

Routing

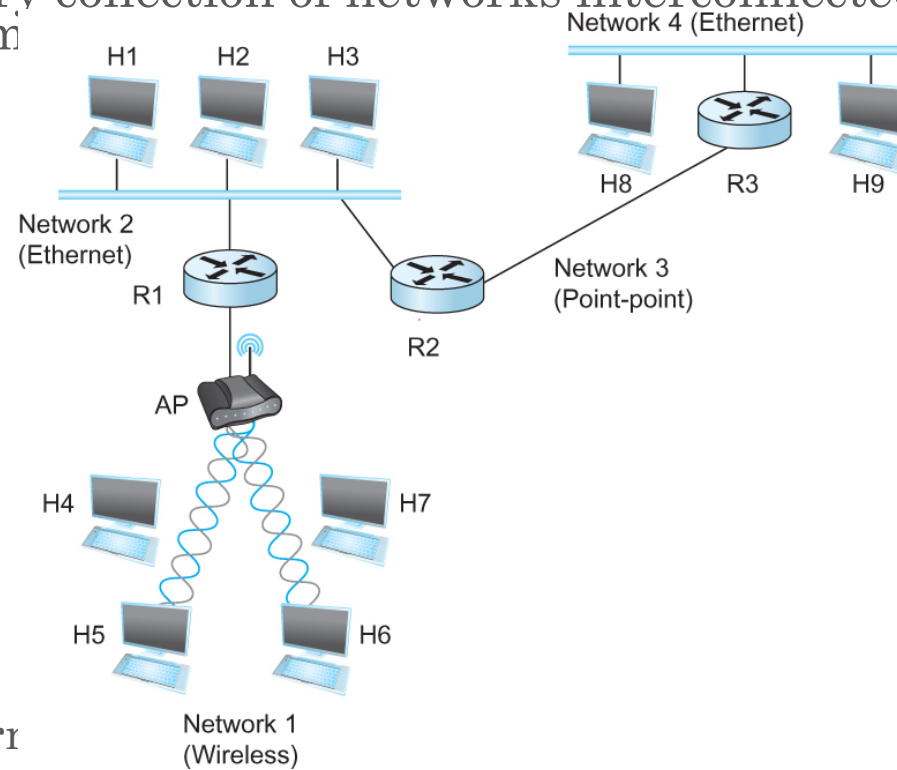
Outline

Algorithms

Scalability

Internetworking

- What is internetworking
 - An arbitrary collection of networks interconnected to provide some service

