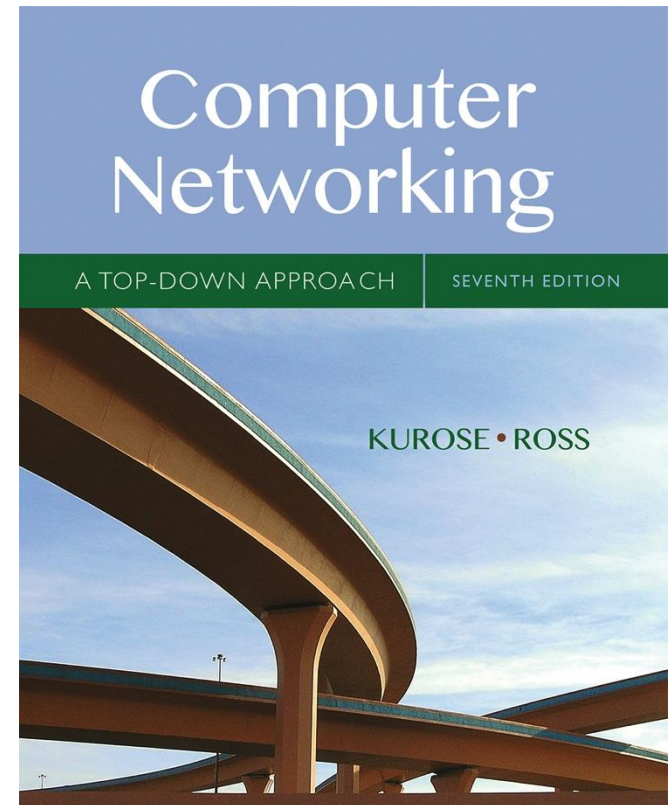# Chapter 4
# Network Layer:
# The Data Plane

## 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

*Computer Networking: A Top Down Approach*

7th edition
Jim Kurose, Keith Ross
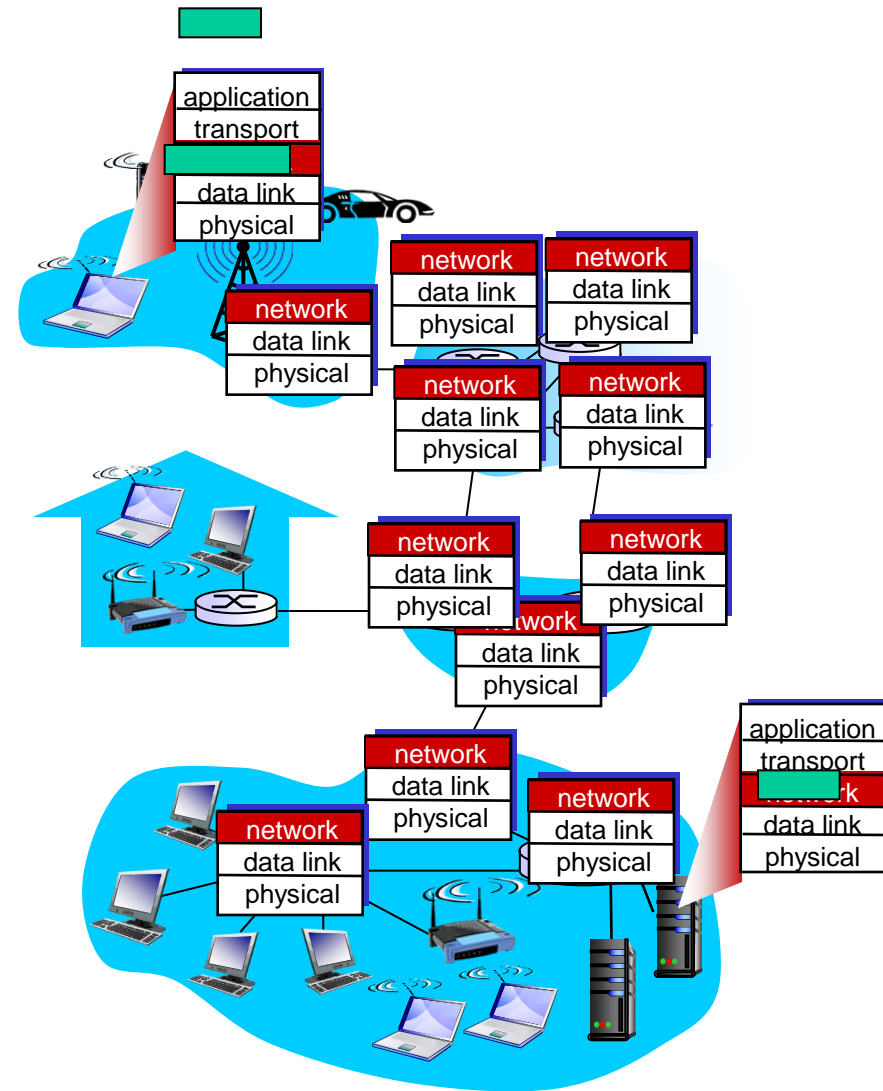Pearson/Addison Wesley
April 2016

