

# Chapter 4

## Network Layer: The Data Plane

---

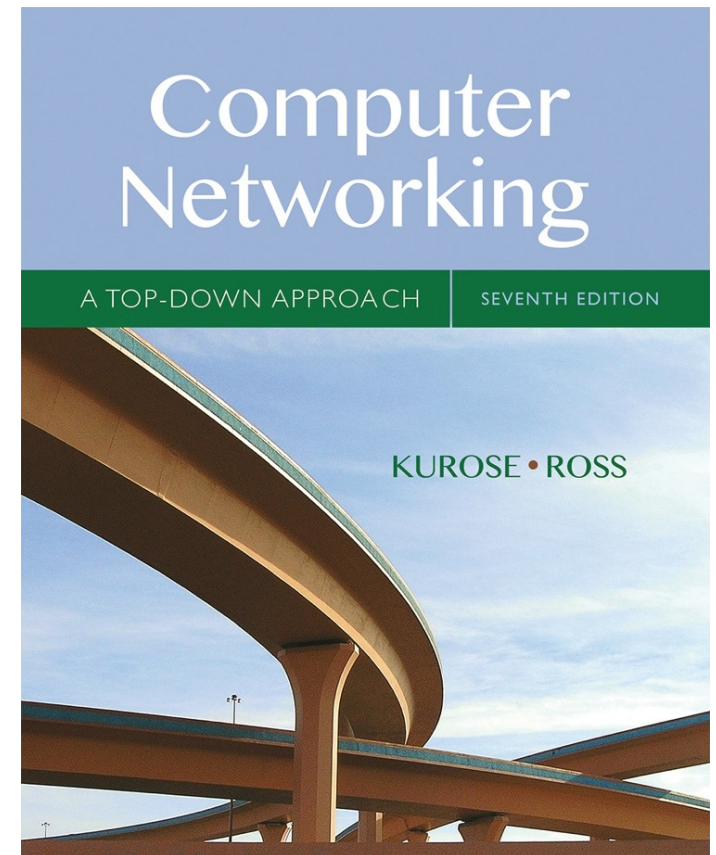
### A note on the use of these Powerpoint slides:

We're making these slides freely available to all (faculty, students, readers). They're in PowerPoint form so you see the animations; and can add, modify, and delete slides (including this one) and slide content to suit your needs. They obviously represent a *lot* of work on our part. In return for use, we only ask the following:

- If you use these slides (e.g., in a class) that you mention their source (after all, we'd like people to use our book!)
- If you post any slides on a www site, that you note that they are adapted from (or perhaps identical to) our slides, and note our copyright of this material.

Thanks and enjoy! JFK/KWR

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



## *Computer Networking: A Top Down Approach*

7<sup>th</sup> edition

Jim Kurose, Keith Ross

Pearson/Addison Wesley

April 2016

# Chapter 4: Outline

## 4.1 Overview of Network layer

- data plane
- control plane

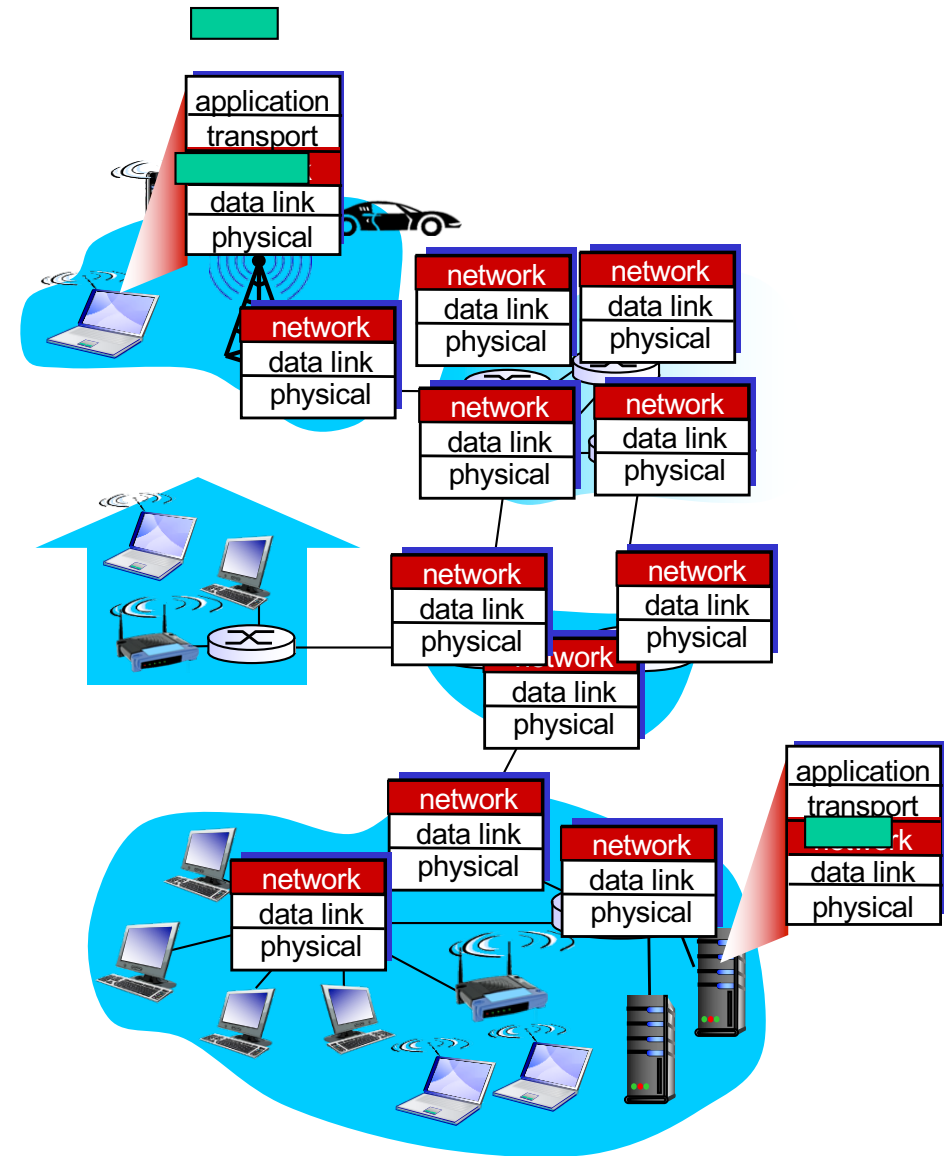
## 4.2 What's inside a router

## 4.3 IP: Internet Protocol

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

# 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



# Two key network-layer functions

## *network-layer functions:*

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

## *analogy: taking a trip*

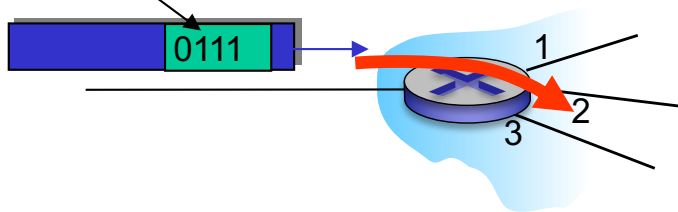
- *forwarding*: process of getting through single interchange
- *routing*: process of planning trip from source to destination