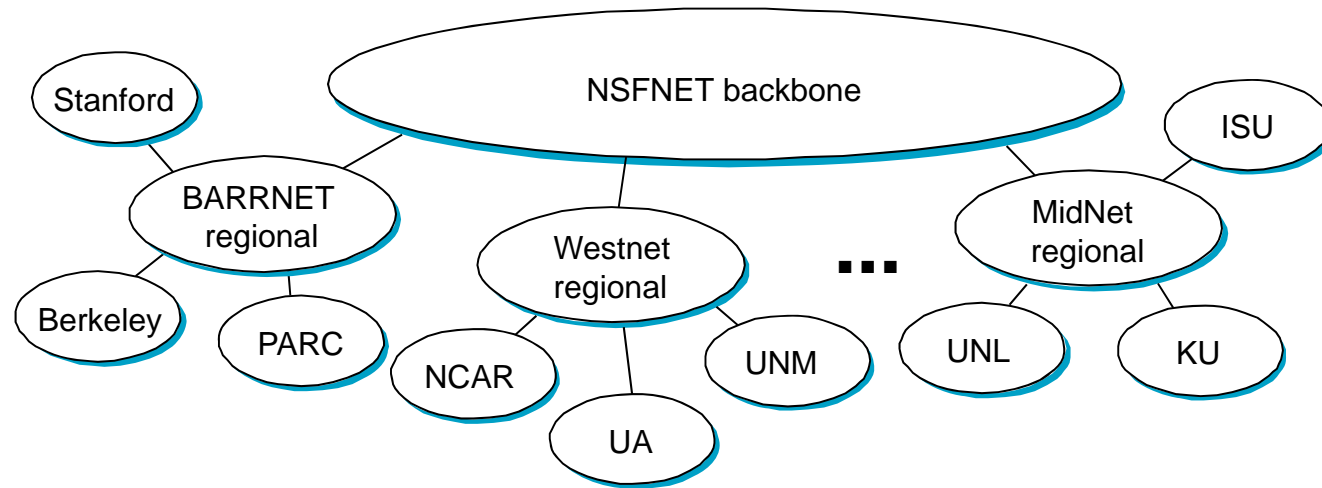


A simple internetworking setup with routers

represents

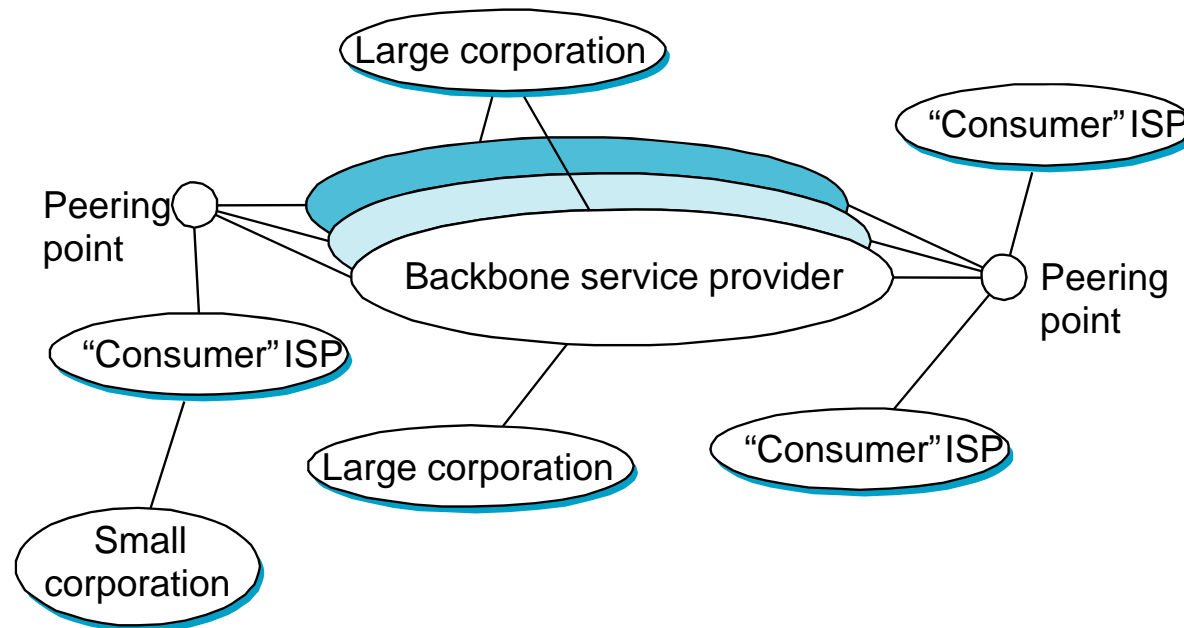
Internet Structure

Recent Past

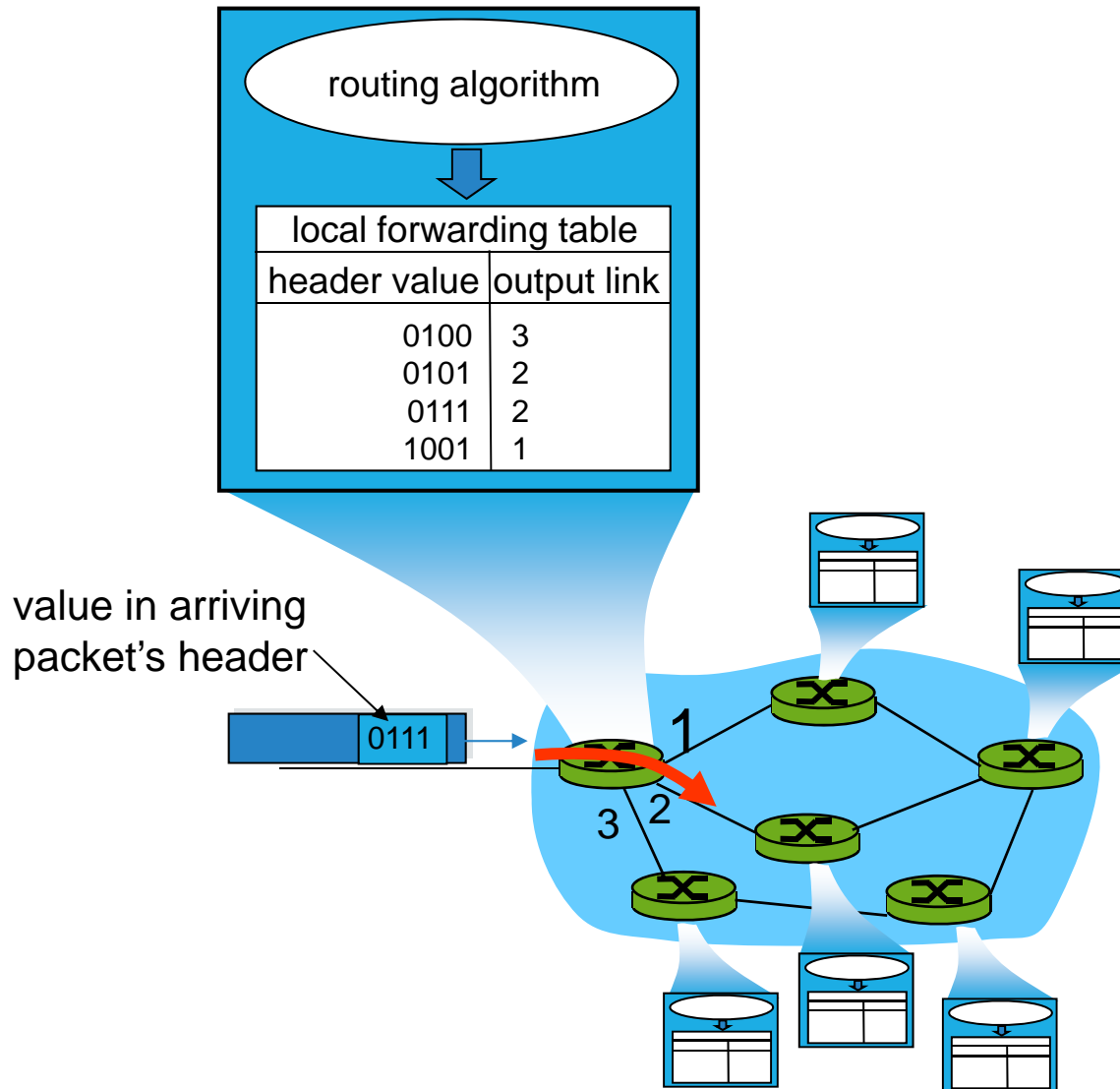


Internet Structure

Today



Routing vs Forwarding



Routing vs Forwarding

- Forwarding table VS Routing table
 - Forwarding table
 - Used when a packet is being forwarded and so must contain enough information to accomplish the forwarding function
 - A row in the forwarding table contains the mapping from a network number to an outgoing interface and some MAC information, such as Ethernet Address of the next hop
 - Routing table
 - Built by the routing algorithm as a precursor to build the forwarding table
 - Generally contains mapping from network numbers to next hops

Forwarding Algorithm

```
D = destination IP address

for each entry (SubnetNum, SubnetMask, NextHop)

    D1 = SubnetMask & D

    if D1 = SubnetNum

        if NextHop is an interface

            deliver datagram directly to D

        else

            deliver datagram to NextHop
```

