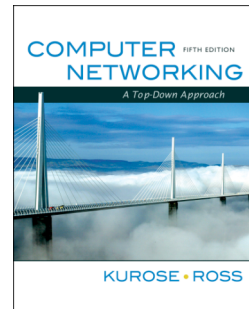


# Introduction to Computer Networking

Guy Leduc

## Chapter 4 Network Layer



*Computer Networking:  
A Top Down Approach,  
5<sup>th</sup> edition.*

Jim Kurose, Keith Ross  
Addison-Wesley, April  
2009.

© From Computer Networking, by Kurose&Ross

Network Layer 4-1

## Chapter 4: Network Layer

### Chapter goals:

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

© From Computer Networking, by Kurose&Ross

Network Layer 4-2

## Chapter 4: Network Layer

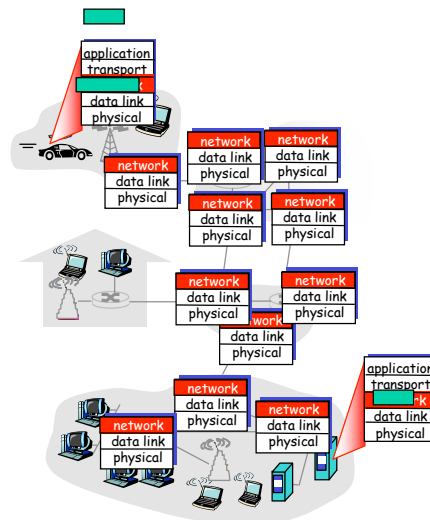
- ❑ 4.1 Introduction
- ❑ 4.2 Virtual circuit and datagram networks
- ❑ 4.3 What's inside a router
- ❑ 4.4 IP: Internet Protocol
  - Datagram format
  - IPv4 addressing
  - ICMP
  - IPv6
- ❑ 4.5 Routing algorithms
  - Link state
  - Distance Vector
  - Hierarchical routing
- ❑ 4.6 Routing in the Internet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

© From Computer Networking, by Kurose&Ross

Network Layer 4-3

### 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



© From Computer Networking, by Kurose&Ross

Network Layer 4-4

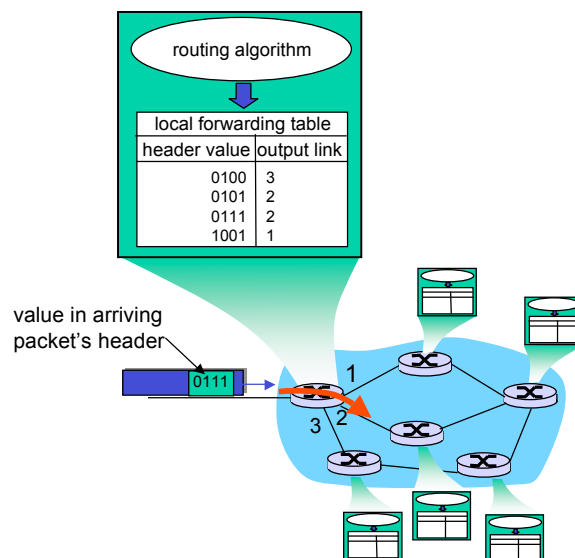
## Two Key Network-Layer Functions

- ❑ *forwarding*: move packets from router's input to appropriate router output
  - ❑ *routing*: determine route taken by packets from source to dest.
    - *routing algorithms*
- analogy:
- ❑ *routing*: process of planning trip from source to dest
  - ❑ *forwarding*: process of getting through single interchange

© From Computer Networking, by Kurose&Ross

Network Layer 4-5

## Interplay between routing and forwarding



© From Computer Networking, by Kurose&Ross

Network Layer 4-6

## Connection setup

- ❑ 3<sup>rd</sup> important function in some network architectures:
  - MPLS, ATM, frame relay, X.25
- ❑ before datagrams flow, two end hosts *and* intervening routers establish virtual connection
  - routers get involved
- ❑ network vs transport layer connection service:
  - **network**: between two hosts (may also involve intervening routers in case of VCs)
  - **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

© From Computer Networking, by Kurose&Ross

Network Layer 4-9

## 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

© From Computer Networking, by Kurose&Ross

Network Layer 4-10

## 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

© From Computer Networking, by Kurose&Ross

Network Layer 4-11

## 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)

© From Computer Networking, by Kurose&Ross

Network Layer 4-12

## VC implementation

a VC consists of:

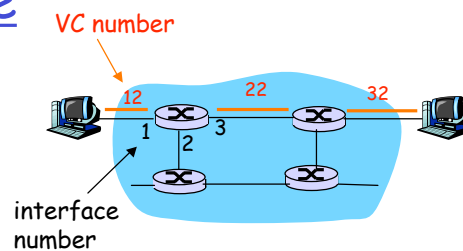
1. path from source to destination
  2. VC numbers, one number for each link along path
  3. entries in forwarding tables in routers along path
- ❑ packet belonging to VC carries VC number (rather than dest address)
  - ❑ VC number can be changed on each link.
    - New VC number comes from forwarding table

© From Computer Networking, by Kurose&Ross

Network Layer 4-13

## Forwarding table

Forwarding table in northwest router:



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

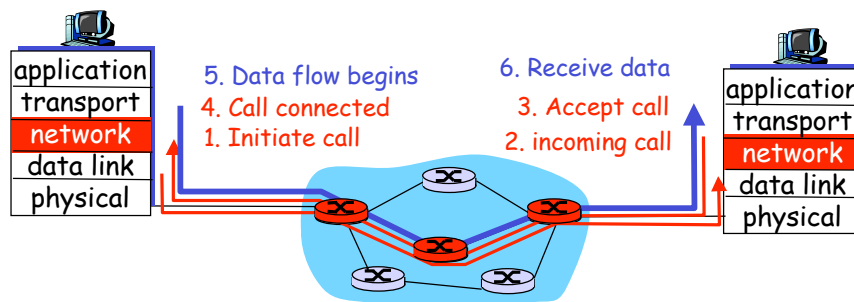
**Routers maintain connection state information!**

© From Computer Networking, by Kurose&Ross

Network Layer 4-14

## Virtual circuits: signaling protocols

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

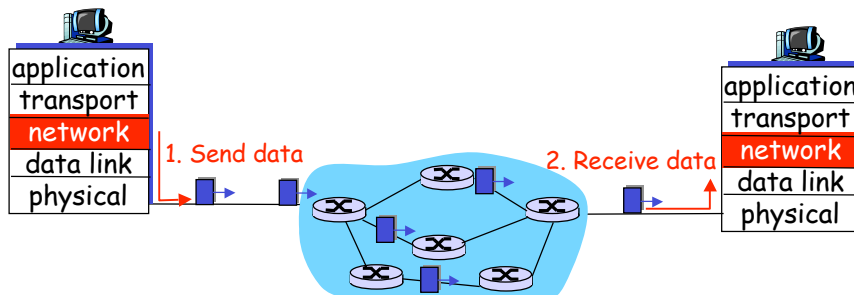


© From Computer Networking, by Kurose&Ross

Network Layer 4-15

## Datagram networks

- ❑ 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



© From Computer Networking, by Kurose&Ross

Network Layer 4-16



## 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

© From Computer Networking, by Kurose&Ross

Network Layer 4-17

## 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 00010000 10101010

Which interface?

© From Computer Networking, by Kurose&Ross

Network Layer 4-18

## 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 (VC)

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

© From Computer Networking, by Kurose&Ross

Network Layer 4-19

## 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

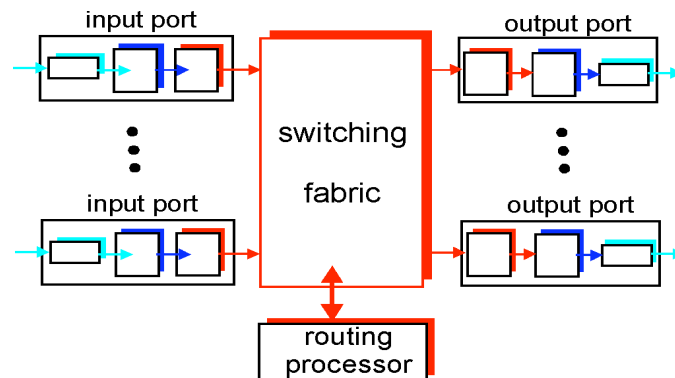
© From Computer Networking, by Kurose&Ross

Network Layer 4-20

## Router Architecture Overview

Two key router functions:

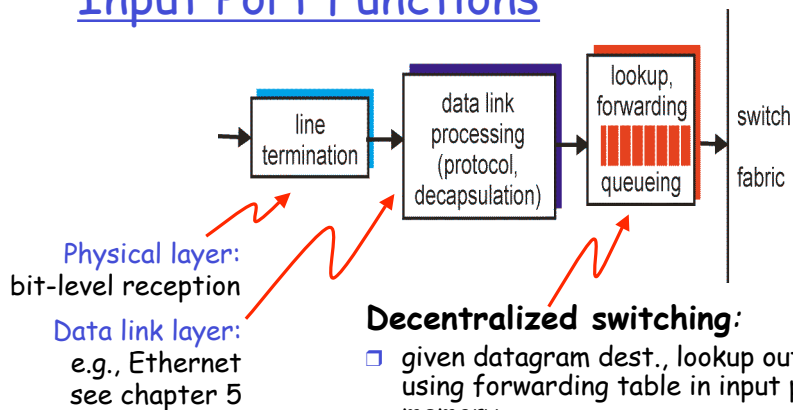
- run routing algorithms/protocol (RIP, OSPF, BGP)
- *forwarding* datagrams from incoming to outgoing link



© From Computer Networking, by Kurose&Ross

Network Layer 4-21

## Input Port Functions



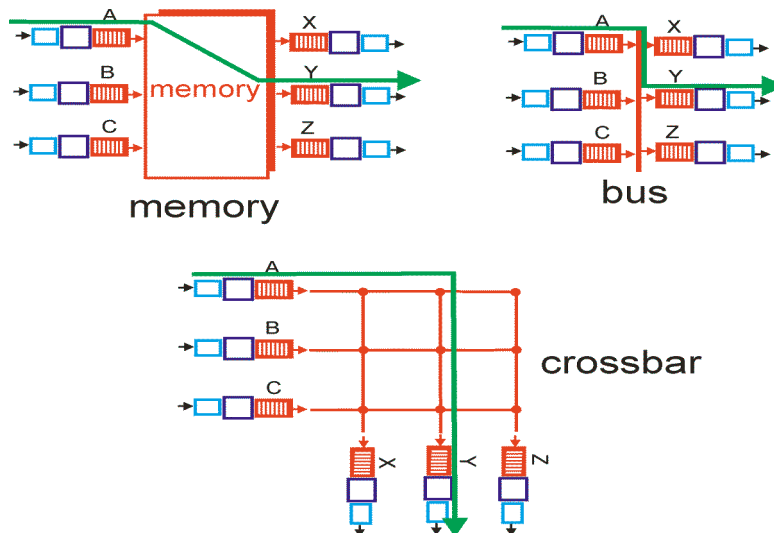
### **Decentralized switching:**

- 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

© From Computer Networking, by Kurose&Ross

Network Layer 4-22

## Three types of switching fabrics



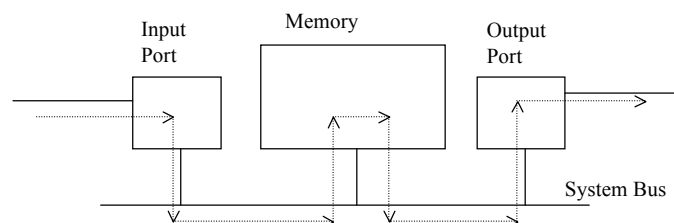
© From Computer Networking, by Kurose&Ross