# 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

values in arriving packet header

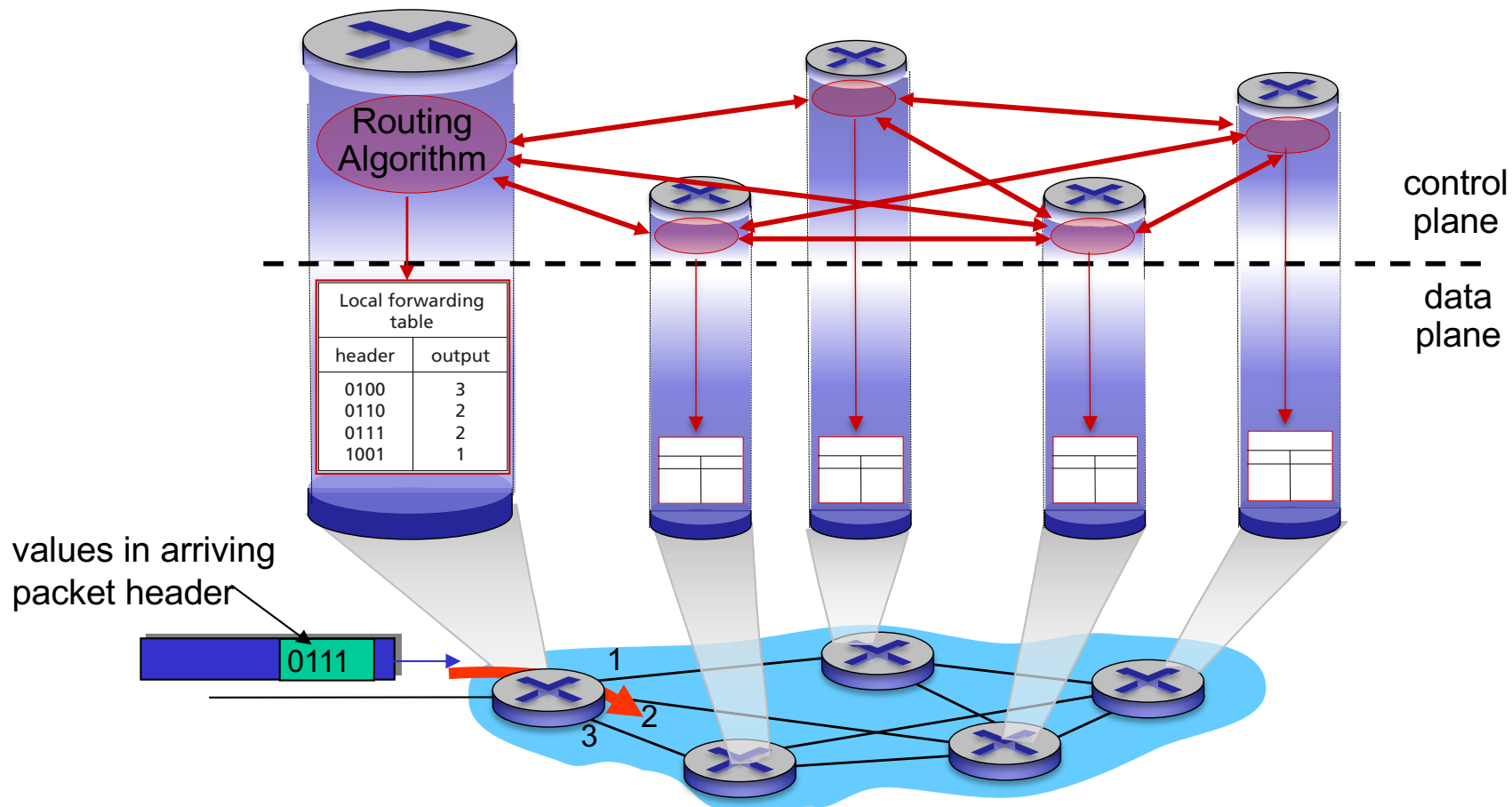


## *Control plane*

- 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

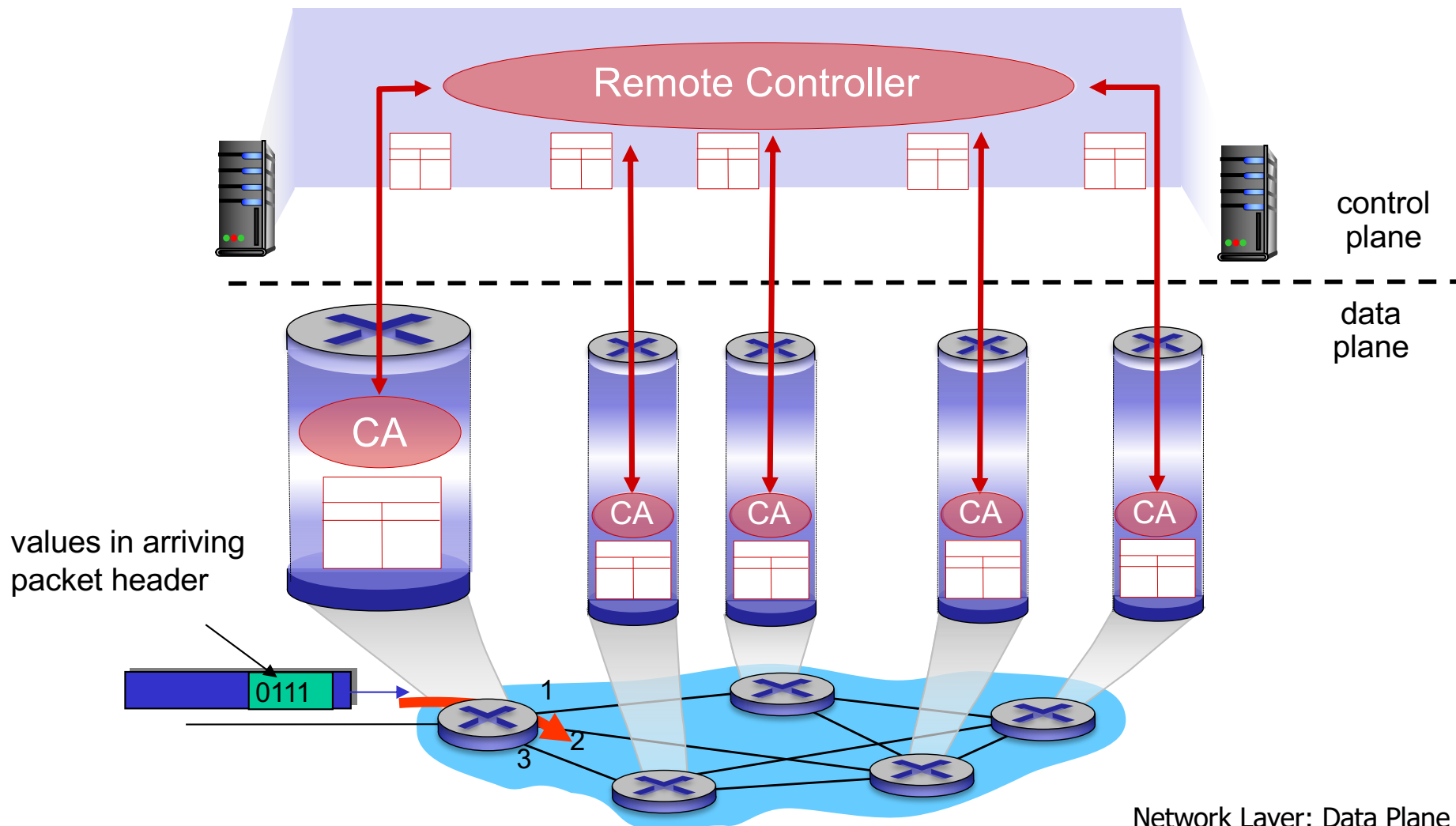
# Per-router control plane

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



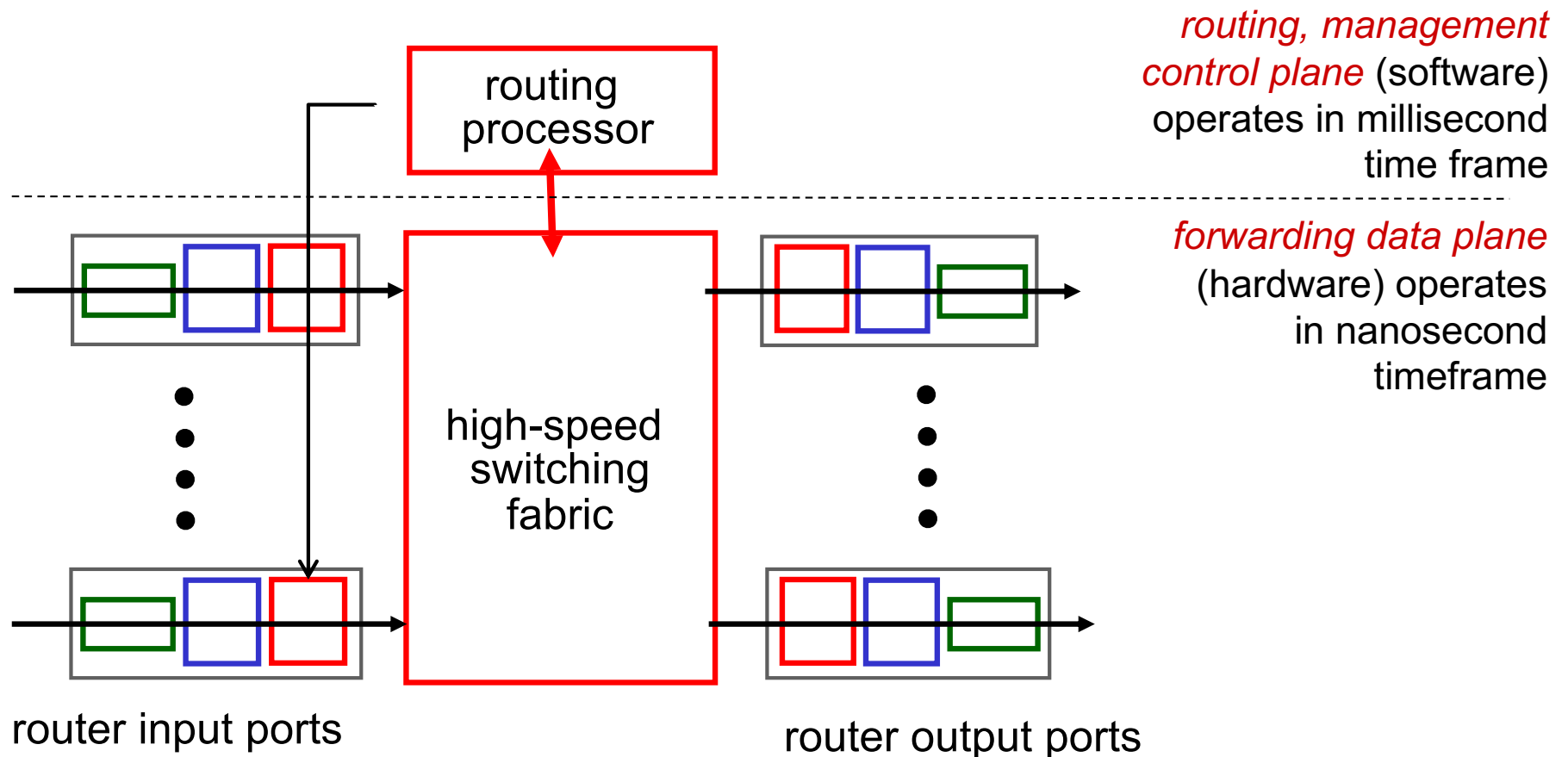
# Logically centralized control plane

A distinct (typically remote) controller interacts with local control agents (CAs)



# Router architecture overview

- high-level view of generic router architecture:





# Destination-based forwarding

*forwarding table*

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?

# Longest prefix matching

## *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 00011000 **10101010**

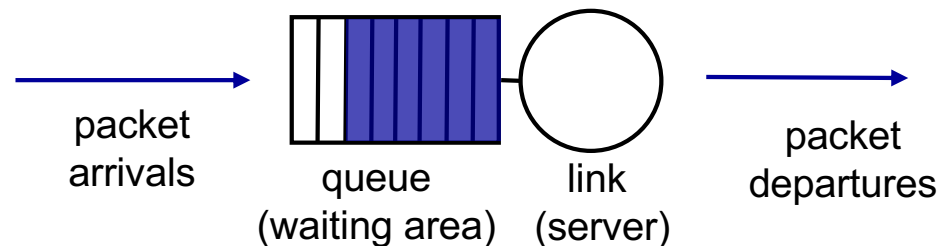
