

The Network Layer

Goal

- To provide the transport layer the service of getting the packets from source to destination hosts.
 - May travel **multiple** hops.
 - The lowest layer that deals with *end-to-end transmission*.

The Network Layer - Four major issues

- Services to Transport Layer
- Routing
 - “to make 'appropriate' routes for packets with the consideration of meeting performance requirements for the network and the user.”
- Congestion Control
 - To avoid an overload of packets at one or more switching nodes (e.g., gateways)
- Internetworking
 - To provide transparent communication between source and destination machines across multiple physical networks.

Chapter 4: Network Layer

Chapter goals:

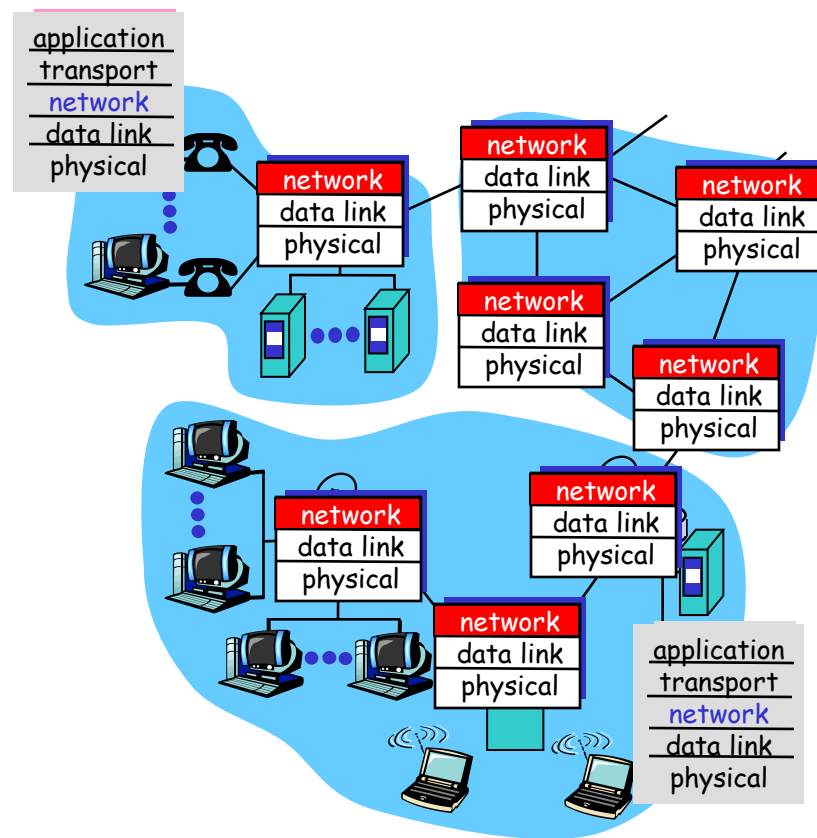
- understand principles behind network layer services:
 - routing (path selection)
 - dealing with scale
 - how a router works
 - advanced topics: IPv6, mobility
- instantiation and implementation in the Internet

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 functions

- 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



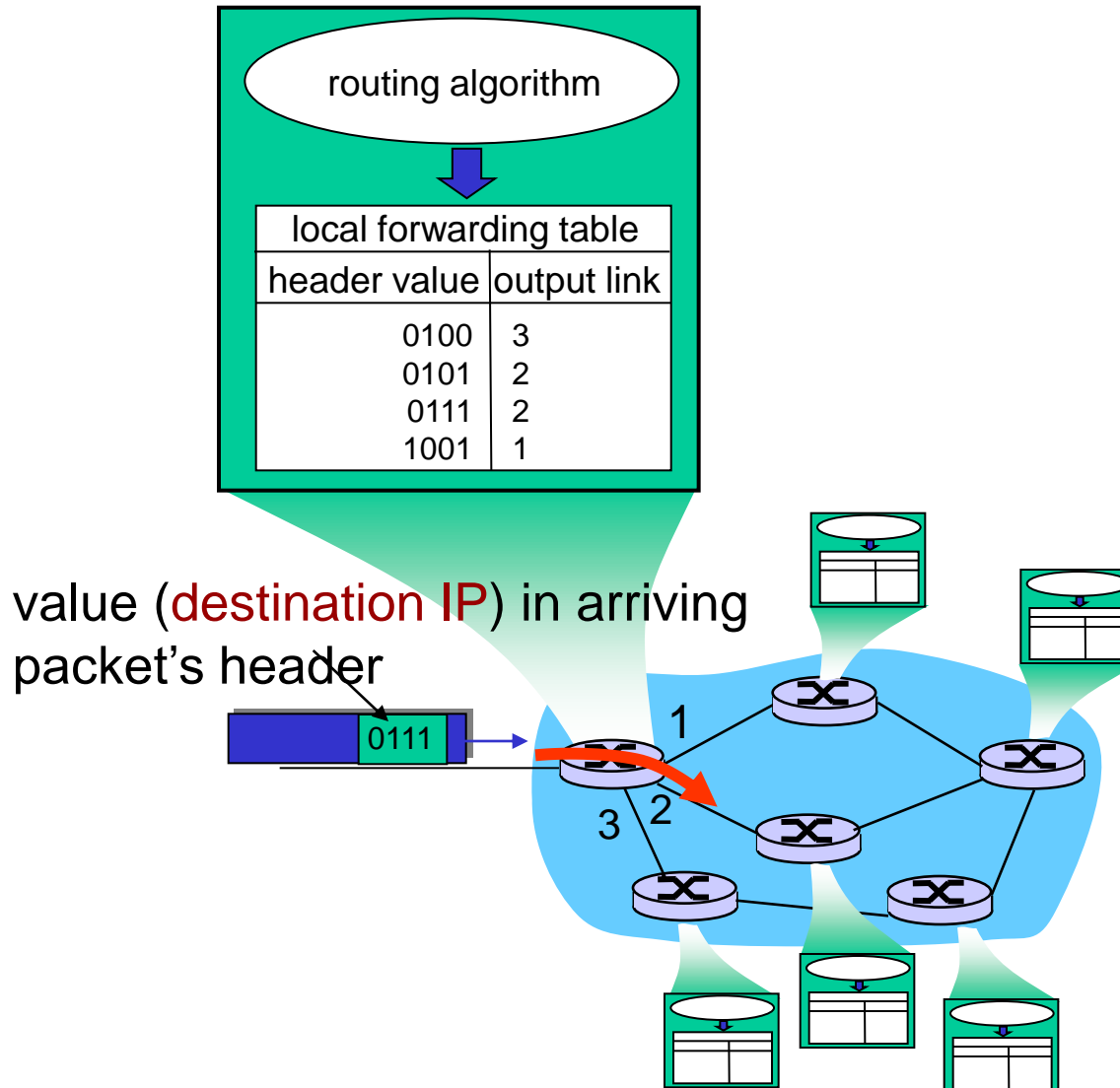
Two Key Network-Layer Functions

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

analogy:

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

Interplay between routing and forwarding



Connection setup

- 3rd 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 and transport layer connection service:
 - **Network**: between two hosts
 - **Transport**: between two processes

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

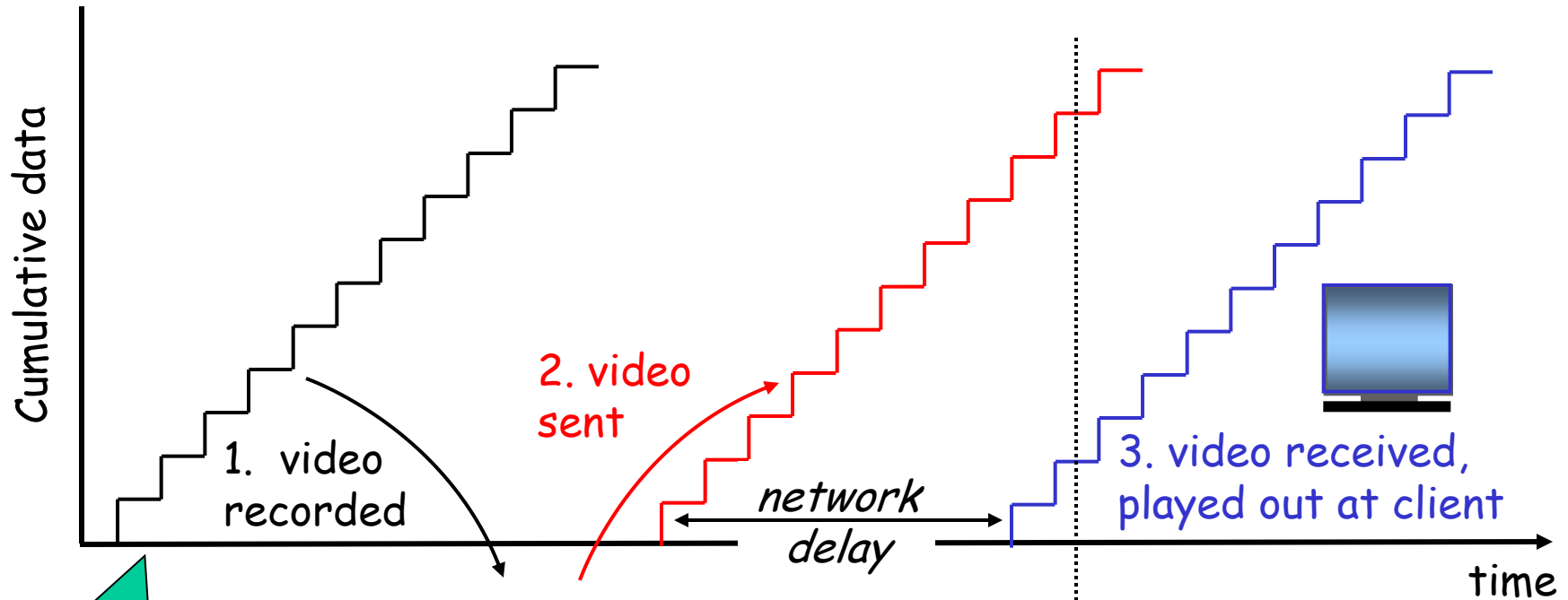
Network layer service models

Network Architecture	Service Model	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

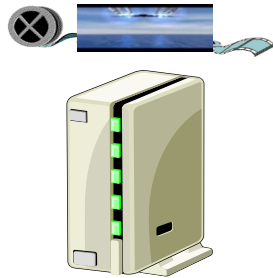
- Internet model being extended: IntServ, DiffServ
 - Chapter 6



Streaming Stored Multimedia: ideal case

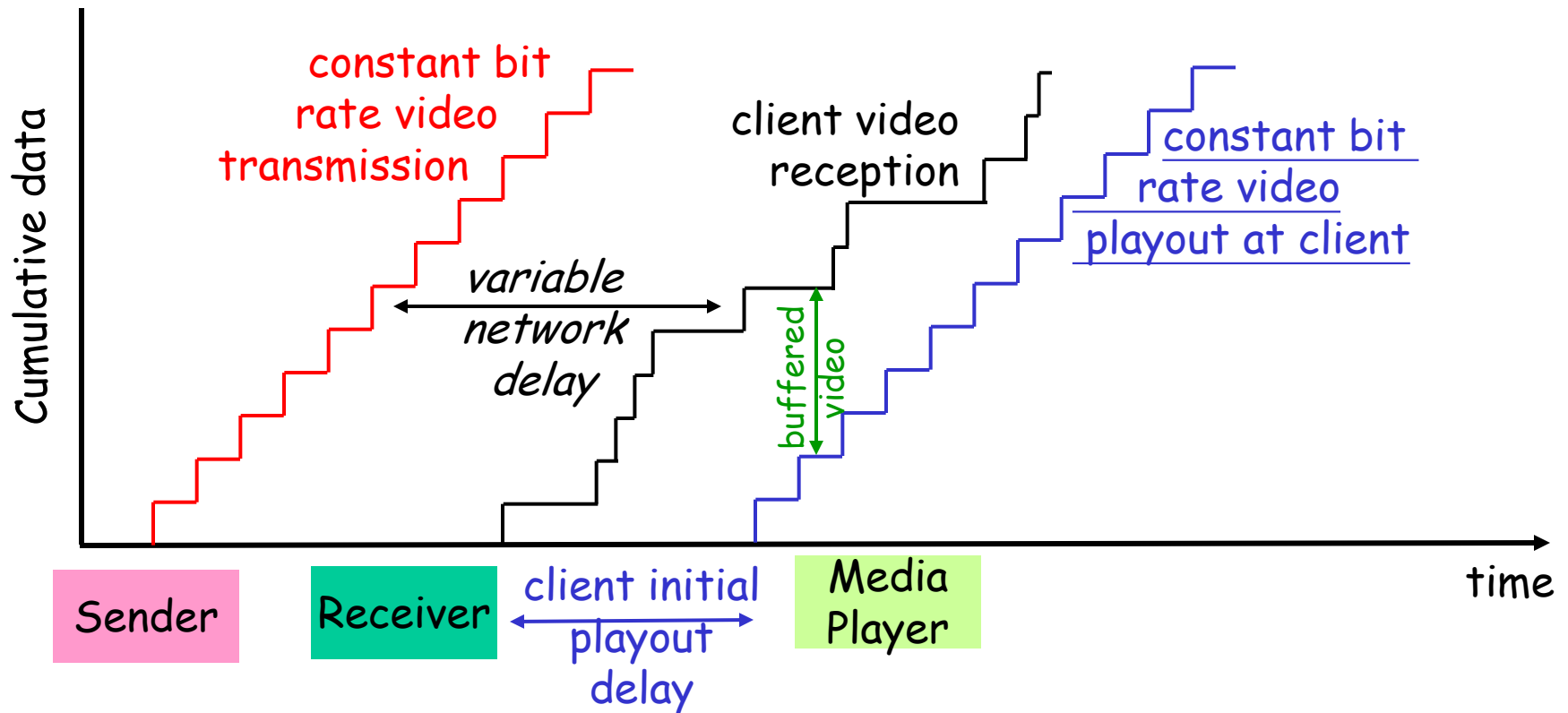


Constant rate
data (frame)
Generation
(e.g., 20
frames/sec)

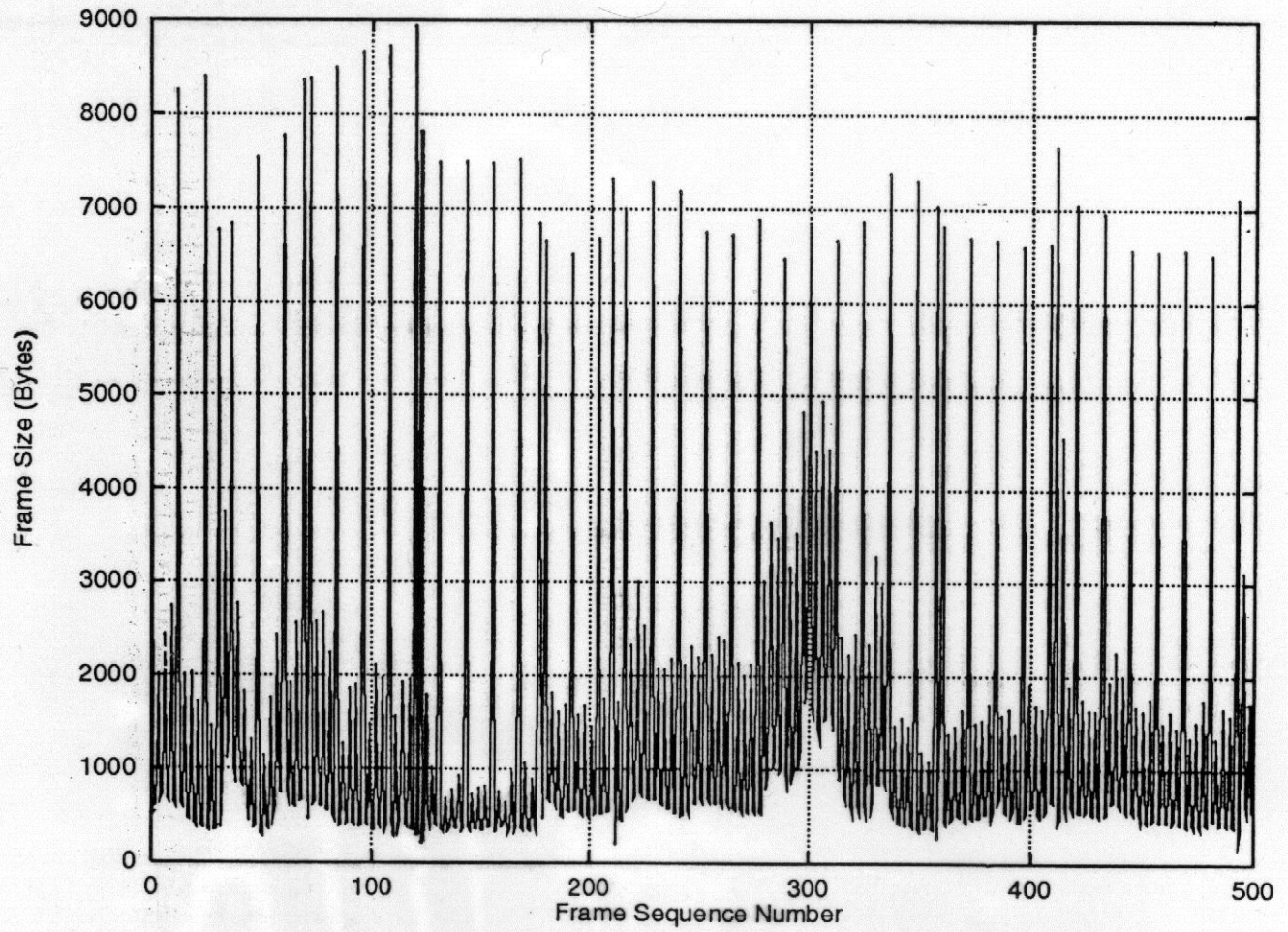


streaming: at this time, client
playing out early part of video,
while server still sending later
part of video

Streaming Multimedia: Client Buffering



- Client-side buffering: playout delay compensates for network-added delay and delay jitters.



"Star War"



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 network core

Virtual circuits

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

- call setup, teardown for each call *before* data can flow
- each packet carries **VC identifier** (NOT destination host ID)
- *every* router on source-dest path maintains “**state**” for each passing connection
- *link, router resources (bandwidth, buffers) may be allocated to VC*
 - to get circuit-like performance.

