

Network Layer: IP and SDN

Oct 29, 2024

Min Suk Kang

Associate Professor

School of Computing/Graduate School of Information Security



Network layer: “data plane” roadmap

- Network layer: overview
 - data plane
 - control plane
 - What's inside a router
 - input ports, switching, output ports
 - buffer management, scheduling
 - IP: the Internet Protocol
 - datagram format
 - addressing
 - network address translation
 - IPv6
- Generalized Forwarding, SDN
 - Match+action
 - OpenFlow: match+action in action
 - Middleboxes



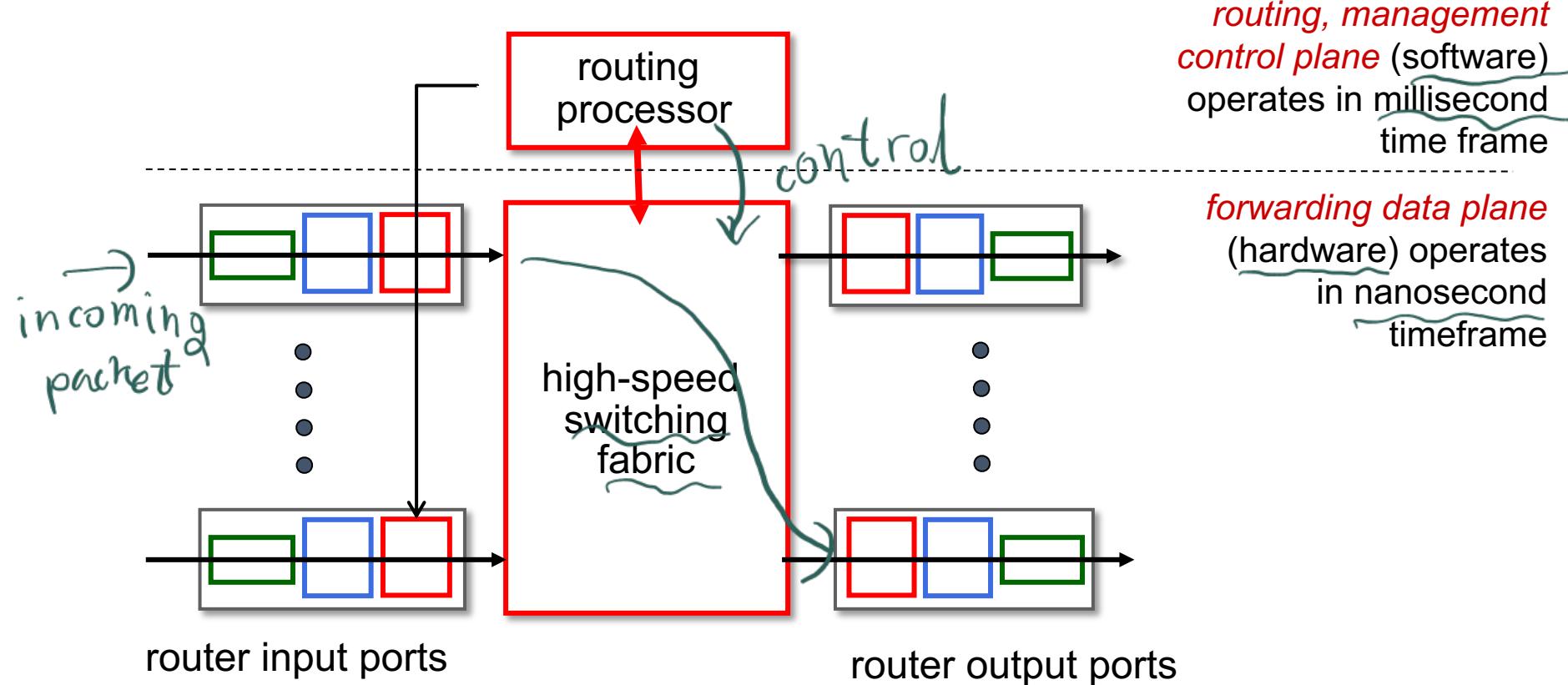
Router architecture overview



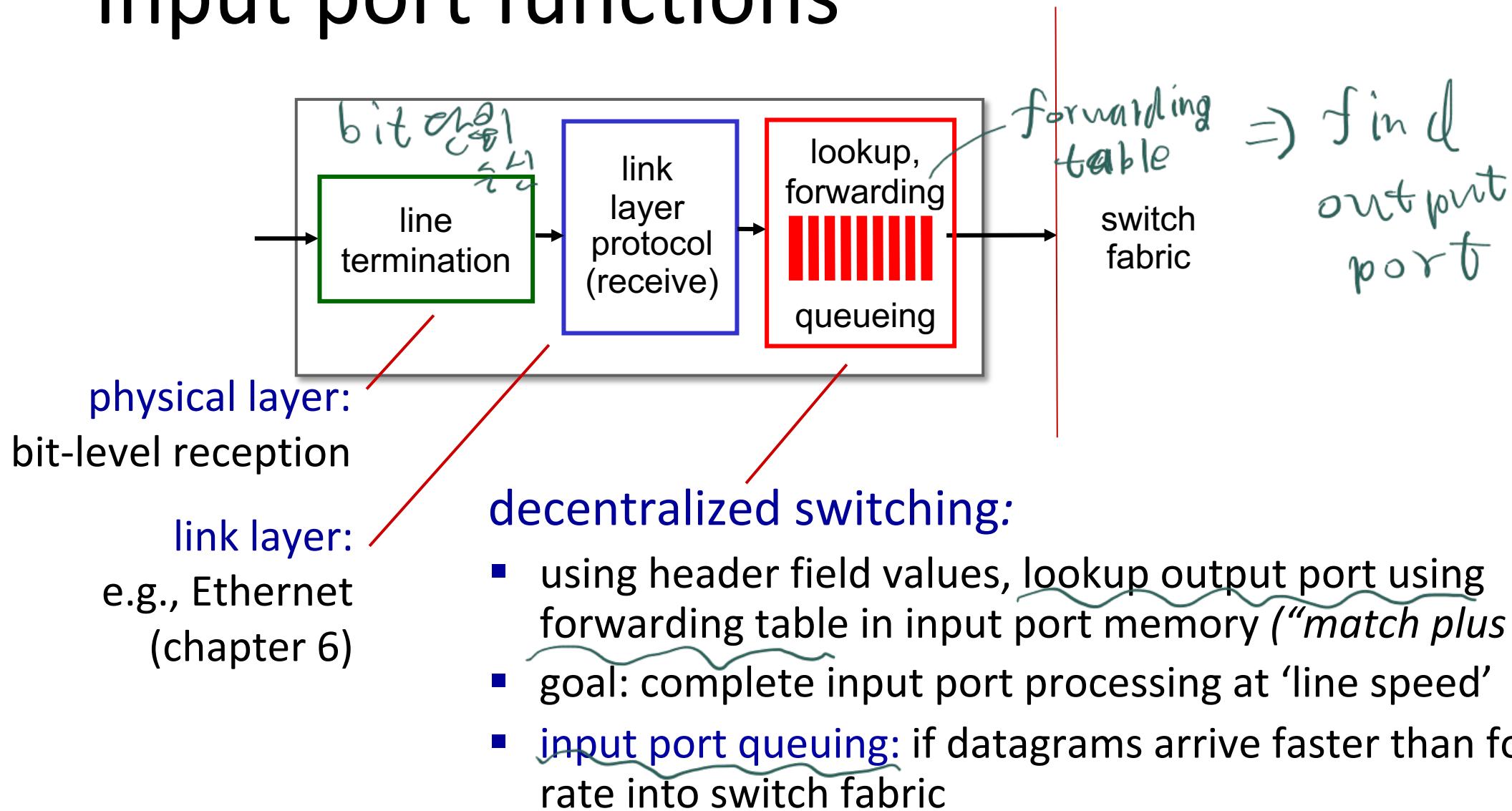
Adobe Stock | #224845489

Router architecture overview

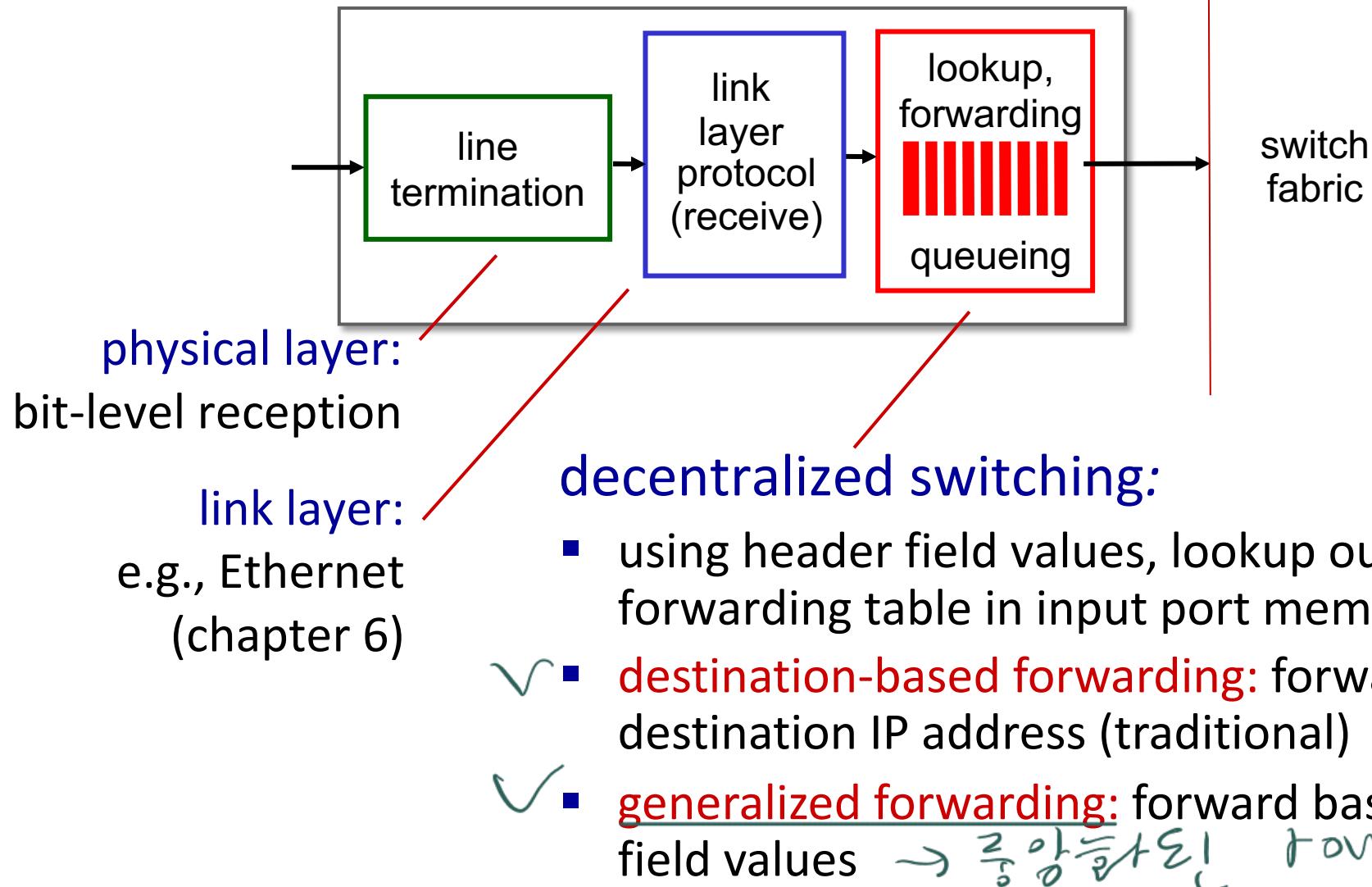
high-level view of generic router architecture:



Input port functions



Input port functions



Destination-based forwarding

bit 91 မြန်။ output link ၇၂၃၁
forwarding table

Destination Address Range	Link Interface
11001000 00010111 00010000 00000000 through 00010111 11111111	range 0 includes 0
11001000 00010111 00010000 00000100 through	range 1 includes 1
11001000 00010111 00010000 00000111 00011000 00000000	range 2 includes 2
11001000 00010111 00011000 11111111	range 3 includes 3
range	
11001000 00010111 00011001 00000000 through	2
11001000 00010111 00011111 11111111	3
otherwise 32 bit address	3

ex) 00010000
00000010
0 or 3?

range 0
range 1
range 2
range 3
separates
range 0
range 1
range 2
range 3
separates
range 0
range 1
range 2
range 3

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

Longest prefix matching

longest prefix match

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:

11001000 00010111 00010110 10100001 which interface?
11001000 00010111 00011000 10101010 which interface?

* (마지막 bit)
로 매칭

C1 가 prefix
(접두사)
선택 = 1

Longest prefix matching

longest prefix match

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 1 00011*** *****	2
otherwise	3

examples:

11001000 00010111 00010110 10100001 which interface?

11001000 00010111 00011000 10101010 which interface?

Longest prefix matching

longest prefix match

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

match!

examples:

11001000	00010111	00010110	10100001	which interface?
11001000	00010111	00011000	10101010	which interface?

Longest prefix matching

longest prefix match

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

match!

examples:

11001000	00010111	00010110	10100001	which interface?
11001000	00010111	00011000	10101010	which interface?

Longest prefix matching

Simple & fast

- 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 \Rightarrow 주소를 주면 한 번에 결과를 찾을 수 있다
 - Cisco Catalyst: ~1M routing table entries in TCAM

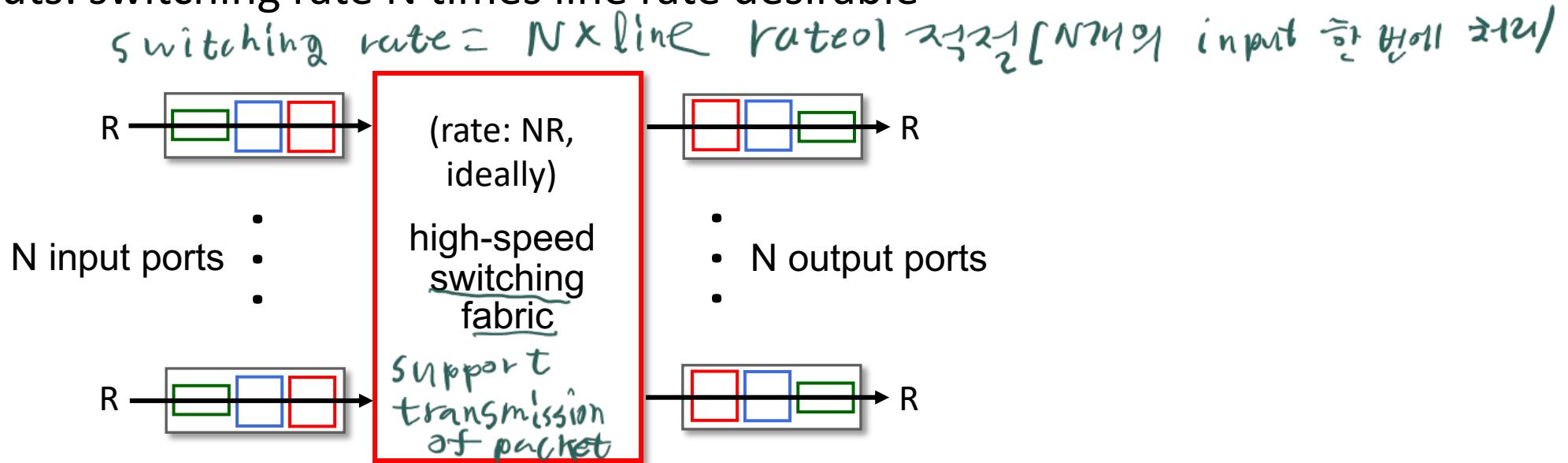
↳ 네트워크
스위치
터미널

3진 - 0, 1, x
0 0 1
0 1 1
1 0 1
1 1 1
기본주소
증가주소
전체주소
검색어 탐색하여
일관성 여부에 반응
(속도↑)

Switching fabrics

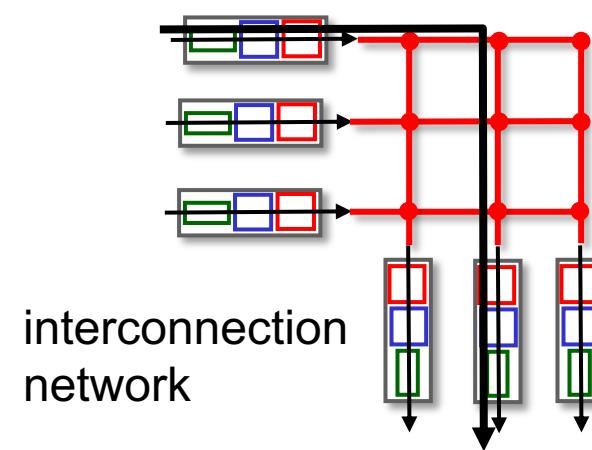
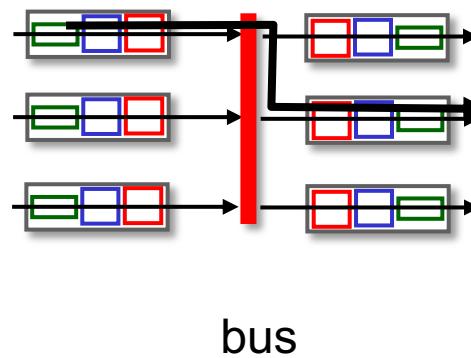
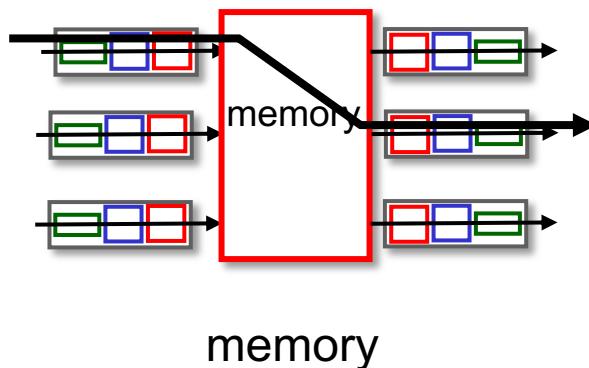
2121-2

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



Switching fabrics

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

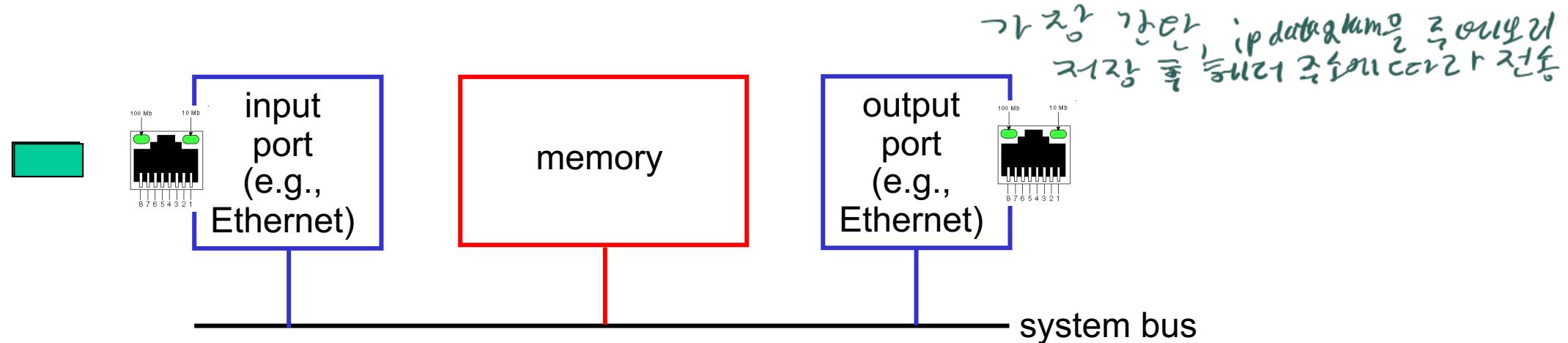


Switching via memory

skip
↓

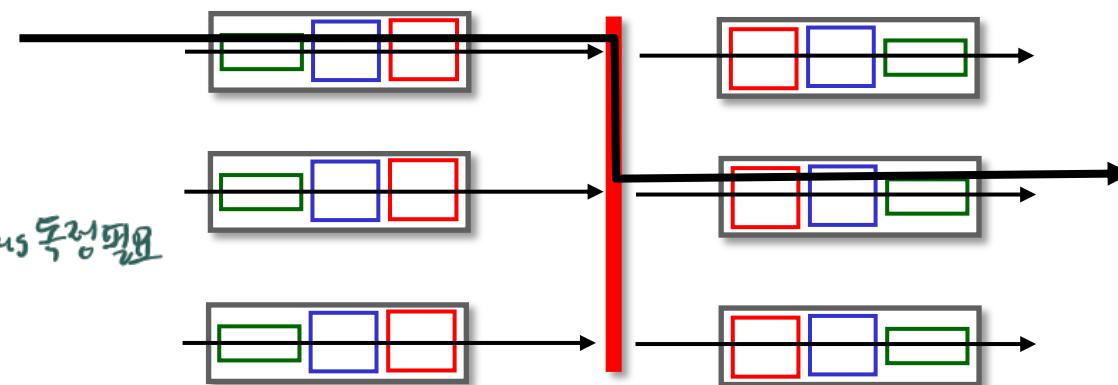
first generation routers:

- 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)



Switching via a bus

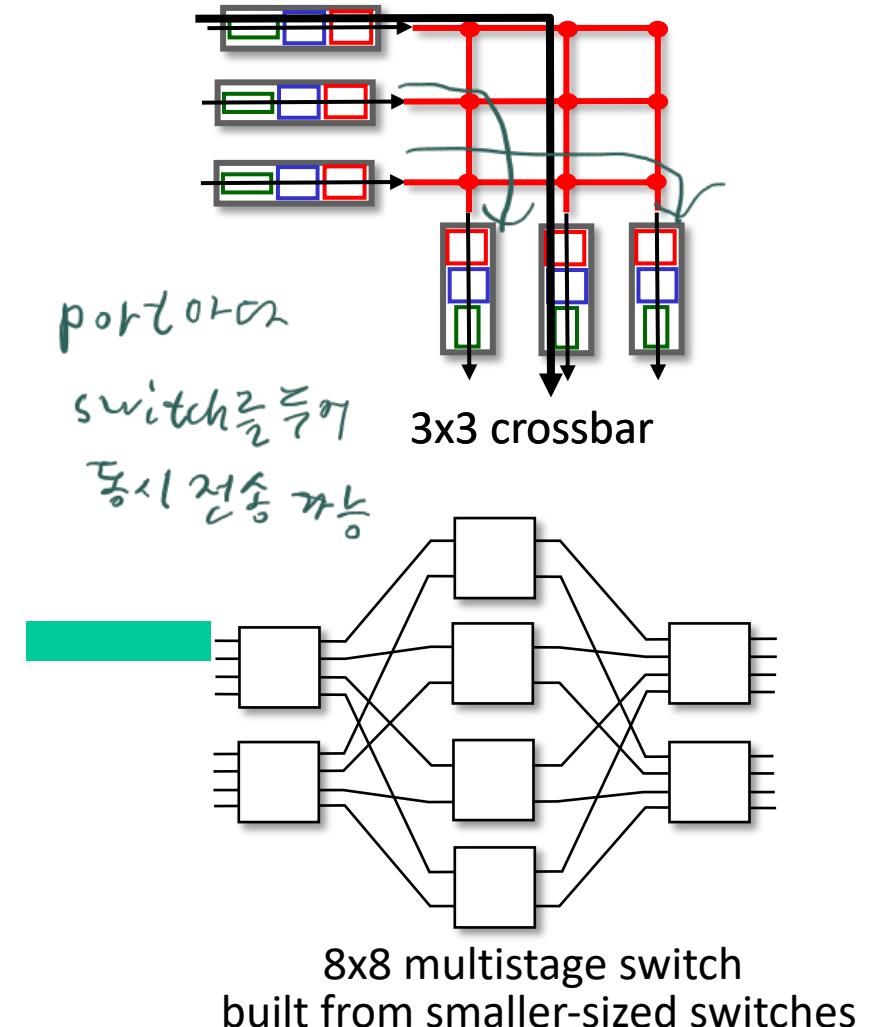
- datagram from input port memory to output port memory via a shared bus
- *bus contention*: switching speed limited by bus bandwidth
- 32 Gbps bus, Cisco 5600: sufficient speed for access routers