which interface?

# 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 up ~1M routing table entries in TCAM

# Scheduling mechanisms

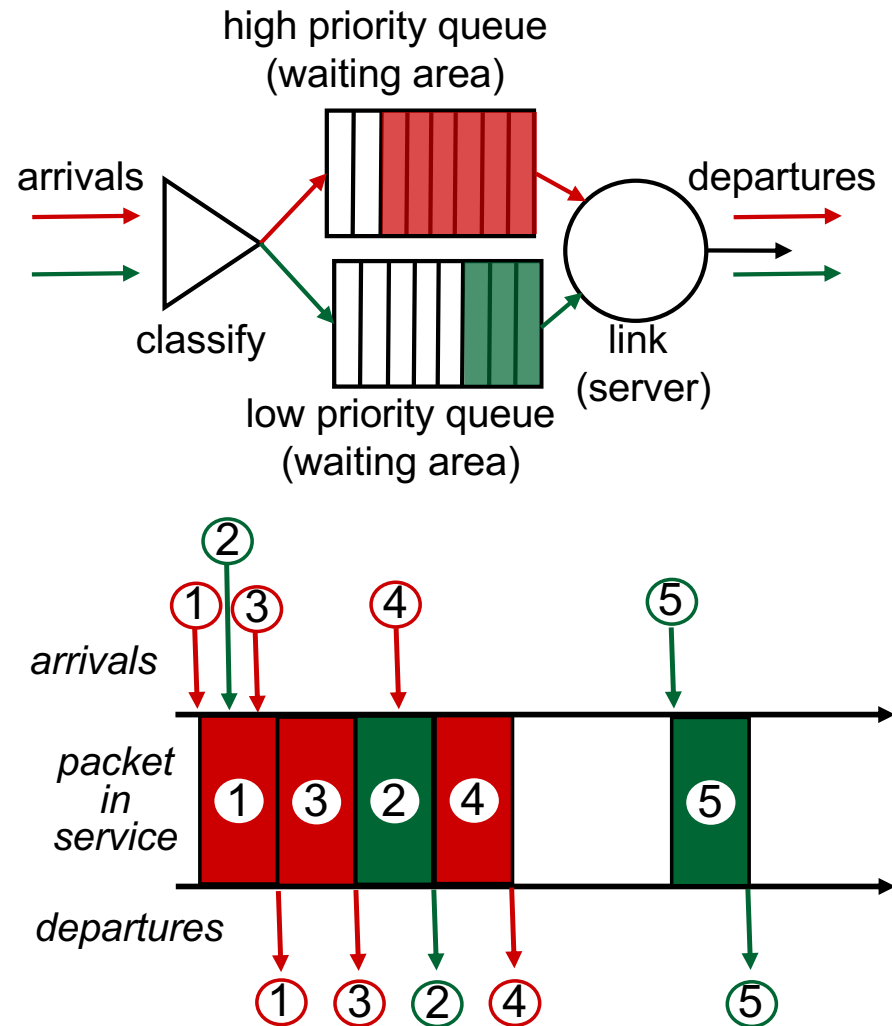
- *scheduling*: choose next packet to send on link
- *FIFO (first in first out) scheduling*: send in order of arrival to queue
  - real-world example?
  - *discard policy*: if packet arrives to full queue: who to discard?
    - *tail drop*: drop arriving packet
    - *priority*: drop/remove on priority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



# Scheduling policies: priority

*priority scheduling*: send highest priority queued packet

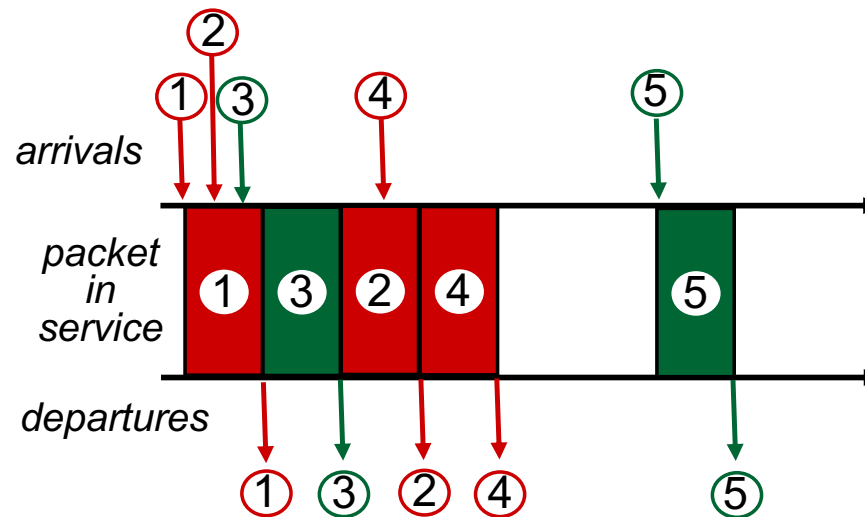
- multiple *classes*, with different priorities
  - class may depend on marking or other header info, e.g., IP source/destination, port numbers, etc.
  - real world example?