Virtual circuits

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

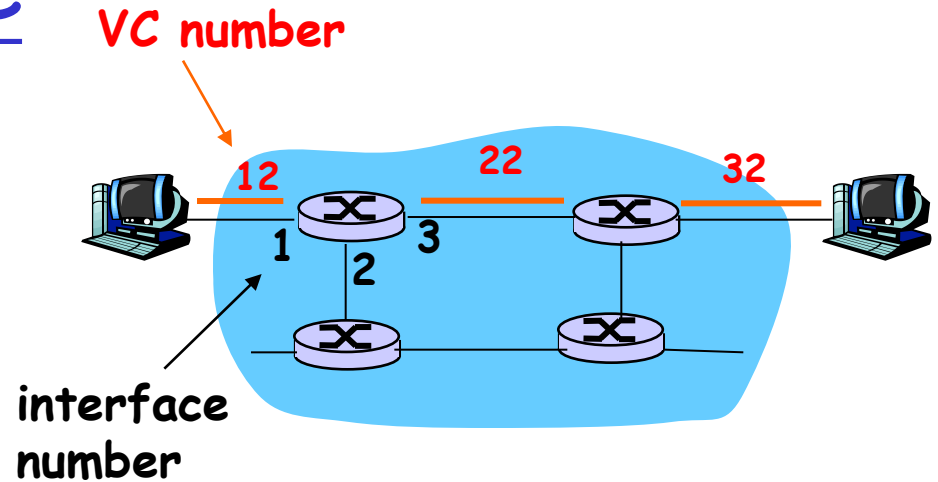
- performance-wise
 - network actions along source-to-dest path
-
- call setup, teardown for each call *before* data can flow
 - each packet carries VC identifier (not destination host address)
 - *every* router on source-dest path maintains “state” for each passing connection
 - link, router resources (bandwidth, buffers) may be *allocated* to VC (dedicated resources = predictable service)

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 VC number (rather than dest address)
 - VC number can be changed on each link.
 - New VC number comes from forwarding table

Forwarding table



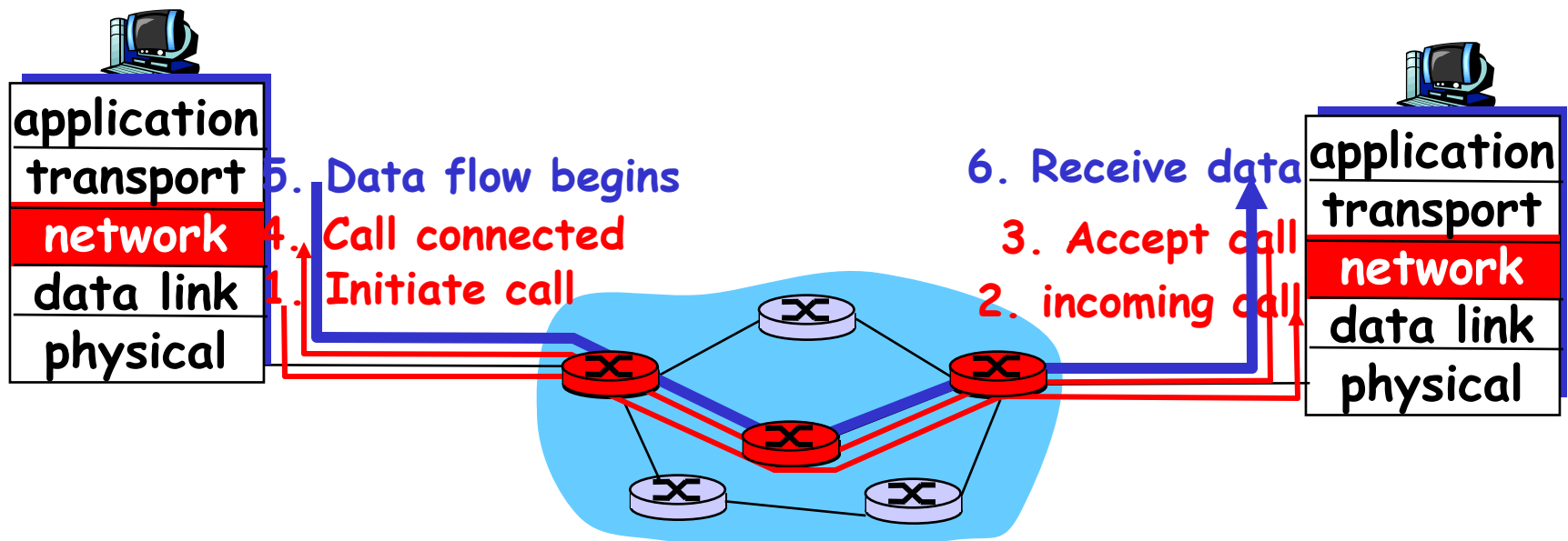
Forwarding table in northwest router:

Incoming interface	Incoming VC #	Outgoing interface	Outgoing VC #
1	12		3
2	63		1
3	7		2
1	97		3
...

Routers maintain connection state information!

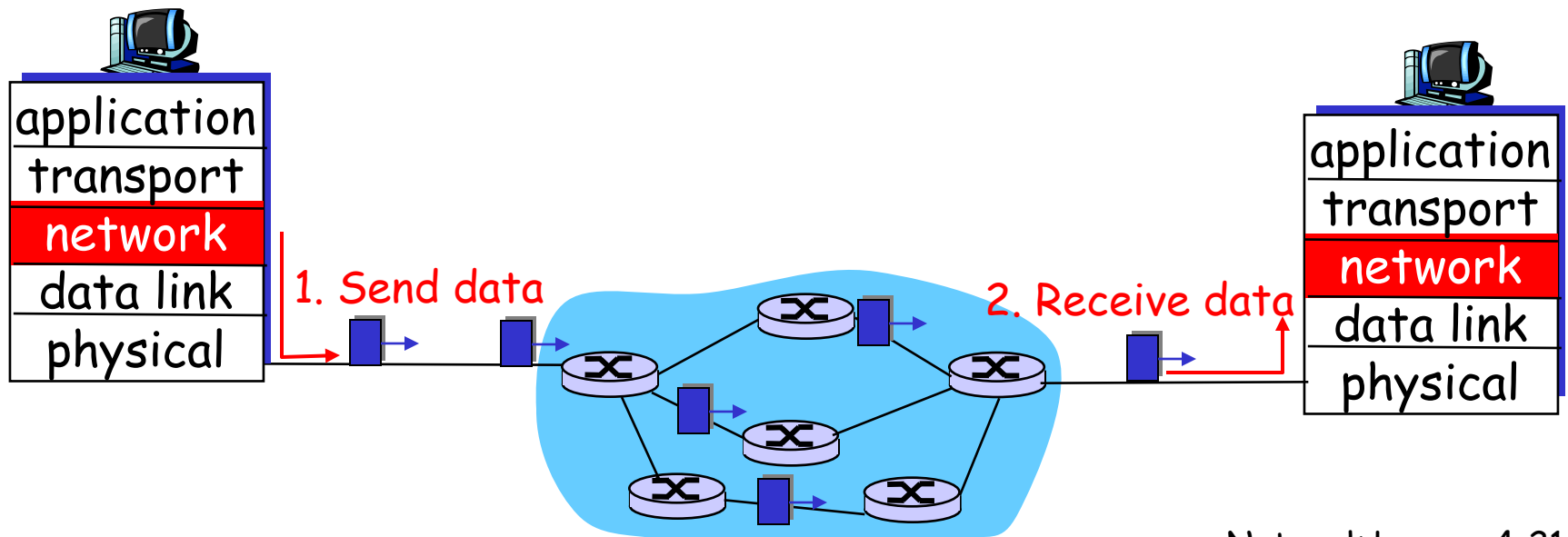
Virtual circuits: signaling protocols

- used to setup, maintain teardown VC
- used in ATM, frame-relay, X.25
- not used in today's Internet



Datagram networks: the Internet model

- **NO** call setup at network layer
- routers: **NO** state 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 (datagram)

- 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 (virtual circuit)

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

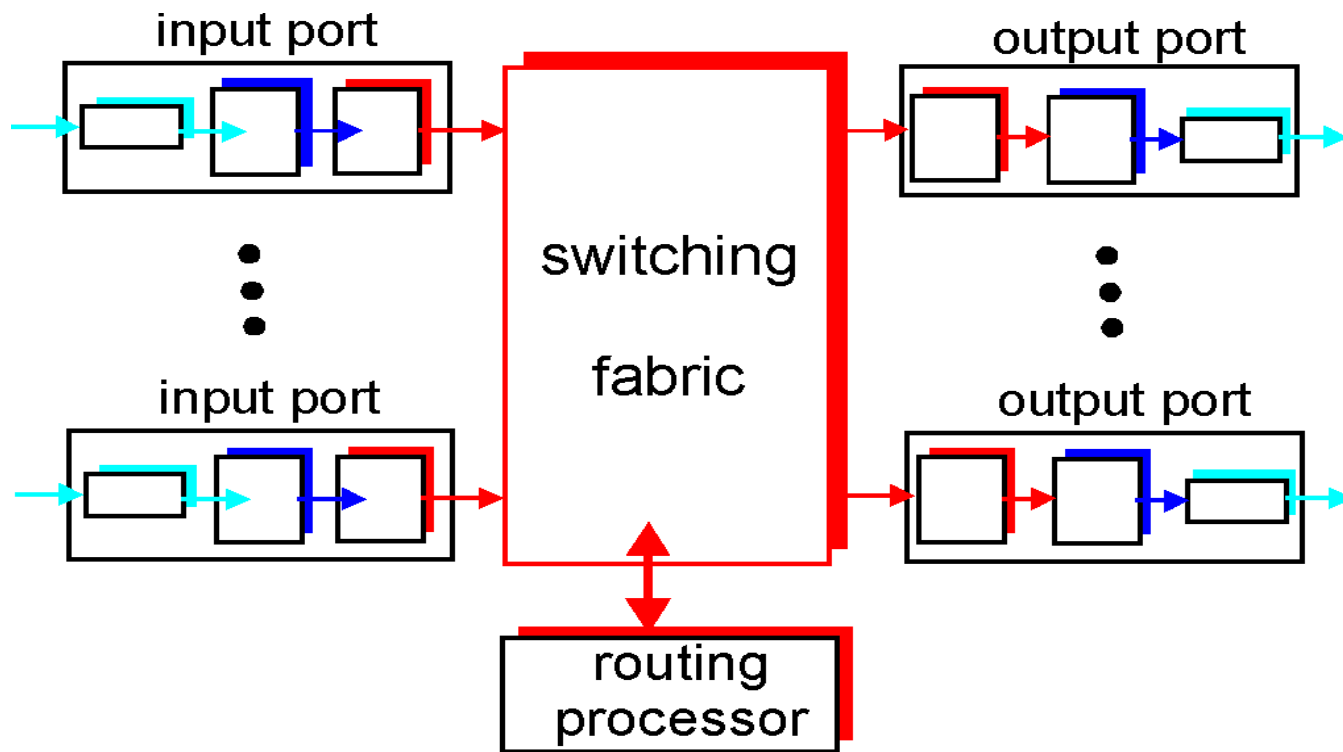
Chapter 4: Network Layer

- 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

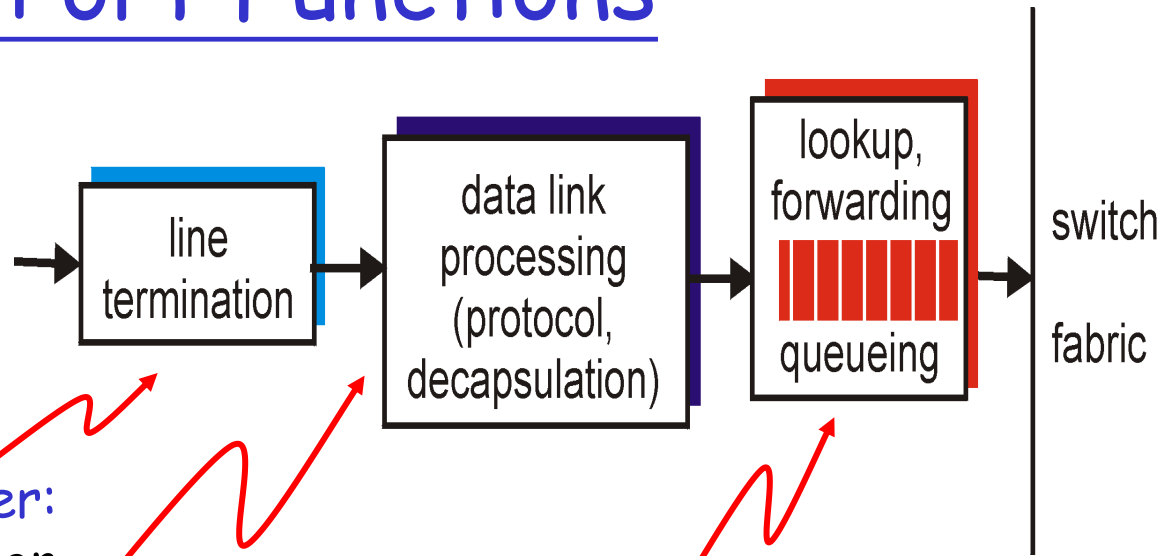
Router Architecture Overview

Two key router functions:

- run routing algorithms/protocol (RIP, OSPF, BGP) - to build routing table
- *switching* datagrams from incoming to outgoing link - packet forwarding by looking up routing table



Input Port Functions



Physical layer:
bit-level reception

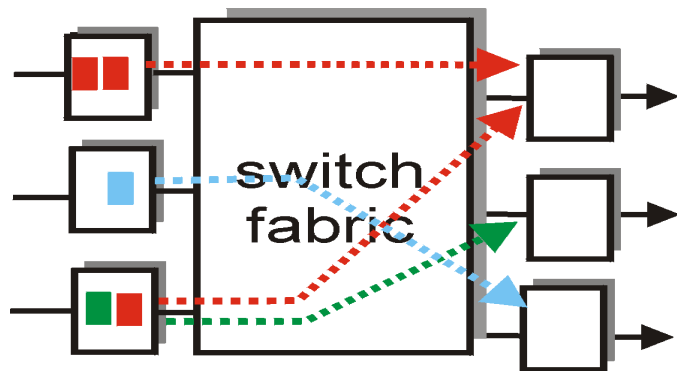
Data link layer:
e.g., Ethernet
see chapter 5

Decentralized switching:

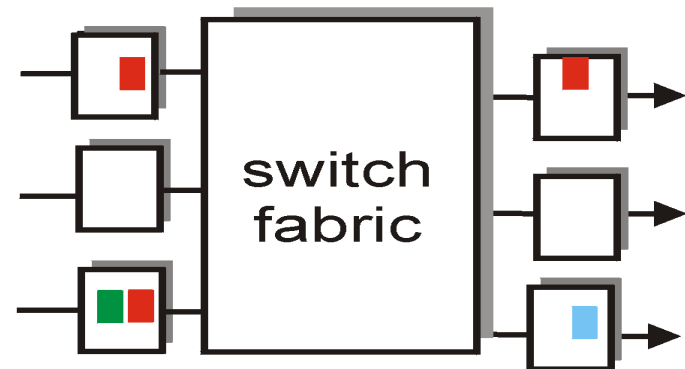
- given datagram dest., lookup output port using routing table in input port memory
- goal: complete input port processing at 'line speed'
- queuing: if datagrams arrive faster than forwarding rate into switch fabric

Input Port Queuing

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

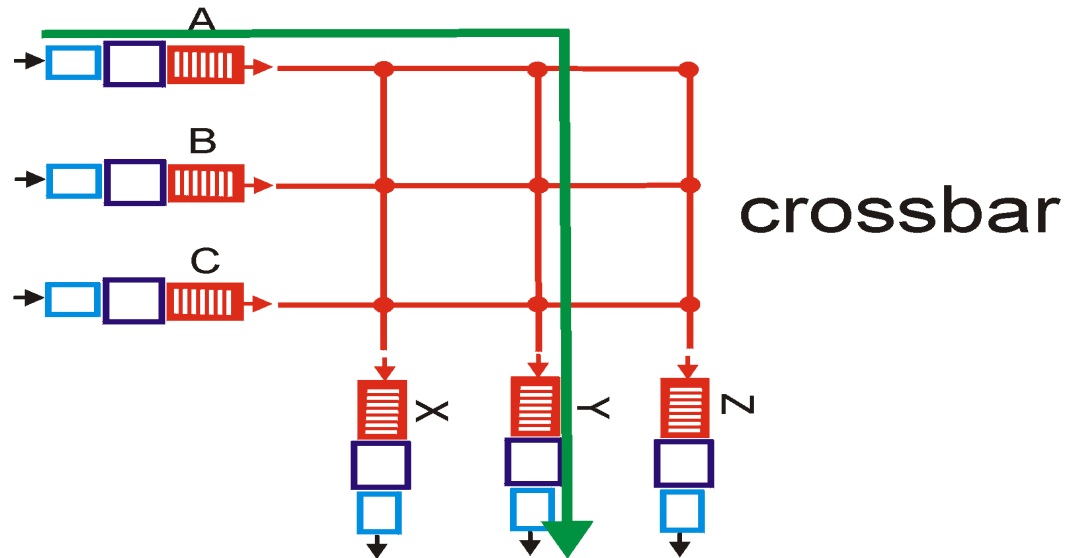
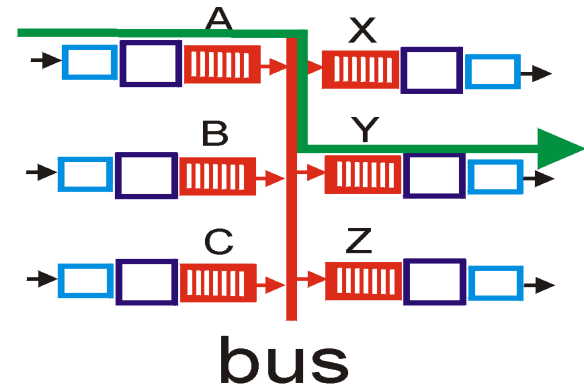
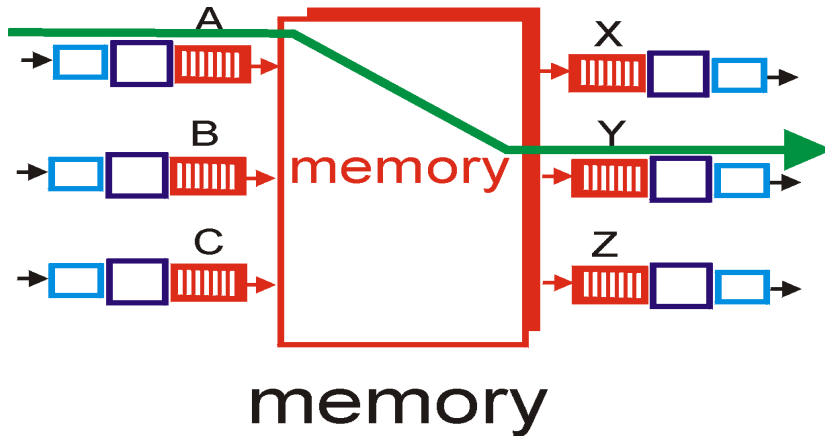


output port contention at time t - only one red packet can be transferred



green packet experiences HOL blocking

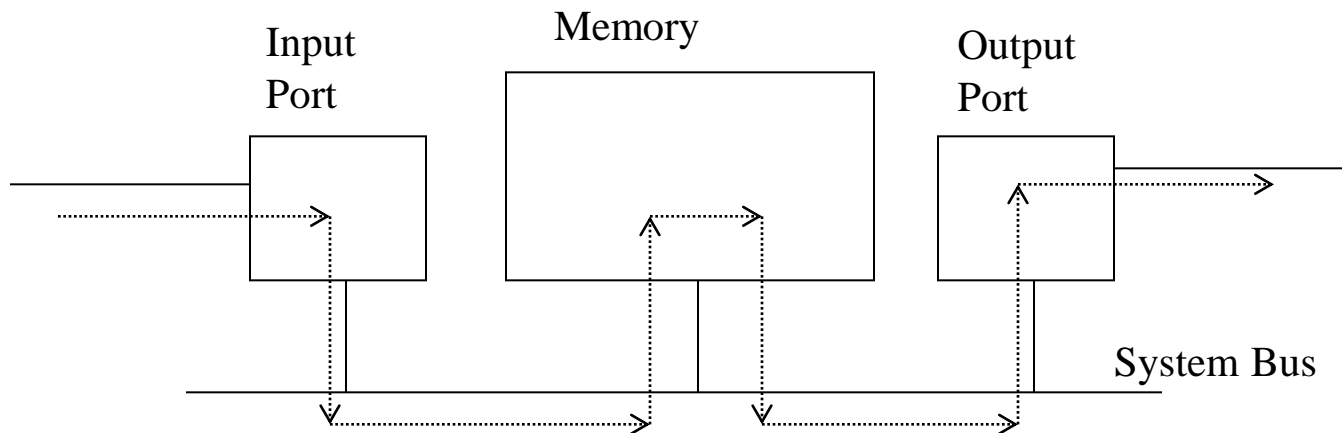
Three types of switching fabrics



Switching via Memory

First generation routers:

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

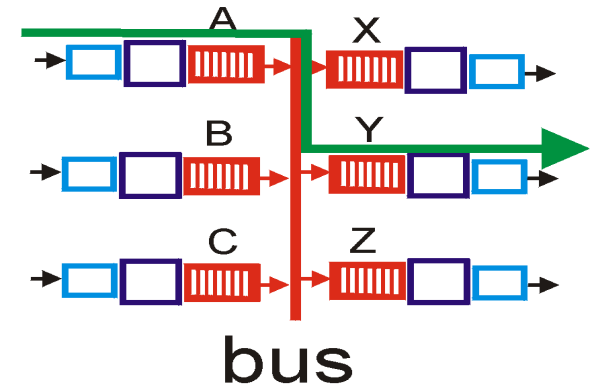


Modern routers:

- input port processor performs lookup, copy into memory
- Cisco Catalyst 6500

Switching via a Bus

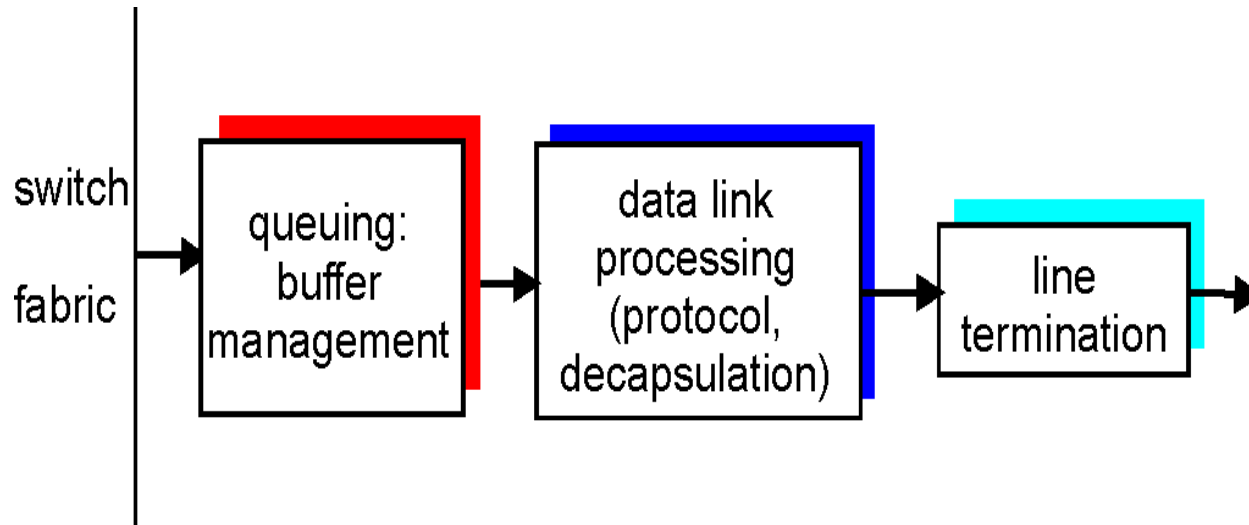
- Datagram from input port memory to output port memory via a **shared bus**
- **bus contention**: switching speed limited by bus bandwidth
- 1 Gbps bus, Cisco 1900: sufficient speed for access and enterprise routers (not regional or backbone)



Switching via an Interconnection Network

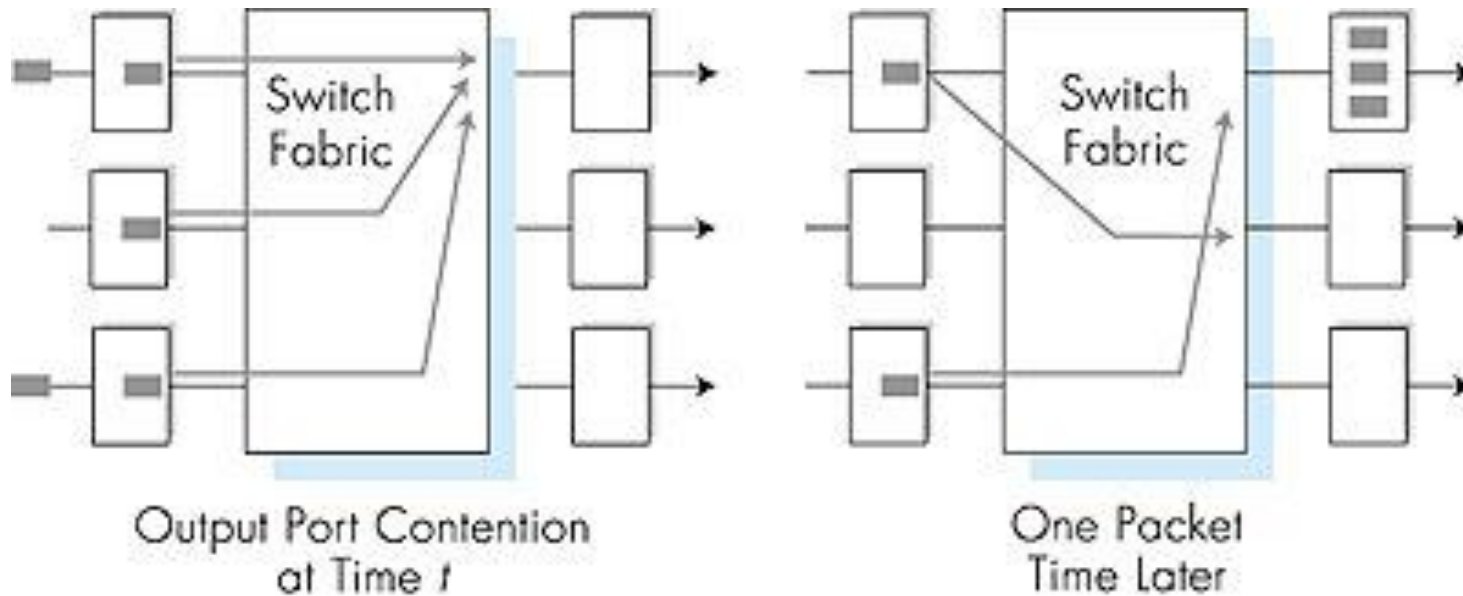
- Overcome bus bandwidth limitations
- Banyan networks, other interconnection nets initially developed to connect processors in multiprocessor
- Advanced design: fragmenting datagram into fixed length cells, switch cells through the fabric.
- Cisco 12000: switches Gbps through the interconnection network

Output Ports



- *Buffering* required when datagrams arrive from fabric faster than the transmission rate
- *Scheduling discipline* chooses among queued datagrams for transmission
- *queueing (delay) and loss due to output port buffer overflow!*

Output port queueing



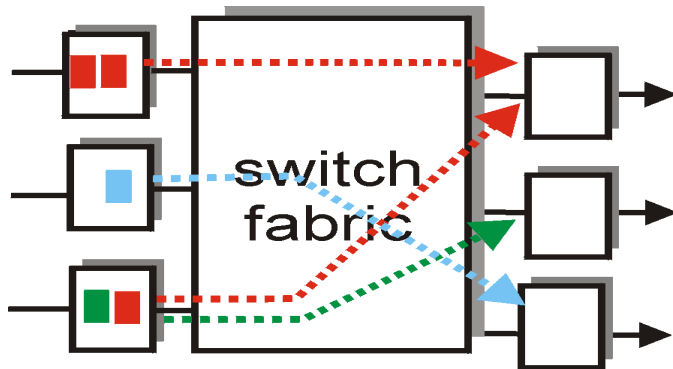
- buffering when arrival rate via switch exceeds output line speed
- *queueing (delay) and loss due to output port buffer overflow!*

How much buffering?

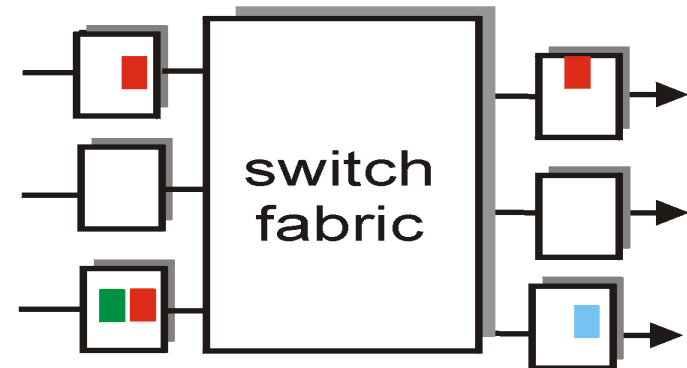
- RFC 3439 rule of thumb: average buffering equal to "typical" RTT (say 250 msec) times link capacity C
 - e.g., $C = 10$ Gps link: 2.5 Gbit buffer
- Recent recommendation: with N flows, buffering equal to $\frac{RTT \cdot C}{\sqrt{N}}$

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

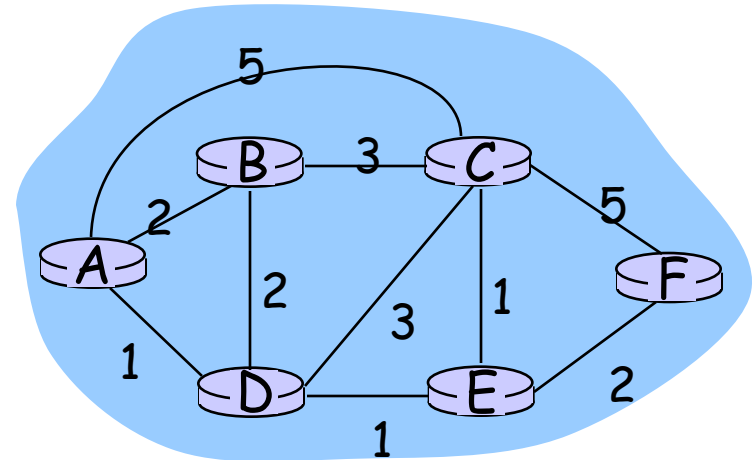
Routing

Routing protocol

Goal: determine “good” path (sequence of routers) thru network from source to dest.

Graph abstraction for routing algorithms:

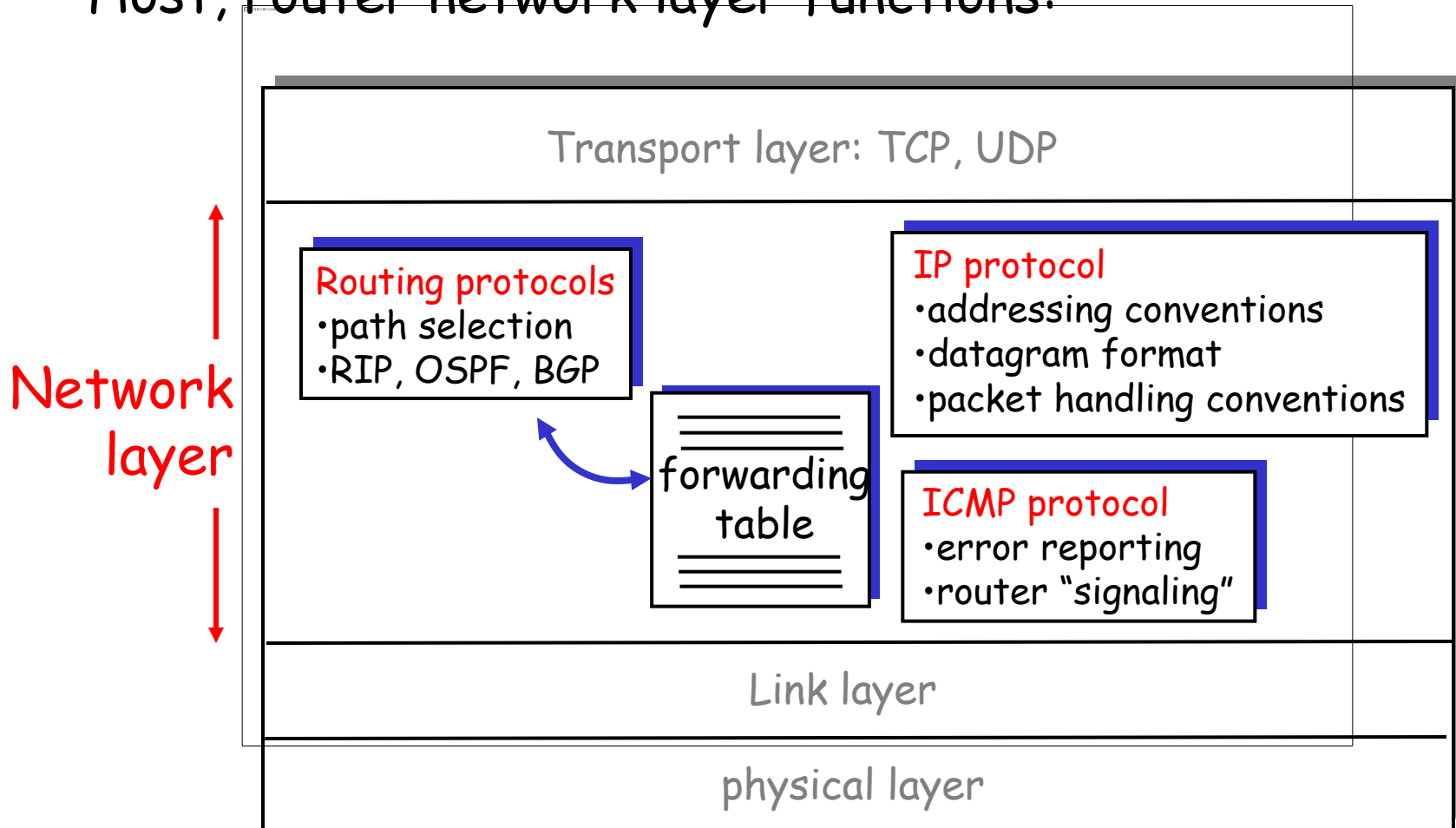
- graph **nodes** are **routers**
- graph **edges** are **physical links**
- link **cost**: delay, \$ cost, or congestion level



- “good” path:
 - typically means minimum cost path
 - other def’s possible

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

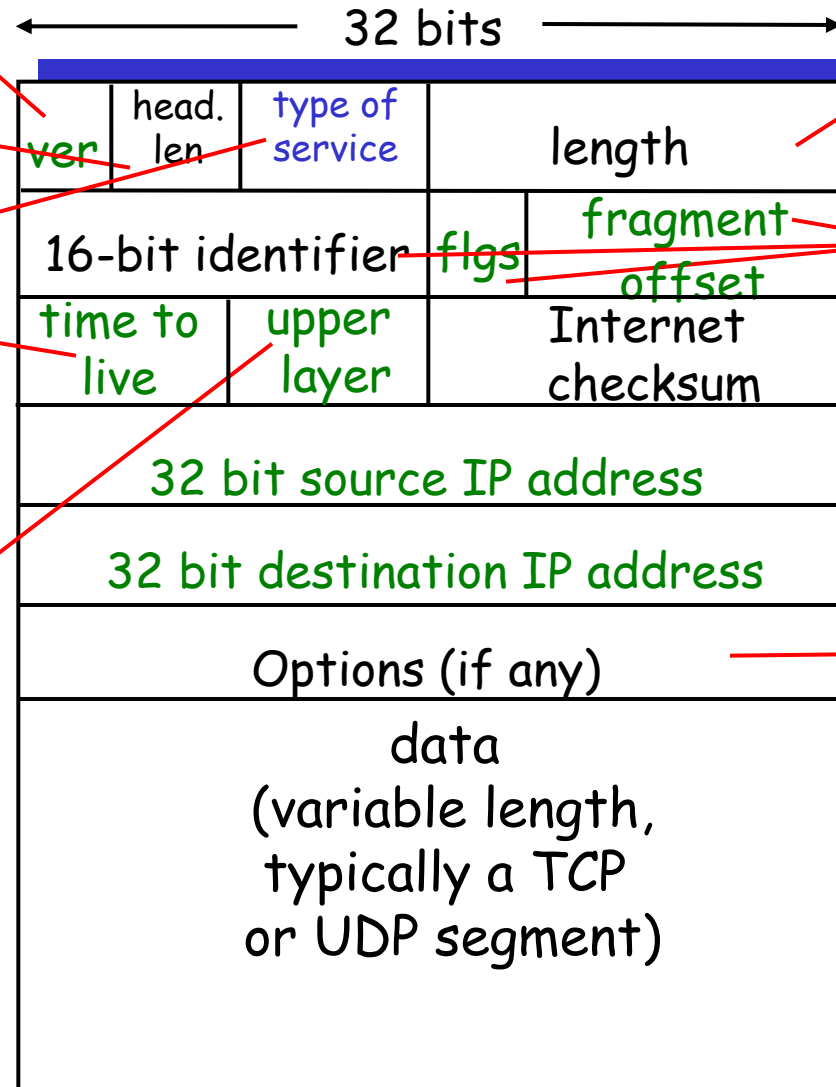
IP protocol version number

header length (bytes)

“type” of data

max number remaining hops (decremented at each router)

upper layer protocol to deliver payload to



total datagram length (bytes)

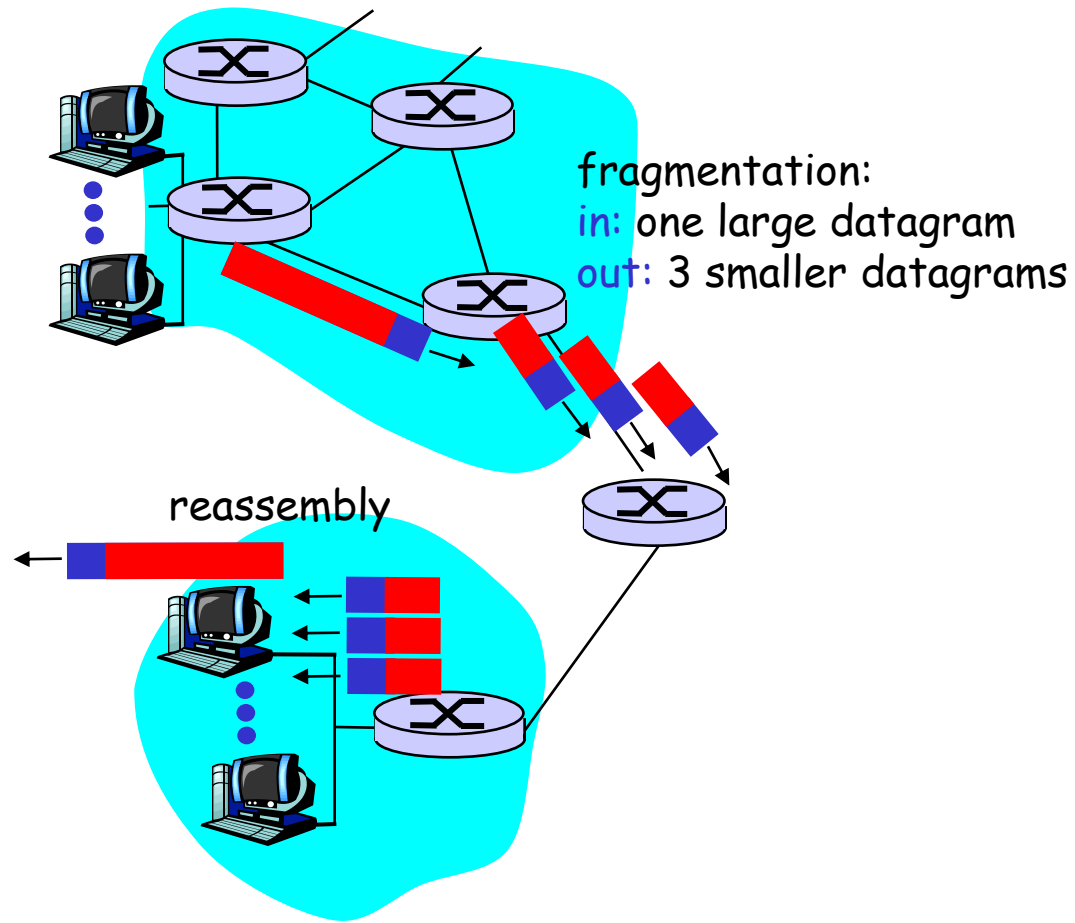
for fragmentation/reassembly

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

- how much overhead with TCP?
- 20 bytes of TCP
 - 20 bytes of IP
 - = 40 bytes + app layer overhead

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

	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	=1480	

	length	ID	fragflag	offset	
	=1040	=x	=0	=2960	

$$\begin{aligned} 4000 &= 20 + 3980 \\ &= (20 + 1480) + (20 + 1480) \\ &\quad + (20 + 1020) \end{aligned}$$