부록 2
'input & output
port buffer 사용'으로,
전송시 bus斗争 필요

Switching via interconnection network

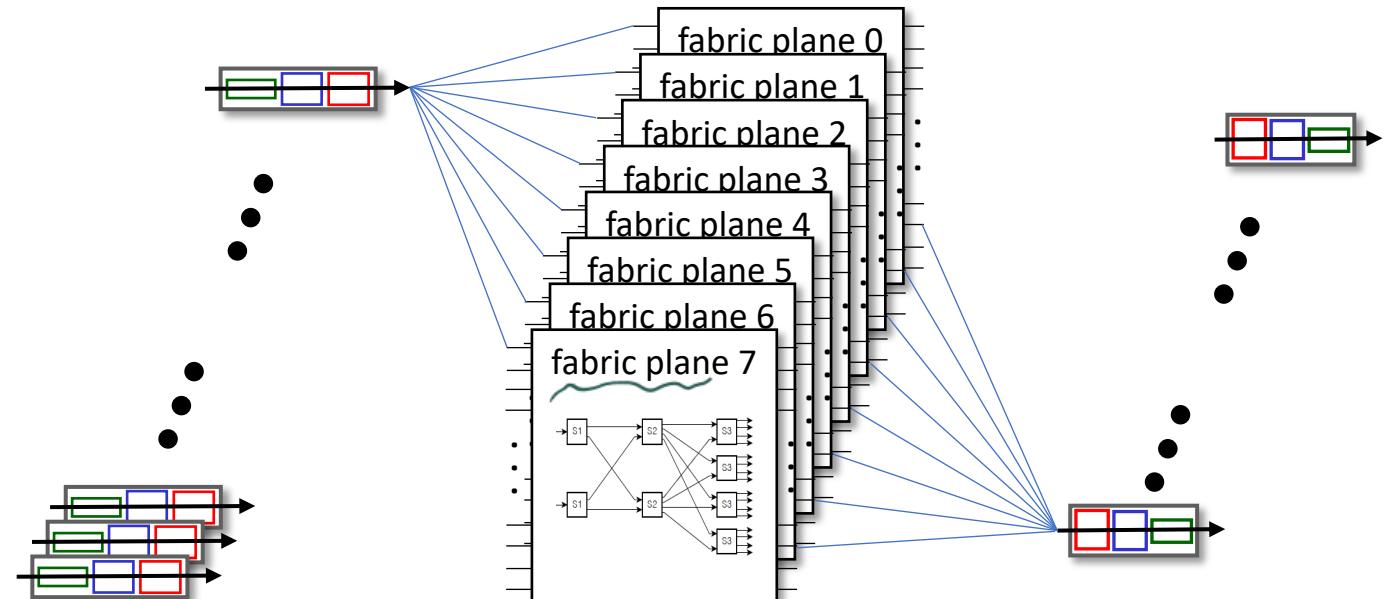
- Crossbar, Clos networks, other interconnection nets initially developed to connect processors in multiprocessor
- multistage switch: $n \times n$ switch from multiple stages of smaller switches
- exploiting parallelism:
• fragment datagram into fixed length cells on entry
• switch cells through the fabric, reassemble datagram at exit



Switching via interconnection network

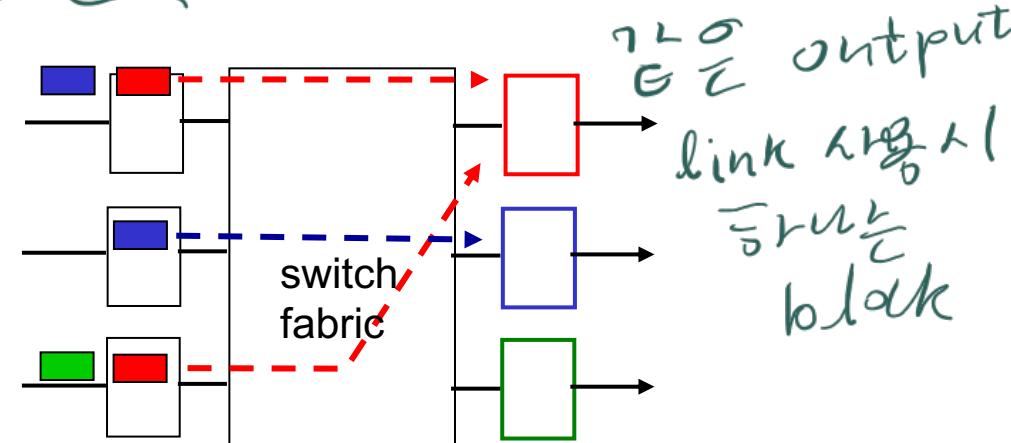
4
skip

- scaling, using multiple switching “planes” in parallel:
 - speedup, scaleup via parallelism
- Cisco CRS router:
 - basic unit: 8 switching planes
 - each plane: 3-stage interconnection network
 - up to 100's Tbps switching capacity

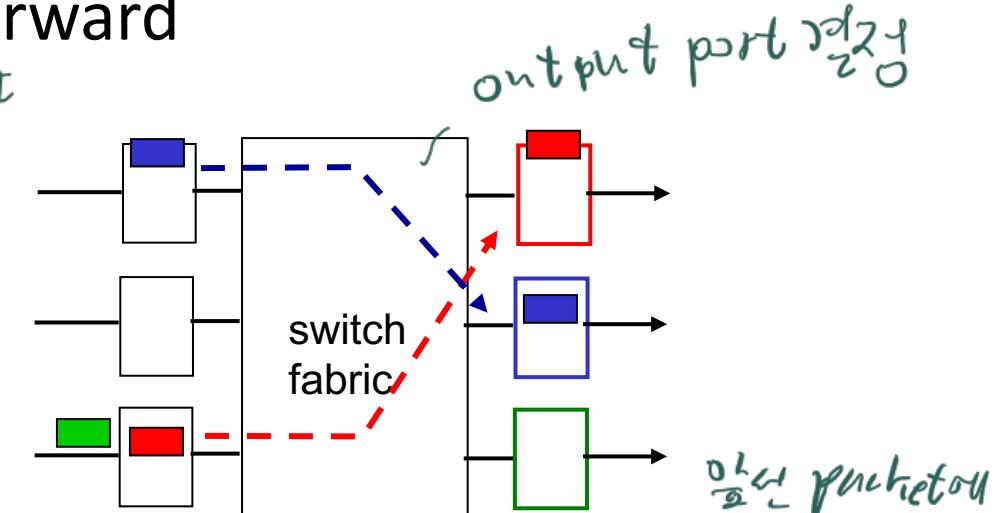


Input port queuing

- If switch 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

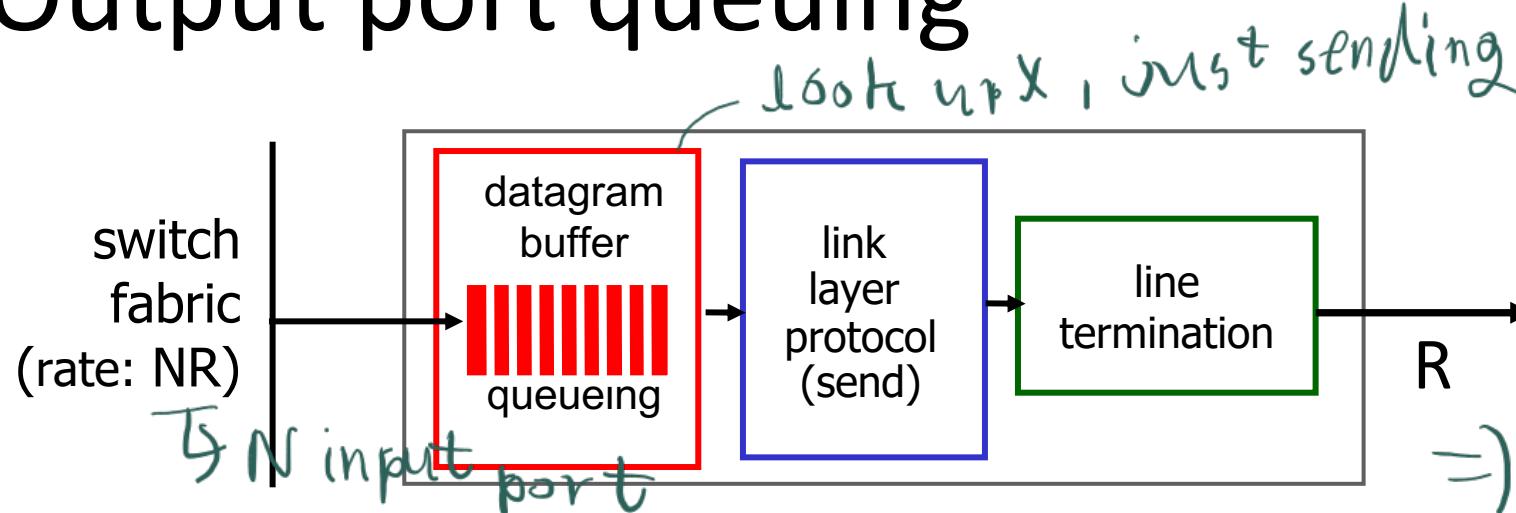


output port contention: only one red datagram can be transferred. lower red packet is **blocked**



one packet time later: green packet experiences HOL blocking

Output port queuing



This is a really important slide

- **Buffering** required when datagrams arrive from fabric faster than link transmission rate. **Drop policy:** which datagrams to drop if no free buffers?

buffer 대기열 drop 발생

tail drop
(\Rightarrow packet drop)

middle drop
(\Rightarrow middle queue drop)

:

- **Scheduling discipline** chooses among queued datagrams for transmission

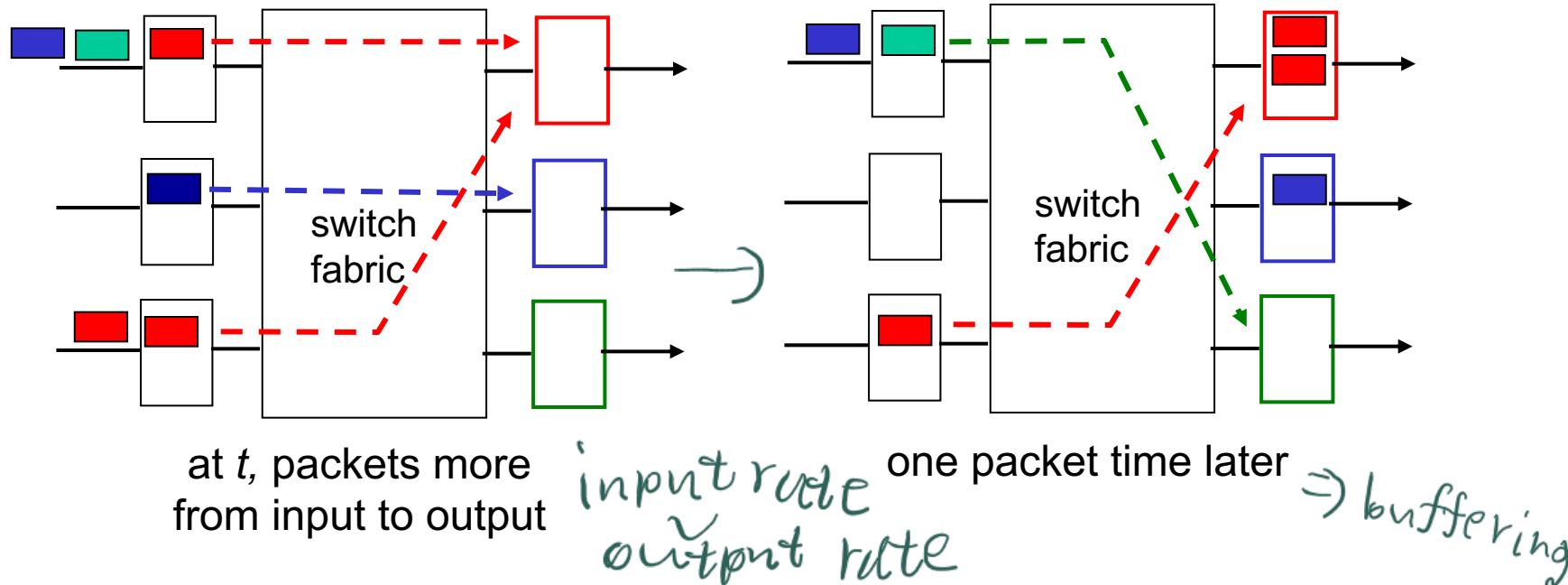
\rightarrow dequeue 순서 정하기

$CNR > R$)

Datagrams can be lost due to congestion, lack of buffers

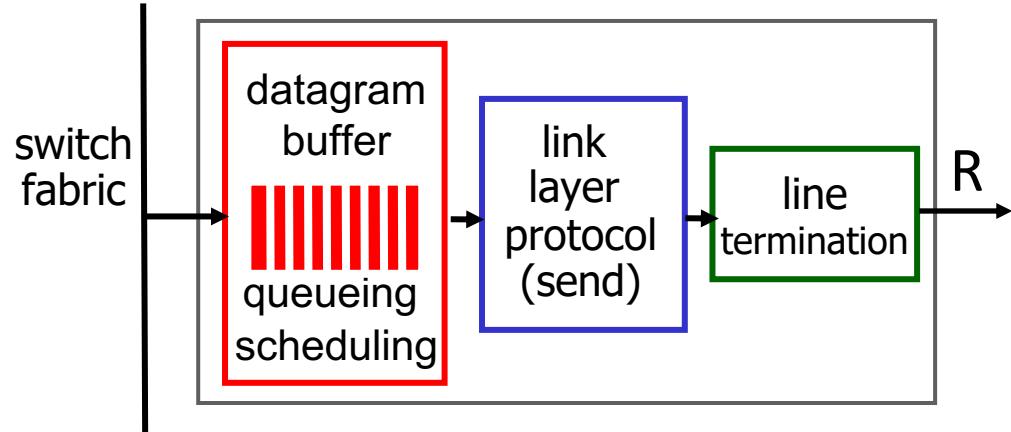
Priority scheduling – who gets best performance, network neutrality

Output port queuing

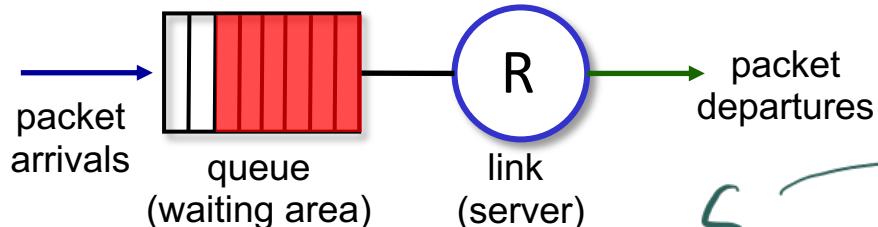


- buffering when arrival rate via switch exceeds output line speed
- *queueing (delay) and loss due to output port buffer overflow!*

Buffer Management



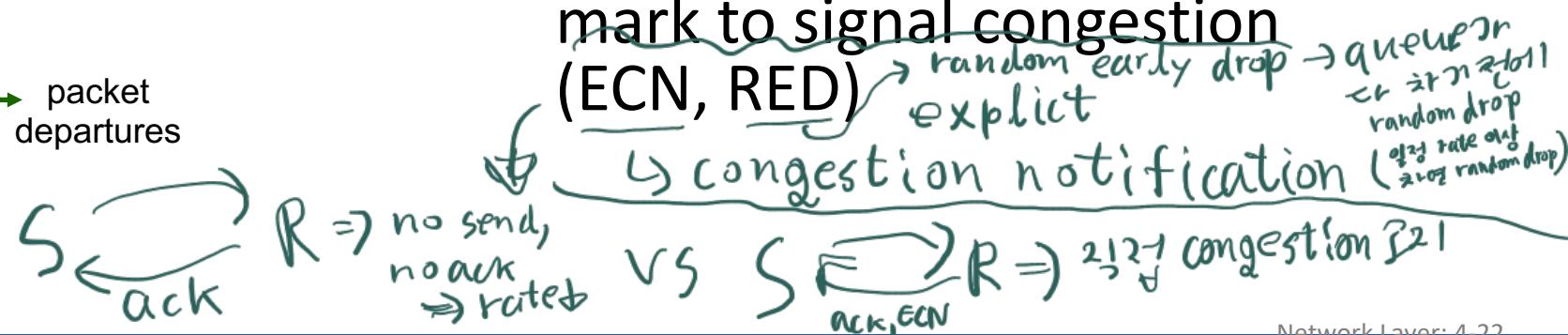
Abstraction: queue



buffer management:

- **drop:** which packet to add, drop when buffers are full
 - **tail drop:** drop arriving packet → 마지막 도착
 - **priority:** drop/remove on priority basis → 우선순위 처리

- **marking:** which packets to mark to signal congestion (ECN, RED)



Packet Scheduling: FCFS

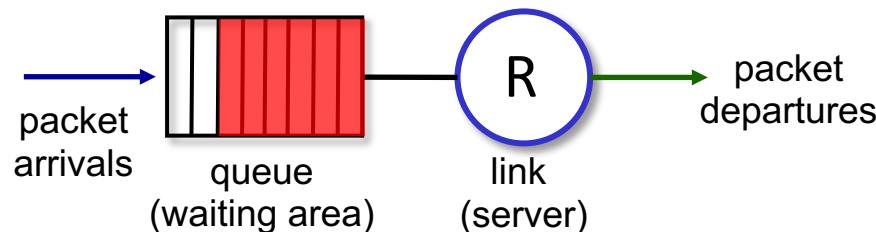
packet scheduling: deciding which packet to send next on link

- first come, first served =FCFS
- priority
- round robin
- weighted fair queueing

FCFS: packets transmitted in order of arrival to output port → 먼저 들어온 순서대로

- also known as: First-in-first-out (FIFO)
- real world examples?

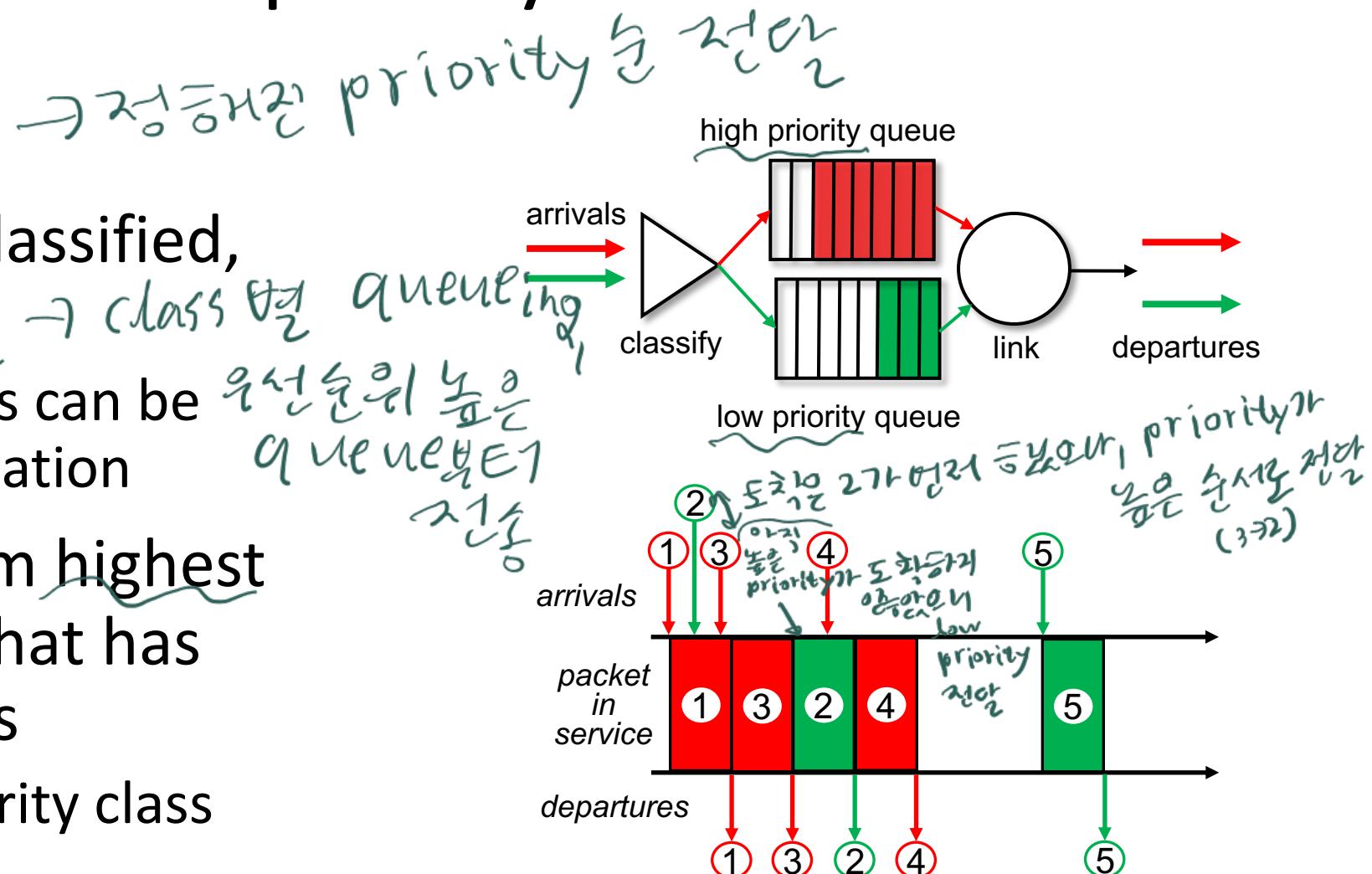
Abstraction: queue



Scheduling policies: priority

Priority scheduling:

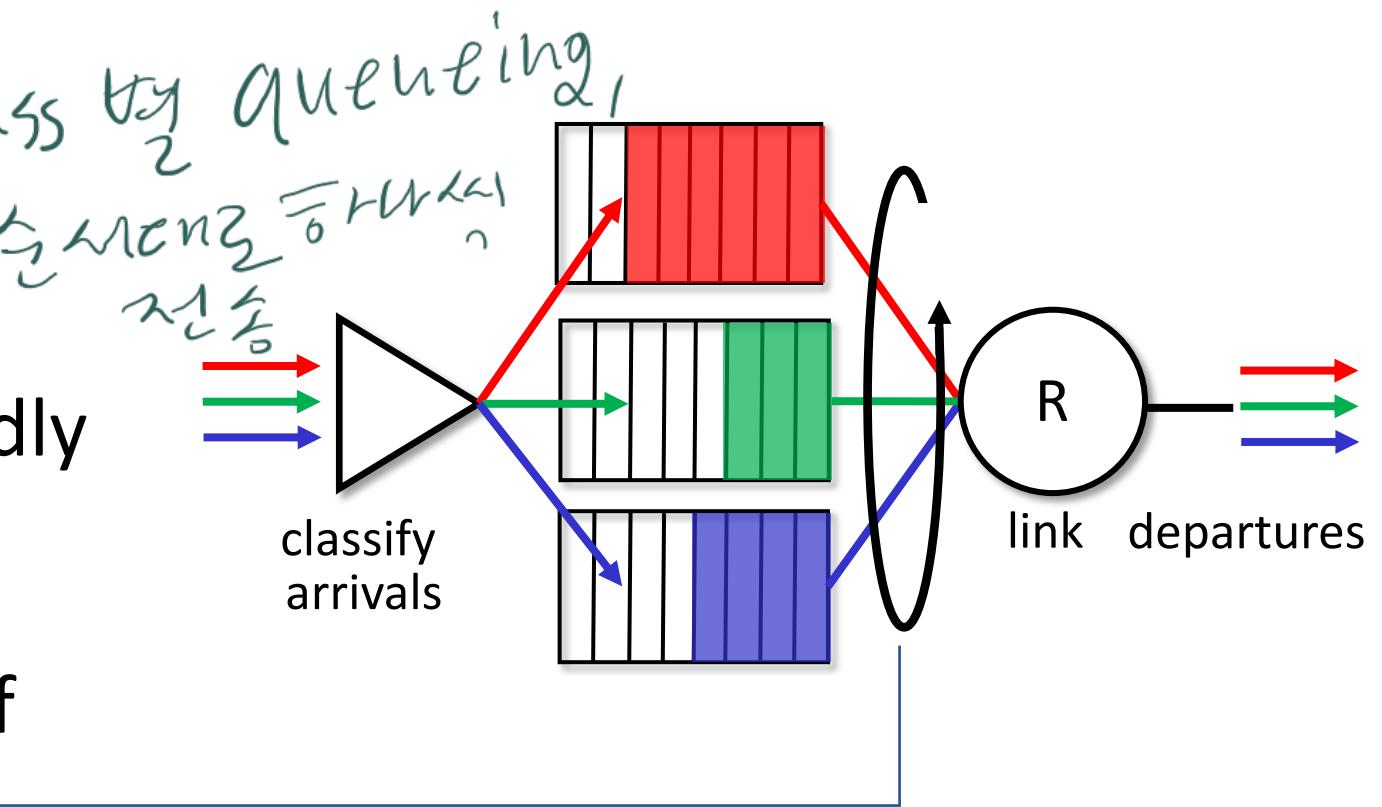
- arriving traffic classified, queued by class
 - any header fields can be used for classification
- send packet from highest priority queue that has buffered packets
 - FCFS within priority class



