

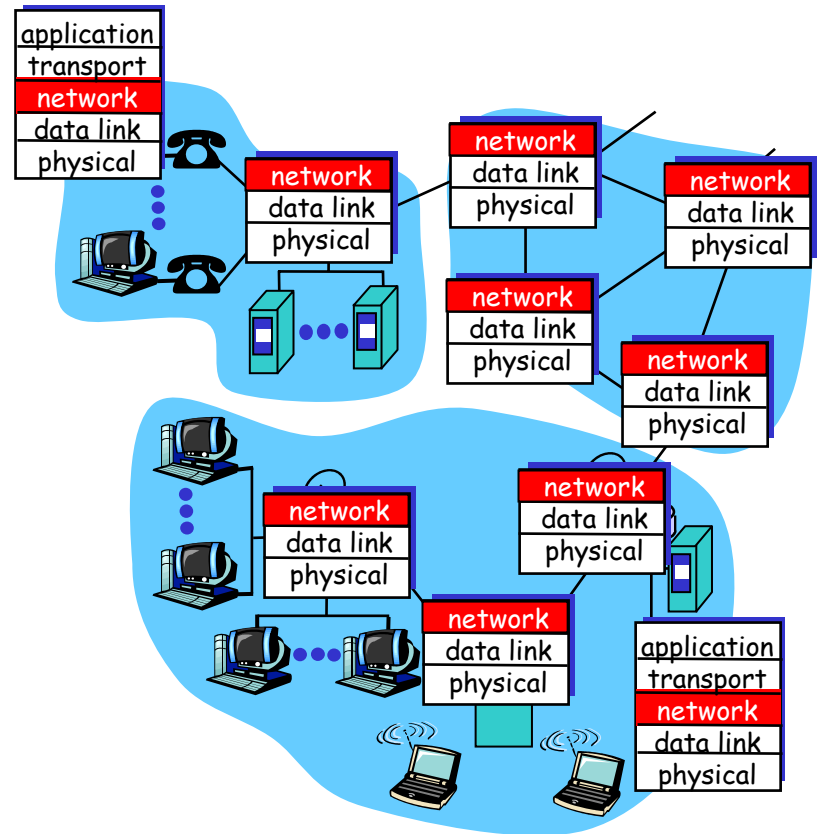
# Chapter 4: Network Layer Part I

(last revised 22/03/05)

- ❑ 4.1 Introduction
- ❑ 4.2 Virtual circuit and datagram networks
- ❑ 4.3 What's inside a router
- ❑ 4.4 IP: Internet Protocol
  - Datagram format
  - IPv4 addressing
  - ICMP
  - IPv6
- ❑ 4.5 Routing algorithms
  - Link state
  - Distance Vector
  - Hierarchical routing
- ❑ 4.6 Routing in the Internet
  - RIP
  - OSPF
  - BGP
- ❑ 4.7 Broadcast and multicast routing (maybe)

# Network layer

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



# Key Network-Layer Functions

□ *forwarding*: move packets from router's input to appropriate router output

□ *routing*: determine route taken by packets from source to dest.

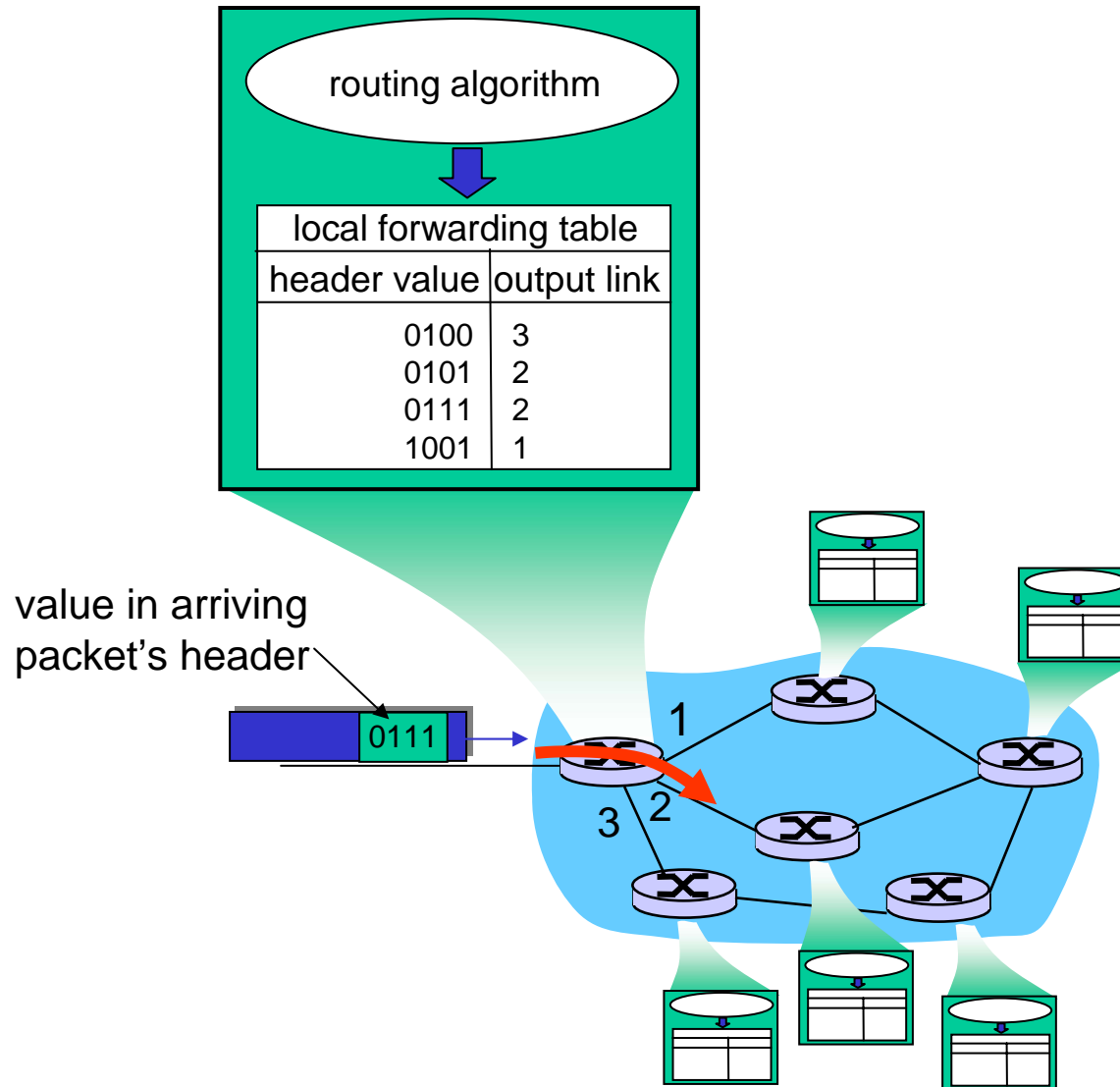
○ *Routing algorithms*

analogy:

□ *routing*: process of planning trip from source to dest

□ *forwarding*: process of getting through single interchange

# Interplay between routing and forwarding



- ❑ 3<sup>rd</sup> important function in *some* network architectures:
  - ATM, frame relay, X.25
  - Before datagrams flow, two hosts and intervening routers establish virtual connection
  - Routers get involved
  
- ❑ Network vs. transport layer cnctn service:
  - **Network:** between two hosts
  - **Transport:** between two processes

# Network service model

Q: What *service model* for "channel" transporting packets from sender to receiver?

- service abstraction
- ☐ guaranteed bandwidth?
  - ☐ preservation of inter-packet timing (no jitter)?
  - ☐ loss-free delivery?
  - ☐ in-order delivery?
  - ☐ congestion feedback to sender?

The most important abstraction provided by network layer:

virtual circuit  
or  
datagram?

# Network layer service models:

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

# Chapter 4: Network Layer

- ❑ 4.1 Introduction
- ❑ 4.2 Virtual circuit and datagram networks
- ❑ 4.3 What's inside a router
- ❑ 4.4 IP: Internet Protocol
  - Datagram format
  - IPv4 addressing
  - ICMP
  - IPv6
- ❑ 4.5 Routing algorithms
  - Link state
  - Distance Vector
  - Hierarchical routing
- ❑ 4.6 Routing in the Internet
  - RIP
  - OSPF
  - BGP
- ❑ 4.7 Broadcast and multicast routing



# Network layer connection and connection-less service

- ❑ Datagram network provides network-layer connectionless service
- ❑ VC network provides network-layer connection service
- ❑ Analogous to the transport-layer services, but:
  - **Service:** host-to-host
  - **No choice:** network provides one or the other
  - **Implementation:** in the core

# Virtual circuits

“source-to-dest path behaves much like telephone circuit”

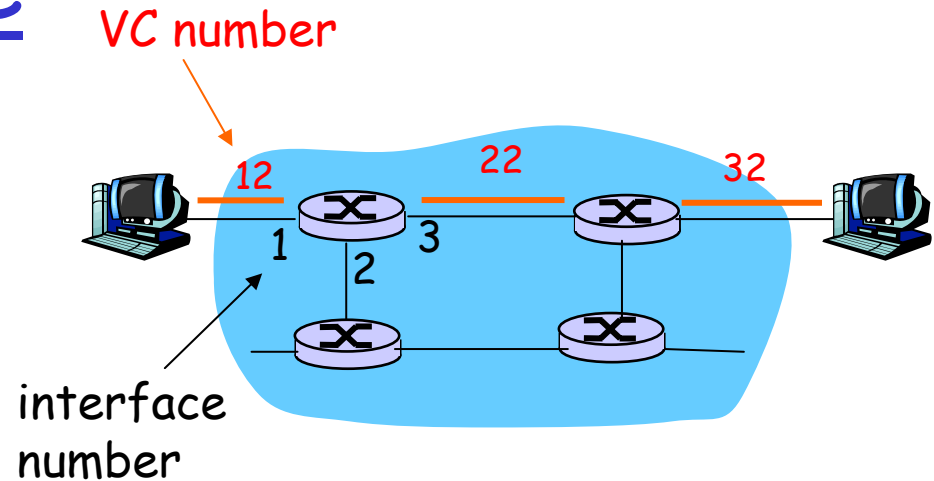
- performance-wise
  - network actions along source-to-dest path
- 
- ❑ call setup, teardown for each call *before* data can flow
  - ❑ each packet carries VC identifier (not destination host OD)
  - ❑ *every* router on source-dest paths maintain the “state” for each passing connection
    - transport-layer connection only involved two end systems
  - ❑ link, router resources (bandwidth, buffers) may be *allocated* to VC
    - to get circuit-like performance

# VC implementation

A VC consists of:

1. Path from source to destination
  2. VC numbers, one number for each link along path
  3. Entries in forwarding tables in routers along path
- ❑ Packet belonging to VC carries a VC number.
  - ❑ VC number must be changed on each link.
    - New VC number comes from forwarding table

