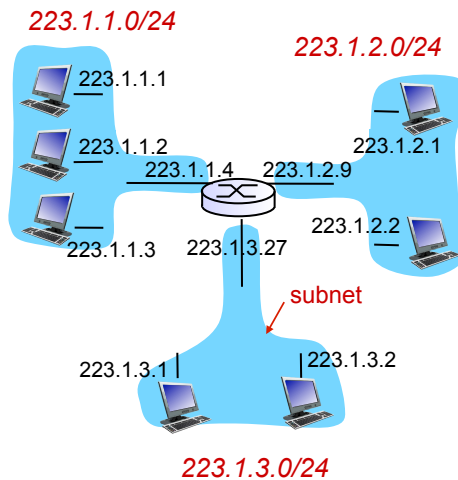
network consisting of 3 subnets

Network Layer: Data Plane 4-38

## Subnets (2)

### *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*

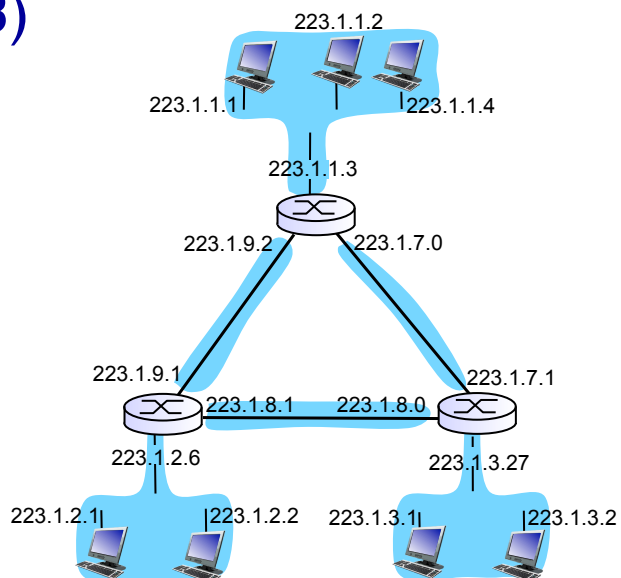


subnet mask: /24

Network Layer: Data Plane 4-39

## Subnets (3)

how many?

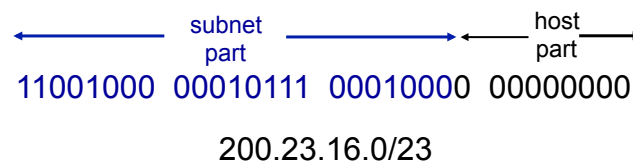


Network Layer: Data Plane 4-40

## IP addressing: CIDR

### CIDR: Classless InterDomain Routing

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



Network Layer: Data Plane 4-41

## IP addressing: how to get one?

**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

Network Layer: Data Plane 4-42

## IP addresses: how to get one? (2)

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

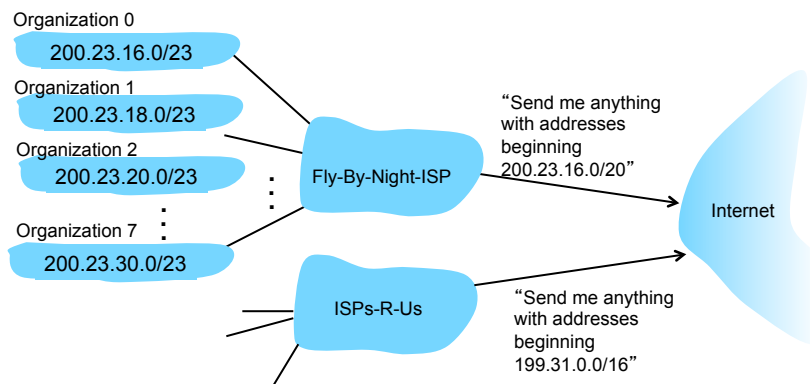
**A:** gets allocated portion of its provider ISP's address space