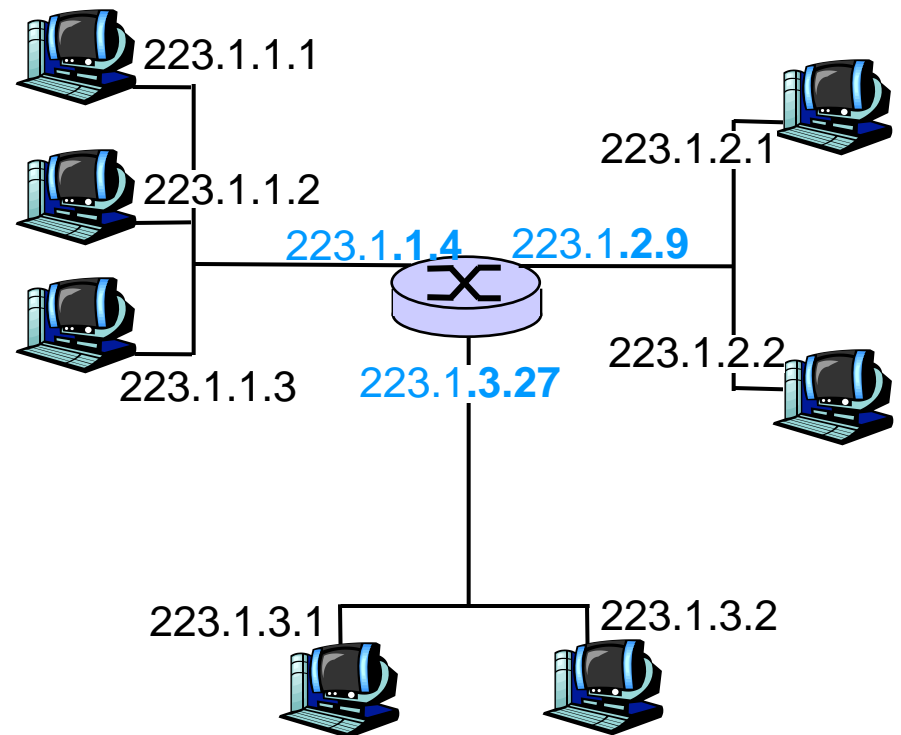
Chapter 4: Network Layer

- 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

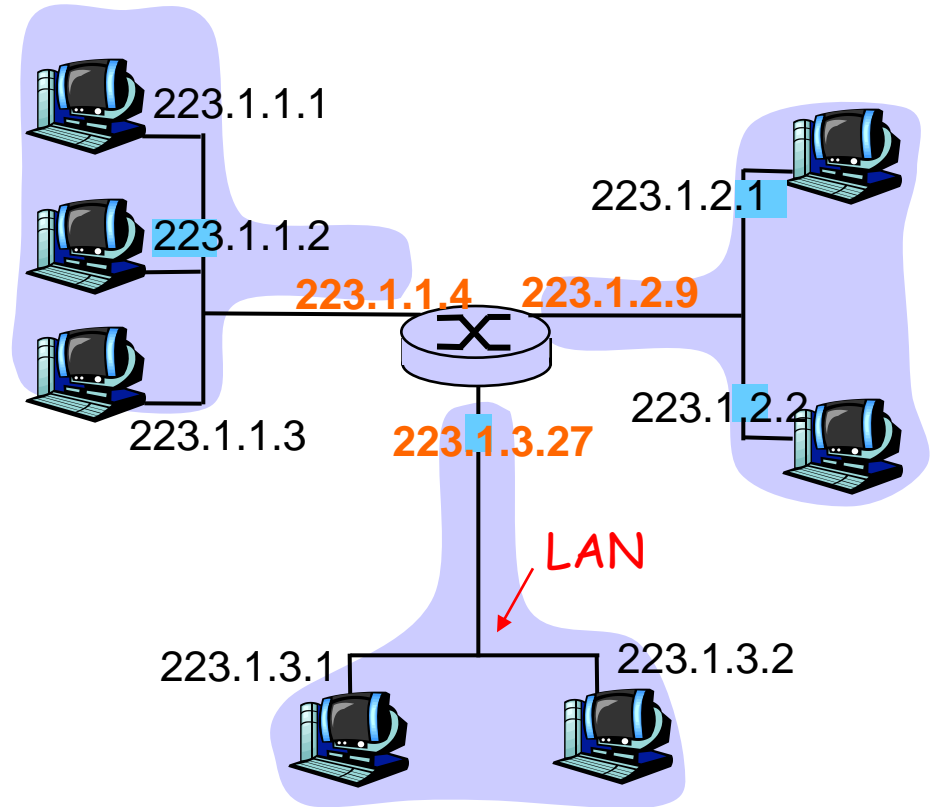
- **IP address:** 32-bit identifier for host, router *interface*
- **interface:** connection between host/router and physical link
 - routers typically have multiple interfaces
 - host may have multiple interfaces
 - One or more IP addresses may be associated with an interface

$$223.1.1.1 = \underbrace{11011111}_{223} \underbrace{00000001}_{1} \underbrace{00000001}_{1} \underbrace{00000001}_{1}$$



IP Addressing

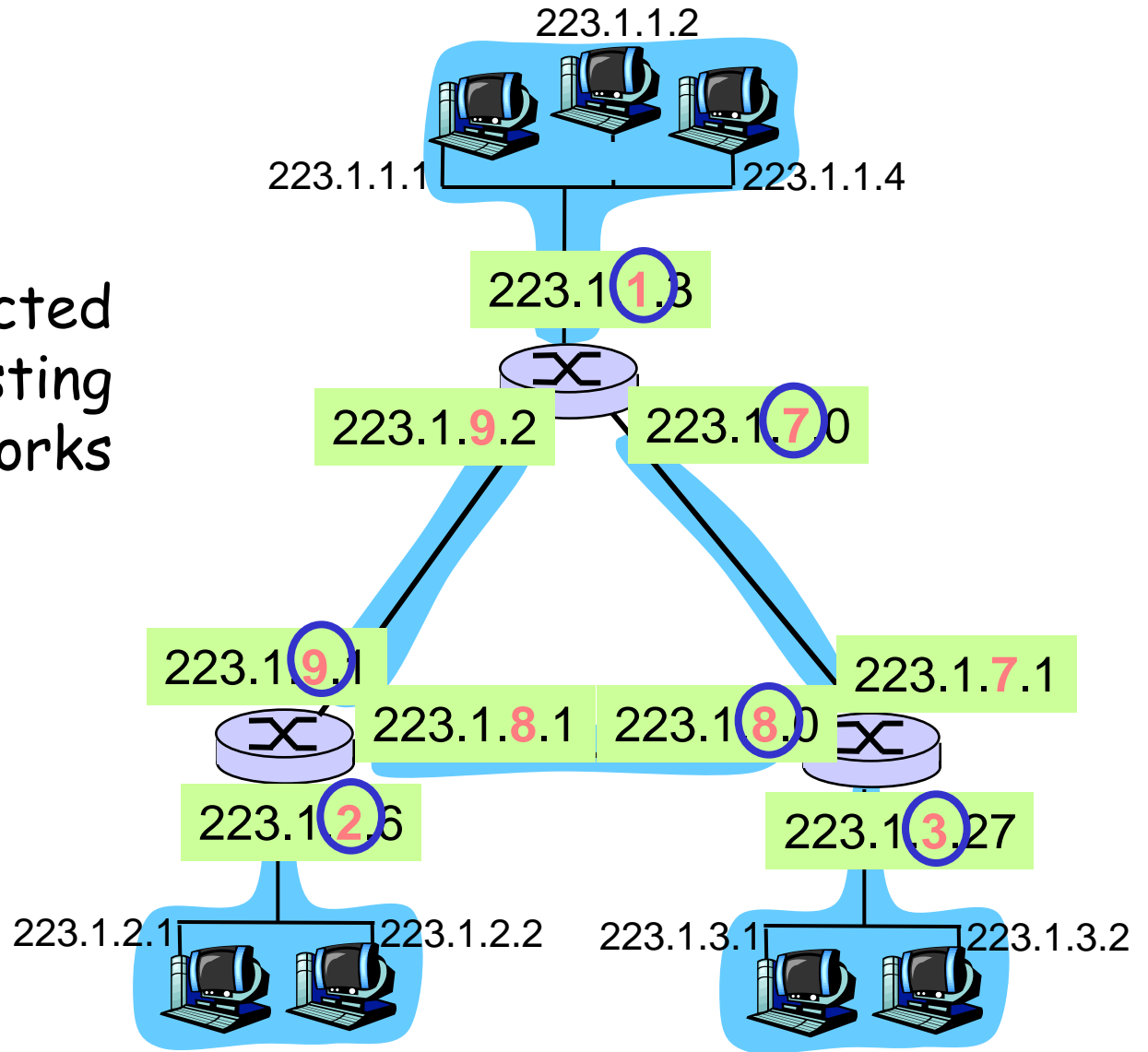
- IP address:
 - network part (high order bits)
 - host part (low order bits)
- *What's a subnet ?*
(from IP address perspective)
 - device interfaces with **same** network part of IP address
 - can physically reach each other without intervening router



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

Subnets

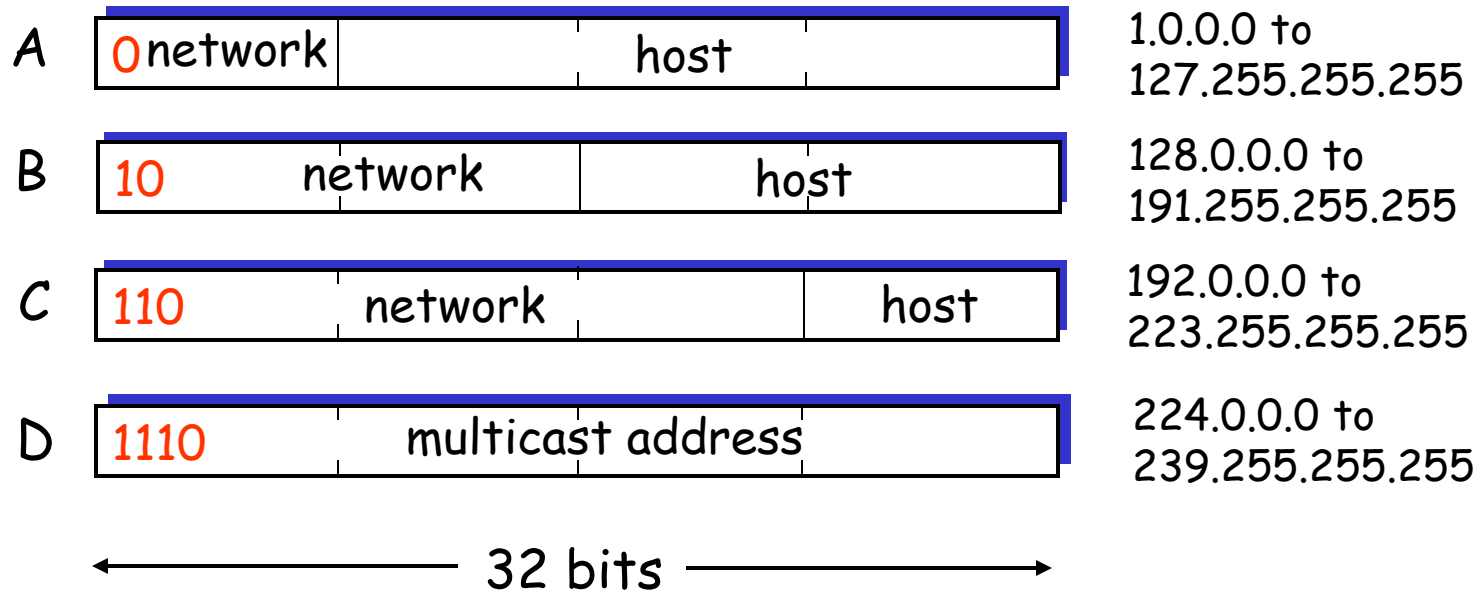
Interconnected system consisting of six networks



IP Addresses

“Class-full” addressing:

class



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 as server
 - "plug-and-play"(more shortly)

DHCP: Dynamic Host Configuration Protocol

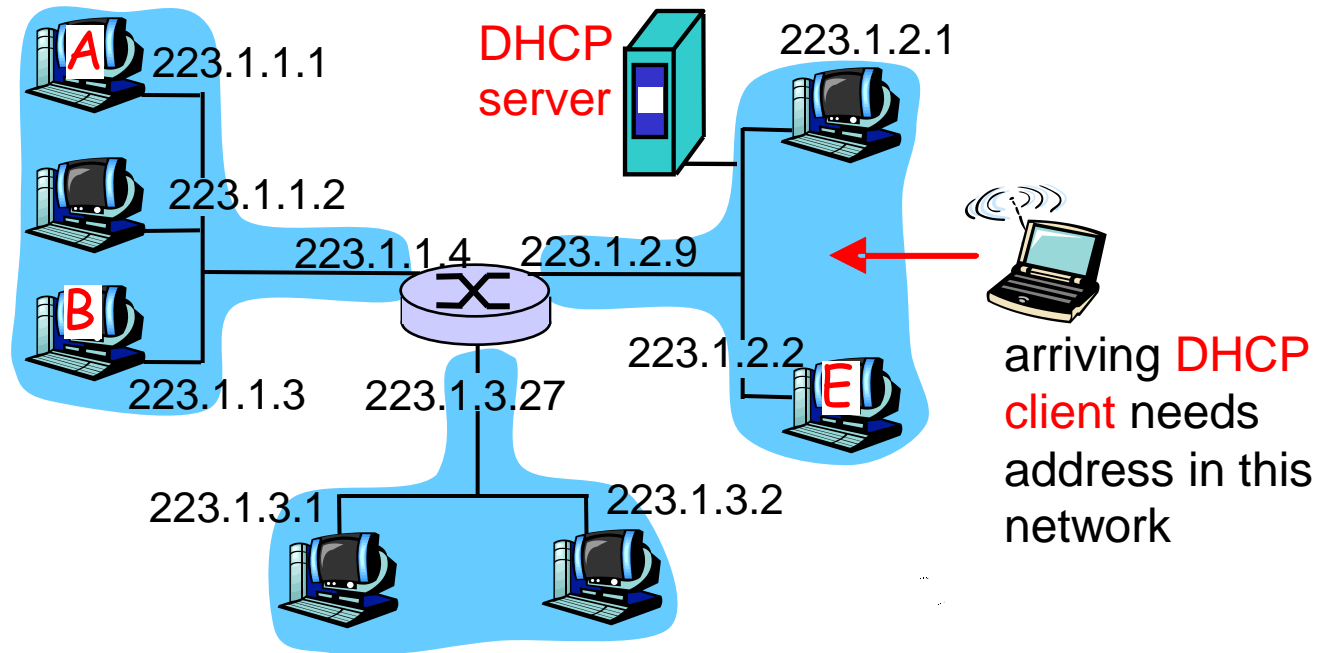
Goal: allow host to *dynamically* obtain its IP address from network *server when it joins network*

- Can *renew* its lease on address in use
- *Allows reuse of addresses* (only hold address while connected and "on")
- Support for *mobile users* who want to join network (more shortly)

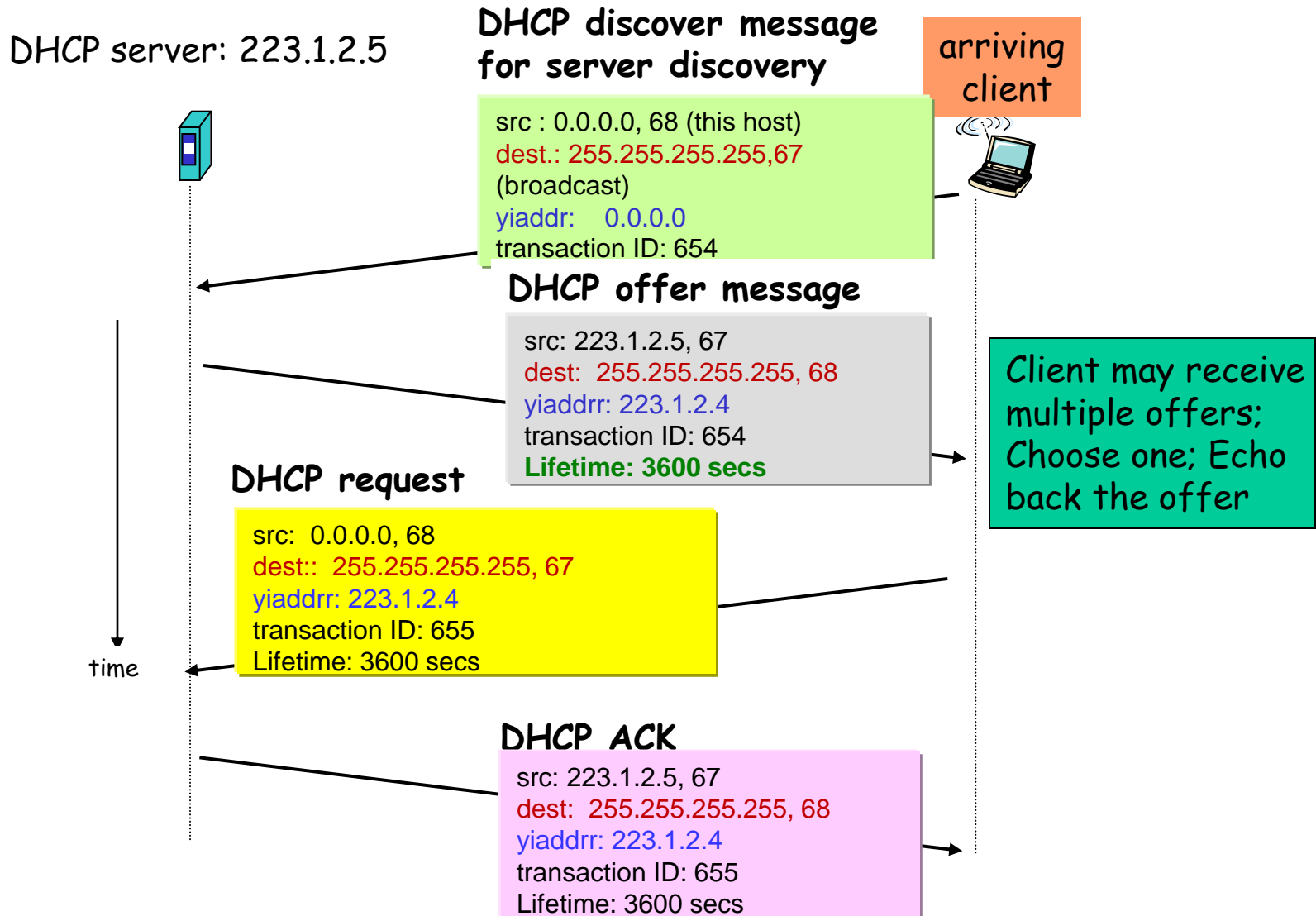
DHCP overview: (on top of UDP)

- host broadcasts "*DHCP discover*" msg
- DHCP server responds with "*DHCP offer*" msg
- host requests IP address: "*DHCP request*" msg
- DHCP server sends address: "*DHCP ack*" msg

DHCP client-server scenario



DHCP client-server scenario



IP addresses: how to get one?

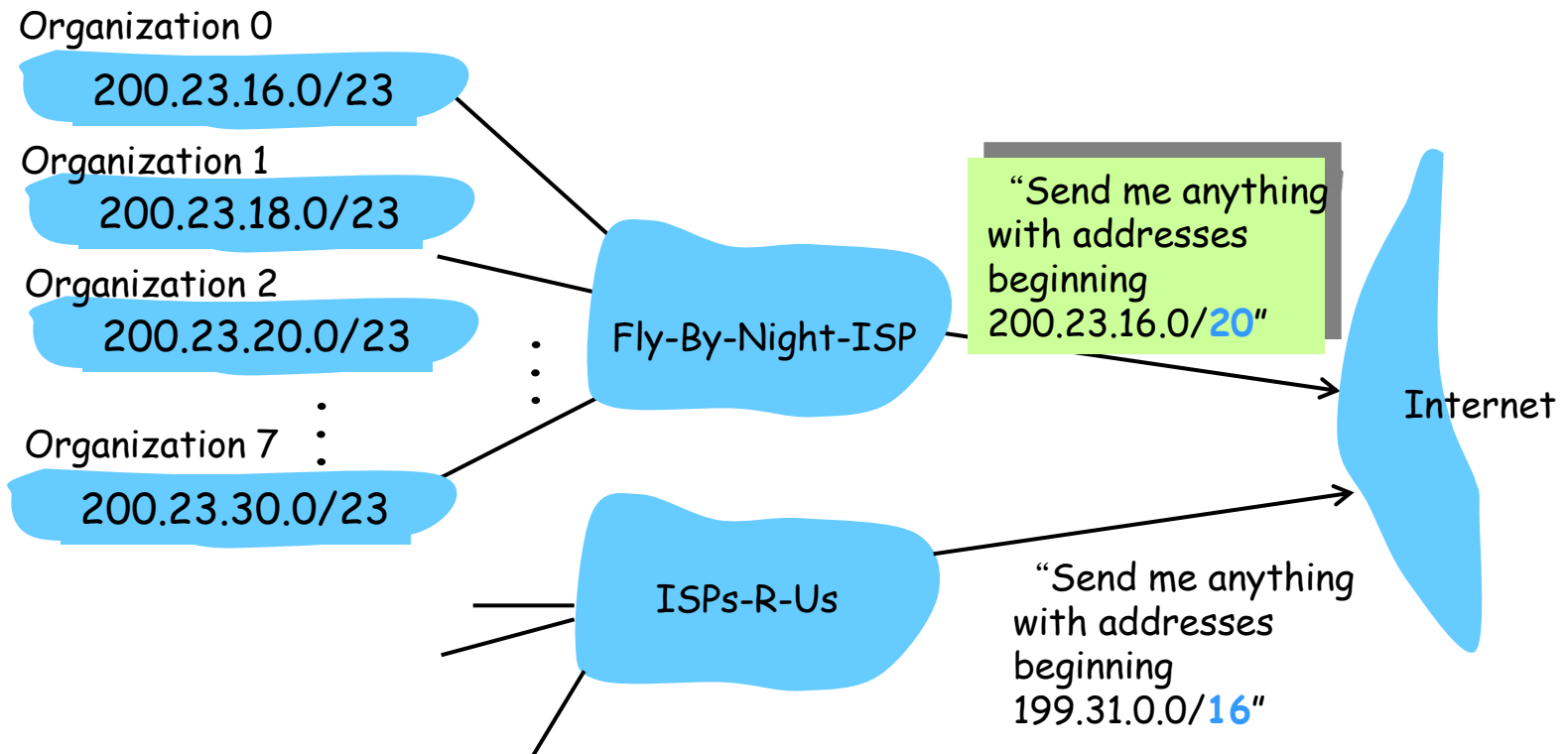
Q: How does *network* get network part of IP addr?