## Forwarding table



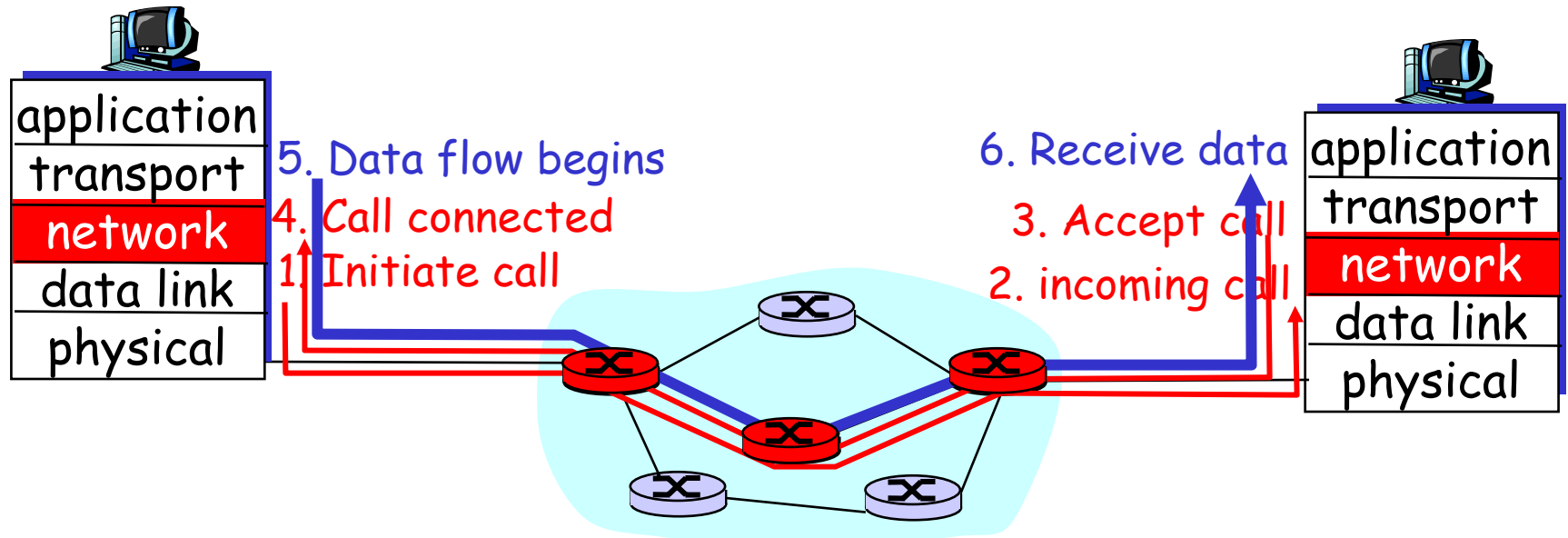
Forwarding table in upper-left router:

Incoming interface	Incoming VC #	Outgoing interface	Outgoing VC #
1	12	2	22
2	63	1	18
3	7	2	17
1	97	3	87
...	...	...	...

## Routers maintain connection state information!

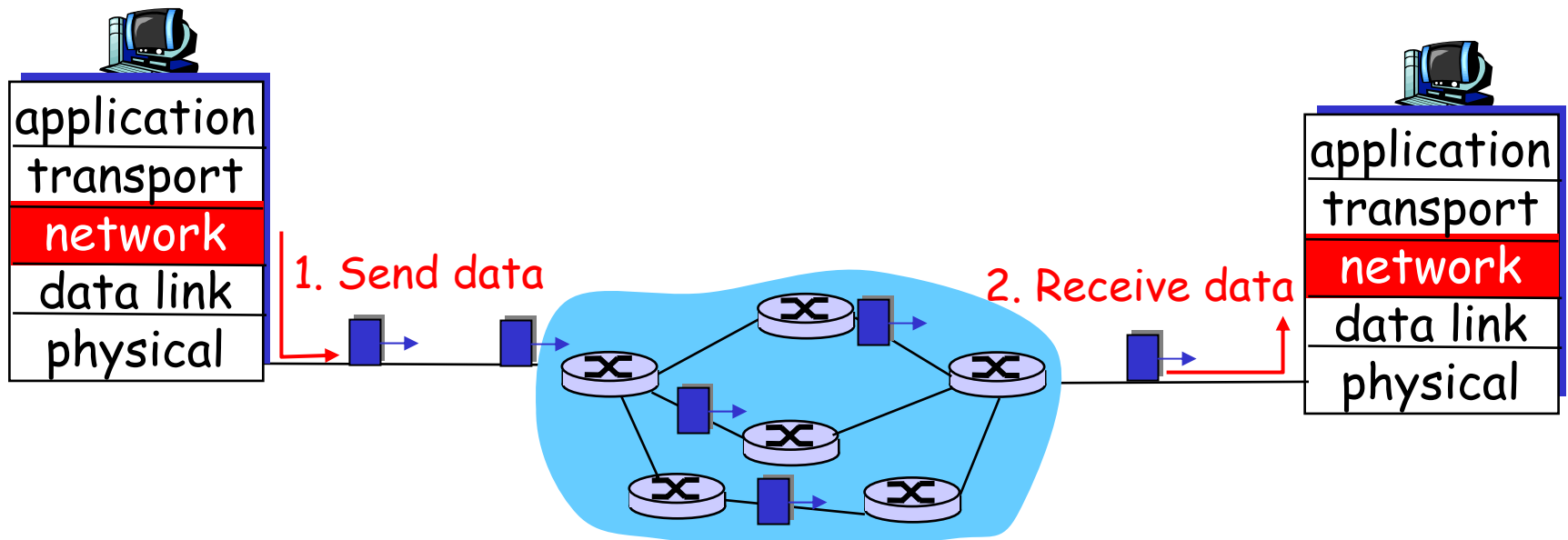
# Virtual circuits: signaling protocols

- used to setup, maintain teardown VC
- connection setup process still needs routing much as in the Internet
- used in ATM, frame-relay, X.25
- not used in today's Internet



# Datagram networks

- ❑ no call setup at network layer
- ❑ routers: no state info about end-to-end connections
  - no network-level concept of "connection"
- ❑ packets forwarded using destination host address
  - packets between same source-dest pair may take different paths



# Forwarding table

4 billion  
possible entries

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

# Longest prefix matching

<u>Prefix Match</u>	<u>Link Interface</u>
11001000 00010111 00010	0
11001000 00010111 00011000	1
11001000 00010111 00011	2
otherwise	3

## Examples

DA: 11001000 00010111 00010110 10100001

Which interface?

DA: 11001000 00010111 00011000 10101010

Which interface?



# Datagram or VC network: why?

## Internet

- ❑ data exchange among computers
  - “elastic” service, no strict timing req.
- ❑ “smart” end systems (computers)
  - can adapt, perform control, error recovery
  - simple inside network, complexity at “edge”
- ❑ many link types
  - different characteristics
  - uniform service difficult

## ATM

- ❑ evolved from telephony
- ❑ human conversation:
  - strict timing, reliability requirements
  - need for guaranteed service
- ❑ “dumb” end systems
  - telephones
  - complexity inside network

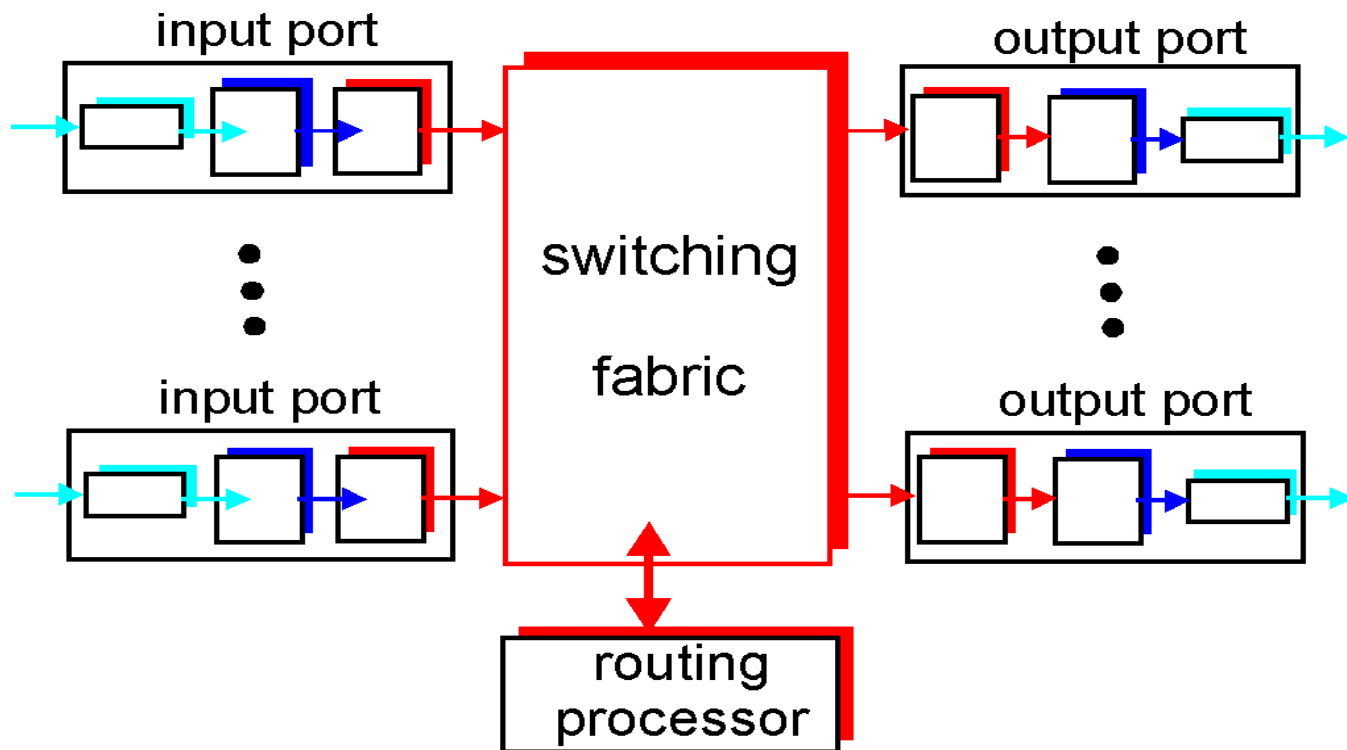
# Chapter 4: Network Layer

- ❑ 4.1 Introduction
- ❑ 4.2 Virtual circuit and datagram networks
- ❑ 4.3 What's inside a router
- ❑ 4.4 IP: Internet Protocol
  - Datagram format
  - IPv4 addressing
  - ICMP
  - IPv6
- ❑ 4.5 Routing algorithms
  - Link state
  - Distance Vector
  - Hierarchical routing
- ❑ 4.6 Routing in the Internet
  - RIP
  - OSPF
  - BGP
- ❑ 4.7 Broadcast and multicast routing