# Chapter 4: outline

# 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

- transports segments 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

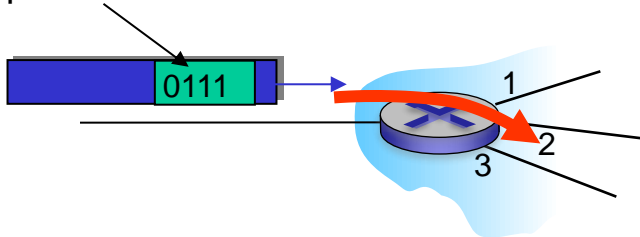# Two key network-layer functions

*network-layer functions:*

- *forwarding:* move packets from router's input to appropriate router output → Data Plane (chapter#4)

- *routing:* determine route taken by packets from source to destination →Control Plane (Chapter#5)
  - *routing algorithms*

# Network layer: data plane, control plane

## Data plane (chapter#4)

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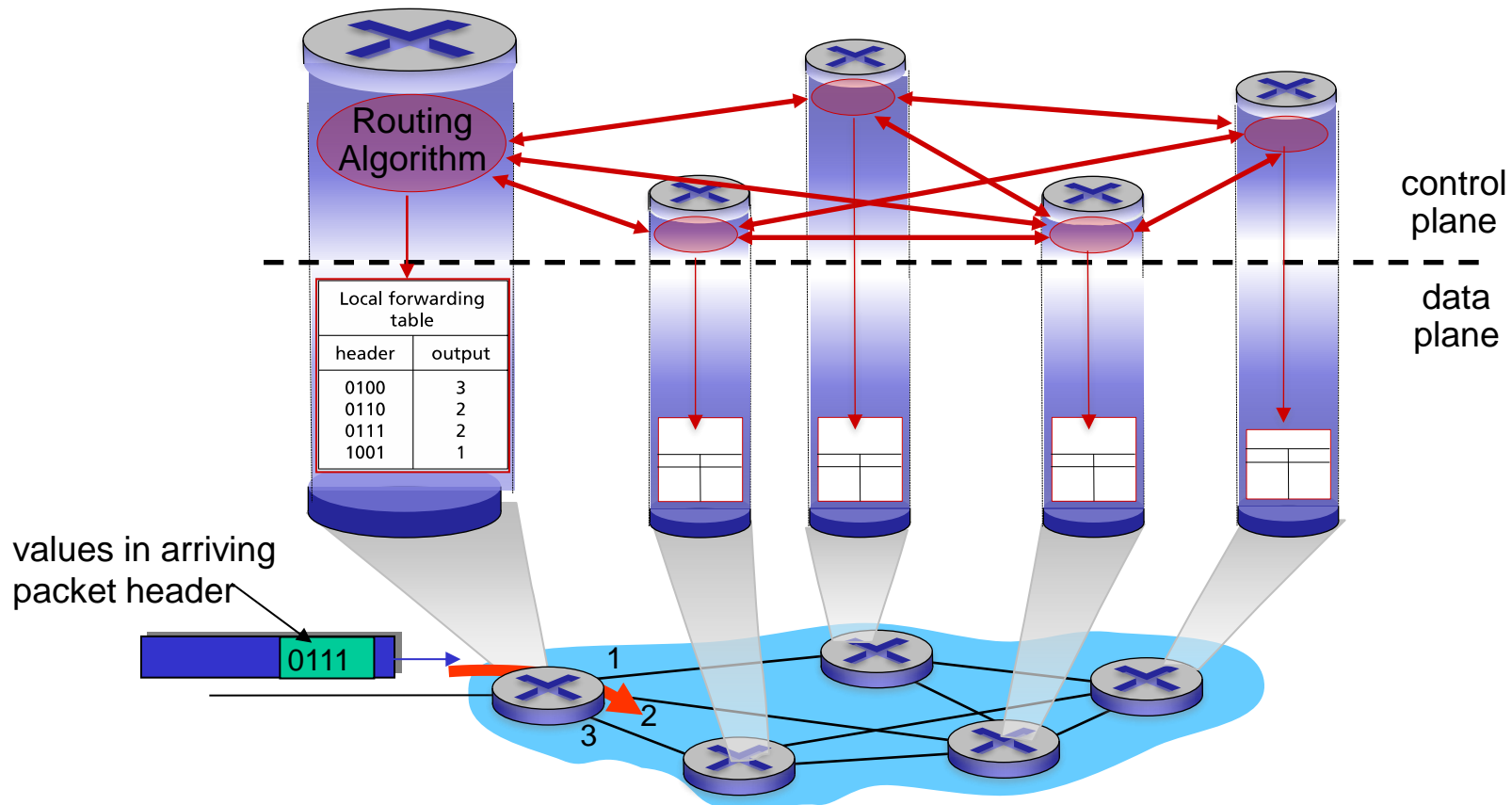
values in arriving
packet header



## Control plane (chapter#5)

- 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



control plane

data plane

Routing Algorithm

Local forwarding table

| header | output |
|--------|--------|
| 0100 | 3 |
| 0110 | 2 |
| 0111 | 2 |
| 1001 | 1 |

values in arriving packet header

0111

1

3    2

# Network service model

*Q:* What possible services network layer could provide?

*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

However, Internet's network layer provides a single service, known as *best-effort service*.