Network Layer 4-23

## Switching Via Memory

### First generation routers:

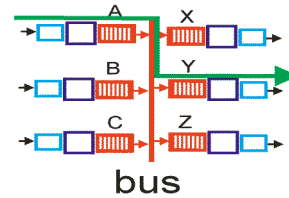
- ❑ traditional computers with switching under direct control of CPU
- ❑ packet copied to system's memory
- ❑ speed limited by memory bandwidth (2 bus crossings per datagram)



© From Computer Networking, by Kurose&Ross

Network Layer 4-24

## 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

© From Computer Networking, by Kurose&Ross

Network Layer 4-25

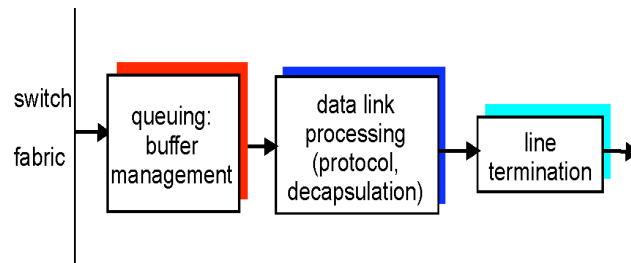
## Switching Via An Interconnection Network

- ❑ overcome bus bandwidth limitations
- ❑ Banyan networks, other interconnection nets initially developed to connect processors in multiprocessor architectures
- ❑ advanced design: fragmenting datagram into fixed length cells, switch cells through the fabric.
- ❑ Cisco 12000: switches 60 Gbps through the interconnection network

© From Computer Networking, by Kurose&Ross

Network Layer 4-26

## Output Ports

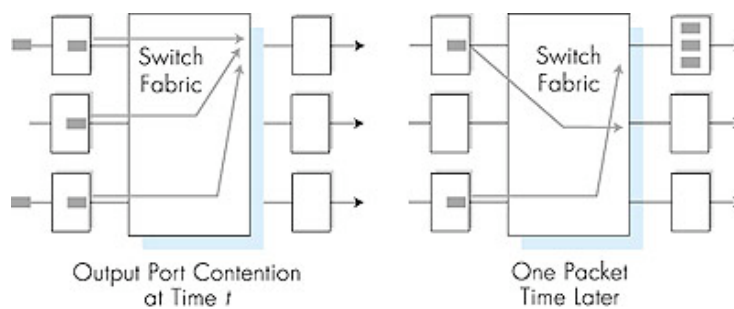


- ❑ *Buffering* required when datagrams arrive from fabric faster than the transmission rate
- ❑ *Scheduling discipline* chooses among queued datagrams for transmission

© From Computer Networking, by Kurose&Ross

Network Layer 4-27

## Output port queueing



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

© From Computer Networking, by Kurose&Ross

Network Layer 4-28

## How much buffering?

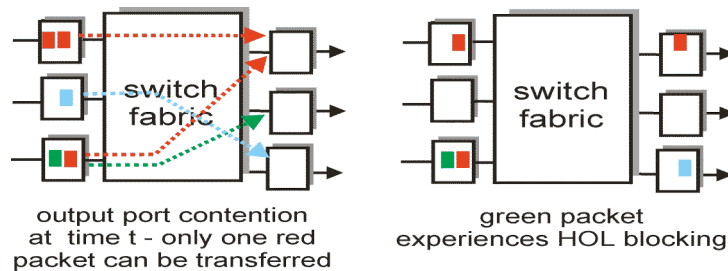
- RFC 3439 rule of thumb: average buffering equal to "typical" RTT (say 250 msec) times link capacity  $C$ 
  - e.g.,  $C = 10$  Gbps link: 2.5 Gbit buffer
- Recent recommendation: with  $N$  TCP flows, buffering equal to  $\frac{RTT \cdot C}{\sqrt{N}}$

© From Computer Networking, by Kurose&Ross

Network Layer 4-29

## Input Port Queuing

- Fabric slower than input ports combined  $\rightarrow$  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!*



© From Computer Networking, by Kurose&Ross

Network Layer 4-30

## Chapter 4: Network Layer

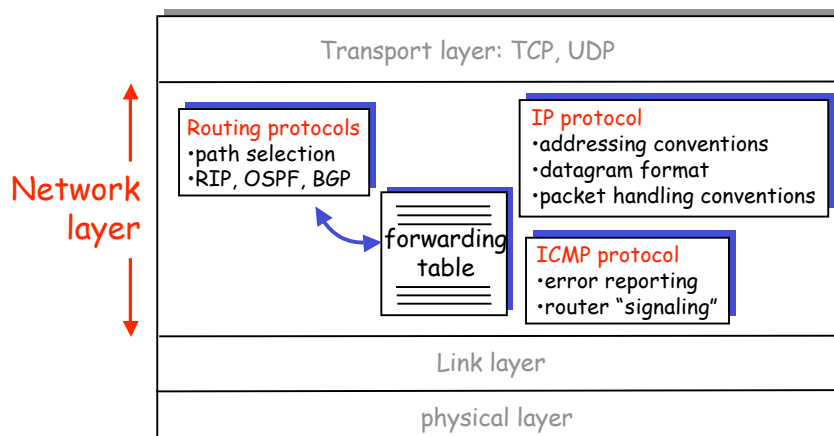
- 4.1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
- 4.4 IP: Internet Protocol
  - Datagram format
  - IPv4 addressing
  - ICMP
  - IPv6
- 4.5 Routing algorithms
  - Link state
  - Distance Vector
  - Hierarchical routing
- 4.6 Routing in the Internet
  - RIP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

© From Computer Networking, by Kurose&Ross

Network Layer 4-31

## The Internet Network layer

Host, router network layer functions:



© From Computer Networking, by Kurose&Ross

Network Layer 4-32



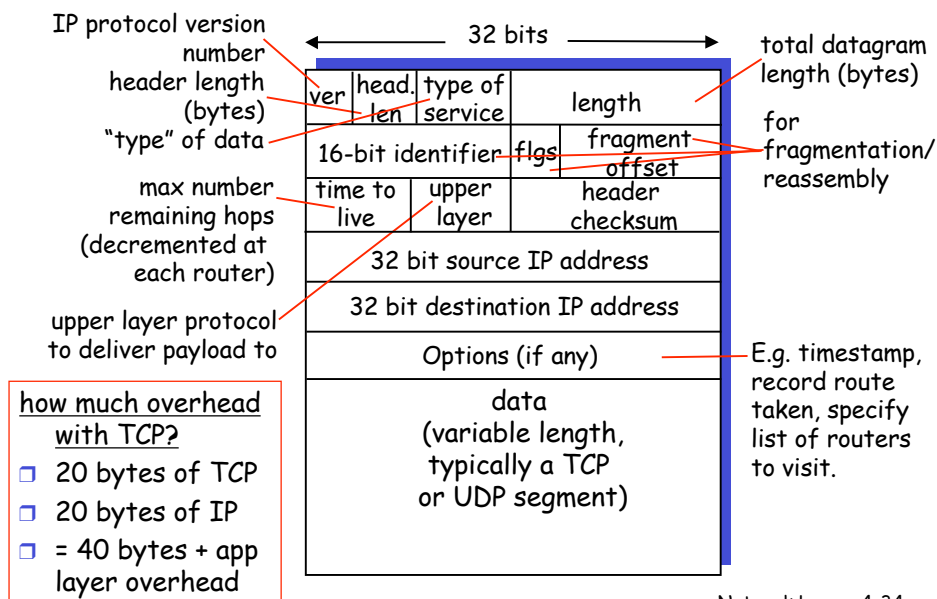
## Chapter 4: Network Layer

- 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

© From Computer Networking, by Kurose&Ross

Network Layer 4-33

### IP datagram format

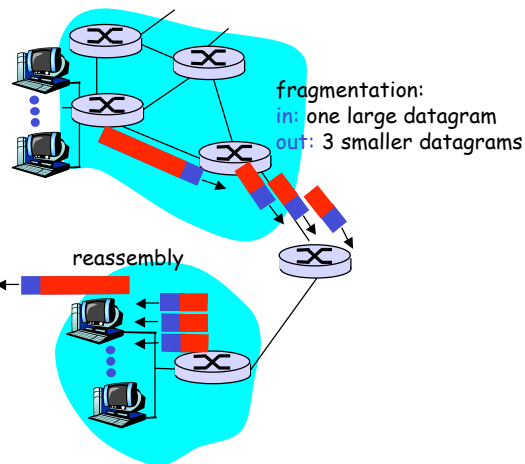


© From Computer Networking, by Kurose&Ross

Network Layer 4-34

## 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



© From Computer Networking, by Kurose&Ross

Network Layer 4-35

## 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
=1500	=x	=1	=185
=1040	=x	=0	=370

© From Computer Networking, by Kurose&Ross

Network Layer 4-36

## Avoiding fragmentation

- ❑ To avoid fragmentation, the source must know the minimal MTU of the path
- ❑ Path MTU discovery (trial and error)
  - Send an IP packet with the "Don't fragment flag" set
  - Routers may be forced to discard the packet
  - If the source receives an ICMP error message (see later), it tries again with a size smaller than the MTU indicated in the ICMP packet
- ❑ Drawback
  - Relies on routers properly returning ICMP error message
  - Also, congestion could discard ICMP messages
  - Also, the route may change afterwards
  - So fragmentation can happen anyway

Network Layer 4-37

## Chapter 4: Network Layer

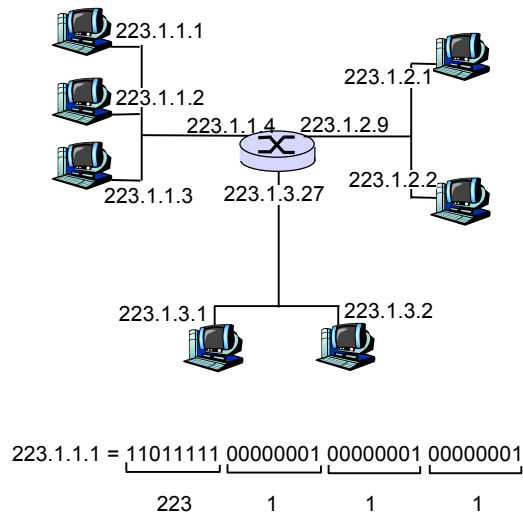
- |  |   |
|--|---|
| <ul style="list-style-type: none"><li>❑ 4.1 Introduction</li><li>❑ 4.2 Virtual circuit and datagram networks</li><li>❑ 4.3 What's inside a router</li><li>❑ 4.4 IP: Internet Protocol<ul style="list-style-type: none"><li>○ Datagram format</li><li>○ IPv4 addressing</li><li>○ ICMP</li><li>○ IPv6</li></ul></li></ul> | <ul style="list-style-type: none"><li>❑ 4.5 Routing algorithms<ul style="list-style-type: none"><li>○ Link state</li><li>○ Distance Vector</li><li>○ Hierarchical routing</li></ul></li><li>❑ 4.6 Routing in the Internet<ul style="list-style-type: none"><li>○ RIP</li><li>○ OSPF</li><li>○ BGP</li></ul></li></ul> |
|--|---|