Scheduling policies: round robin

Round Robin (RR) scheduling: \rightarrow fairness

- arriving traffic classified, queued by class \rightarrow class by queuing,
 - any header fields can be used for classification
- server cyclically, repeatedly scans class queues, sending one complete packet from each class (if available) in turn



r, g, b 순서로 반복하여 흑자내수작업

Scheduling policies: weighted fair queueing

Weighted Fair Queuing (WFQ):

→ RR + priority
(weight)

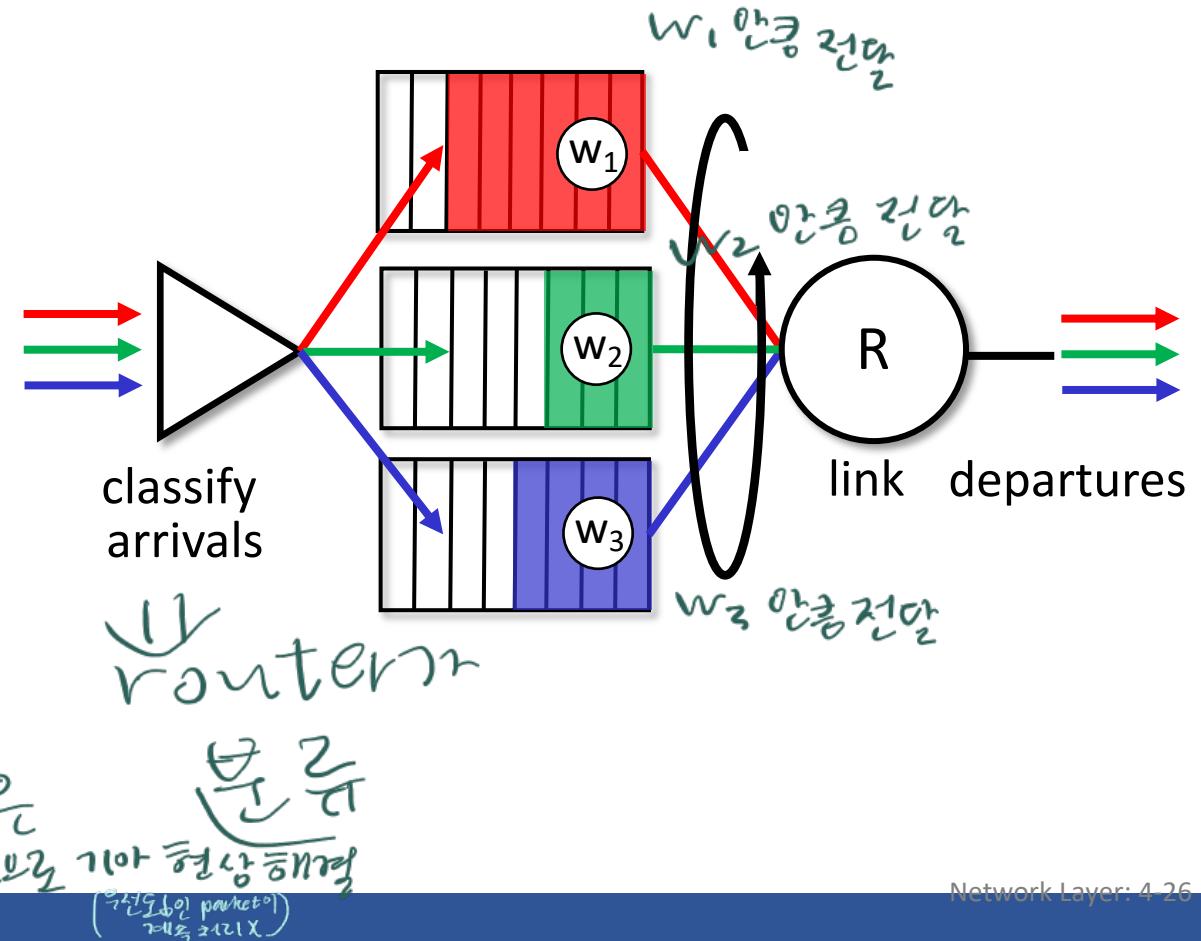
- generalized Round Robin
- each class, i , has weight, w_i ,
and gets weighted amount
of service in each cycle:

$$\frac{w_i}{\sum_j w_j}$$

각 cycle에서
가장 많은 큐 전달
 $\rightarrow w \rightarrow 전송량 \uparrow$

- minimum bandwidth
guarantee (per-traffic-class)

전체 대역폭 중 일정 비율은
보장되어야 함



Sidebar: Network Neutrality

What is network neutrality?

- *technical*: how an ISP should share/allocation its resources
 - packet scheduling, buffer management are the *mechanisms*
- *social, economic principles*
 - protecting free speech
 - encouraging innovation, competition
- enforced *legal* rules and policies

Different countries have different “takes” on network neutrality

Sidebar: Network Neutrality

2015 US FCC *Order on Protecting and Promoting an Open Internet*: three “clear, bright line” rules:

- no blocking ... “shall not block lawful content, applications, services, or non-harmful devices, subject to reasonable network management.”
- no throttling ... “shall not impair or degrade lawful Internet traffic on the basis of Internet content, application, or service, or use of a non-harmful device, subject to reasonable network management.”
- no paid prioritization. ... “shall not engage in paid prioritization”

2017 roll back \Rightarrow rules being weakened

ISP: telecommunications or information service?

Is an ISP a “telecommunications service” or an “information service” provider?

internet & cellular network

- the answer *really* matters from a regulatory standpoint! is “very” different

(영어성어의 다른 뜻과 차이)

US Telecommunication Act of 1934 and 1996:

- **Title II:** imposes “common carrier duties” on **telecommunications services**: reasonable rates, non-discrimination and requires regulation
- **Title I:** applies to **information services**:
 - no common carrier duties (not regulated) → 규제가 없음
 - but grants FCC authority “... as may be necessary in the execution of its functions” → 필요한 경우 FCC(미국항통신위원회)가 할 수 있음

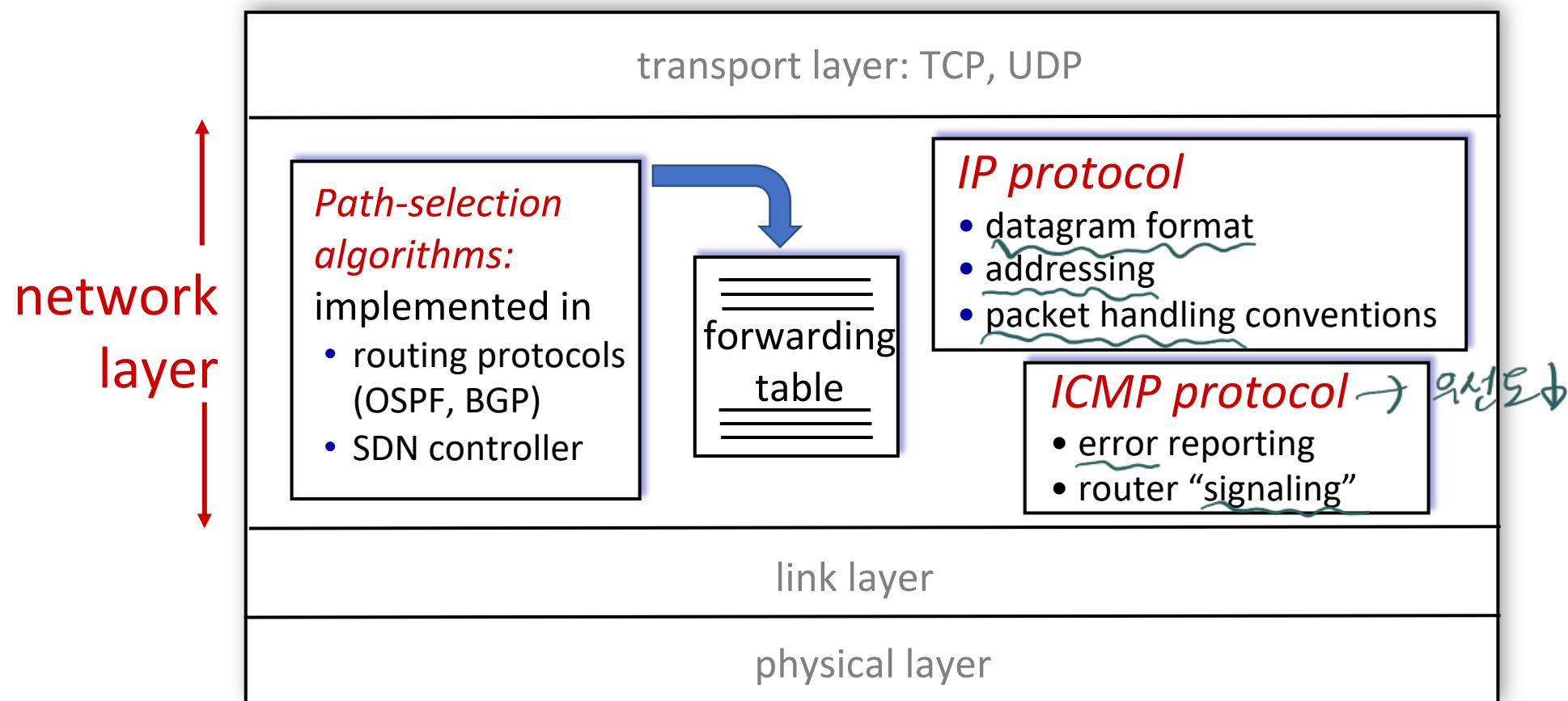
Network layer: “data plane” roadmap

- Network layer: overview
 - data plane
 - control plane
- What's inside a router
 - input ports, switching, output ports
 - buffer management, scheduling
- IP: the Internet Protocol
 - datagram format
 - addressing
 - network address translation
 - IPv6
- Generalized Forwarding, SDN
 - match+action
 - OpenFlow: match+action in action
- Middleboxes



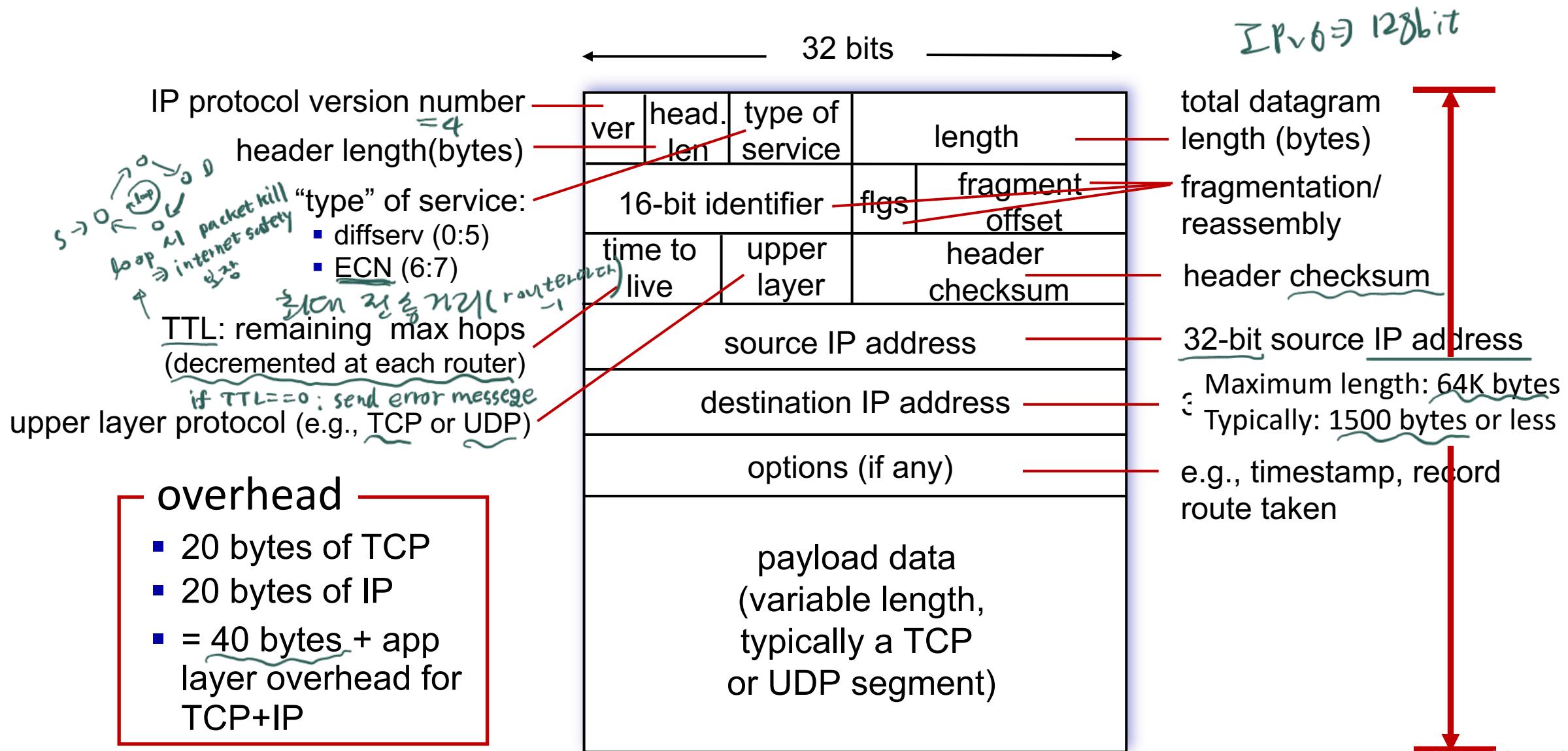
Network Layer: Internet

host, router network layer functions:



IP Datagram format

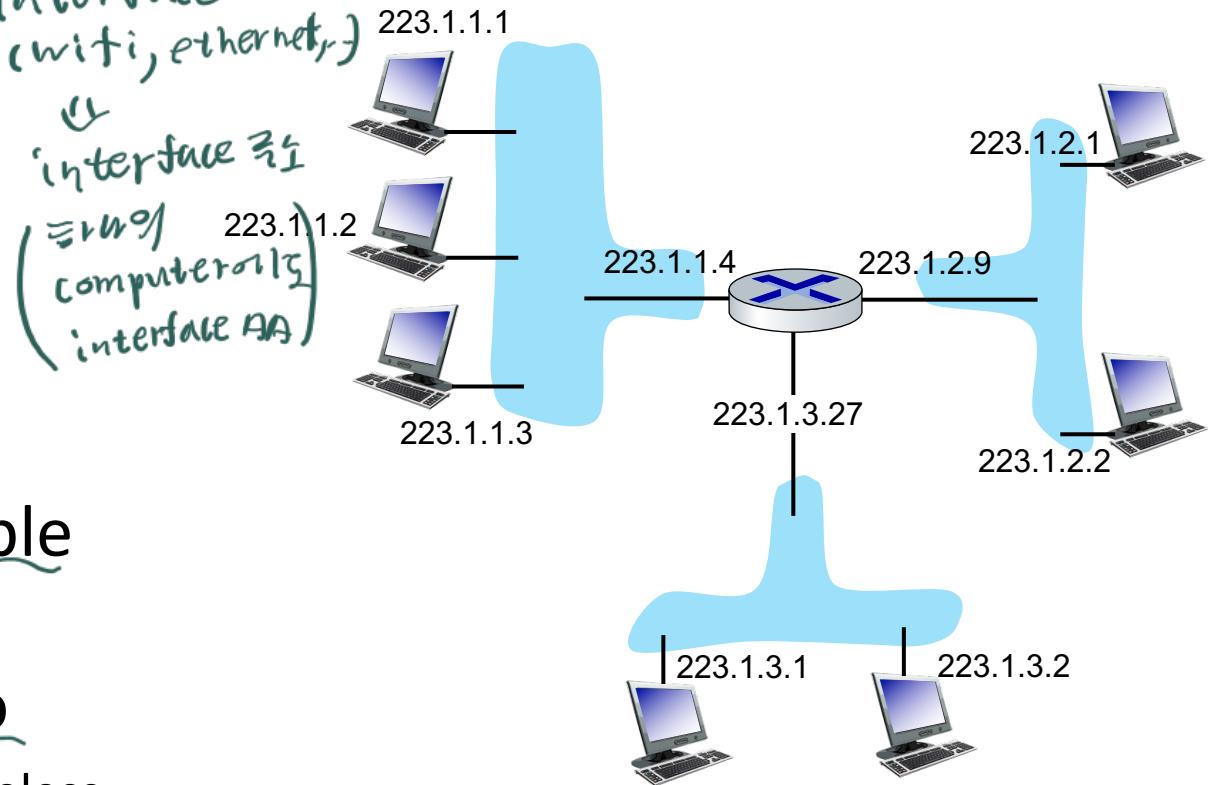
→ version 4 (IPv4)
↔ 32 bit



IP addressing: introduction

*address of "one" interface
(wifii, ethernetr)*

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



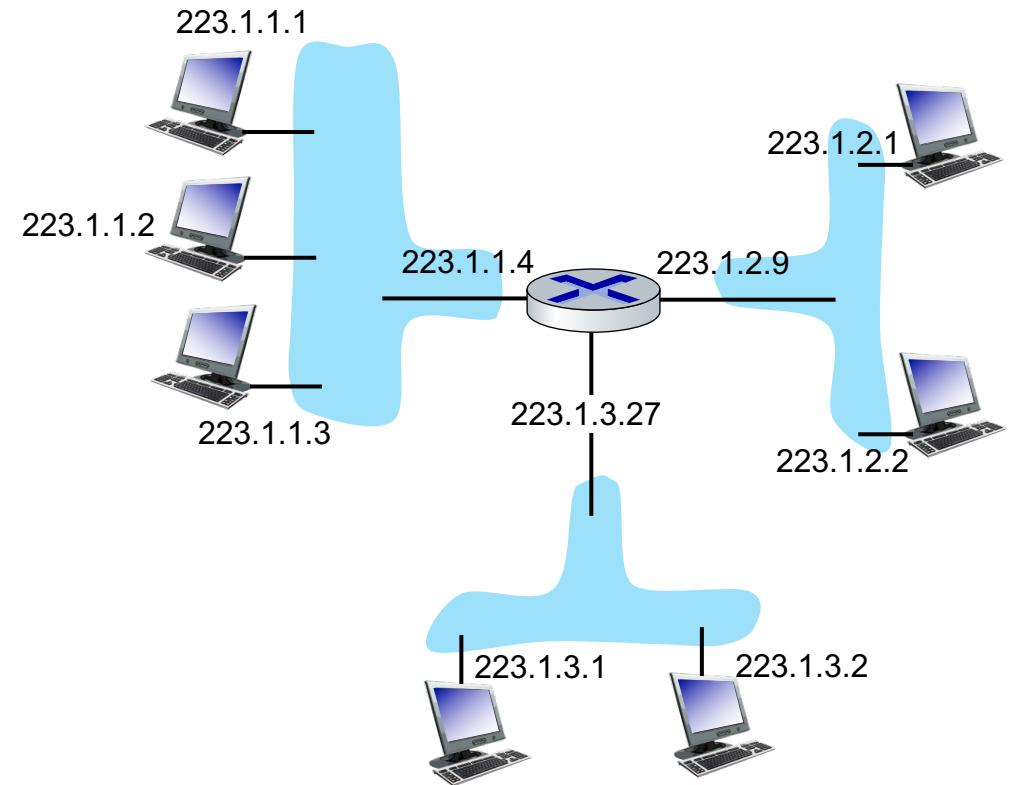
dotted-decimal IP address notation:

223.1.1.1 = $\begin{array}{cccc} 11011111 & 00000001 & 00000001 & 00000001 \end{array}$

223 1 1 1

IP addressing: introduction

- **IP address:** 32-bit identifier associated with each host or 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)



dotted-decimal IP address notation:

223.1.1.1 = 11011111 00000001 00000001 00000001

223 1 1 1

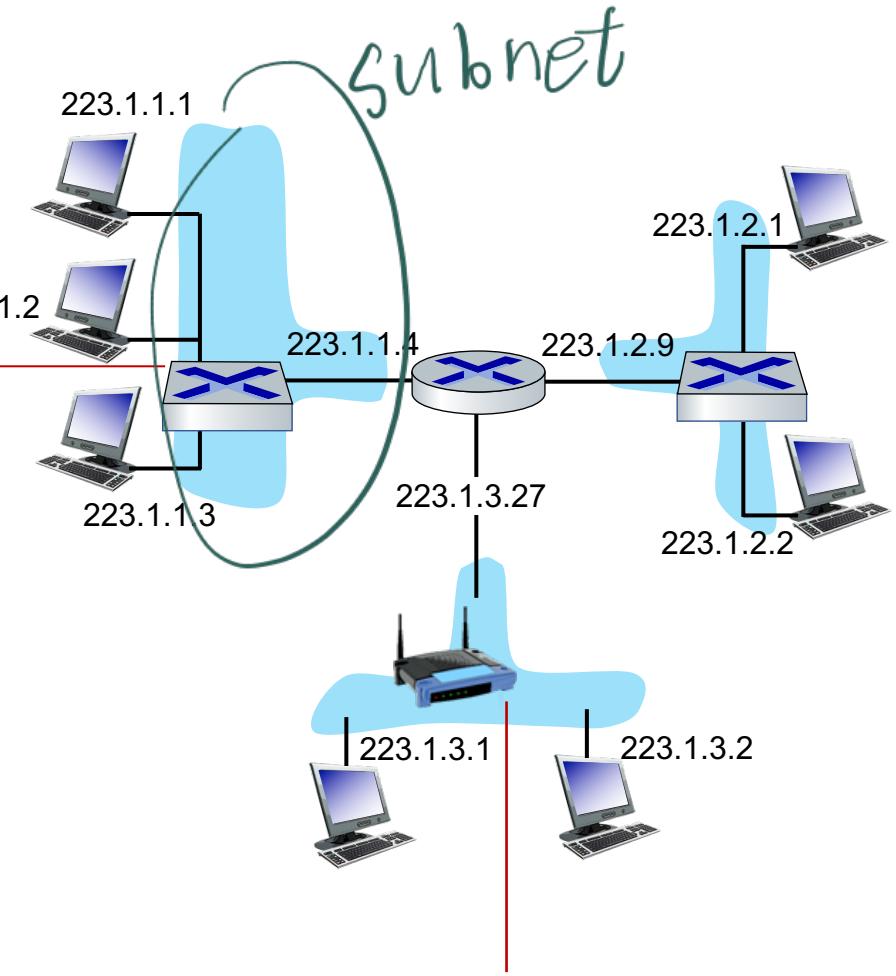
IP addressing: introduction

Q: how are interfaces actually connected?