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

Network Layer: Data Plane 4-43

## Hierarchical addressing: route aggregation

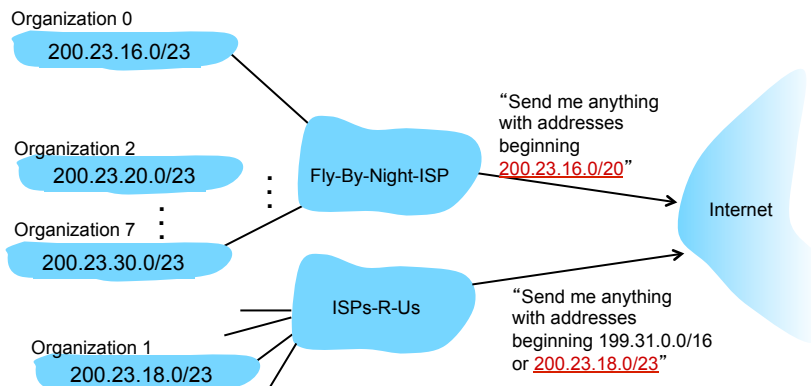
hierarchical addressing allows efficient advertisement of routing information:



Network Layer: Data Plane 4-44

## Hierarchical addressing: more specific routes

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



Network Layer: Data Plane 4-45

## IP addresses: how to get one? (3)

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

- hard-coded by system admin in a file
  - 
  -
- **DHCP:** **D**ynamic **H**ost **C**onfiguration **P**rotocol:  
dynamically get address from as server
  - "plug-and-play"

Network Layer: Data Plane 4-46

## DHCP: Dynamic Host Configuration Protocol

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

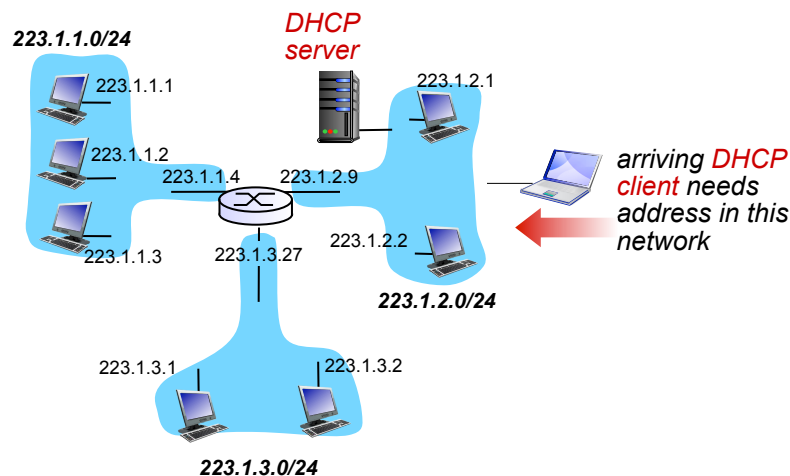
- can renew its lease on address in use
- allows reuse of addresses (only hold address while connected/“on”)
- support for mobile users who want to join network (more shortly)