© From Computer Networking, by Kurose&Ross

Network Layer 4-38

## 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 interface
  - IP addresses associated with each interface

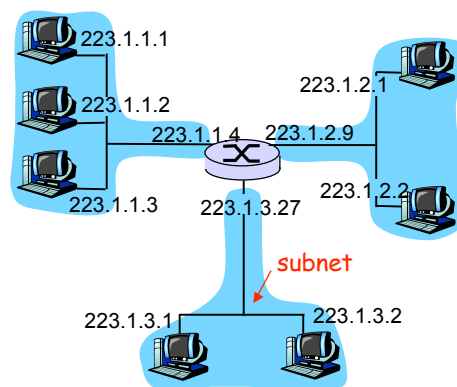


© From Computer Networking, by Kurose&Ross

Network Layer 4-39

## 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

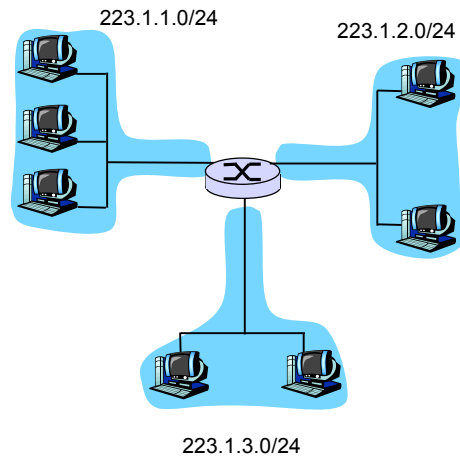
© From Computer Networking, by Kurose&Ross

Network Layer 4-40

# 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**.



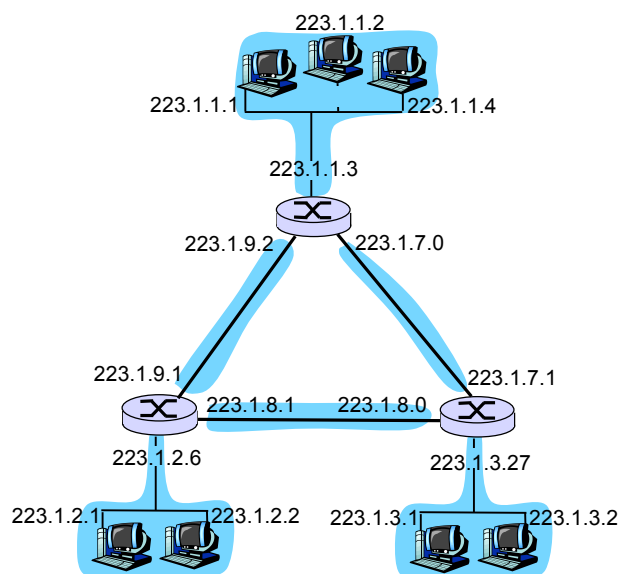
Subnet mask: /24

© From Computer Networking, by Kurose&Ross

Network Layer 4-41

# Subnets

How many?



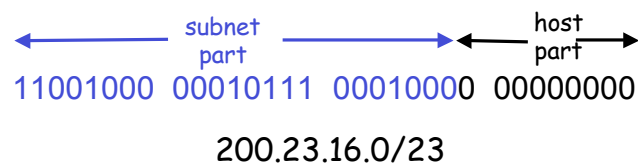
© From Computer Networking, by Kurose&Ross

Network Layer 4-42

## 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



© From Computer Networking, by Kurose&Ross

Network Layer 4-43

## 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 a server
  - "plug-and-play"

© From Computer Networking, by Kurose&Ross

Network Layer 4-44

## 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)

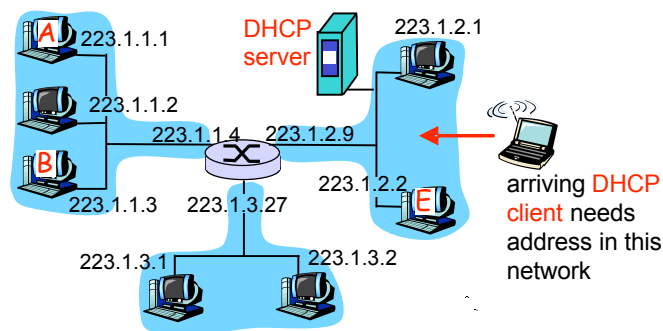
DHCP overview:

- 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

© From Computer Networking, by Kurose&Ross

Network Layer 4-45

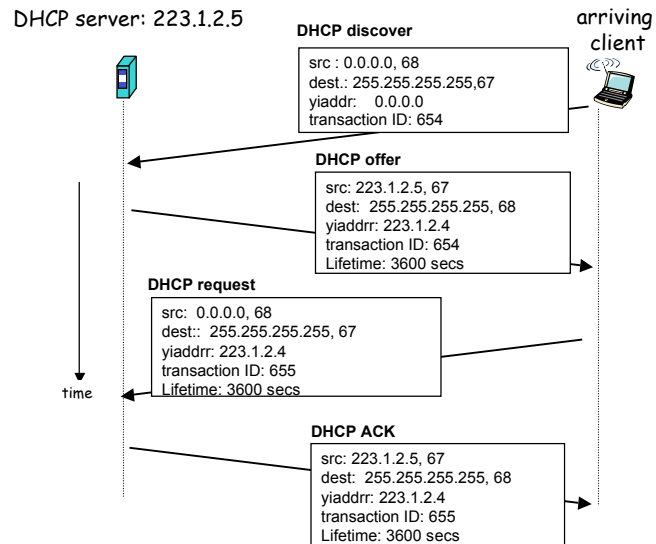
## DHCP client-server scenario



© From Computer Networking, by Kurose&Ross

Network Layer 4-46

## DHCP client-server scenario



© From Computer Networking, by Kurose&Ross

Network Layer 4-47

## DHCP: more than IP address

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

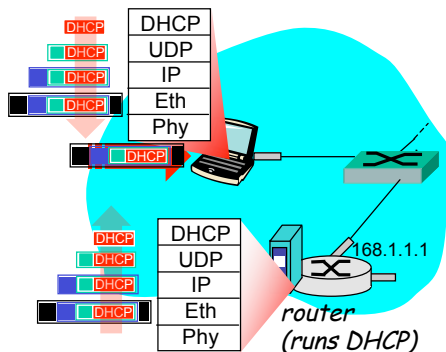
- address of first-hop router for client
- name and IP address of DNS server
- network mask (indicating network versus host portion of address)

© From Computer Networking, by Kurose&Ross

Network Layer 4-48



## DHCP: example

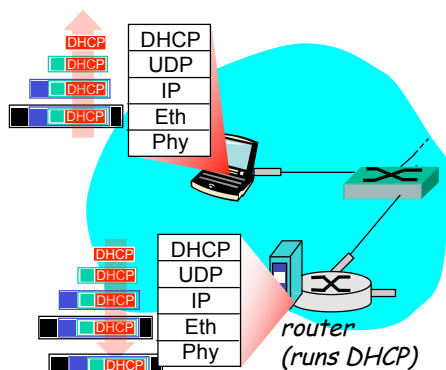


- connecting laptop needs its IP address, addr of first-hop router, addr of DNS server: use DHCP
- DHCP request encapsulated in UDP, encapsulated in IP, encapsulated in Ethernet
- Ethernet frame broadcast (dest: FFFFFFFFFF) on LAN, received at router running DHCP server
- Ethernet demux'ed to IP demux'ed, UDP demux'ed to DHCP

© From Computer Networking, by Kurose&Ross

Network Layer 4-49

## DHCP: example



- DHCP 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, demux'ing up to DHCP at client
- client now knows its IP address, name and IP address of DNS server, IP address of its first-hop router

© From Computer Networking, by Kurose&Ross

Network Layer 4-50

## IP addresses: how to get one?

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

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

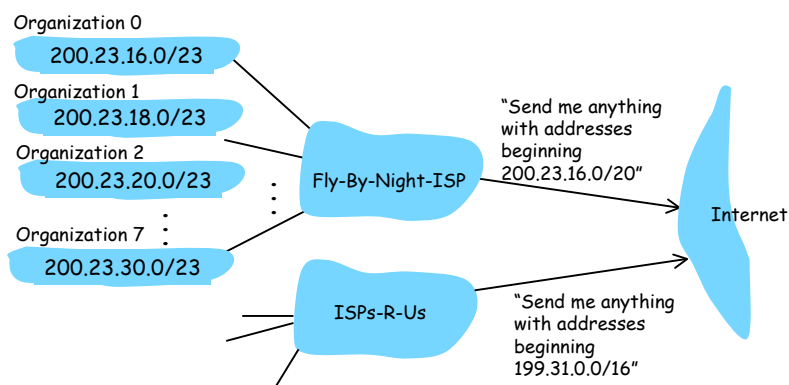
ISP's block	11001000	00010111	00010000	00000000	200.23.16.0/20
Organization 0	11001000	00010111	00010000	00000000	200.23.16.0/23
Organization 1	11001000	00010111	00010010	00000000	200.23.18.0/23
Organization 2	11001000	00010111	00010100	00000000	200.23.20.0/23
...	....	....	....	....	....
Organization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

© From Computer Networking, by Kurose&Ross

Network Layer 4-51

## Hierarchical addressing: route aggregation

Hierarchical addressing allows efficient advertisement of routing information:

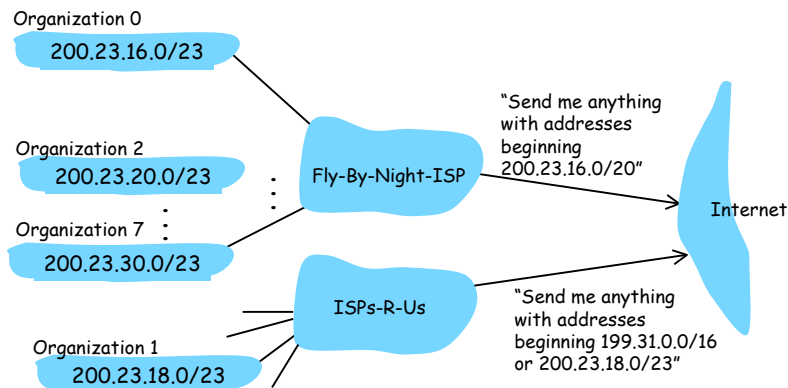


© From Computer Networking, by Kurose&Ross

Network Layer 4-52

## Hierarchical addressing: more specific routes

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

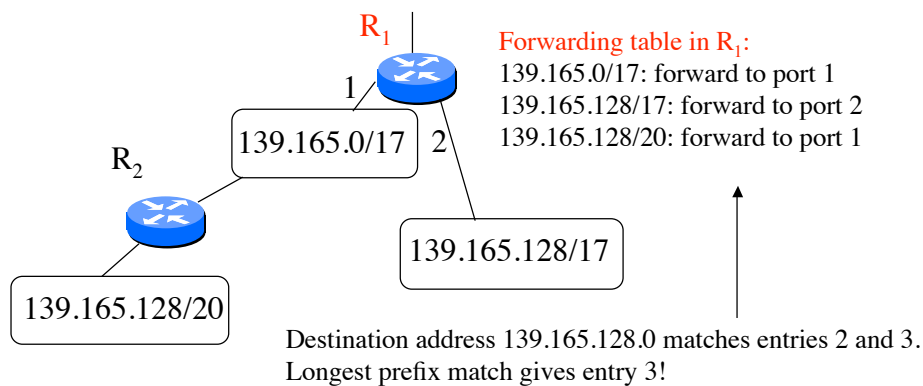


© From Computer Networking, by Kurose&Ross

Network Layer 4-53

## Longest prefix match - Another ex.

- Suppose a network has the address range 139.165.0.0/16
- It is first split into 139.165.0.0/17 and 139.165.128.0/17
- Then **part of** 139.165.128.0/17 (namely 139.165.128.0/20) is reallocated elsewhere.