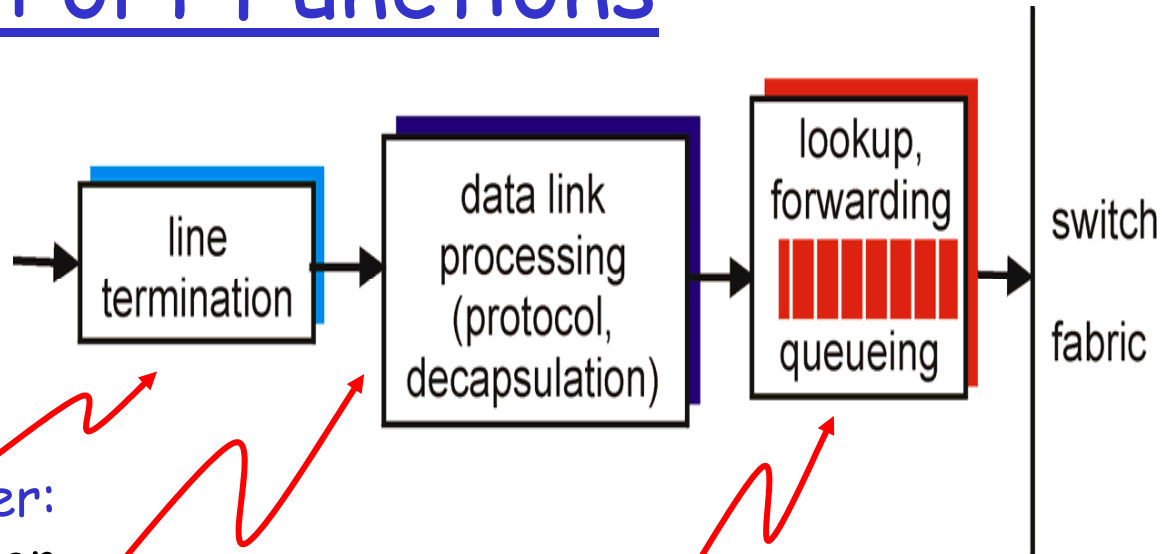
# Router Architecture Overview

Two key router functions:

- run routing algorithms/protocol (RIP, OSPF, BGP)
- *forwarding* datagrams from incoming to outgoing link



# Input Port Functions



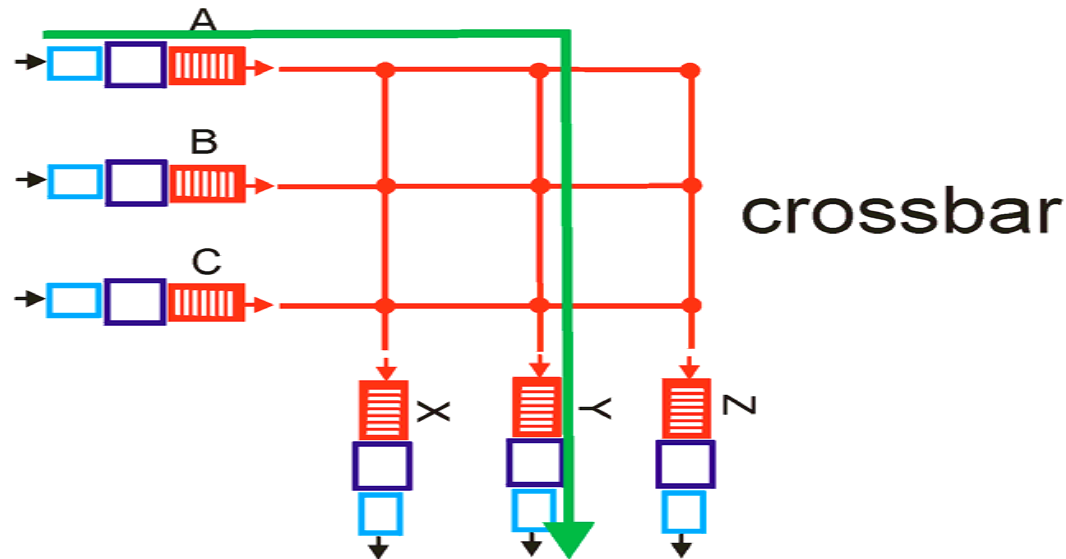
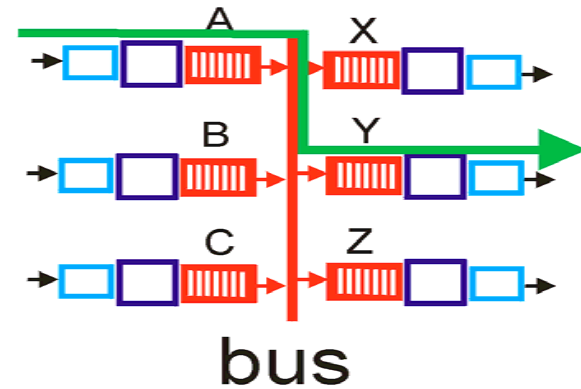
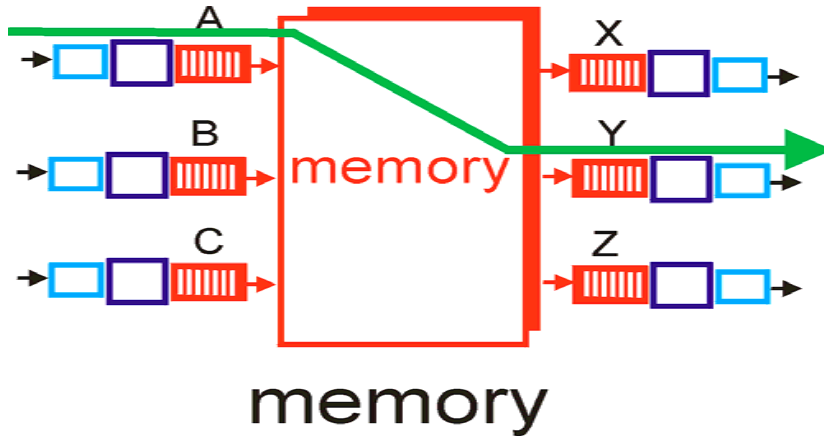
Physical layer:  
bit-level reception

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

## **Decentralized switching:**

- given datagram dest., lookup output port using forwarding table in input port memory
- goal: complete input port processing at 'line speed'
- queuing: if datagrams arrive faster than forwarding rate into switch fabric

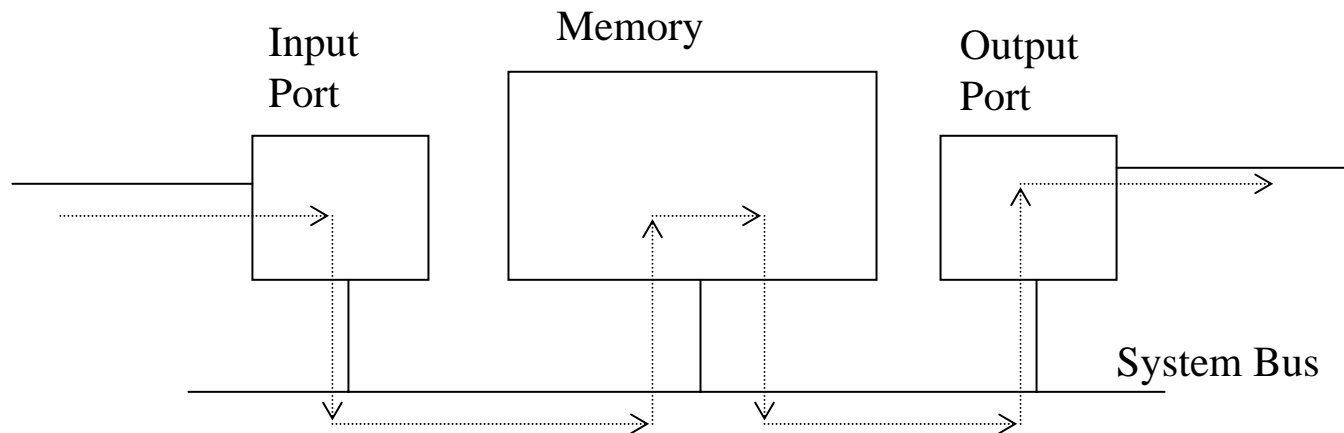
# Three types of switching fabrics



# Switching Via Memory

## First generation routers:

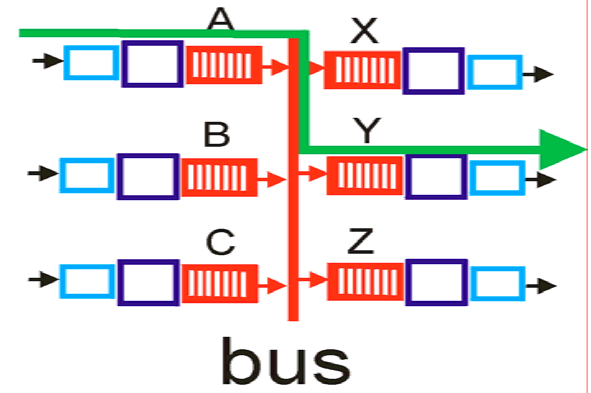
- ❑ packet copied by system's (single) CPU
- ❑ speed limited by memory bandwidth (2 bus crossings per datagram)



## Modern routers:

- ❑ input port processor performs lookup, copy into memory
- ❑ Cisco Catalyst 8500

# Switching Via a Bus



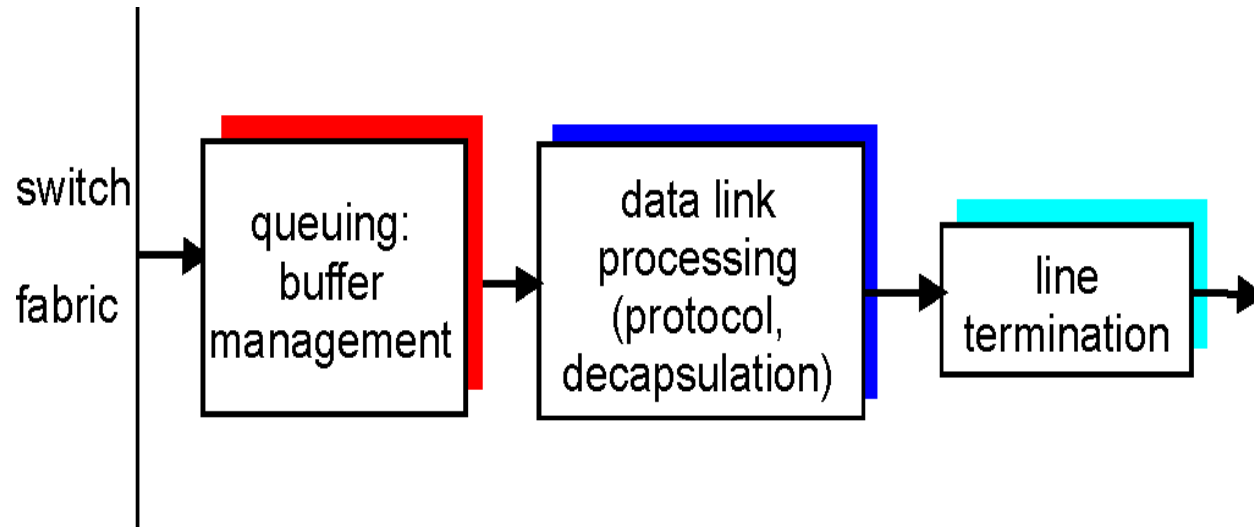
- ❑ datagram from input port memory to output port memory via a shared bus
- ❑ **bus contention:** switching speed limited by bus bandwidth
- ❑ 1 Gbps bus, Cisco 1900: sufficient speed for access and enterprise routers (not regional or backbone)