# Scheduling policies: still more

## *Round Robin (RR) scheduling:*

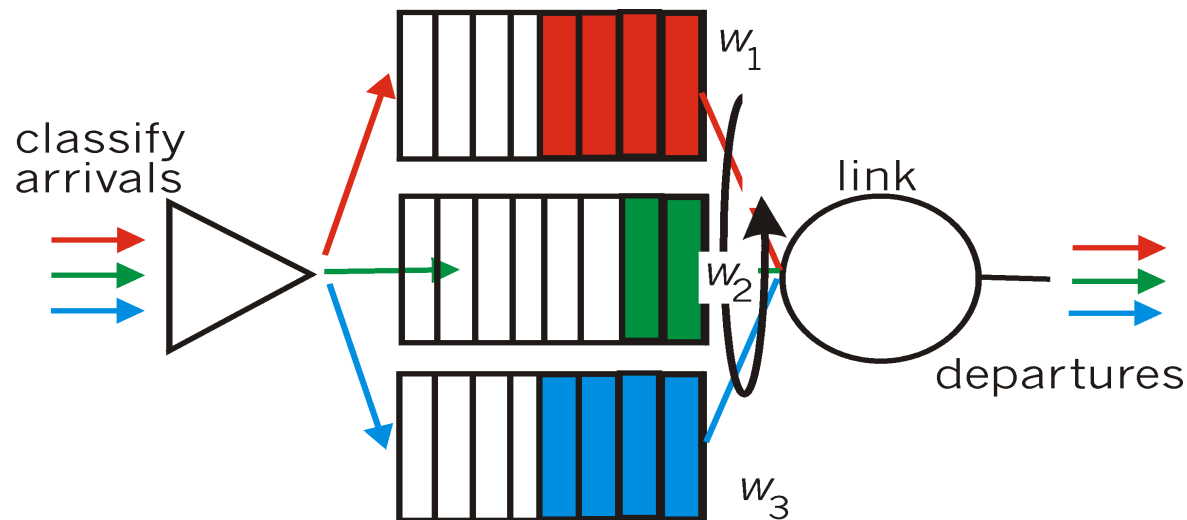
- multiple classes
- cyclically scan class queues, sending one complete packet from each class (if available)
- real world example?



# Scheduling policies: still more

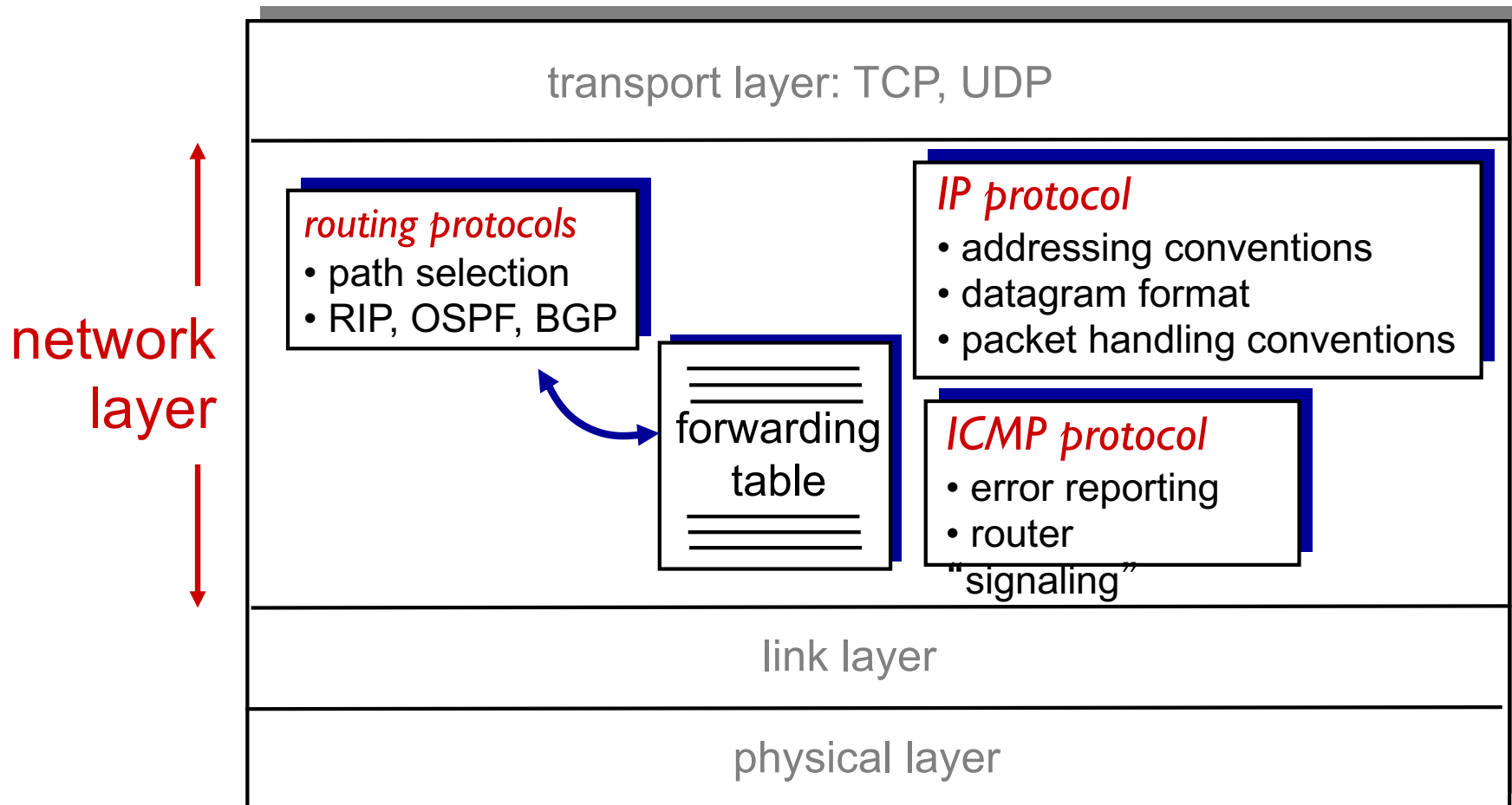
## *Weighted Fair Queuing (WFQ):*

- generalized Round Robin
- each class gets weighted amount of service in each cycle
- real-world example?



# The Internet network layer

host, router network layer functions:





# IP datagram format

IP protocol version  
number

header length  
(4 bytes)

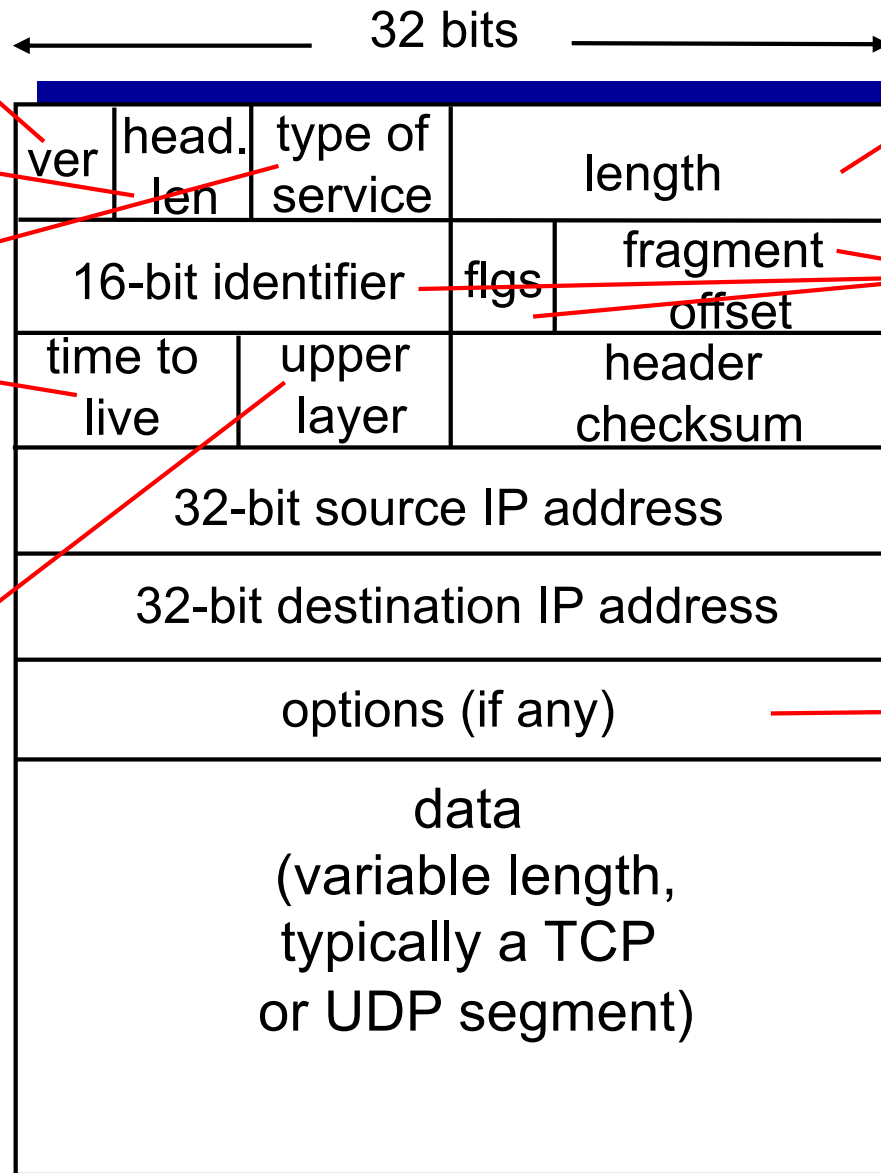
“type” of data

max number  
remaining hops  
(decremented at  
each router)