Network Layer 4-54

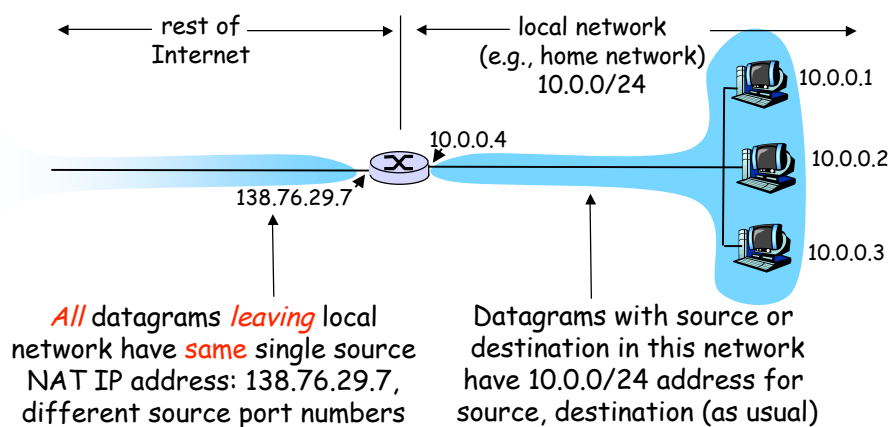
## IP addressing: the last word...

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

**A:** **ICANN**: Internet Corporation for Assigned Names and Numbers

- 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:
  - 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).

© From Computer Networking, by Kurose&Ross

Network Layer 4-57

## NAT: Network Address Translation

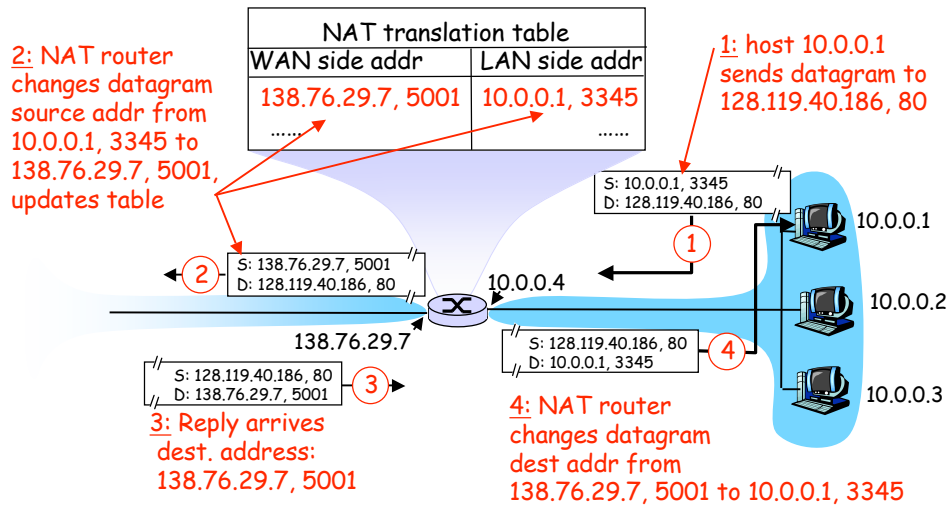
**Implementation:** NAT router must:

- *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

© From Computer Networking, by Kurose&Ross

Network Layer 4-58

## NAT: Network Address Translation



© From Computer Networking, by Kurose&Ross

Network Layer 4-59

## NAT: Network Address Translation

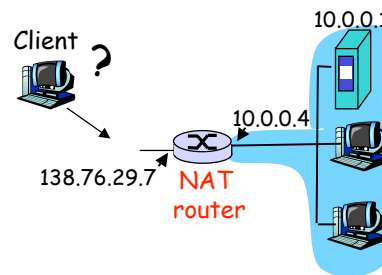
- ❑ 16-bit port-number field:
  - 60,000 simultaneous connections with a single LAN-side address!
- ❑ NAT is controversial:
  - routers should only process up to layer 3
  - violates end-to-end argument
    - NAT possibility must be taken into account by app designers, eg, P2P applications
  - address shortage should instead be solved by IPv6

© From Computer Networking, by Kurose&Ross

Network Layer 4-60

## 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 25000

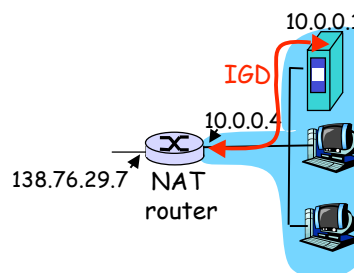


© From Computer Networking, by Kurose&Ross

Network Layer 4-61

## 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)



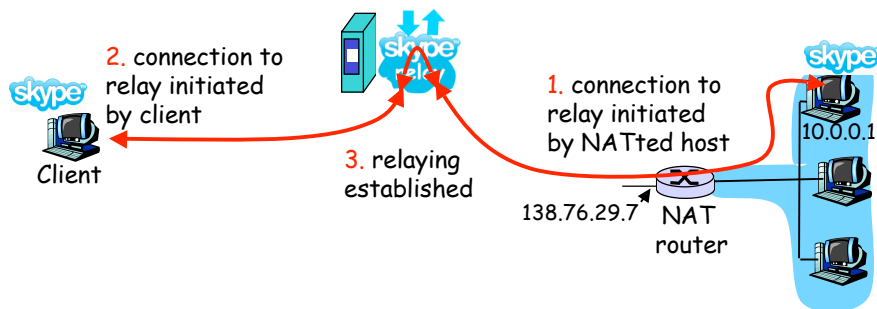
i.e., automate static NAT port map configuration

© From Computer Networking, by Kurose&Ross

Network Layer 4-62

## NAT traversal problem

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



© From Computer Networking, by Kurose&Ross

Network Layer 4-63

## 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

© From Computer Networking, by Kurose&Ross

Network Layer 4-64



## 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
- | Type | Code | description                                   |
|------|------|---|
| 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                                 |

© From Computer Networking, by Kurose&Ross

Network Layer 4-65

## 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 "port unreachable" packet (type 3, code 3)
  - ❑ When source gets this ICMP, stops.

© From Computer Networking, by Kurose&Ross

Network Layer 4-66

## 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

© From Computer Networking, by Kurose&Ross

Network Layer 4-67

## 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

© From Computer Networking, by Kurose&Ross

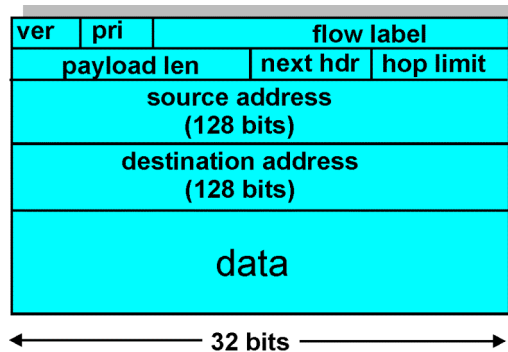
Network Layer 4-68

## 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



© From Computer Networking, by Kurose&Ross

Network Layer 4-69

## Other Changes from IPv4

- ❑ *(Header) 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

© From Computer Networking, by Kurose&Ross

Network Layer 4-70

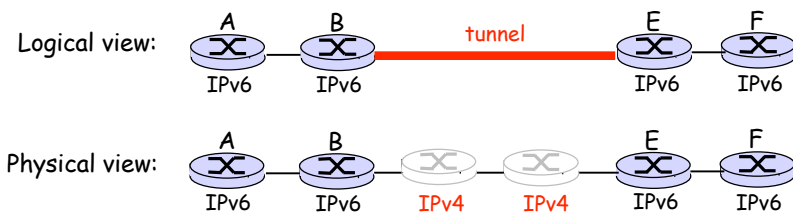
## Transition From IPv4 To IPv6

- ❑ Not all routers can be upgraded simultaneously
  - no "flag days"
  - How will the network operate with mixed IPv4 and IPv6 routers?
- ❑ **Tunneling:** IPv6 carried as payload in IPv4 datagram among IPv4 routers

