

Week 4 Network Layer



These slides are modified from the slides made available by Kurose and Ross.

Computer Networking: A Top Down Approach Featuring the Internet, 2nd edition.
Jim Kurose, Keith Ross
Addison-Wesley, July 2002.

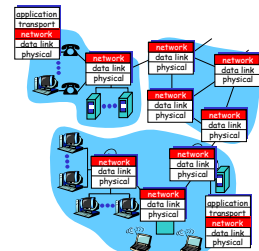
Network Layer 4-1

Network layer functions

- transport packet from sending to receiving hosts
- network layer protocols in every host, router

three important functions:

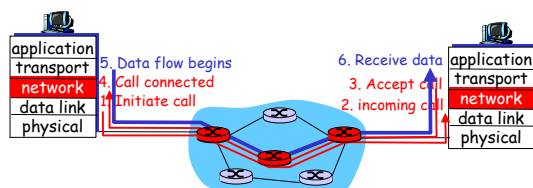
- path determination:** route taken by packets from source to dest. *Routing algorithms*
- forwarding:** move packets from router's input to appropriate router output
- call setup:** some network architectures require router call setup along path before data flows



Network Layer 4-2

Virtual circuits: signaling protocols

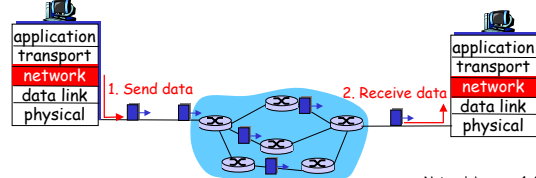
- used to setup, maintain teardown VC
- used in ATM, frame-relay, X.25
- not used in today's Internet**



Network Layer 4-3

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



Network Layer 4-4

Datagram or VC network: why?

Internet

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

ATM

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

Network Layer 4-5

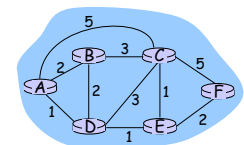
Routing

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 definitions possible

Network Layer 4-6

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

Network Layer 4-7

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 v
- N : set of nodes whose least cost path definitively known

Network Layer 4-8

Dijkstra's Algorithm

1 Initialization:

- $N = \{A\}$
- for all nodes v
- if v adjacent to A
- then $D(v) = c(A,v)$
- else $D(v) = \text{infinity}$

8 Loop

- find w not in N such that $D(w)$ is a minimum
- add w to N
- update $D(v)$ for all v adjacent to w and not in N:
 $D(v) = \min(D(v), D(w) + c(w,v))$
- /* new cost to v is either old cost to v or known shortest path cost to w plus cost from w to v */
- until all nodes in N

Network Layer 4-9

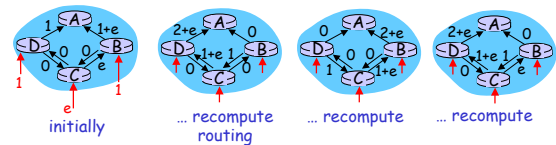
Dijkstra's algorithm, discussion

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:

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



Network Layer 4-10

Distance Table Routing Algorithm

iterative:

- continues until no nodes exchange info.
- self-terminating: no "signal" to stop

asynchronous:

- nodes need not exchange info/iterate in lock step!

distributed:

- each node communicates only with directly-attached neighbors

Distance Table data structure

- each node has its own
- row for each possible destination
- column for each directly-attached neighbor to node
- example: in node X, for dest. Y via neighbor Z:

$$D^X(Y,Z) = \text{distance from X to Y, via Z as next hop} \\ = c(X,Z) + \min_w \{D^Z(Y,w)\}$$

Network Layer 4-11

Distance Vector Routing: overview

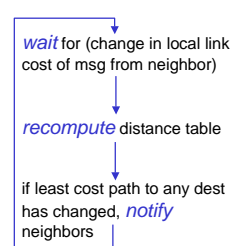
Iterative, asynchronous:

- each local iteration caused by:
 - local link cost change
 - message from neighbor: its least cost path change from neighbor

Distributed:

- each node notifies neighbors only when its least cost path to any destination changes
 - neighbors then notify their neighbors if necessary

Each node:



Network Layer 4-12

Distance Vector Algorithm:

At all nodes, X:

- 1 Initialization:
- 2 for all adjacent nodes v:
- 3 $D^X(*,v) = \text{infinity}$ /* the * operator means "for all rows" */
- 4 $D^X(v,v) = c(X,v)$
- 5 for all destinations, y
- 6 send $\min_w D^X(y,w)$ to each neighbor /* w over all X's neighbors */