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>	0001 <u>0000</u>	00000000	200.23.16.0/23
Organization 1	<u>11001000</u>	<u>00010111</u>	0001 <u>0010</u>	00000000	200.23.18.0/23
Organization 2	<u>11001000</u>	<u>00010111</u>	0001 <u>0100</u>	00000000	200.23.20.0/23
...
Organization 7	<u>11001000</u>	<u>00010111</u>	0001 <u>1110</u>	00000000	200.23.30.0/23

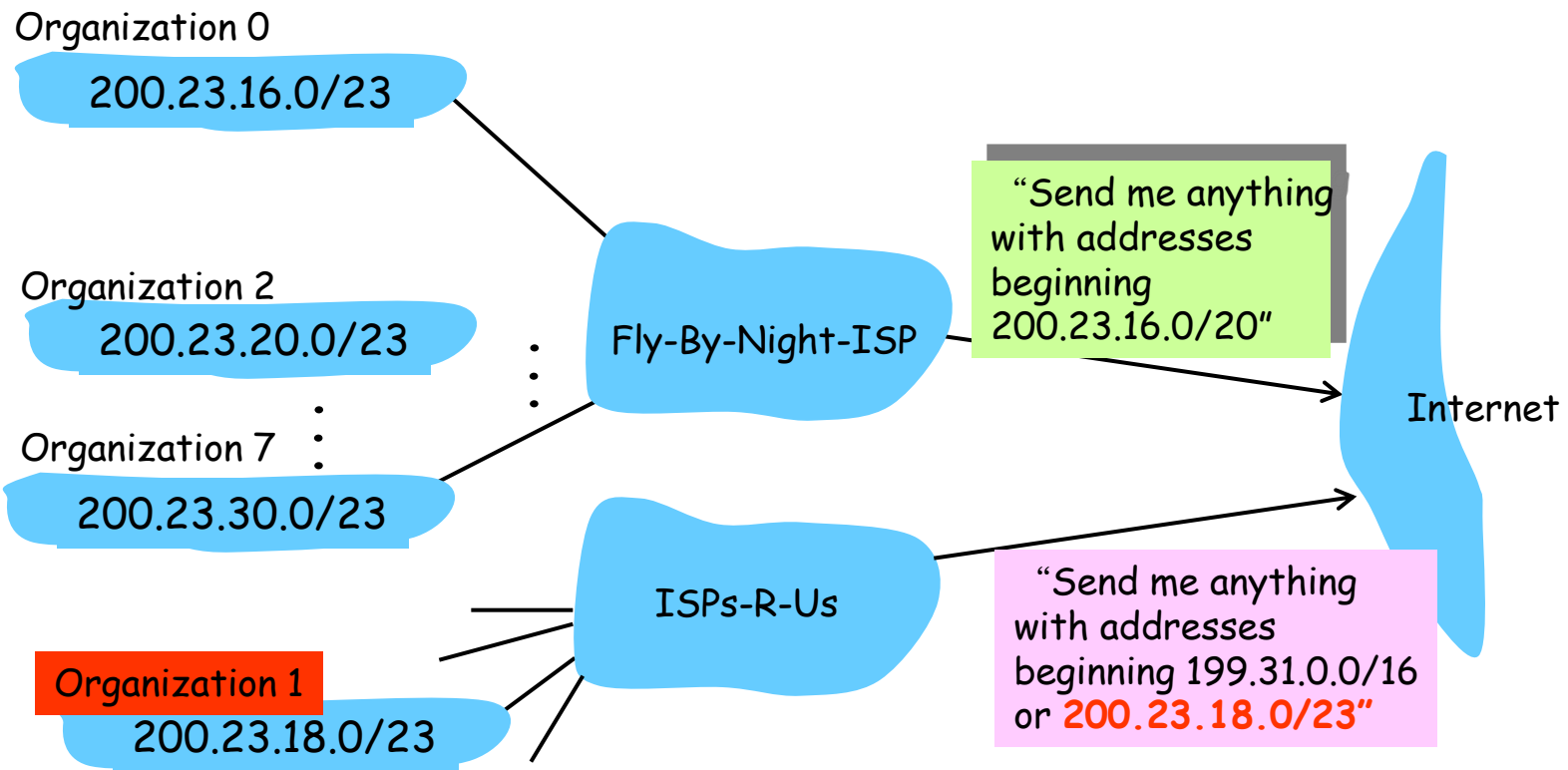
Hierarchical addressing: route aggregation

Hierarchical addressing allows efficient advertisement of routing information:



Hierarchical addressing: Longest Prefix Match

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



Routing Table - Longest Prefix Match

Routing table

200.23.16.0/20
199.31.0.0/16
200.23.18.0/23

✓

✓



IP addressing: the last word...

Q: How does an ISP get block of addresses?

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

- allocates addresses
 - manages DNS
 - assigns domain names, resolves disputes
-
- **APEC: APNIC**
 - **TWNIC**

Connecting a LAN to the Internet using Private IP Addresses - NAT (network address translation)

Motivation

- How to do private IP addressing?
- How Private IP networks connect to the Internet?

Connecting Private LANs to the Internet

- *An Internet account isn't expensive nowadays!*
- A **person** may have a whole (Ethernet) network of computers at home
 - you want to access the Internet from each of those machines
 - *one machine*, which has the modem, will be the "*middle-man*" for the other machines

Connecting Private LANs to the Internet (cont'd)

- A **school** has a couple of computers and wants to connect to the Internet at low cost.
- A **small office** wants to connect to the Internet
- **Dynamic IP address assignment**
 - An **ISP** usually buys a smaller **block of IP addresses**
 - When a subscriber calls in, he/she receives one of the IP addresses out of this block during the connection setup negotiation process
 - Subscribers *don't* know the IP address in advance

Connecting Private LANs to the Internet (cont'd)

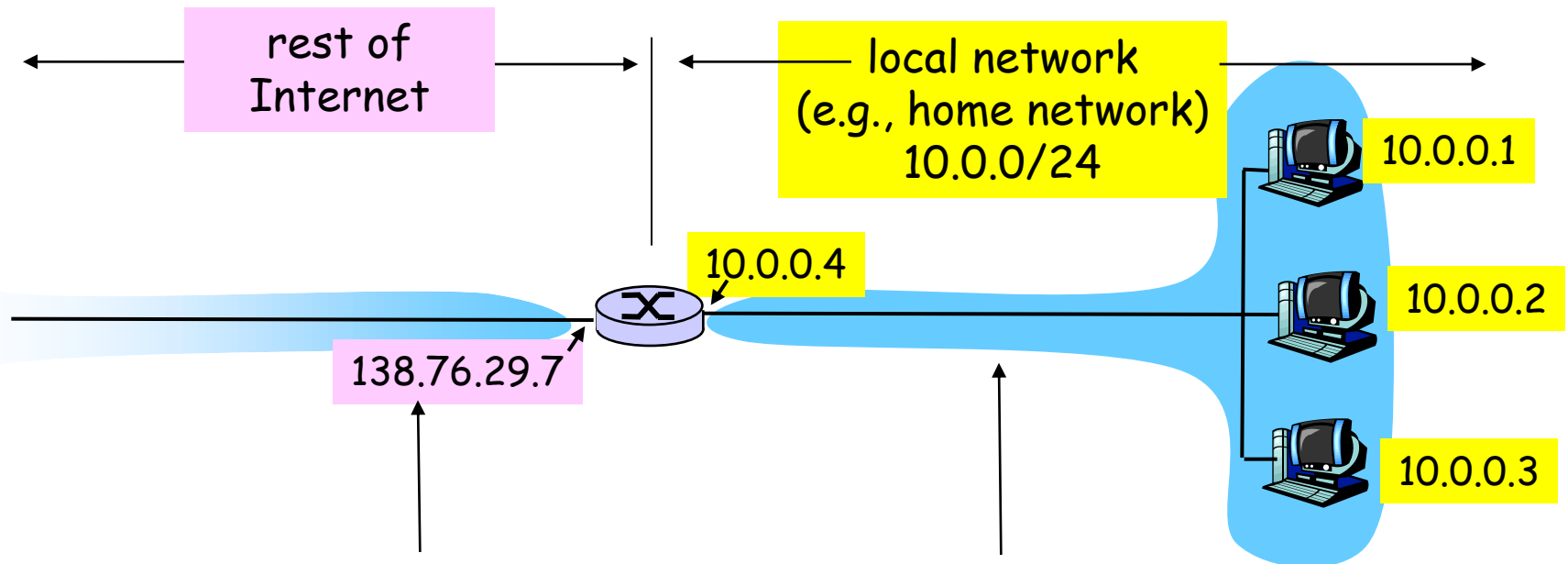
■ Static IP addresses

- some ISPs also offer fixed addresses, i.e. every time you dial-in you get the same IP address.
- e.g. if you want to run servers

Different Approaches to Connecting a LAN to the Internet

- Serial port sharing
- Routing
- Proxy servers
- IP Masquerading

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:
 - 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**).

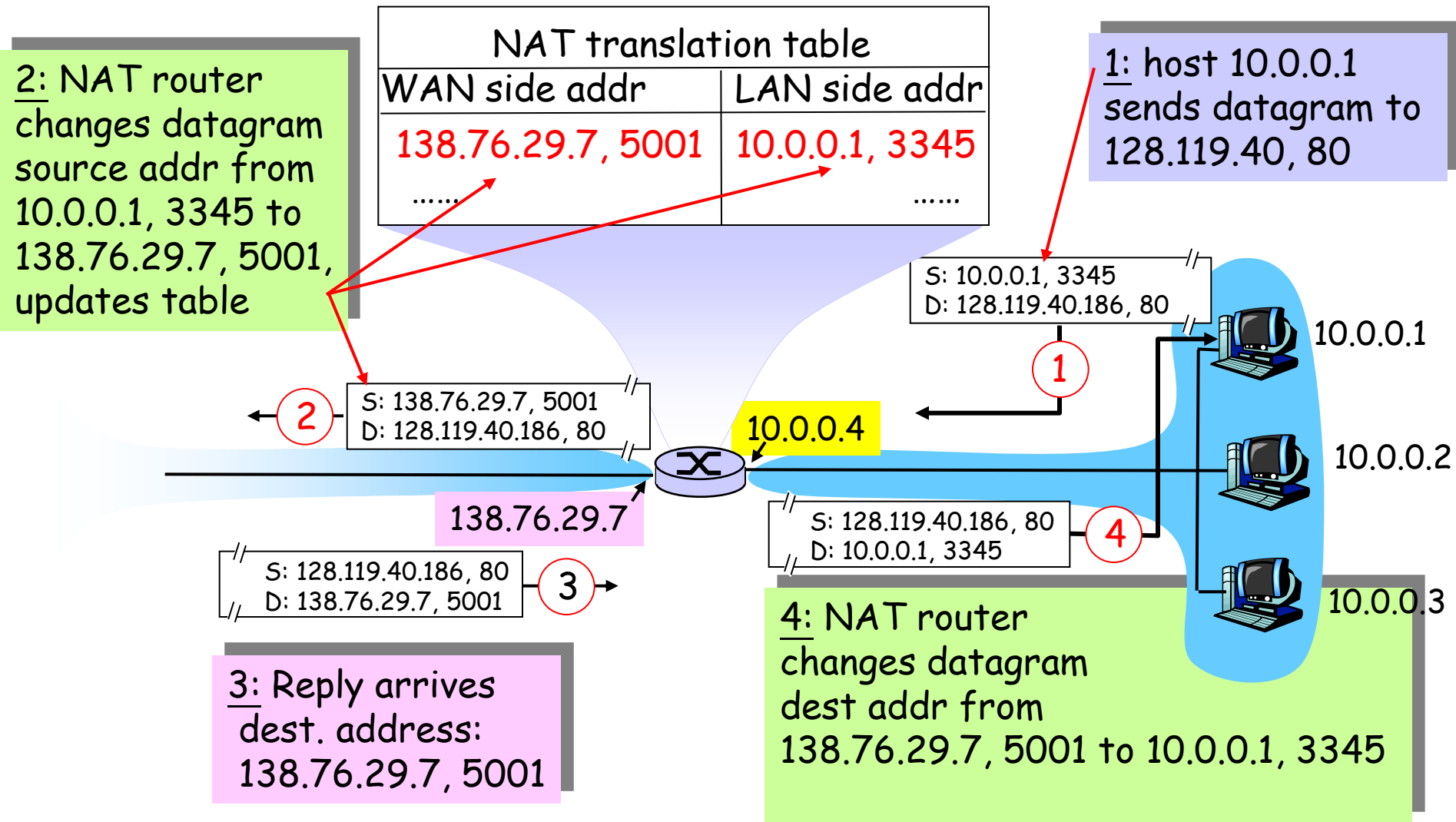
This public IP address could be dynamically allocated via DHCP.

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

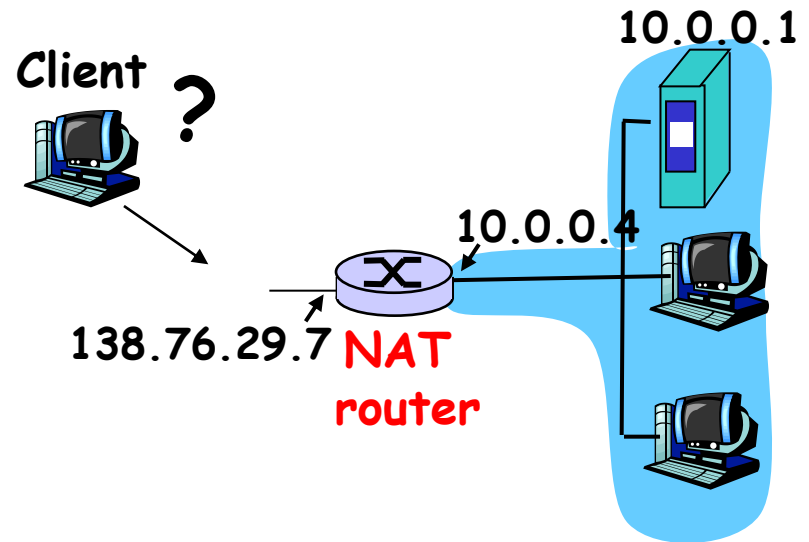


NAT: Network Address Translation

- 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, e.g., P2P applications
 - address shortage should instead be solved by IPv6

NAT traversal problem

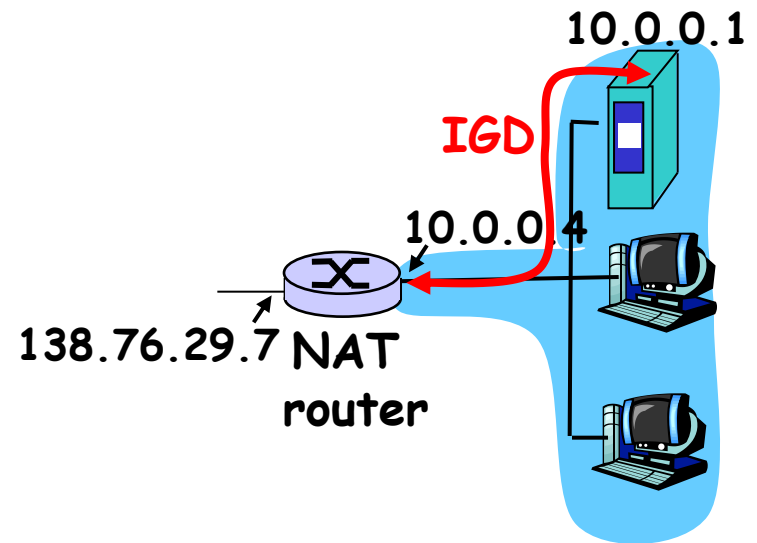
- client wants to connect to server with address 10.0.0.1
 - server address 10.0.0.1 local to LAN (client can't use it as destination addr)
 - only one externally visible NATted address: 138.76.29.7
- solution 1: statically configure NAT to forward incoming connection requests at given port to server
 - e.g., (138.76.29.7, port 2500) always forwarded to 10.0.0.1 port 2500



NAT traversal problem

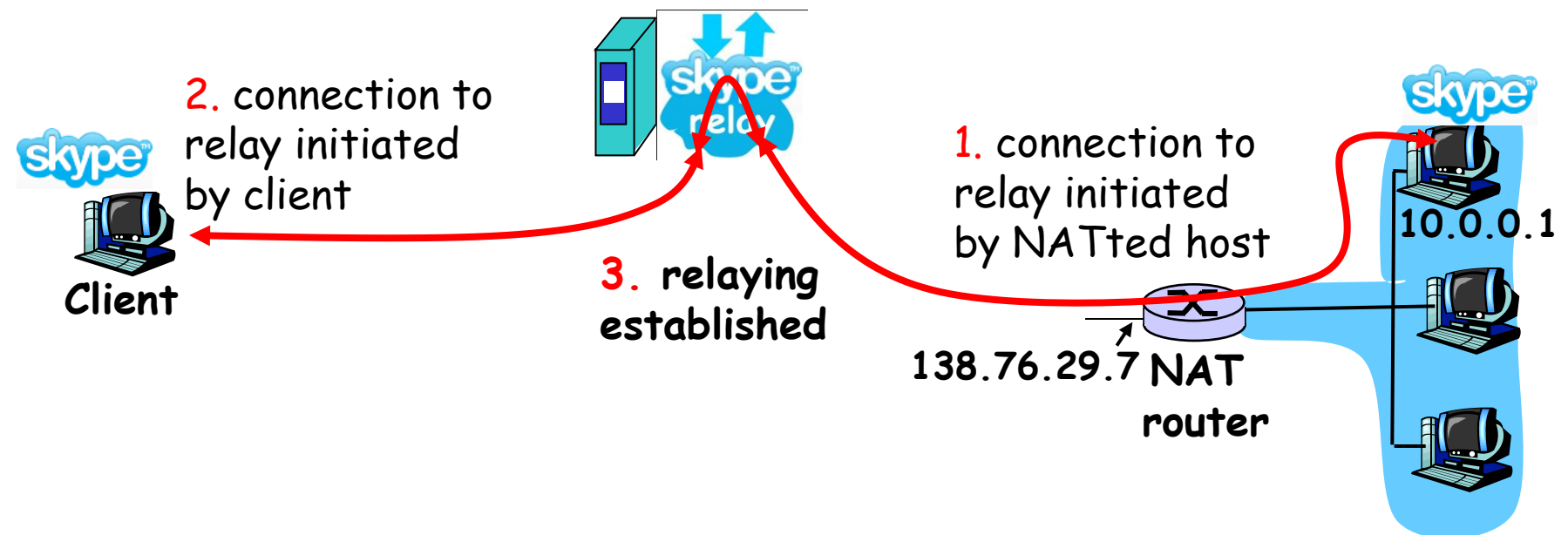
- solution 2: Universal Plug and Play (UPnP) Internet Gateway Device (IGD) Protocol. Allows NATted host to:
 - ❖ learn public IP address (138.76.29.7)
 - ❖ add/remove port mappings (with lease times)

i.e., automate static NAT port map configuration



NAT traversal problem

- solution 3: relaying (used in Skype)
 - NATed client establishes connection to relay
 - External client connects to relay
 - relay bridges packets between to connections



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

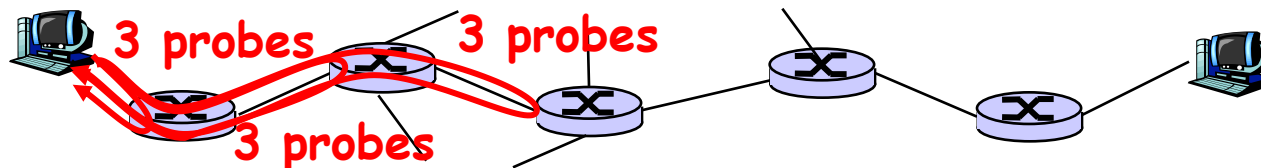
ICMP: Internet Control Message Protocol

- used by hosts, routers, gateways 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 n^{th} datagram arrives to n^{th} 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 "protocol unreachable" packet (type 3, code 3)**
 - When source gets this ICMP, stops.