### **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

Network Layer: Data Plane 4-47

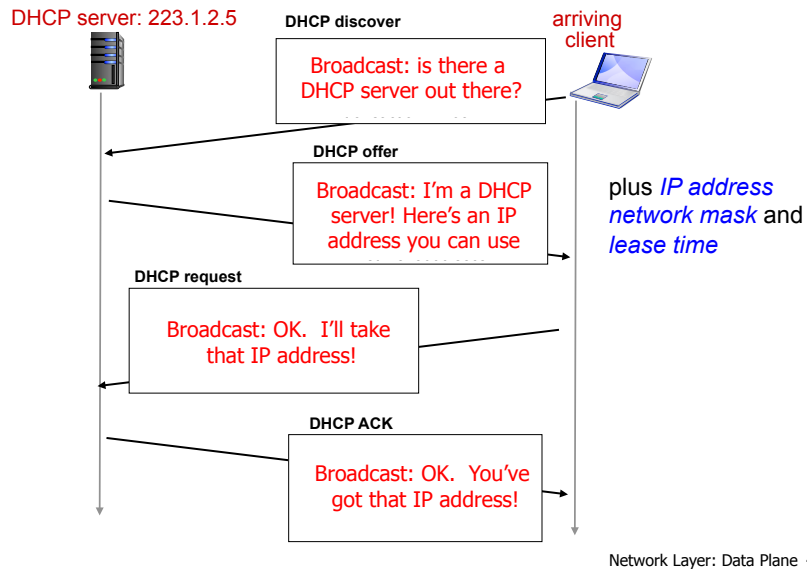
## DHCP client-server scenario



Network Layer: Data Plane 4-48



## DHCP client-server scenario



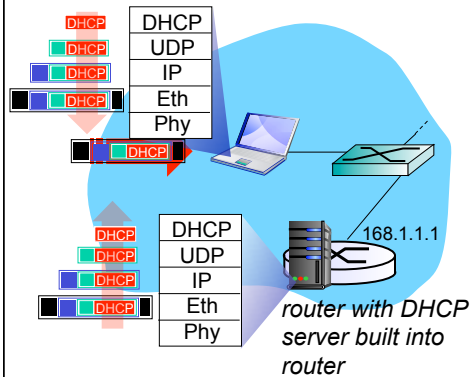
## DHCP: more than IP addresses

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

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

Network Layer: Data Plane 4-50

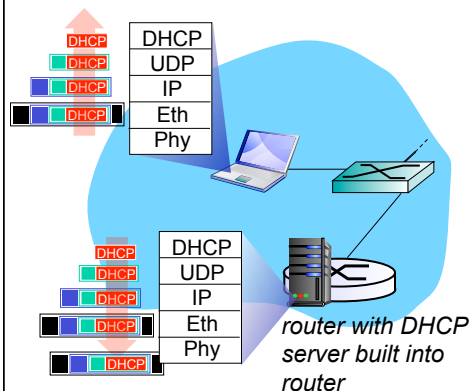
## DHCP: example



- 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: FFFFFFFF) on LAN, received at router running DHCP server
- Ethernet demuxed to IP demuxed, UDP demuxed to DHCP

Network Layer: Data Plane 4-51

## DHCP: example



- DHCP 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 DNS server, IP address of its first-hop router

Network Layer: Data Plane 4-52

## 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  
 .....

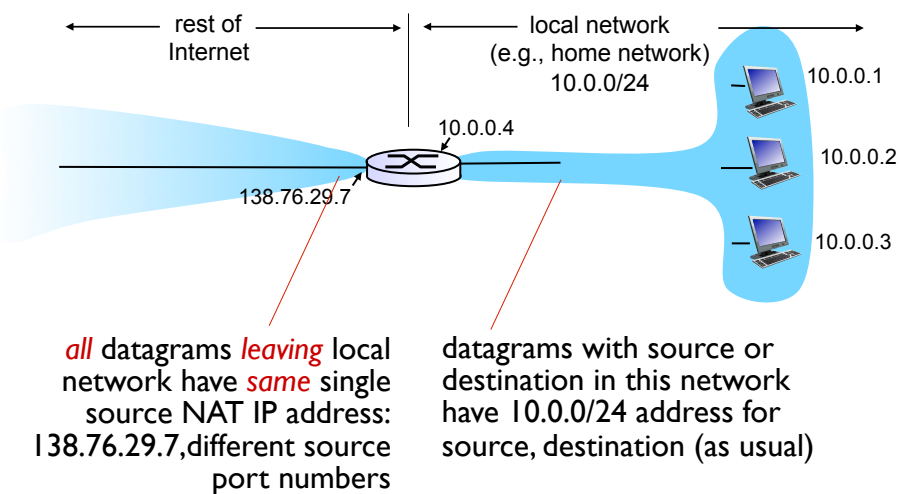
request

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

Network Layer: Data Plane 4-53

## NAT: network address translation



Network Layer: Data Plane 4-54

## NAT: motivation

*motivation:* 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
- devices inside local net not explicitly addressable, visible by outside world (a security plus)

