

# Chapter 4

## Network Layer:

## The Data Plane

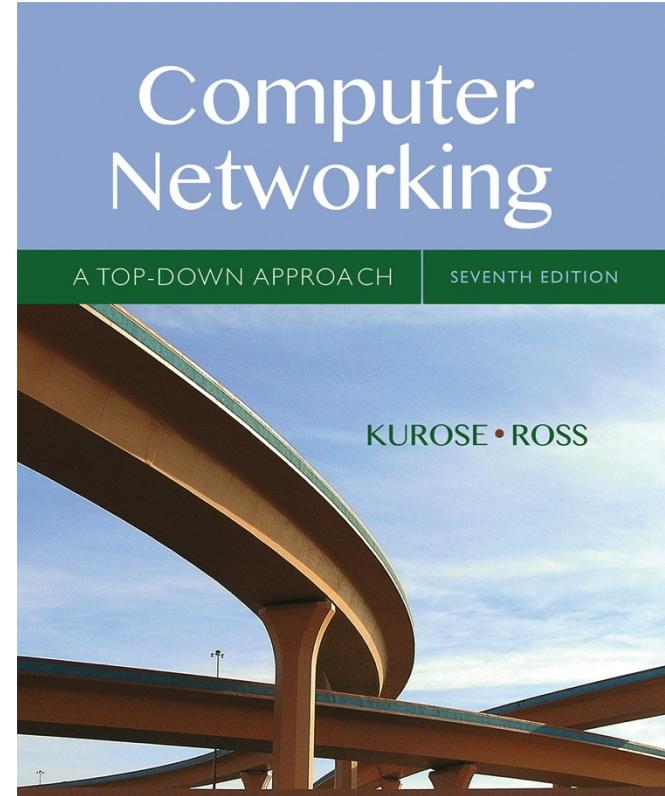
### A note on the use of these Powerpoint slides:

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 (NAT)
- IPv6

- DATA PLANE: parte della rete che si occupa di smistare i dati
- CONTROL PLANE: " definisce algoritmi che calcolano i cammini di trasferimento dei pacchetti

## 4.4 Generalized Forward and SDN (Software Define Network)

- match  
→ gestione dei pacchetti a livello software
- action  
→ il pacchetto lo indirizza il software
- OpenFlow examples of match-plus-action in action

LIVELLO RETE

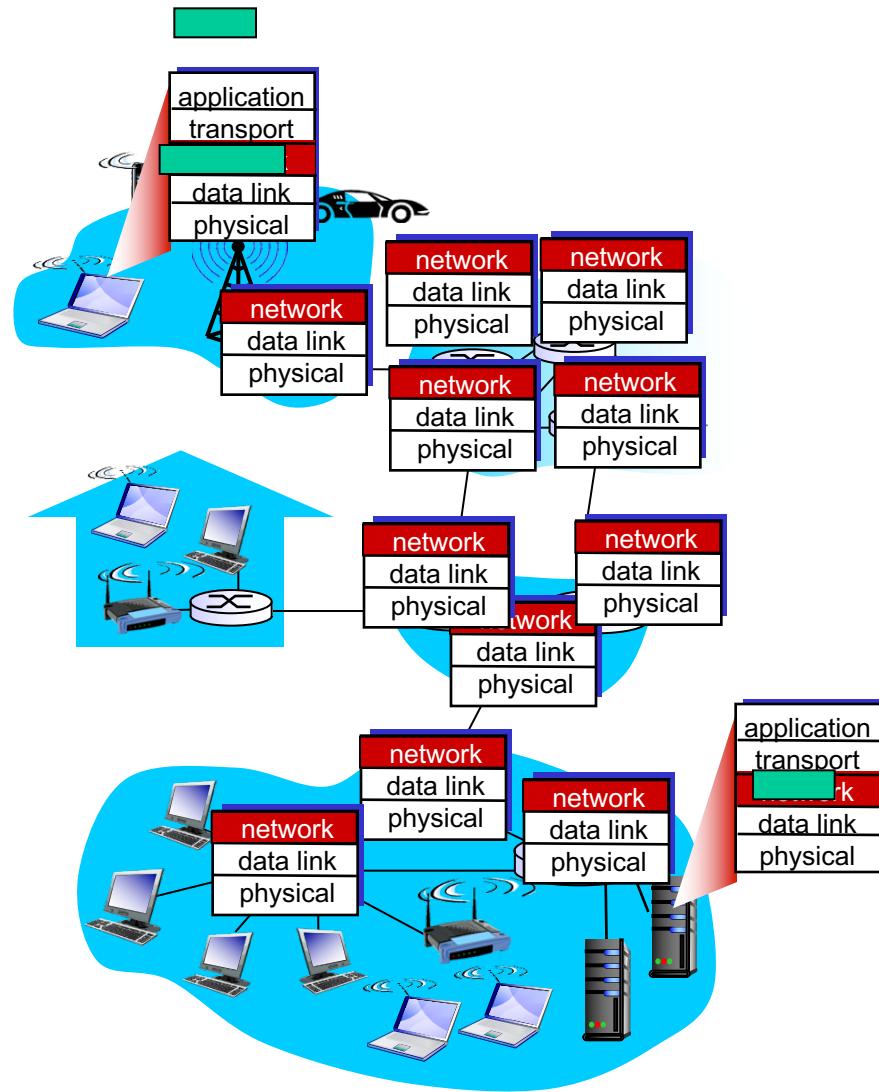
# Chapter 4: network layer

*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

- 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*

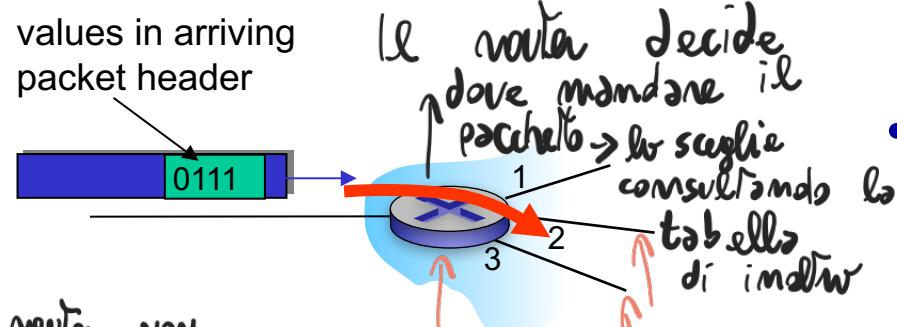
*analogy: taking a trip*

- **forwarding:** process of getting through single interchange
  - **routing:** process of planning trip from source to destination
- [operazione del calcolo delle Tabelle di inoltro → calcolo del percorso]*

# 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



N.B. Il router non  
sempre ha  
la tabella di  
indirizzamento

## 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

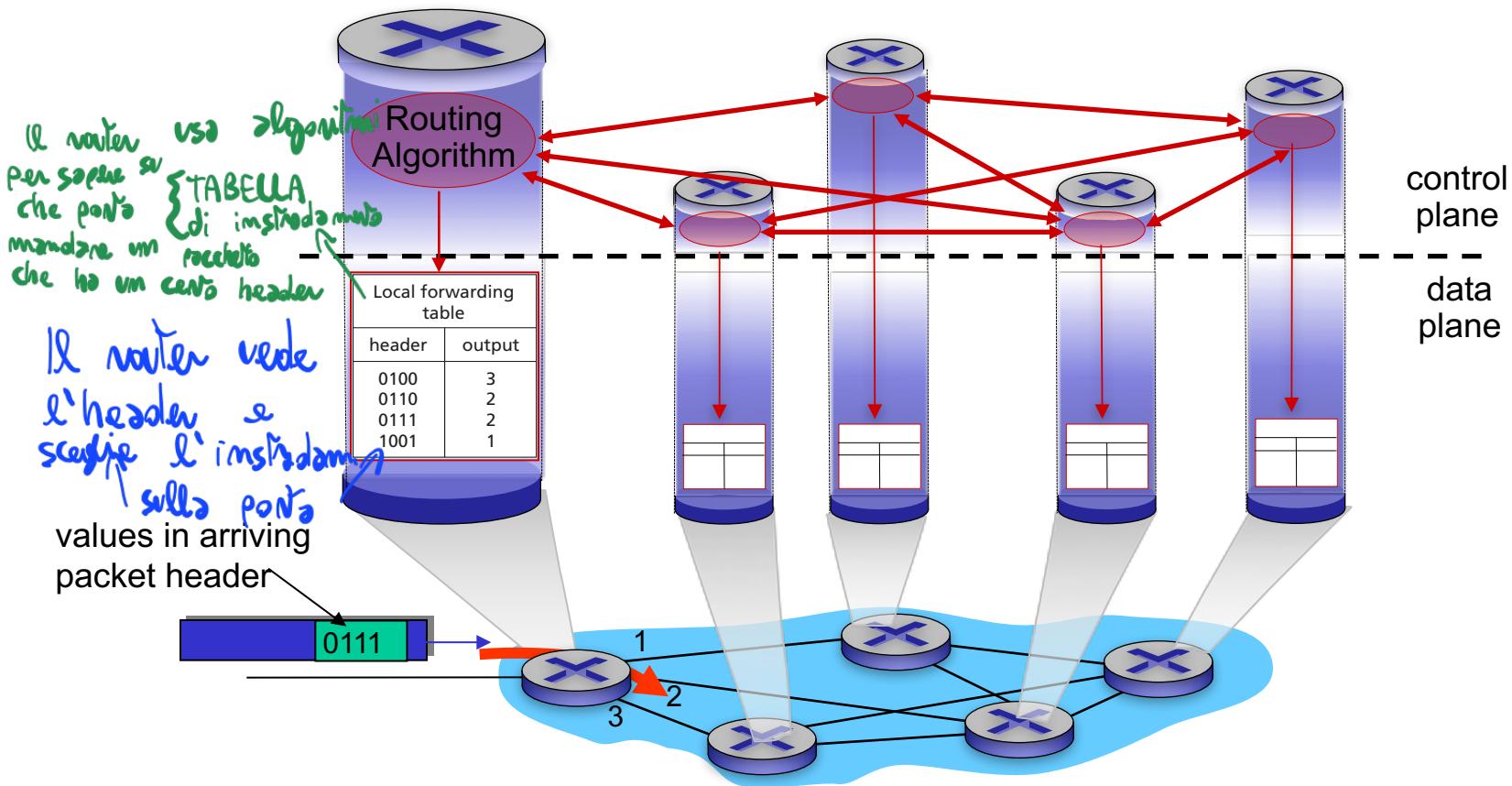
I router possono avere più porte di uscita a seconda di quante interfacce di rete sono implementate

# Per-router control plane

1º modo

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

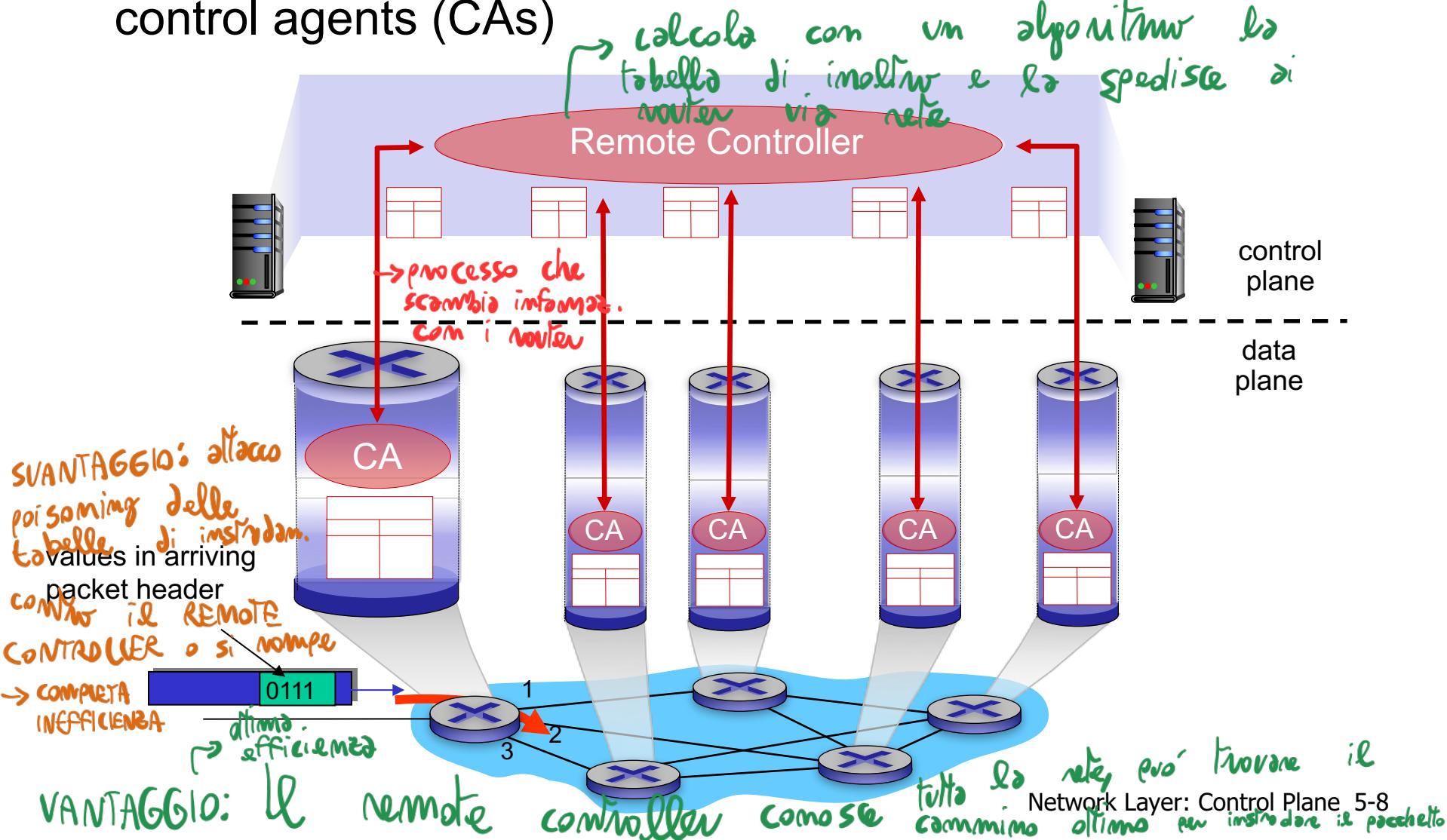
→ trova applicazione in contesti di reti dinamiche (reti veicolari)



# Logically centralized control plane

2° modo : remote controller

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



# 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 service models:

Network Architecture	Service Model	Bandwidth	Loss	Order	Timing	Congestion feedback
F, schifo, posso fare meglio!	Internet	best effort <i>no garanzie di banda, di ordine, di ritardo,</i>	none <i>pendola pacchetti</i>	no <i>ordine dei pacchetti</i>	no <i>arrivo dei pacchetti in un certo tempo</i>	no (inferred via loss)
Asynchronous transfer model	ATM CBR	constant bit rate → flusso	constant bit rate <i>costante di bit usata per la voce</i>	yes	yes	yes
desire alcuni bit molte	ATM VBR	VARIABLE BIT RATE <i>per i video</i>	guaranteed rate	yes	yes	no congestion
non c'è detto che i bit arrivino in modo sincrono	ATM ABR	ASYNCHRONOUS BIT RATE	guaranteed minimum	no	yes	no congestion
simmetria	ATM UBR	none	no	yes	no	no
<p>Internet usa un buffer per garantire la consegna di file multimediali (che con ATM è garantita) pacchetti in anticipo rispetto allo stato nel buffer, che viene sovraccarico e riempito</p>						