© From Computer Networking, by Kurose&Ross

Network Layer 4-71

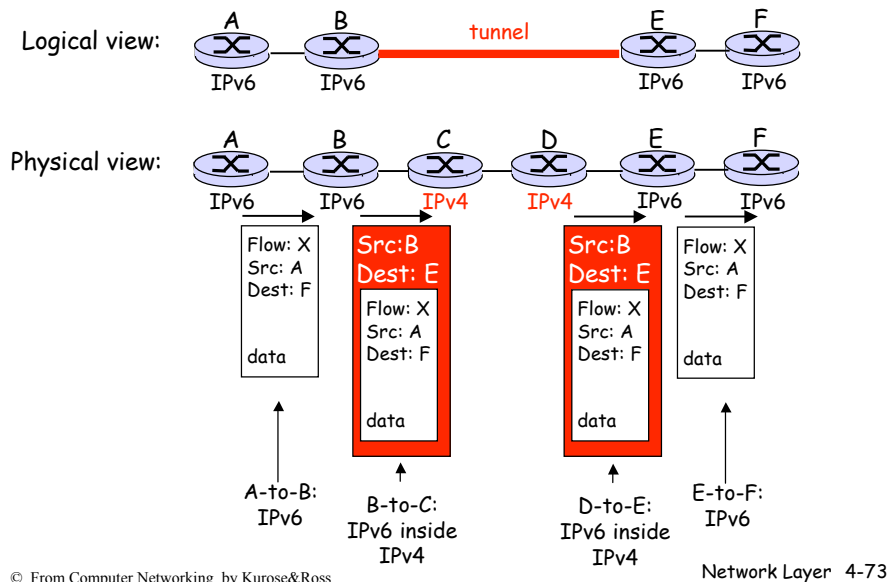
## Tunneling



© From Computer Networking, by Kurose&Ross

Network Layer 4-72

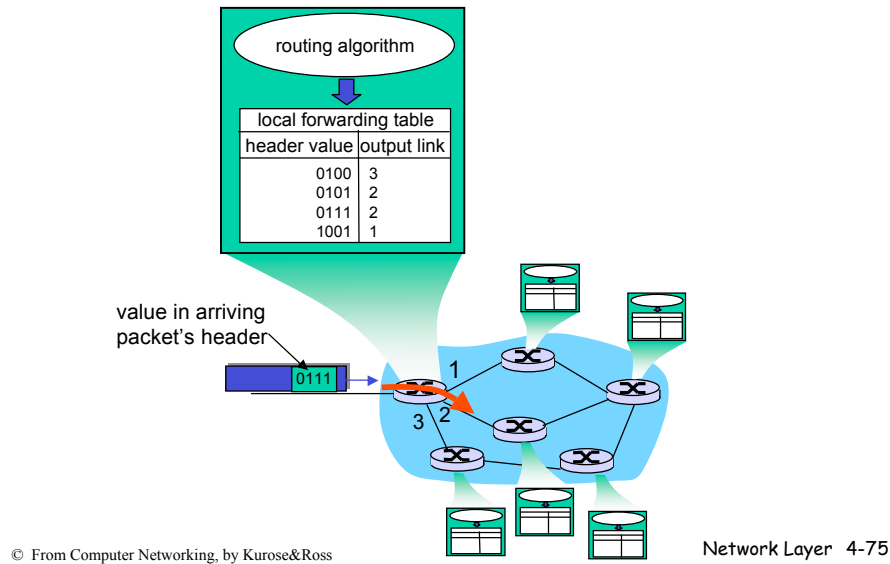
## Tunneling



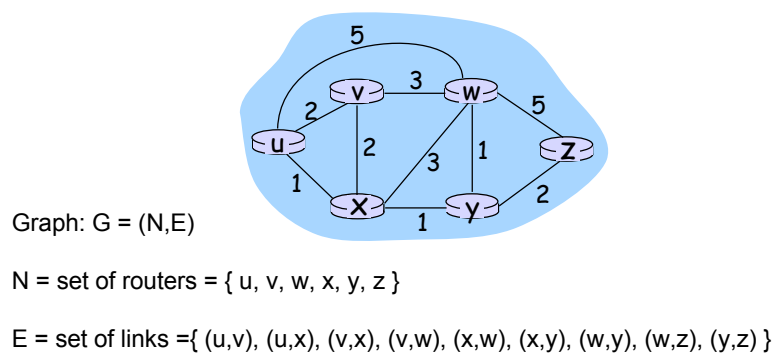
## Chapter 4: Network Layer

- ❑ 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

## Interplay between routing, forwarding



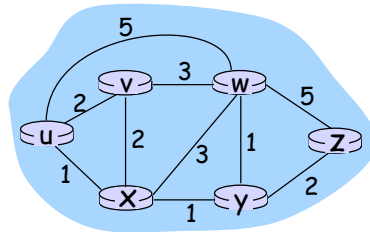
## Graph abstraction



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



•  $c(x,x') = \text{cost of link } (x,x')$

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

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

© From Computer Networking, by Kurose&Ross

Network Layer 4-77

## How to set link costs?

- ❑ To achieve minimum hop routing
  - Set all link costs to 1
  - Will minimize (average) link load (and node processing)
    - See next slide
  - Does not necessarily minimize delay, nor congestion!
- ❑ Changing link costs will surely change the least-cost paths!
- ❑ Link costs can be engineered to optimize the network to some extent
  - But this usually requires to know the traffic matrix (TM)
  - What is a TM?
    - For every pair of nodes  $(i,j)$ ,  $TM(i,j)$  is the amount of traffic entering the network at node  $i$  and exiting the network at node  $j$

Network Layer 4-78

## Minimum hop routing minimises the average link **load** (for any TM)

$$Score = Avg\_link\_load = \frac{\sum_{i \in links} load_i}{N}$$

- Minimizing the average link load is equivalent to minimizing the sum of all the link loads. So, remove denominator N from the score.
- Routing a new flow of rate R along a given path P will increase the score:

$$Score\_increase = \sum_{i \in P} R = R \times nb\_hops(P)$$

- Therefore, minimizing the average link load is equivalent to minimizing the number of hops of each flow
- So to achieve this, each link will simply get the static metric = 1

Network Layer 4-79

## InvCap routing minimises the average link **utilisation** (for any TM)

$$Score = Avg\_link\_util = \frac{\sum_{i \in links} util_i}{N} = \frac{\sum_{i \in links} \frac{load_i}{capacity_i}}{N}$$

- Equivalent to minimizing the sum of all the link **utilisations**
- Routing a new flow of rate R along a given path P will increase the score:

$$Score\_increase = \sum_{i \in P} \frac{R}{C_i} = R \times \sum_{i \in P} \frac{1}{C_i}$$

- Therefore, minimizing the average link utilisation is equivalent to finding the path that minimizes

$$\sum_{i \in P} \frac{1}{C_i}$$

- So, to achieve this, each link will simply get the static metric =  $1/C_i$
- InvCap = metric is the inverse of the capacity

Network Layer 4-80



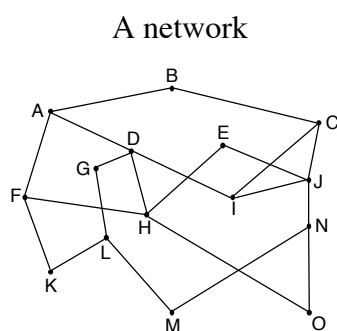
## Other possible metrics

- ❑ Link delay metric
  - Minimizes delay
  - But delay has several components:
    - Propagation delay
    - Transmission delay ( $= \text{packet\_size} / \text{link\_capacity}$ )
    - Queuing delay (variable, depends on the load, difficult to take into account)
- ❑ Administrative link cost (weight)
  - Any link metric computed so as to optimise a given score
  - For example to better balance the load
    - but traffic matrix dependent!
- ❑ Basically any summable quantity
  - Summable = cost of a path is the sum of the costs of the links

Network Layer 4-81

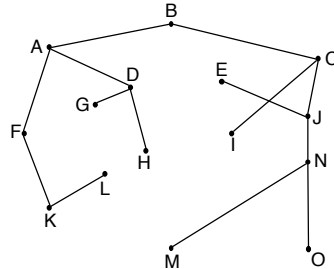
## Optimality principle

- ❑ If a router J is on the optimal path from router I to router K, then the optimal path from J to K also falls along the same route
- ❑ Consequence: the set of optimal routes from all sources to a destination form a tree rooted at the destination
- ❑ Similarly: the set of optimal routes from one source to all destinations form a tree rooted at the source



From Computer Networks, by Tanenbaum © Prentice Hall

A sink tree for router B based on the hop count metric



Network Layer 4-82

## 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

© From Computer Networking, by Kurose&Ross

Network Layer 4-83

## 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

© From Computer Networking, by Kurose&Ross

Network Layer 4-84

## A Link-State Routing Algorithm

### Principle

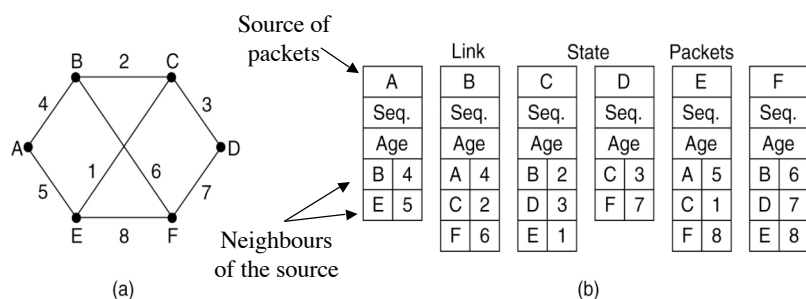
- "link state broadcast"
  - See next slides
  - So, all nodes have the same info
- Every node computes least-cost paths to all other nodes
  - It uses Dijkstra's algorithm (see later)
  - This gives **forwarding table** for that node

© From Computer Networking, by Kurose&Ross

Network Layer 4-85

## Building link state packets

- Link State Packets are composed of
  - the **source node**, a sequence number and an age (see later)
  - a distance vector **limited to the neighbours**

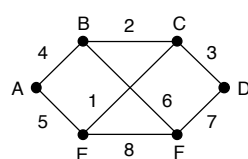


From Computer Networks, by Tanenbaum © Prentice Hall

Network Layer 4-86

## Distributing Link State Packets

- ❑ Packets are flooded selectively
  - Packets are not forwarded on the lines they arrived on
  - Duplicate (or older) packets are detected by the sequence number
- ❑ Packets are acknowledged



Packet received from  
D via C and F

The packet buffer for router B

Source	Seq.	Age	Send flags			ACK flags			Data
			A	C	F	A	C	F	
A	21	60	0	1	1	1	0	0	
F	21	60	1	1	0	0	0	1	
E	21	59	0	1	0	1	0	1	
C	20	60	1	0	1	0	1	0	
D	21	59	1	0	0	0	1	1	

Clearly, routers do not forward the received packets immediately but put them for a short while in the **packet buffer (holding area)**

From Computer Networks, by Tanenbaum © Prentice Hall

Network Layer 4-87

## Potential problems

- ❑ What if the sequence number wraps around?
  - **Solution:** choose 32-bit sequence number
  - Needs 137 years to wrap around if one link state packet is sent every second (in practice one packet per e.g. 10 seconds)
- ❑ What if a router crashes?
  - It restarts with sequence number 0 and its packets are ignored until the sequence number reaches the previous value
  - **Solution:** the age field is decremented by 1 every second and the entry removed when it hits 0
- ❑ What if a sequence number is corrupted?
  - Same consequence, same solution

© From Computer Networking, by Kurose&Ross

Network Layer 4-88

## 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 **forwarding table** for that node
- iterative: after k iterations, know least-cost path to k destinations

### Notation:

- $c(x,y)$ : link cost from node x to y;  $= \infty$  if not direct neighbors
- $D(v)$ : current value of cost of path from source to destination v
- $p(v)$ : predecessor node along path from source to v
- $N'$ : set of nodes whose least-cost path definitively known

© From Computer Networking, by Kurose&Ross

Network Layer 4-89

## Dijkstra's Algorithm

(as executed in node u)

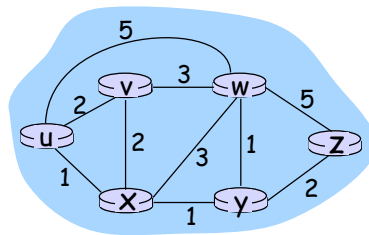
```
1 Initialization:
2   $N' = \{u\}$ 
3  for all nodes v
4    if v adjacent to u
5      then  $D(v) = c(u,v)$ 
6    else  $D(v) = \infty$ 
7
8 Loop
9  find w not in  $N'$  such that  $D(w)$  is a minimum
10 add w to  $N'$ 
11 update  $D(v)$  for all v adjacent to w and not in  $N'$  :
12    $D(v) = \min( D(v), D(w) + c(w,v) )$ 
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  $N'$ 
```