A: we'll learn about that in chapters 6, 7

For now: don't need to worry about how one interface is connected to another (with no intervening router)

A: wired Ethernet interfaces connected by Ethernet switches



A: wireless WiFi interfaces connected by WiFi base station

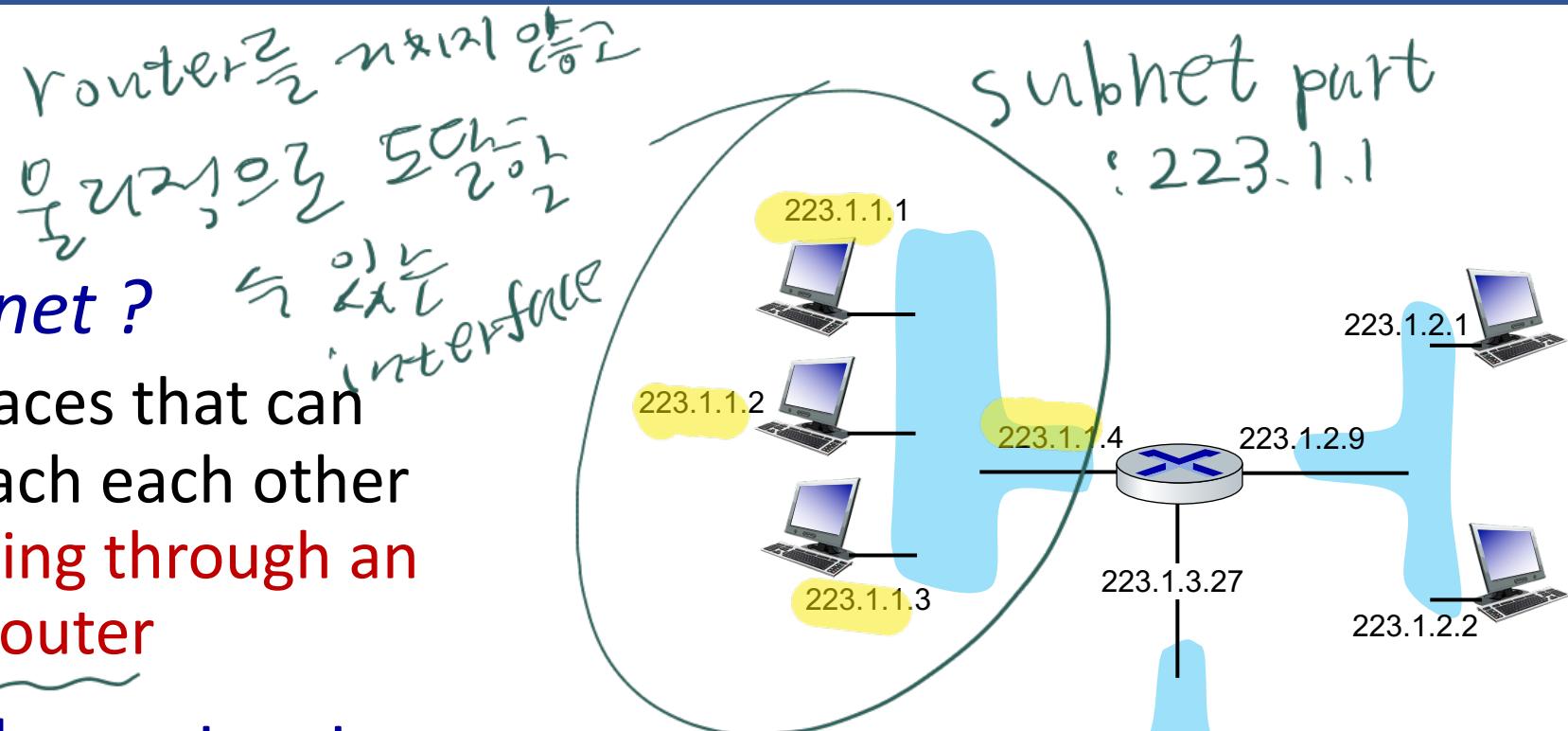
Subnets

■ 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

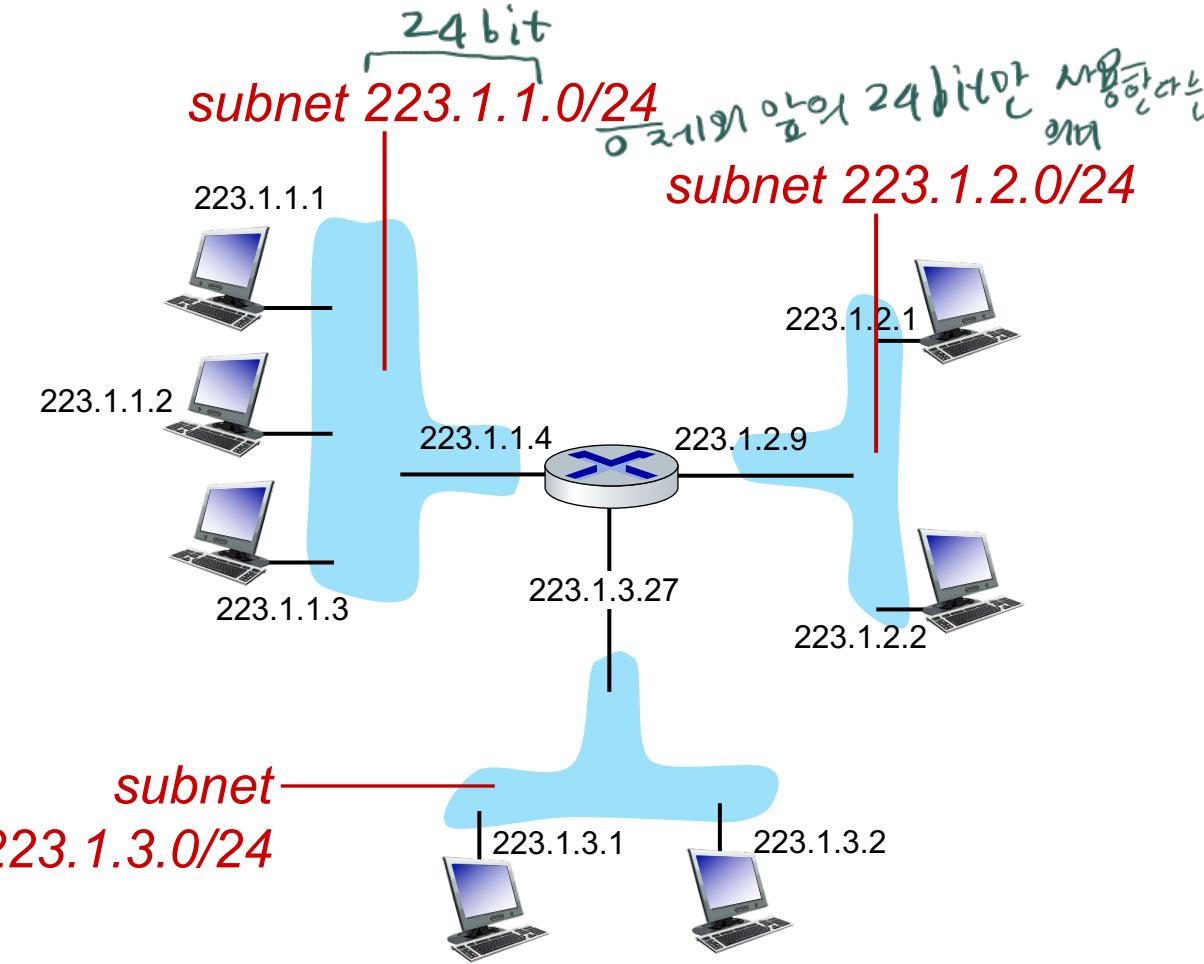
같은 subnet
인가요
devices
subnet part
인가요

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*

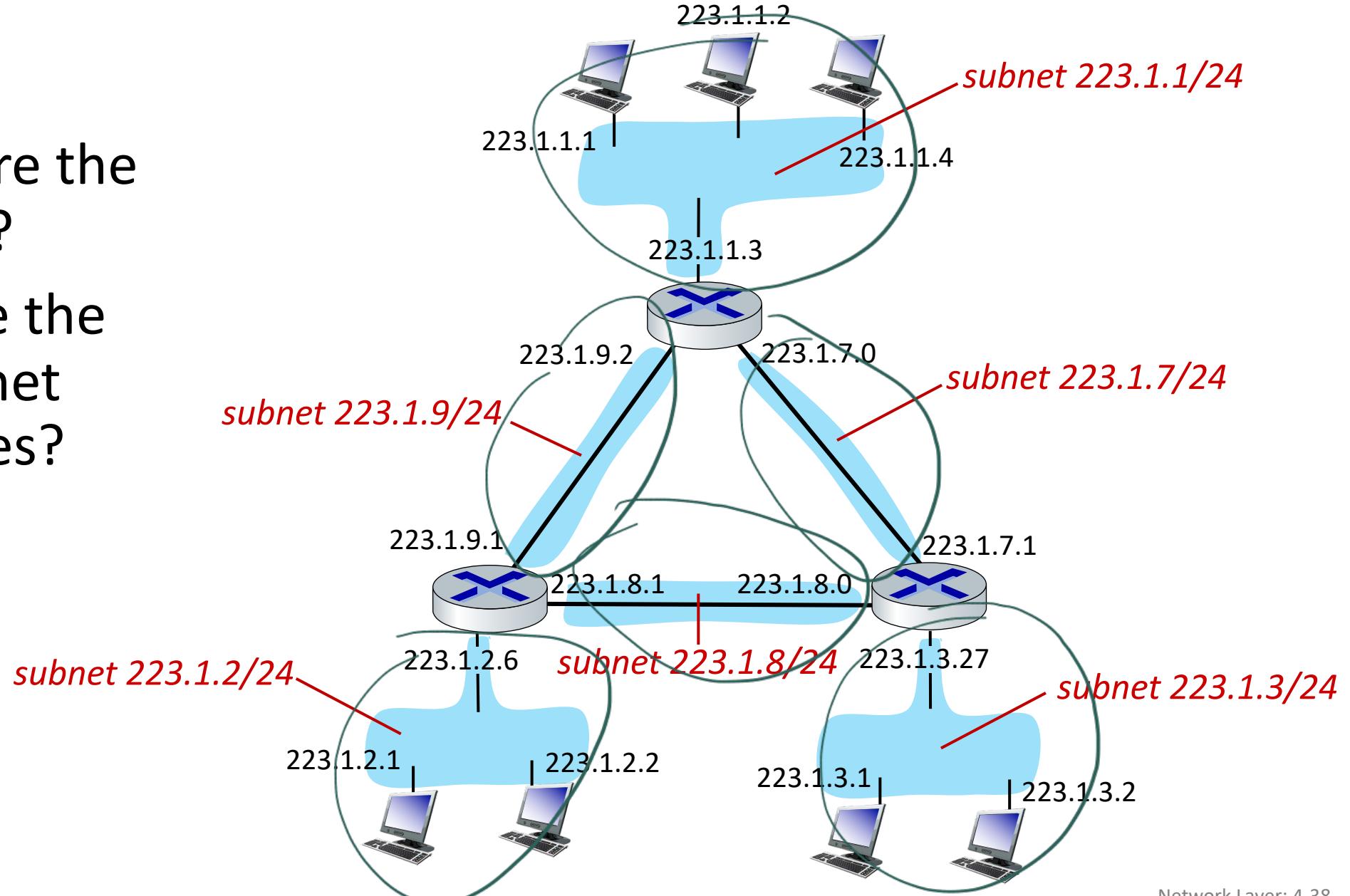
↳ host/router 한곳에
interface를 뚫고 네트워크
구성되는 \Rightarrow subnet



subnet mask: /24
(high-order 24 bits: subnet part of IP address)

Subnets

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



IP: address + identity

First name

Last name

Address

Apartment, suite, etc.

City

State/province

Country

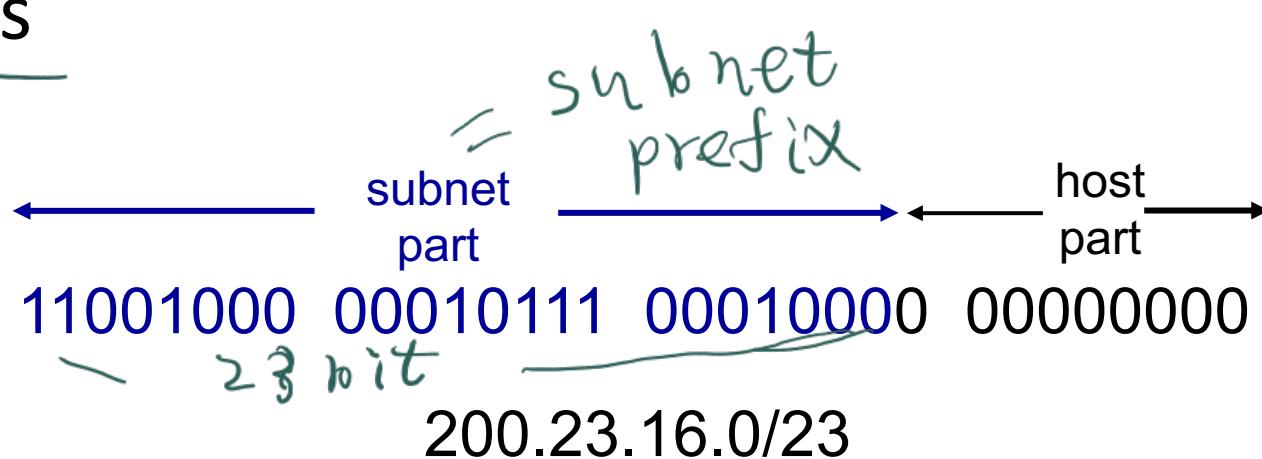
ZIP/postal code

IP addressing: CIDR

subnet 域名

CIDR: Classless InterDomain Routing (pronounced “cider”)

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



IP addresses: how to get one?

That's actually **two** questions:

$143.248.*.*$
= $143.248.0.0/16$

1. Q: How does a *host* get IP address within its *network* (host part of address)?
2. Q: How does a *network* get IP address for *itself* (network part of address)

How does *host* get IP address?

- hard-coded by sysadmin in config file (e.g., /etc/rc.config in UNIX)
- DHCP: Dynamic Host Configuration Protocol: dynamically get address from server → 호스트가 서버에게 IP 할당
 - “plug-and-play”

DHCP: Dynamic Host Configuration Protocol

skip (wiki)
↓

설명

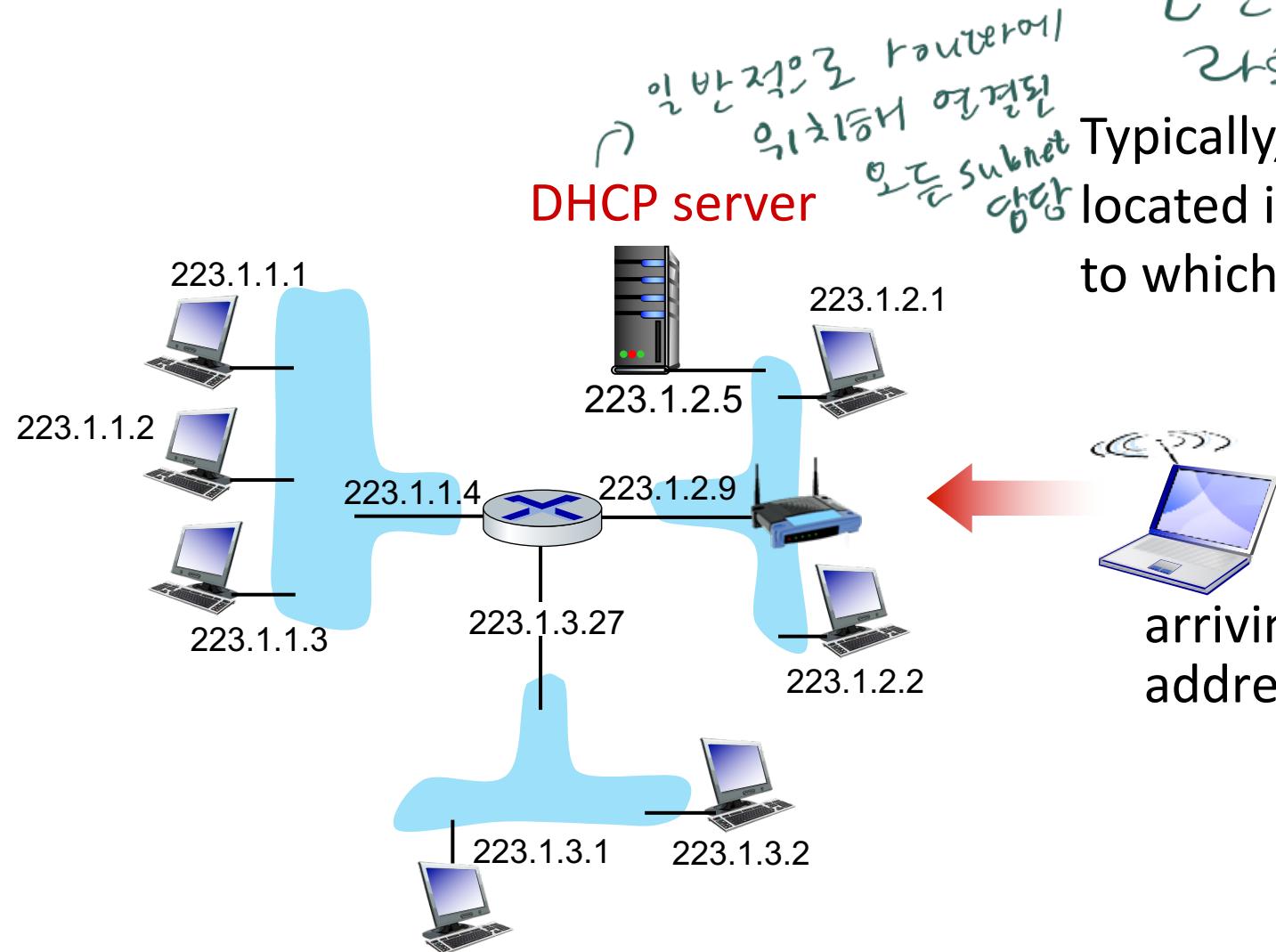
goal: host *dynamically* obtains IP address from network server when it "joins" network → host가 네트워크에 접속시 ip 자동 할당

- can renew its lease on address in use
- allows reuse of addresses (only hold address while connected/on)
- support for mobile users who join/leave network

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

DHCP client-server scenario



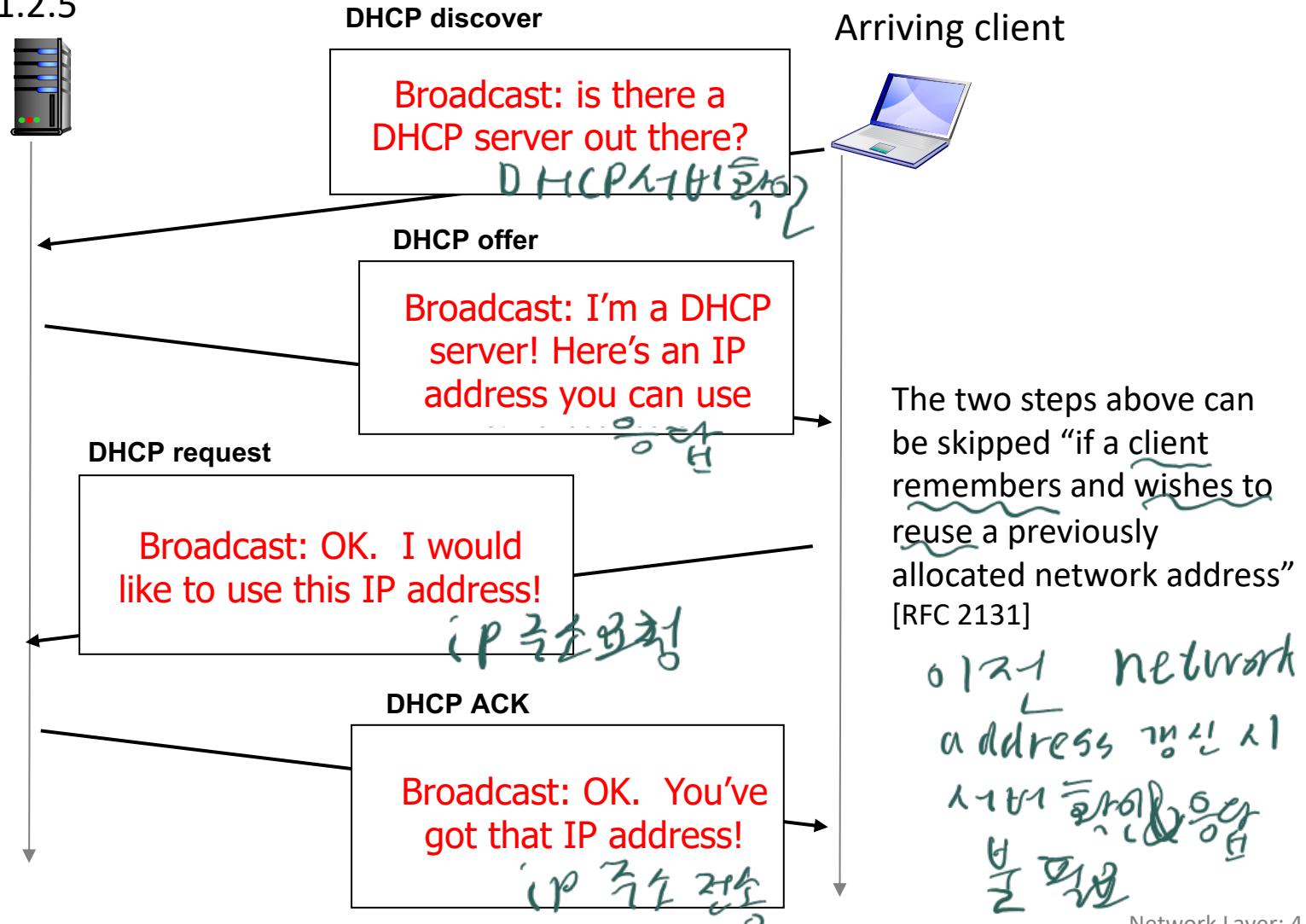
일반적으로, DHCP서버는
라우터에 위치한 모든 subnet
를 지원합니다.