# 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

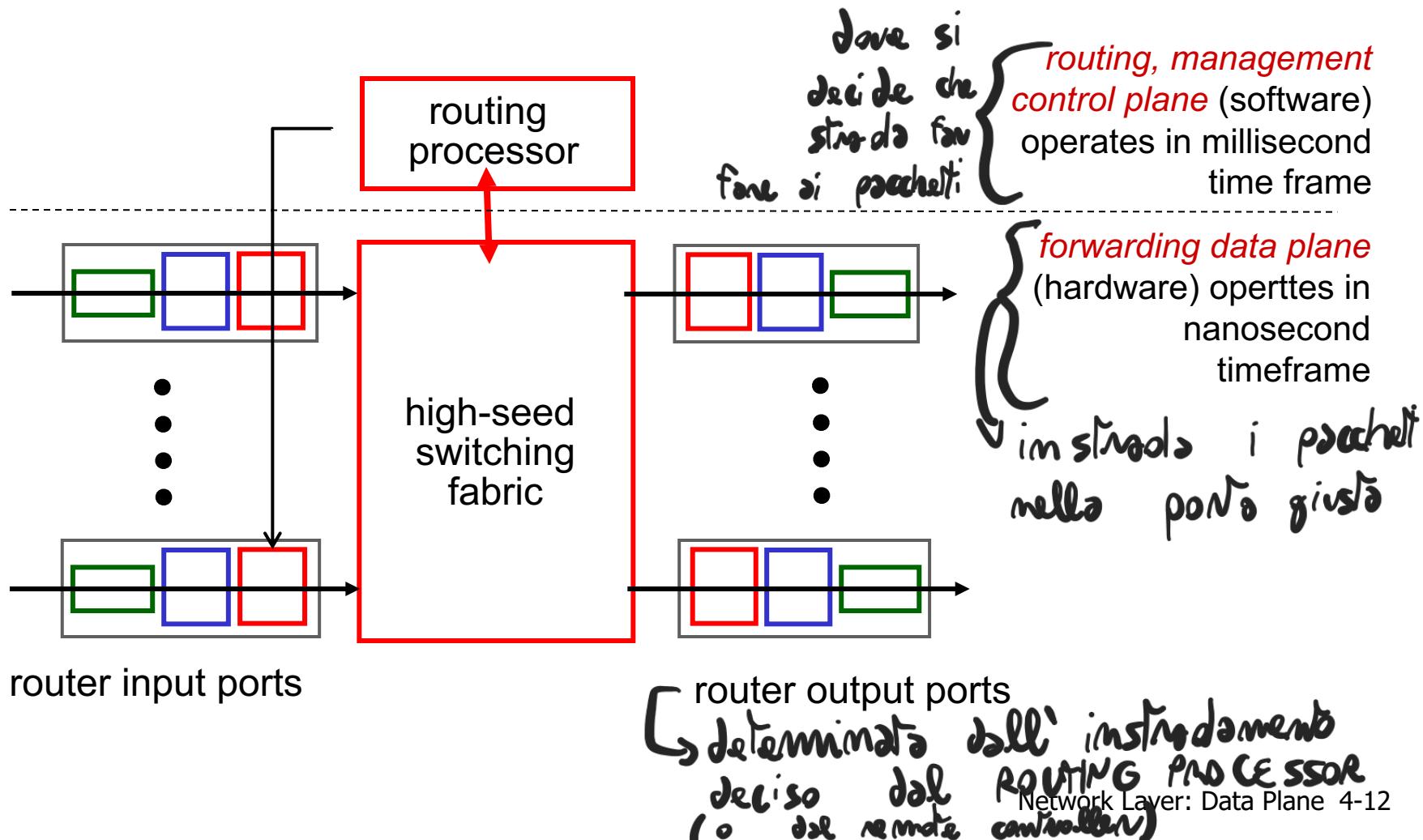
- datagram format
- fragmentation
- IPv4 addressing
- network address translation
- IPv6

## 4.4 Generalized Forward and SDN

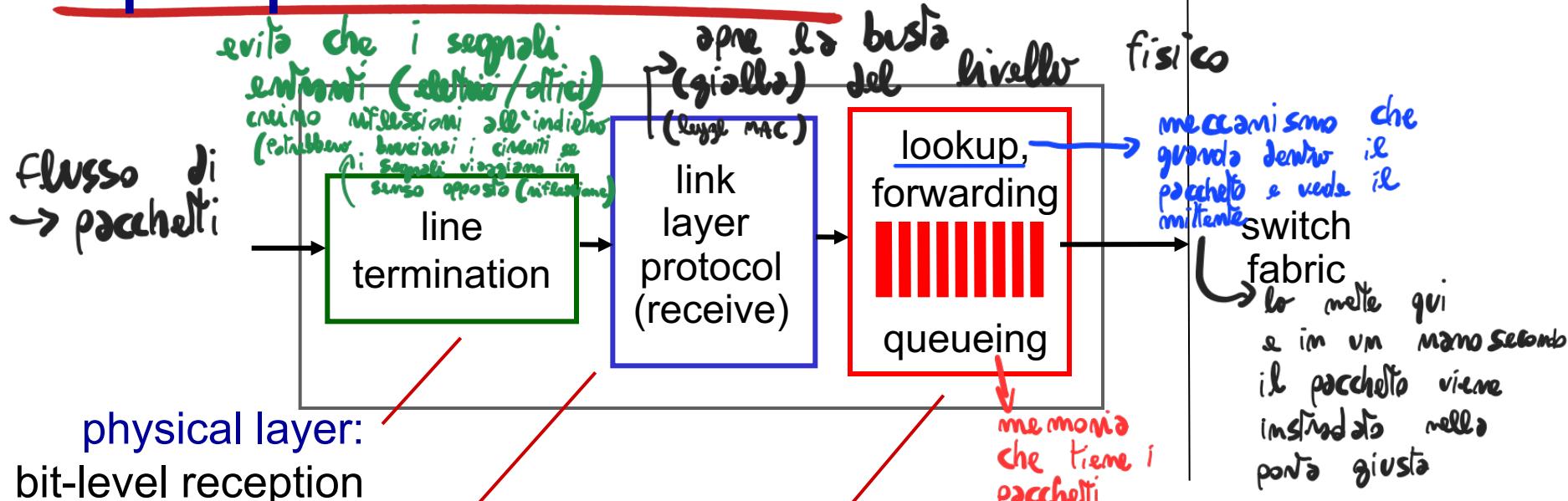
- match
- action
- OpenFlow examples of match-plus-action in action

# Router architecture overview

- high-level view of generic router architecture:



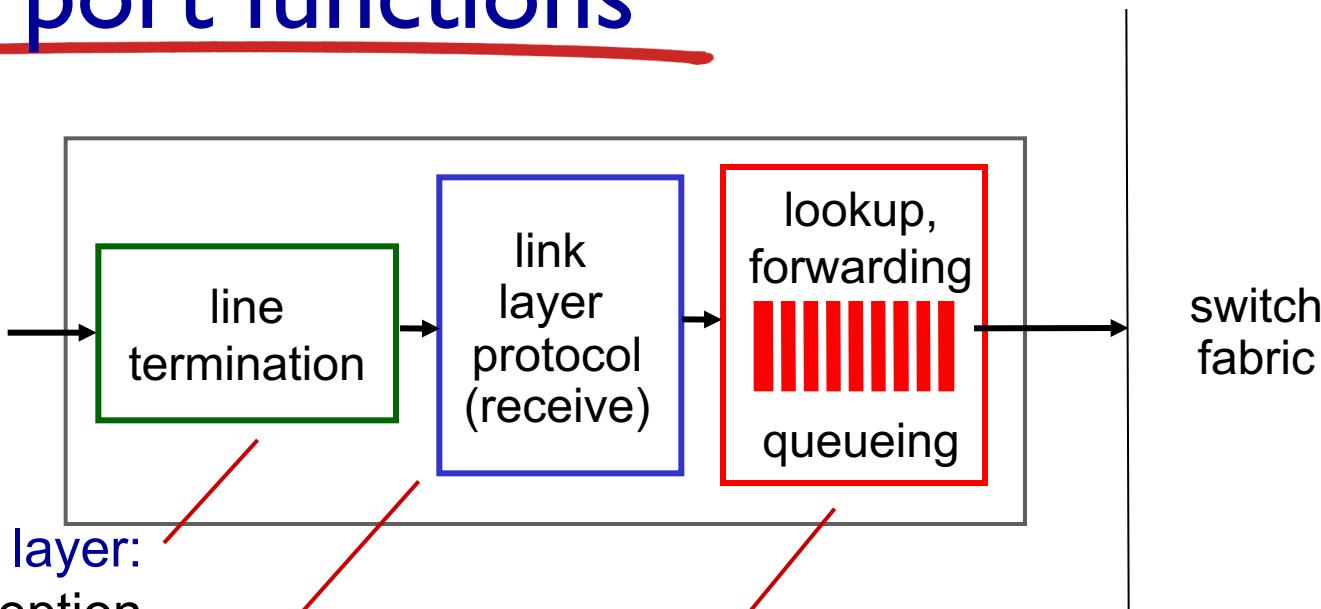
# Input port functions



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

# 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

# Destination-based forwarding

Questa tabella di indirizzamento viene scelta dall' algoritmo di indirizzamento **forwarding table**

Destination Address Range	Link Interface
11001000 00010111 00010000 00000000 through 11001000 00010111 00010111 11111111 indirizzati sulla porta 0	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.

TABELLA DI INOLTRO

Destination Address Range <i>indirizzo</i>	Link interface
11001000 00010111 00010*** ******	0
11001000 00010111 00011000 *	1
11001000 00010111 00011*** ******	2
otherwise	3

prefissi:  
...  
...  
...  
...

examples:

DA: 11001000 00010111 00010110 10100001

DA: 11001000 00010111 00011000 10101010

which interface?  $\Rightarrow 0$

which interface?  $1 \text{ o } 2?$

$\Rightarrow 2$   
(prefisso più lungo matched)

# Longest prefix matching

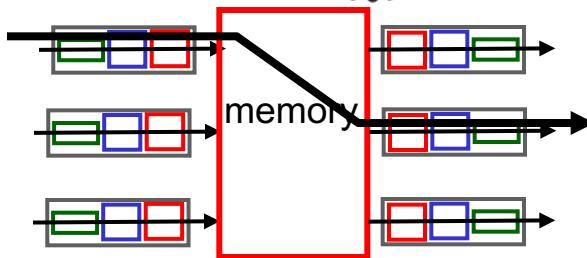
- we'll see *why* longest prefix matching is used shortly, when we study addressing
- longest prefix matching: often performed using ternary content addressable memories (TCAMs)
  - *content addressable*: present address to TCAM: retrieve address in one clock cycle, regardless of table size
  - Cisco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

# Switching fabrics

richiede un ciclo di clock  
per

↳ blocco che trasferisce da INPUT ad OUTPUT

- transfer packet from input buffer to appropriate output buffer
- 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 types of switching fabrics → come le realizzo?  
*voglio : velocità  
costo : basso*

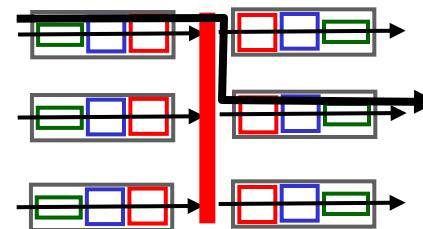


memory

si salva in una RAM  
il pacchetto

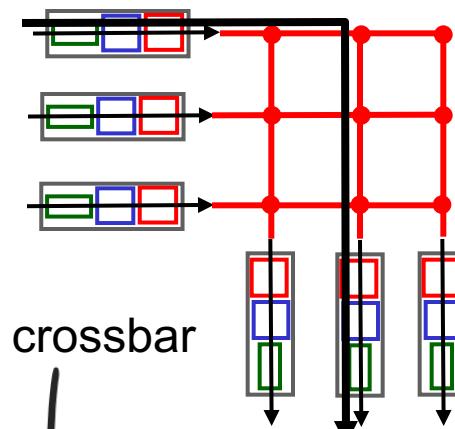
Svantaggio: ciclo di clock

aggiuntivo per leggere la memoria



bus

Svantaggio: puo' gestire un solo  
flusso alla volta



crossbar

→ si può decidere il percorso  
di più flussi → cammini paralleli  
Network Layer: Data Plane 4-18

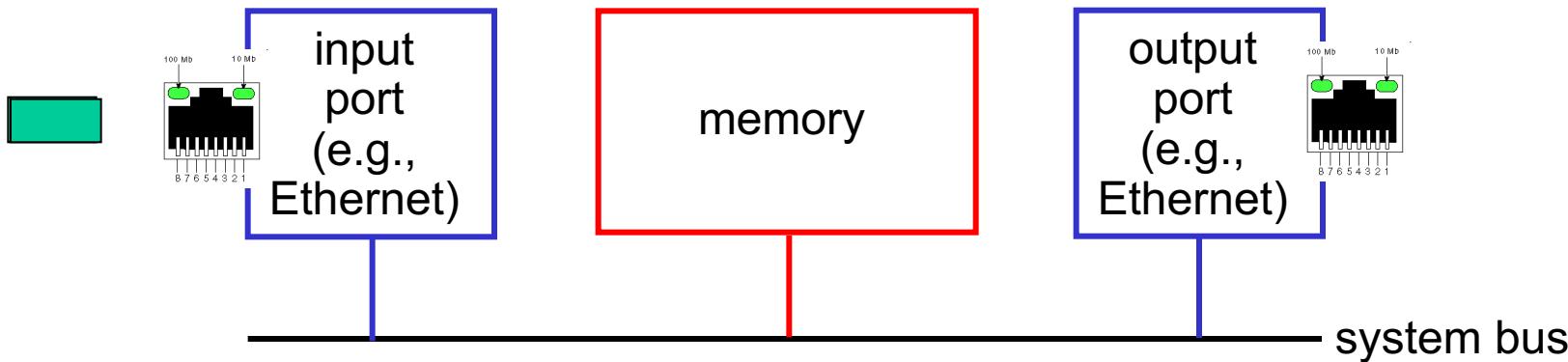
Vantaggio: cammini paralleli → più flussi

RETE SU CHIP: Architettura di rete che spostano dati da una scheda madre ad un'altra

## Switching via memory

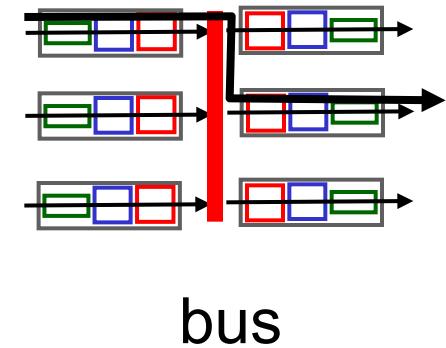
*first generation routers:*

- 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)