“Real” Internet delays and routes


traceroute: gaia.cs.umass.edu to www.eurecom.fr

Three delay measurements from
gaia.cs.umass.edu to cs-
gw.cs.umass.edu



1 cs-gw (128.119.240.254) 1 ms 1 ms 2 ms
2 border1-rt-fa5-1-0.gw.umass.edu (128.119.3.145) 1 ms 1 ms 2 ms
3 cht-vbns.gw.umass.edu (128.119.3.130) 6 ms 5 ms 5 ms
4 jn1-at1-0-0-19.wor.vbns.net (204.147.132.129) 16 ms 11 ms 13 ms
5 jn1-so7-0-0-0.wae.vbns.net (204.147.136.136) 21 ms 18 ms 18 ms
6 abilene-vbns.abilene.ucaid.edu (198.32.11.9) 22 ms 18 ms 22 ms
7 nycm-wash.abilene.ucaid.edu (198.32.8.46) 22 ms 22 ms 22 ms
8 62.40.103.253 (62.40.103.253) 104 ms 109 ms 106 ms
9 de2-1.de1.de.geant.net (62.40.96.129) 109 ms 102 ms 104 ms
10 de.fr1.fr.geant.net (62.40.96.50) 113 ms 121 ms 114 ms
11 renater-gw.fr1.fr.geant.net (62.40.103.54) 112 ms 114 ms 112 ms
12 nio-n2.cssi.renater.fr (193.51.206.13) 111 ms 114 ms 116 ms
13 nice.cssi.renater.fr (195.220.98.102) 123 ms 125 ms 124 ms
14 r3t2-nice.cssi.renater.fr (195.220.98.110) 126 ms 126 ms 124 ms
15 eurecom-valbonne.r3t2.ft.net (193.48.50.54) 135 ms 128 ms 133 ms
16 194.214.211.25 (194.214.211.25) 126 ms 128 ms 126 ms
17 * * *
18 * * *
19 fantasia.eurecom.fr (193.55.113.142) 132 ms 128 ms 136 ms

trans-oceanic link



* means no response (probe lost, router not replying)

Chapter 4: Network Layer

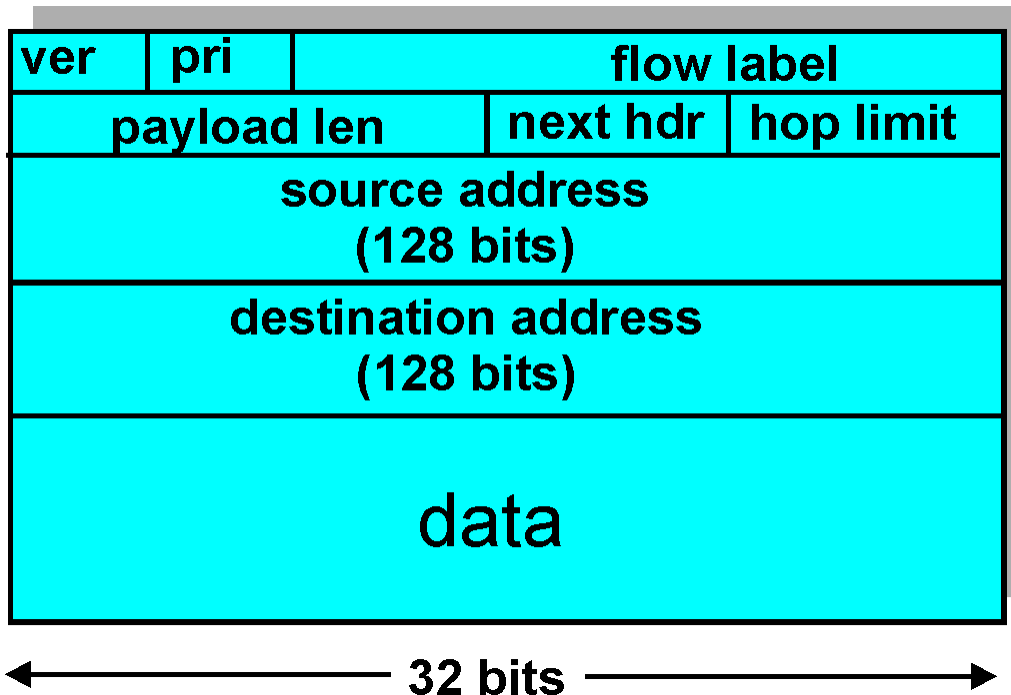
- 4.1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
- 4.4 IP: Internet Protocol
 - Datagram format
 - IPv4 addressing
 - ICMP
 - IPv6
- 4.5 Routing algorithms
 - Link state
 - Distance Vector
 - Hierarchical routing
- 4.6 Routing in the Internet
 - RIP
 - OSPF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



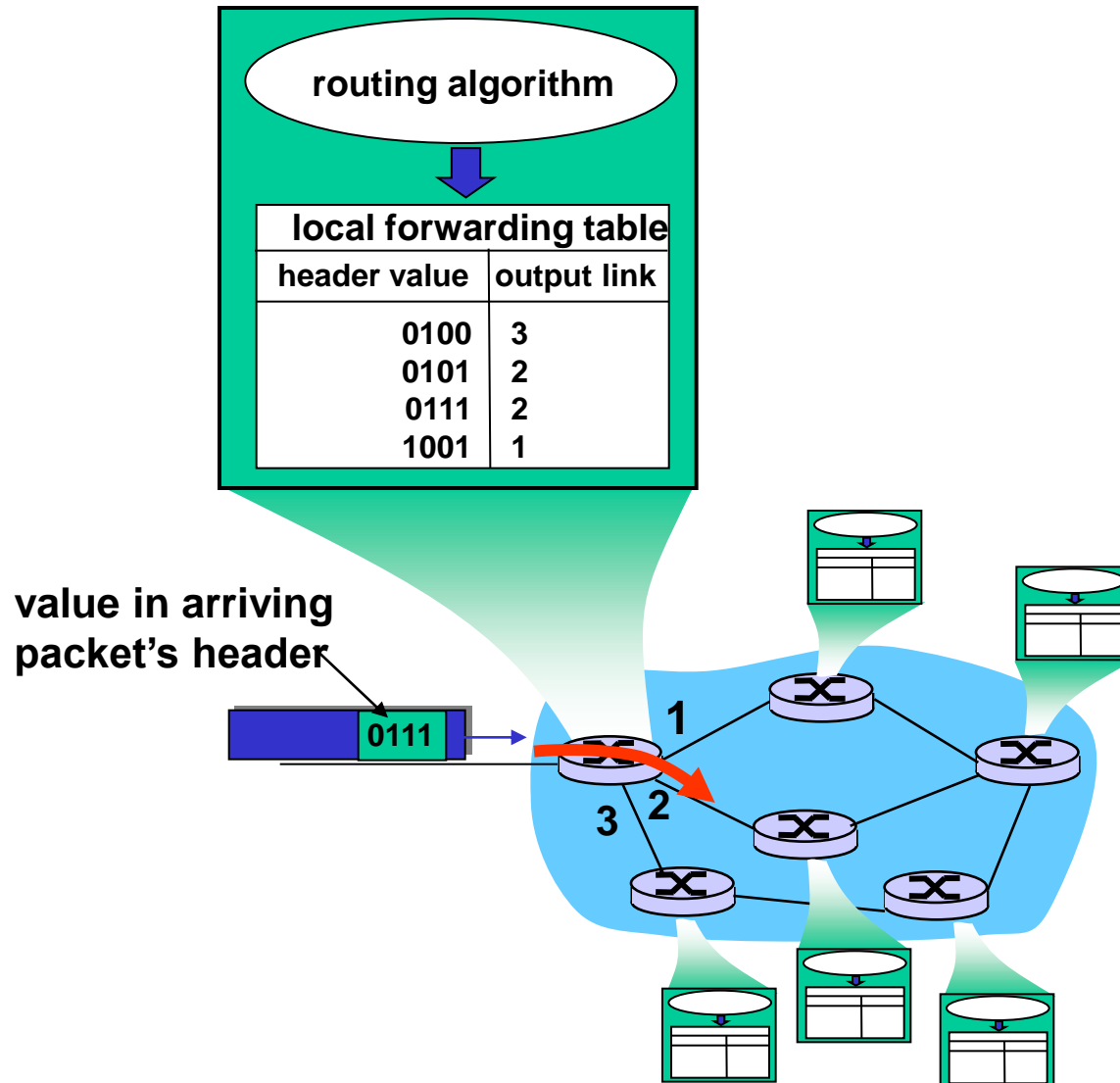
Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
- 4.4 IP: Internet Protocol
 - IPv4 addressing
 - Moving a datagram from source to destination
 - Datagram format
 - IP fragmentation
 - ICMP: Internet Control Message Protocol
 - DHCP: Dynamic Host Configuration Protocol
 - NAT: Network Address Translation
- 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

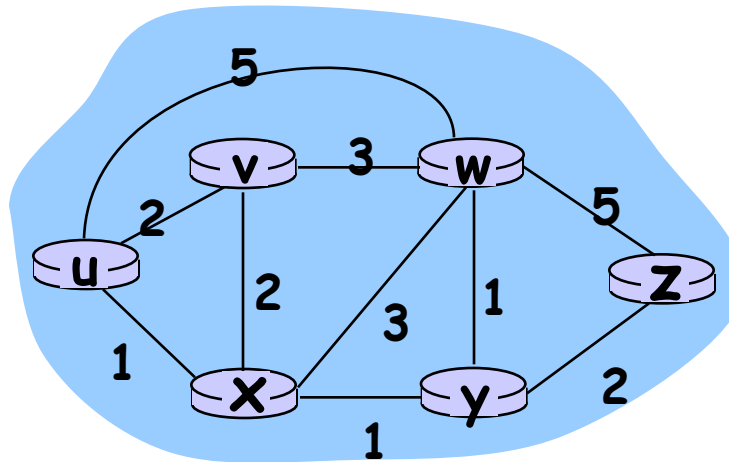
Chapter 4: Network Layer

- 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

Interplay between routing, forwarding



Graph abstraction



Graph: $G = (N,E)$

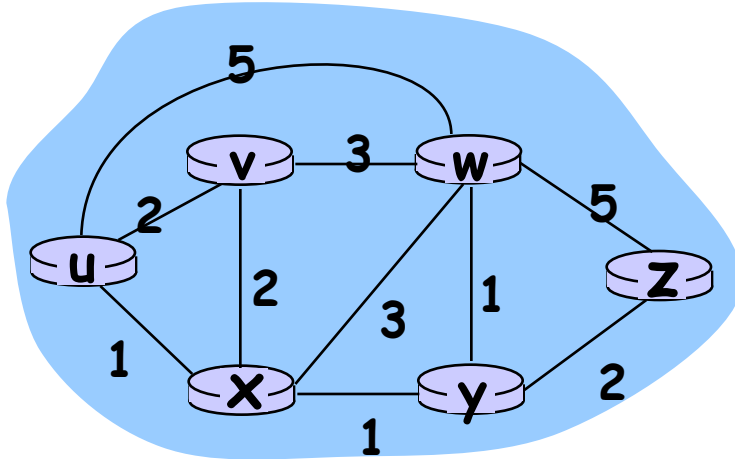
$N =$ set of routers = $\{ u, v, w, x, y, z \}$

$E =$ set of links = $\{ (u,v), (u,x), (v,x), (v,w), (x,w), (x,y), (w,y), (w,z), (y,z) \}$

Remark: Graph abstraction is useful in other network contexts

Example: P2P, where N is set of peers and E is set of TCP connections

Graph abstraction: costs



• $c(x,x')$ = cost of link (x,x')

- e.g., $c(w,z) = 5$

• cost could always be 1, or inversely related to bandwidth, or inversely related to congestion

Cost of path $(x_1, x_2, x_3, \dots, x_p) = c(x_1, x_2) + c(x_2, x_3) + \dots + c(x_{p-1}, x_p)$

Question: What's the least-cost path between u and z ?

Routing algorithm: algorithm that finds least-cost path

Routing Algorithm classification

Global or decentralized information?

Global:

- all routers have **complete topology, link cost info**
- "link state" algorithms

Decentralized:

- router knows **physically-connected neighbors**, link costs to neighbors
- iterative process of computation, exchange of info with neighbors
- "distance vector" algorithms

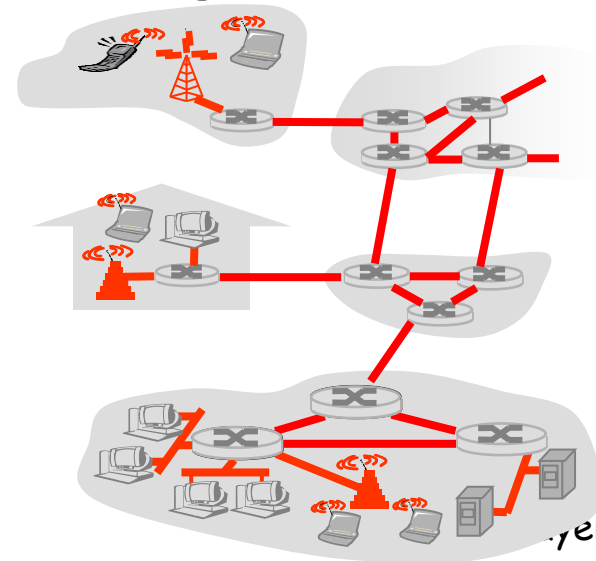
Static or dynamic?

Static:

- routes change slowly over time

Dynamic:

- routes change more quickly
 - periodic update
 - in response to link cost changes



Flooding

- An extreme form of **isolated routing**
- Algorithm
 - Every received incoming packet is sent to out on every outgoing link except the one it arrived on.
- Issue
 - How to terminate packet flooding?

Flooding (cont'd)

■ Method#1

- Use a hop count which is attached to each packet and is decremented at each hop visited.
- A packet is discarded when the counter reaches zero.
- The counter is often initialized to the worst-case length of the path from a source to destination, i.e. the longest path length of the network.

Flooding (cont'd)

■ Method#2

- Each packet is assigned a sequence number by the source node.
- Each switching node maintains a table about which sequence number originating from which source node have already been received.
- The table is consulted for each received packet.
- A duplicate packet is discarded, otherwise the packet's information is entered the table and is copied and sent over all the outgoing link except the arriving one.

Flooding (cont'd)

■ Drawbacks

- Generate enormous amount of duplicate packets and use a large amount of bandwidth.
- Not practical in most applications. (but now is commonly used in some mobile wireless applications - simple!)

Flooding (cont'd)

- Useful for applications that requires
 - **robust** delivery, e.g., in military
 - **concurrent** data delivery to all nodes, e.g., in network database update.
- Often used as a performance metric against other routing algorithms
- Variations, e.g.,
 - Selective flooding

A Link-State Routing Algorithm

Dijkstra's algorithm

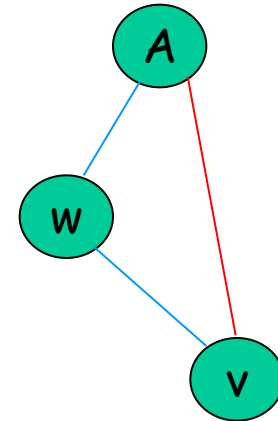
- net topology, link costs known to all nodes
 - accomplished via "link state broadcast"
 - all nodes have same info
- computes least cost paths from one node ("source") to all other nodes
 - gives routing table for that node
- iterative: after k iterations, know least cost path to k dest.'s

Notation:

- $c(i,j)$: link cost from node i to j . cost infinite if not direct neighbors
- $D(v)$: current value of cost of path from source to dest. v
- $p(v)$: predecessor node along path from source to v , that is next to v
- N : set of nodes whose least cost path definitively known

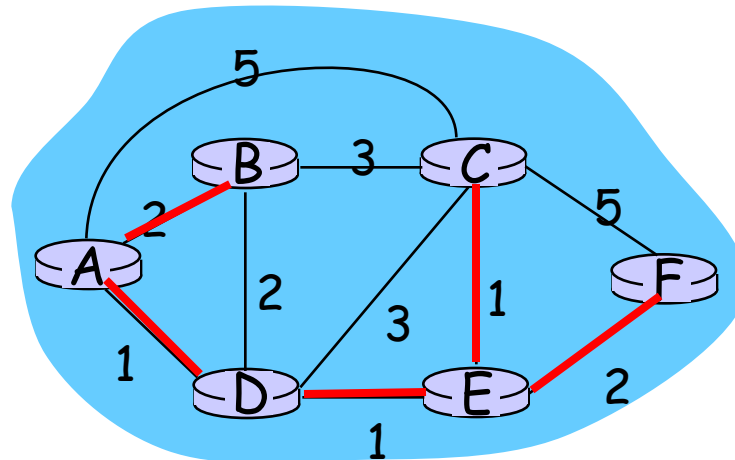
Dijkstra's Algorithm (shortest paths for all src-dst pairs)

- 1 **Initialization:**
- 2 $N = \{A\}$
- 3 for all nodes v
- 4 if v adjacent to A
- 5 then $D(v) = c(A,v)$ // link cost from A to v .
- 6 else $D(v) = \text{infinity}$
- 7
- 8 **Loop**
- 9 find w not in N such that $D(w)$ is a minimum
- 10 add w to N
- 11 update $D(v)$ for all v adjacent to w and not in N :
- 12 $D(v) = \min(D(v), D(w) + c(w,v))$
- 13 /* new cost to v is either old cost to v or known
- 14 shortest path cost to w plus cost from w to v */
- 15 **until all nodes in N**



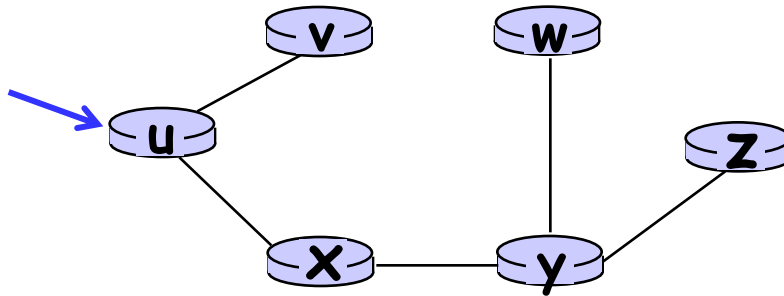
Dijkstra's algorithm: example

Step	start N	D(B),p(B)	D(C),p(C)	D(D),p(D)	D(E),p(E)	D(F),p(F)
→ 0	A	2,A	5,A	1,A	infinity	infinity
→ 1	AD	2,A	4,D		2,D	infinity
→ 2	ADE	2,A	3,E			4,E
→ 3	ADEB		3,E			4,E
→ 4	ADEBC					4,E
5	ADEBCF					



Dijkstra's algorithm: example (2)

Resulting shortest-path tree from u:



Resulting forwarding table in u:

destination	link
v	(u,v)
x	(u,x)
y	(u,x)
w	(u,x)
z	(u,x)

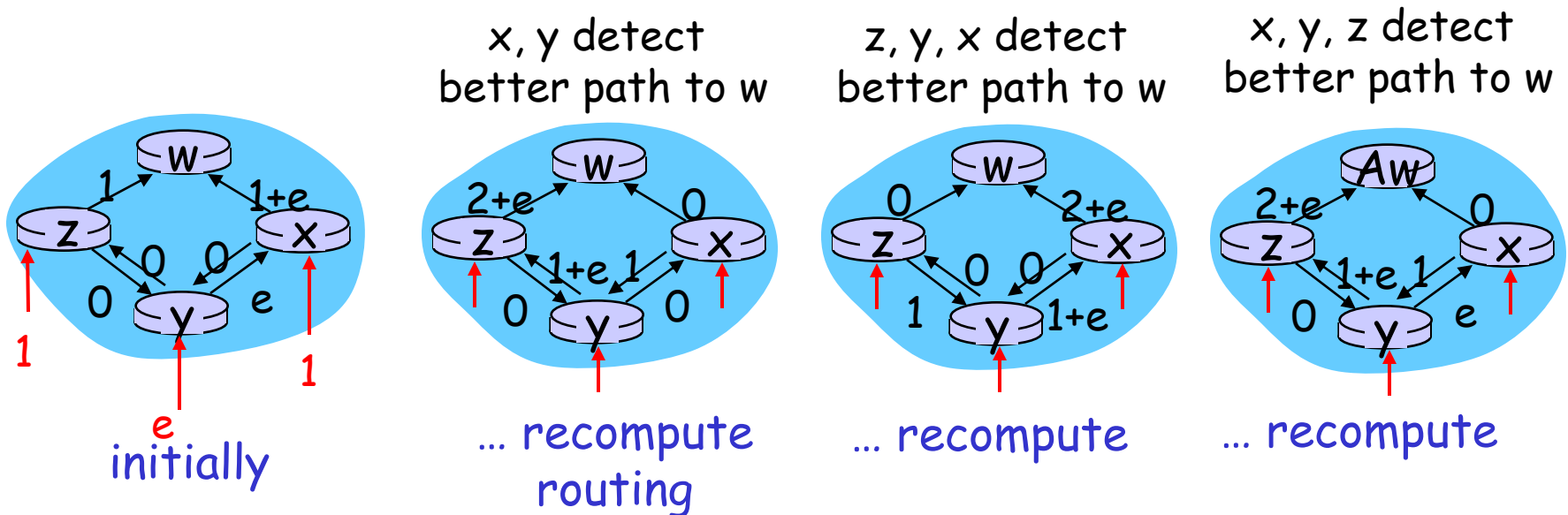
Dijkstra's algorithm, discussion

Algorithm complexity: n nodes

- each iteration: need to check all nodes, w, not in N
- $n*(n+1)/2$ comparisons: $O(n^2)$
- more efficient implementations possible: $O(n \log n)$ (using heap)