© From Computer Networking, by Kurose&Ross

Network Layer 4-90

## Dijkstra's algorithm: example

Step	N'	D(v),p(v)	D(w),p(w)	D(x),p(x)	D(y),p(y)	D(z),p(z)
0	u	2,u	5,u	1,u	$\infty$	$\infty$
1	ux	2,u	4,x		2,x	$\infty$
2	uxy	2,u	3,y			4,y
3	uxyv		3,y			4,y
4	uxyvw					4,y
5	uxyvwz					

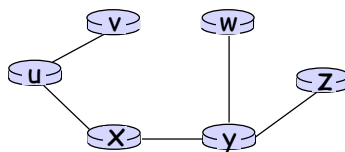


© From Computer Networking, by Kurose&Ross

Network Layer 4-91

## 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)

© From Computer Networking, by Kurose&Ross

Network Layer 4-92

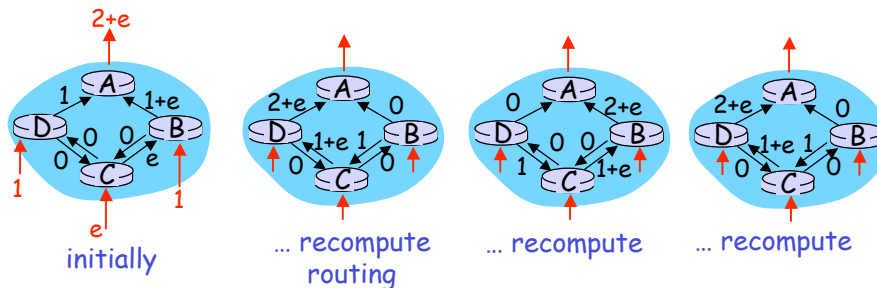
## Dijkstra's algorithm, discussion

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:  $O(n \log n)$

Oscillations possible when link cost are traffic dependent:

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



© From Computer Networking, by Kurose&Ross

Network Layer 4-93

## 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

© From Computer Networking, by Kurose&Ross

Network Layer 4-94

## Distance Vector Algorithm

### Bellman-Ford Equation (dynamic programming)

Define

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

Then

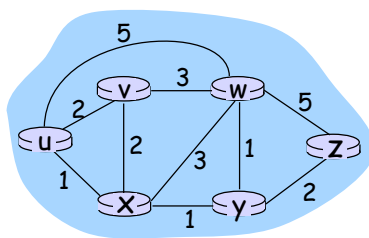
$$d_x(y) = \min_v \{c(x,v) + d_v(y)\}$$

where min is taken over all neighbors  $v$  of  $x$

© From Computer Networking, by Kurose&Ross

Network Layer 4-95

## 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

© From Computer Networking, by Kurose&Ross

Network Layer 4-96



## 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  
 $D_v = [D_v(y): y \in N]$

© From Computer Networking, by Kurose&Ross

Network Layer 4-97

## Distance vector algorithm (4)

### Basic idea:

- Each node periodically 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  $D_x(y)$  converges to the actual least cost  $d_x(y)$

© From Computer Networking, by Kurose&Ross

Network Layer 4-98

## Distance Vector Algorithm (5)

### Iterative, asynchronous:

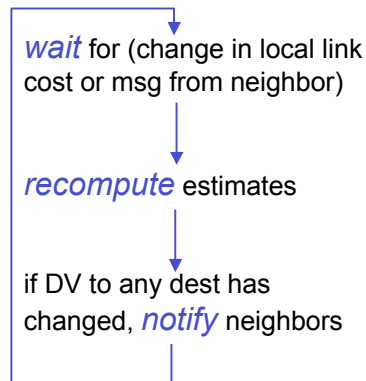
each local iteration caused by:

- local link cost change
- DV update message from neighbor

### Distributed:

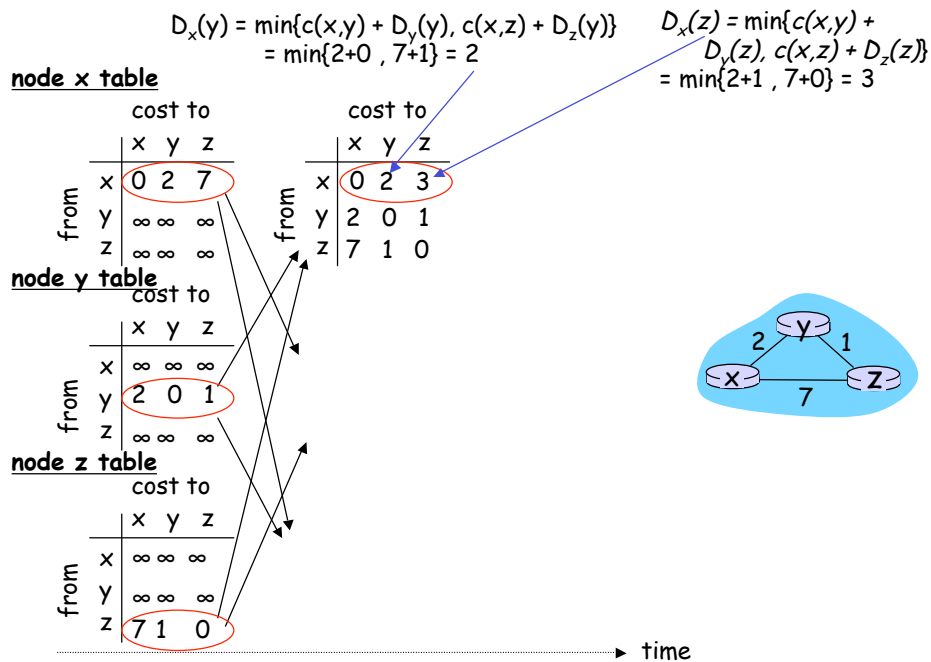
- each node notifies neighbors *only* when its DV changes
  - neighbors then notify their neighbors if necessary

### Each node:



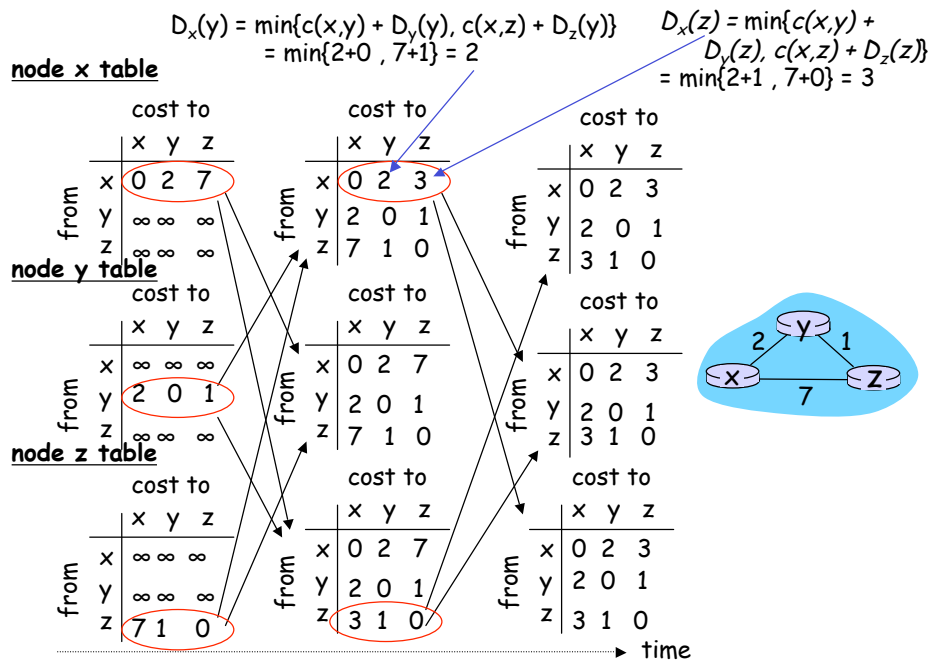
© From Computer Networking, by Kurose&Ross

Network Layer 4-99



© From Computer Networking, by Kurose&Ross

Network Layer 4-100



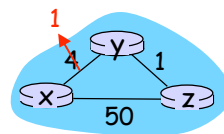
© From Computer Networking, by Kurose&Ross

Network Layer 4-101

## 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



"good news travels fast"

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.

© From Computer Networking, by Kurose&Ross

Network Layer 4-102

## Distance Vector: link cost changes

### Link cost changes:

- good news travels fast
- bad news travels slowly - "count to infinity" problem!

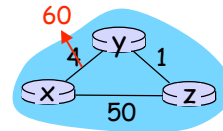
Distance to A  
in nr. of hops

Propagation of bad news  
(link A-B is down,  $c(B,A) = \infty$ )

A	B	C	D	E	
•	•	•	•	•	Initially
	1	2	3	4	Initially
	3	2	3	4	After 1 exchange
	3	4	3	4	After 2 exchanges
	5	4	5	4	After 3 exchanges
	5	6	5	6	After 4 exchanges
	7	6	7	6	After 5 exchanges
	7	8	7	8	After 6 exchanges
	⋮				
	∞	∞	∞	∞	

From Computer Networks, by Tanenbaum © Prentice Hall

- Other example (see book):
  - 44 iterations before algorithm stabilizes



Network Layer 4-103

© From Computer Networking, by Kurose&Ross

## Poisoned reverse

- Also called "split-horizon"
- If C routes through B to get to A :
  - C tells B its (C's) distance to A is infinite
    - (so B won't route to A via C)
- With poisoned reverse, we get (when link A-B goes down):

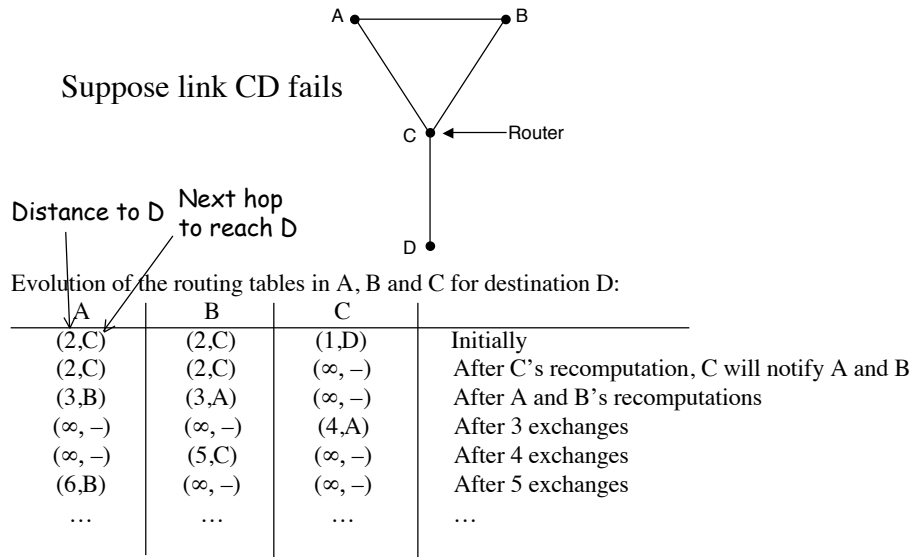
A	B	C	D	E	
•	•	•	•	•	Initially
	1	2	3	4	Initially
	∞	2	3	4	After 1 exchange
	∞	∞	3	4	After 2 exchanges
	∞	∞	∞	4	After 3 exchanges
	∞	∞	∞	∞	After 4 exchanges

- Q: Will this completely solve count to infinity problem?