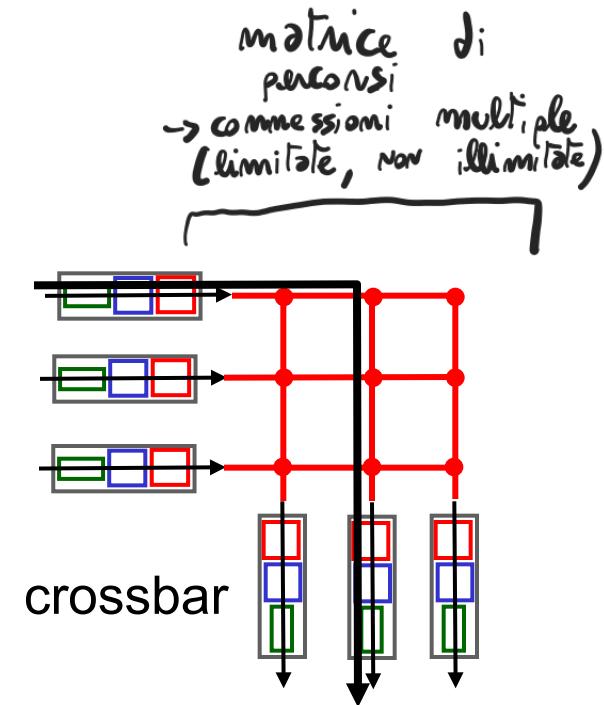
# Switching via a bus

- 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



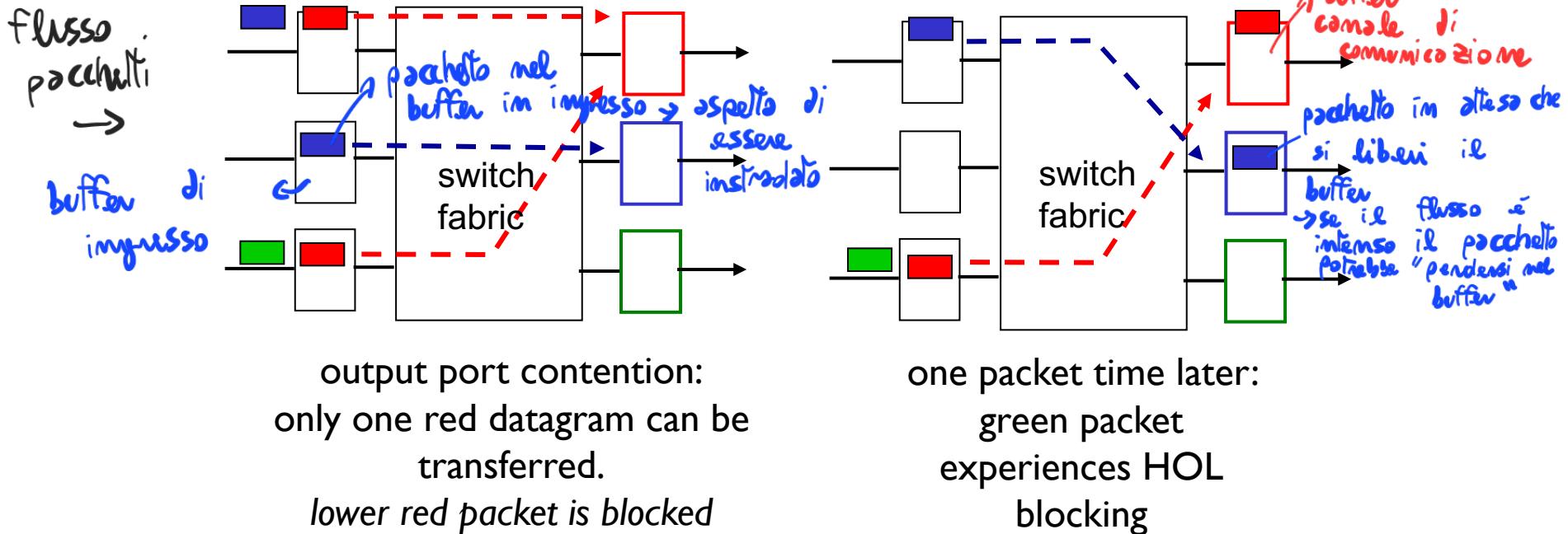
# Switching via interconnection network

- 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



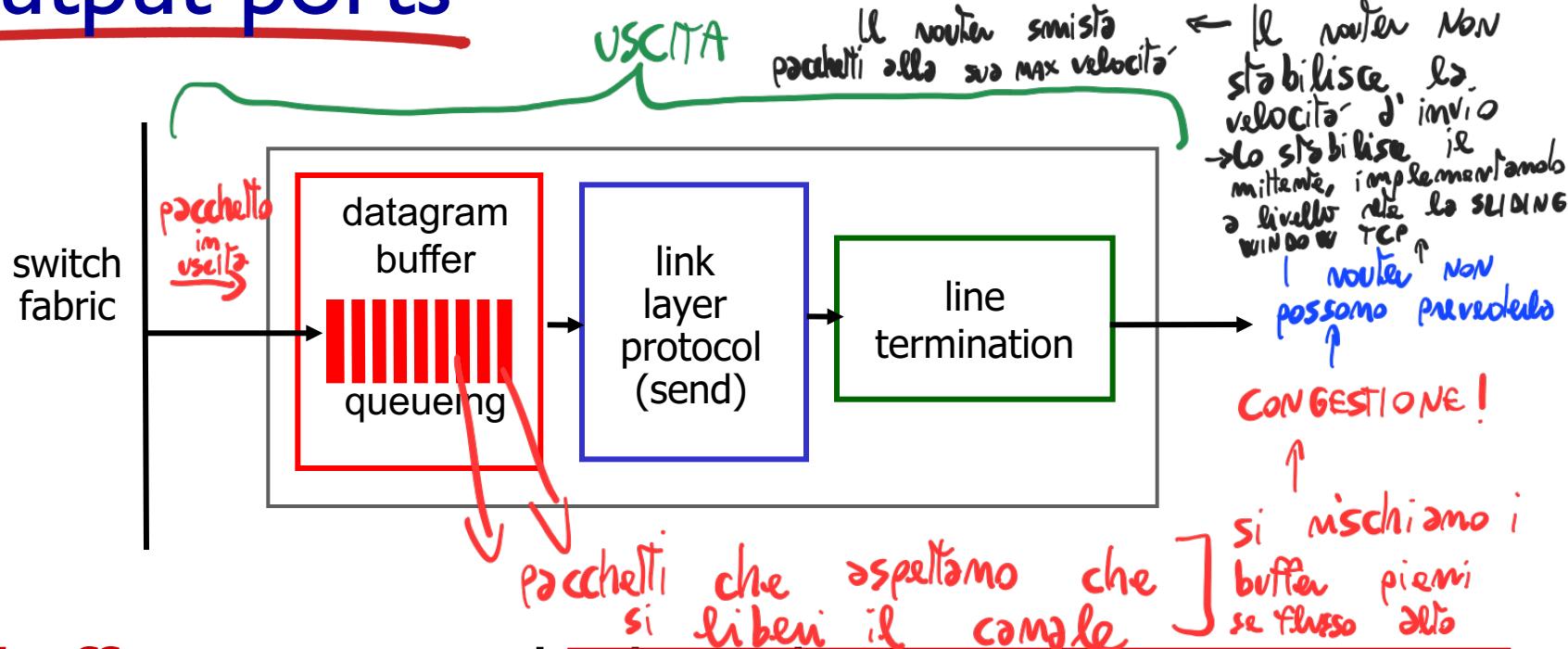
# Input port queuing

- 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



# Output ports

*This slide is HUGELY important!*

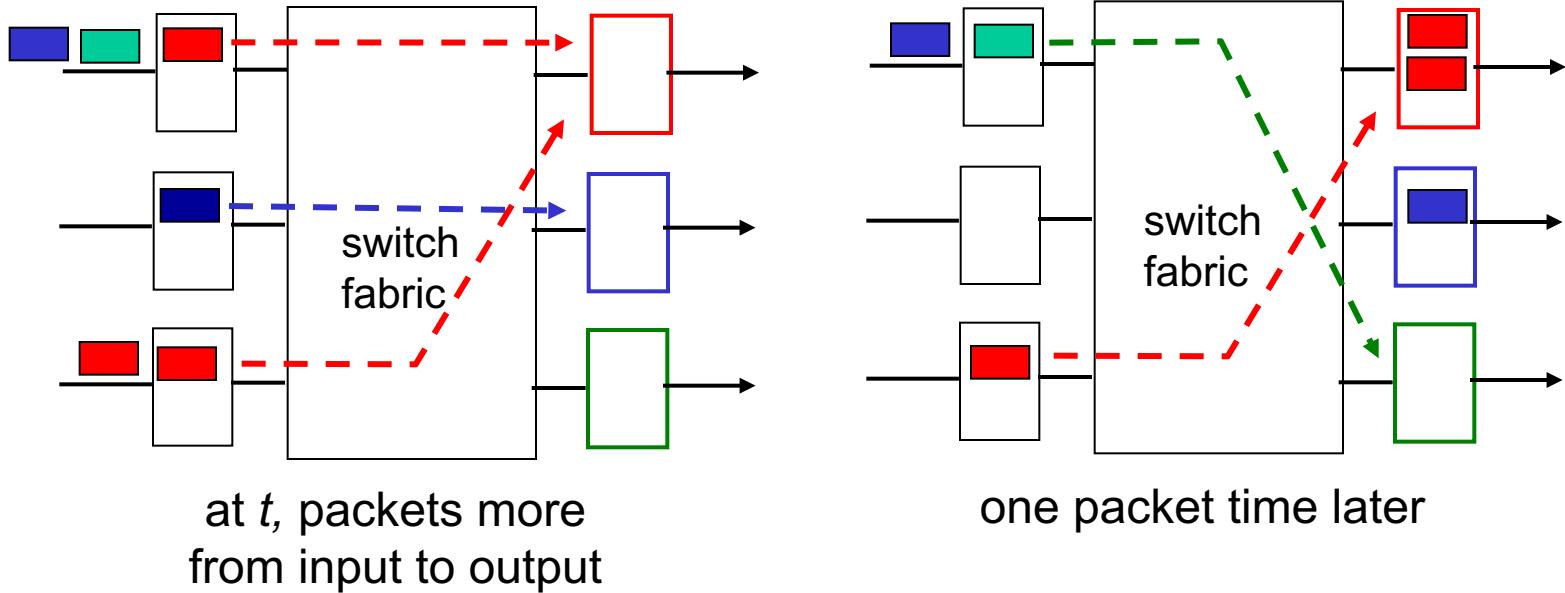


- **buffering** required from fabric faster rate
- **scheduling** datagrams

Datagram (packets) can be lost due to congestion, lack of buffers

Priority scheduling – who gets best performance, network neutrality

# Output port queueing



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

→ NON LA CHIEDE questa slide

# How much buffering? → slide filosofica

Quanto spazio dedicare ai buffer di uscita sulle porte dei router?

- RFC 3439 rule of thumb: average buffering equal to “typical” RTT (say 250 msec) times link

capacity  $C$

• e.g.,  $C = 10 \text{ Gpbs}$  link: 2.5 Gbit buffer

- recent recommendation: with  $N$  flows, buffering equal to

quantità del  
buffer di  
memoria

(minima  
necessaria)

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

→  $N^2$  punti di controllo  
(flussi massimi che passano)  
 $\sqrt{N^2} = N$

statisticamente

ne passano  $\sqrt{N}$   
(è una stima)

Se costa

meno → buffer piccoli  
servizio più lento

Se ne do'  
di più → costa di più

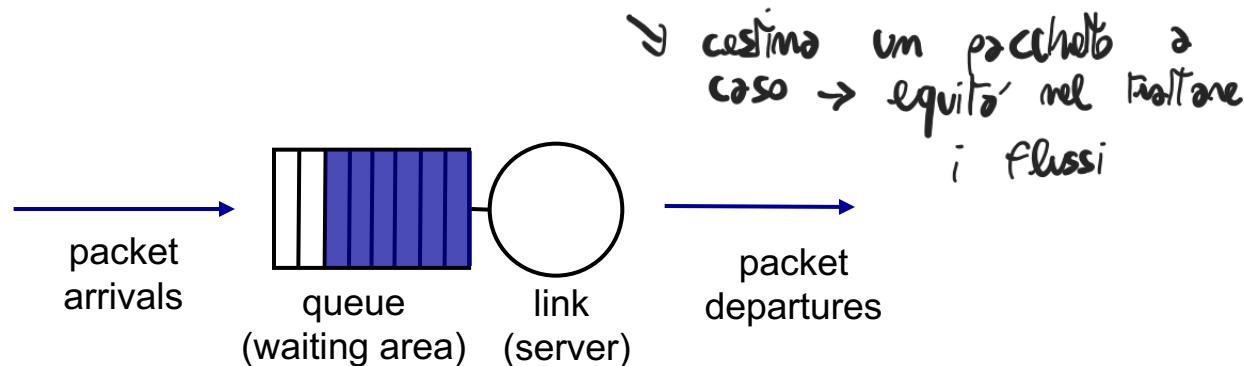
→ oppure ho il buffer pieno, devo togliere un pacchetto chi scelgo? Il router non tiene conto delle priorità dei pacchetti

## Scheduling mechanisms

→ Se ci sono dei pacchetti in uscita che attendono e ho pacchetti meno importanti (mail) e più importanti ("pezzo di voce")

- **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 → elimina l'ultimo pacchetto arrivato (se pieno).
  - **priority:** drop/remove on priority basis → mettezza di priority che decide chi cestimare
  - **random:** drop/remove randomly



NON è il caso né di internet, né dei router → ma nei GARR

## Scheduling policies: priority

I router gestiscono i pacchetti secondo

**priority scheduling:** send livelli di priorità

highest priority queued packet

- multiple classes, with different priorities

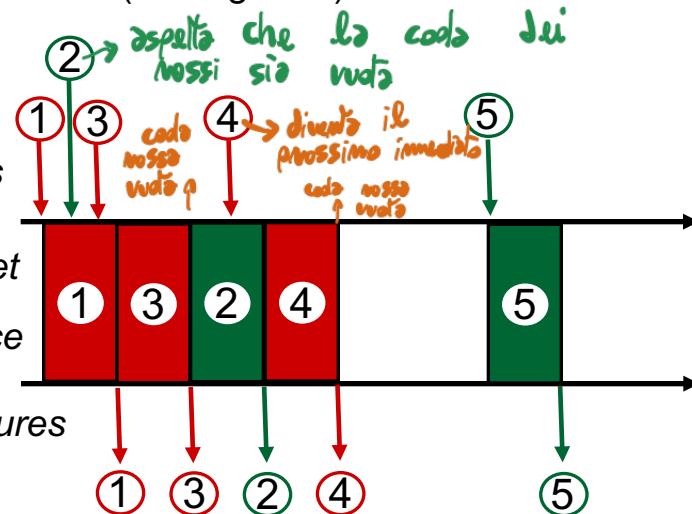
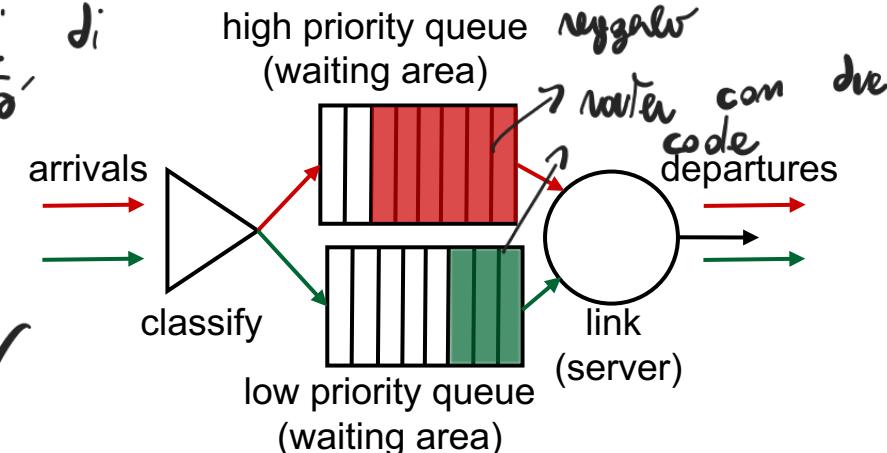
- class may depend on marking or other header info, e.g. IP source/dest, port numbers, etc.

- real world example?

Se i buffer sono pieni

→ cestina quelli con meno priorità

COSTANO MOLTISSIMO  
· i best effort di internet mette a segnale

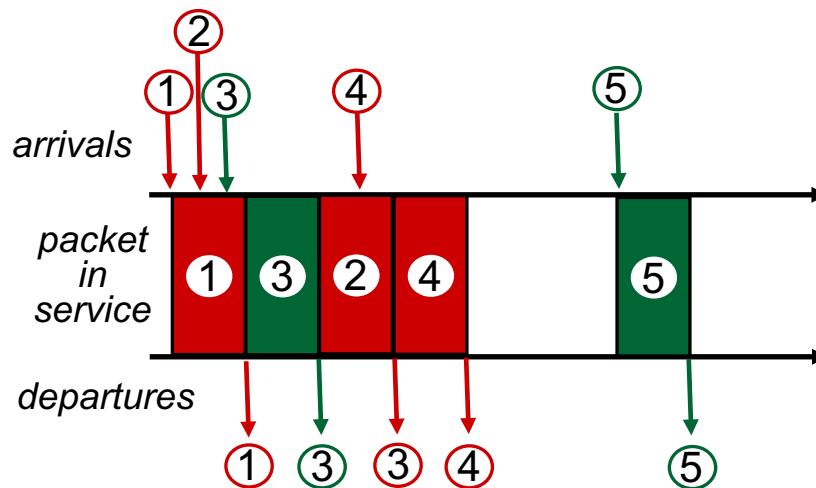


Network Layer: Data Plane 4-27

# Scheduling policies: still more

## *Round Robin (RR) scheduling:*

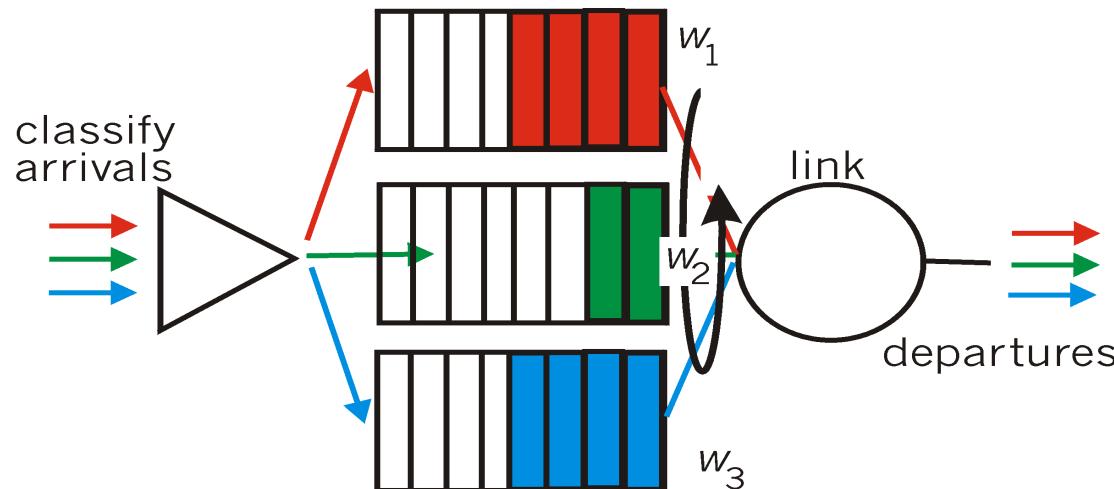
- multiple classes
- cyclically scan class queues, sending one complete packet from each class (if available)
- real world example?