Network Layer 4-13

Distance Vector Algorithm (cont.):

```

8 loop
9   wait (until I see a link cost change to neighbor V
10    or until I receive update from neighbor V)
11
12   if (c(X,V) changes by d)
13     /* change cost to all dest's via neighbor v by d */
14     /* note: d could be positive or negative */
15     for all destinations y:  $D^X(y,V) = D^X(y,V) + d$ 
16
17   else if (update received from V wrt destination Y)
18     /* shortest path from V to some Y has changed */
19     /* V has sent a new value for its  $\min_w DV(Y,w)$  */
20     /* call this received new value is "newval" */
21     for the single destination y:  $D^X(Y,V) = c(X,V) + \text{newval}$ 
22
23   if we have a new  $\min_w D^X(Y,w)$  for any destination Y
24     send new value of  $\min_w D^X(Y,w)$  to all neighbors
25
26 forever
    
```

Network Layer 4-14

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

Network Layer 4-15

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

Network Layer 4-16

Hierarchical Routing

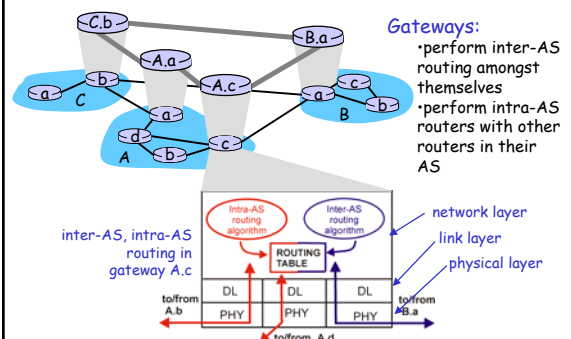
- aggregate routers into regions, "autonomous systems" (AS)
- routers in same AS run same routing protocol
 - "intra-AS" routing protocol
 - routers in different AS can run different intra-AS routing protocol

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

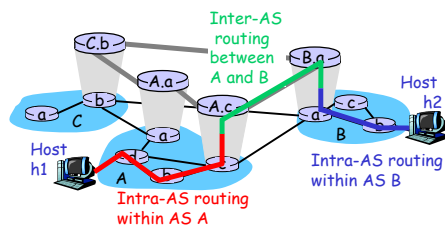
Network Layer 4-17

Intra-AS and Inter-AS routing



Network Layer 4-18

Intra-AS and Inter-AS routing

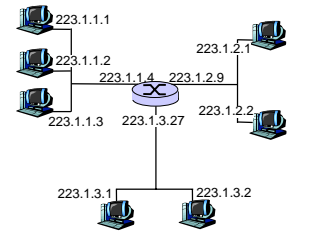


- We'll examine specific inter-AS and intra-AS Internet routing protocols shortly

Network Layer 4-19

IP Addressing: introduction

- IP address: 32-bit identifier for host, router *interface*
- *interface*: connection between host/router and physical link
 - router's typically have multiple interfaces
 - host may have multiple interfaces
 - IP addresses associated with each interface

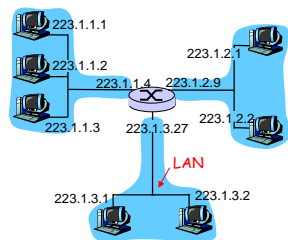


223.1.1.1 = 11011111 00000001 00000001 00000001
223 1 1 1

Network Layer 4-20

IP Addressing

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



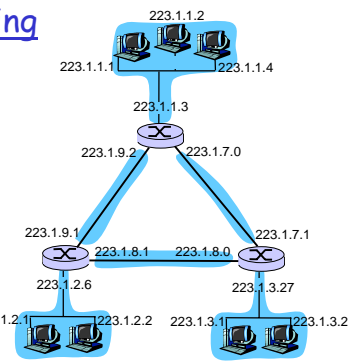
network consisting of 3 IP networks
(for IP addresses starting with 223,
first 24 bits are network address)

Network Layer 4-21

IP Addressing

- How to find the networks?
- Detach each interface from router, host
- create "islands of isolated networks"

Interconnected
system consisting
of six networks



Network Layer 4-22

IP Addresses

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

"class-full" addressing:

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

← 32 bits →

Network Layer 4-23

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

← network part → host part →
11001000 00010111 00010000 00000000
200.23.16.0/23

Network Layer 4-24

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

Network Layer 4-25

IP addresses: how to get one?