# Chapter 4: outline

# Router architecture overview

- high-level view of generic router architecture:



*routing, management control plane* operates in millisecond time frame

*forwarding data plane* operates in nanosecond timeframe

routing processor

high-speed switching fabric

router input ports

router output ports

# Input port functions



physical layer:
bit-level reception

data link layer:
 e.g., Ethernet
 see chapter 6

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

# Destination-based forwarding

| Destination Address Range | Link Interface |
|---|---|
| 11001000 00010111 00010000 00000000<br>through<br>11001000 00010111 00010111 11111111 | 0 |
| 11001000 00010111 00011000 00000000<br>through<br>11001000 00010111 00011000 11111111 | 1 |
| 11001000 00010111 00011001 00000000<br>through<br>11001000 00010111 00011111 11111111 | 2 |
| otherwise | 3 |

# Longest prefix matching

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

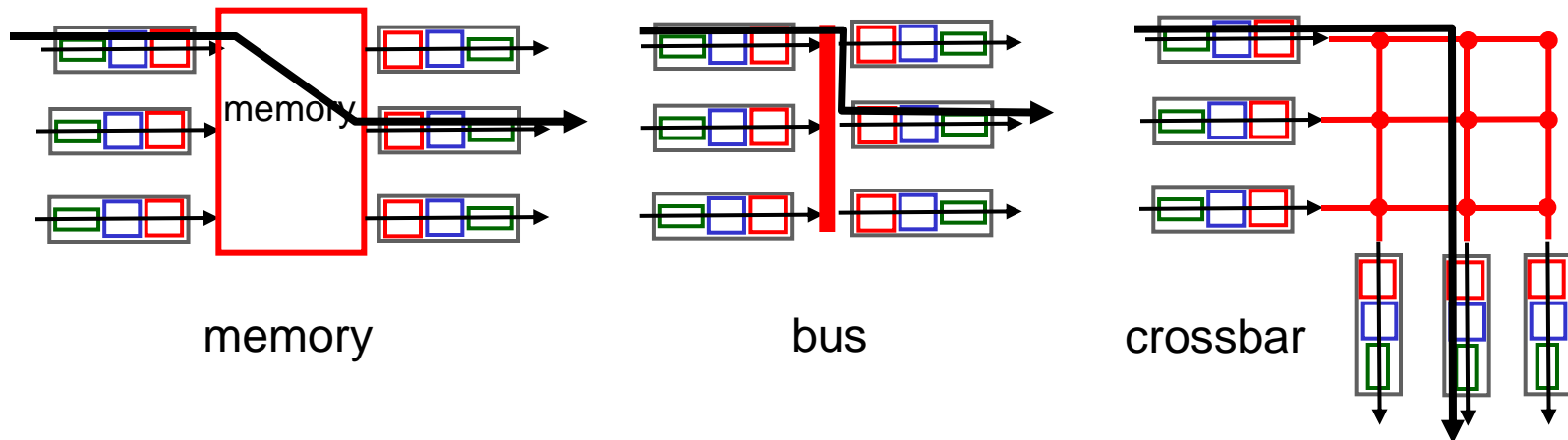| Destination Address Range | Link interface |
|---|---|
| 11001000 00010111 00010*** ******** | 0 |
| 11001000 00010111 00011000 ******** | 1 |
| 11001000 00010111 00011*** ******** | 2 |
| otherwise | 3 |

# Longest prefix matching

- **longest prefix matching:** often performed using ternary content addressable memories (TCAMs)
  - Cisco Catalyst: holds up to ~1M routing table entries in TCAM

# Switching fabrics

- Switching fabric:
  - Connects input ports to its output ports
  - forwards packet from input buffer to appropriate output buffer
- switching rate, $R_{switch}$ : rate at which packets can be transferred from inputs to outputs
  - often measured as multiple of input/output line rate denoted as $R_{line}$
  - N inputs: switching rate $R_{switch}$ , N times line rate is desirable
- three types of switching fabrics

(1) Switching via memory (2) via Bus (3) via interconnection network



memory                    bus                    crossbar

# Switching via memory

*first generation routers:*

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

input
port
(e.g.,
Ethernet)

CPU/memory

output
port
(e.g.,
Ethernet)

system bus

# Switching via a bus

- Only one 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

bus

# Switching via interconnection network

- overcome bus bandwidth limitations
- Allows datagrams to be transferred from multiple input ports to multiple output ports simultaneously
- Cisco 12000 series switches use a crossbar switching network

crossbar

# Input port queuing

- fabric slower than input ports combined, $R_{switch} < N \times R_{line}$
  - 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



output port contention:
only one red datagram can be
transferred.
*lower red packet is blocked*

one packet time later:
green packet
experiences HOL
blocking

# Output ports: Buffering and Queueing



- *buffering* required when datagrams arrive from fabric faster than the transmission rate

- *scheduling discipline* chooses among queued datagrams for transmission

- *queueing (delay) and loss due to output port buffer overflow!*

# 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

packet arrivals

queue
(waiting area)

link
(server)

packet departures

# 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.
  - real world example?



high priority queue (waiting area)

arrivals

classify

low priority queue (waiting area)

departures

link (server)

arrivals

packet in service

departures

# 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?

# Chapter 4: outline

# IP datagram format

32 bits

IP protocol version number

header length (bytes)