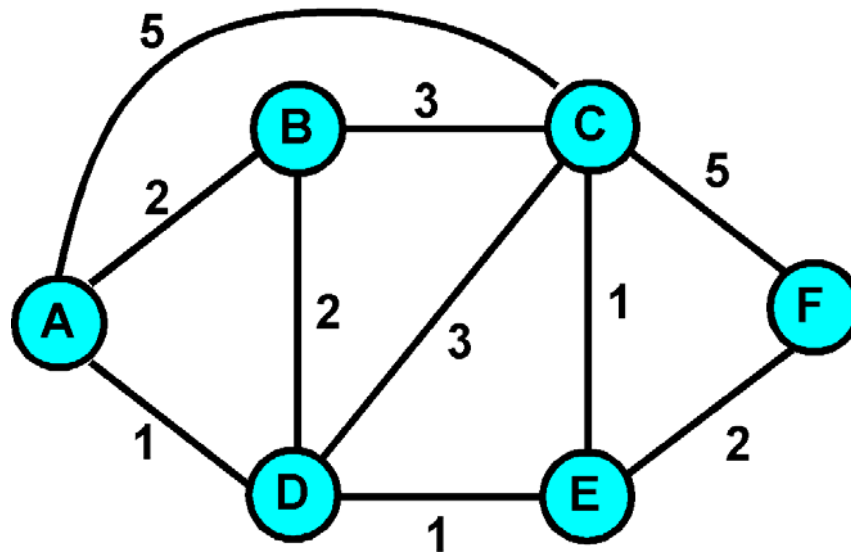
- Use a default router if nothing matches
- Not necessary for all 1s in subnet mask to be contiguous
- Can put multiple subnets on one physical network
- Subnets not visible from the rest of the Internet

Routing Principles

- Routing: delivering a packet to its destination on the best possible path
- Routing steps:
 - (a) determine node network address
 - (b) compute/construct the path
 - (c) forward the packet to destination
- Here, we will focus on (b) - routing alg. For path computation

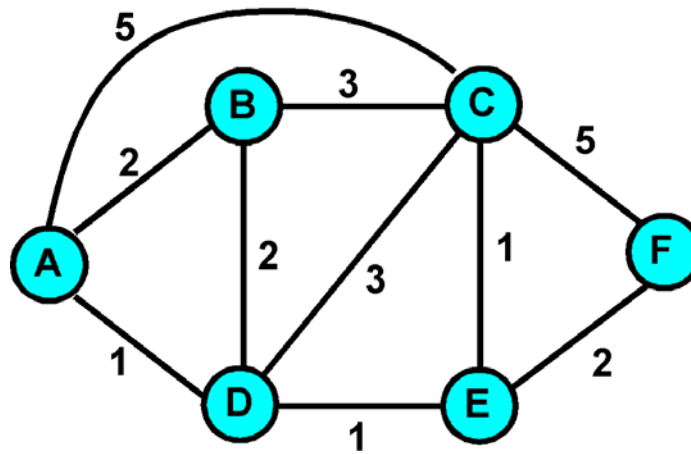
Routing Alg Requirements

- Find path with min delay, cost or other metric
- dynamic reconfiguration after failures/changes
- adaptive load balancing



Graph abstraction

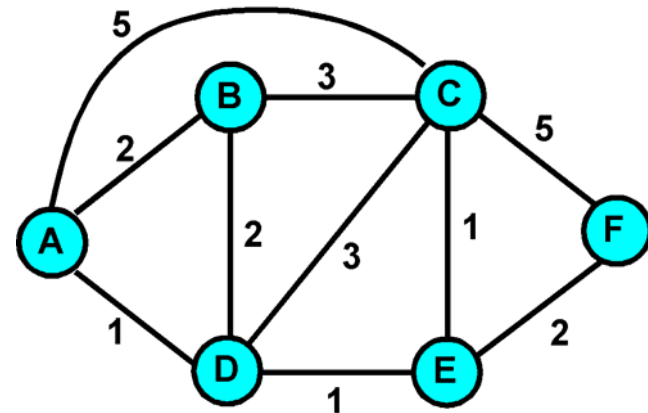
- Graph: $G = (N, E)$
- $N = \text{set of routers} = \{ u, v, w, x, y, z \}$
- $E = \text{set of links} = \{ (u,v), (u,x), (v,x), (v,w), (x,w), (x,y), (w,y), (w,z), (y,z) \}$