# Switching Via An Interconnection Network

- ❑ overcome bus bandwidth limitations
- ❑ Banyan networks, 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 Gbps through the interconnection network

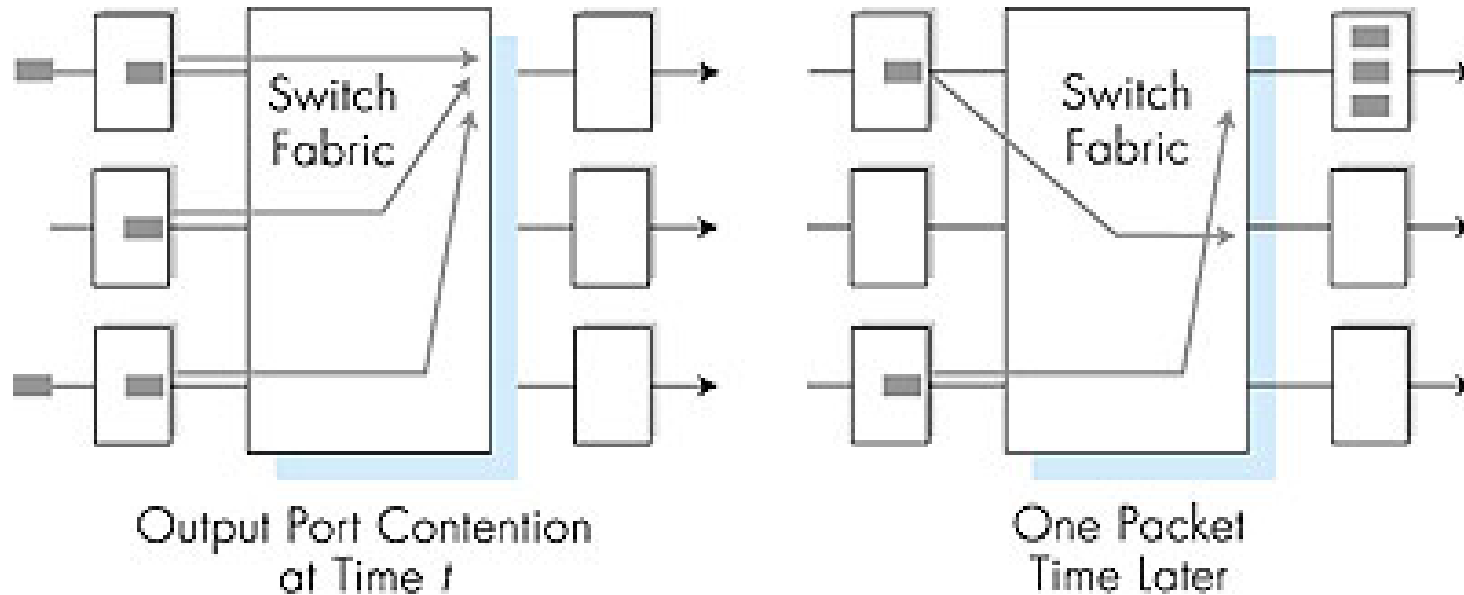


# Output Ports



- ❑ *Buffering* required when datagrams arrive from fabric faster than the transmission rate
- ❑ *Scheduling discipline* chooses among queued datagrams for transmission

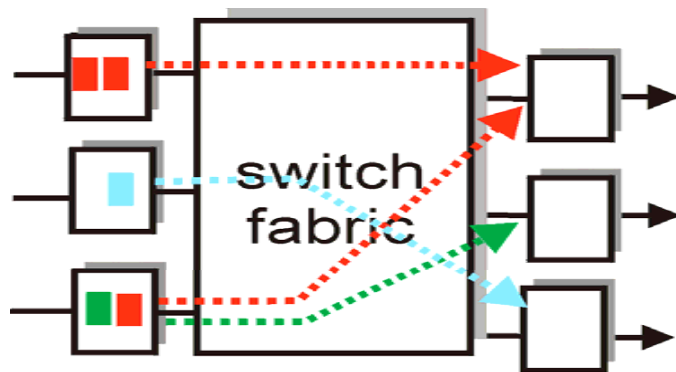
# Output port queueing



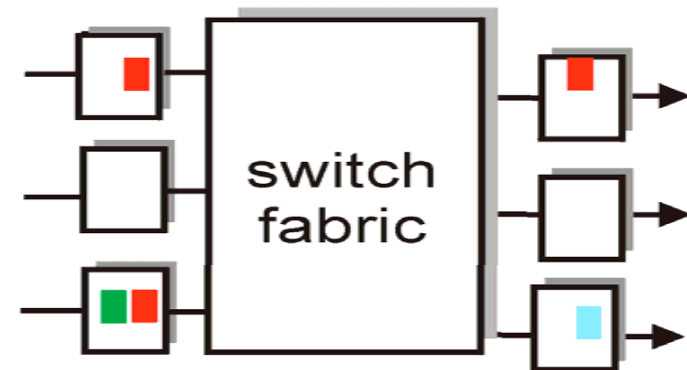
- buffering when arrival rate via switch exceeds output line speed
- *queueing (delay) and loss due to output port buffer overflow!*

# Input Port Queuing

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



output port contention  
at time t - only one red  
packet can be transferred



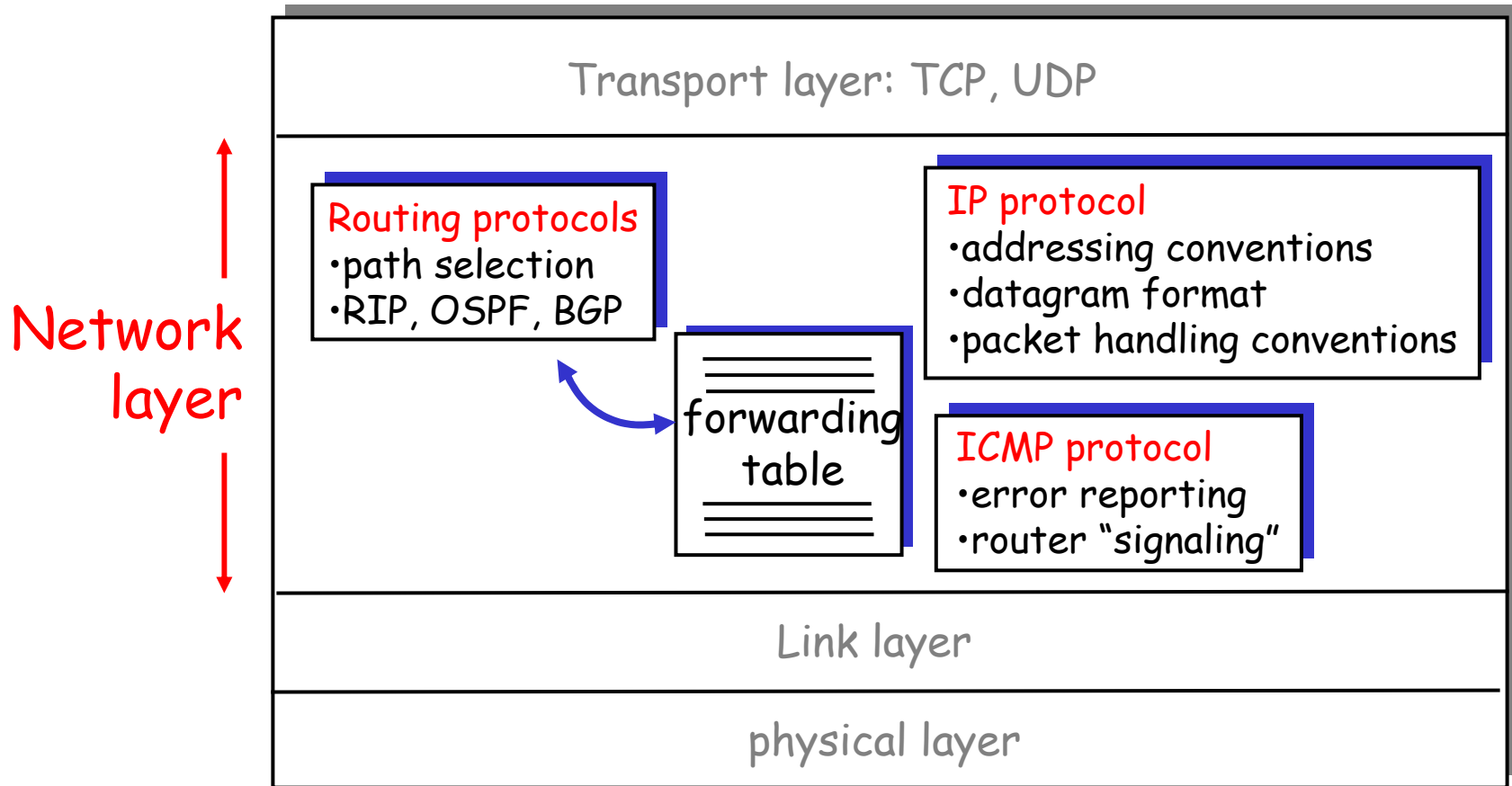
green packet  
experiences HOL blocking

# Chapter 4: Network Layer

- ❑ 4.1 Introduction
- ❑ 4.2 Virtual circuit and datagram networks
- ❑ 4.3 What's inside a router
- ❑ 4.4 IP: Internet Protocol
  - Datagram format
  - IPv4 addressing
  - ICMP
  - IPv6
- ❑ 4.5 Routing algorithms
  - Link state
  - Distance Vector
  - Hierarchical routing
- ❑ 4.6 Routing in the Internet
  - RIP
  - OSPF
  - BGP
- ❑ 4.7 Broadcast and multicast routing

# The Internet Network layer

Host, router network layer functions:



# Chapter 4: Network Layer

- ❑ 4.1 Introduction
- ❑ 4.2 Virtual circuit and datagram networks
- ❑ 4.3 What's inside a router
- ❑ 4.4 IP: Internet Protocol
  - Datagram format
  - IPv4 addressing
  - ICMP
  - IPv6
- ❑ 4.5 Routing algorithms
  - Link state
  - Distance Vector
  - Hierarchical routing
- ❑ 4.6 Routing in the Internet
  - RIP
  - OSPF
  - BGP
- ❑ 4.7 Broadcast and multicast routing