© From Computer Networking, by Kurose&Ross

Network Layer 4-104

## Poisoned reverse is not a panacea



From Computer Networks, by Tanenbaum © Prentice Hall

Network Layer 4-105

## Comparison of LS and DV algorithms

### Message complexity

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

### Speed of Convergence

- **LS:**  $O(n \log n)$  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 propagates through network

© From Computer Networking, by Kurose&Ross

Network Layer 4-106

## 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

© From Computer Networking, by Kurose&Ross

Network Layer 4-107

## Hierarchical Routing

Our routing study thus far - idealization

- ❑ all routers identical
  - ❑ 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

© From Computer Networking, by Kurose&Ross

Network Layer 4-108

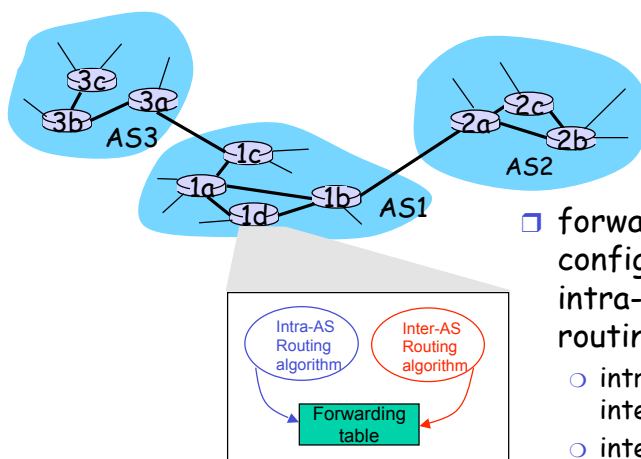
## Hierarchical Routing

- 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 router
- Direct link to router in another AS

© From Computer Networking, by Kurose&Ross

Network Layer 4-109

## Interconnected ASes



- forwarding table configured by both intra- and inter-AS routing algorithm
  - intra-AS sets entries for internal destinations
  - inter-AS & intra-AS set entries for external destinations

© From Computer Networking, by Kurose&Ross

Network Layer 4-110

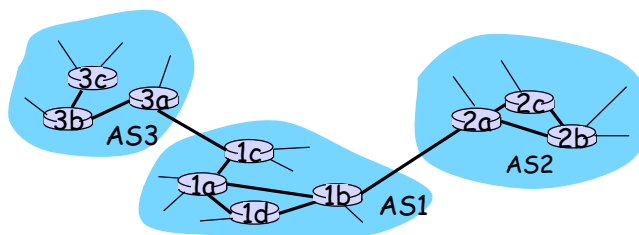
## Inter-AS tasks

- suppose router in AS1 receives datagram whose dest is outside of AS1
  - router should forward packet to gateway router, but which one?

### AS1 must:

1. learn which dests reachable through AS2, which through AS3
2. propagate this reachability info to all routers in AS1

Job of inter-AS routing!

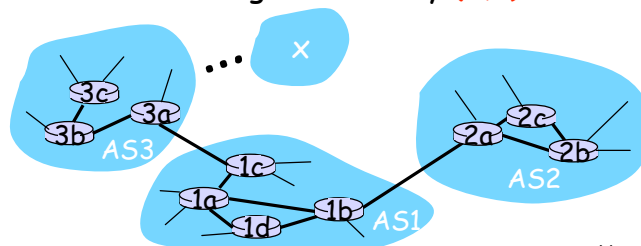


© From Computer Networking, by Kurose&Ross

Network Layer 4-111

## Example: Setting forwarding table in router 1d

- suppose AS1 learns (via inter-AS protocol) that subnet x reachable via AS3 (gateway 1c) but not via AS2.
- inter-AS protocol propagates reachability info to all internal routers.
- router 1d determines from intra-AS routing info that its interface I is on the least cost path to 1c.
  - installs forwarding table entry (x, I)



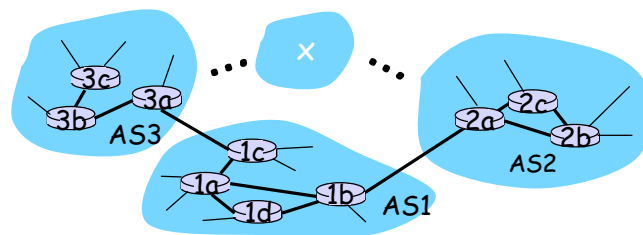
© From Computer Networking, by Kurose&Ross

Network Layer 4-112



## Example: Choosing among multiple ASes

- ❑ now suppose AS1 learns from inter-AS protocol that subnet **x** is reachable from AS3 and from AS2.
- ❑ to configure forwarding table, router 1d must determine towards which gateway it should forward packets for dest **x**.
  - this is also job of inter-AS routing protocol!

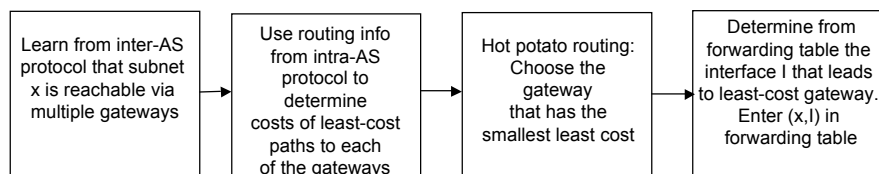


© From Computer Networking, by Kurose&Ross

Network Layer 4-113

## Example: Choosing among multiple ASes

- ❑ now suppose AS1 learns from inter-AS protocol that subnet **x** is reachable from AS3 and from AS2.
- ❑ to configure forwarding table, router 1d must determine towards which gateway it should forward packets for dest **x**.
  - this is also job of inter-AS routing protocol!
- ❑ **hot potato routing**: send packet towards closest of two routers.



© From Computer Networking, by Kurose&Ross

Network Layer 4-114

## 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

© From Computer Networking, by Kurose&Ross

Network Layer 4-115

## 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
    - recommended by IETF
  - IS-IS: Intermediate System to Intermediate System
    - standardized by ISO
  - IGRP: Interior Gateway Routing Protocol
    - Cisco proprietary

© From Computer Networking, by Kurose&Ross

Network Layer 4-116

## Chapter 4: Network Layer

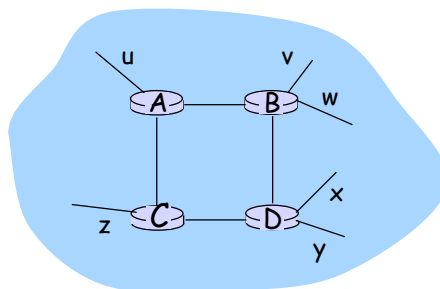
- 4.1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
- 4.4 IP: Internet Protocol
  - Datagram format
  - IPv4 addressing
  - ICMP
  - IPv6
- 4.5 Routing algorithms
  - Link state
  - Distance Vector
  - Hierarchical routing
- 4.6 Routing in the Internet
  - RIP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

© From Computer Networking, by Kurose&Ross

Network Layer 4-117

## 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:

destination	hops
u	1
v	2
w	2
x	3
y	3
z	2

© From Computer Networking, by Kurose&Ross

Network Layer 4-118

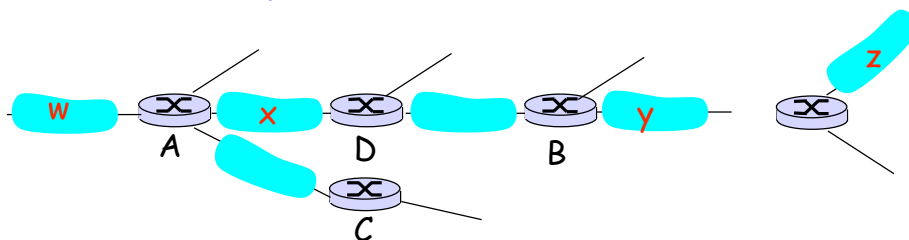
## RIP advertisements

- *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

© From Computer Networking, by Kurose&Ross

Network Layer 4-119

## RIP: Example



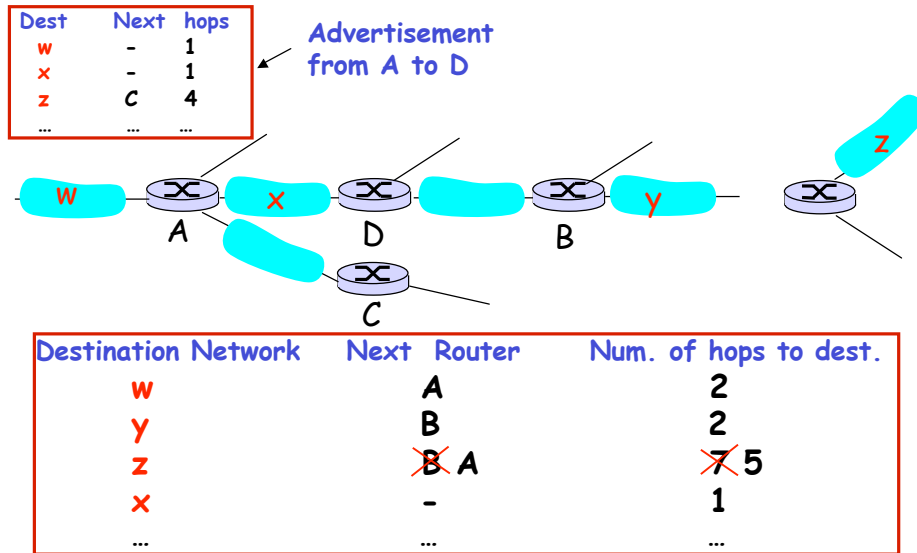
Destination Network	Next Router	Num. of hops to dest.
<b>w</b>	<b>A</b>	<b>2</b>
<b>y</b>	<b>B</b>	<b>2</b>
<b>z</b>	<b>B</b>	<b>7</b>
<b>x</b>	<b>-</b>	<b>1</b>
...	...	...