Network Layer: Data Plane 4-55

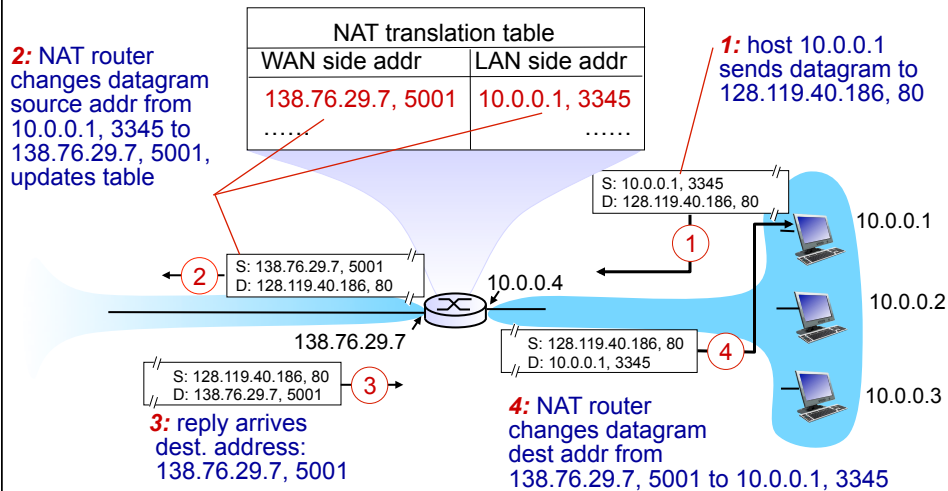
## NAT: implementation

*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

Network Layer: Data Plane 4-56

## NAT: example



\* Check out the online interactive exercises for more examples: [http://gaia.cs.umass.edu/kurose\\_ross/interactive/](http://gaia.cs.umass.edu/kurose_ross/interactive/)

Network Layer: Data Plane 4-57

## NAT: other issues

- 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
  - address shortage should be solved by IPv6
  - violates end-to-end argument
    - NAT possibility must be taken into account by app designers, e.g., P2P applications
  - NAT traversal: what if client wants to connect to server behind NAT?
    - well-known port numbers
    - peers in P2P acting as servers

Network Layer: Data Plane 4-58

## Chapter 4: outline

### 4.1 Overview of Network layer

- data plane
- control plane

### 4.2 What's inside a router

### 4.3 IP: Internet Protocol

- datagram format
- fragmentation
- IPv4 addressing
- network address translation
- **IPv6**

### 4.4 Generalized Forward and SDN

- match
- action
- OpenFlow examples of match-plus-action in action

Network Layer: Data Plane 4-59

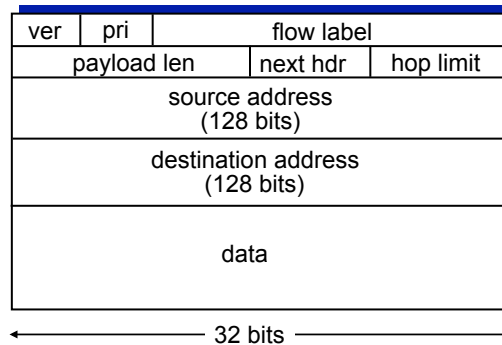
## IPv6: motivation

- *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
    - ICMP message "Packet Too Big"