# IP datagram format

IP protocol version  
number

header length  
(bytes)

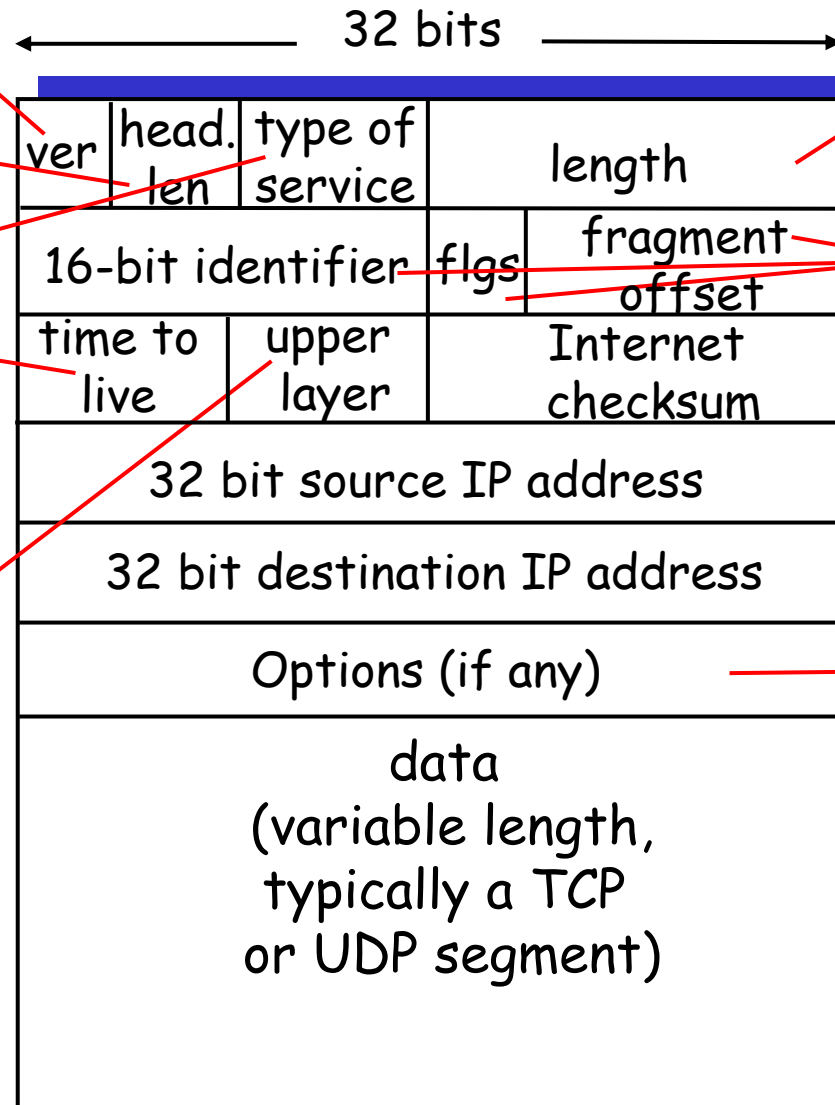
"type" of data

max number  
remaining hops  
(decremented at  
each router)

upper layer protocol  
to deliver payload to

## how much overhead with TCP?

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



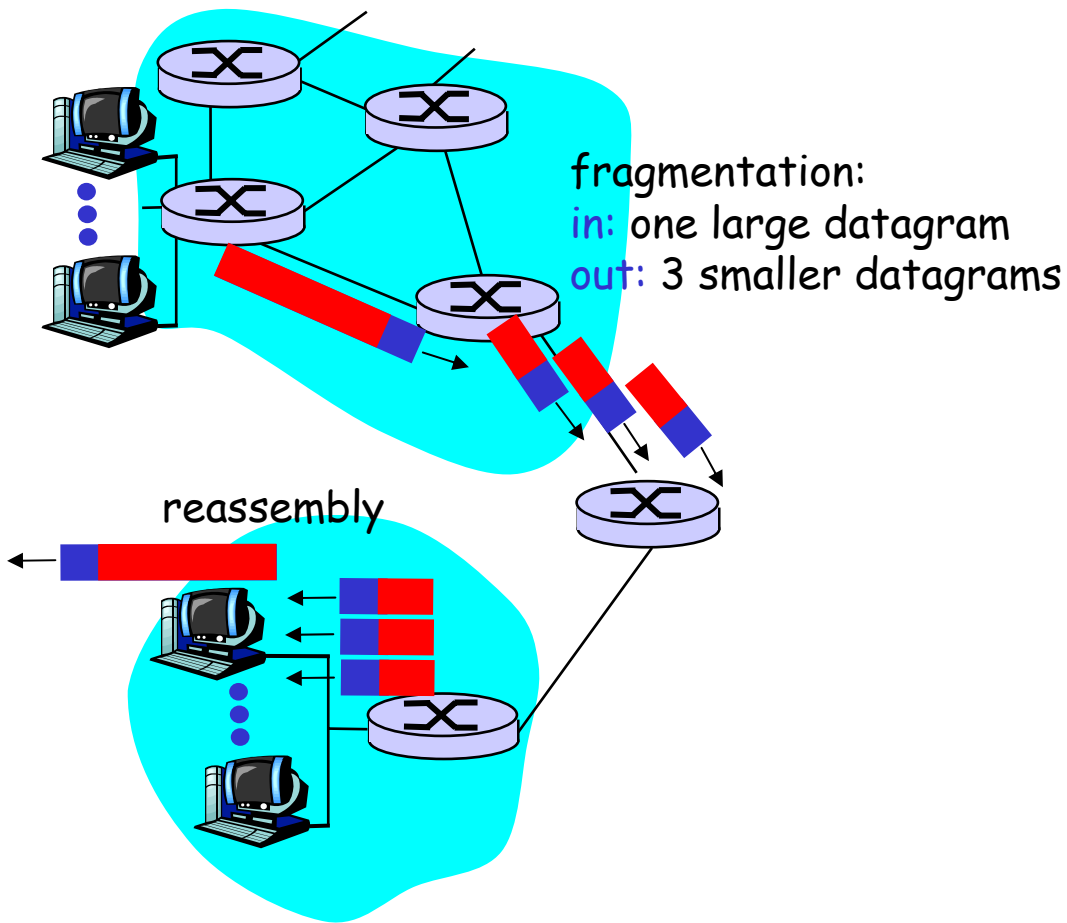
total datagram  
length (bytes)

for  
fragmentation/  
reassembly

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

# IP Fragmentation & Reassembly

- ❑ network links have MTU (max.transfer size) - largest possible link-level frame.
  - different link types, different MTUs
- ❑ large IP datagram divided ("fragmented") within net
  - one datagram becomes several datagrams
  - "reassembled" only at final destination
  - IP header bits used to identify, order related fragments





# IP Fragmentation and Reassembly

## Example

- ❑ 4000 byte datagram
- ❑ MTU = 1500 bytes

1480 bytes in data field

offset =  
 $1480/8$

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

One large datagram becomes several smaller datagrams

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

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

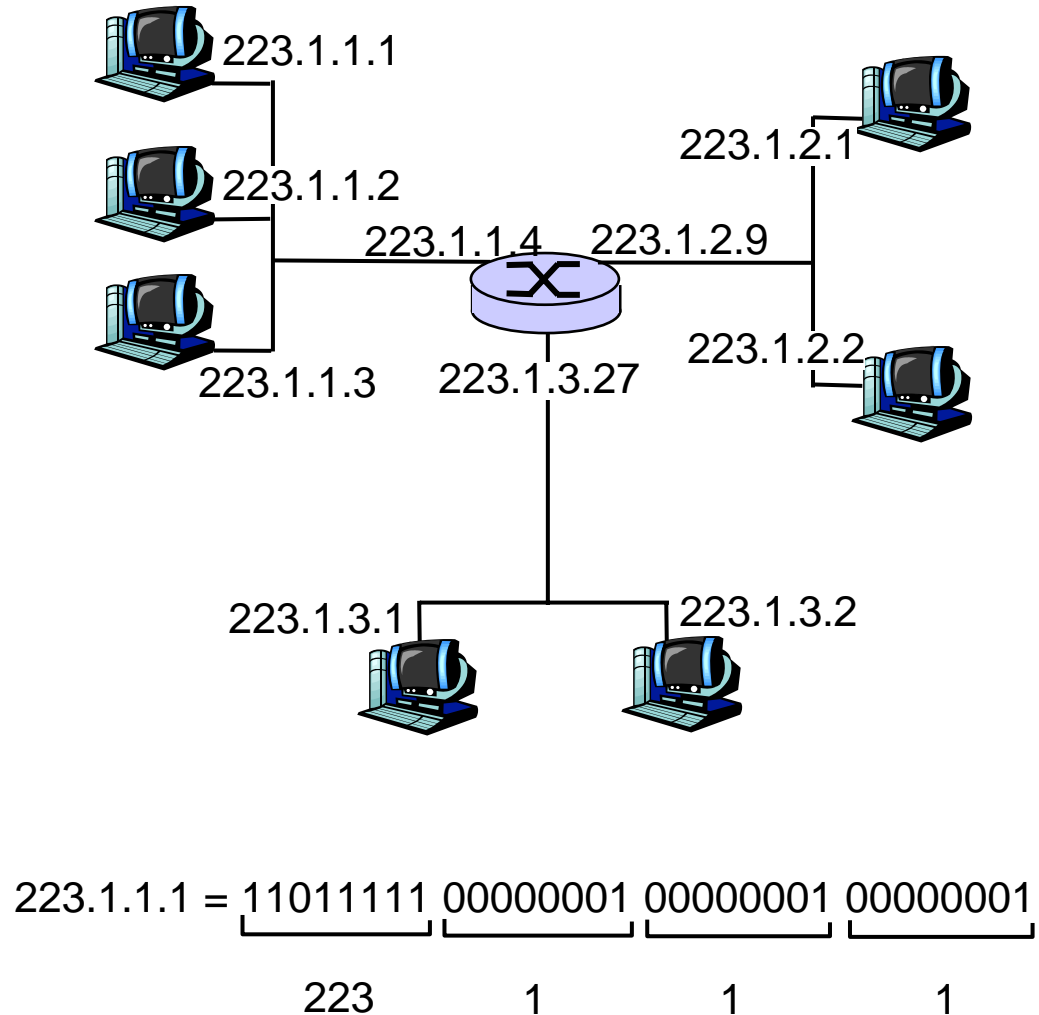
	length	ID	fragflag	offset
	=1040	=x	=0	=370