upper layer protocol  
to deliver payload to

*how much overhead?*

- ❖ 20 bytes of TCP
- ❖ 20 bytes of IP
- ❖ = 40 bytes + app layer overhead



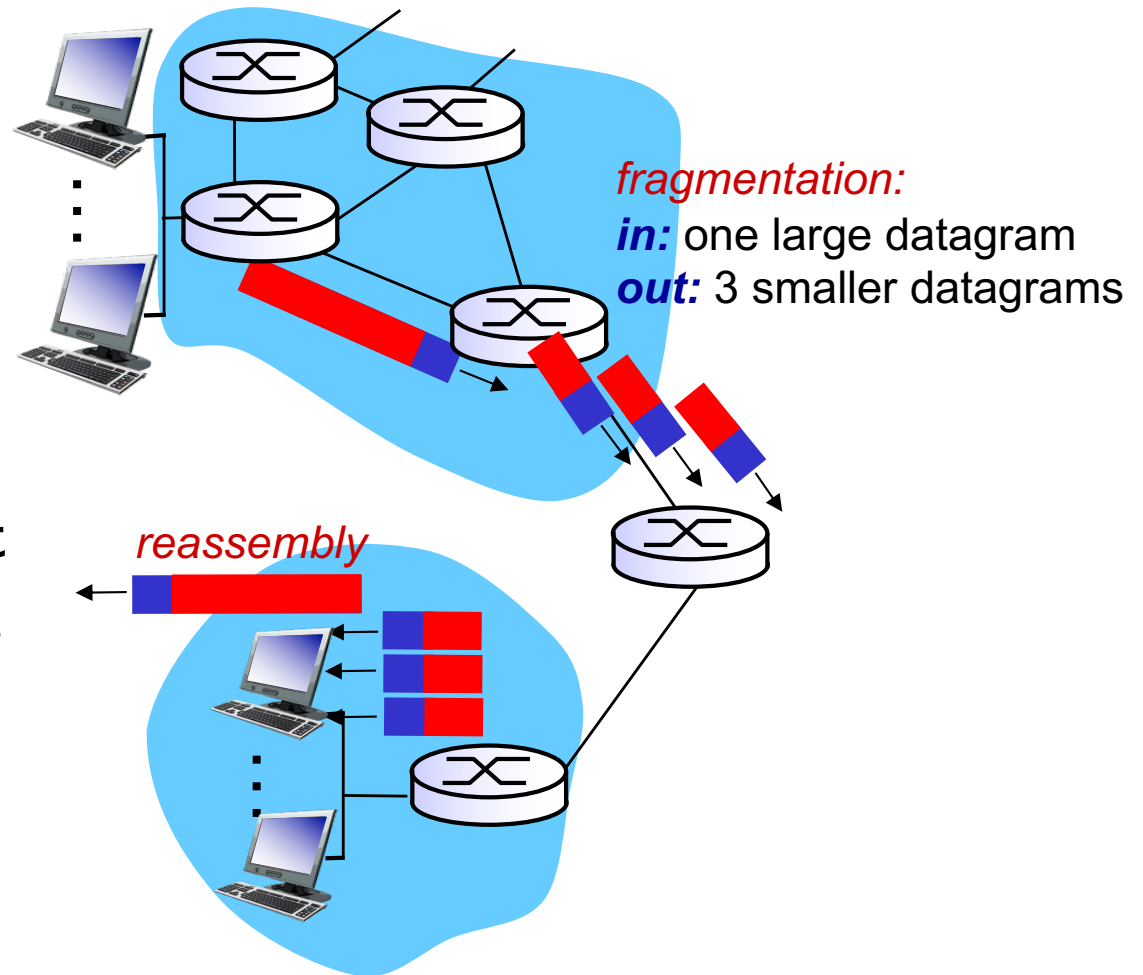
total datagram  
length (bytes)

for  
fragmentation/  
reassembly

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

# 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



# IP fragmentation, reassembly

*example:*

- ❖ 4000-byte datagram
- ❖ MTU = 1500 bytes

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

*one large datagram becomes  
several smaller datagrams*

1480 bytes in  
data field

offset =  
 $1480/8$

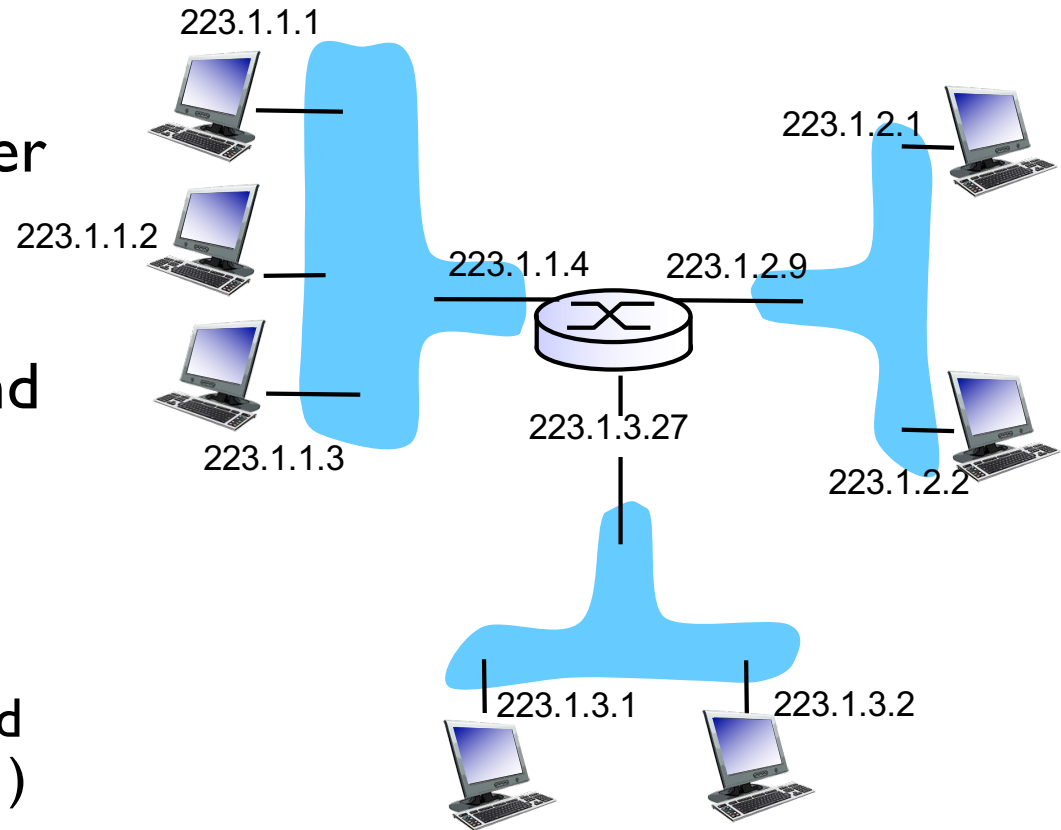
	length	ID	fragflag	offset	
	=1500	=x	=1	=0	

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

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

# IP addressing: introduction

- **IP address:** 32-bit identifier for host, router interface
- **interface:** connection between host/router and physical link
  - router's typically have multiple interfaces
  - host typically has one or two interfaces (e.g., wired Ethernet, wireless 802.11)
- **IP addresses associated with each interface**



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

# Binary Numbers

- 00000000 in binary = 0 in decimal
- 00000001 in binary = 1 in decimal
- 00000010 in binary = 2 in decimal
- ...
- 11111110 in binary = 254 in decimal
- 11111111 in binary = 255 in decimal

# IP Address

- Written in dotted-decimal format for humans to read
  - Binary = 11000000 10101000 00000001 00000001
  - Dotted decimal format = 192.168.1.1
- Two parts specified by subnet mask
  - Network address
  - Host address
- Example
  - 192.168.2.26 IP address with 255.255.255.0 subnet mask
  - Network portion = 192.168.2.0
  - Host portion = 26
  - How many host IP address in the network 192.168.2.0 (with 255.255.255.0 subnet mask)?
    - 192.168.2.0 - 192.168.2.255 (total 256 IP addresses)
    - 192.168.2.26 is one of them

# Special IP Addresses

- Host bits are **all zero** describes network
  - 192.168.2.0/24 refers to the network
- Host bits **all ones** refers to directed broadcast
  - 192.168.2.255/24
- Broadcast
  - 255.255.255.255
- Multicast range
  - 224-239.0.0.0
- Private IP ranges (not supposed to be routed)
  - 10.0.0.0 - 10.255.255.255
  - 172.16.0.0 - 172.31.255.255
  - 192.168.0.0 - 192.168.255.255
- Loopback - 127.X.Y.Z like 127.0.0.1 (also called home)