Oscillations possible:

- e.g., link cost = amount of carried traffic



Distance Vector

Distance Vector Algorithm

Bellman-Ford Equation (dynamic programming)

Define

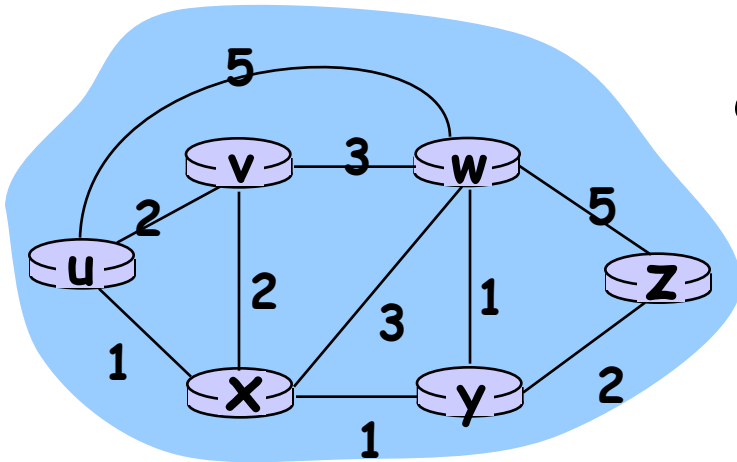
$d_x(y) :=$ cost of least-cost path from x to y

Then

$$d_x(y) = \min_v \{c(x,v) + d_v(y)\}$$

where min is taken over all neighbors v of x

Bellman-Ford example



Clearly, $d_v(z) = 5$, $d_x(z) = 3$, $d_w(z) = 3$

B-F equation says:

$$\begin{aligned}d_u(z) &= \min \{ c(u,v) + d_v(z), \\ &\quad c(u,x) + d_x(z), \\ &\quad c(u,w) + d_w(z) \} \\ &= \min \{ 2 + 5, \\ &\quad 1 + 3, \\ &\quad 5 + 3 \} = 4\end{aligned}$$

Node that achieves minimum is next hop in shortest path → forwarding table

Distance Vector Algorithm

- $D_x(y)$ = estimate of least cost from x to y
- Node x knows cost to each neighbor v :
 $c(x,v)$
- Node x maintains distance vector $D_x = [D_x(y): y \in N]$
- Node x also maintains its neighbors' distance vectors
 - For each neighbor v , x maintains $D_v = [D_v(y): y \in N]$

Distance vector algorithm (4)

Basic idea:

- From time-to-time, each node sends its own distance vector estimate to neighbors
- Asynchronous
- When a node x receives new DV estimate from neighbor, it updates its own DV using B-F equation:

$$D_x(y) \leftarrow \min_v \{c(x, v) + D_v(y)\} \quad \text{for each node } y \in N$$

- Under minor, natural conditions, the estimate $D_x(y)$ converge to the actual least cost $d_x(y)$

Distance Vector Algorithm (5)

Iterative, asynchronous:

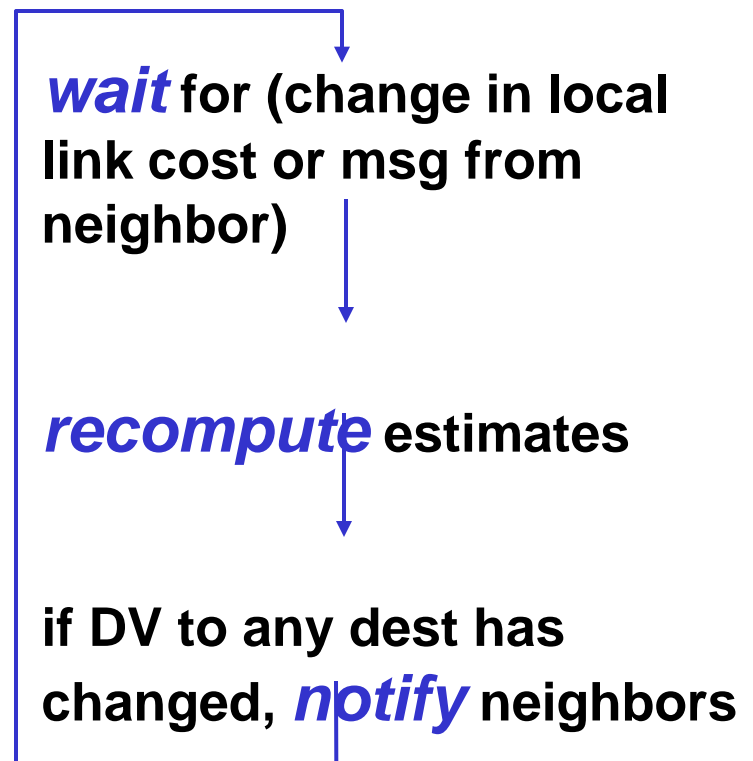
each local iteration caused by:

- local link cost change
- DV update message from neighbor

Distributed:

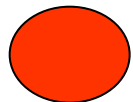
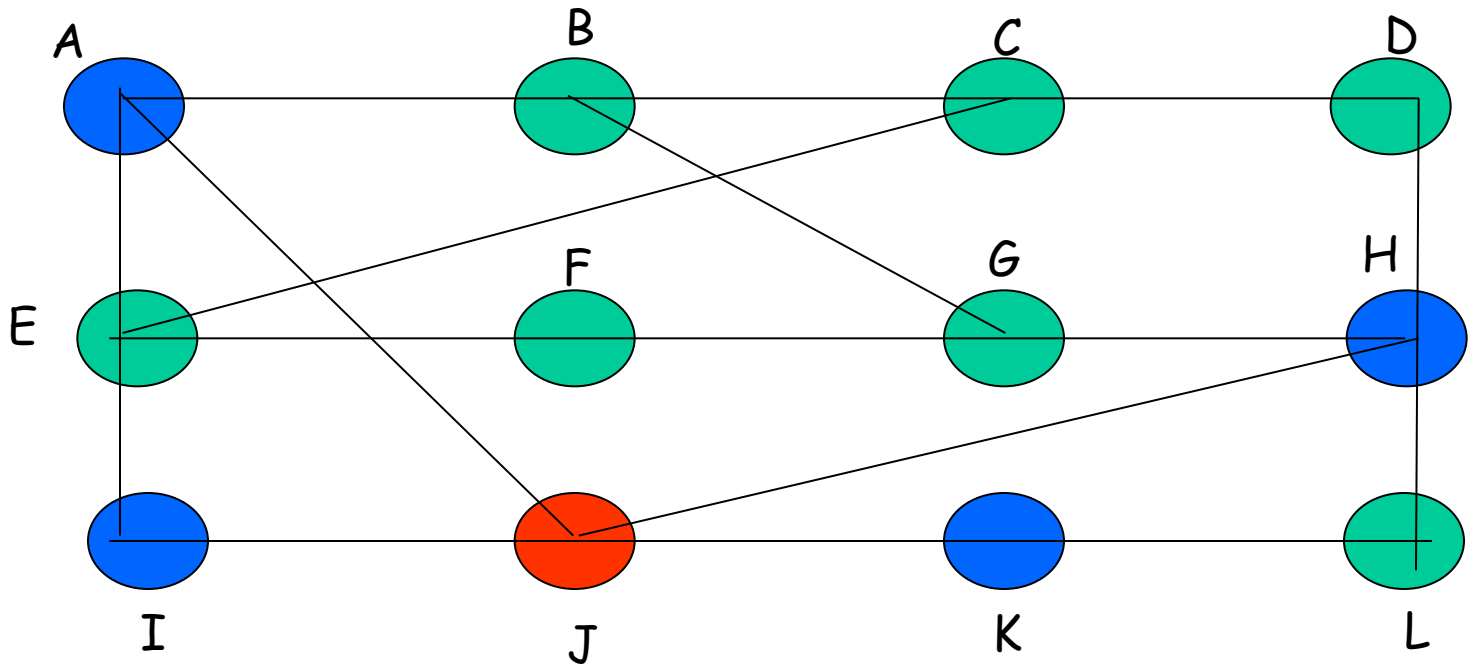
- each node notifies neighbors *only* when its DV changes
 - neighbors then notify their neighbors if necessary

Each node:

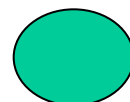


Example Topology

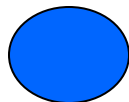
	J
A	8
B	∞
C	∞
D	∞
E	∞
F	∞
G	∞
H	12
I	10
J	0
K	6
L	∞



: J



others



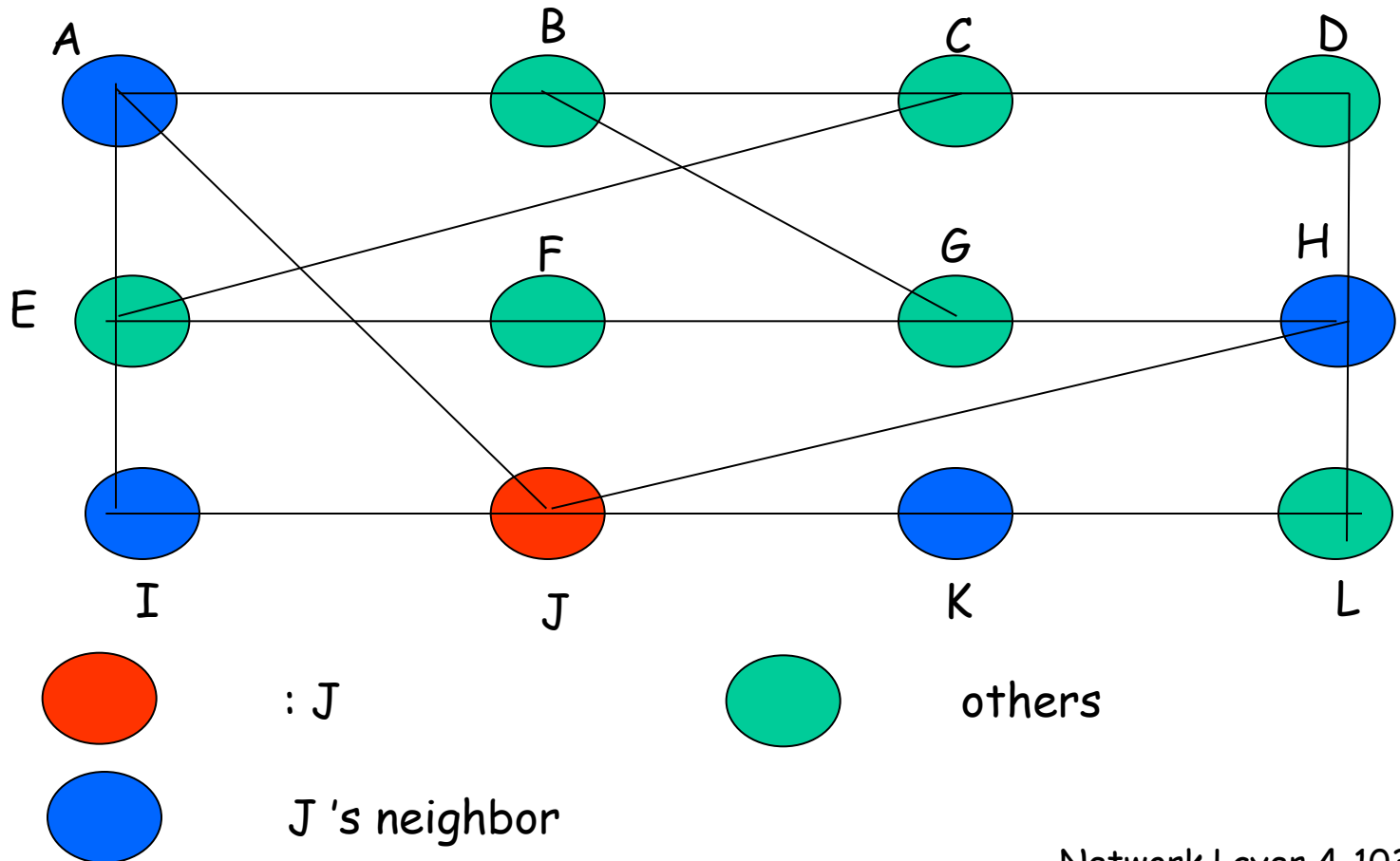
J's neighbor

Example Topology

JA delay
=8

A's distance vector to J

A	0
B	12
C	∞
D	∞
E	14
F	∞
G	∞
H	∞
I	∞
J	8
K	∞
L	∞

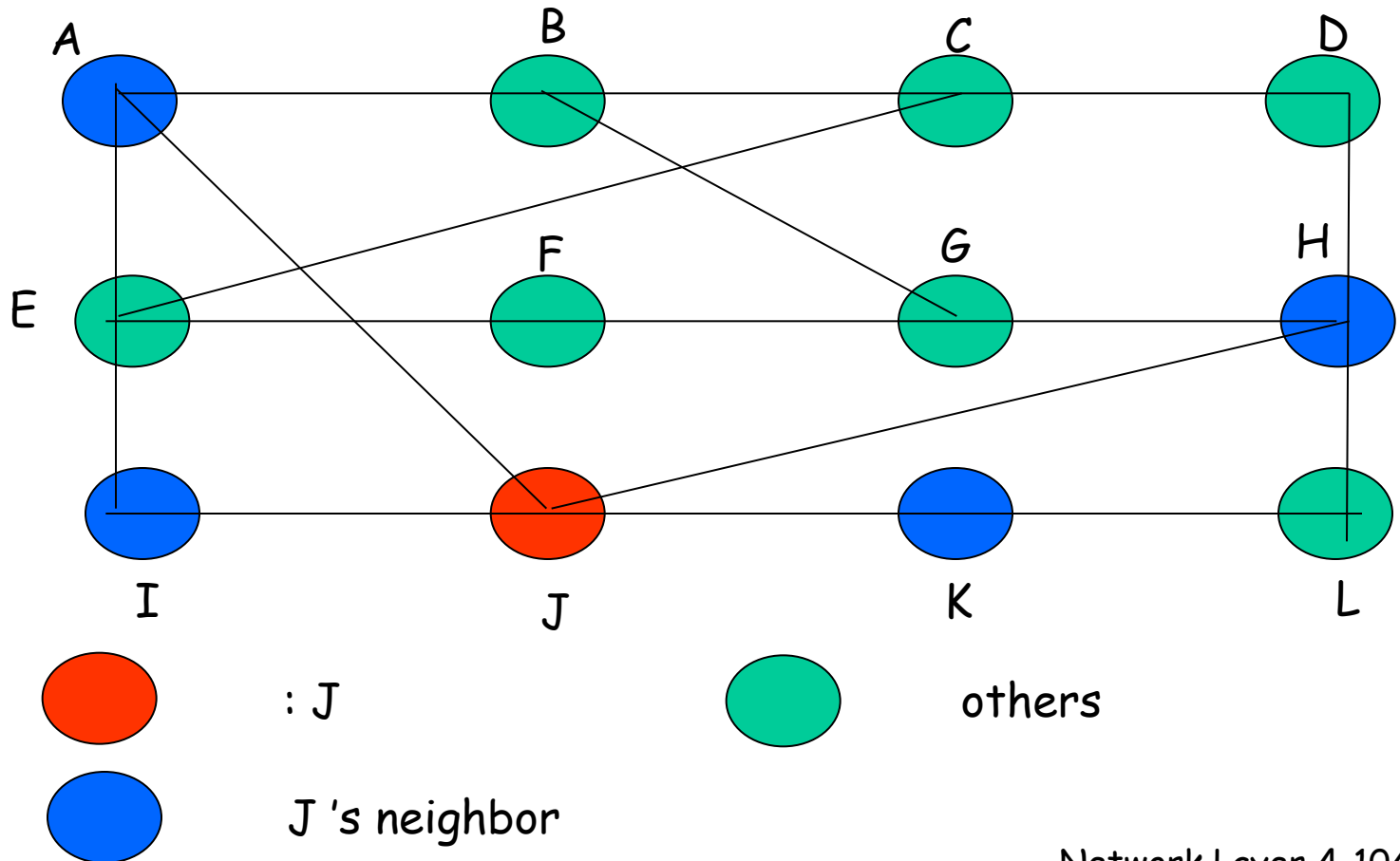


Example Topology

JH delay
=12

H's distance vector to J

A	∞
B	∞
C	∞
D	17
E	∞
F	∞
G	22
H	0
I	∞
J	12
K	∞
L	9



JA delay
=8

JH delay
=12

JI delay
=10

JK delay
=6

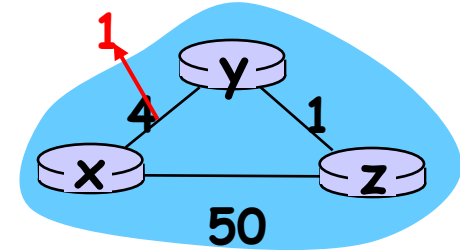
	J	A	H	I	K	I
A	8	0	∞	∞	∞	8
B	∞	12	∞	∞	∞	20 (A)
C	∞	∞	∞	∞	∞	∞
D	∞	∞	17	∞	∞	29 (H)
E	∞	14	∞	9	∞	19 (I)
F	∞	∞	∞	∞	∞	∞
G	∞	∞	22	∞	∞	34 (H)
H	12	∞	0	∞	∞	12
I	10	∞	∞	0	∞	10
J	0	8	12	10	6	0
K	6	∞	∞	∞	0	6
L	∞	∞	9	∞	30	21 (H)

Network Layer 4-105

Distance Vector: link cost changes

Link cost changes:

- ❑ node detects local link cost change
- ❑ updates routing info, recalculates distance vector
- ❑ if DV changes, notify neighbors



At time t_0 , y detects the link-cost change, updates its DV, and informs its neighbors.

At time t_1 , z receives the update from y and updates its table. It computes a new least cost to x and sends its neighbors its DV.

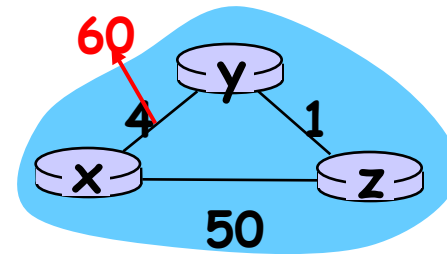
At time t_2 , y receives z 's update and updates its distance table. y 's least costs do not change and hence y does *not* send any message to z .

“good
news
travels
fast”

Distance Vector: link cost changes

Link cost changes:

- ❑ good news travels fast
- ❑ bad news travels slow - "count to infinity" problem!
- ❑ 44 iterations before algorithm stabilizes: see text



Poisoned reverse:

- ❑ If Z routes through Y to get to X :
 - Z tells Y its (Z's) distance to X is infinite (so Y won't route to X via Z)
- ❑ will this completely solve count to infinity problem?

Comparison of LS and DV algorithms

Message complexity

- LS: with n nodes, E links, $O(nE)$ msgs sent each
- DV: exchange between neighbors only
 - convergence time varies

Speed of Convergence

- LS: $O(n^2)$ algorithm requires $O(nE)$ msgs
 - may have oscillations
- DV: convergence time varies
 - may be routing loops
 - count-to-infinity problem

Robustness: what happens if router malfunctions?