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



Routing

- For a simple network, we can calculate all shortest paths and load them into some nonvolatile storage on each node.
- Such a static approach has several shortcomings
 - It does not deal with node or link failures
 - It does not consider the addition of new nodes or links
 - It implies that edge costs cannot change
- What is the solution?
 - Need a distributed and dynamic protocol
 - Two main classes of protocols
 - Distance Vector
 - Link State

Routing Algorithm classification

Global or decentralized?

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

Routing

(a)		
Prefix/Length	Next Hop	
18/8	171.69.245.10	

(b)		
Prefix/Length	Interface	MAC Address
18/8	if0	8:0:2b:e4:b:1:2

Example rows from (a) routing and (b) forwarding tables

Metrics

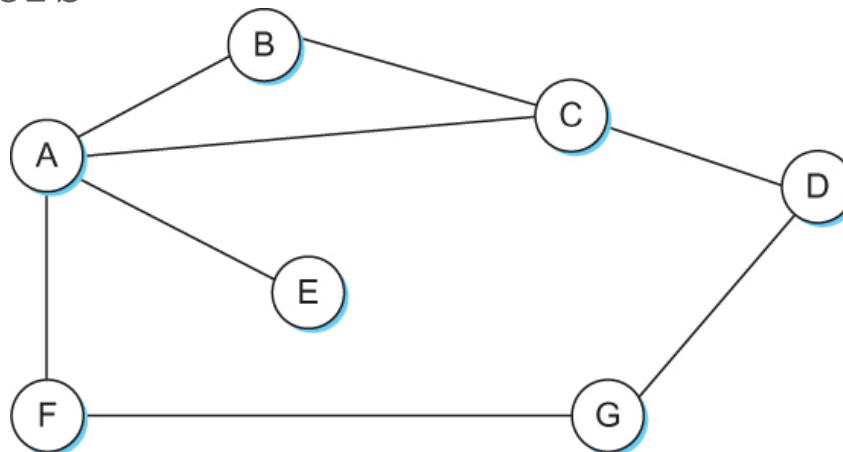
- Original ARPANET metric
 - measures number of packets queued on each link
 - took neither latency or bandwidth into consideration
- New ARPANET metric
 - stamp each incoming packet with its arrival time (**AT**)
 - record departure time (**DT**)
 - when link-level ACK arrives, compute
$$\text{Delay} = (\text{DT} - \text{AT}) + \text{Transmit} + \text{Latency}$$
 - if timeout, reset **DT** to departure time for retransmission
 - link cost = average delay over some time period
- Fine Tuning
 - compressed dynamic range
 - replaced **Delay** with link utilization

Route Propagation

- Know a smarter router
 - hosts know local router
 - local routers know site routers
 - site routers know core router
 - core routers know everything
- Autonomous System (AS)
 - corresponds to an administrative domain
 - examples: University, company, backbone network
 - assign each AS a 16-bit number
- Two-level route propagation hierarchy
 - interior gateway protocol (each AS selects its own)
 - exterior gateway protocol (Internet-wide standard)

Distance Vector

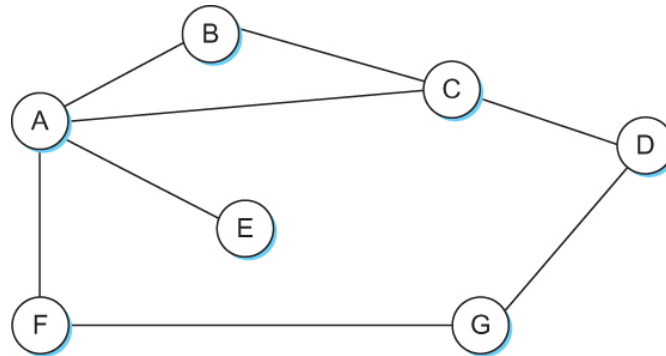
- Each node constructs a one dimensional array (a vector) containing the “distances” (costs) to all other nodes and distributes that vector to its immediate neighbors
- Starting assumption is that each node knows the cost of the link to each of its directly connected neighbors