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

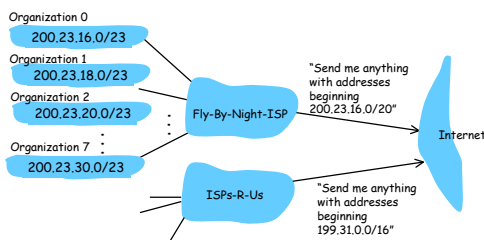
A: gets allocated portion of its provider ISP's address space

ISP's block	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

Network Layer 4-26

Hierarchical addressing: route aggregation

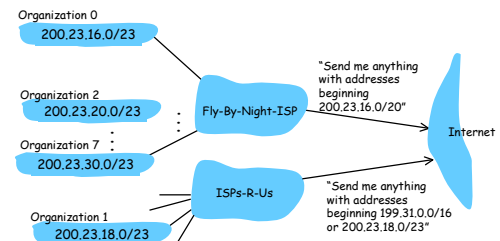
Hierarchical addressing allows efficient advertisement of routing information:



Network Layer 4-27

Hierarchical addressing: more specific routes

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



Network Layer 4-28

IP addressing: the last word...

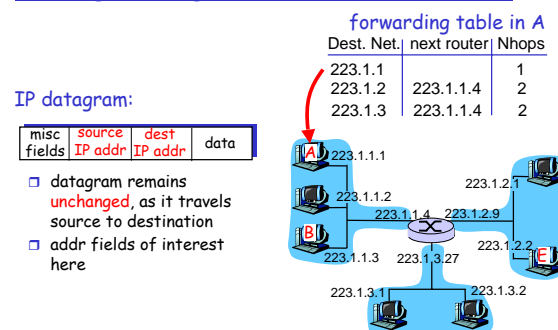
Q: How does an ISP get block of addresses?

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

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

Network Layer 4-29

Getting a datagram from source to dest.



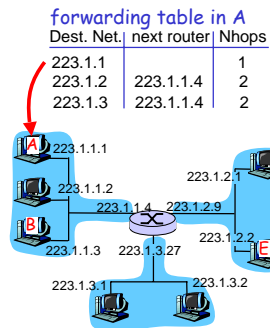
Network Layer 4-30

Getting a datagram from source to dest.

misc	223.1.1.1	223.1.1.3	data
------	-----------	-----------	------

Starting at A, send IP datagram addressed to B:

- look up net. address of B in forwarding table
- find B is on same net. as A
- link layer will send datagram directly to B inside link-layer frame
 - B and A are directly connected



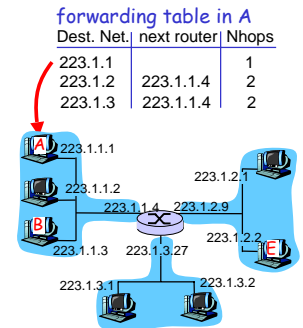
Network Layer 4-31

Getting a datagram from source to dest.

misc	223.1.1.1	223.1.2.3	data
------	-----------	-----------	------

Starting at A, dest. E:

- look up network address of E in forwarding table
 - A, E not directly attached
- routing table: next hop router to E is 223.1.1.4
- link layer sends datagram to router 223.1.1.4 inside link-layer frame
 - datagram arrives at 223.1.1.4
 - continued.....



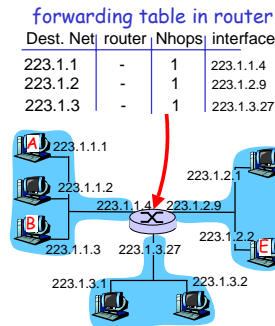
Network Layer 4-32

Getting a datagram from source to dest.

misc	223.1.1.1	223.1.2.3	data
------	-----------	-----------	------

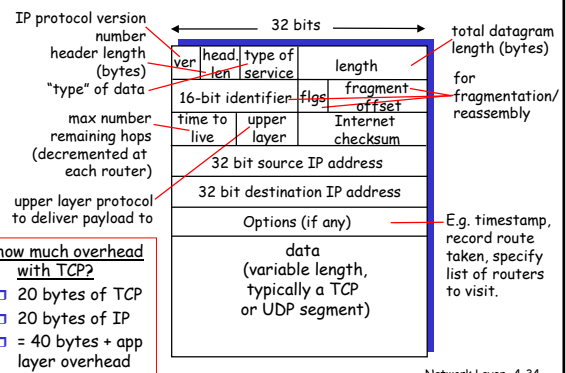
Arriving at 223.1.4, destined for 223.1.2.2

- look up network address of E in router's forwarding table
 - E on same network as router's interface 223.1.2.9
 - router, E directly attached
- link layer sends datagram to 223.1.2.2 inside link-layer frame via interface 223.1.2.9
- datagram arrives at 223.1.2.2!!! (hooray!)