"type" of data (real-time/non-rea time)

max number remaining hops (decremented at each router)

upper layer protocol to deliver payload to

| ver | head. len | type of service | length |
|---|---|---|---|
| 16-bit identifier | | flgs | fragment offset |
| time to live | upper layer Protocol | header checksum | |
| 32 bit source IP address | | | |
| 32 bit destination IP address | | | |
| options (if any) | | | |
| data (variable length, typically a TCP or UDP segment) | | | |

total datagram length (bytes)

for fragmentation/ reassembly

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

*how much overhead?*
- ❖ 20 bytes of TCP
- ❖ 20 bytes of IP
- ❖ = 40 bytes

# IP fragmentation, reassembly

- network links have MTU (max. transfer unit) - 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 (receiver host)
  - IP header bits used to identify, order related fragments

*fragmentation:*
*in:* one large datagram
*out:* 3 smaller datagrams

*reassembly*

# IP fragmentation, reassembly

*example:*

- ❖ 4000 byte datagram including IP header
- ❖ MTU = 1500 bytes including IP header

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

*one large datagram becomes several smaller datagrams*

1480 bytes in data field

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

offset = 1480/8

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

offset = 2960/8

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

# IP fragmentation, reassembly

*example:*

❖ 8060 byte datagram
   including IP header

❖ MTU = 1500 bytes
   including IP header

| | length =8060 | ID =x | fragflag =0 | offset =0 | | |
|---|---|---|---|---|---|---|

# Chapter 4: outline
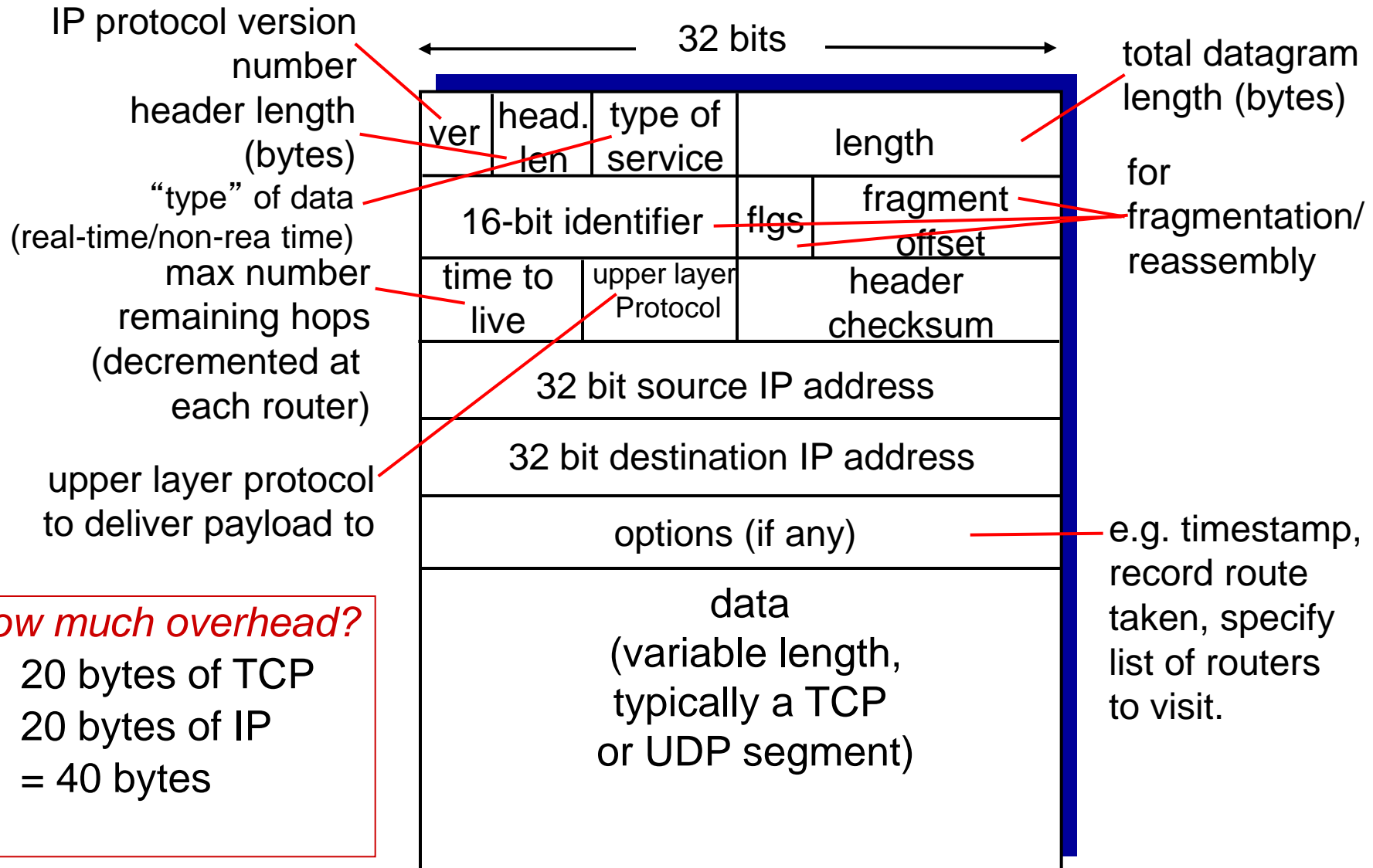
4.1 Overview of Network layer
- data plane
- control plane
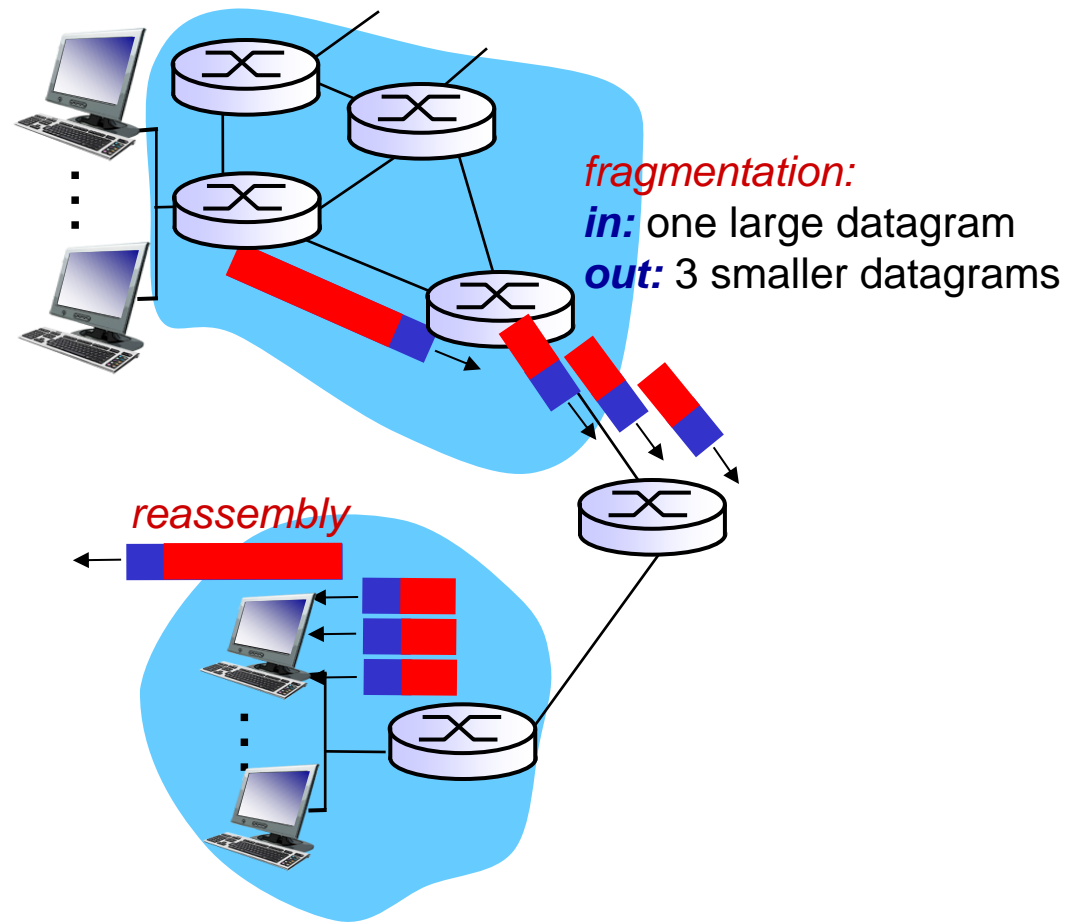
4.2 What's inside a router