Typically, DHCP server will be co-located in router, serving all subnets to which router is attached

arriving **DHCP client** needs address in this network

DHCP client-server scenario

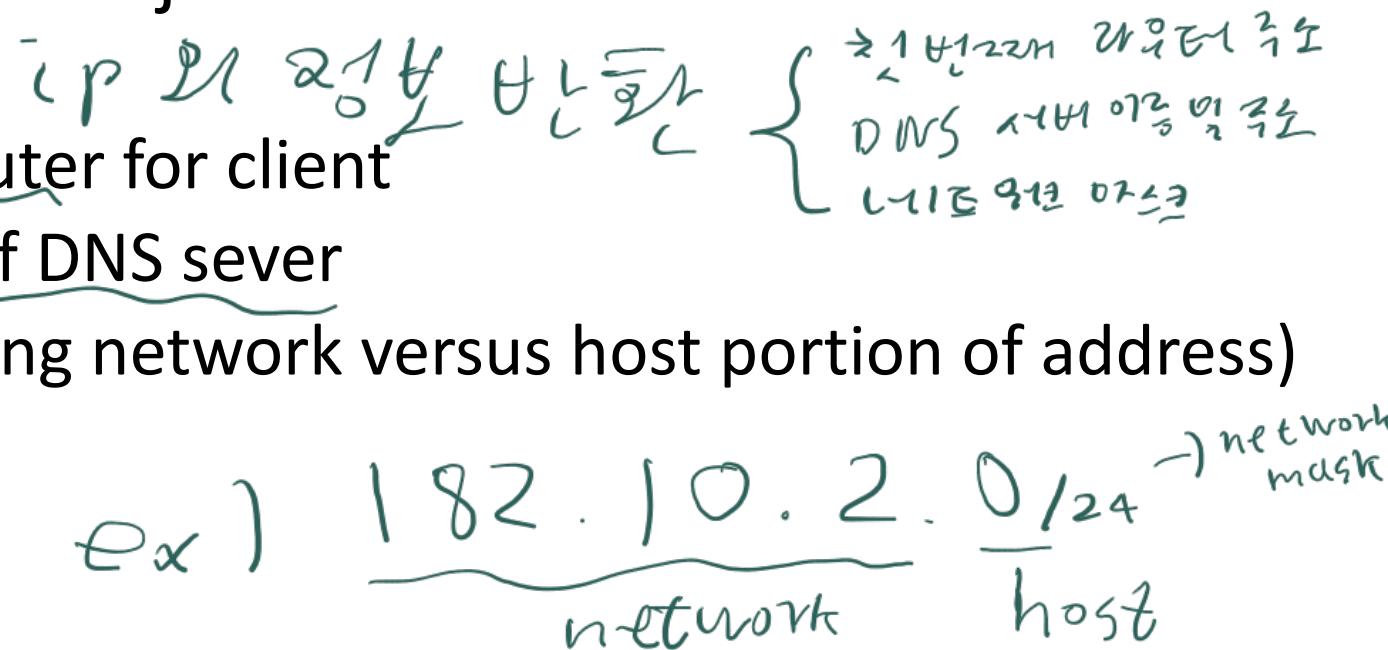
DHCP server: 223.1.2.5



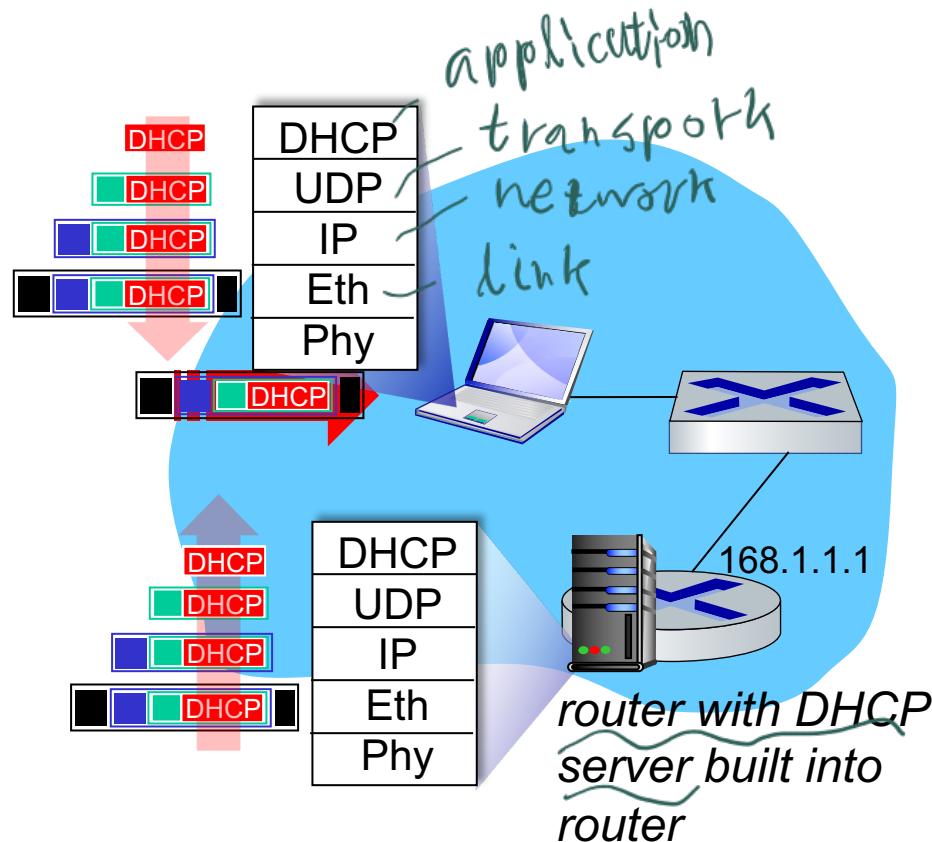
DHCP: more than IP addresses

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

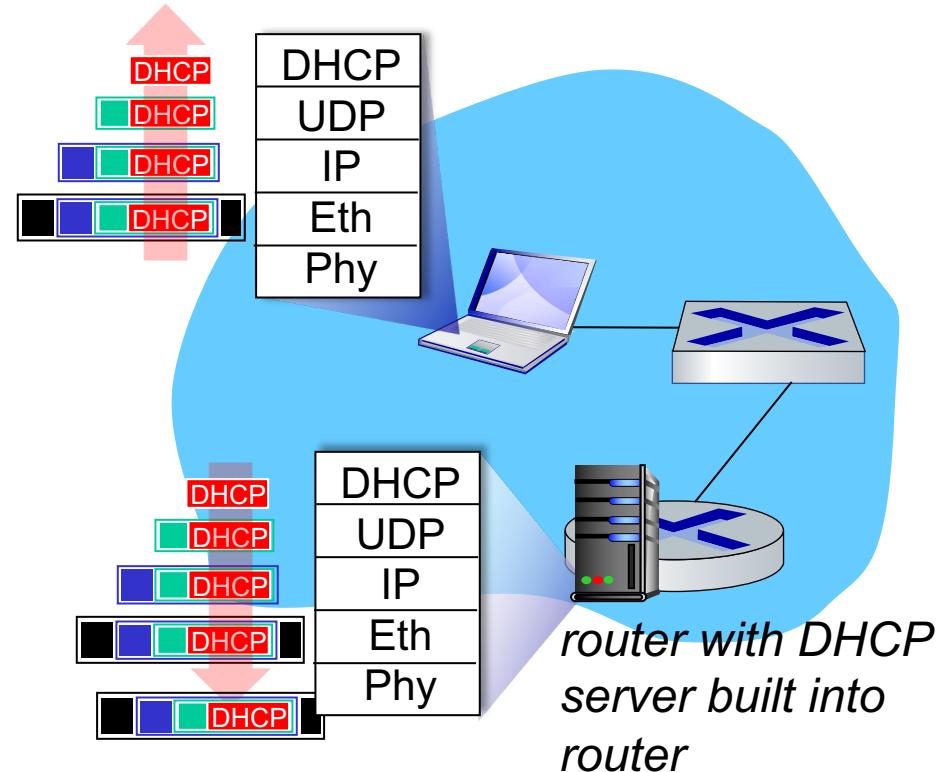


host 네트워크 연결

⇒ DHCP에서 IP, 첫 라우터 주소, DNS 서버 주소
확인

- Connecting laptop will use DHCP to get IP address, address of first-hop router, address of DNS server.
- DHCP REQUEST message encapsulated
 - 1. in UDP, encapsulated in IP, encapsulated in Ethernet **DHCP REQUEST** 양호화
 - 2. Ethernet frame broadcast (dest: FFFFFFFFFFFF) on LAN, received at router running DHCP server **frame** 방송
 - 3. Ethernet demux'ed to IP demux'ed, UDP demux'ed to DHCP **DHCP REQUEST** 복호화

DHCP: example



4.
 - DCP server formulates DHCP ACK containing client's IP address, IP address of first-hop router for client, name & IP address of DNS server

5.
 - encapsulated DHCP server reply forwarded to client, demuxing up to DHCP at client

6.
 - client now knows its IP address, name and IP address of DNS server, IP address of its first-hop router

IP addresses: how to get one?

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

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

ISP's block 11001000 00010111 00010000 00000000 200.23.16.0/20

\Rightarrow ISP space = subnet part

ISP can then allocate out its address space in 8 blocks:

9 bit \Rightarrow host part

Organization 0 11001000 00010111 00010000 00000000 200.23.16.0/23

Organization 1 11001000 00010111 00010010 00000000 200.23.18.0/23

Organization 2 11001000 00010111 00010100 00000000 200.23.20.0/23

...

.....

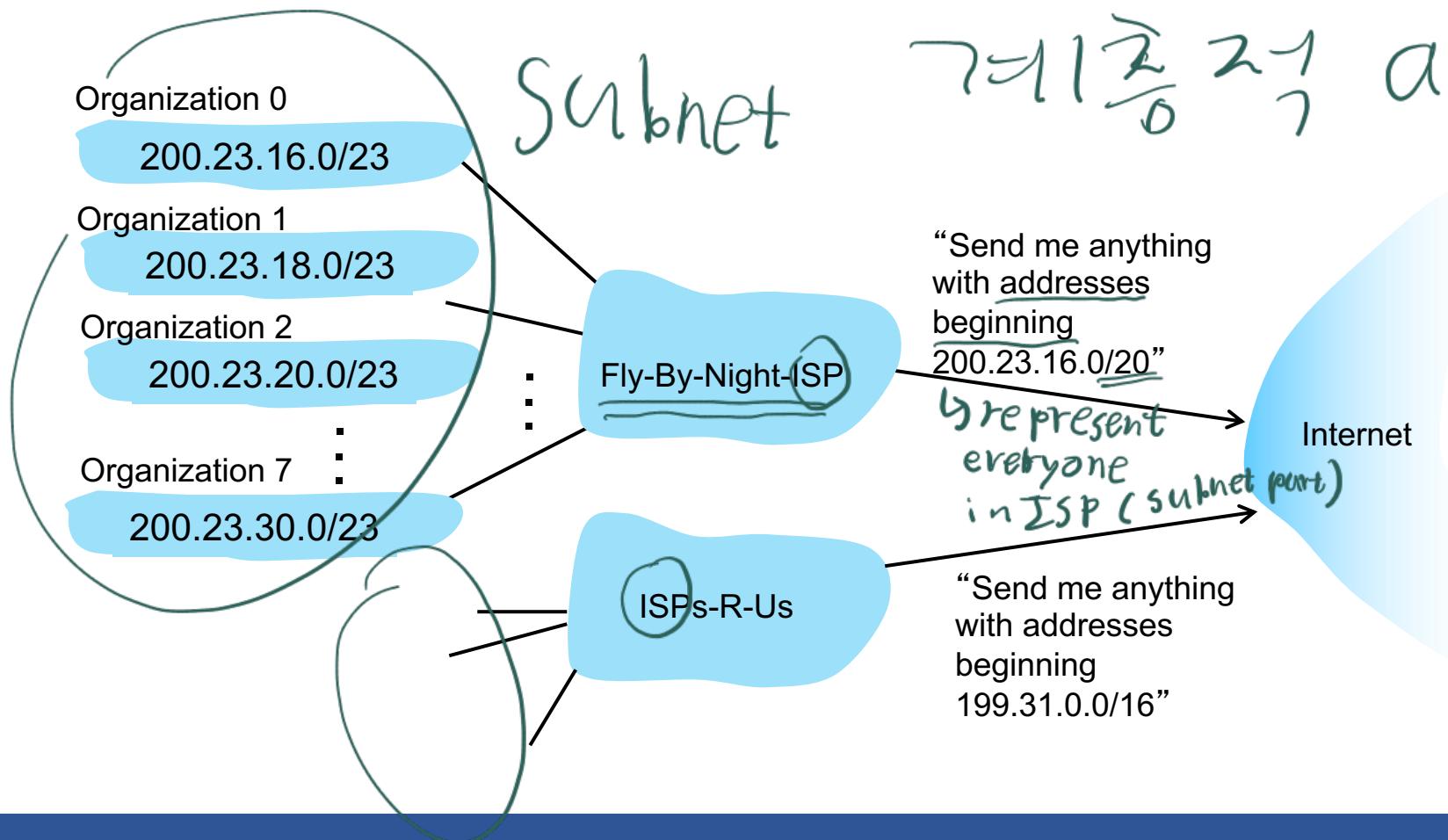
.....

....

Organization 7 11001000 00010111 00011110 00000000 200.23.30.0/23

Hierarchical addressing: route aggregation

hierarchical addressing allows efficient advertisement of routing information:



7-11 $\frac{2}{0}$ 21 addressing

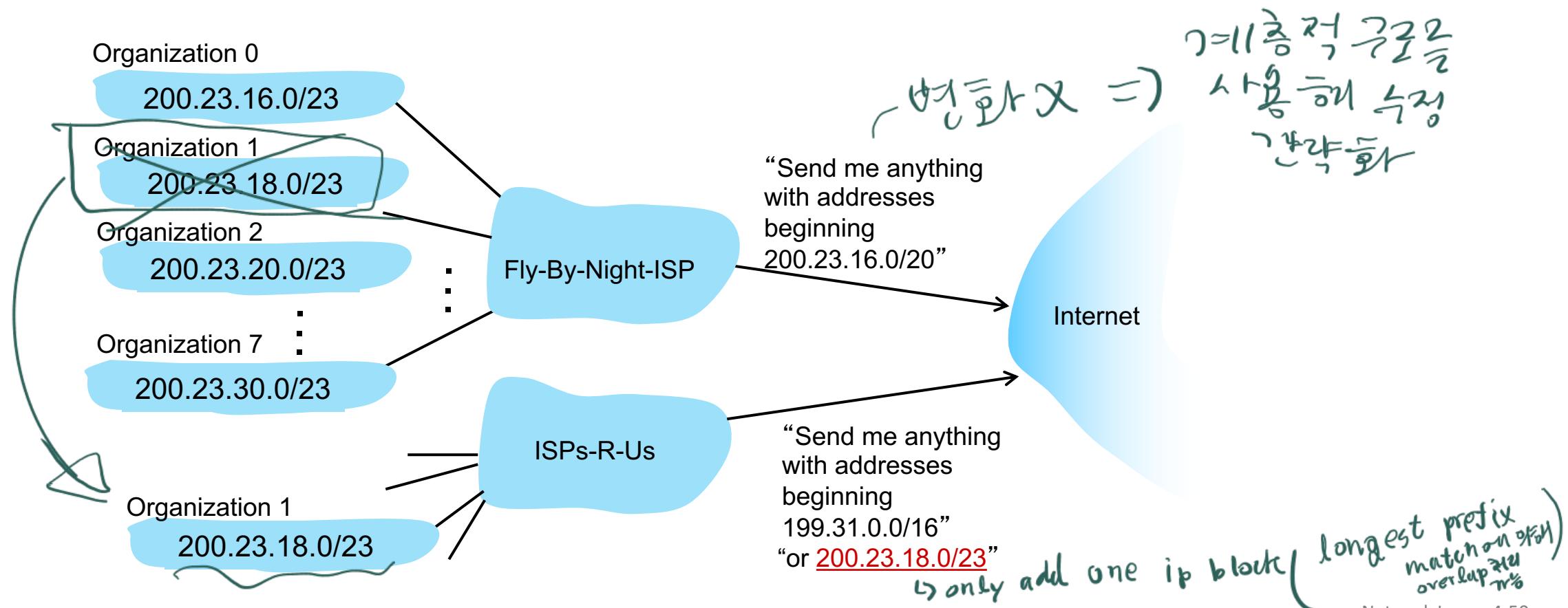
"Send me anything
with addresses
beginning
200.23.16.0/20"

↳ represent
everyone
in ISP (subnet part)

"Send me anything
with addresses
beginning
199.31.0.0/16"

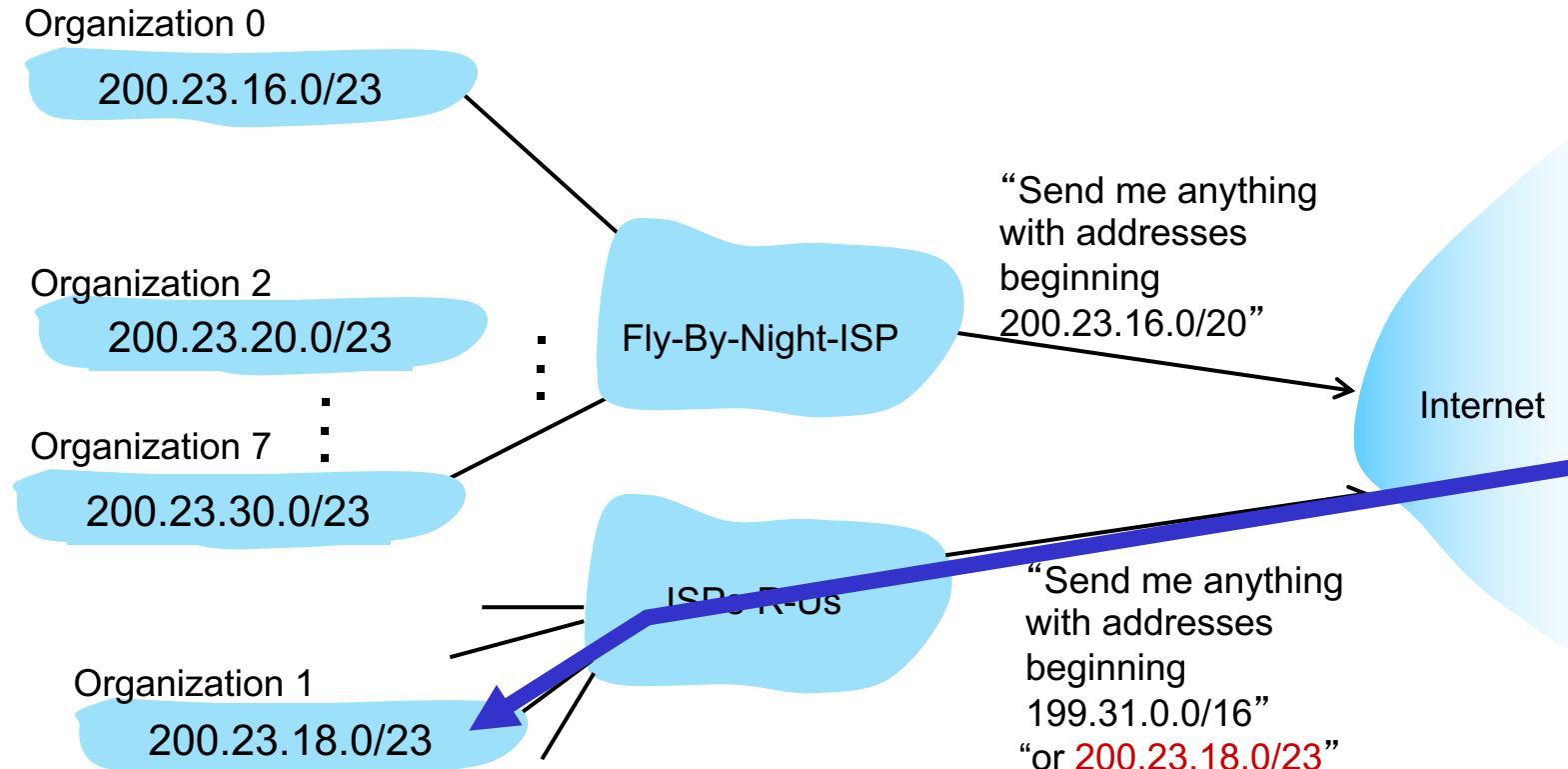
Hierarchical addressing: more specific routes

- Organization 1 moves from Fly-By-Night-ISP to ISPs-R-Us
- ISPs-R-Us now advertises a more specific route to Organization 1



Hierarchical addressing: more specific routes

- Organization 1 moves from Fly-By-Night-ISP to ISPs-R-Us
- ISPs-R-Us now advertises a more specific route to Organization 1



IP addressing: last words ...

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

TLD 포함 DNS root 영역
2^n-1

Q: are there enough 32-bit IP addresses? → 여전히 IPv4 주소 부족
not enough

- ICANN allocated last chunk of IPv4 addresses to RRs in 2011
- NAT (next) helps IPv4 address space exhaustion NAT로 실 ip 부족, 중복
- 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)

IP as identity

- Use IP as identities on the Internet

- E-commerce: to remember users (along with cookies)

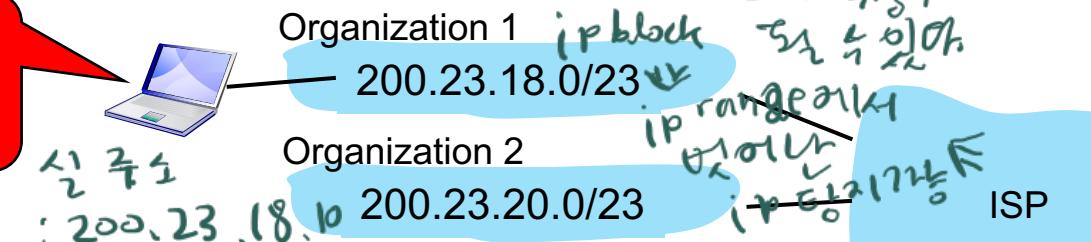
- Firewall: to identify misbehaving hosts

- ...

- Is it okay to use IP as identity?

- IP spoofing: what would happen?

Send this:
srcIP=100.10.10.25,
dstIP=142.250.206.206



ip address IP주소
⇒ identity를 사용할 수 있음
but 실제 사용하는 사용자

→ ban certain

ip 주소는 본인의 IP
가능한 대로
보안 차이 약화

Internet

IP spoofing:

- One of the open problems in the Internet (but why?)