# Chapter 4: Network Layer

- ❑ 4.1 Introduction
- ❑ 4.2 Virtual circuit and datagram networks
- ❑ 4.3 What's inside a router
- ❑ 4.4 IP: Internet Protocol
  - Datagram format
  - IPv4 addressing
  - ICMP
  - IPv6
- ❑ 4.5 Routing algorithms
  - Link state
  - Distance Vector
  - Hierarchical routing
- ❑ 4.6 Routing in the Internet
  - RIP
  - OSPF
  - BGP
- ❑ 4.7 Broadcast and multicast routing

# 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 may have multiple interfaces
  - IP addresses associated with each *interface*, not with *host or router*



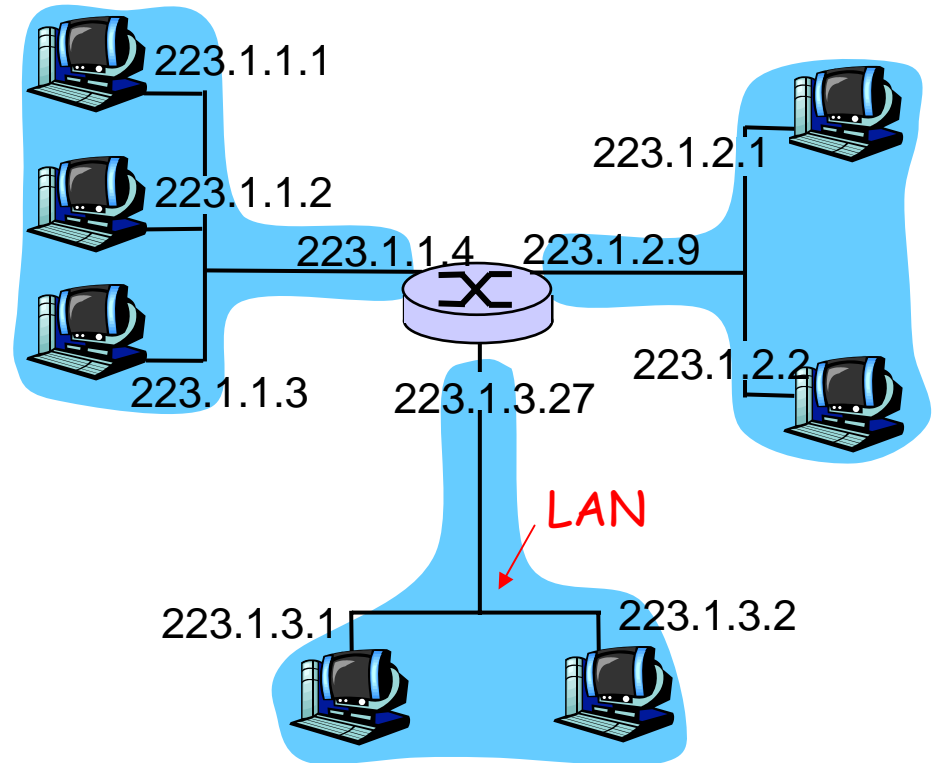
# Subnets

## ❑ 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

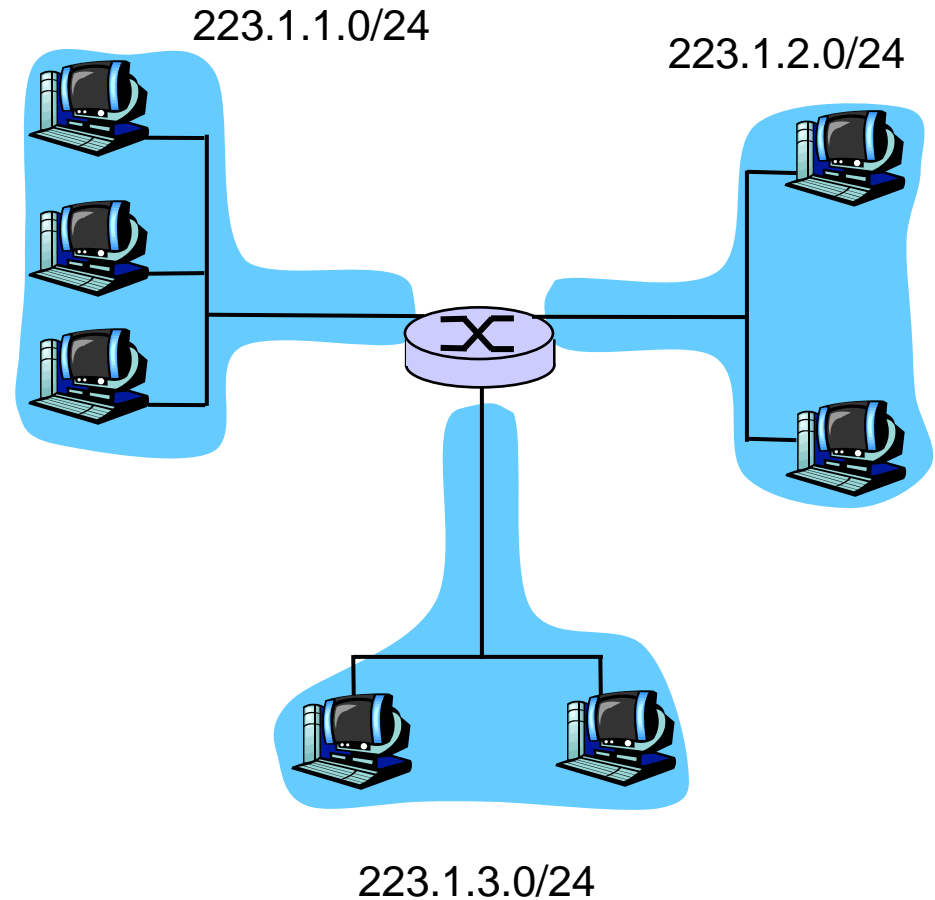


network consisting of 3 subnets  
(for IP addresses starting with 223,  
first 24 bits are network address)

# Subnets

## Recipe

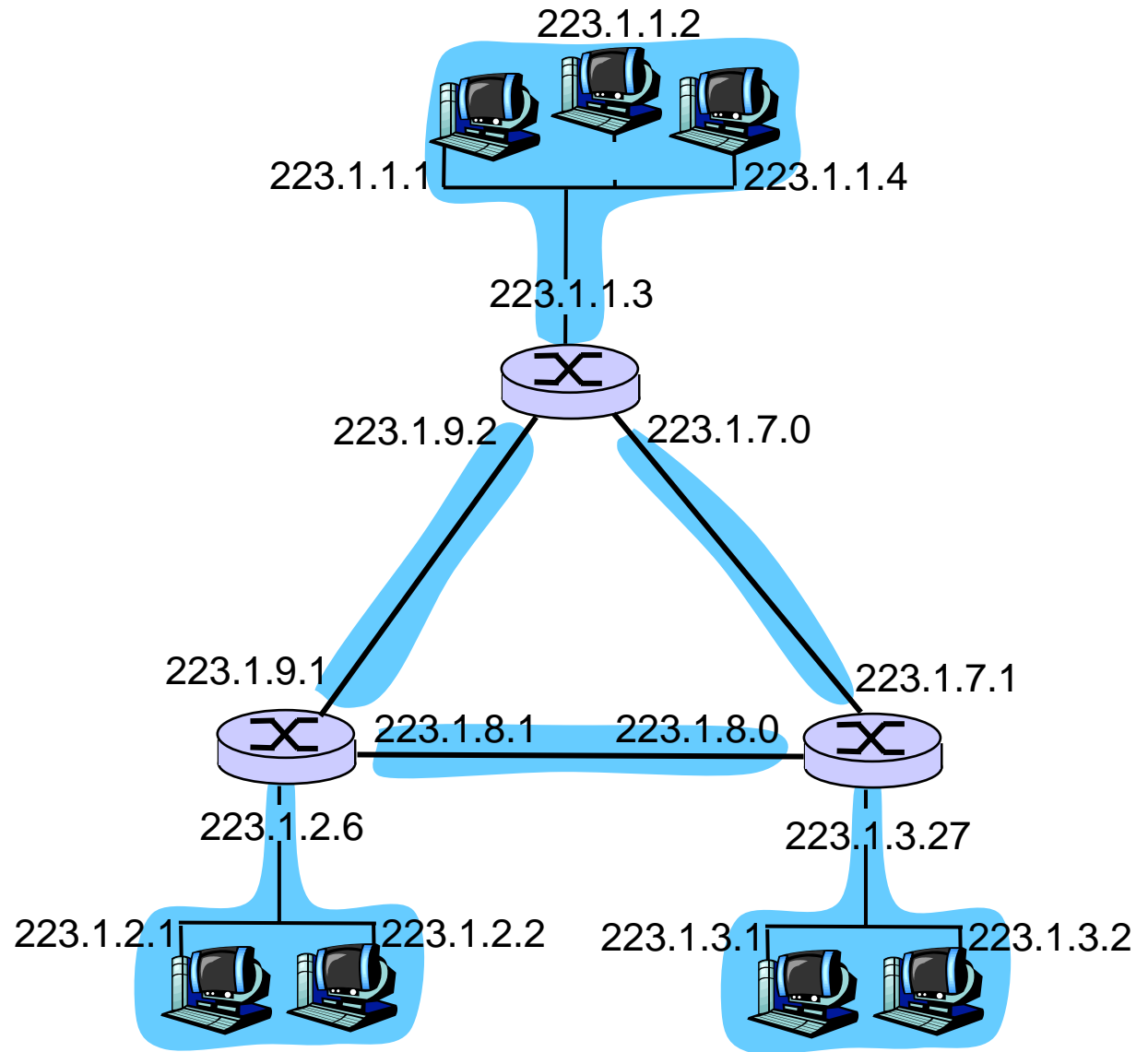
- ❑ To determine the subnets, detach each interface from its host or router, creating islands of isolated networks. Each isolated network is called a **subnet**.



Subnet mask: /24

# Subnets

How many?



# IP Addresses

given notion of "network", let's re-examine IP addresses:

"class-full" addressing:

class

A	0	network		host		1.0.0.0 to 127.255.255.255
B	10		network		host	128.0.0.0 to 191.255.255.255
C	110		network		host	192.0.0.0 to 223.255.255.255
D	1110		multicast address			224.0.0.0 to 239.255.255.255

← 32 bits →

# IP addressing: CIDR

## □ classful addressing:

- inefficient use of address space, address space exhaustion
- e.g., class B net allocated enough addresses for 65K hosts, even if only 2K hosts in that network

## □ CIDR: Classless InterDomain Routing

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



200.23.16.0/23



# IP addresses: how to get one?

Q: How does *host* get IP address?