# How many valid host IP addresses can I use?

---

- 192.168.2.0/24
  - 192.168.2.1-254
  - Why?
    - 192.168.2.0 reserved to specify network
    - 192.168.2.255 reserved for directed broadcast



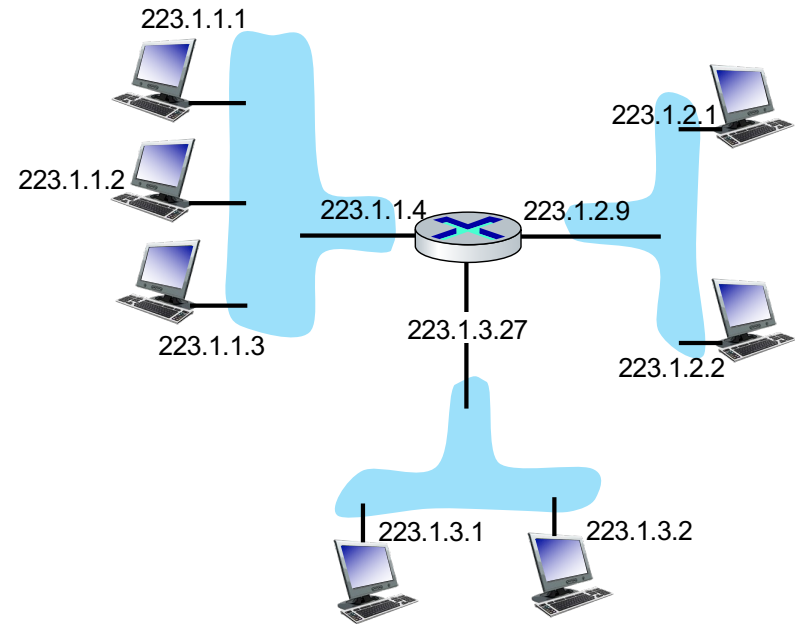
# Subnets

## ■ *What's a subnet ?*

- device interfaces that can physically reach each other **without passing through an intervening router**

## ■ IP addresses have structure:

- **subnet part:** devices in same subnet have common high order bits
- **host part: remaining** low order bits

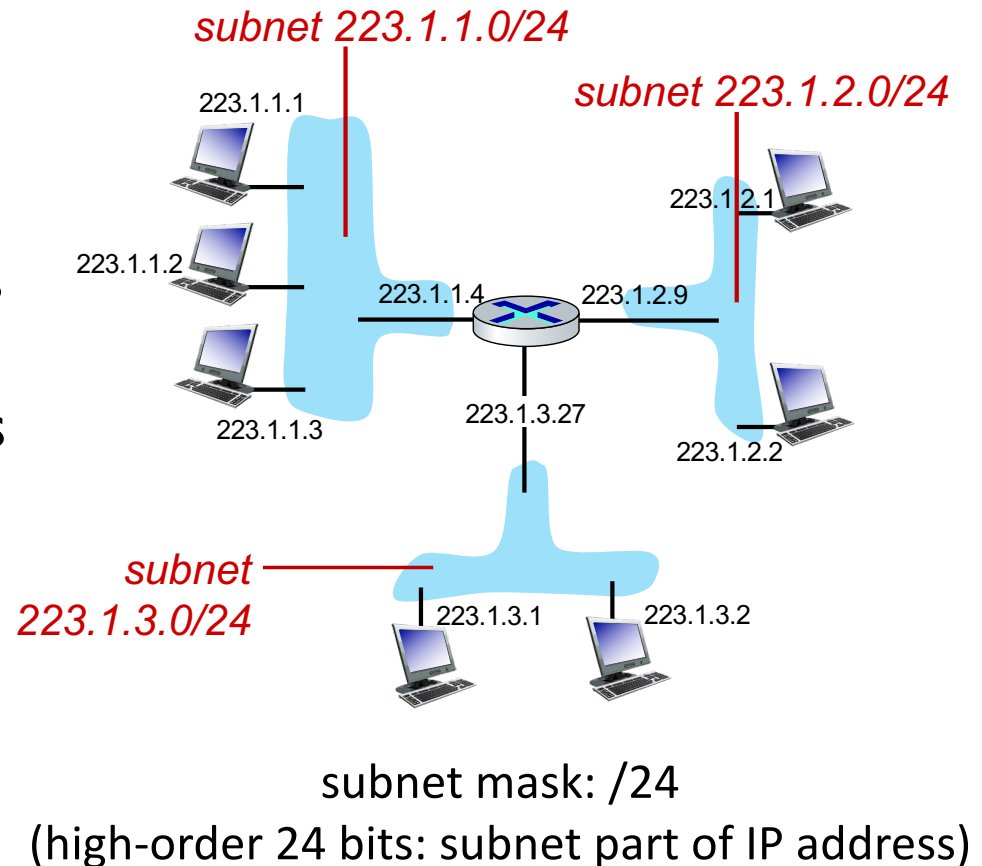


network consisting of 3 subnets

# Subnets

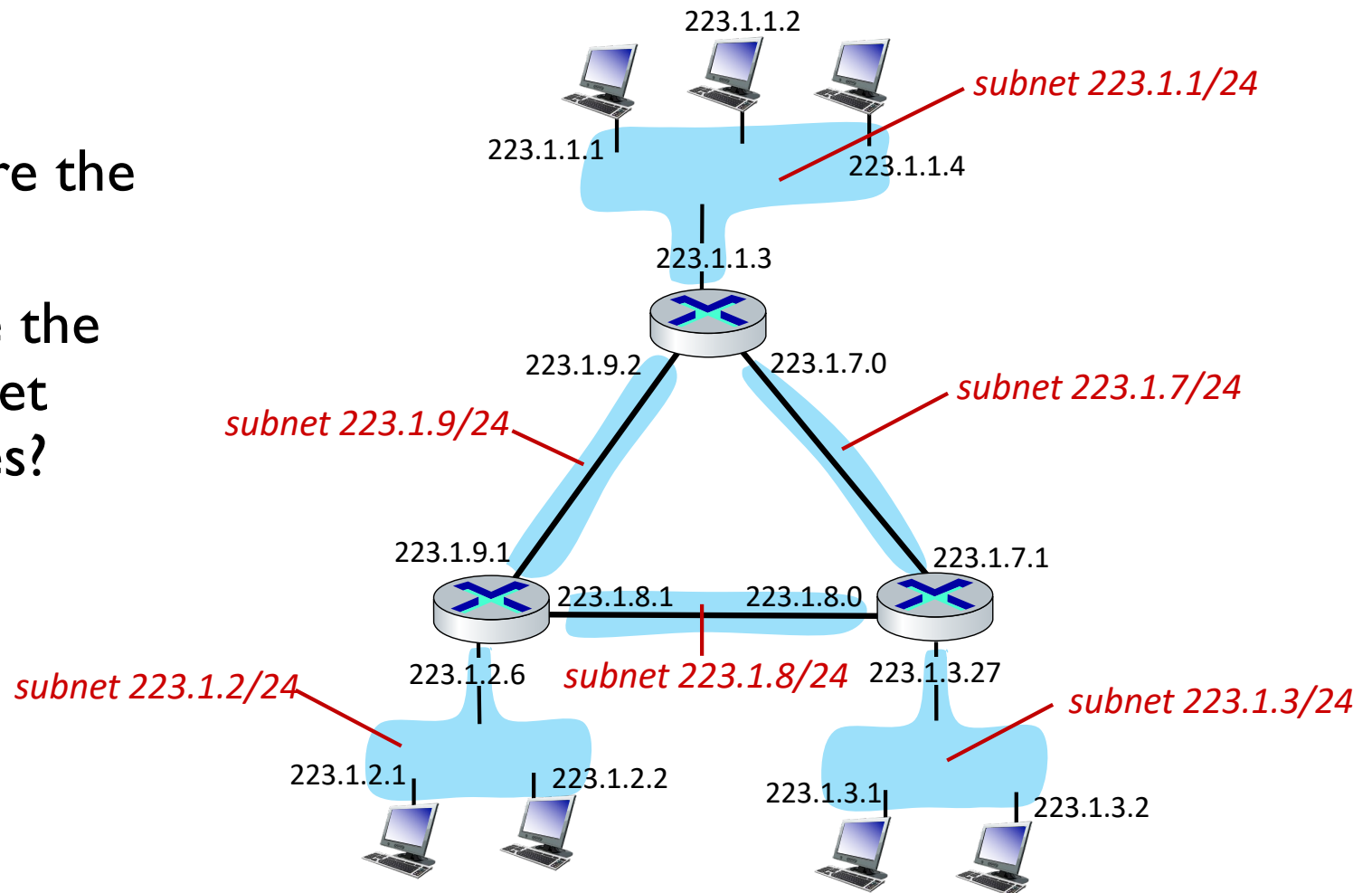
## *Recipe for defining subnets:*

- detach each interface from its host or router, creating “islands” of isolated networks
- each isolated network is called a **subnet**



# Subnets

- where are the subnets?
- what are the /24 subnet addresses?