Distance Vector

- Each node maintains a set of triples
 - (Destination, Cost, NextHop)
- Directly connected neighbors exchange updates
 - periodically (on the order of several seconds)
 - whenever table changes (called *triggered* update)
- Each update is a list of pairs:
 - (Destination, Cost)
- Update local table if receive a “better” route
 - smaller cost
 - came from next-hop
- Refresh existing routes; delete if they time out

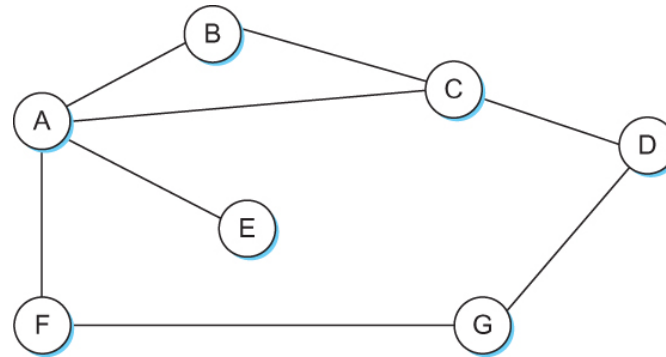
Distance Vector



Information Stored at Node	Distance to Reach Node						
	A	B	C	D	E	F	G
A	0	1	1	∞	1	1	∞
B	1	0	1	∞	∞	∞	∞
C	1	1	0	1	∞	∞	∞
D	∞	∞	1	0	∞	∞	1
E	1	∞	∞	∞	0	∞	∞
F	1	∞	∞	∞	∞	0	1
G	∞	∞	∞	1	∞	1	0

Initial distances stored at each node (global view)

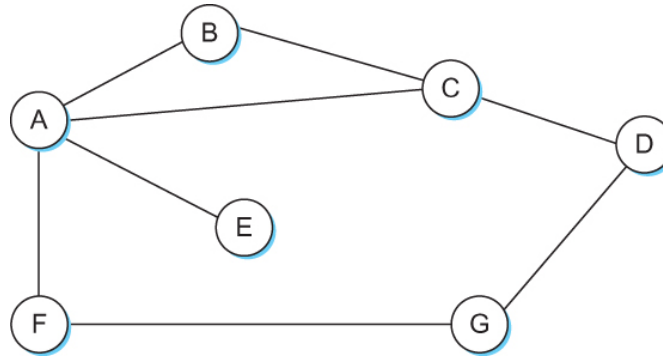
Distance Vector



Destination	Cost	NextHop
B	1	B
C	1	C
D	∞	—
E	1	E
F	1	F
G	∞	—

Initial routing table at node A

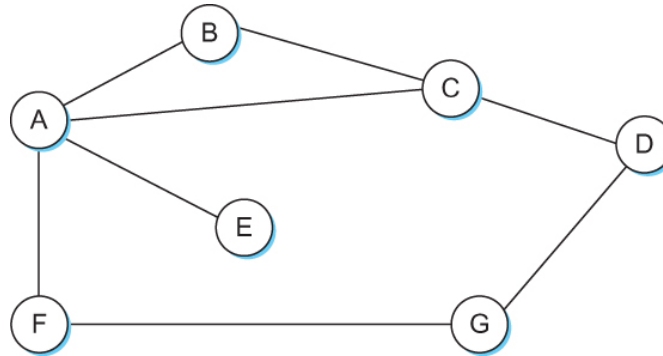
Distance Vector



Destination	Cost	NextHop
B	1	B
C	1	C
D	2	C
E	1	E
F	1	F
G	2	F

Final routing table at node A

Distance Vector



Information Stored at Node	Distance to Reach Node						
	A	B	C	D	E	F	G
A	0	1	1	2	1	1	2
B	1	0	1	2	2	2	3
C	1	1	0	1	2	2	2
D	2	2	1	0	3	2	1
E	1	2	2	3	0	2	3
F	1	2	2	2	2	0	1
G	2	3	2	1	3	1	0