# 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

## 4.1 Overview of Network layer

- data plane
- control plane

## 4.2 What's inside a router

## 4.3 IP: Internet Protocol

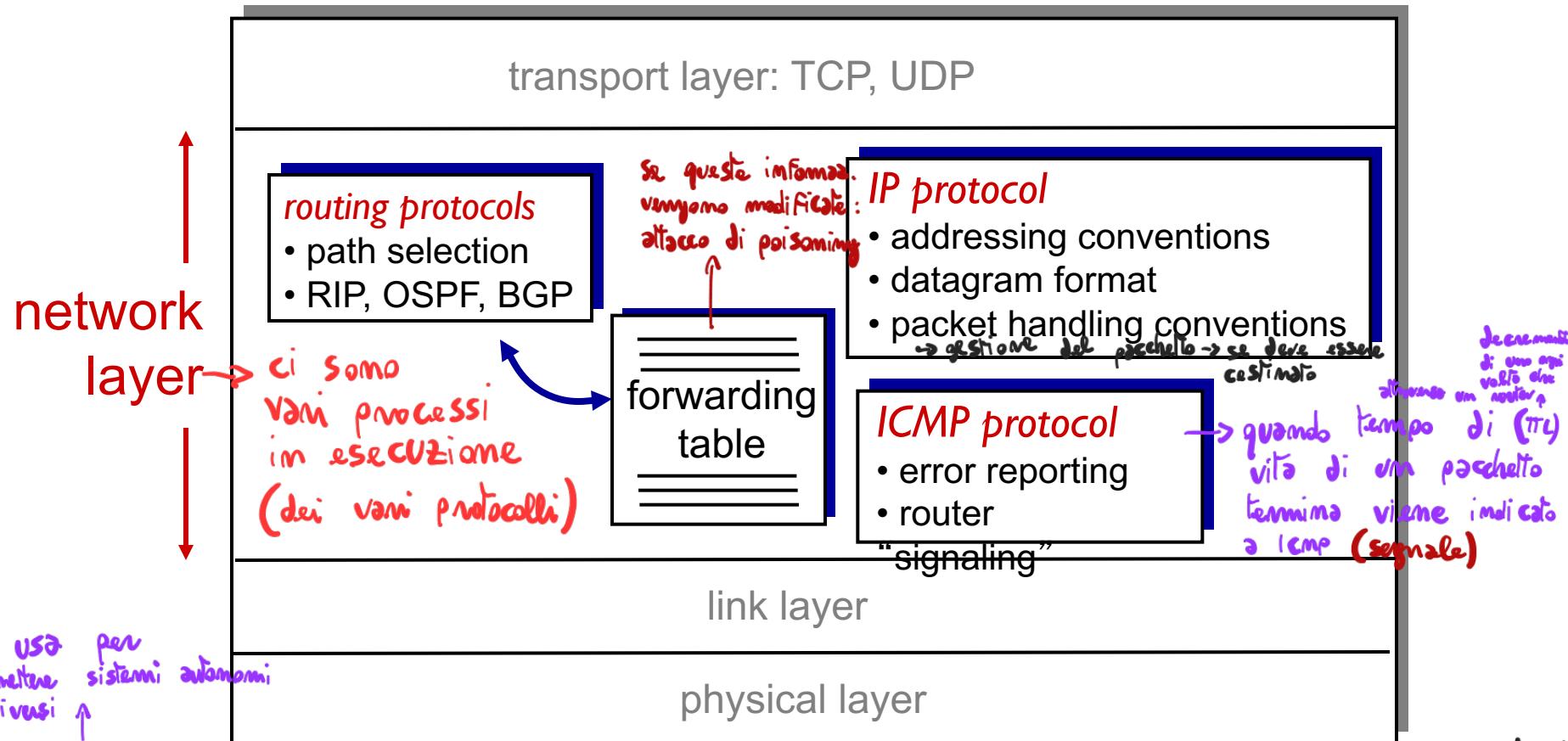
- datagram format
- fragmentation
- IPv4 addressing
- network address translation
- IPv6

## 4.4 Generalized Forward and SDN

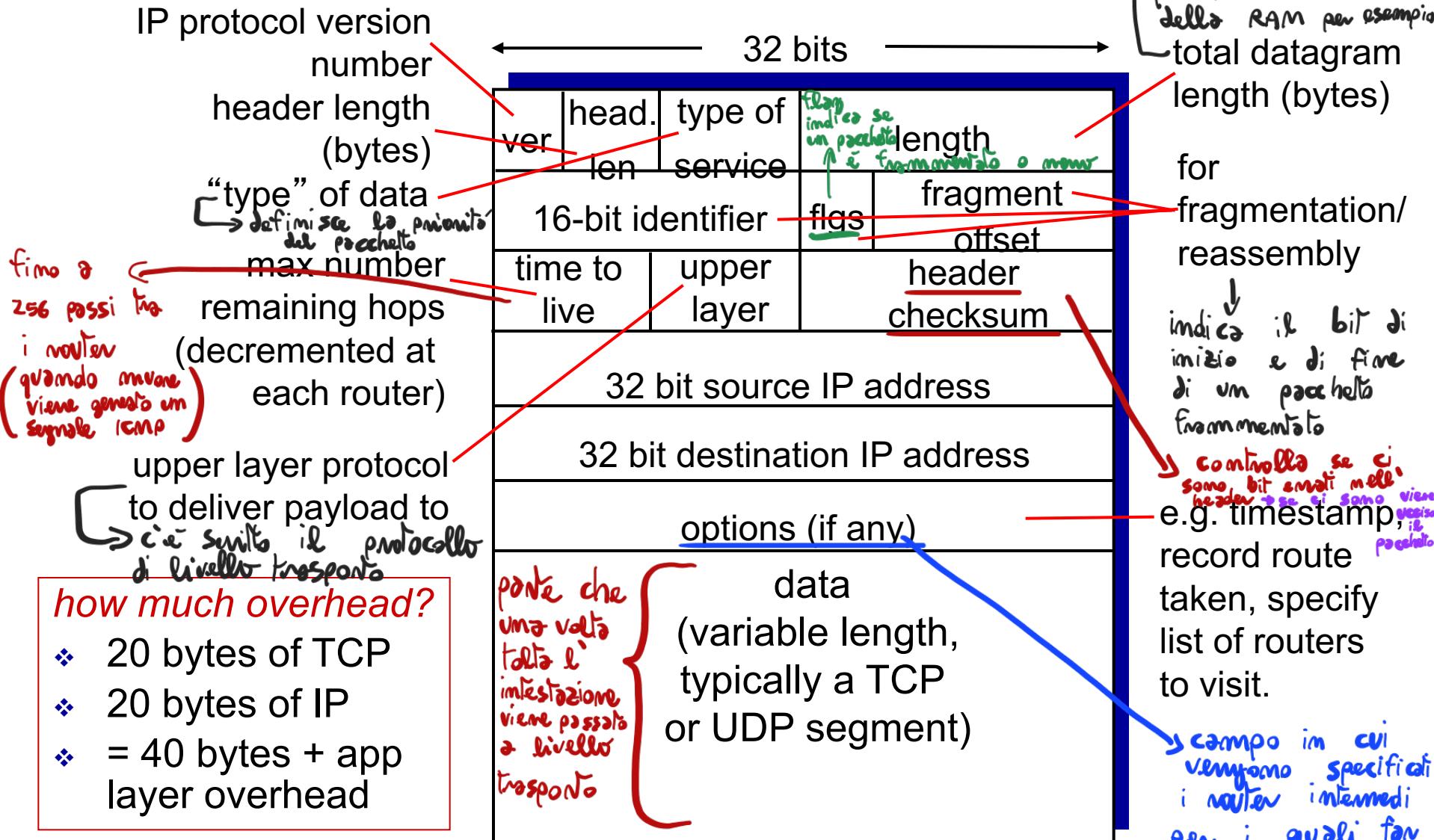
- match
- action
- OpenFlow examples of match-plus-action in action

# The Internet network layer

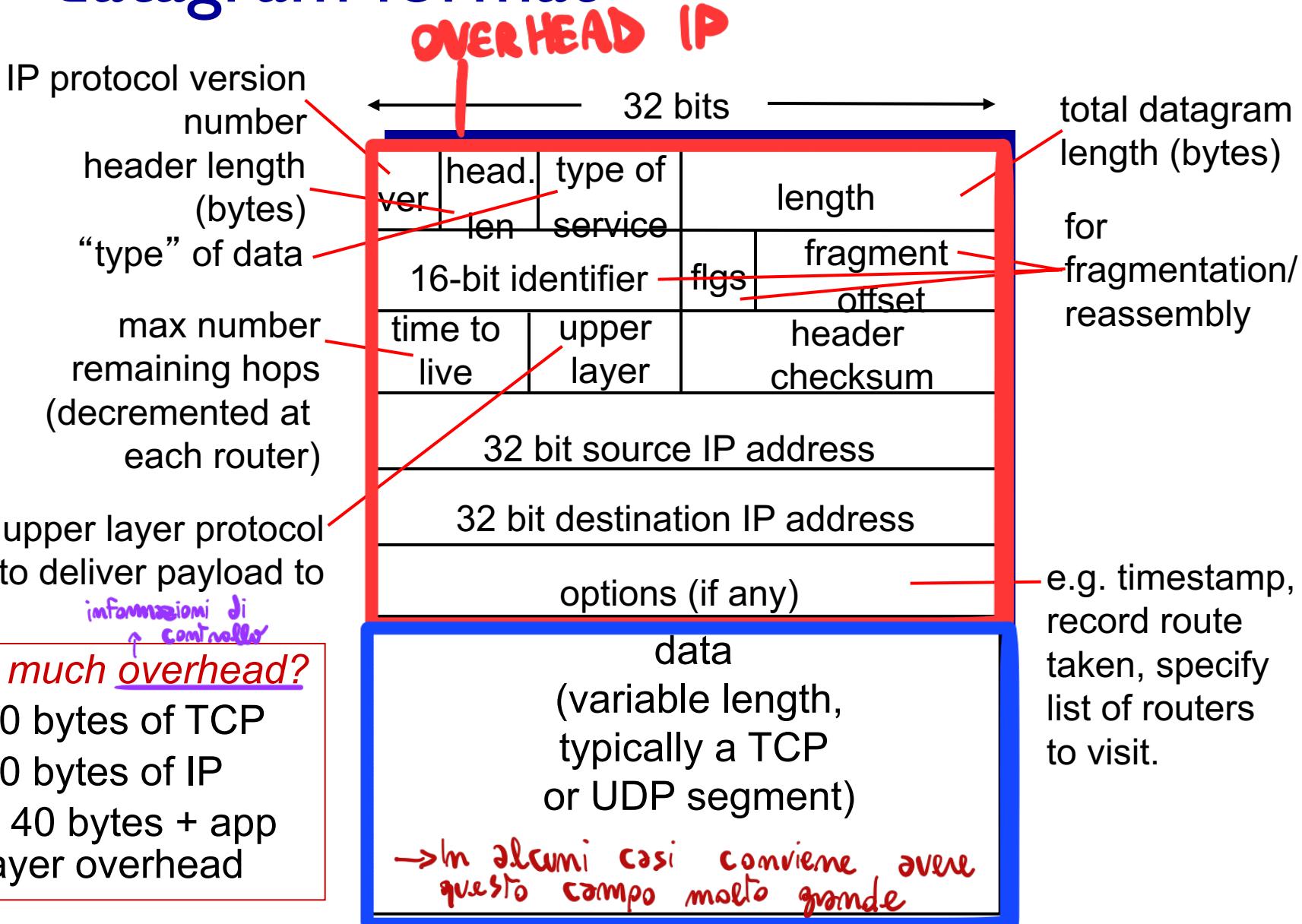
host, router network layer functions:



# IP datagram format

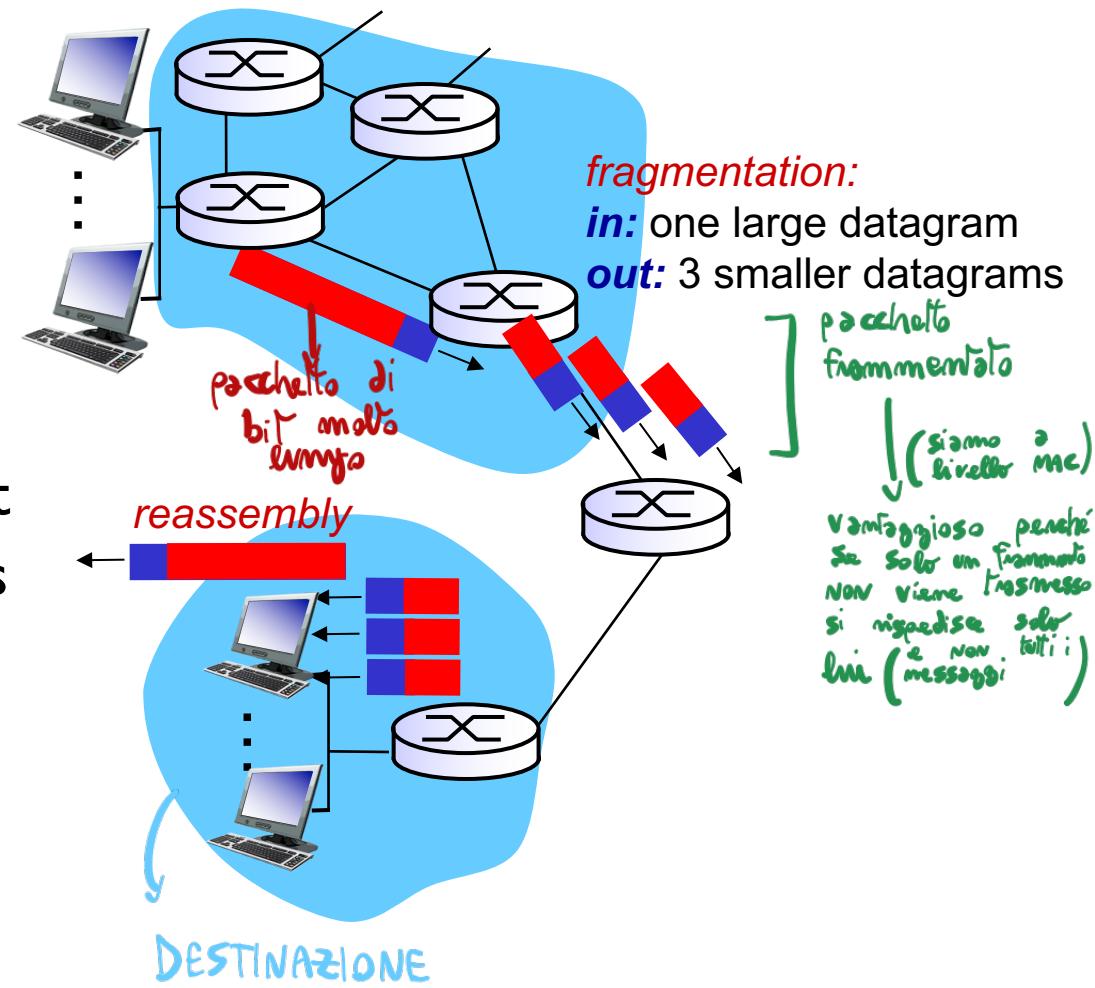


# IP datagram format



# 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

→ Lo fa il router che frammenta  
 (È il router a destinazione che riassumba  
 il messaggio)

## example:

- ❖ 4000 byte datagram → dimensione pacchetto
- ❖ MTU = 1500 bytes

Maximum Transfer Unit → indica la dimensione massima del pacchetto da trasmettere per non

avere errori di trasmissione

1480 bytes in data field → + 20 bytes di overhead

$$\text{offset} = \\ 1480/8$$

	length	ID	<u>fragflag</u>	offset	
	=4000	=x	=0	=0	

fragmentation flag

one large datagram becomes several smaller datagrams  
 → 3 frammenti

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

primo byte di dati del pacchetto

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

// ID del pacchetto "grande"

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

primo byte di questo pacchetto

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

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

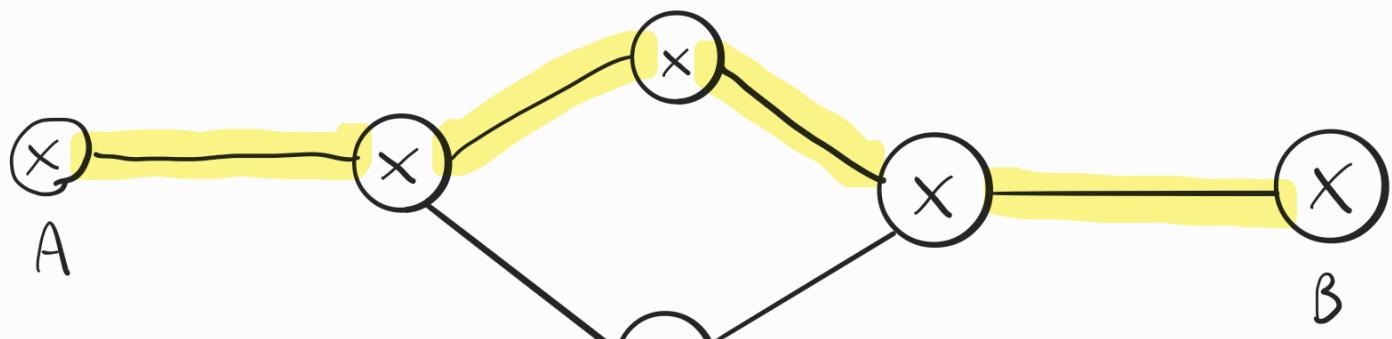
primo byte

indica che dopo questo frammento non deve più esserci l'ultimo

## PROTO COLLO ATM:

- Protocollo alternativo a IP
- La dimensione del campo dati è COSTANTE (In IP è variabile) → 48 byte (dovuto a dispute fra vari Paesi)
- Nell'header ci sono 5 byte → sono un'etichetta di flusso

- C'è una tabella che mappa mittente e destinatario con una ETICHETTA di FLUSSO riservata / "prenotata" prima delle risorse di rete
  - Viene definito primo il cammino di instradamento dei pacchetti riservando queste quantità di risorse di rete prima di inviare i pacchetti
- Vantaggi:  
possibile perché riserva prima le risorse necessarie
- { → posso stabilire un numero minimo di bit (riserva prima le risorse)  
→ Definisco un tempo di arrivo dei pacchetti.



Ricorda i circuiti a commutazione

(INTERNET NON FUNZIONA COSÌ!!)