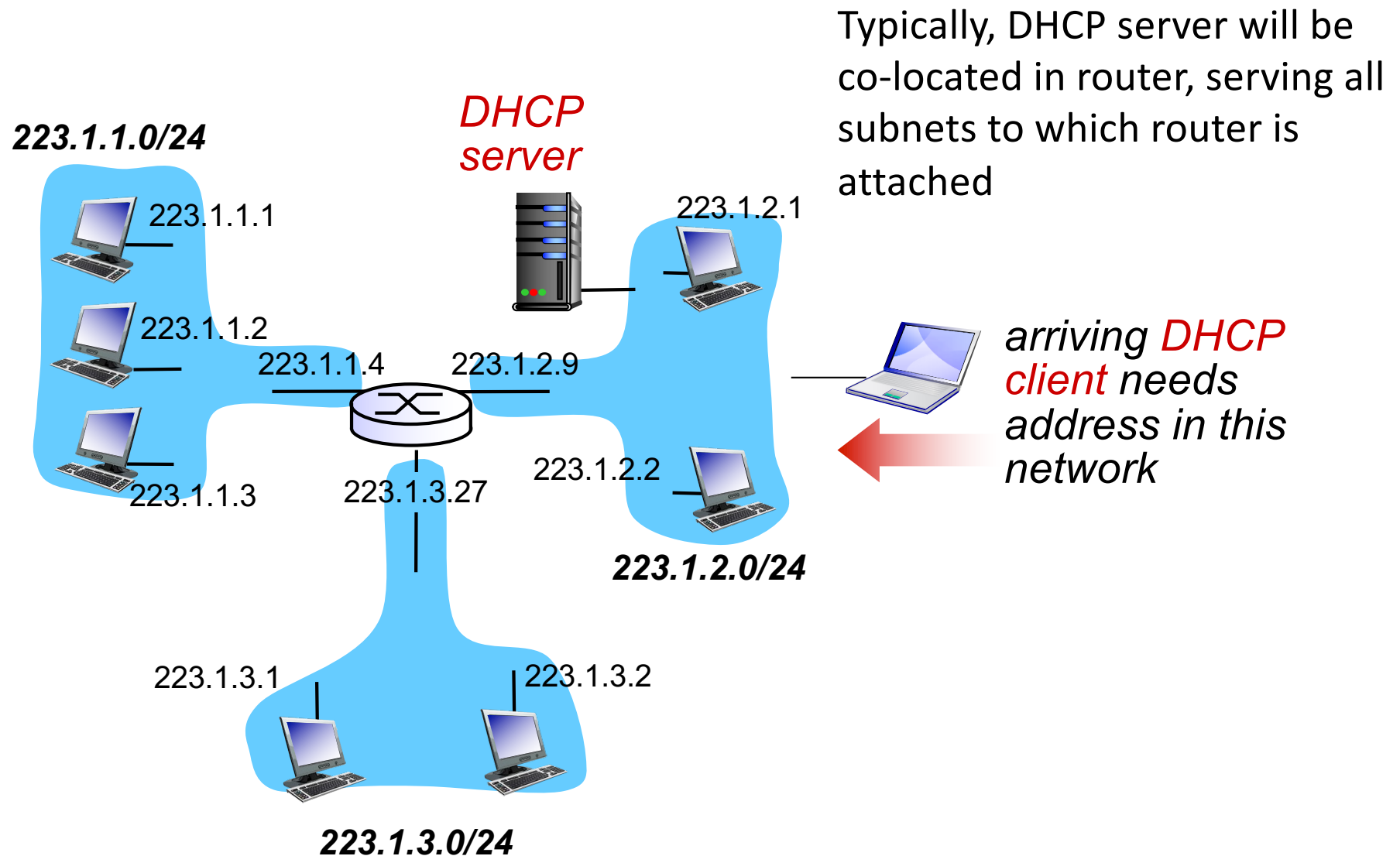


# IP addresses: how to get one?

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

- hard-coded by system admin in a file
  - Windows: control-panel->network->configuration  
->tcp/ip->properties
  - UNIX: /etc/rc.config
- **DHCP: D**ynamic **H**ost **C**onfiguration **P**rotocol:  
dynamically get address from as server
  - “plug-and-play”
  - allow host to *dynamically* obtain its IP address  
from network server when it joins network

# DHCP client-server scenario



# DHCP client-server scenario

DHCP server: 223.1.2.5

DHCP discover

arriving  
client



Broadcast: is there a  
DHCP server out there?

DHCP offer

Broadcast: I'm a DHCP  
server! Here's an IP  
address you can use

DHCP request

Broadcast: OK. I'll take  
that IP address!

DHCP ACK

Broadcast: OK. You've  
got that IP address!

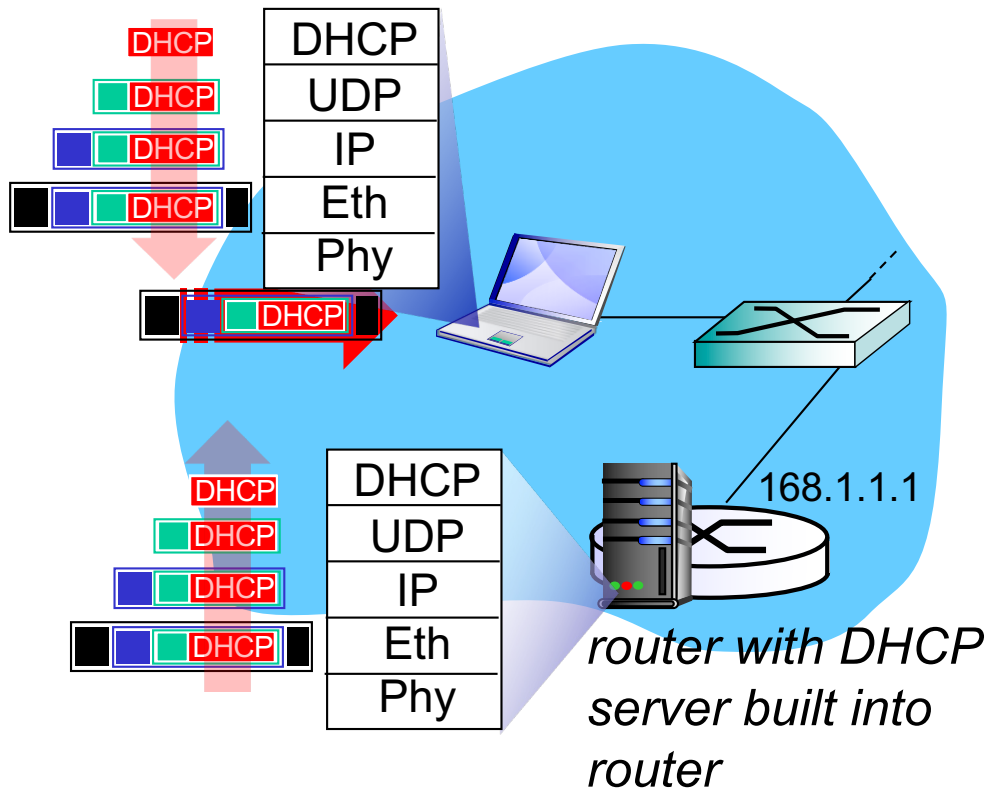
The two steps above can  
be skipped "if a client  
remembers and wishes to  
reuse a previously  
allocated network address"  
[RFC 2131]