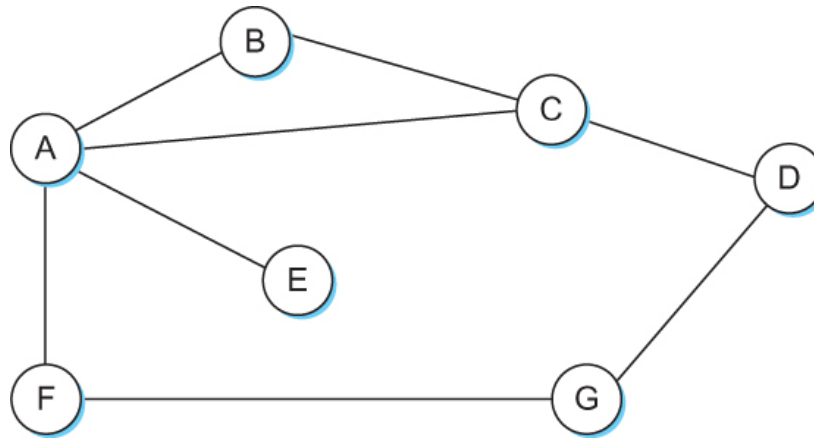
Final distances stored at each node (global view)

Distance Vector

- The distance vector routing algorithm is sometimes called as Bellman-Ford algorithm
- Every T seconds each router sends its table to its neighbor each each router then updates its table based on the new information
- Problems include fast response to good new and slow response to bad news. Also too many messages to update

Distance Vector

- When a node detects a link failure
 - F detects that link to G has failed
 - F sets distance to G to infinity and sends update to A
 - A sets distance to G to infinity since it uses F to reach G
 - A receives periodic update from C with 2-hop path to G
 - A sets distance to G to 3 and sends update to F
 - F decides it can reach G in 4 hops via A



Distance Vector

- Slightly different circumstances can prevent the network from stabilizing
 - Suppose the link from A to E goes down
 - In the next round of updates, A advertises a distance of infinity to E, but B and C advertise a distance of 2 to E
 - Depending on the exact timing of events, the following might happen
 - Node B, upon hearing that E can be reached in 2 hops from C, concludes that it can reach E in 3 hops and advertises this to A
 - Node A concludes that it can reach E in 4 hops and advertises this to C
 - Node C concludes that it can reach E in 5 hops; and so on.
 - This cycle stops only when the distances reach some number that is large enough to be considered infinite
 - Count-to-infinity problem

Count-to-infinity Problem

- Use some relatively small number as an approximation of infinity
- For example, the maximum number of hops to get across a certain network is never going to be more than 16
- One technique to improve the time to stabilize routing is called split horizon
 - When a node sends a routing update to its neighbors, it does not send those routes it learned from each neighbor back to that neighbor
 - For example, if B has the route (E, 2, A) in its table, then it knows it must have learned this route from A, and so whenever B sends a routing update to A, it does not include the route (E, 2) in that update

Count-to-infinity Problem

- In a stronger version of split horizon, called split horizon with poison reverse
 - B actually sends that back route to A, but it puts negative information in the route to ensure that A will not eventually use B to get to E
 - For example, B sends the route (E, ∞) to A

Routing Loops

- Example 1
 - F detects that link to G has failed
 - F sets distance to G to infinity and sends update to A
 - A sets distance to G to infinity since it uses F to reach G
 - A receives periodic update from C with 2-hop path to G
 - A sets distance to G to 3 and sends update to F
 - F decides it can reach G in 4 hops via A
- Example 2
 - link from A to E fails
 - A advertises distance of infinity to E
 - B and C advertise a distance of 2 to E
 - B decides it can reach E in 3 hops; advertises this to A
 - A decides it can reach E in 4 hops; advertises this to C
 - C decides that it can reach E in 5 hops...