Routing table in D

© From Computer Networking, by Kurose&Ross

Network Layer 4-120

## RIP: Example



© From Computer Networking, by Kurose&Ross

Routing table in D

Network Layer 4-121

## RIP: Link Failure and Recovery

If no advertisement heard after 180 sec -->  
neighbor/link declared dead

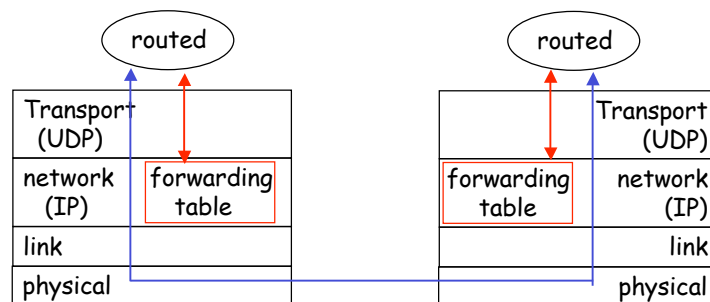
- routes via neighbor invalidated
- new advertisements sent to neighbors
- neighbors in turn send out new advertisements (if tables changed)
- link failure info quickly (?) propagates to entire net
- *poison reverse* used to prevent ping-pong loops (infinite distance = 16 hops)

© From Computer Networking, by Kurose&Ross

Network Layer 4-122

## RIP Table processing

- ❑ RIP routing tables managed by **application-level** process called route-d (daemon)
- ❑ advertisements sent in UDP packets, periodically repeated



© From Computer Networking, by Kurose&Ross

Network Layer 4-123

## 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

© From Computer Networking, by Kurose&Ross

Network Layer 4-124

## 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)

© From Computer Networking, by Kurose&Ross

Network Layer 4-125

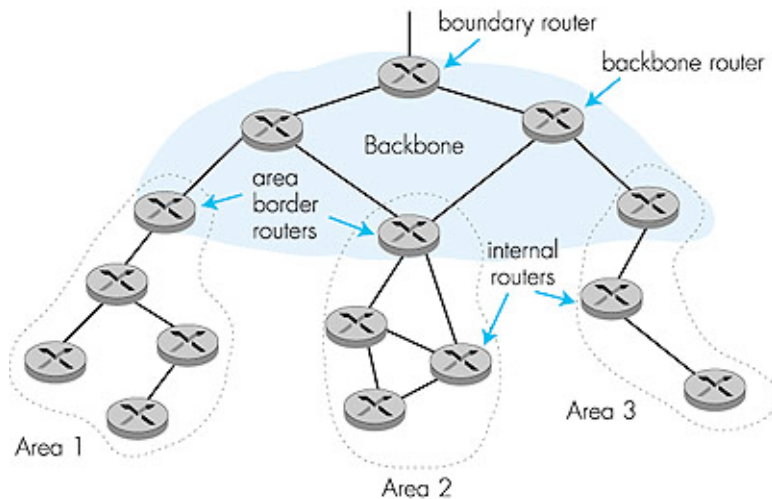
## OSPF "advanced" features (not in RIP)

- ❑ **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.

© From Computer Networking, by Kurose&Ross

Network Layer 4-126

## Hierarchical OSPF



© From Computer Networking, by Kurose&Ross

Network Layer 4-127

## Hierarchical OSPF

- ❑ **two-level hierarchy:** local area, backbone.
  - Link-state advertisements only in area
  - each node has detailed area topology; only knows direction (shortest path) to nets in other areas
- ❑ **area border routers:** "summarize" distances to nets in own area, advertise to other Area Border routers
- ❑ **backbone routers:** run OSPF routing limited to backbone
- ❑ **boundary routers:** connect to other AS's

© From Computer Networking, by Kurose&Ross

Network Layer 4-128



## 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

© From Computer Networking, by Kurose&Ross

Network Layer 4-129

## Internet inter-AS routing: BGP

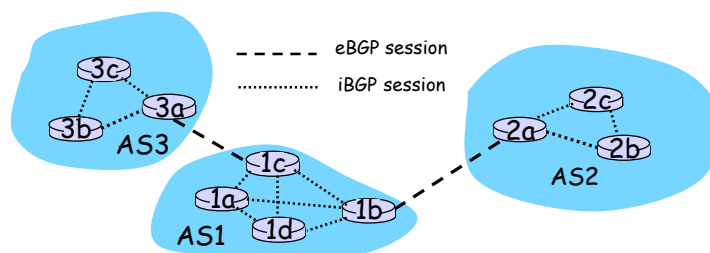
- BGP (Border Gateway Protocol): *the de facto standard*
- BGP provides each AS a means to:
  1. Obtain subnet reachability information from neighboring ASs
  2. Propagate reachability information to all AS-internal routers
  3. Determine "good" routes to subnets based on reachability information and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

© From Computer Networking, by Kurose&Ross

Network Layer 4-130

## 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 prefix to AS1:
  - AS2 **promises** it will forward any datagram addressed towards that prefix
  - AS2 can aggregate prefixes in its advertisement

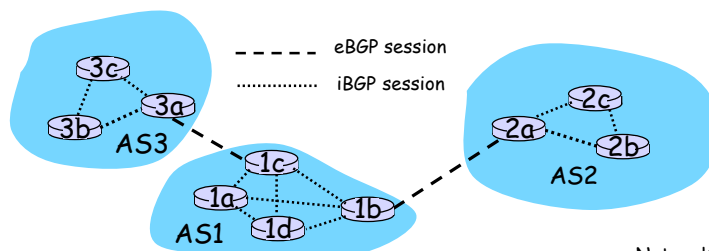


© From Computer Networking, by Kurose&Ross

Network Layer 4-131

## Distributing reachability info

- ❑ using eBGP session between 3a and 1c, AS3 sends prefix reachability info to AS1
  - 1c can then use iBGP to distribute new prefix info to all routers in AS1
  - 1b can then re-advertise new reachability info to AS2 over 1b-to-2a eBGP session
- ❑ when router learns of new prefix, it creates entry for this prefix in its forwarding table



© From Computer Networking, by Kurose&Ross

Network Layer 4-132

## Path attributes & BGP routes

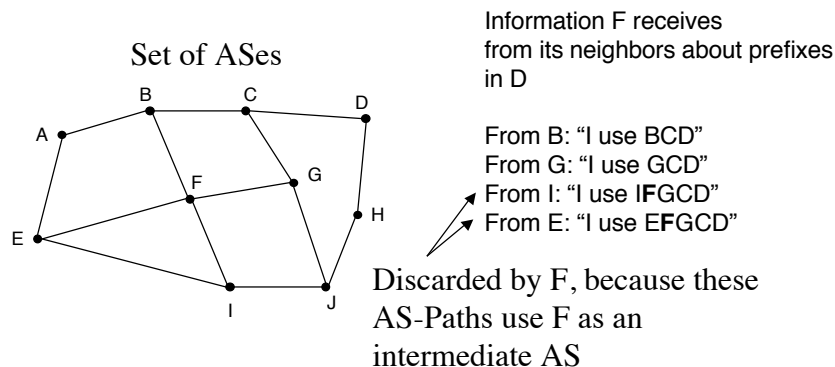
- ❑ advertised prefix includes BGP attributes
  - prefix + attributes = "route"
- ❑ two important attributes:
  - **AS-PATH**: contains ASs through which prefix advertisement has passed: e.g, AS 67, AS 17
    - BGP is a path vector protocol (AS-Path is propagated)
    - In a DV protocol, only the distance is propagated
  - **NEXT-HOP**: indicates specific internal-AS router to next-hop AS (may be multiple links from current AS to next-hop-AS)
- ❑ when gateway router receives route advertisement, uses **import policy** to accept/decline

© From Computer Networking, by Kurose&Ross

Network Layer 4-133

## AS-Paths avoid loops

- ❑ If an AS sees itself in the AS-Path advertised by a neighbor AS, it discards it, otherwise it would create a loop
- ❑ More powerful than poisoned reverse (in DV)
  - Made possible by the presence of the AS-Path in the advertisements



From Computer Networks, by Tanenbaum © Prentice Hall

Network Layer 4-134

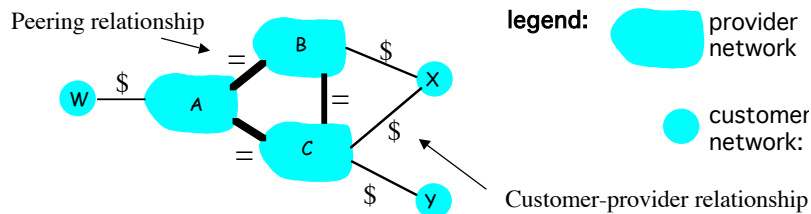
## BGP route selection

- ❑ router may learn about more than 1 route to some prefix. Router must select route.
- ❑ elimination rules:
  1. local preference value attribute: policy decision
  2. shortest AS-PATH
  3. closest NEXT-HOP router: hot potato routing
  4. additional criteria

## BGP messages

- ❑ BGP messages exchanged using TCP
- ❑ BGP messages:
  - **OPEN**: opens TCP connection to peer and authenticates sender
  - **UPDATE**: advertises new path (or withdraws old)
  - **KEEPALIVE** keeps connection alive in absence of UPDATES; also ACKs OPEN request
  - **NOTIFICATION**: reports errors in previous msg; also used to close connection

## BGP routing policy

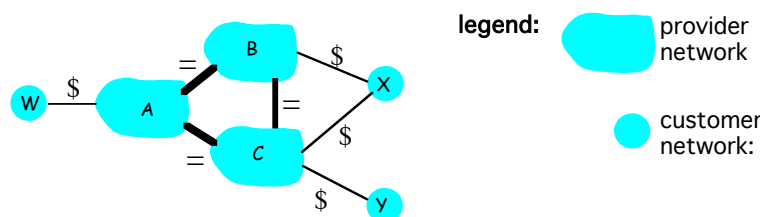


- A,B,C are **provider networks**
- X,W,Y are customers (of provider networks)
- X is **dual-homed**: attached to two networks
  - X does not want to route from B via X to C
  - ... so X will not advertise to B a route to C

© From Computer Networking, by Kurose&Ross

Network Layer 4-137

## BGP routing policy (2)



- A advertises to B the path AW
- B advertises to X the path BAW
- Should B advertise to C the path BAW?
  - No way! B gets no "revenue" for routing CBAW since neither W nor C are B's customers
  - B wants to force C to route to W via A
  - B wants to route **only** to/from its customers!

© From Computer Networking, by Kurose&Ross

Network Layer 4-138

## Why different Intra- and Inter-AS routing ?

### Policy:

- ❑ Inter-AS: admin wants control over how its traffic is routed, who routes through its net
- ❑ Intra-AS: single admin, so no policy decisions needed

### Scale:

- ❑ hierarchical routing saves table size, reduced update traffic

### Performance:

- ❑ Intra-AS: can focus on performance
- ❑ Inter-AS: policy may dominate over performance

## Chapter 4: Network Layer

- ❑ Virtual circuit and datagram networks
- ❑ What's inside a router
- ❑ IP: Internet Protocol
  - Datagram format
    - Fragmentation
  - IPv4 addressing
    - Subnets, CIDR
    - DHCP
    - NAT
  - ICMP
  - IPv6
- ❑ Routing algorithms
  - Intra-domain
    - LS - Link State (OSPF)
      - Dijkstra
    - DV - Distance Vector (RIP)
      - Bellman-Ford
  - Inter-domain routing
    - Hierarchical routing
    - Path Vector
    - Policies
    - Hot-potato
    - BGP