# DHCP: more than IP addresses

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

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

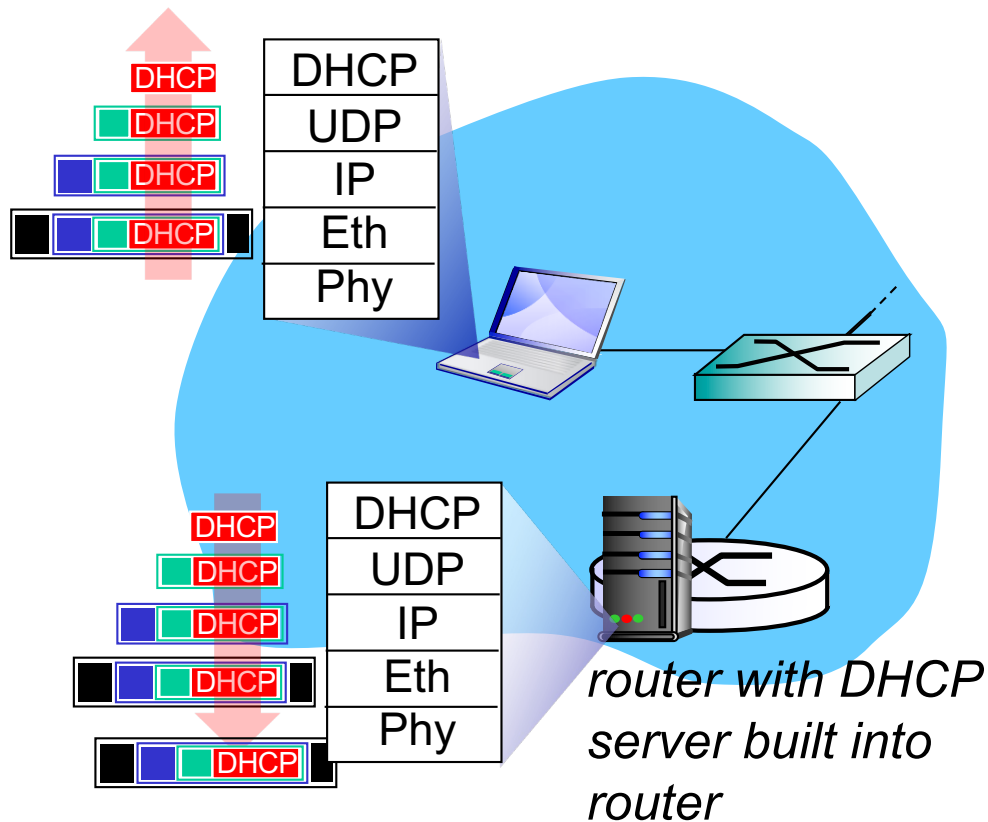
# DHCP: example



- 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



# DHCP: example



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

# IP addressing: the last word...

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

**A:** **ICANN:** Internet Corporation for Assigned Names and Numbers <http://www.icann.org/>

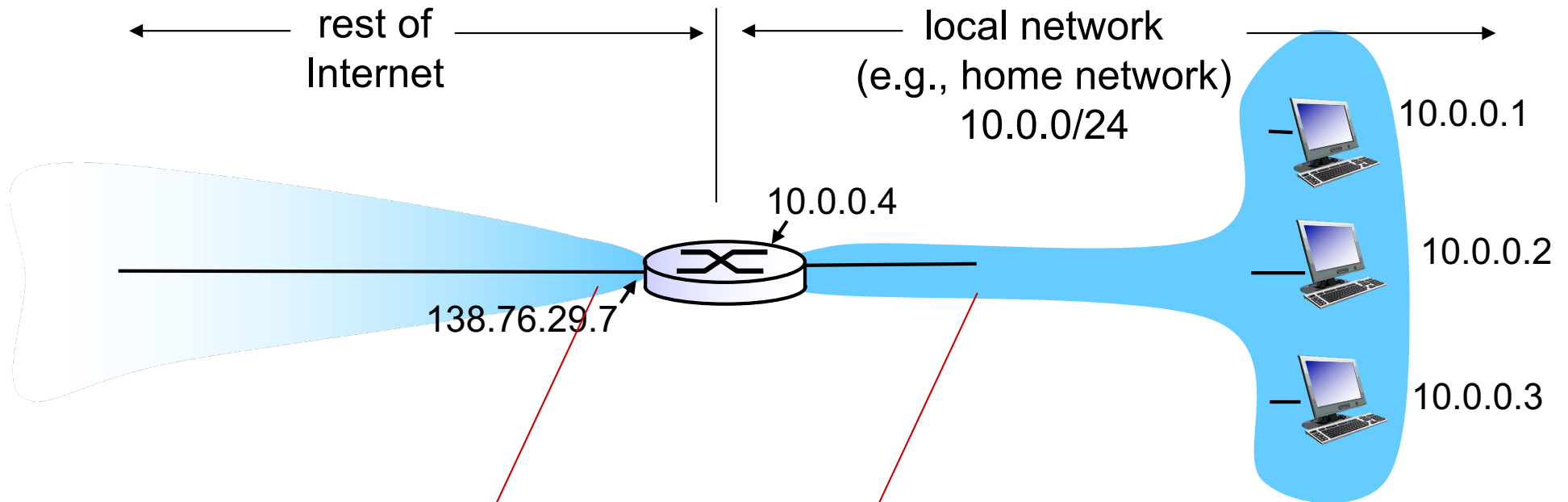
- allocates IP addresses, through **5 regional registries (RRs)** (who may then allocate to local registries)
- manages DNS root zone, including delegation of individual TLD (.com, .edu , ... ) management

**Q:** are there enough 32-bit IP addresses?

- ICANN allocated last chunk of IPv4 addresses to RRs in 2011
- NAT (next) helps IPv4 address space exhaustion
- IPv6 has 128-bit address space

"Who the hell knew how much address space we needed?" Vint Cerf (reflecting on decision to make IPv4 address 32 bits long)

# 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

Local network uses just one IP address as far as outside world is concerned. **Why?**

- 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)

# NAT: network address translation

*implementation:* NAT router must:

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

Q: When the remote host responds, what is the destination address and port#?

A: (NAT IP address, new port #)

- *remember (in NAT translation table)* every translation pair  
(source IP address, port #) to (NAT IP address, new port #)
- *incoming datagrams: replace* destination fields of every incoming datagram

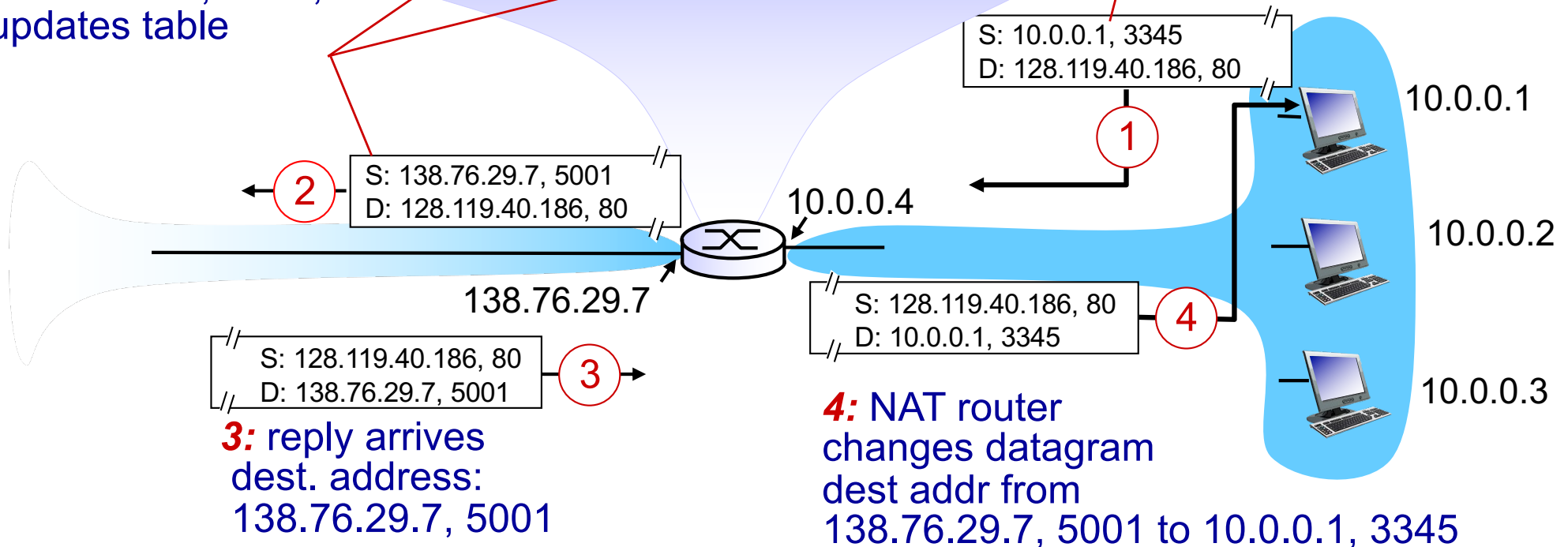
(NAT IP address, new port #) to (source IP address, port #)  
stored in NAT table

# 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



# NAT: network address translation

- 16-bit port-number field:
  - 60,000 simultaneous connections with a single LAN-side address!
- NAT is controversial, why?
  - 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?

# 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

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

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