4.3 IP: Internet Protocol
- datagram format
- fragmentation
- IPv4 addressing
- network address translation
- IPv6

# 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 active interface
- *IP addresses associated with each interface*

| Address Class | Bit Pattern of First Byte | First Byte Decimal Range | Host Assignment Range in Dotted Decimal |
|---|---|---|---|
| A | 0xxxxxxx | 1 to 127 | 1.0.0.1 to 126.255.255.254 |
| B | 10xxxxxx | 128 to 191 | 128.0.0.1 to 191.255.255.255.254 |
| C | 110xxxxx | 192 to 223 | 192.0.0.1 to 223.255.255.254 |
| D | 1110xxxx | 224 to 239 | 224.0.0.1 to 239.255.255.254 |
| E | 11110xxx | 240 to 255 | 240.0.0.1 to 255.255.255.255 |

# 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 active interface

- *IP addresses associated with each interface*

223.1.1.1

223.1.1.2

223.1.1.3

223.1.1.4    223.1.2.9

223.1.2.1

223.1.2.2

223.1.3.27

223.1.3.1    223.1.3.2

223.1.1.1 = 11011111 00000001 00000001 00000001

223            1            1            1

# 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*

- each isolated network is called a *subnet*

*223.1.1.0/24*

*223.1.2.0/24*

223.1.1.1

223.1.1.2

223.1.1.4   223.1.2.9

223.1.2.1

223.1.2.2

223.1.1.3   223.1.3.27

subnet

223.1.3.1    223.1.3.2

*223.1.3.0/24*

subnet mask: /24

# IP addressing: CIDR

CIDR: Classless InterDomain Routing
- subnet portion of address of arbitrary length
- address format: a.b.c.d/x, where x is # of bits in subnet portion of address



| subnet part | | host part |
|---|---|---|
| 11001000  00010111  00010000 | | 00000000 |