- I router decidono l'instradamento volto per volto
- È un servizio best effort ("fa il meglio che può")

# 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

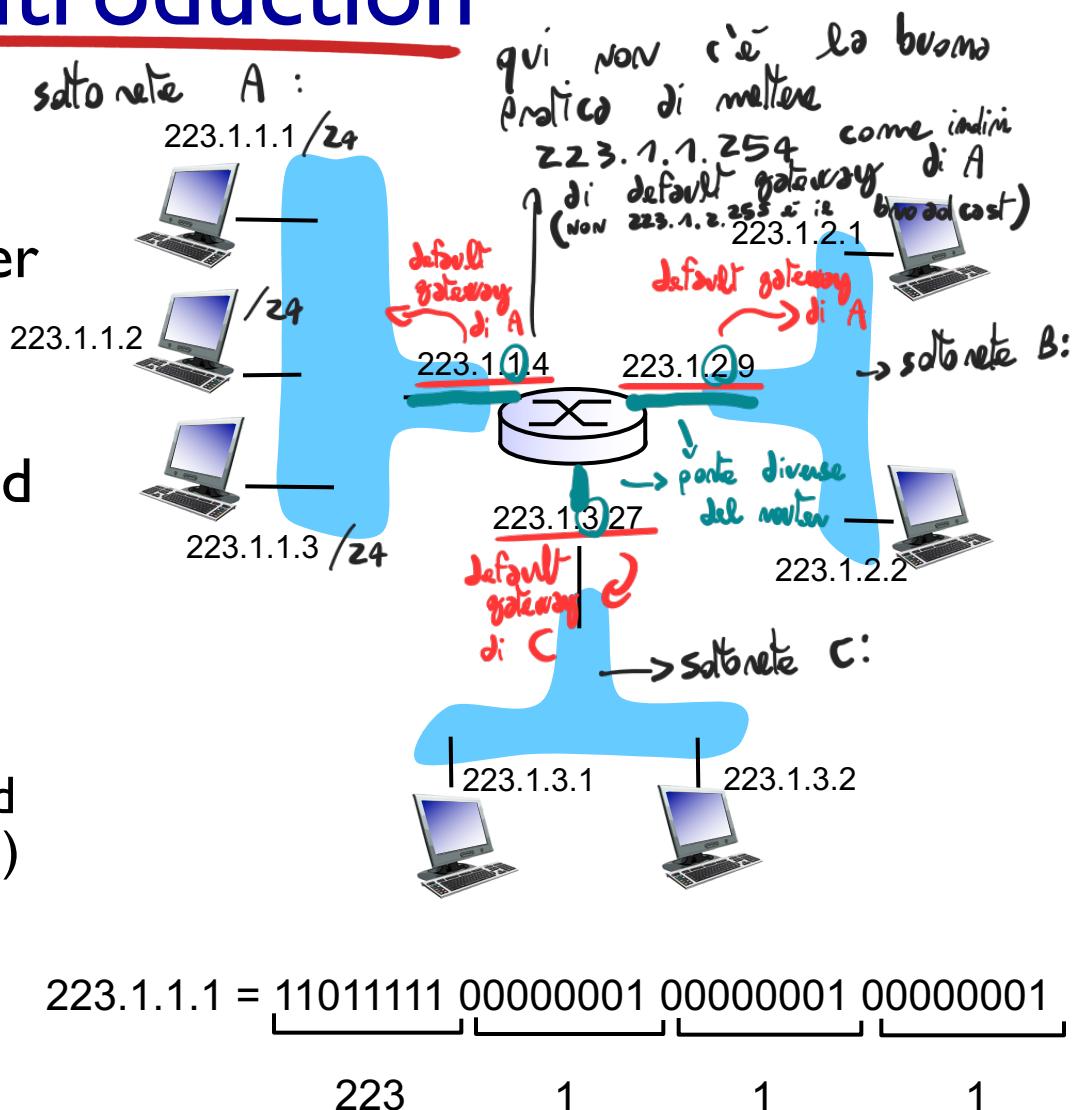
- datagram format
- fragmentation
- IPv4 addressing
- network address translation
- IPv6

## 4.4 Generalized Forward and SDN

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

# 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**



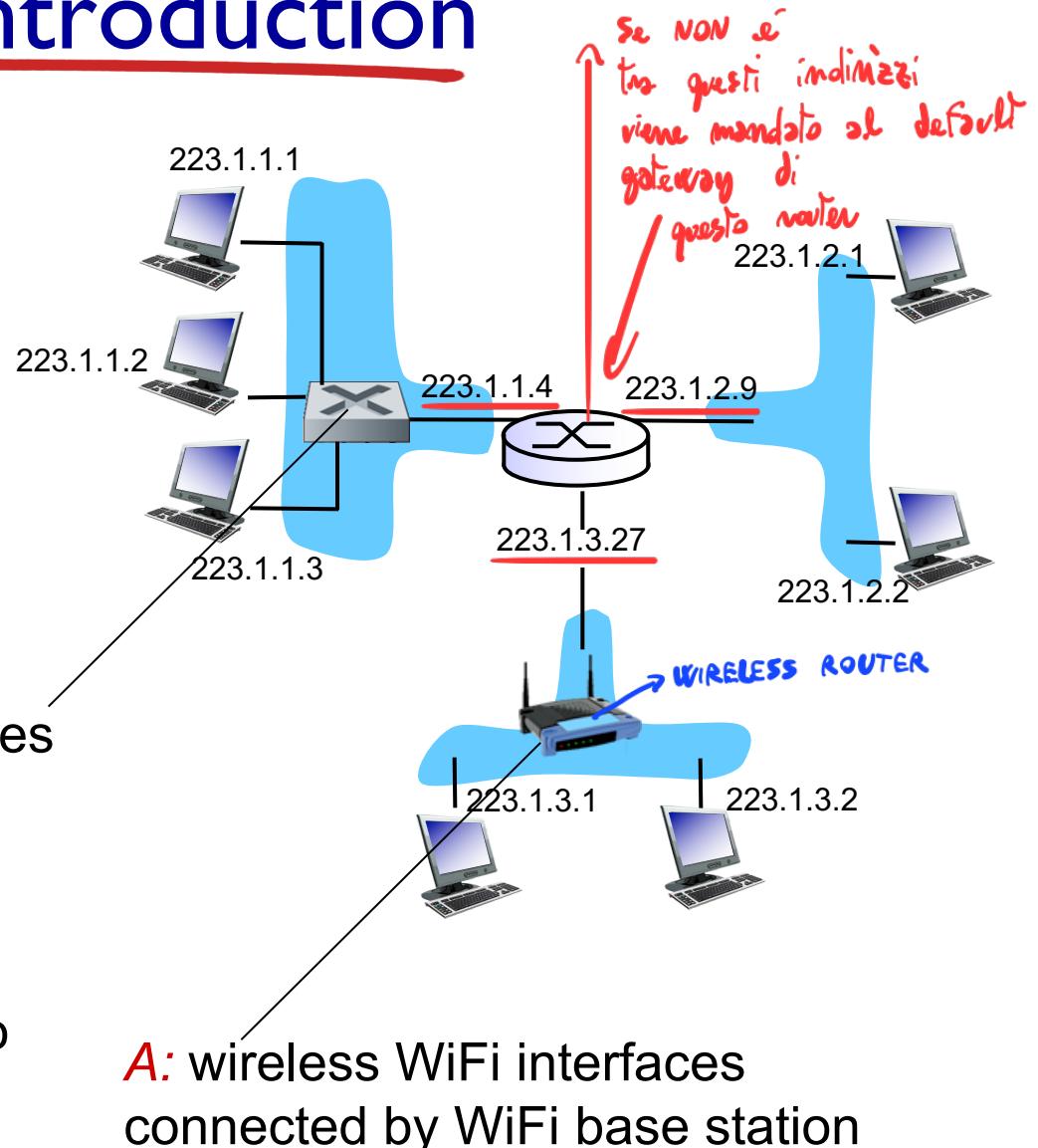
# IP addressing: introduction

*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)



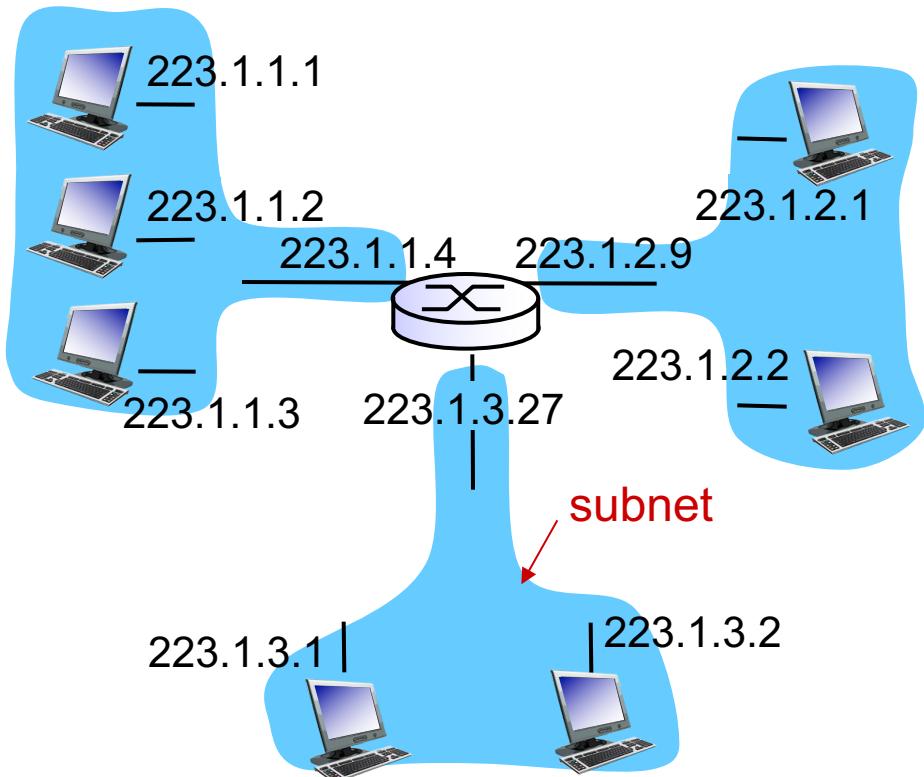
# Subnets

## ■ 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*

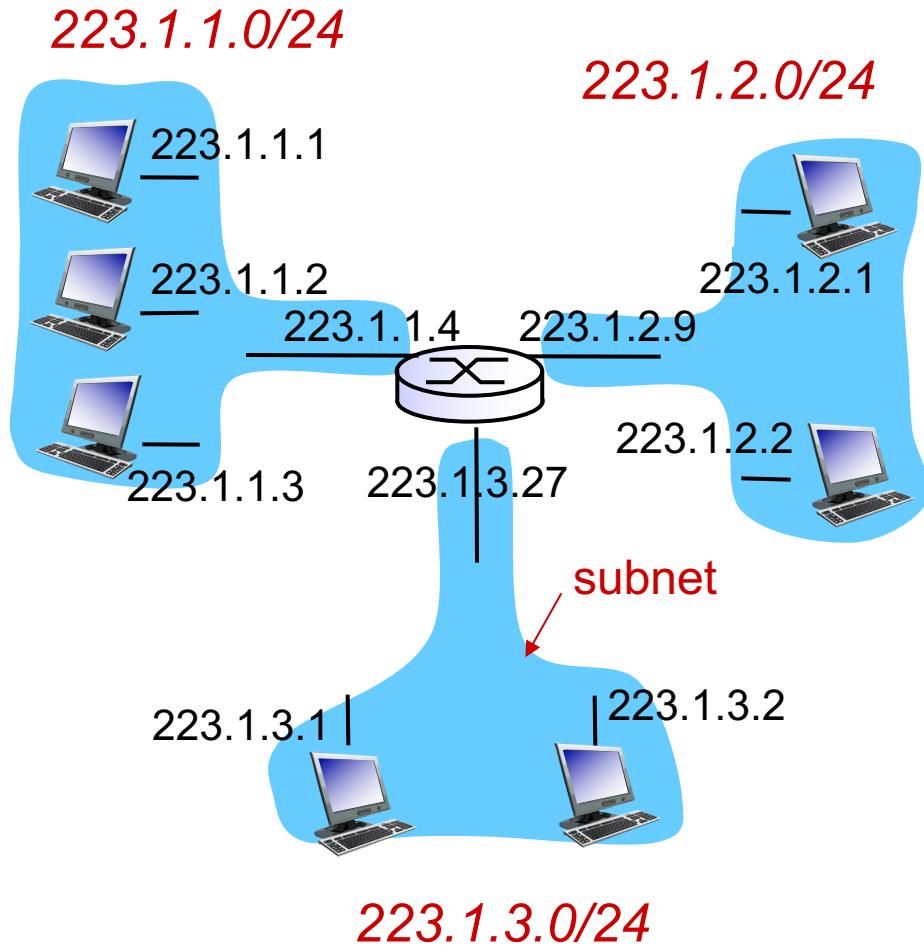


network consisting of 3 subnets

# Subnets

## *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

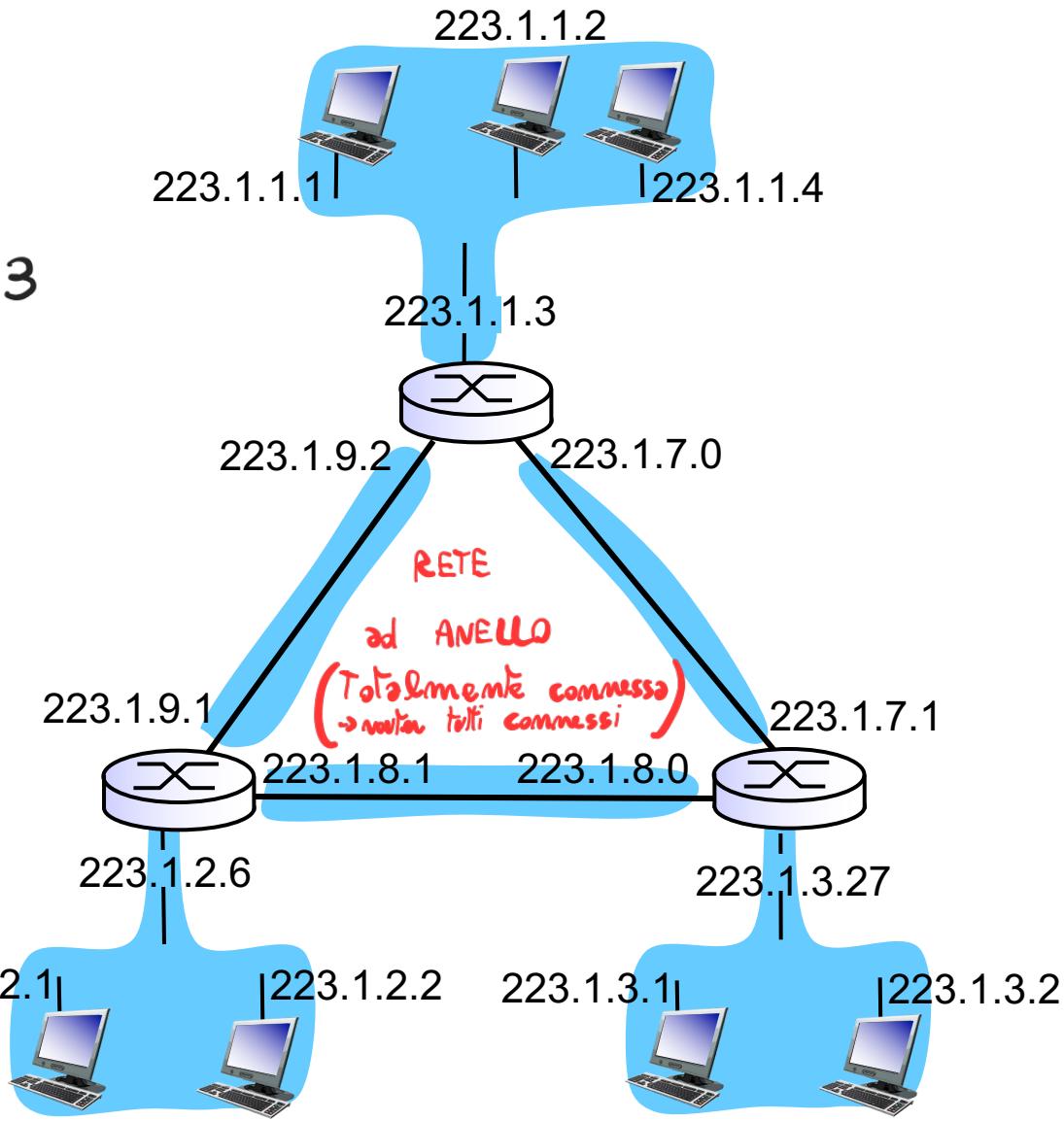
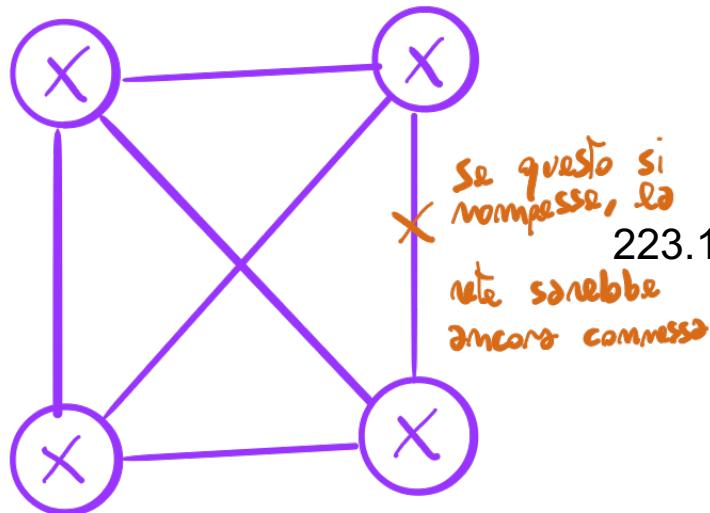
# Subnets

3 sottonet di rete tutto diversi

how many?

- Quante ne ci sono? 3  
(sono di classe C)

Un altro esempio di rete ad anello totalmente connessa:



# IP addressing: CIDR

## CIDR: Classless InterDomain Routing

- 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?