- hard-coded by system admin in a file
  - Wintel: control-panel->network->configuration->tcp/ip->properties
  - UNIX: /etc/rc.config
- **DHCP: Dynamic Host Configuration Protocol:**  
dynamically get address from server
  - "plug-and-play"(more later)

# IP addresses: how to get one?

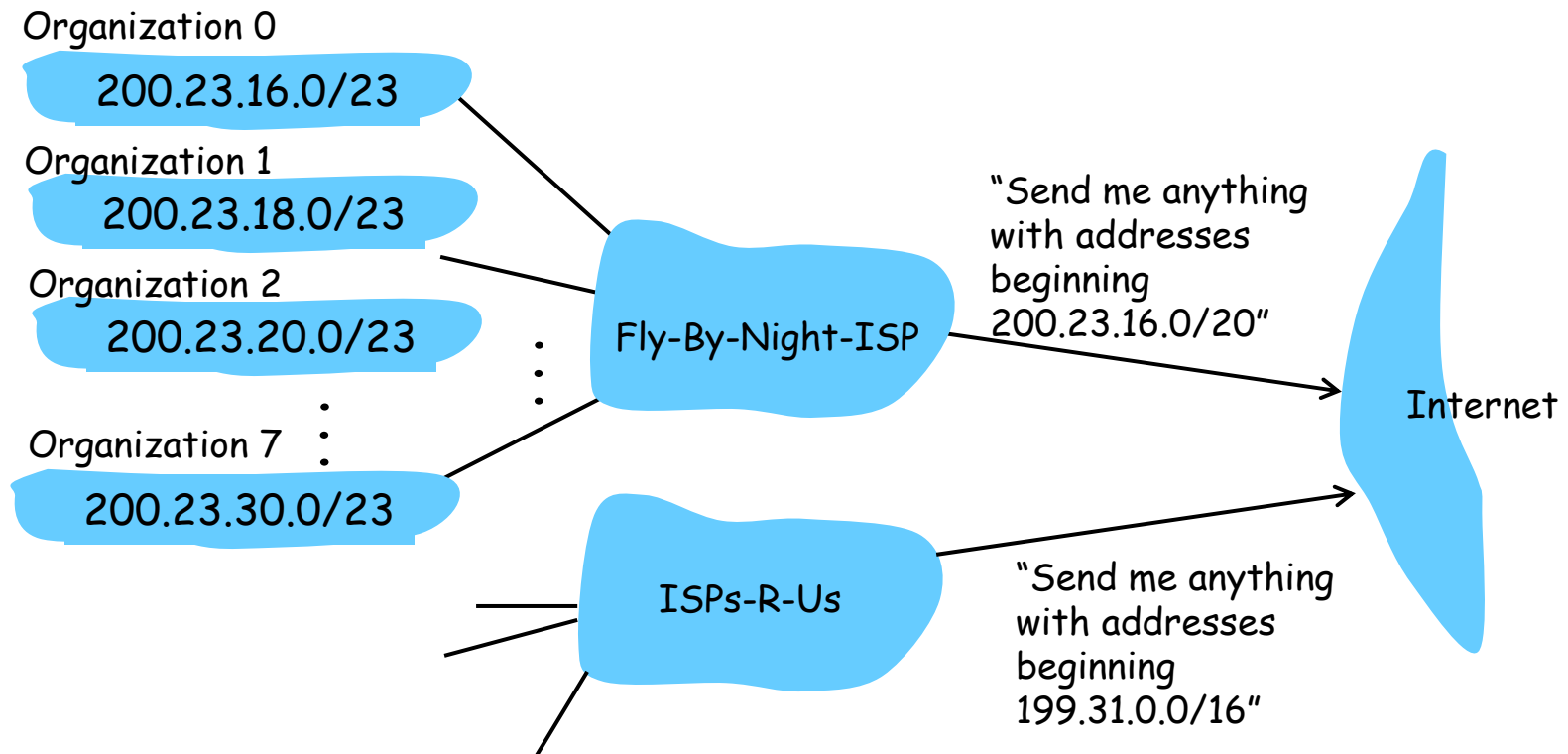
Q: How does *network* get subnet part of IP addr?

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

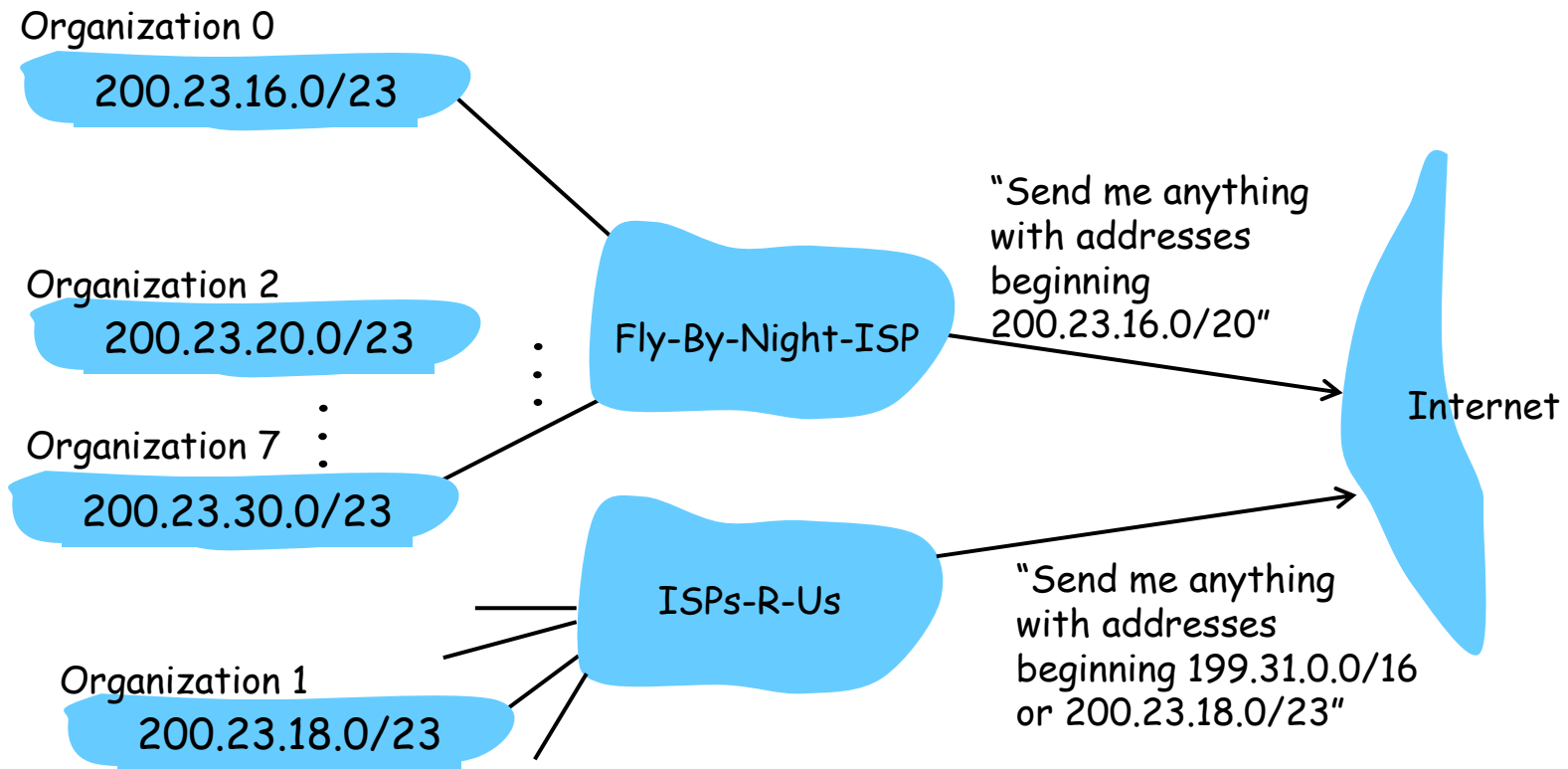
# Hierarchical addressing: route aggregation

Hierarchical addressing allows efficient advertisement of routing information:



# Hierarchical addressing: more specific routes

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



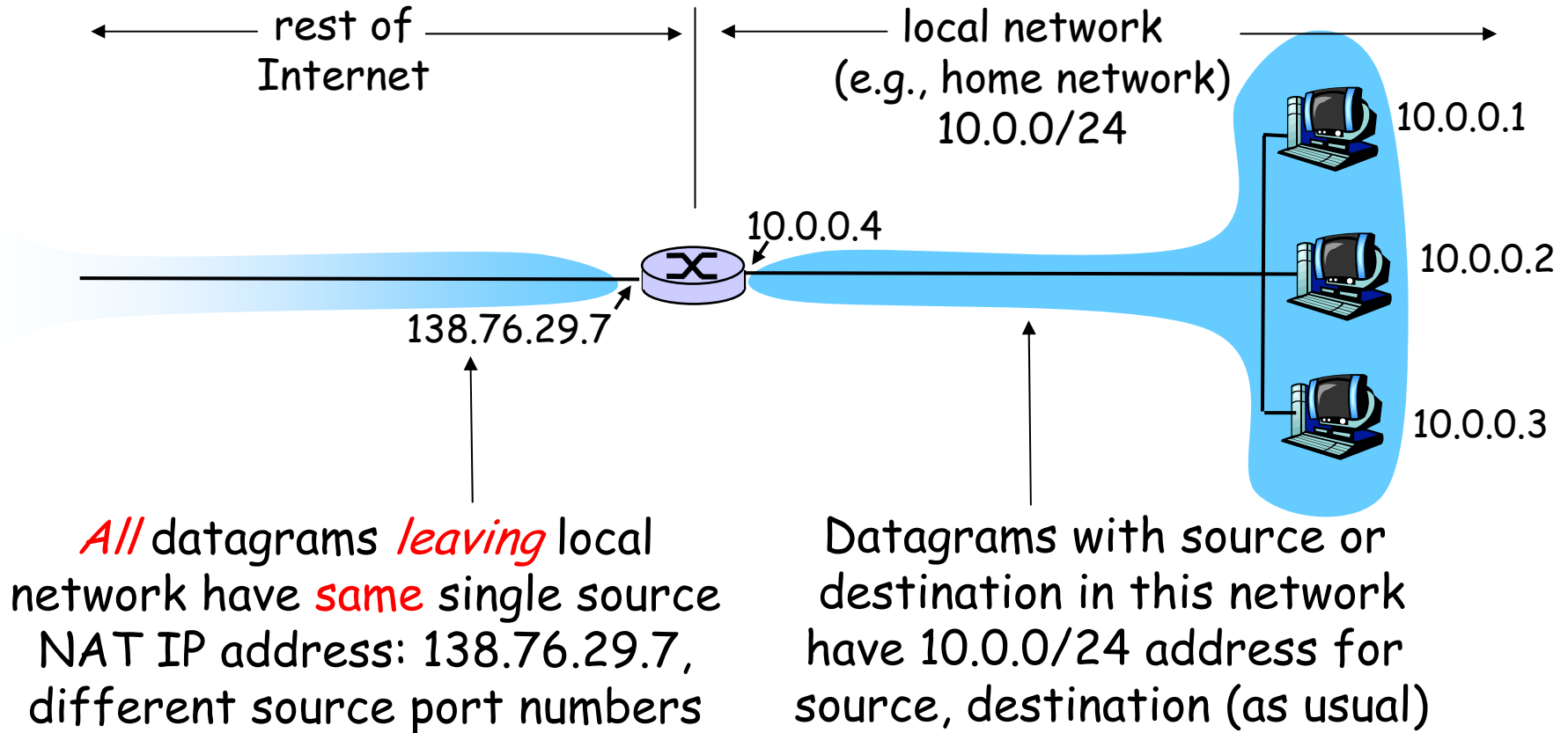
## IP addressing: the last word...