Network Layer 4-33

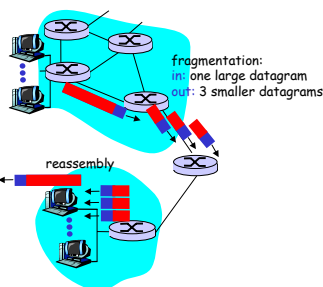
IP datagram format



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

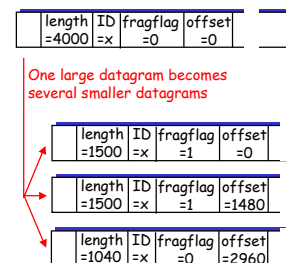


Network Layer 4-35

IP Fragmentation and Reassembly

Example

- 4000 byte datagram
- MTU = 1500 bytes



Network Layer 4-36

DHCP: Dynamic Host Configuration Protocol

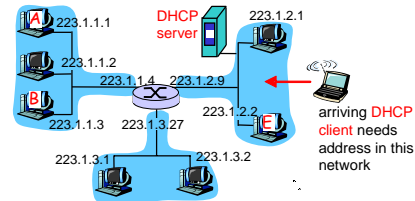
Goal: allow host to *dynamically* obtain its IP address from network server when it joins network
 Can renew its lease on address in use
 Allows reuse of addresses (only hold address while connected an "on")
 Support for mobile users who want to join network (more shortly)

DHCP overview:

- 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

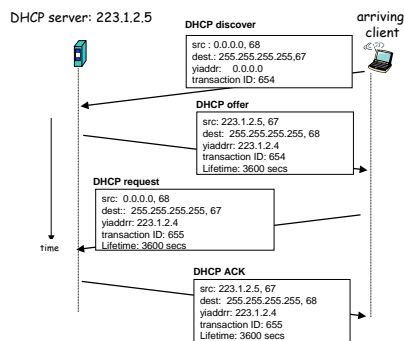
Network Layer 4-37

DHCP client-server scenario



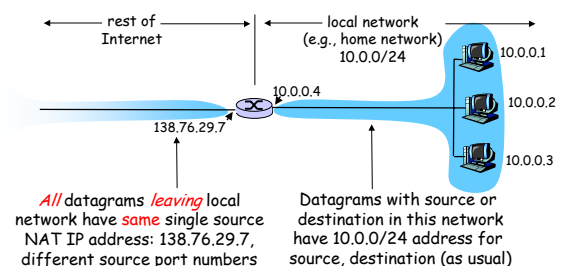
Network Layer 4-38

DHCP client-server scenario



Layer 4-39

NAT: Network Address Translation



Network Layer 4-40

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).

Network Layer 4-41

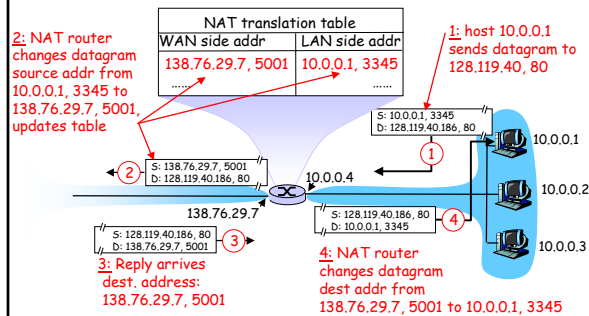
NAT: Network Address Translation

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

Network Layer 4-42

NAT: Network Address Translation



Network Layer 4-43

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

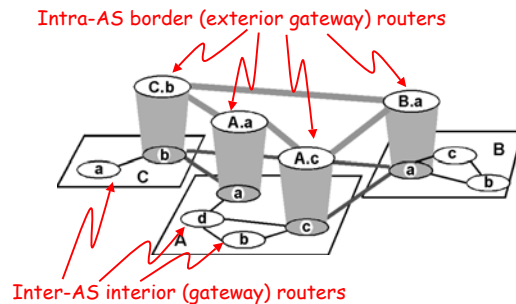
Network Layer 4-44

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

Network Layer 4-45

Internet AS Hierarchy



Network Layer 4-46

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)

Network Layer 4-47

RIP (Routing Information Protocol)

- Distance vector algorithm
- Included in BSD-UNIX Distribution in 1982
- Distance metric: # of hops (max = 15 hops)
- Distance vectors: exchanged among neighbors every 30 sec via Response Message (also called **advertisement**)
- Each advertisement: list of up to 25 destination nets within AS