- Difficult to incentivize ISPs

- All or nothing problem

ip spoofing은
악용 가능 X
가능한 대로
ISP 차이
보안 차이 약화
단일 규제 어려움

→ 전송 주소
: 100.10.10.25
srcIP를 변경해 IP속성
→ srcIP를 변경해 IP속성
이제는 체크를 해보면 알겠지
spoofing 가능

Network layer: “data plane” roadmap

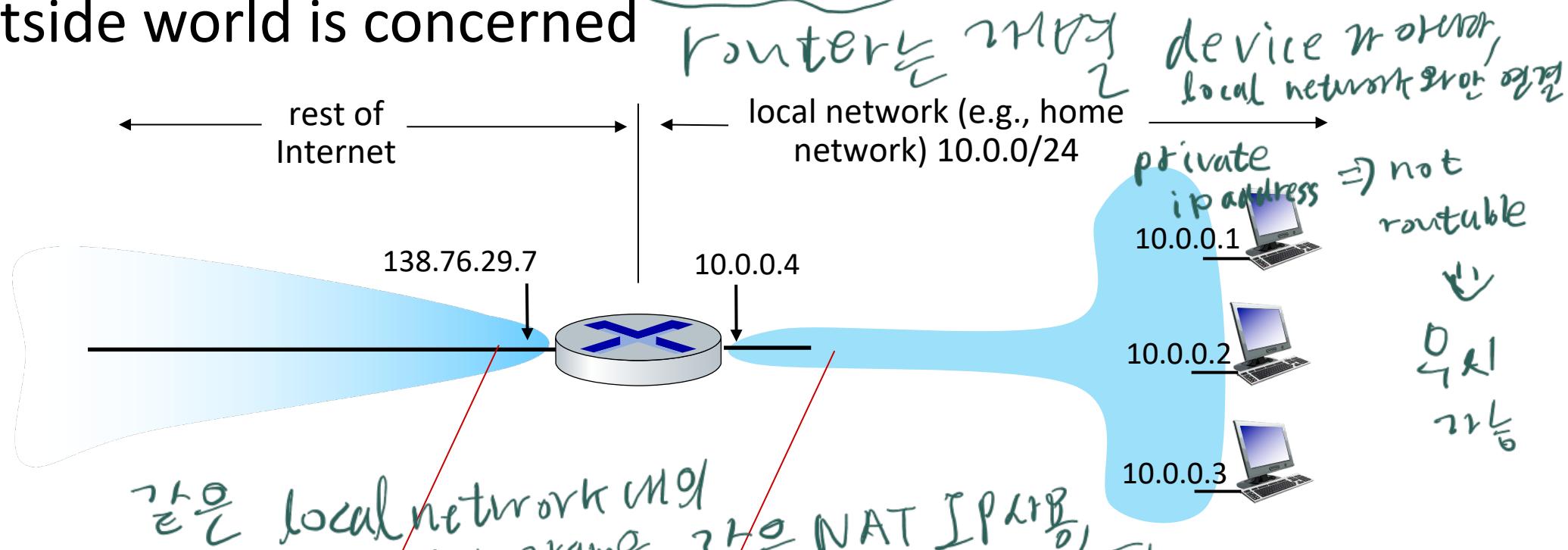
- Network layer: overview
 - data plane
 - control plane
- What's inside a router
 - input ports, switching, output ports
 - buffer management, scheduling
- IP: the Internet Protocol
 - datagram format
 - addressing
 - network address translation
 - IPv6
- Generalized Forwarding, SDN
 - match+action
 - OpenFlow: match+action in action
- Middleboxes



NAT: network address translation

사용 가능한
ip 주소

NAT: all devices in local network share just **one** IPv4 address as far as outside world is concerned



all datagrams *leaving* local network have *same* source NAT IP address: 138.76.29.7, but *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

- all devices in local network have 32-bit addresses in a “private” IP address space (10/8, 172.16/12, 192.168/16 prefixes) that can only be used in local network → local network 만에서만 통용되는 private ip 발급
- advantages:
 - just one IP address needed from provider ISP for all devices ISP당 1개
 - can change addresses of host in local network without notifying outside world → 외부 변경 없이 local network 내 ip 변경 안
 - can change ISP without changing addresses of devices in local network → device ip 주소 변경 없이 ISP 주소 변경 가능 (일반화 변경 x)
 - security: devices inside local net not directly addressable, visible by outside world
외부에서 local net 내부로 직접 연결 불가 있어

NAT: network address translation

네트워크
주소
번역

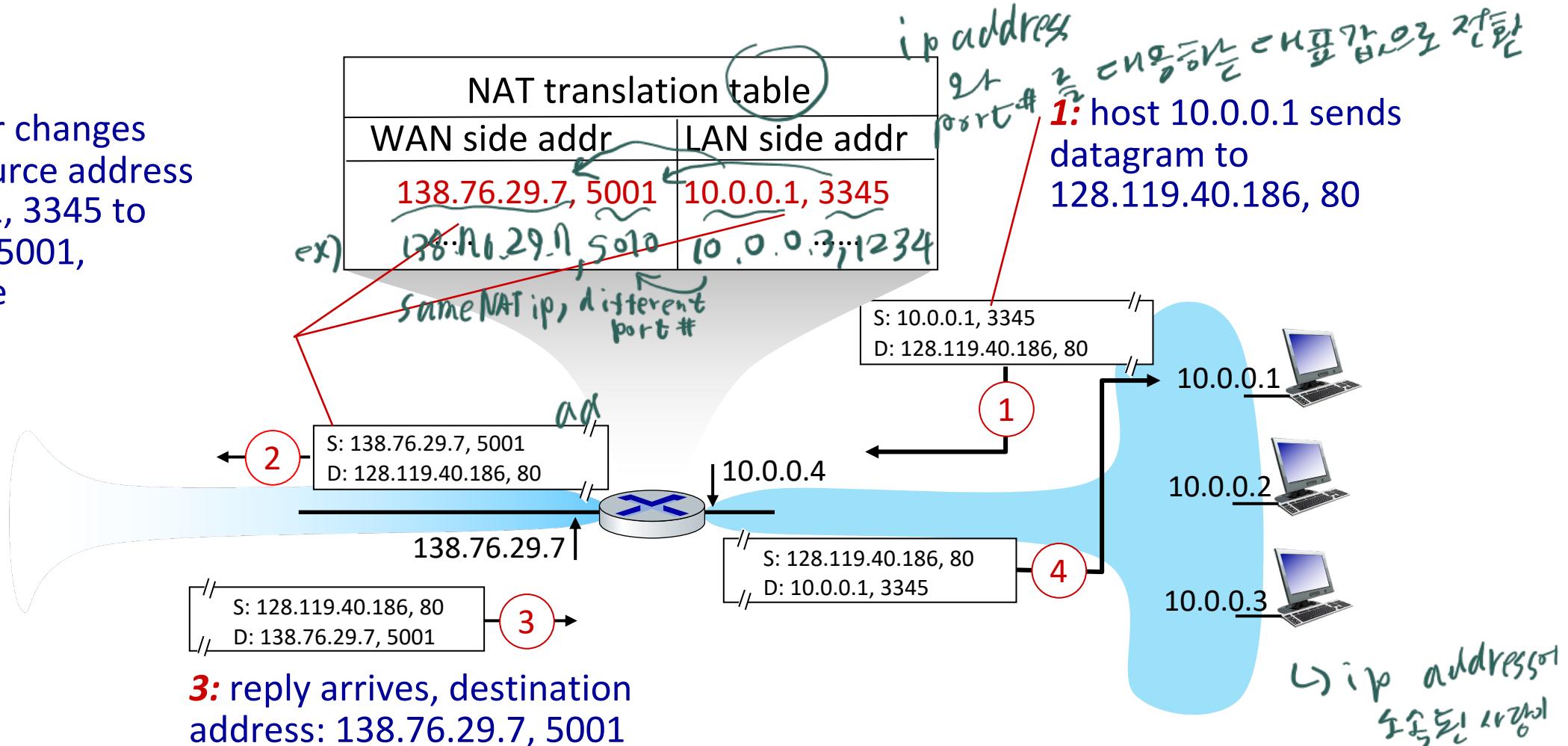
implementation: NAT router must (transparently):

- 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 address → 원격 서버나 NAT IP로
- 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 destination fields of every incoming datagram with corresponding (source IP address, port #) stored in NAT table

NAT: network address translation

→ IPv6로 네트워크 접속
일시적 대체주소

- 2:** NAT router changes datagram source address from 10.0.0.1, 3345 to 138.76.29.7, 5001, updates table



NAT: network address translation

논란이 있는

- NAT has been controversial:

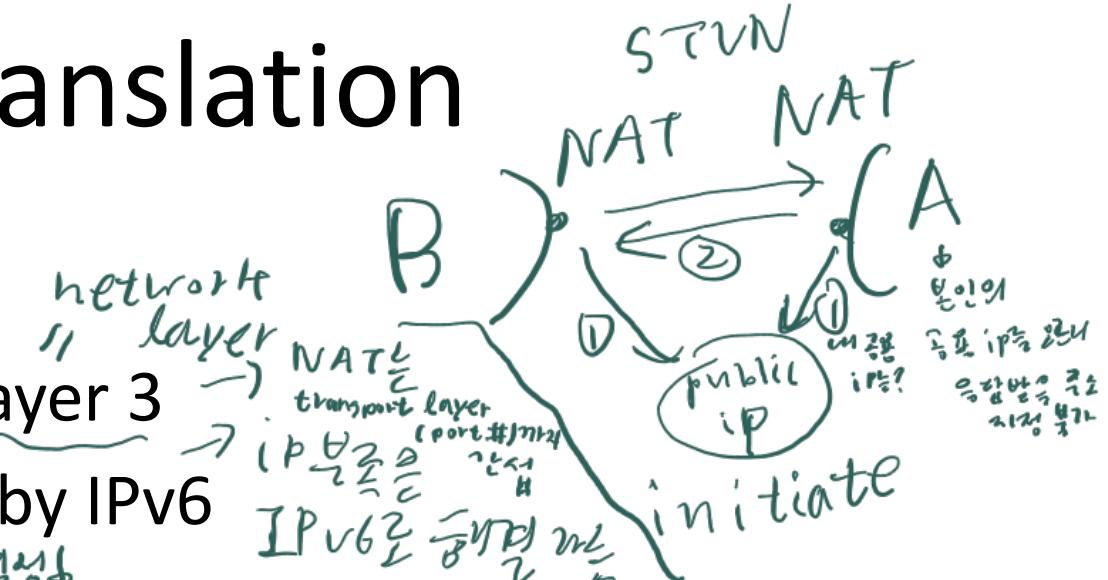
- routers “should” only process up to layer 3
- address “shortage” should be solved by IPv6
→ 중간 단계인 network layer에서 정부 설정 → 투명성
- violates end-to-end argument (port # manipulation by network-layer device)
- NAT traversal: what if client wants to connect to server behind NAT? hawk

- but NAT is here to stay:

- extensively used in home and institutional nets, 4G/5G cellular nets

↳ ipv6로의 전환은 대체로

업데이트가 필요해 NAT를
유지하는 설정



↳ AT 외부에서 연결하는 경우,
직접 연결이 불가능해 public ip로 연결해
서로의
ip를
얻어
통신하고

public ip는 아우나 접속할 수 없어
공격에 취약

IPv6: motivation

skip ↴

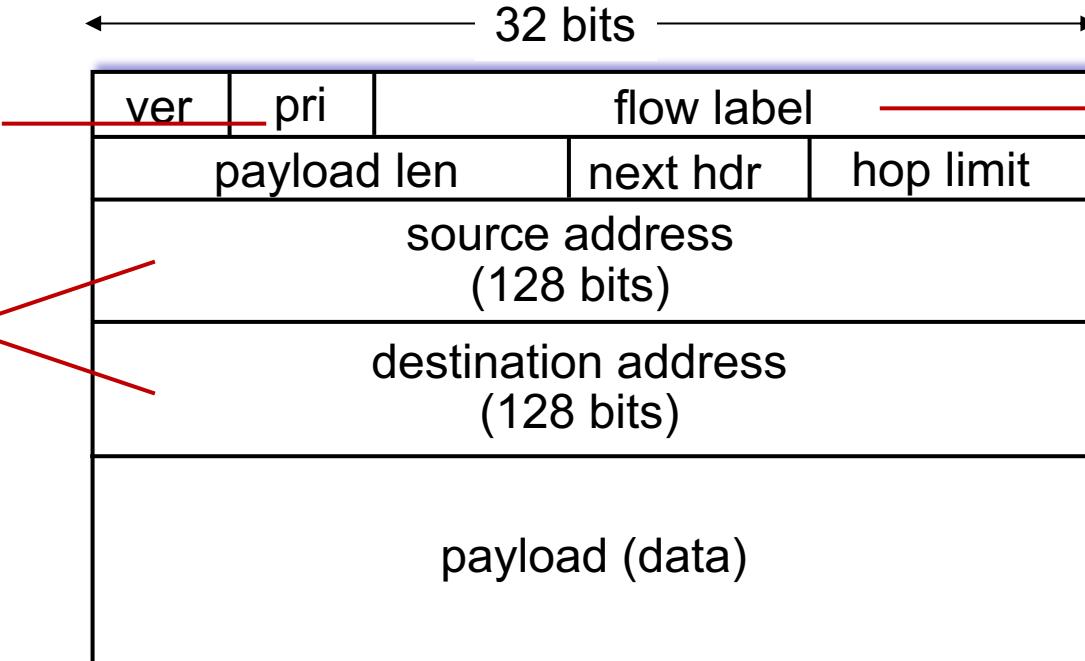
- **initial motivation:** 32-bit IPv4 address space would be completely allocated → 32-bit 주소 공간이 완전히 차운다.
- additional motivation:
 - speed processing/forwarding: 40-byte fixed length header
 - enable different network-layer treatment of “flows”

→ 흐름을
traffic
여러개 처리

IPv6 datagram format

priority: identify priority among datagrams in flow

128-bit
IPv6 addresses



flow label: identify datagrams in same "flow." (concept of "flow" not well defined).

What's missing (compared with IPv4):

- no checksum (to speed processing at routers) → checksum 제거함
- no fragmentation/reassembly → packet 단편화 및 재조립 X → 라우터 부담
- no options (available as upper-layer, next-header protocol at router) → option은 상위 layer에서 처리

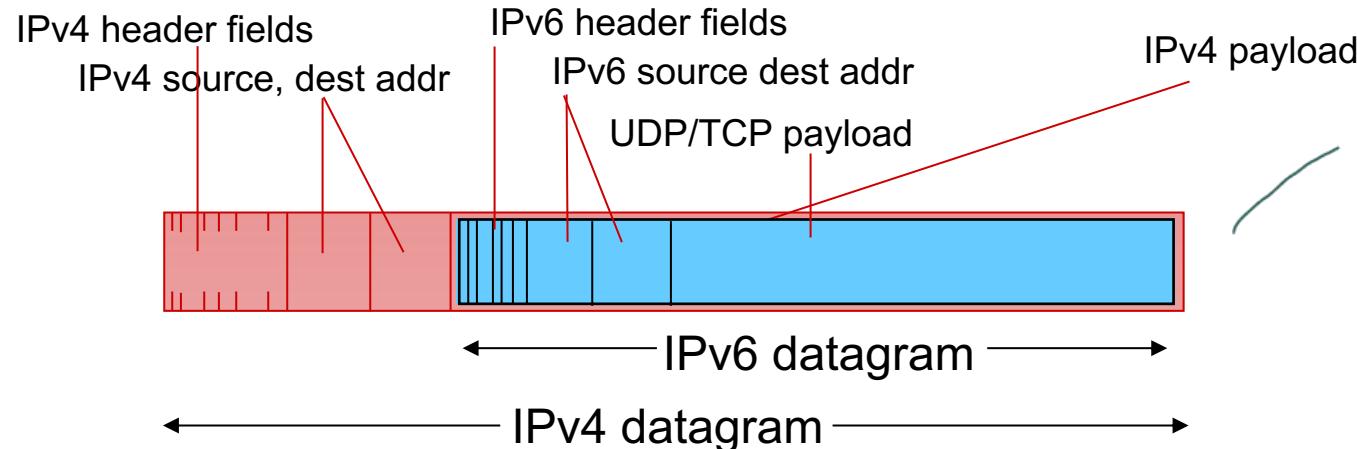
Transition from IPv4 to IPv6

모든 라우터 동시에 업그레이드 불가

- not all routers can be upgraded simultaneously
 - no “flag days”
 - how will network operate with mixed IPv4 and IPv6 routers?

→ 전화로 바꾸면서 혼용할
방법?

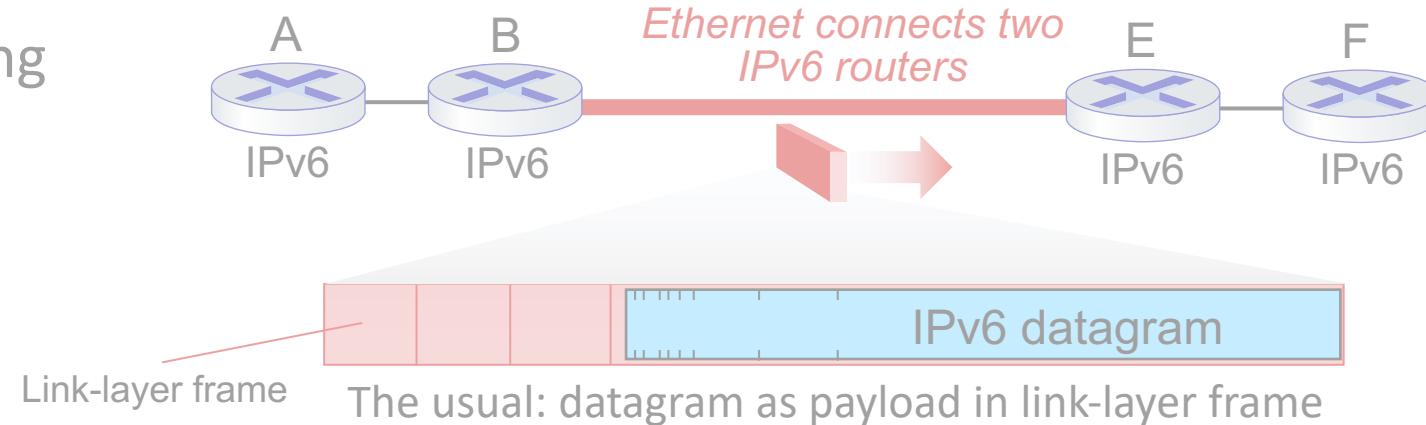
- **tunneling:** IPv6 datagram carried as *payload* in IPv4 datagram among IPv4 routers (“packet within a packet”)
 - tunneling used extensively in other contexts (4G/5G)



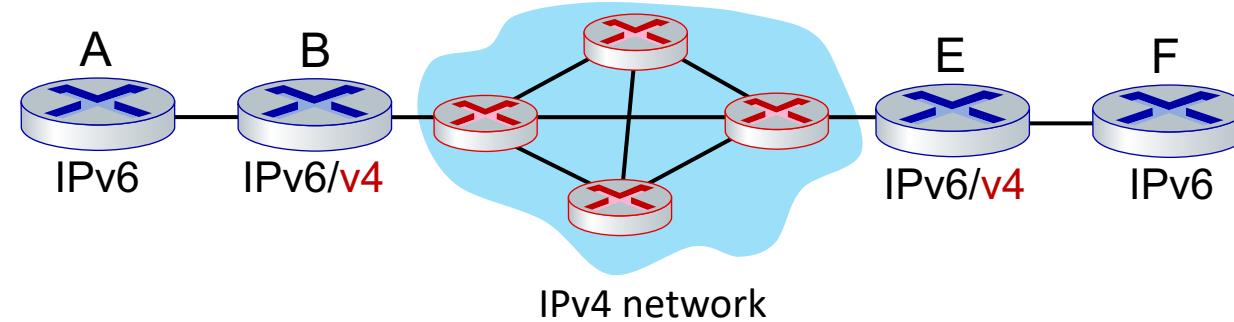
ip v4 datagram
내부
ip v6 datagram
외부

Tunneling and encapsulation

Ethernet connecting
two IPv6 routers:

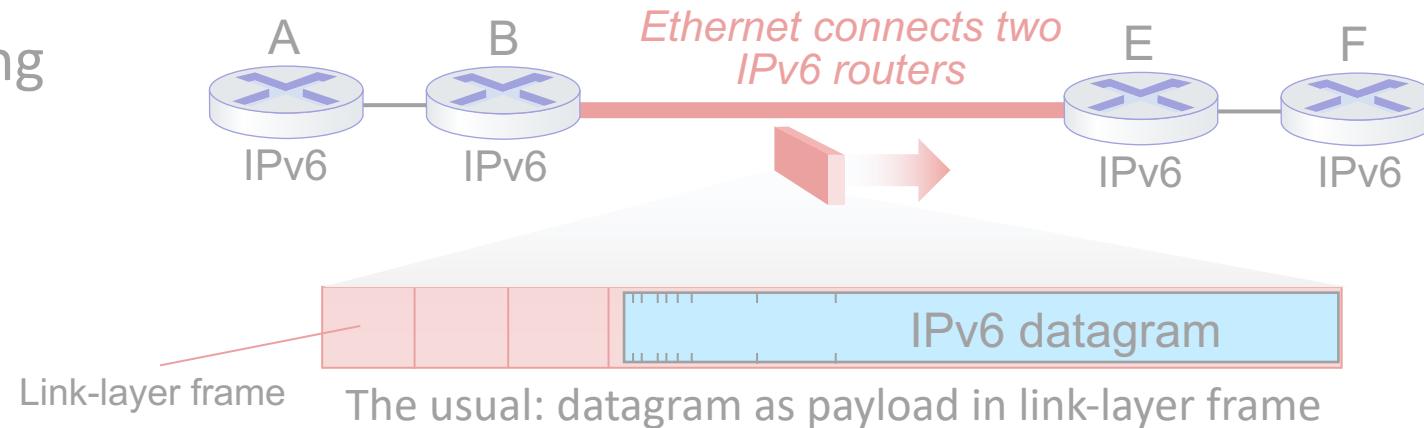


IPv4 network
connecting two
IPv6 routers

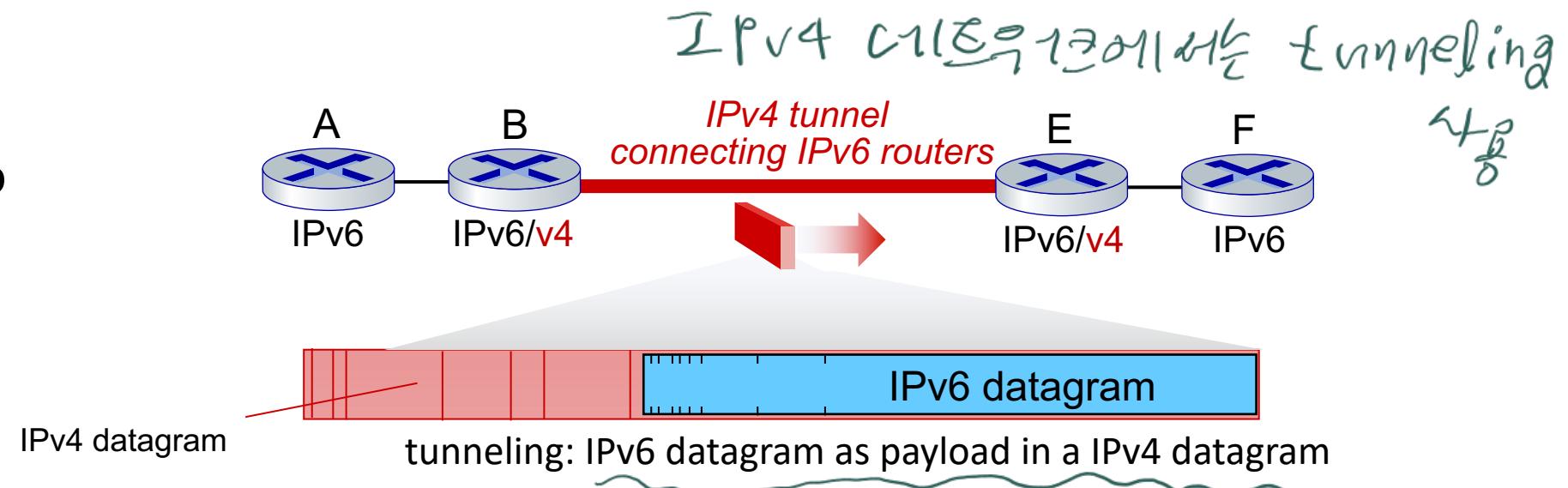


Tunneling and encapsulation

Ethernet connecting
two IPv6 routers:



IPv4 tunnel
connecting two
IPv6 routers



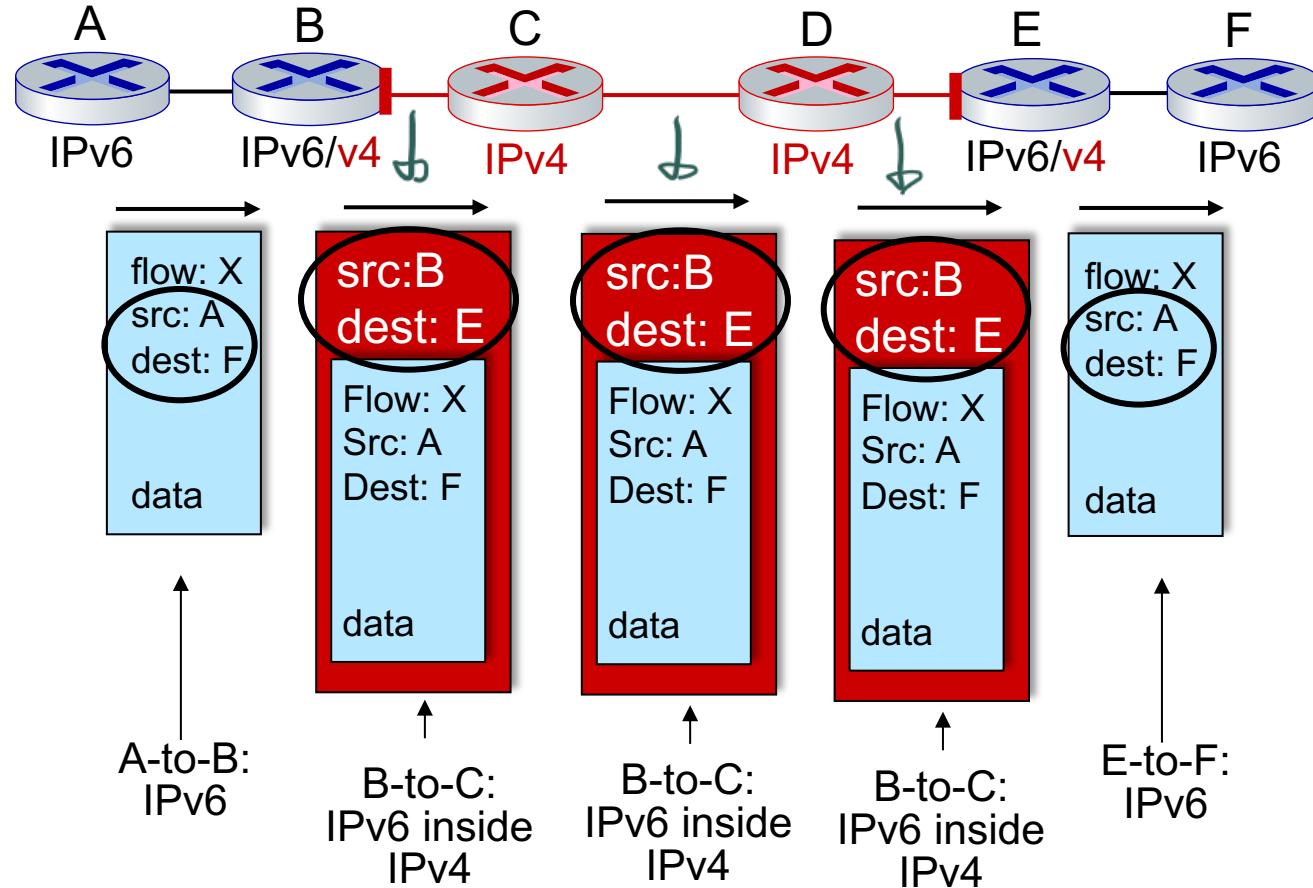
Tunneling

logical view:



physical view:

Note source and destination addresses!