Q: How does an ISP get block of addresses?

A: **ICANN**: Internet **C**orporation for **A**ssigned  
**N**ames and **N**umbers

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

# NAT: Network Address Translation



# NAT: Network Address Translation

- **Motivation:** local network uses just one IP address as far as outside world is concerned:
  - no need to be allocated range of addresses from ISP:
    - just one IP address is used 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).

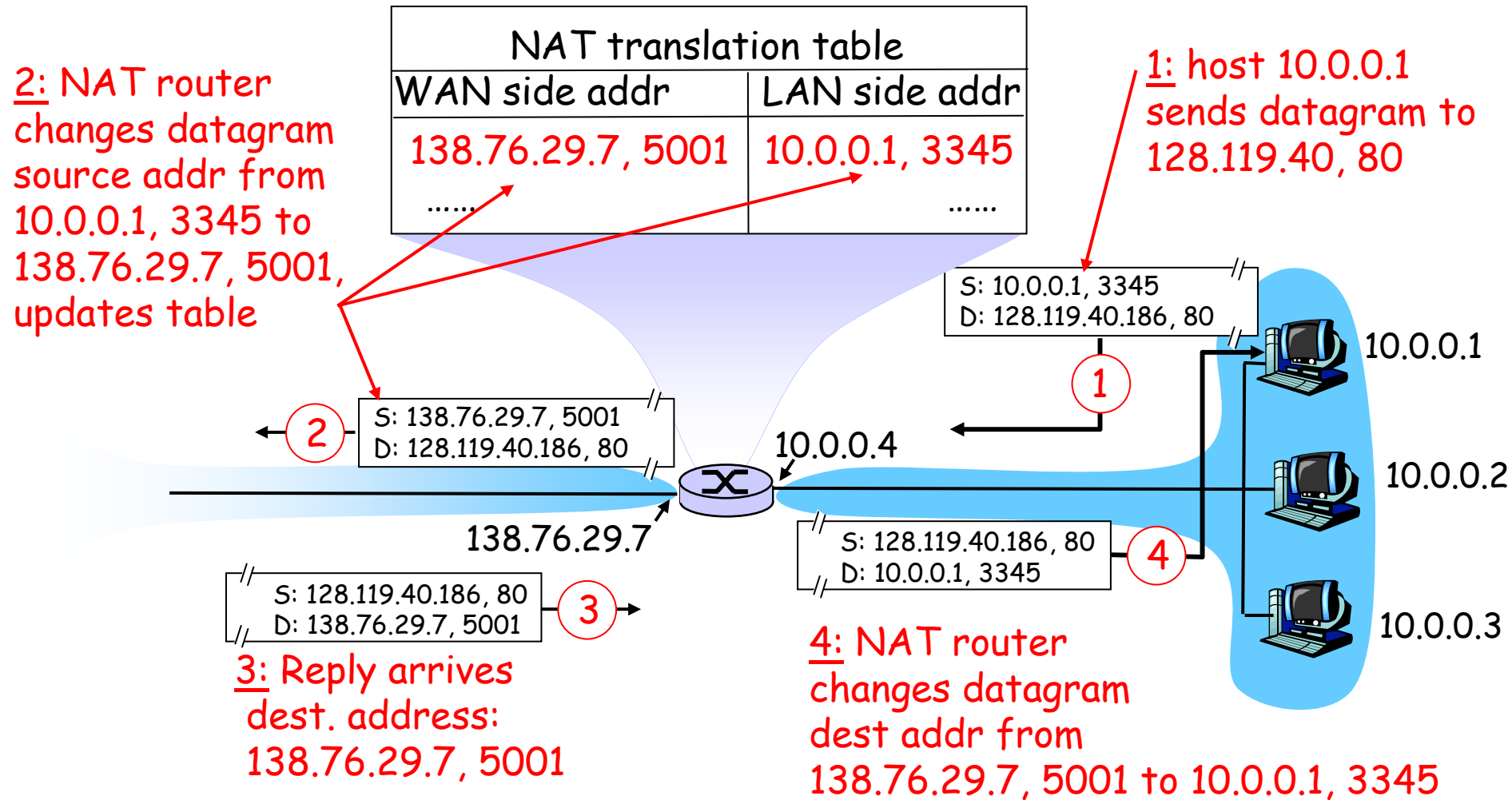
# NAT: Network Address Translation

**Implementation:** NAT router must:

- *in 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
- *in 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



# 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
  - violates end-to-end argument
    - NAT possibility must be taken into account by app designers, eg, P2P applications
  - address shortage should instead be solved by IPv6

# Chapter 4: Network Layer

- ❑ 4.1 Introduction
- ❑ 4.2 Virtual circuit and datagram networks
- ❑ 4.3 What's inside a router
- ❑ 4.4 IP: Internet Protocol
  - Datagram format
  - IPv4 addressing
  - ICMP
  - IPv6
- ❑ 4.5 Routing algorithms
  - Link state
  - Distance Vector
  - Hierarchical routing
- ❑ 4.6 Routing in the Internet
  - RIP
  - OSPF
  - BGP
- ❑ 4.7 Broadcast and multicast routing

# ICMP: Internet Control Message Protocol

- ❑ used by hosts & routers to communicate network-level information
  - error reporting: unreachable host, network, port, protocol
  - echo request/reply (used by ping)
- ❑ network-layer "above" IP:
  - ICMP msgs carried in IP datagrams
- ❑ **ICMP message:** type, code plus first 8 bytes of IP datagram causing error

<u>Type</u>	<u>Code</u>	<u>description</u>
0	0	echo reply (ping)
3	0	dest. network unreachable
3	1	dest host unreachable
3	2	dest protocol unreachable
3	3	dest port unreachable
3	6	dest network unknown
3	7	dest host unknown
4	0	source quench (congestion control - not used)
8	0	echo request (ping)
9	0	route advertisement
10	0	router discovery
11	0	TTL expired
12	0	bad IP header

# Traceroute and ICMP

- ❑ Source sends series of UDP segments to dest
  - First has TTL =1
  - Second has TTL=2, etc.
  - Unlikely port number
- ❑ When nth datagram arrives to nth router:
  - Router discards datagram
  - And sends to source an ICMP message (type 11, code 0)
  - Message includes name of router & IP address

- ❑ When ICMP message arrives, source calculates RTT
- ❑ Traceroute does this 3 times

## Stopping criterion

- ❑ UDP segment eventually arrives at destination host
- ❑ Destination returns ICMP "host unreachable" packet (type 3, code 3)
- ❑ When source gets this ICMP, stops.

# Chapter 4: Network Layer

- ❑ 4.1 Introduction
- ❑ 4.2 Virtual circuit and datagram networks
- ❑ 4.3 What's inside a router
- ❑ 4.4 IP: Internet Protocol
  - Datagram format
  - IPv4 addressing
  - ICMP
  - IPv6
- ❑ 4.5 Routing algorithms
  - Link state
  - Distance Vector
  - Hierarchical routing
- ❑ 4.6 Routing in the Internet
  - RIP
  - OSPF
  - BGP
- ❑ 4.7 Broadcast and multicast routing

# IPv6

- ❑ **Initial motivation:** 32-bit address space soon to be completely allocated.
  - ❑ **Additional motivation:**
    - header format helps speed processing/forwarding
    - header changes to facilitate QoS
- IPv6 datagram format:**
- fixed-length 40 byte header
  - no fragmentation allowed

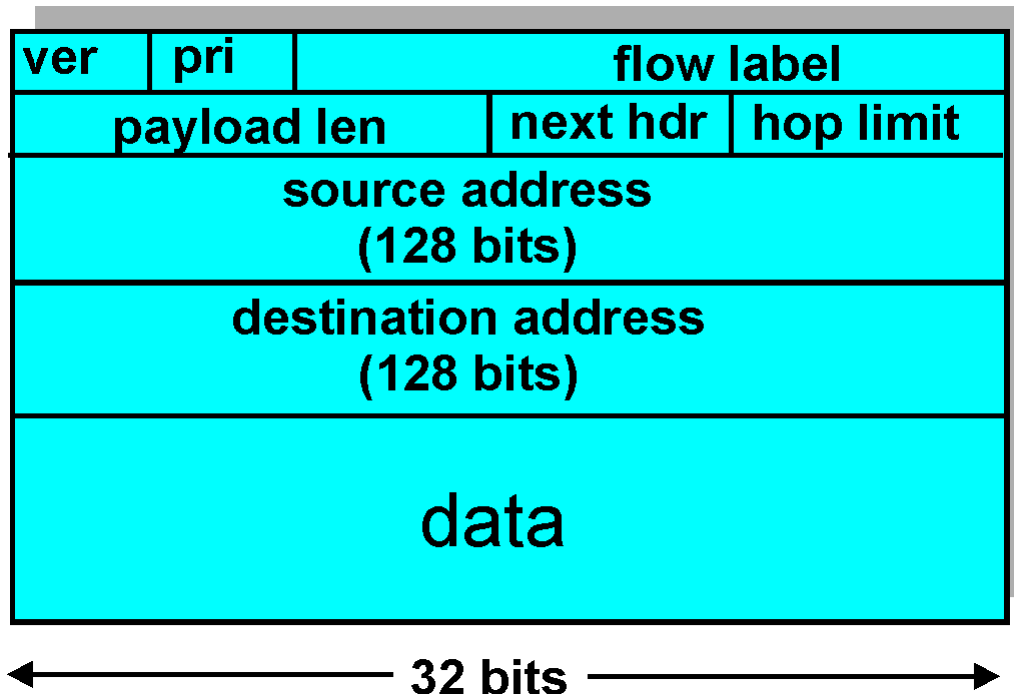
# IPv6 Header (Cont)

*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





# Other Changes from IPv4

- ❑ *Checksum*: removed entirely to reduce processing time at each hop
- ❑ *Options*: allowed, but outside of header, indicated by "Next Header" field
- ❑ *ICMPv6*: new version of ICMP
  - additional message types, e.g. "Packet Too Big"
  - multicast group management functions

# Transition From IPv4 To IPv6

- ❑ Not all routers can be upgraded simultaneous
  - no “flag days”
  - How will the network operate with mixed IPv4 and IPv6 routers?
- ❑ Two proposed approaches:
  - *Dual Stack*: some routers with dual stack (v6, v4) can “translate” between formats
  - *Tunneling*: IPv6 carried as payload in IPv4 datagram among IPv4 routers

# Tunneling