Network Layer 4-48

RIP: Link Failure and Recovery

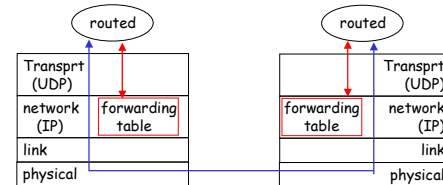
If no advertisement heard after 180 sec --> neighbor/link declared dead

- 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)

Network Layer 4-49

RIP Table processing

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



Network Layer 4-50

RIP Table example (continued)

Router: giroflee.eurocom.fr

Destination	Gateway	Flags	Ref	Use	Interface
127.0.0.1	127.0.0.1	UH	0	26492	lo0
192.168.2.	192.168.2.5	U	2	13	fa0
193.55.114.	193.55.114.6	U	3	58503	le0
192.168.3.	192.168.3.5	U	2	25	qaa0
224.0.0.0	193.55.114.6	U	3	0	le0
default	193.55.114.129	UG	0	143454	

- Three attached class C networks (LANs)
- Router only knows routes to attached LANs
- Default router used to "go up"
- Route multicast address: 224.0.0.0
- Loopback interface (for debugging)

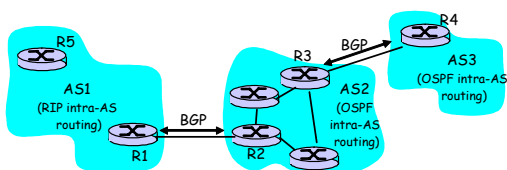
Network Layer 4-51

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)

Network Layer 4-52

Inter-AS routing in the Internet: BGP



Network Layer 4-53

Internet inter-AS routing: BGP

- **BGP (Border Gateway Protocol)**: the de facto standard
- **Path Vector** protocol:
 - similar to Distance Vector protocol
 - each Border Gateway broadcast to neighbors (peers) *entire path* (i.e., sequence of AS's) to destination
 - BGP routes to networks (ASs), not individual hosts
 - E.g., Gateway X may send its path to dest. Z:

Path (X,Z) = X,Y1,Y2,Y3,...,Z

Network Layer 4-54

Internet inter-AS routing: BGP

- Suppose:** gateway X send its path to peer gateway W
- W may or may not select path offered by X
 - cost, policy (don't route via competitors AS), loop prevention reasons.
 - If W selects path advertised by X, then:
 - Path (W,Z) = w, Path (X,Z)
 - Note: X can control incoming traffic by controlling its route advertisements to peers:
 - e.g., don't want to route traffic to Z → don't advertise any routes to Z

Network Layer 4-55

Why different Intra- and Inter-AS routing ?

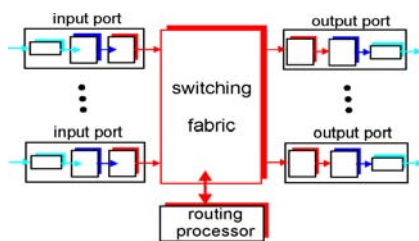
- Policy:**
- Inter-AS: admin wants control over how its traffic routed, who routes through its net.
 - Intra-AS: single admin, so no policy decisions needed
- Scale:**
- hierarchical routing saves table size, reduced update traffic
- Performance:**
- Intra-AS: can focus on performance
 - Inter-AS: policy may dominate over performance

Network Layer 4-56

Router Architecture Overview

Two key router functions:

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



Network Layer 4-57

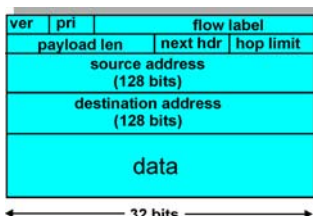
IPv6

- **Initial motivation:** 32-bit address space completely allocated by 2008.
- **Additional motivation:**
 - header format helps speed processing/forwarding
 - header changes to facilitate QoS
 - new "anycast" address: route to "best" of several replicated servers
- **IPv6 datagram format:**
 - fixed-length 40 byte header
 - no fragmentation allowed

Network Layer 4-58

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



Network Layer 4-59

Other Changes from IPv4

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

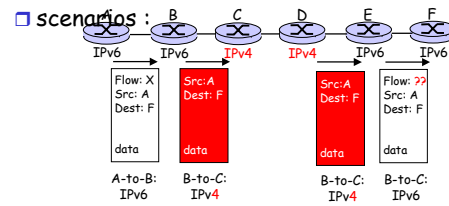
Network Layer 4-60

Transition From IPv4 To IPv6

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

Network Layer 4-61

Dual Stack Approach



Network Layer 4-62