200.23.16.0/23

# 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 | 11001000  00010111  0001<u>0000</u> | 00000000 | 200.23.16.0/20 |
| | | | |
| Organization 0 | <u>11001000  00010111  0001**000**0</u> | 00000000 | 200.23.16.0/23 |
| Organization 1 | <u>11001000  00010111  0001**001**0</u> | 00000000 | 200.23.18.0/23 |
| Organization 2 | <u>11001000  00010111  0001**010**0</u> | 00000000 | 200.23.20.0/23 |
| ... | ..... | .... | .... |
| Organization 7 | <u>11001000  00010111  0001**111**0</u> | 00000000 | 200.23.30.0/23 |

# Hierarchical addressing: route aggregation

hierarchical addressing allows efficient advertisement of routing information:

Organization 0
200.23.16.0/23

Organization 1
200.23.18.0/23

Organization 2
200.23.20.0/23

Organization 7
200.23.30.0/23

ISP-A

ISP-B

"Announcement: Send me anything with addresses beginning 200.23.16.0/20"

"Announcement: Send me anything with addresses beginning 199.31.0.0/16"

Internet

# Hierarchical addressing: more specific routes

Organization now moved under ISP-B
ISP-B has a more *specific* route to Organization 1

Organization 0
200.23.16.0/23

Organization 2
200.23.20.0/23

Organization 7
200.23.30.0/23

ISP-A

"**Announcement:** Send me anything with addresses beginning 200.23.16.0/20"

Internet

ISP-B

Organization 1
200.23.18.0/23

"**Announcement:** Send me anything with addresses beginning 199.31.0.0/16 or 200.23.18.0/23"

| ISP-A | 11001000 00010111 00010000 00000000 | 200.23.16.0/20 |
|---|---|---|
| Organization 1 | 11001000 00010111 00010010 00000000 | 200.23.18.0/23 |

# 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 addresses
- manages DNS
- assigns domain names, resolves disputes

# NAT: network address translation

rest of Internet

local network (e.g., home network)
10.0.0/24

10.0.0.1

10.0.0.4

10.0.0.2

138.76.29.7

10.0.0.3

*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

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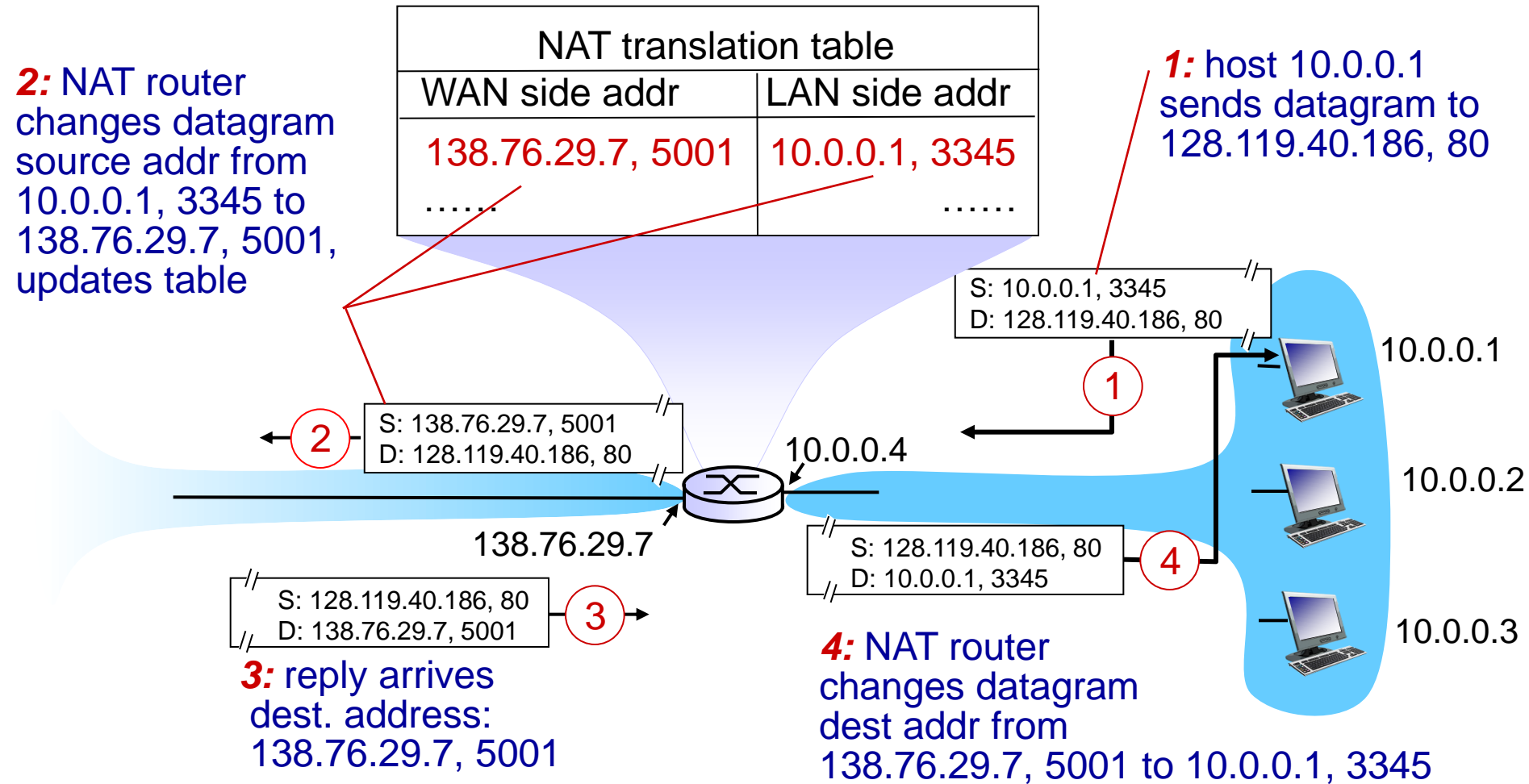
- range of addresses not needed from ISP:  just one IP address for all devices
- can change addresses of devices in local network without notifying outside world
- can change ISP without changing addresses of devices in local network
- devices inside local net not explicitly addressable, visible by outside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