**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: Dynamic Host Configuration Protocol:** dynamically get address from as server
  - “plug-and-play”

# DHCP: Dynamic Host Configuration Protocol

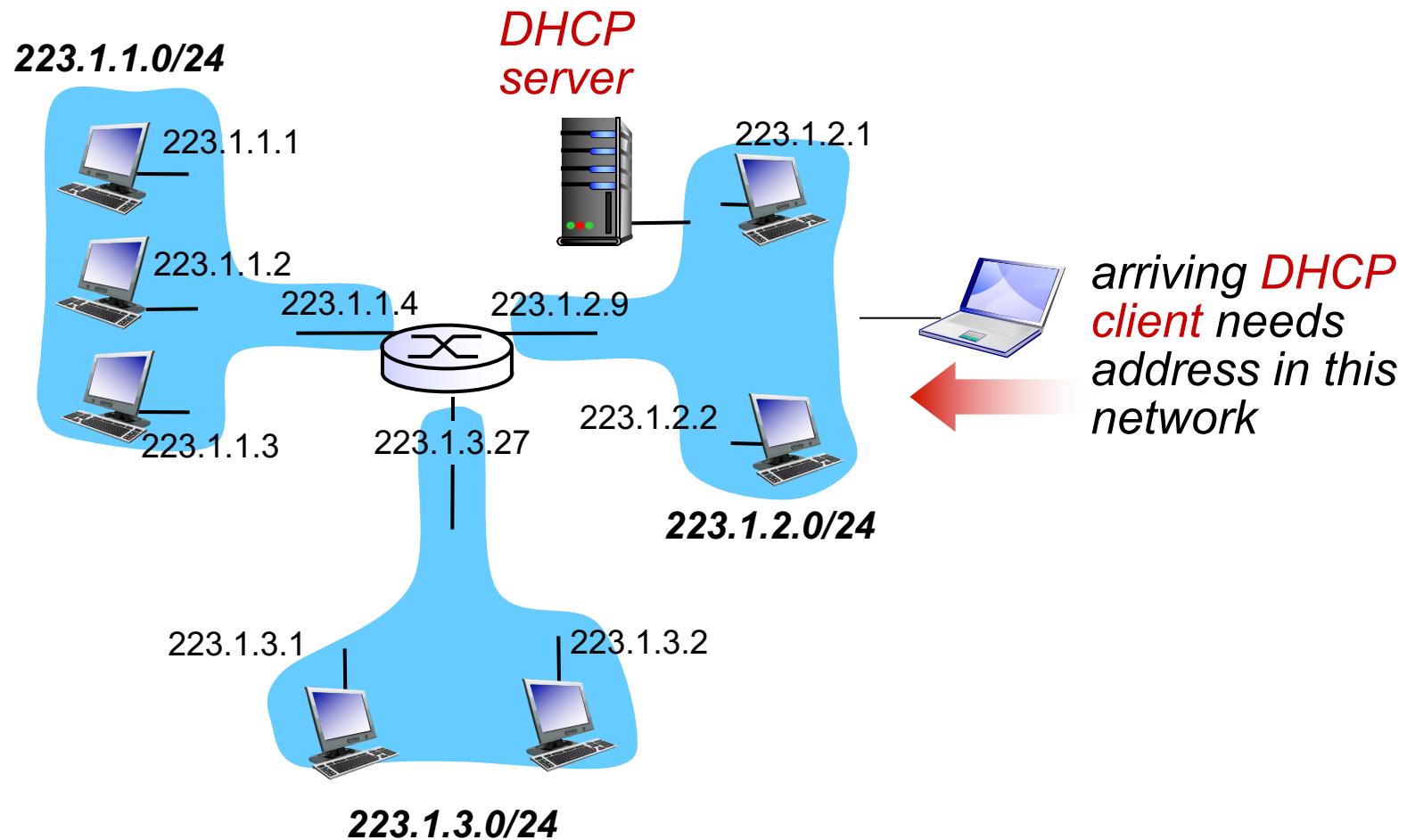
**goal:** allow host to *dynamically* obtain its IP address from network server when it joins network

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

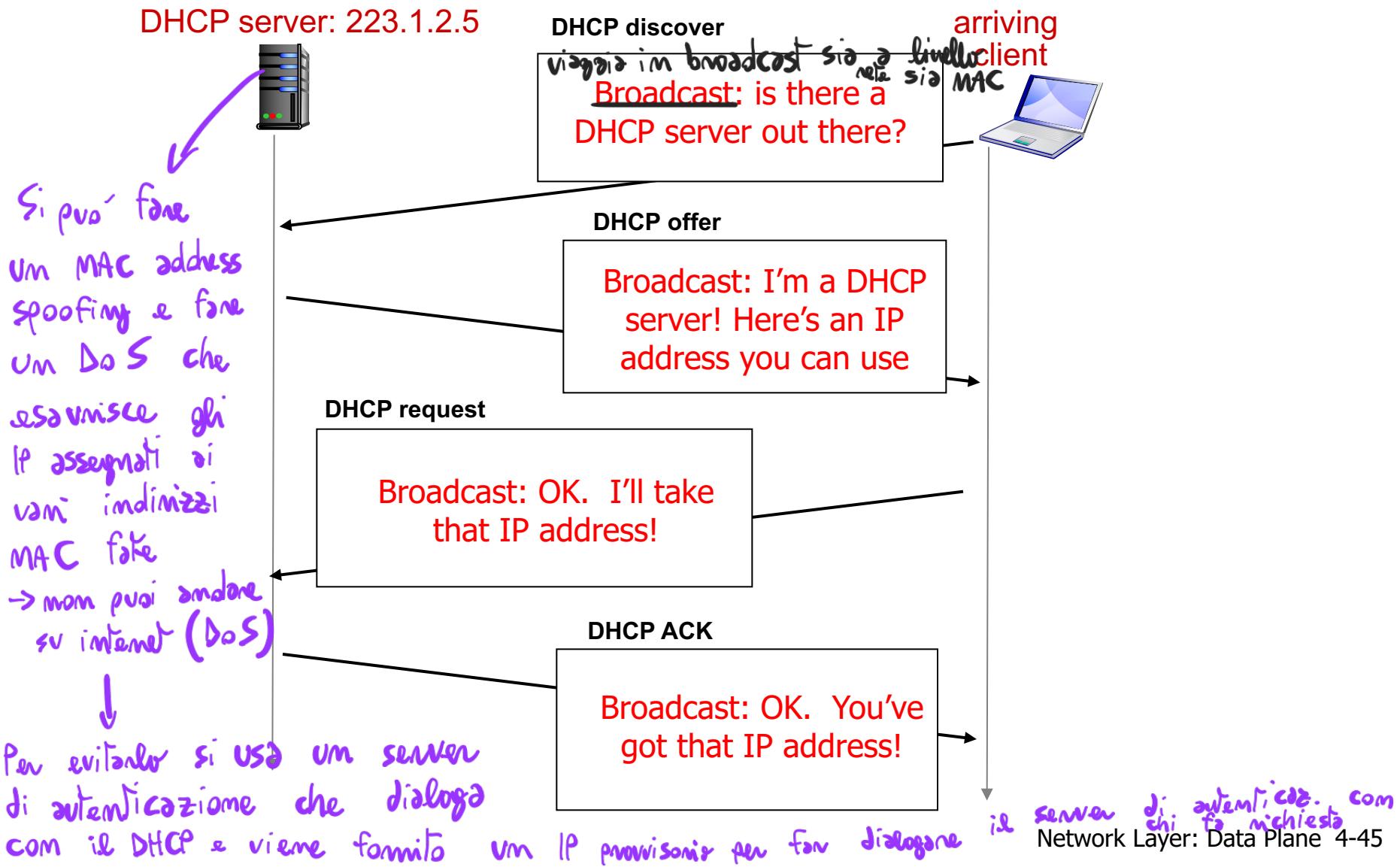
**DHCP overview:**

- 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

# DHCP client-server scenario



# DHCP client-server scenario

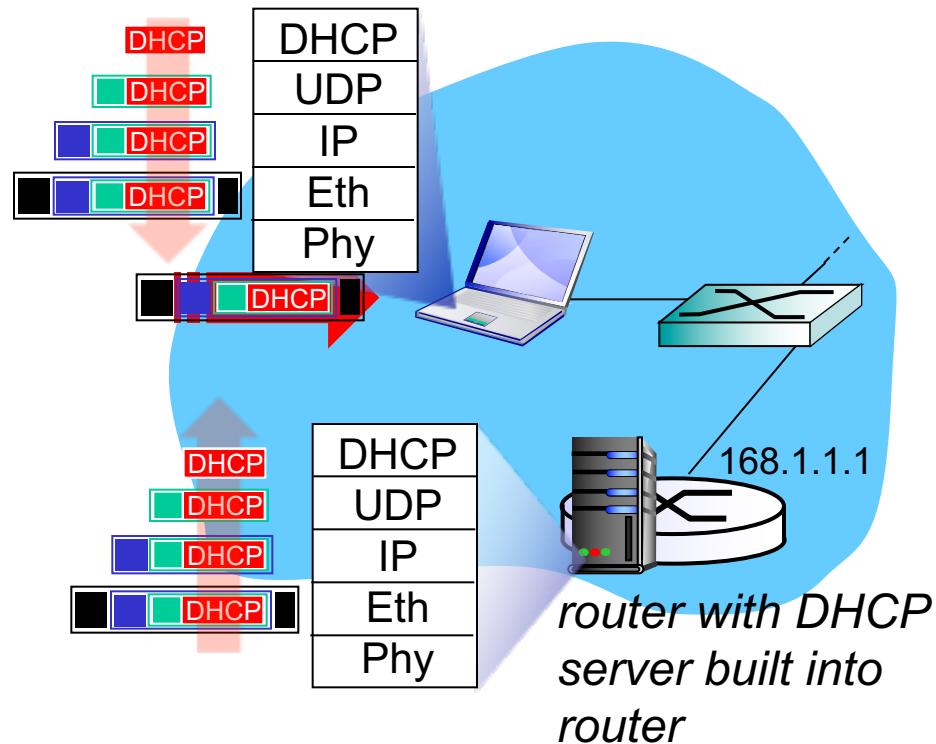


# DHCP: more than IP addresses

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

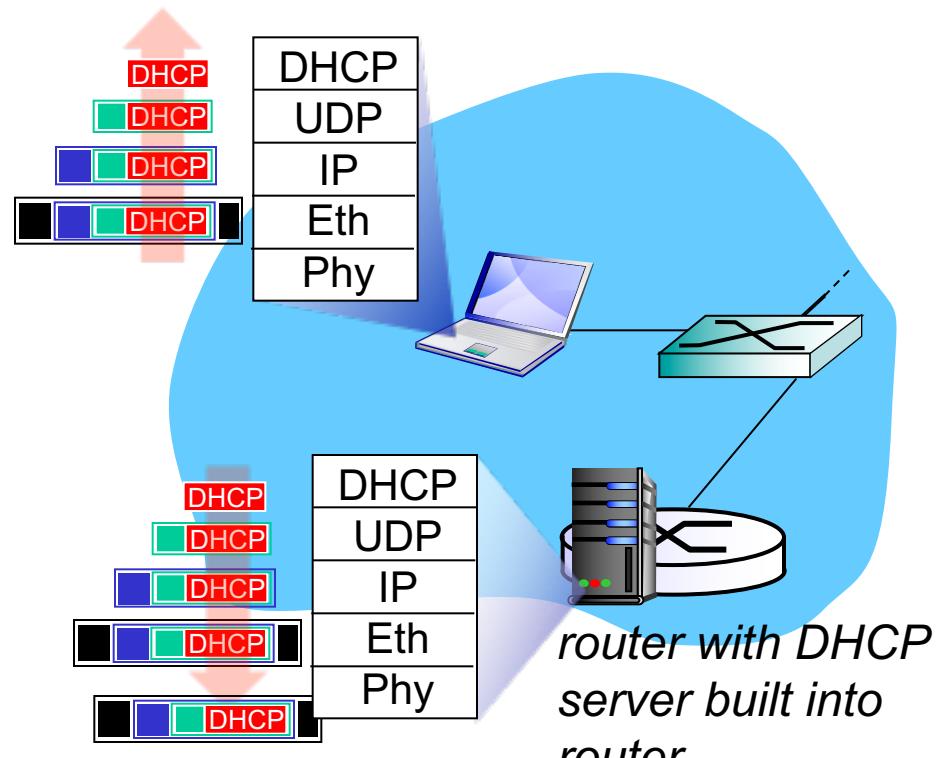
- address of first-hop router for client
- name and IP address of DNS sever
- network mask (indicating network versus host portion of address)

# 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: FFFFFFFFFFFF) 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 DSN server, IP address of its first-hop router

# 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

reply

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."**

# IP addresses: how to get one?

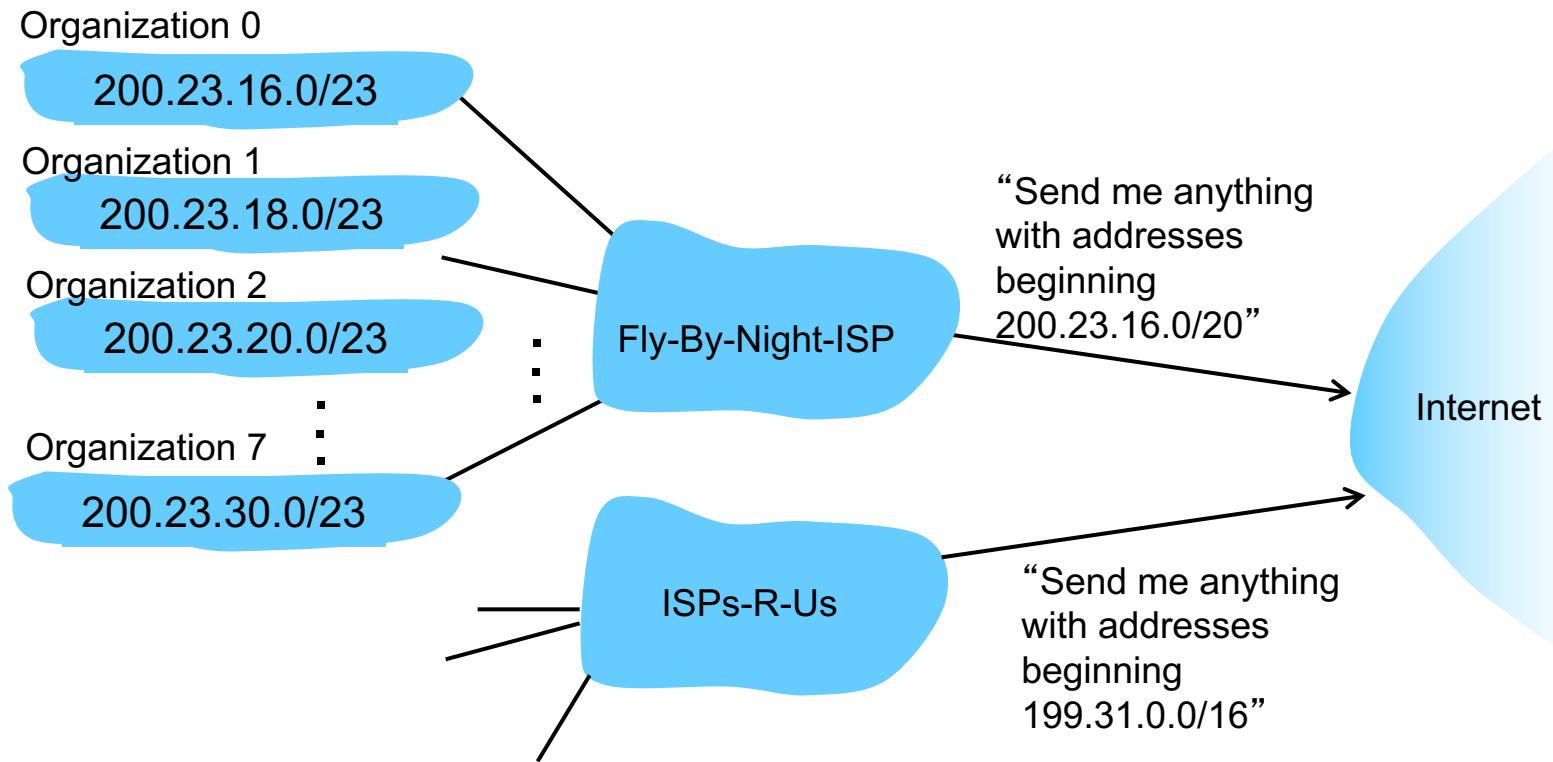
**Q:** how does *network* get subnet part of IP addr?