LS:

- node can advertise incorrect *link* cost
- each node computes only its *own* table

DV:

- DV node can advertise incorrect *path* cost
- each node's table used by others
 - error propagate thru network

Hierarchical Routing

Hierarchical Routing

Our routing study thus far - idealization

- all routers identical (in performance (bps, pps(packet per second), purpose(access, edge, core routers, etc.)
- network "flat"

... *not true in practice*

scale: with 200 million destinations:

- can't store all dest's in routing tables!
- routing table exchange would swamp links!

administrative autonomy

- internet = network of networks
- each network admin may want to control routing in its own network

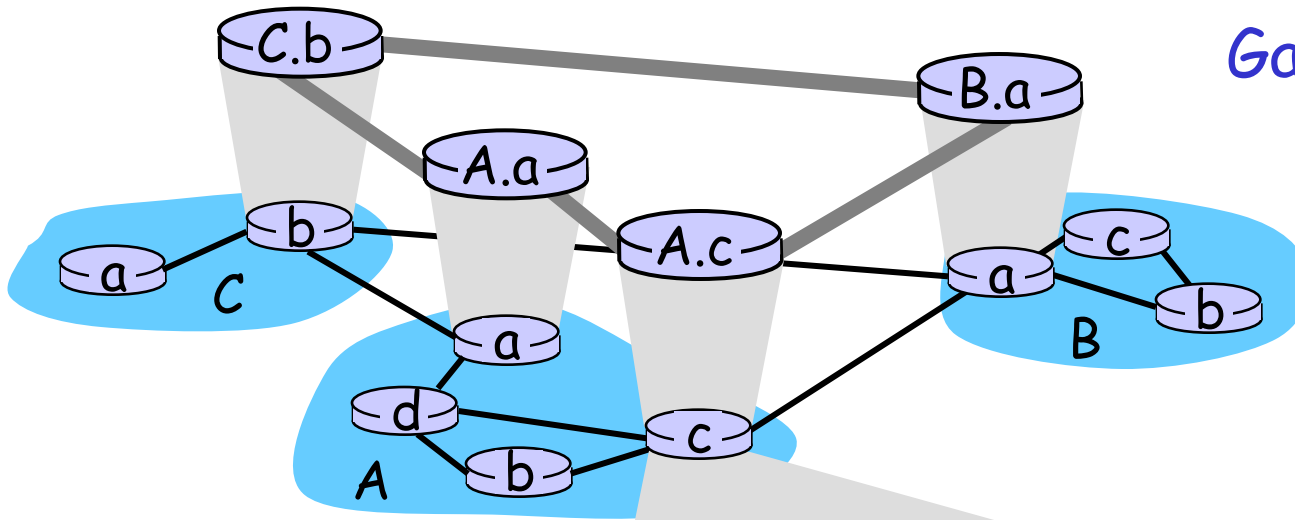
Hierarchical Routing

- aggregate routers into regions, "Autonomous Systems" (AS)
- routers in same AS run same routing protocol
 - "intra-AS" routing protocol
 - routers in different AS can run different intra-AS routing protocol

gateway routers

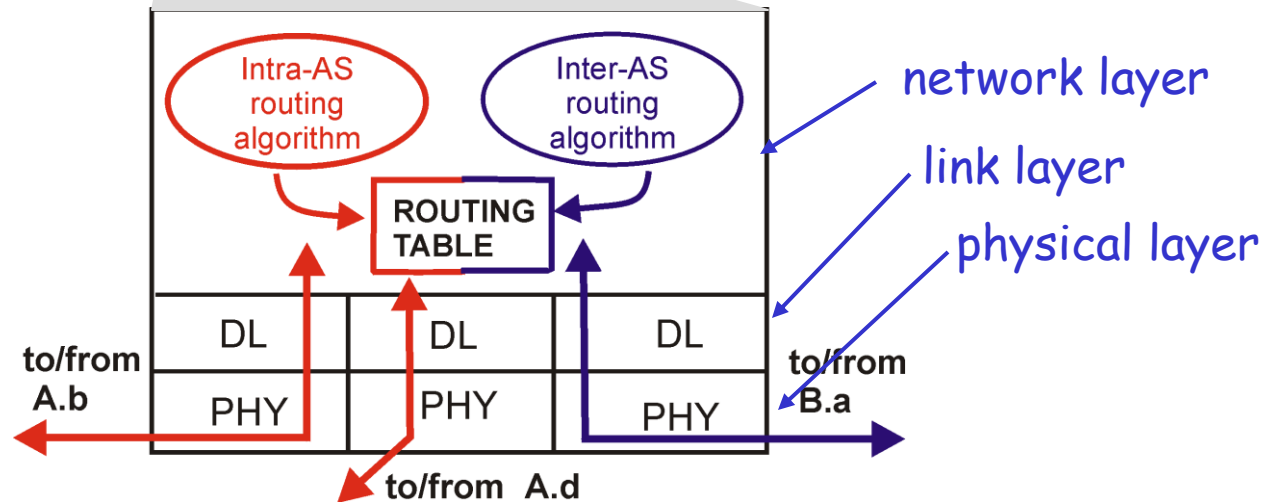
- special routers in AS
- run intra-AS routing protocol with all other routers in AS
- *also* responsible for routing to destinations outside AS
 - run *inter-AS routing protocol* with other gateway routers

Intra-AS and Inter-AS routing

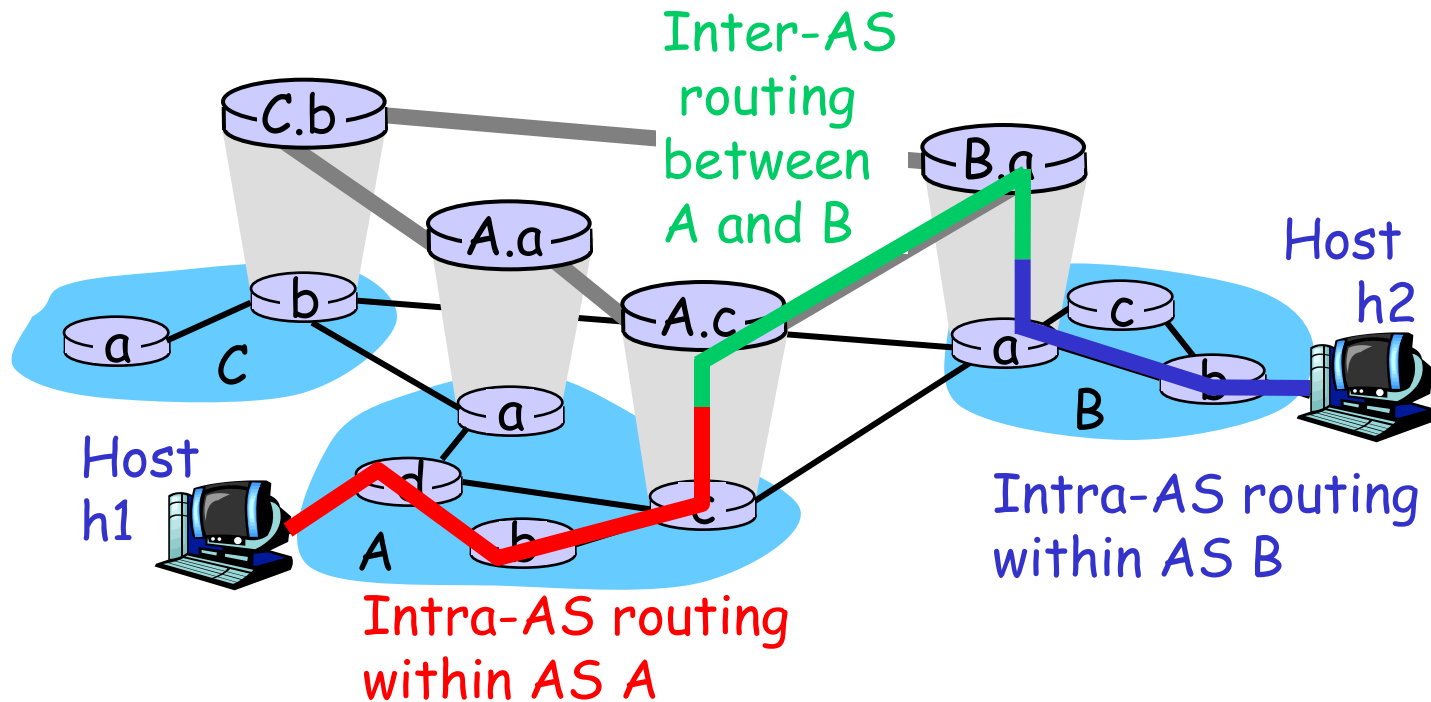


Gateways:

- perform inter-AS routing amongst themselves
- perform intra-AS routing with other routers in their AS



Intra-AS and Inter-AS routing



- We'll examine specific inter-AS and intra-AS Internet routing protocols shortly

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

Routing in the Internet

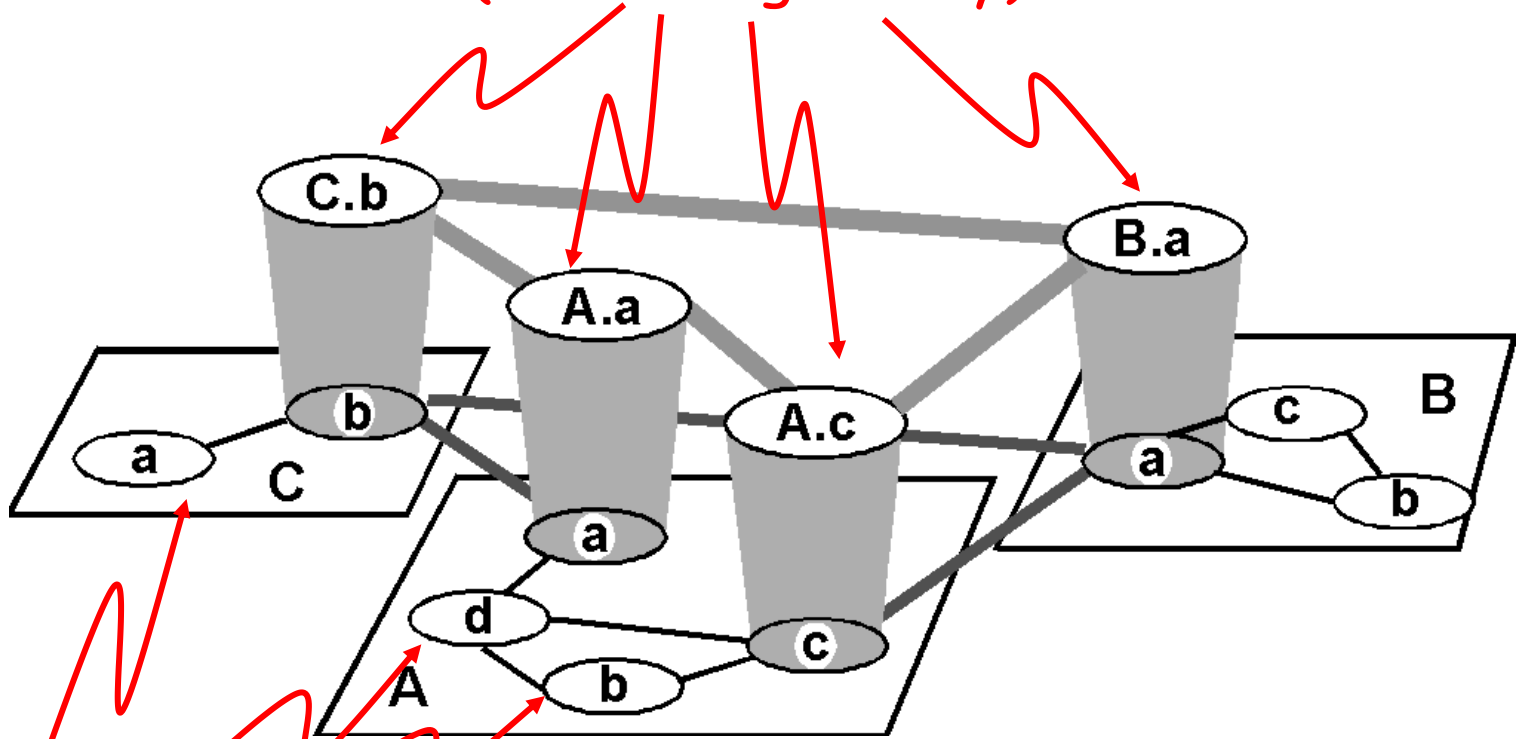
- The Global Internet consists of **Autonomous Systems (AS)** interconnected with each other:
 - **Stub AS**: small corporation: one connection to other AS's
 - **Multihomed AS**: large corporation (no transit): multiple connections to other AS's
 - **Transit AS**: provider, hooking many AS's together

Routing in the Internet (cont'd)

- **Two-level routing:**
 - **Intra-AS:** administrator responsible for choice of routing algorithm within network
 - **Inter-AS:** unique standard for inter-AS routing:
BGP

Internet AS Hierarchy

Inter-AS border (exterior gateway) routers



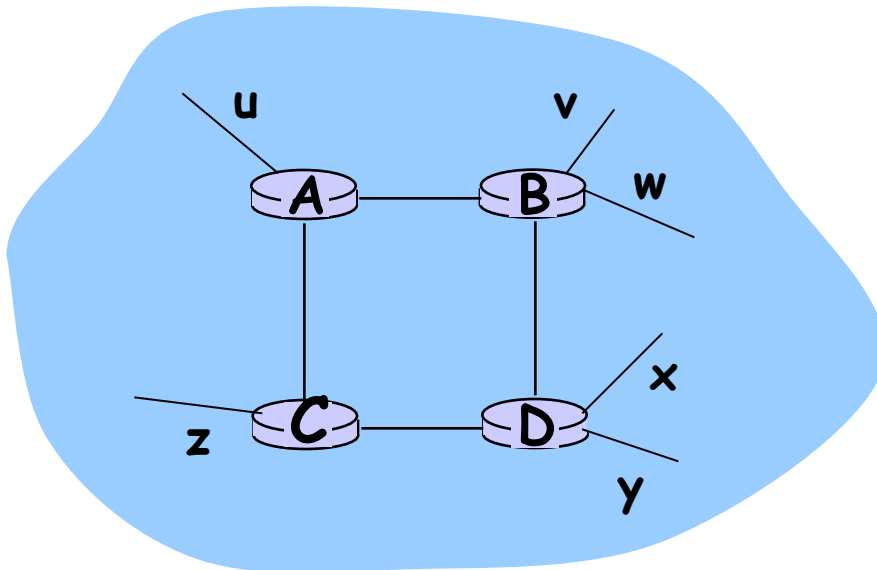
Intra-AS interior (gateway) routers

Intra-AS Routing

- Also known as **Interior Gateway Protocols (IGP)**
- Most common Intra-AS routing protocols:
 - **RIP**: Routing Information Protocol
 - **OSPF**: Open Shortest Path First
 - **IGRP**: Interior Gateway Routing Protocol (Cisco proprietary)

RIP (Routing Information Protocol)

- distance vector algorithm
- included in BSD-UNIX Distribution in 1982
- distance metric: # of hops (max = 15 hops)



From router A to subnets:

<u>destination</u>	<u>hops</u>
u	1
v	2
w	2
x	3
y	3
z	2

RIP advertisements

- distance vectors: exchanged among neighbors every 30 sec via Response Message (also called **advertisement**)
- each advertisement: list of up to 25 destination subnets within AS

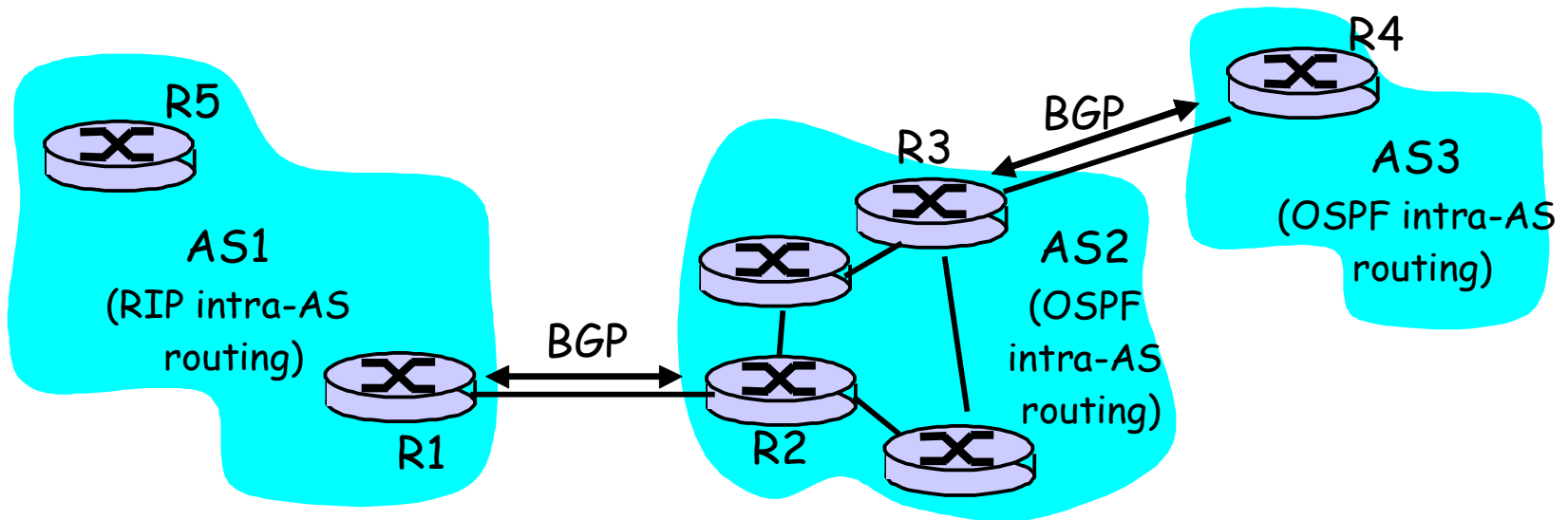
OSPF (Open Shortest Path First)

- “open”: publicly available
- Uses Link State algorithm
 - LS packet dissemination
 - Topology map at each node
 - Route computation using Dijkstra’s algorithm
- OSPF advertisement carries **one entry per neighbor router**
- Advertisements disseminated to **entire AS** (via flooding)
 - **Carried in OSPF messages directly over IP** (rather than TCP or UDP)

OSPF "advanced" features (not in RIP)

- **Security:** all OSPF messages **authenticated** (to prevent malicious intrusion)
- **Multiple** same-cost **paths** allowed (only one path in RIP)
- For each link, multiple cost metrics for different **TOS** (e.g., satellite link cost set "low" for best effort; high for real time)
- Integrated uni- and **multicast** support:
 - Multicast OSPF (MOSPF) uses same topology data base as OSPF
- **Hierarchical** OSPF in large domains.

Inter-AS routing in the Internet: BGP

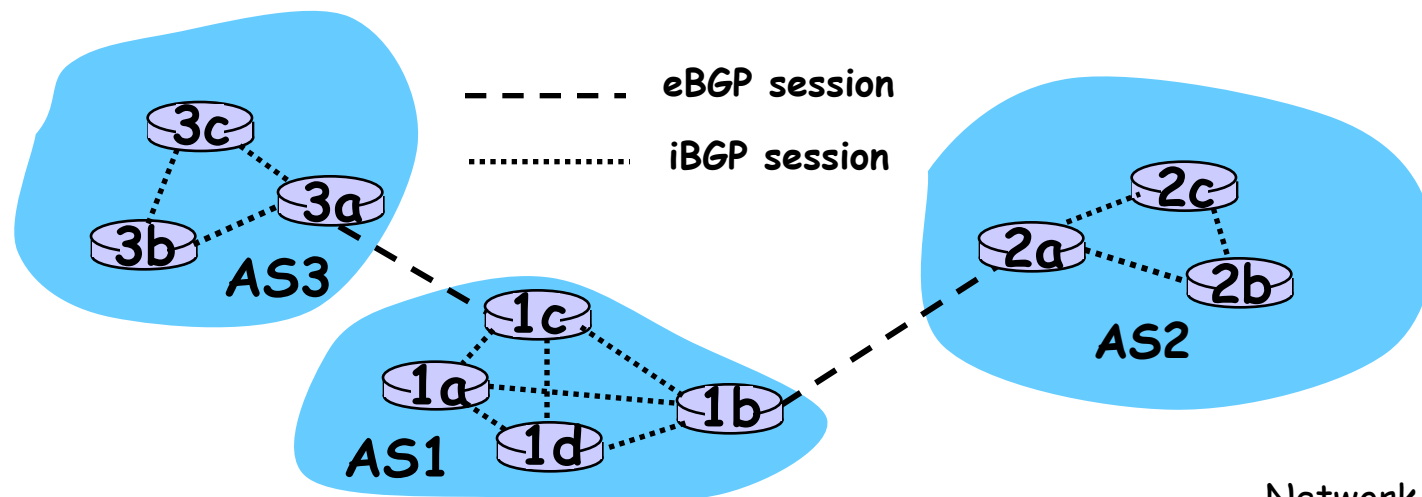


Internet inter-AS routing: BGP

- **BGP (Border Gateway Protocol):** *the de facto standard*
- BGP provides each AS a means to:
 1. Obtain subnet reachability information from neighboring ASs.
 2. Propagate reachability information to all AS-internal routers.
 3. Determine "good" routes to subnets based on reachability information and policy.
- allows subnet to advertise its existence to rest of Internet: *"I am here"*

BGP basics

- pairs of routers (BGP peers) exchange routing info over semi-permanent TCP connections: **BGP sessions**
 - BGP sessions need not correspond to physical links.
- when AS2 advertises a prefix to AS1:
 - AS2 *promises* it will forward datagrams towards that prefix.
 - AS2 can aggregate prefixes in its advertisement



The end. 😊

Homework

- R1, R7, R10, R11, R12, R16, R20
- P10, P18, P19, P25, P32