# NAT: network address translation

*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

# 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

S: 10.0.0.1, 3345
D: 128.119.40.186, 80

① 

10.0.0.1

S: 138.76.29.7, 5001
D: 128.119.40.186, 80

②

10.0.0.4

138.76.29.7

S: 128.119.40.186, 80
D: 10.0.0.1, 3345

④

10.0.0.2

S: 128.119.40.186, 80
D: 138.76.29.7, 5001

③

**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
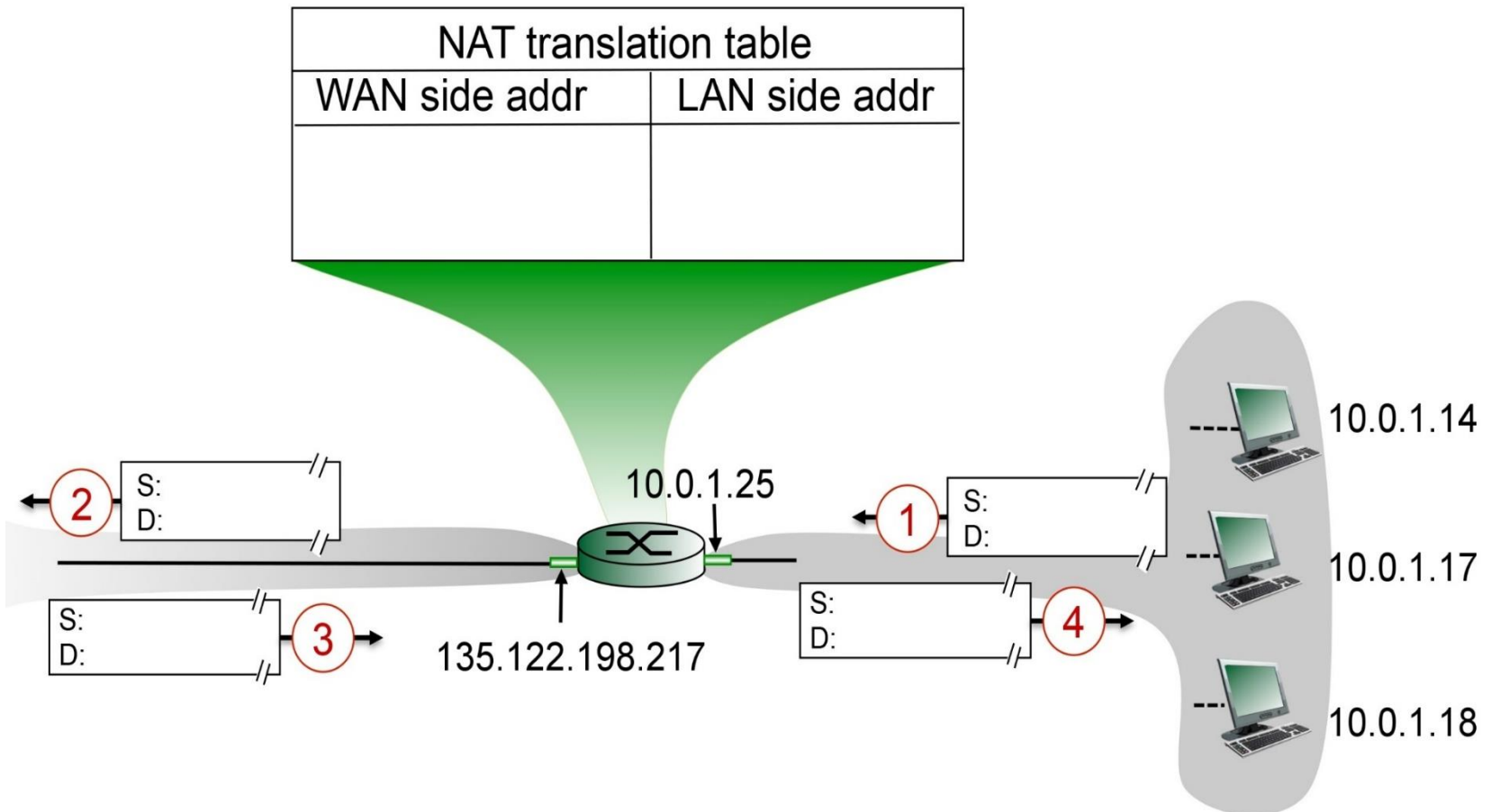
10.0.0.3

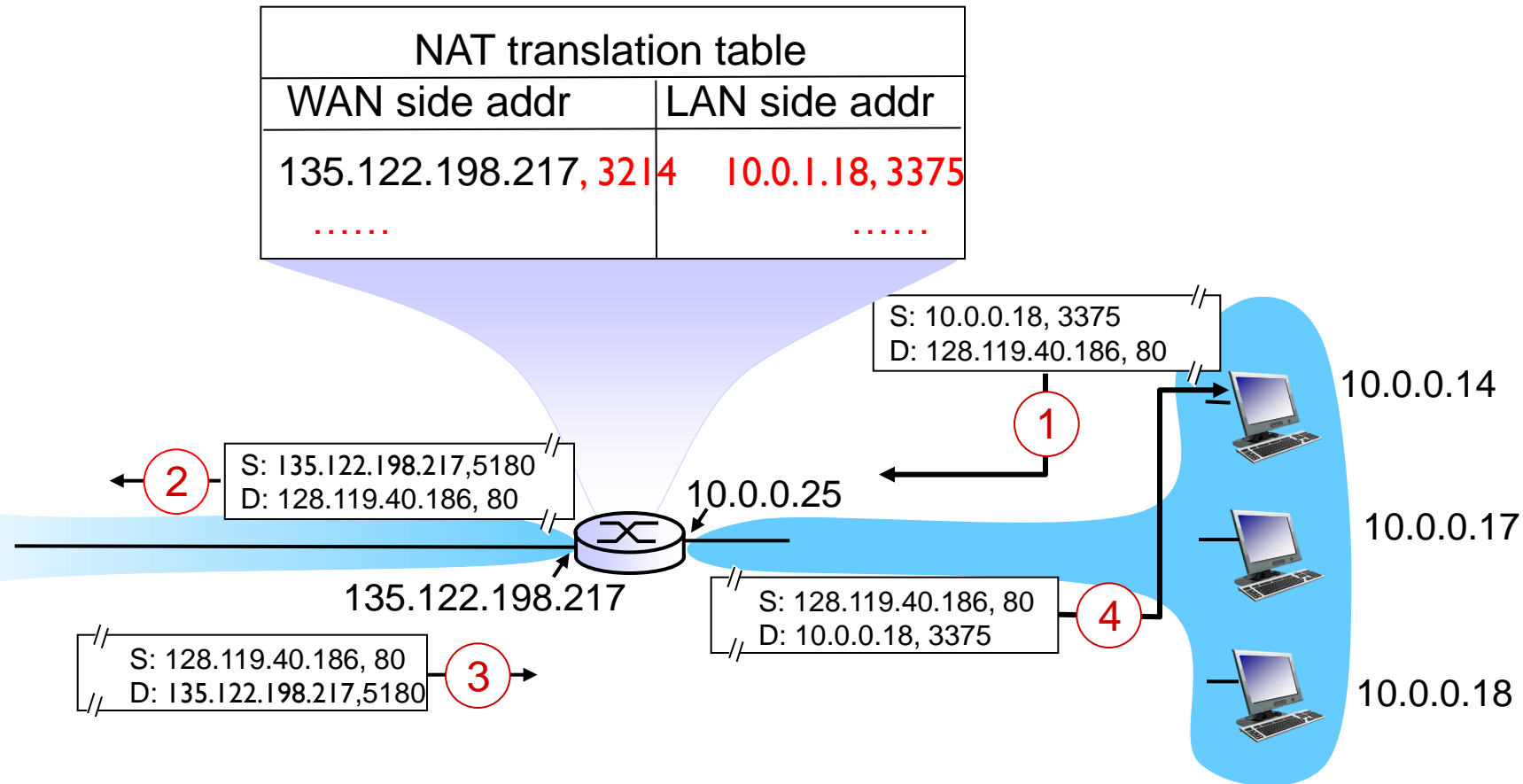# NAT: network address translation

- 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

# NAT: network address translation

- Assume that the computer with IP address 10.0.1.18, port#3375 wants to communicate with a web server, 128.119.40.186, 80.

# NAT: network address translation

| NAT translation table | |
|---|---|
| WAN side addr | LAN side addr |
| 135.122.198.217, 3214 | 10.0.1.18, 3375 |
| …… | …… |

S: 10.0.0.18, 3375
D: 128.119.40.186, 80

1

10.0.0.14

2

S: 135.122.198.217, 5180
D: 128.119.40.186, 80

10.0.0.25

135.122.198.217

10.0.0.17

S: 128.119.40.186, 80
D: 10.0.0.18, 3375

4

3

S: 128.119.40.186, 80
D: 135.122.198.217, 5180

10.0.0.18

# Chapter 4: outline

4.1 Overview of Network layer
- data plane
- control plane

4.2 What's inside a router

<span style="color:red">4.3 IP: Internet Protocol</span>
- datagram format
- fragmentation
- IPv4 addressing
- network address translation

# 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

# 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)