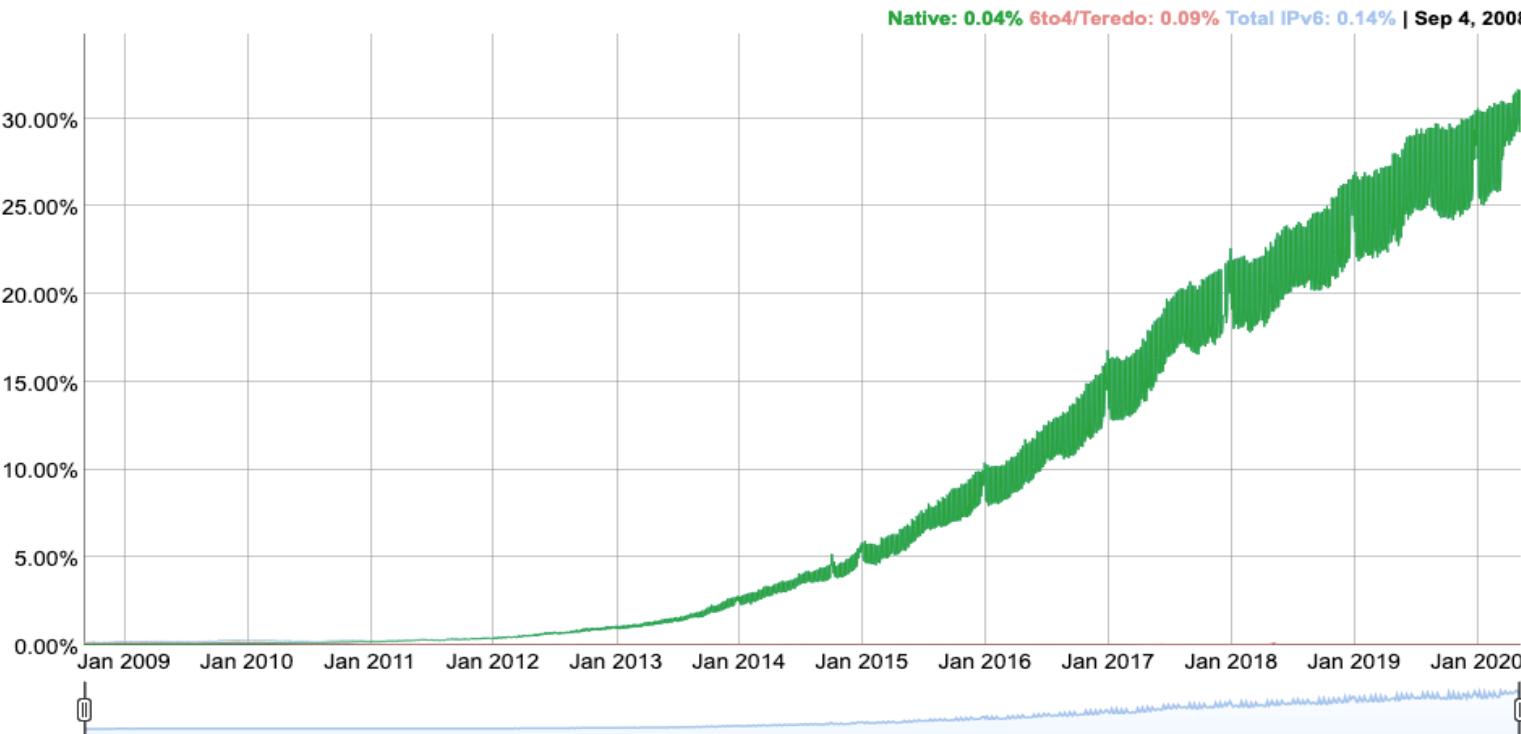
IPv6: adoption

KIEH

- Google¹: ~ 30% of clients access services via IPv6
- NIST: 1/3 of all US government domains are IPv6 capable

IPv6 Adoption

We are continuously measuring the availability of IPv6 connectivity among Google users. The graph shows the percentage of users that access Google over IPv6.



1

<https://www.google.com/intl/en/ipv6/statistics.html>

IPv6: adoption

SKIP

- Google¹: ~ 30% of clients access services via IPv6
- NIST: 1/3 of all US government domains are IPv6 capable
- Long (long!) time for deployment, use
 - 25 years and counting! ~25년 (약 25년)
 - think of application-level changes in last 25 years: WWW, social media, streaming media, gaming, telepresence, ...
 - *Why?*

¹ <https://www.google.com/intl/en/ipv6/statistics.html>

Network layer: “data plane” roadmap

- Network layer: overview
 - data plane
 - control plane
- What's inside a router
 - input ports, switching, output ports
 - buffer management, scheduling
- IP: the Internet Protocol
 - datagram format
 - addressing
 - network address translation
 - IPv6



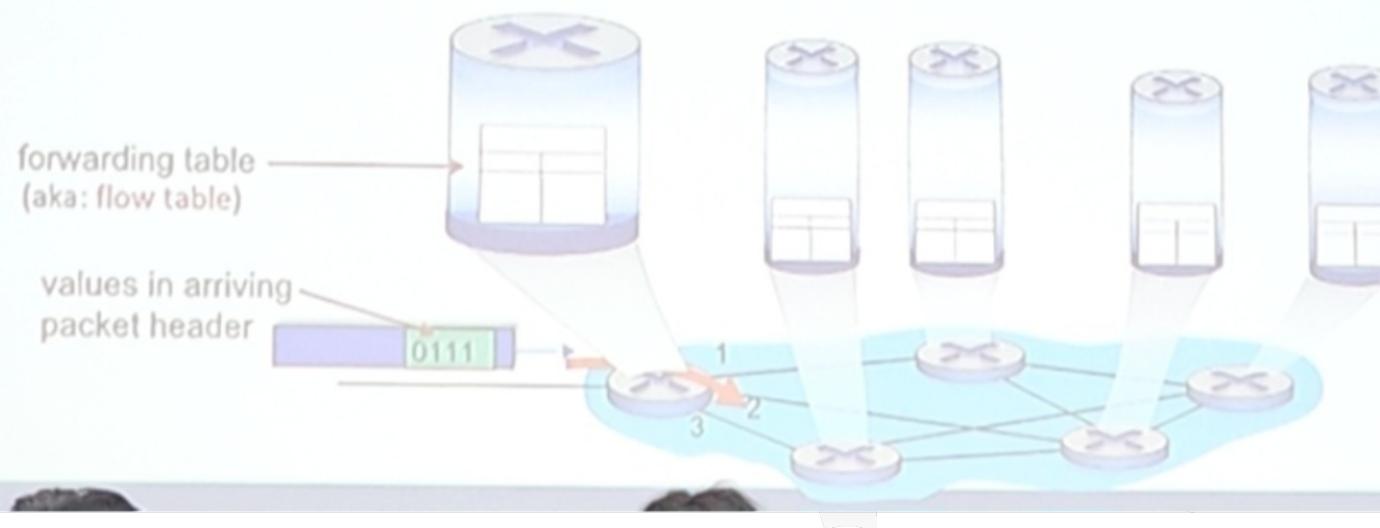
- Generalized Forwarding, SDN
 - Match+action
 - OpenFlow: match+action in action
- Middleboxes

Generalized forwarding: match plus action

Review: each router contains a forwarding table (aka: flow table)

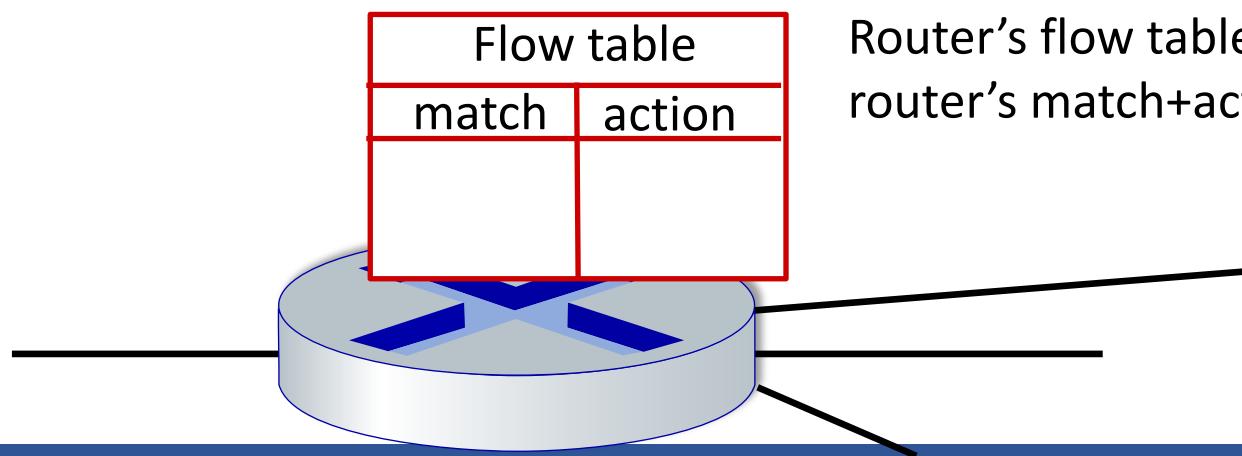
bit ↳ ↳ ↳ action

- “match plus action” abstraction: match bits in arriving packet, take action
 - destination-based forwarding: forward based on dest. IP address
 - generalized forwarding:
 - many header fields can determine action
 - many actions possible: drop/copy/modify/log packet



Flow table abstraction

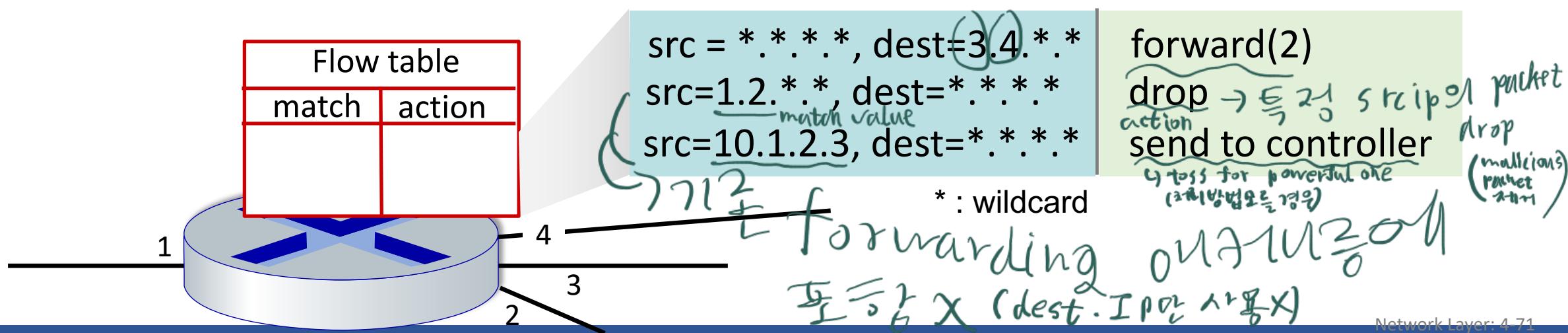
- **flow:** defined by header field values (in link-, network-, transport-layer fields)
- **generalized forwarding:** simple packet-handling rules
 - **match:** pattern values in packet header fields
 - **actions:** for matched packet: drop, forward, modify, matched packet or send matched packet to controller
 다중화하지 않은 큐
 - **priority:** disambiguate overlapping patterns
 - **counters:** #bytes and #packets



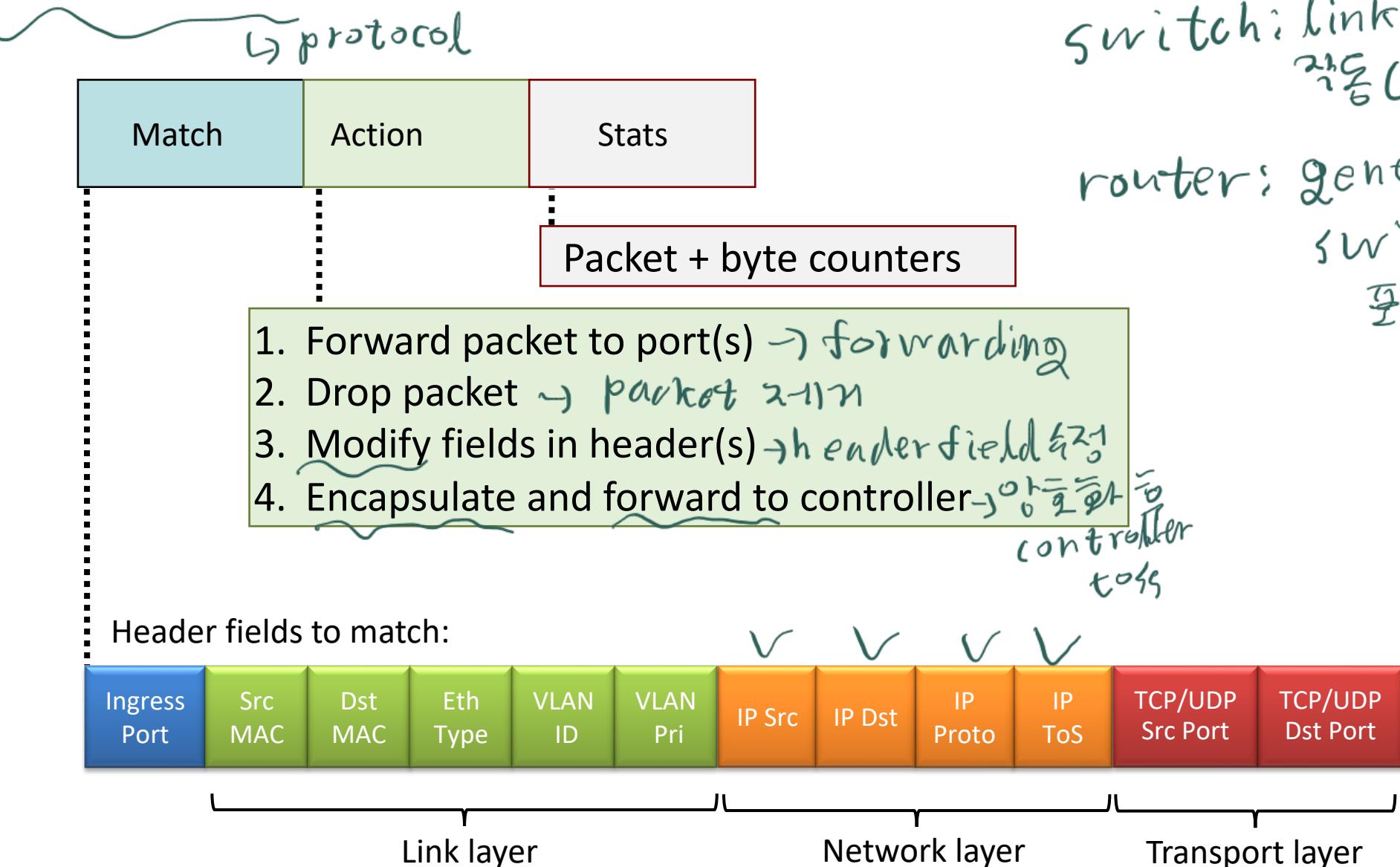
Router's flow table define
router's match+action rules

Flow table abstraction

- **flow:** defined by header fields
- **generalized forwarding:** simple packet-handling rules
 - **match:** pattern values in packet header fields
 - **actions:** for matched packet: drop, forward, modify, matched packet or send matched packet to controller
 - **priority:** disambiguate overlapping patterns
 - **counters:** #bytes and #packets



OpenFlow: flow table entries



switch: link layer or higher
작동 (ip x)

router: general,
switch는,
포함

OpenFlow: examples

Destination-based forwarding:

Switch Port	MAC src	MAC dst	Eth type	VLAN ID	VLAN Pri	IP Src	IP Dst	IP Prot	IP ToS	TCP s-port	TCP d-port	Action
*	*	*	*	*	*	*	51.6.0.8	*	*	*	*	port6

IP datagrams destined to IP address 51.6.0.8 should be forwarded to router output port 6

Firewall:

Switch Port	MAC src	MAC dst	Eth type	VLAN ID	VLAN Pri	IP Src	IP Dst	IP Prot	IP ToS	TCP s-port	TCP d-port	Action
*	*	*	*	*	*	*	*	*	*	*	22	drop

Block (do not forward) all datagrams destined to TCP port 22 (ssh port #)

Switch Port	MAC src	MAC dst	Eth type	VLAN ID	VLAN Pri	IP Src	IP Dst	IP Prot	IP ToS	TCP s-port	TCP d-port	Action
*	*	*	*	*	*	*	128.119.1.1	*	*	*	*	drop

Block (do not forward) all datagrams sent by host 128.119.1.1

OpenFlow: examples

Layer 2 destination-based forwarding:

Switch Port	MAC src	MAC dst	Eth type	VLAN ID	VLAN Pri	IP Src	IP Dst	IP Prot	IP ToS	TCP s-port	TCP d-port	Action
-------------	---------	---------	----------	---------	----------	--------	--------	---------	--------	------------	------------	--------

* * 22:A7:23:
11:E1:02 * * * * * * * * * * * port3
layer 2 frames with destination MAC address 22:A7:23:11:E1:02 should be forwarded to output port 3

OpenFlow abstraction

- **match+action:** abstraction unifies different kinds of devices

측정화를 통해
erof한 종류의
구조(동일)

Router

- *match:* longest destination IP prefix
- *action:* forward out a link

Switch

- *match:* destination MAC address
- *action:* forward or flood

Firewall

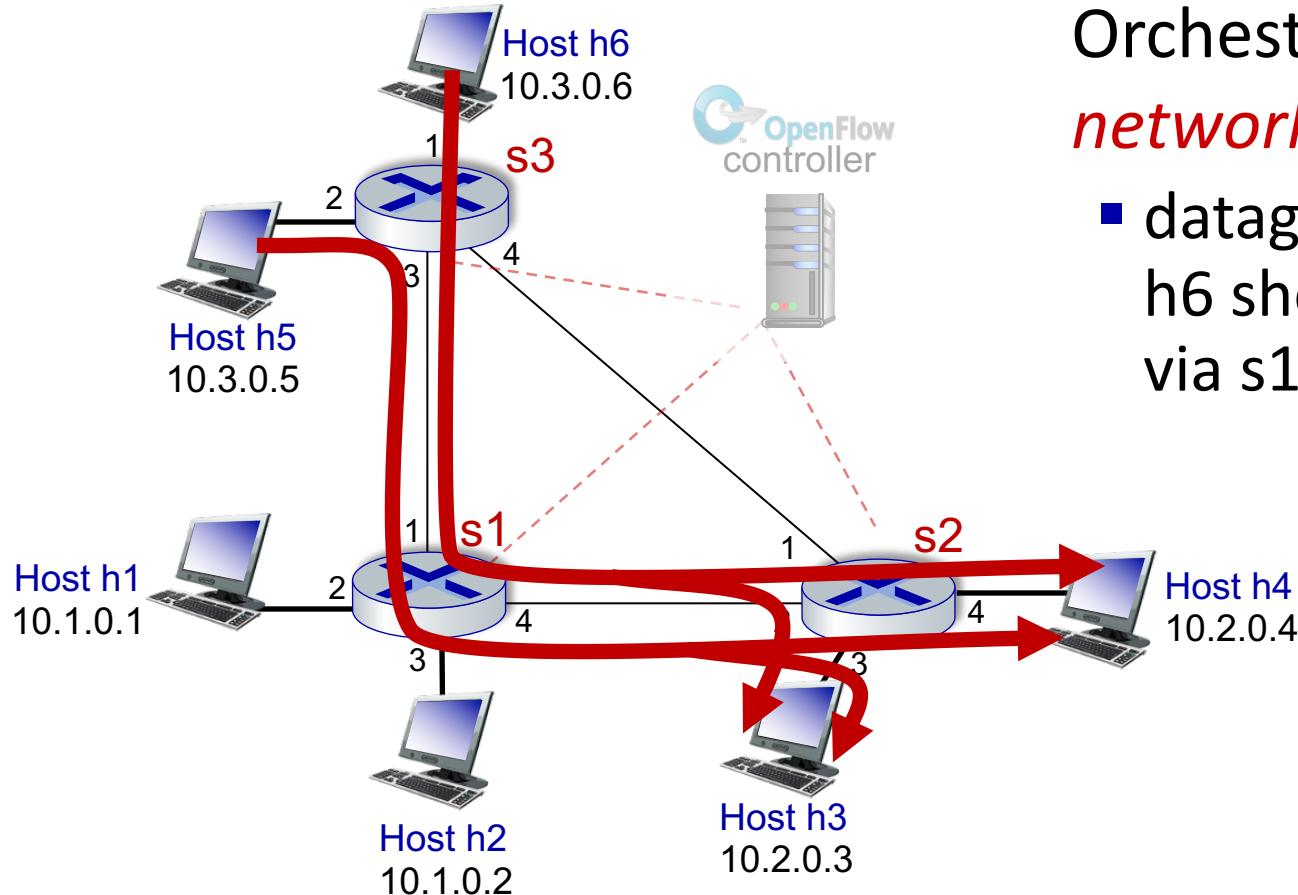
- *match:* IP addresses and TCP/UDP port numbers
- *action:* permit or deny