Network Layer: Data Plane 4-60

## IPv6 datagram format

- 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



Network Layer: Data Plane 4-61

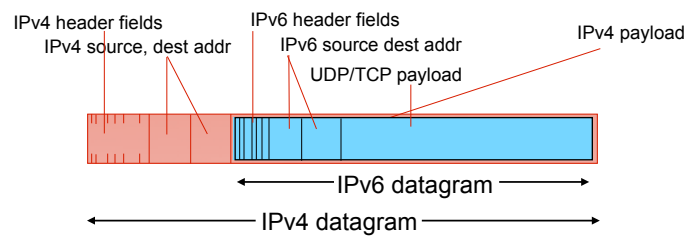
## Other changes from IPv4

- *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

Network Layer: Data Plane 4-62

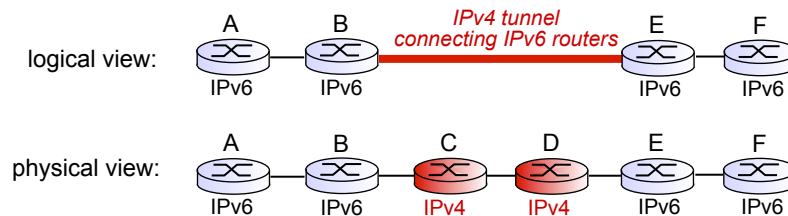
## Transition from IPv4 to IPv6

- 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



Network Layer: Data Plane 4-63

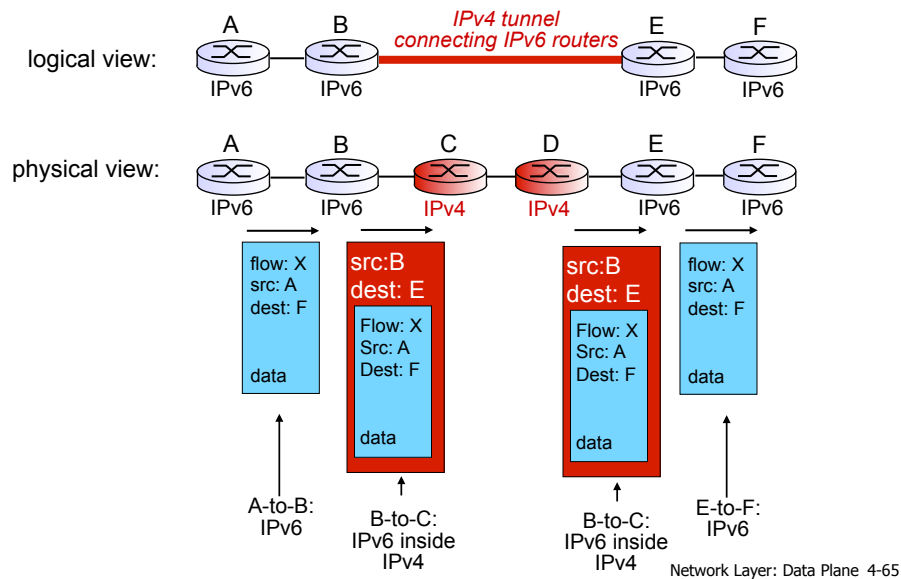
## Tunneling



Network Layer: Data Plane 4-64



# Tunneling

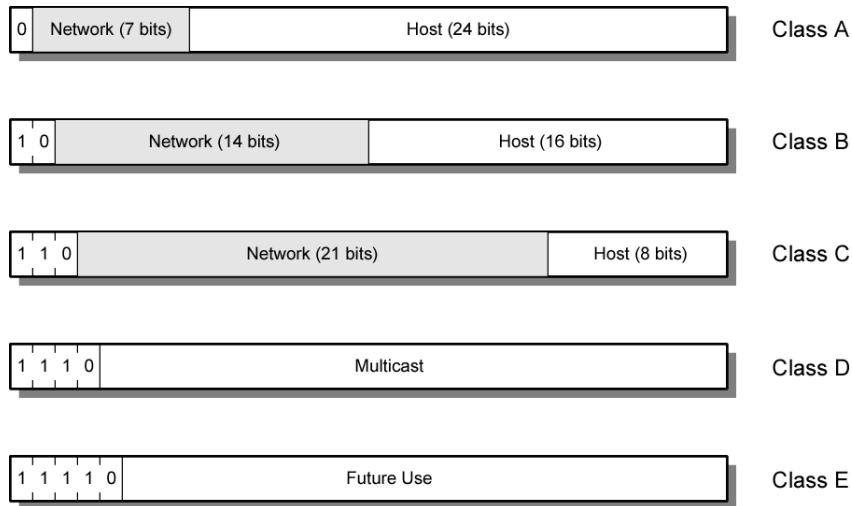


## IPv6: adoption

- Google: 8% of clients access services via IPv6
- NIST: 1/3 of all US government domains are IPv6 capable
- *Long (long!) time for deployment, use*
  - 20 years and counting!
  - think of application-level changes in last 20 years: WWW, Facebook, streaming media, Skype, ...
  - *Why?*

Network Layer: Data Plane 4-66

## Classful IPv4 Address Formats



67

## Chapter 4: done!

### 4.1 Overview of Network layer: data plane and control plane

### 4.2 What's inside a router

### 4.3 IP: Internet Protocol

- datagram format
- fragmentation
- IPv4 addressing
- NAT
- IPv6

### 4.4 Generalized Forward and SDN

- match plus action
- OpenFlow example

**Question:** how do forwarding tables (destination-based forwarding) computed?

**Answer:** by the control plane (next chapter)

Network Layer: Data Plane 4-68