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

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

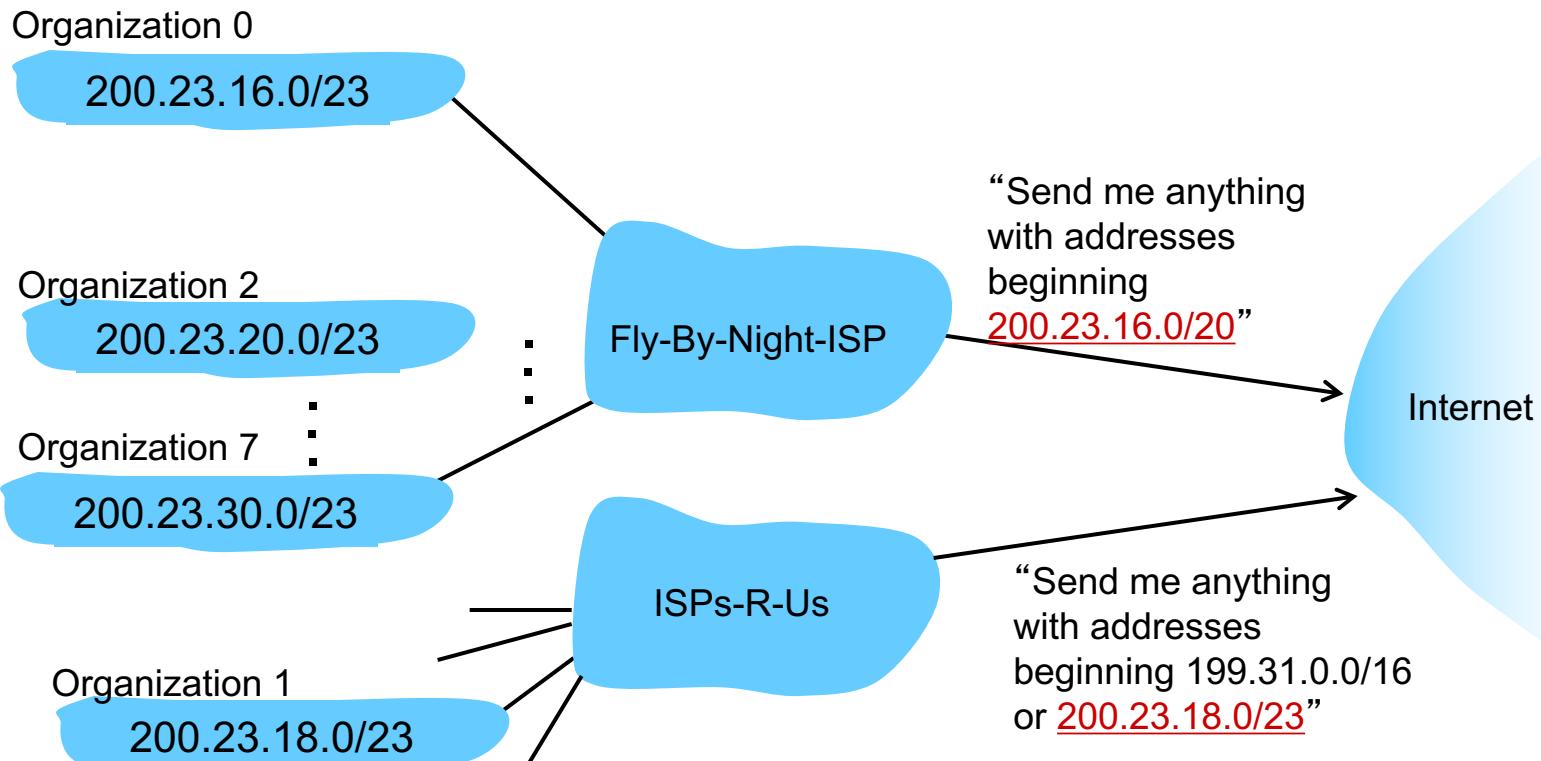
# Hierarchical addressing: route aggregation

hierarchical addressing allows efficient advertisement of routing information:



# Hierarchical addressing: more specific routes

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



# IP addressing: the last word...

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

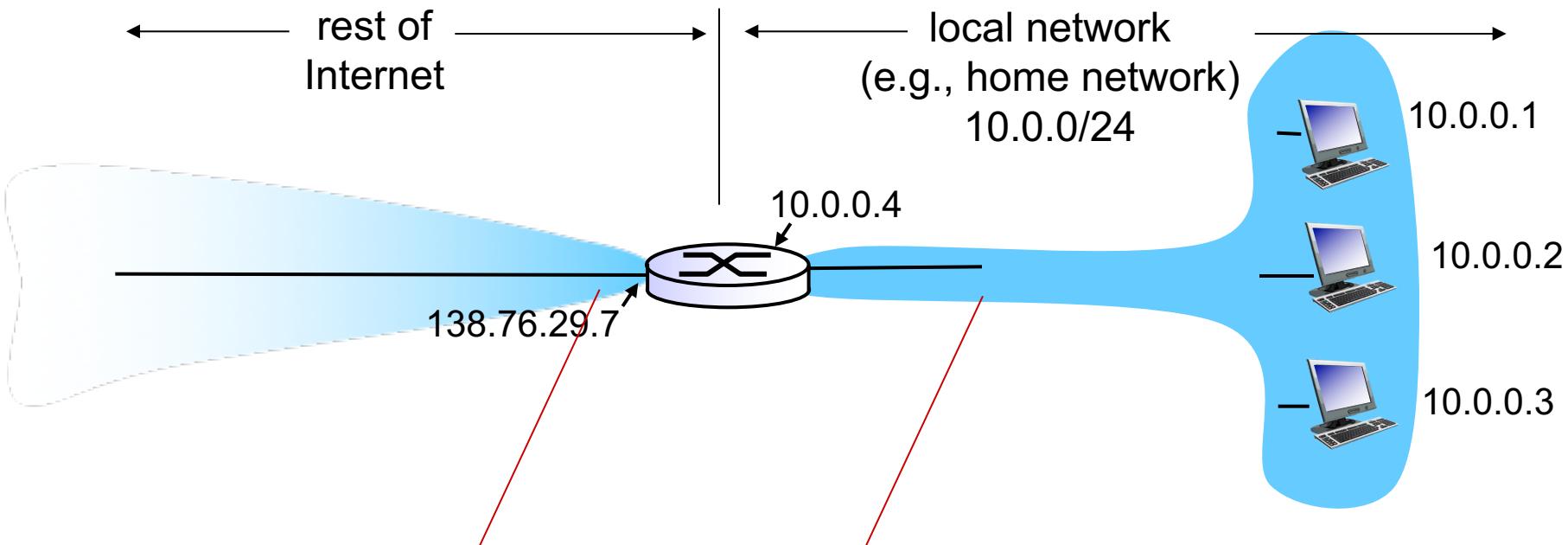
**A:** ICANN: Internet Corporation for Assigned

Names and Numbers <http://www.icann.org/>

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

→ Oggi si compra un solo indirizzo IP e si mette un router NAT → nella rete locale posso fare "quello che voglio"

# 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

***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:

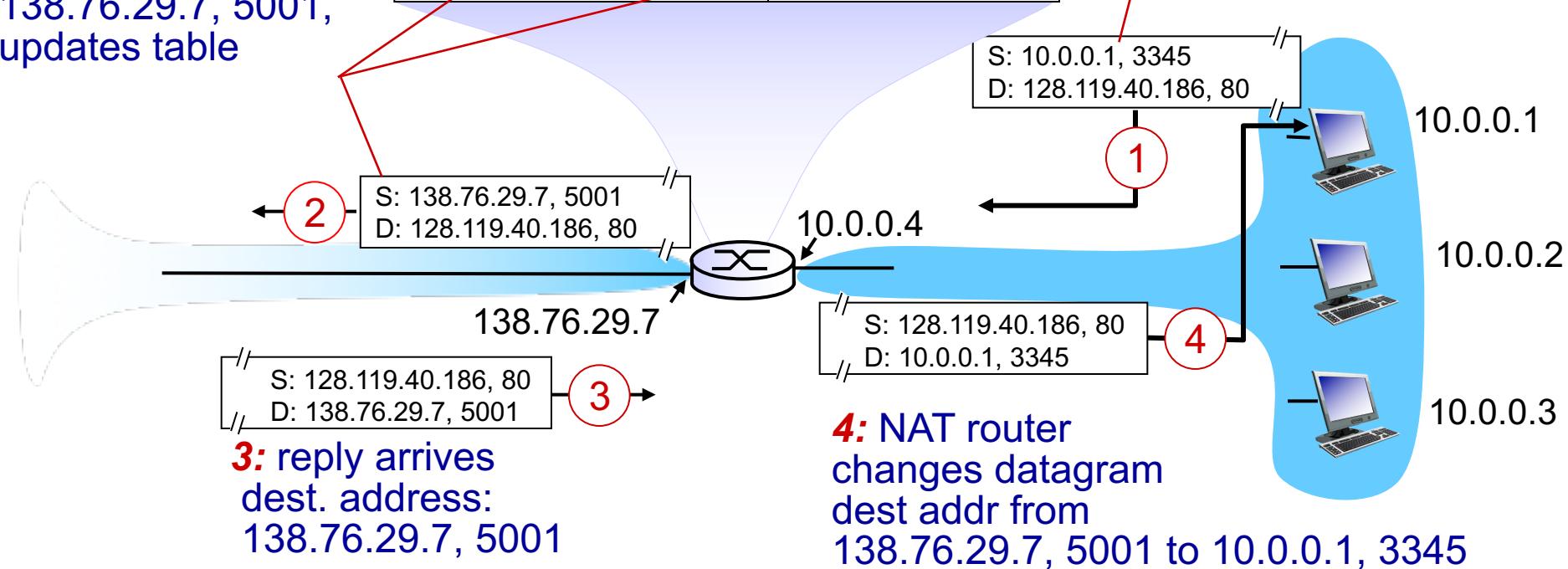
- *outgoing datagrams:* replace (source IP address, port #) of every outgoing datagram to (NAT IP address, new port #)  
    . . . remote clients/servers will respond using (NAT IP address, new port #) as destination addr
- *remember (in NAT translation table)* every (source IP address, port #) to (NAT 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



\* 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/)

# 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
  - 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: 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

# IPv6: motivation

- *initial motivation:* 32-bit address space soon to be completely allocated. → Gli IPv4 sarebbero finiti
- additional motivation:
  - header format helps speed processing/forwarding
  - header changes to facilitate QoS

## *IPv6 datagram format:*

- fixed-length 40 byte header
- no fragmentation allowed → No FRAGMENTAZIONE

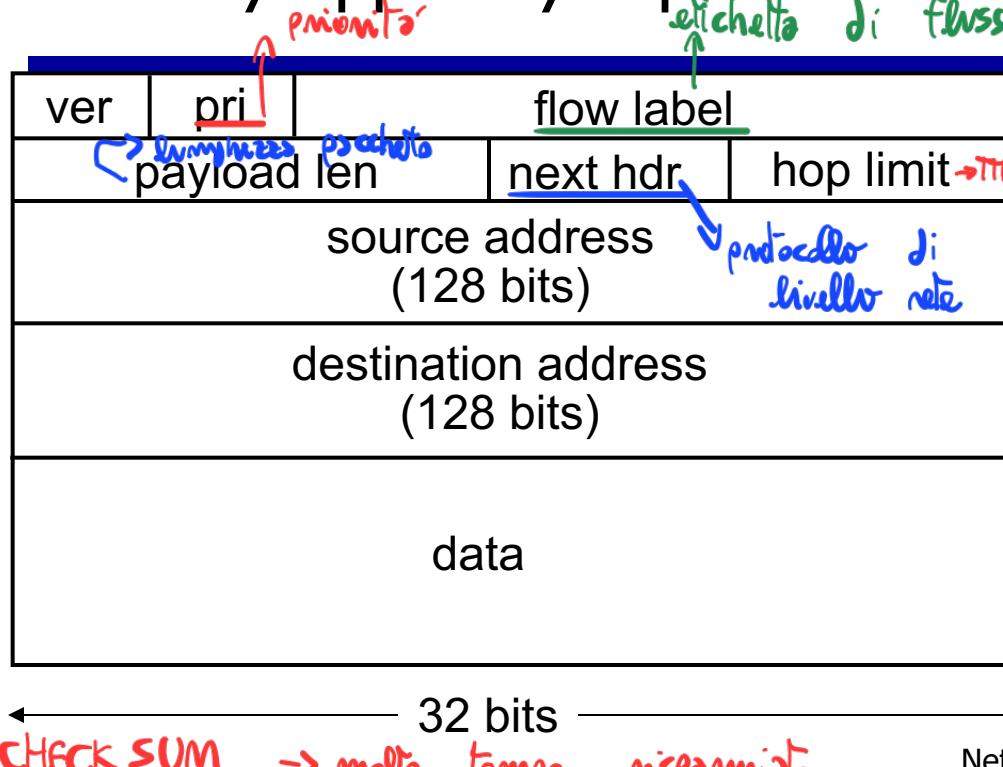
# IPv6 datagram format

*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



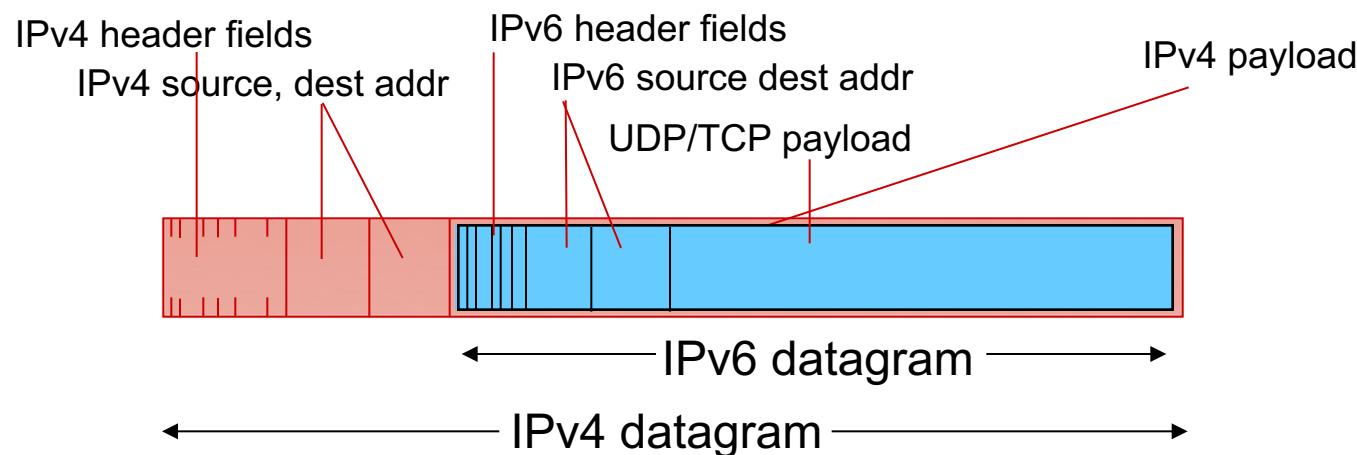
i router IPv6  
hanno una tabella  
che creano e controllano  
l'indirizzo del mittente  
e destinatario

# 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” → ERRORE che può accadere in IPv6
    - multicast group management functions
- c'è in IPv6 poiché manca la frammentazione*
- ↑
- ↓
- pacchetto ucciso  
e modificato al  
mittente con ICMP  
se il mittente che  
fa partire il pacchetto  
e lo invia più piccolo  
non è router

# 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

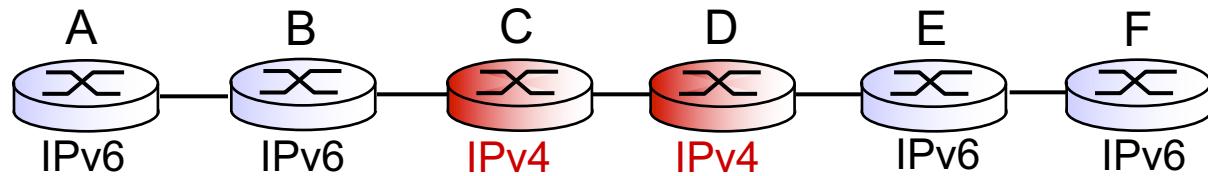


# Tunneling

logical view:

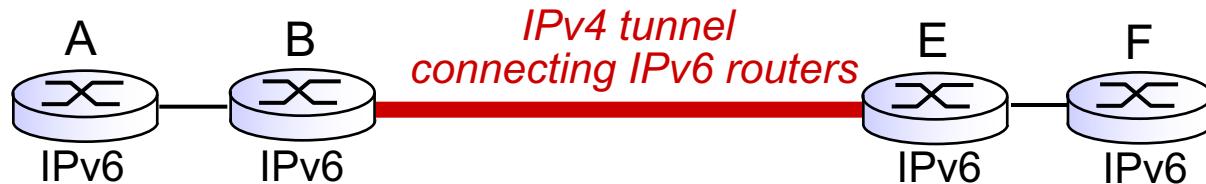


physical view:

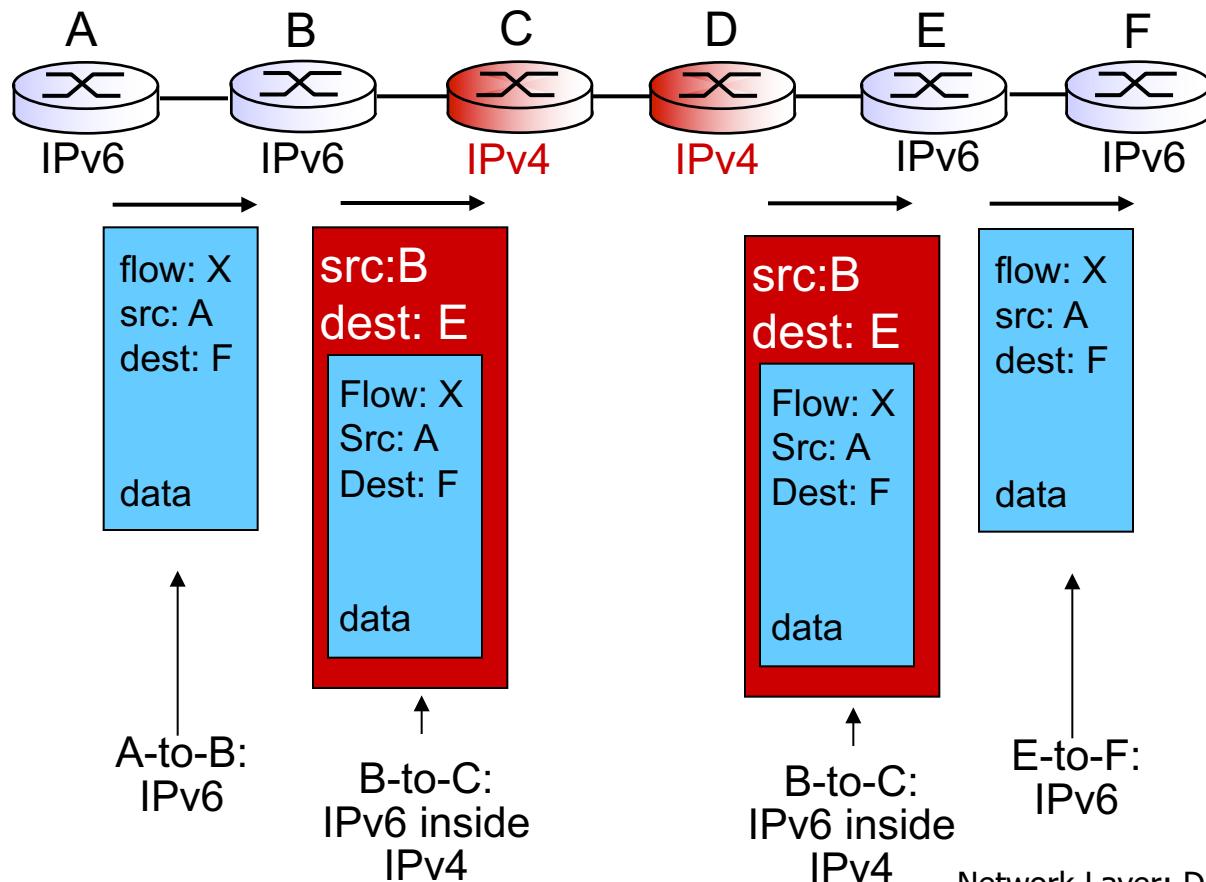


# Tunneling

logical view:



physical view:



# 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?*

# 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

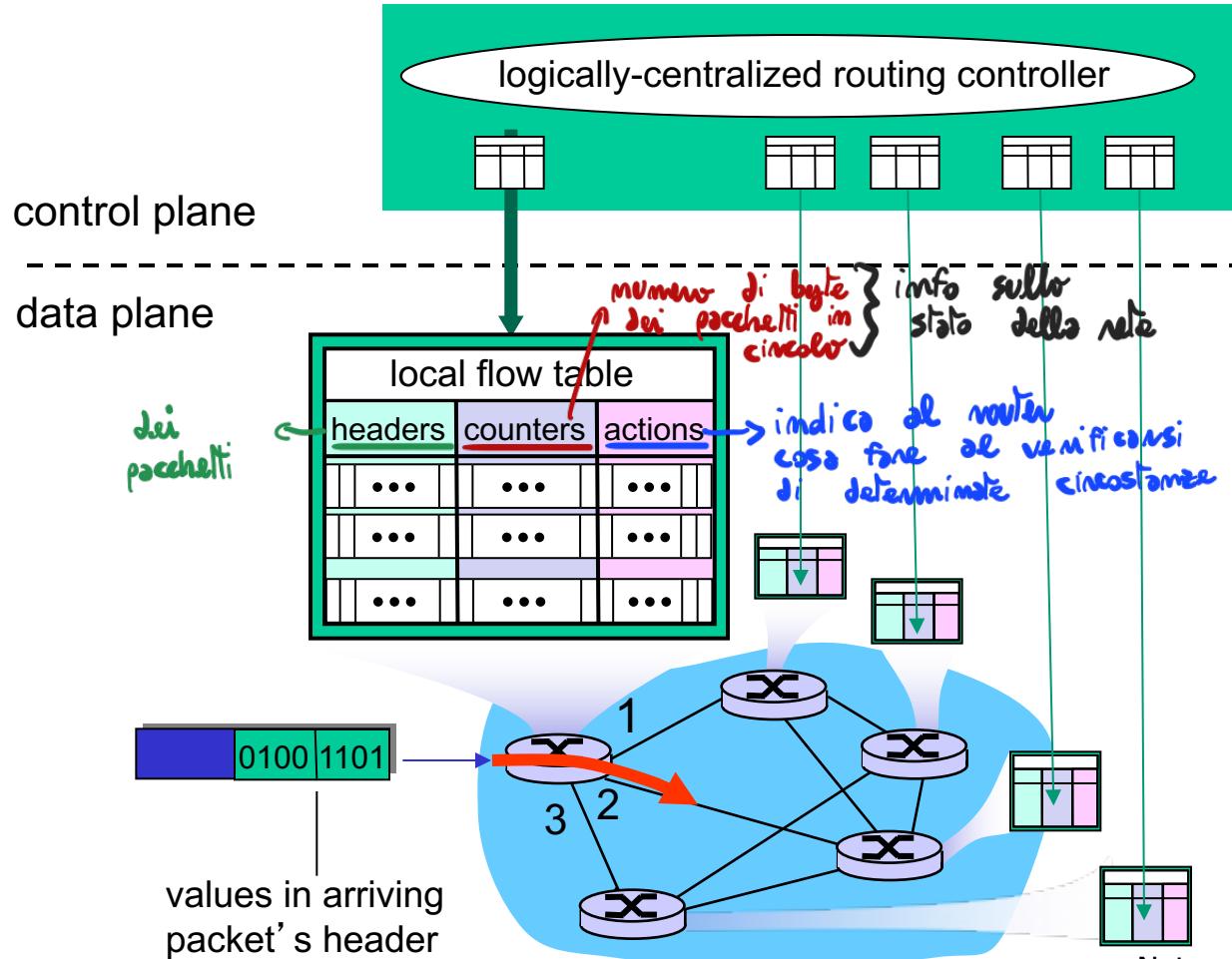
## 4.4 Generalized Forward and SDN

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

# Generalized Forwarding and SDN

Software Defined Network

Each router contains a **flow table** that is computed and distributed by a *logically centralized* routing controller

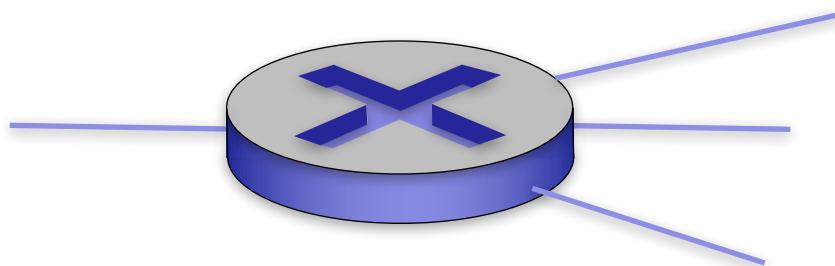


# OpenFlow data plane abstraction

→ identifica lo sequenza dei dati scambiati tra mittente e destinatario

- flow: defined by header fields
- generalized forwarding: simple packet-handling rules
  - *Pattern*: match 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

→ azioni convenzionali per cui se un pacchetto fa match con un certo pattern → si esegue una certa azione



*Flow table in a router (computed and distributed by controller) define router's match+action rules*

# OpenFlow data plane abstraction

- *flow*: defined by header fields
- generalized forwarding: simple packet-handling rules
  - Pattern: match 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

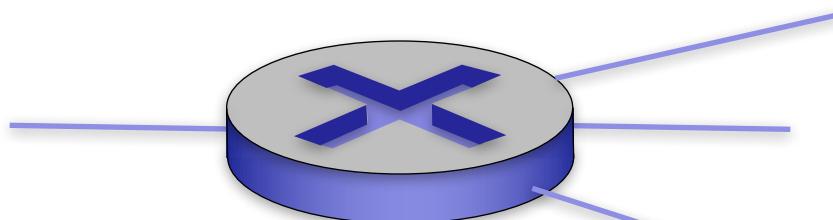
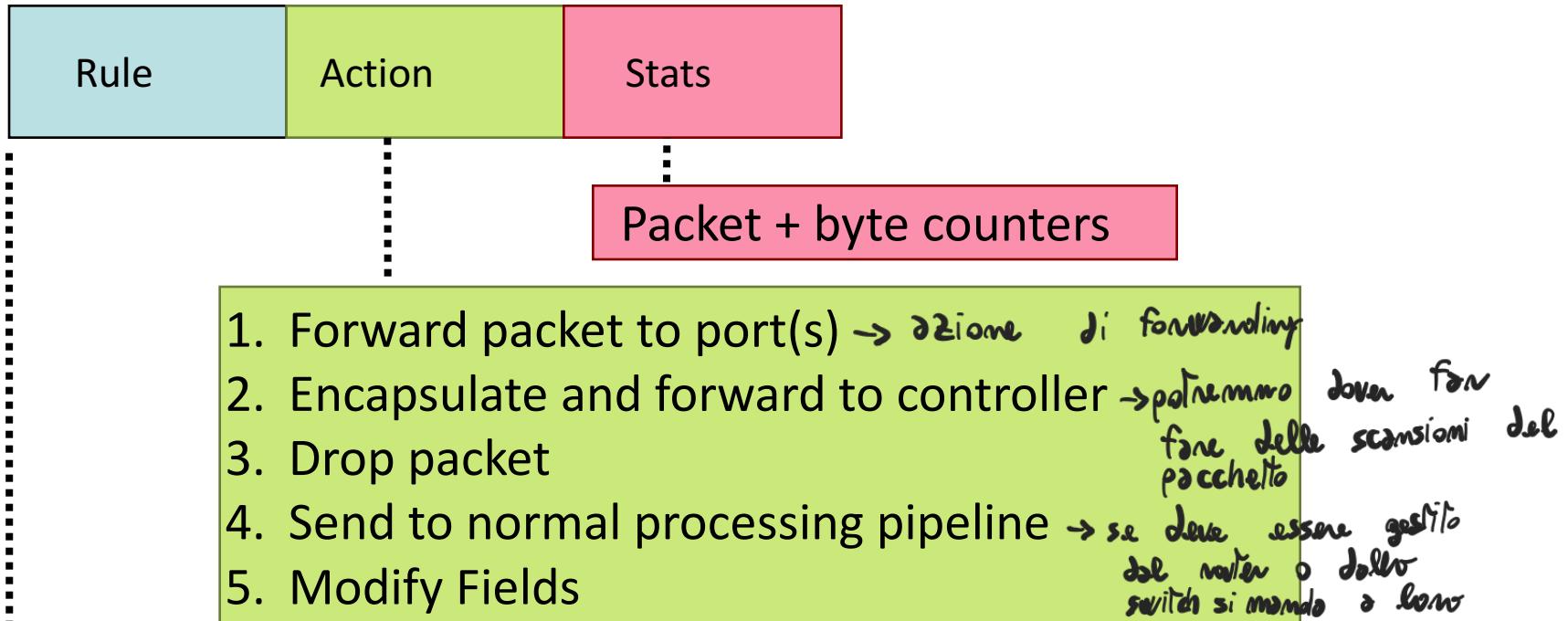


TABELLA DI FLUSSO

- indirizzi IP
- regole di firewalling \* : wildcard
1. src=1.2.\*.\*, dest=3.4.5.\* → drop *se c'è match con l'header di un pacchetto → umazione (drop)*
  2. src = \*.\*.\*.\*, dest=3.4.\*.\* → forward(2) *verso la porta 2 → regola di instidamento (forwarding)*
  3. src=10.1.2.3, dest=\*.\*.\*.\* → send to controller *E> FIREWALL ... E> REMOTE CONTROLLER ...*

# OpenFlow: Flow Table Entries



→ Virtual LAN → può essere usato per separare il traffico

Switch Port	VLAN ID	MAC src	MAC dst	Eth type	IP Src	IP Dst	IP Prot	TCP sport	TCP dport
LINK LAYER → similar MAC									

stessa famiglia di macchine che comunicano  
sulla stessa porta di uno switch → serve uno switch manager

Network layer

Transport layer

# Examples

## Destination-based forwarding:

Sia il  
seguinte:

Switch Port	MAC src	MAC dst	Eth type	VLAN ID	IP Src	IP Dst	IP Prot	TCP sport	TCP dport	Action
-------------	---------	---------	----------	---------	--------	--------	---------	-----------	-----------	--------

**PATTERN** { \*      \*      \*      ...      \*      \*      \*      51.6.0.8 \*      \*      \*      port6      inserire il pachetto  
 ↓  
 se le informaz. NON interessano ciò che è scritto qui  
 del pacchetto fanno  
**MATCH** ↓  
**action**

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	IP Src	IP Dst	IP Prot	TCP sport	TCP dport	Action
*	*	*	*	*	*	*	*	*	22	drop

do not forward (block) all datagrams destined to TCP port 22

Switch Port	MAC src	MAC dst	Eth type	VLAN ID	IP Src	IP Dst	IP Prot	TCP sport	TCP dport	Action
*	*	*	*	*	128.119.1.1	*	*	*	*	drop

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

# Examples

Destination-based layer 2 (switch) forwarding: in una rete locale (LAN)

Switch Port	MAC src	MAC dst	Eth type	VLAN ID	IP Src	IP Dst	IP Prot	TCP sport	TCP dport	Action
*	22:A7:23: 11:E1:02	*	*	*	*	*	*	*	*	port3

! *nella rete layer 2 frames from MAC address 22:A7:23:11:E1:02  
non guardo locale  
should be forwarded to output port 6*

*l'indirizzo IP*

*(sono su una rete locale se sto usando il MAC)*

# OpenFlow abstraction

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

Se ho un software che implementa OpenFlow si può programmare la rete

# OpenFlow example

TABELLA DI FLUSSO SWITCH S3

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

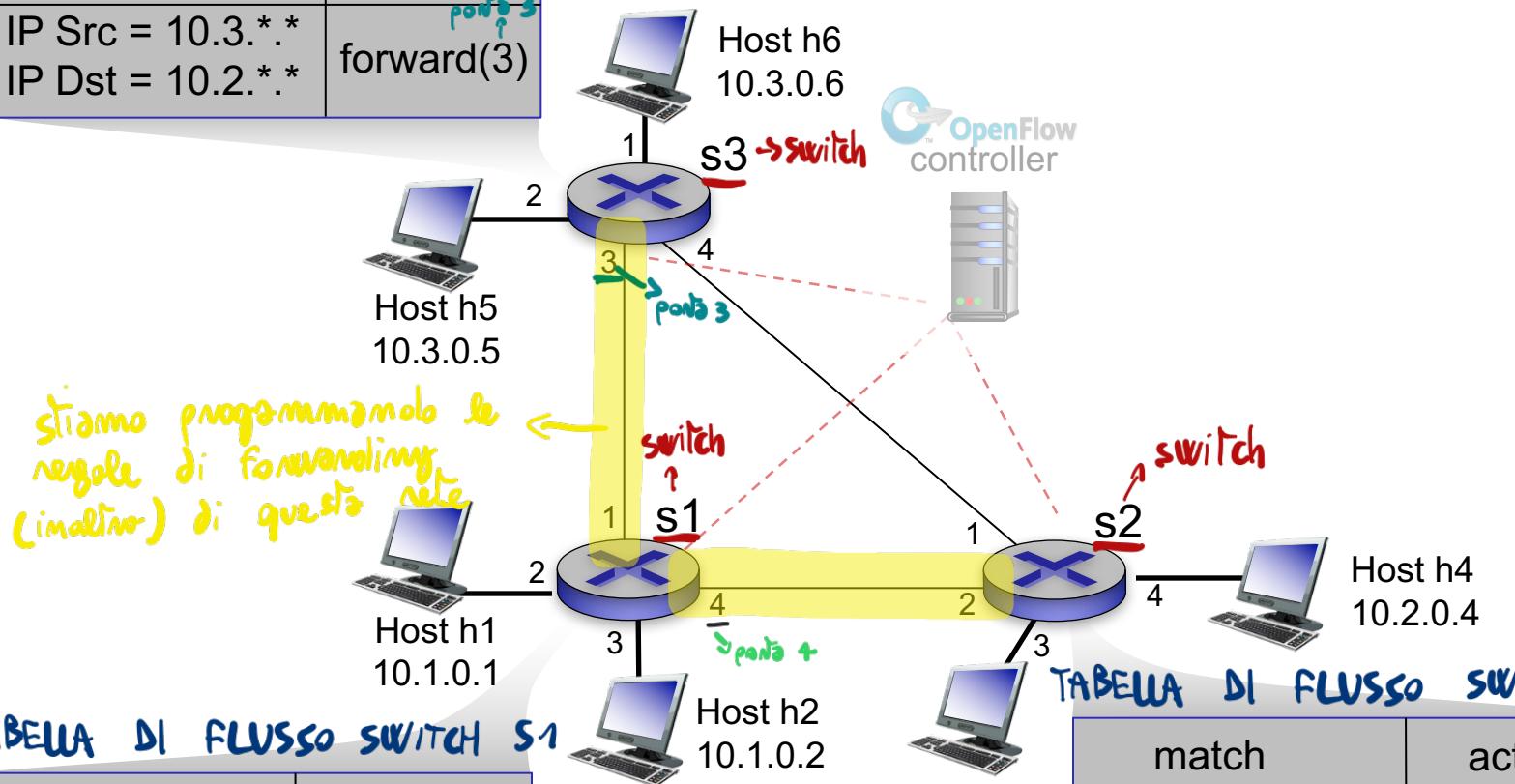


TABELLA DI FLUSSO SWITCH S1

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

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

TABELLA DI FLUSSO SWITCH S2

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

# 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) or flow tables (generalized forwarding) computed?

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