NAT

- *match:* IP address and port
- *action:* rewrite address and port

⇒ match-action으로 전부 표현 가능

OpenFlow example



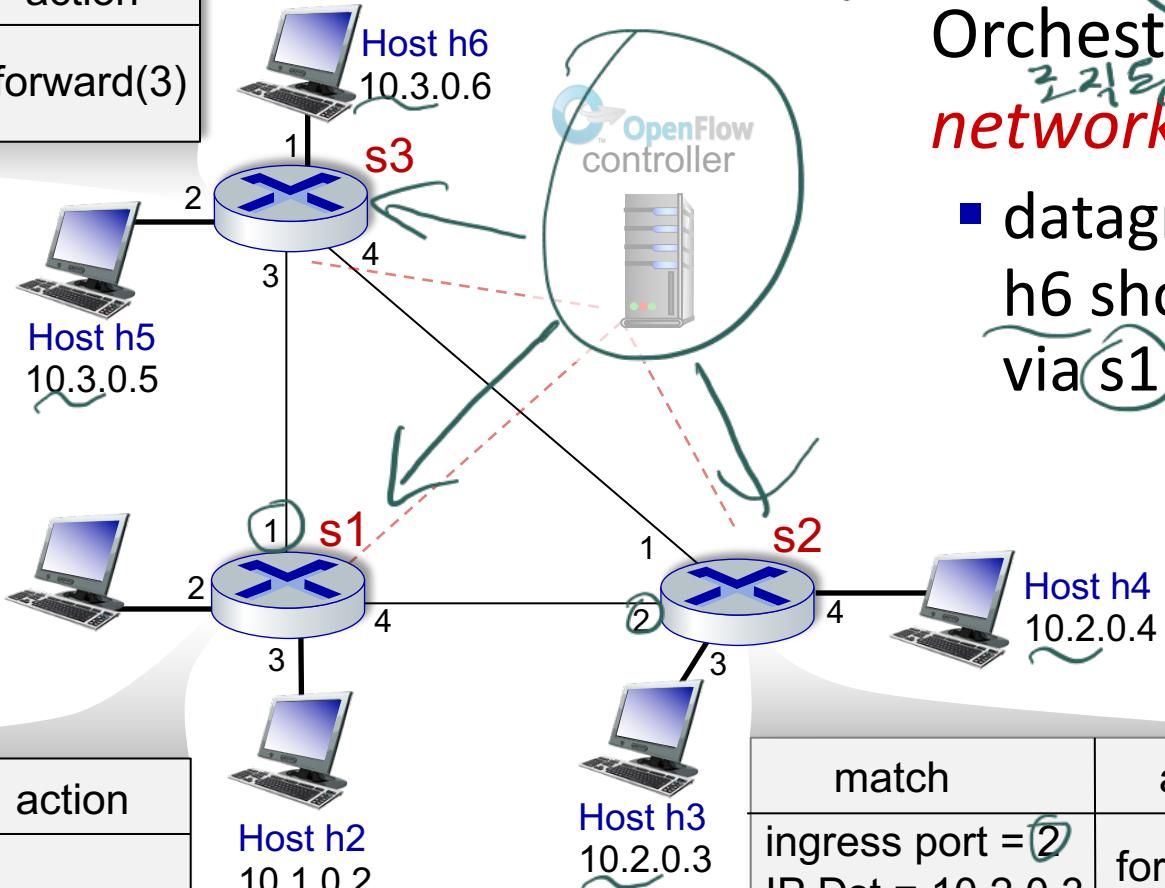
Orchestrated tables can create *network-wide* behavior, e.g.,:

- datagrams from hosts h5 and h6 should be sent to h3 or h4, via s1 and from there to s2

OpenFlow example

match	action
IP Src = 10.3.*.*	
IP Dst = 10.2.*.*	forward(3)

host 5/6
→ host 3/4



match	action
ingress port = 1	
IP Src = 10.3.*.*	
IP Dst = 10.2.*.*	forward(4)

match	action
ingress port = 2 IP Dst = 10.2.0.3	forward(3)
ingress port = 2 IP Dst = 10.2.0.4	forward(4)

write network function with single thread of high-level language \rightarrow 고등 영어로
서버에 네트워크
function을 작성하는 경우

Orchestrated tables can create
network-wide behavior, e.g.,:

- datagrams from hosts h5 and h6 should be sent to h3 or h4, via s1 and from there to s2

Generalized forwarding: summary

- “match plus action” abstraction: match bits in arriving packet header(s) in any layers, take action
 - matching over many fields (link-, network-, transport-layer)
 - local actions: drop, forward, modify, or send matched packet to controller
 - “program” network-wide behaviors
- simple form of “network programmability”
 - programmable, per-packet “processing”
 - *historical roots*: active networking
 - *today*: more generalized programming:
P4 (see p4.org). *flexible action*
→ code in AI₂ assign flow table entry

Network Layer: Data Plane

- Overview of Network Layer
- What's Inside a Router?
- The Internet Protocol: IPv4, Addressing, NAT IPv6
- Generalized Forwarding and SDN
- **Middleboxes**
 - middlebox functions
 - evolution, architectural principles of the Internet
- Summary

COMPSCI 453 **Computer Networks**

Professor Jim Kurose

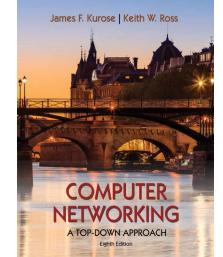
College of Information and Computer Sciences

University of Massachusetts



Class textbook:
*Computer Networking: A Top-
Down Approach* (8th ed.)

J.F. Kurose, K.W. Ross
Pearson, 2020
http://gaia.cs.umass.edu/kurose_ross

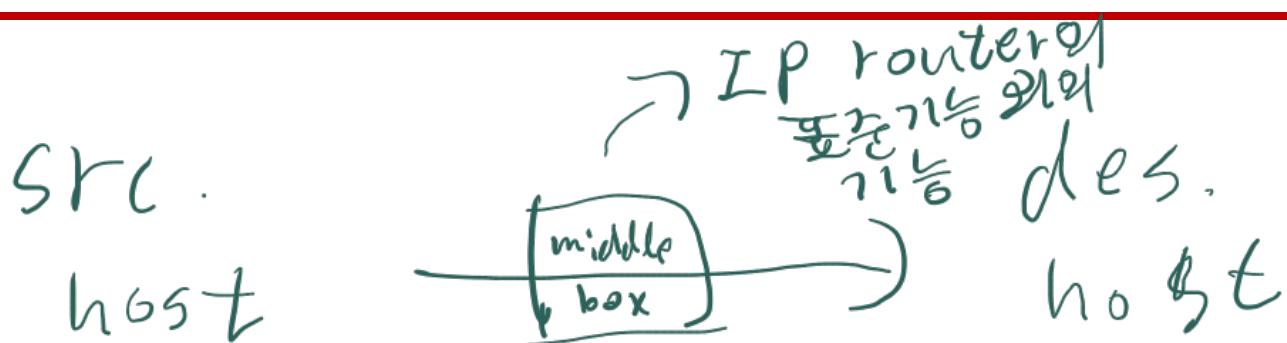


Middleboxes

skip
↓

Middlebox (RFC 3234)

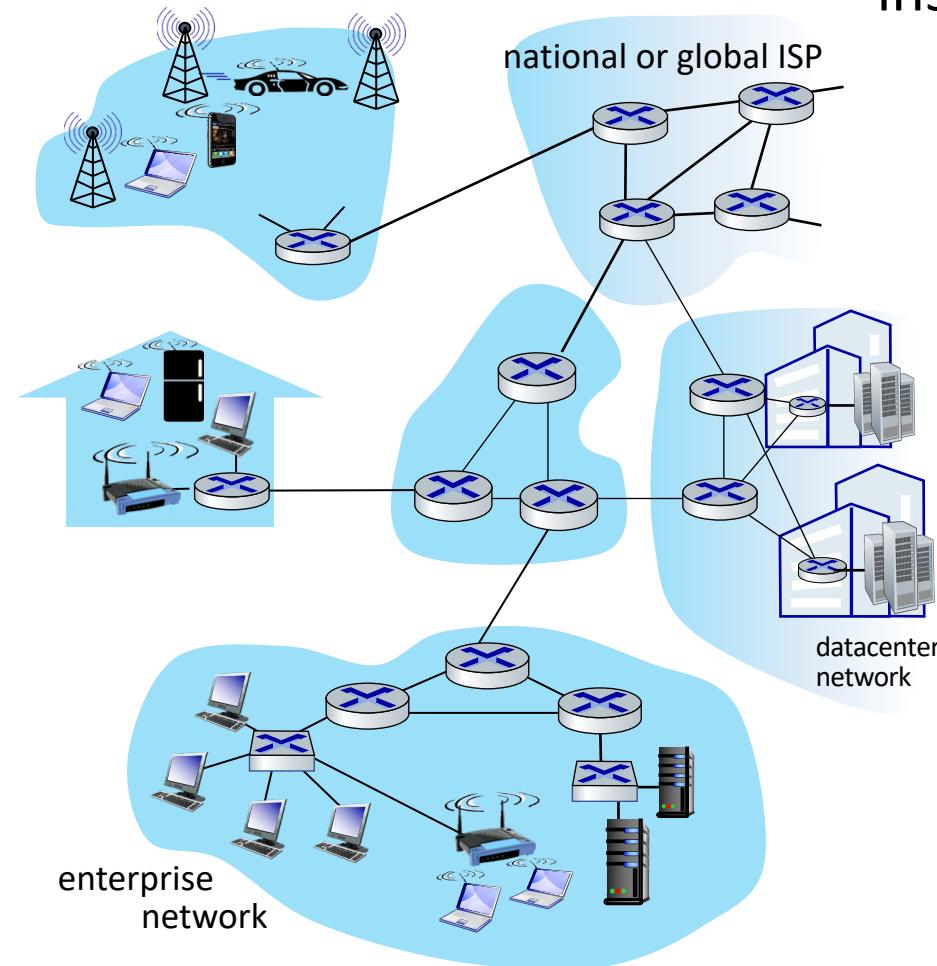
“any intermediary box performing functions apart from normal, standard functions of an IP router on the data path between a source host and destination host”



Middleboxes everywhere!

NAT: home,
cellular,
institutional

Application-specific: service
providers,
institutional,
CDN



Firewalls, IDS: corporate,
institutional, service providers,
ISPs

Load balancers:
corporate, service
provider, data center,
mobile nets

Caches: service
provider, mobile, CDNs

Middleboxes

독점적(폐쇄적) hw 솔루션

- initially: proprietary (closed) hardware solutions
- move towards “whitebox” hardware implementing open API
 - move away from proprietary hardware solutions
 - programmable local actions via match+action
 - move towards innovation/differentiation in software
- SDN: (logically) centralized control and configuration management often in private/public cloud
- network functions virtualization (NFV): programmable services over white box networking, computation, storage

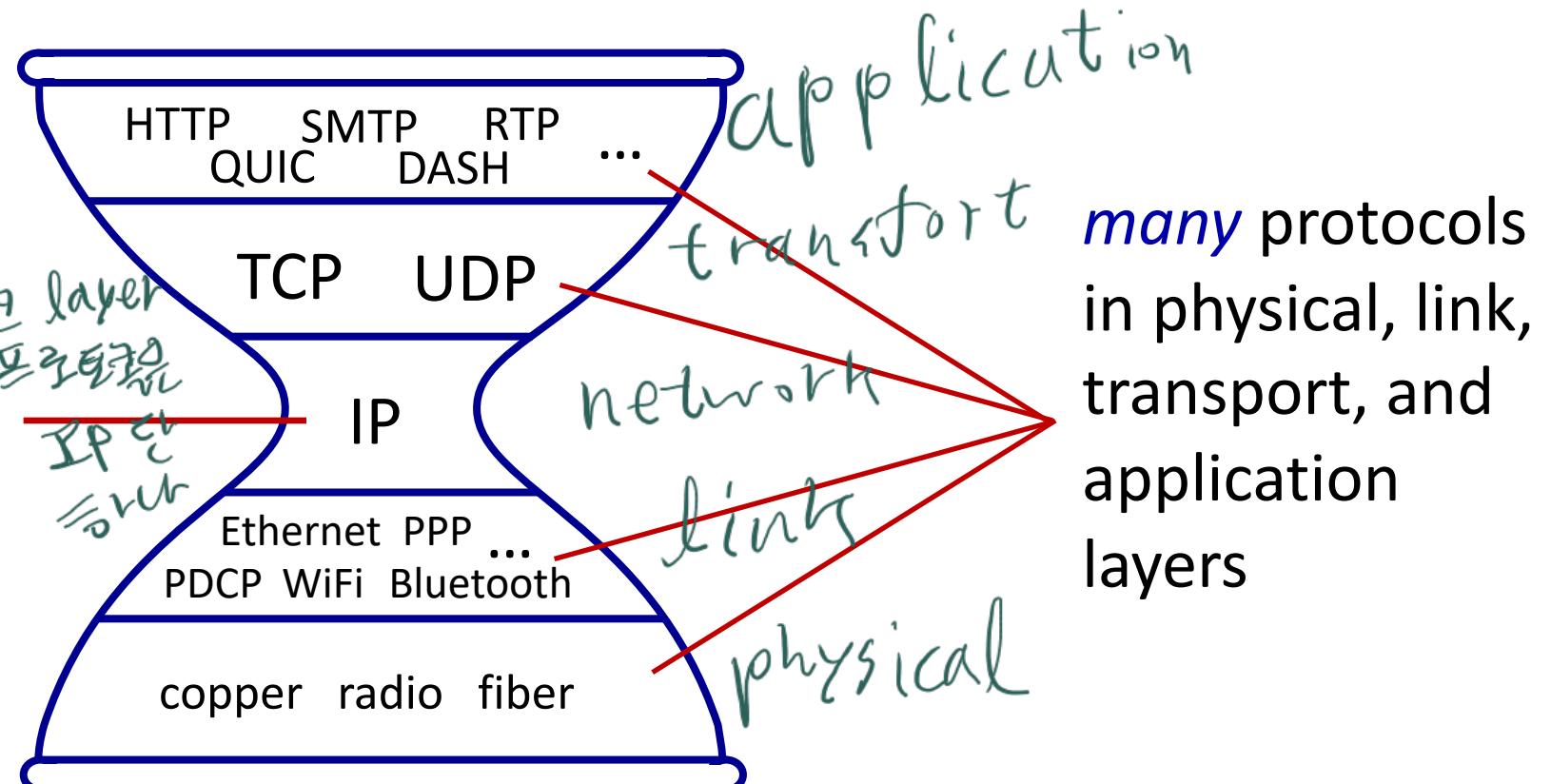
마켓
공개화
중립화

match+action
으로 흐르게
(whitebox)

구성

The IP hourglass

- Internet's "thin waist":
- one network layer protocol: IP
 - must be implemented by every (billions) of Internet-connected devices



The IP hourglass, at middle age

t
skip

Internet's middle age
“love handles”?

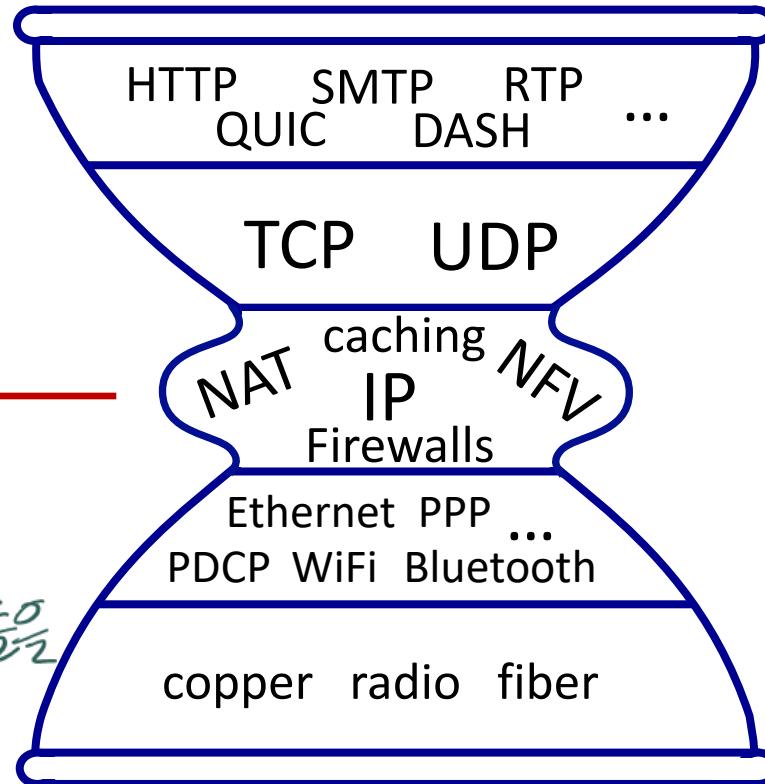
- middleboxes,
operating inside the
network

↳ middlebox 들

추가 학습

network

layer protocol



Architectural Principles of the Internet

↳ most fundamental principle

RFC 1958

"Many members of the Internet community would argue that there is no architecture, but only a tradition, which was not written down for the first 25 years (or at least not by the IAB). However, in very general terms, the community believes that

the goal is connectivity, the tool is the Internet Protocol, and the intelligence is end to end rather than hidden in the network."

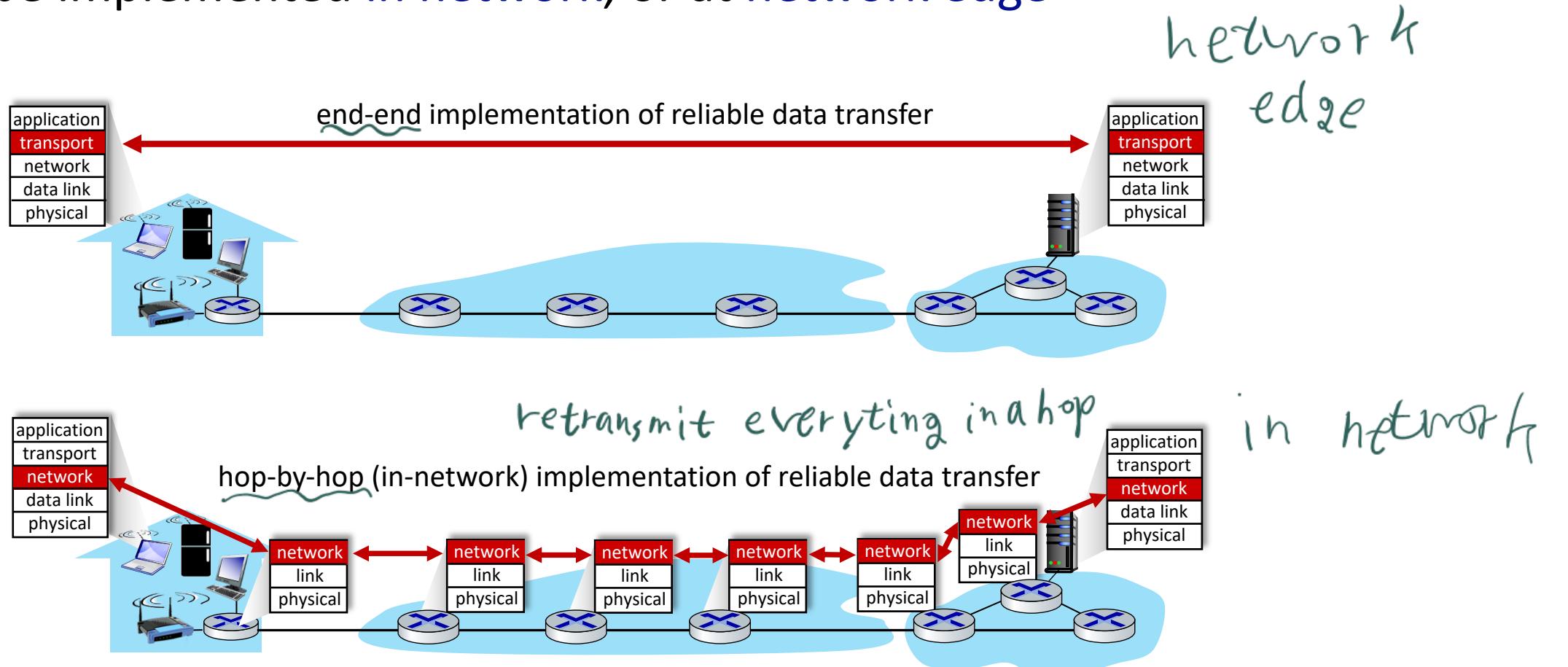
초석

Three cornerstone beliefs:

- simple connectivity → 간단한 연결
- IP protocol: that narrow waist → IP 프로토콜
- intelligence, complexity at network edge → edge 기기
지능적 복잡성

The end-end argument

- some network functionality (e.g., reliable data transfer, congestion) can be implemented in **network**, or at **network edge**



The end-end argument

- some network functionality (e.g., reliable data transfer, congestion) can be implemented in network, or at network edge

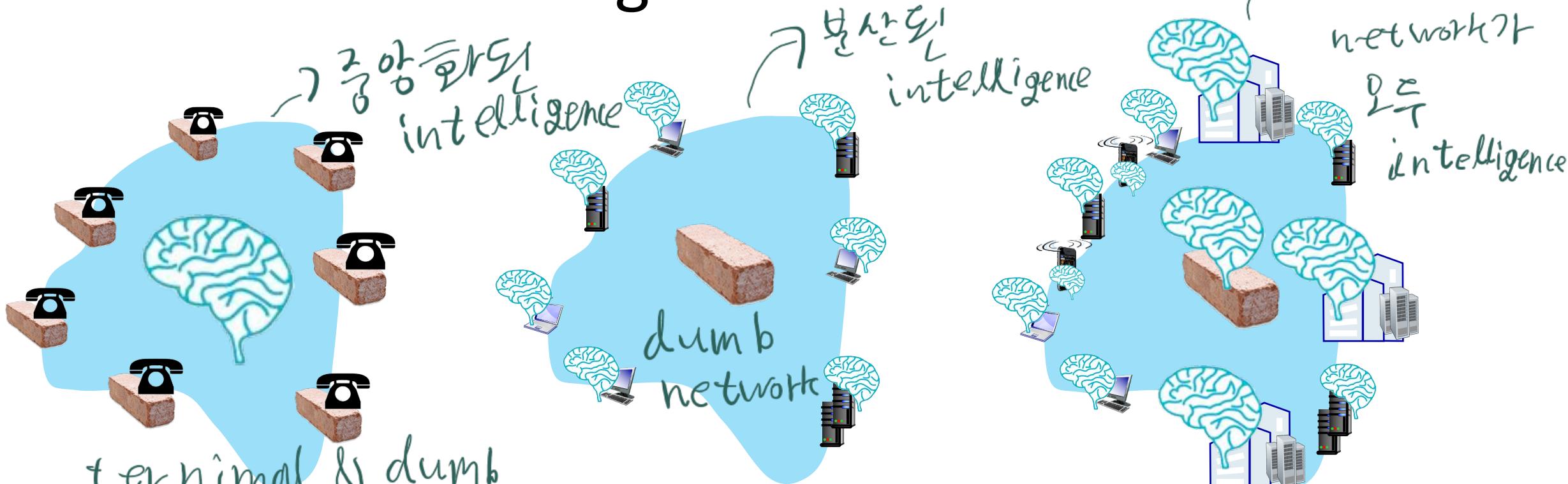
ex) lost 패킷
packetizing

“The function in question can completely and correctly be implemented only with the knowledge and help of the application standing at the end points of the communication system. Therefore, providing that questioned function as a feature of the communication system itself is not possible. (Sometimes an incomplete version of the function provided by the communication system may be useful as a performance enhancement.)

We call this line of reasoning against low-level function implementation the “end-to-end argument.”

문제의 기능은 통신 시스템의 종단점에 있는 애플리케이션의 지식과 도움으로만 완전하고 올바르게 구현될 수 있습니다. 따라서 문제의 기능을 통신 시스템 자체의 기능으로 제공하는 것은 불가능합니다. (때로는 통신 시스템이 제공하는 기능의 불완전한 버전이 성능 향상에 유용할 수 있습니다.) 저희는 저수준 기능 구현에 반대하는 이러한 추론을 "엔드투엔드 논증"이라고 부릅니다.

Where's the intelligence?



20th century phone net:

- intelligence/computing at network switches

Internet (pre-2005)

- intelligence, computing at edge

Internet (post-2005)

- programmable network devices
- intelligence, computing, massive application-level infrastructure at edge

인프라

Chapter 4: done!

- Network layer: overview
- What's inside a router
- IP: the Internet Protocol
- Generalized Forwarding, SDN
- Middleboxes



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

Answer: by the control plane (next chapter)

Network Layer: Data Plane

- Overview of Network Layer
- What's Inside a Router?
- The Internet Protocol: IPv4, Addressing, NAT IPv6
- Generalized Forwarding and SDN
- Middleboxes
 - middlebox functions
 - evolution, architectural principles of the Internet
- Summary

COMPSCI 453 **Computer Networks**

Professor Jim Kurose

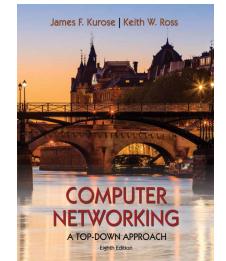
College of Information and Computer Sciences

University of Massachusetts



Class textbook:
*Computer Networking: A Top-
Down Approach (8th ed.)*

J.F. Kurose, K.W. Ross
Pearson, 2020
http://gaia.cs.umass.edu/kurose_ross



Next...

- *Anonymity Networks: Tor, Onion networks, Darknet*