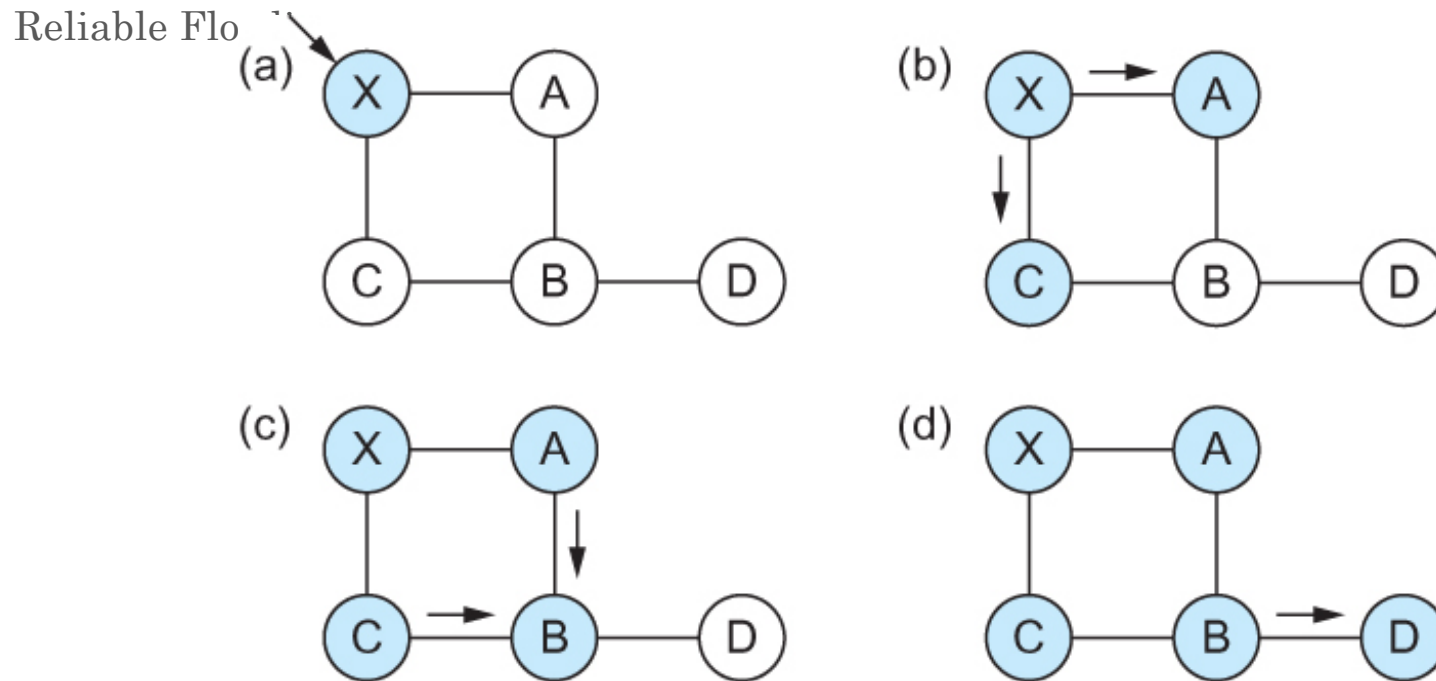
Loop-Breaking Heuristics

- Set infinity to 16
- Split horizon
- Split horizon with poison reverse

Link State Routing

- Strategy: Send to all nodes (not just neighbors) information about directly connected links (not entire routing table).
- Link State Packet (LSP)
 - id of the node that created the LSP
 - cost of link to each directly connected neighbor
 - sequence number (SEQNO)
 - time-to-live (TTL) for this packet
- Reliable Flooding
 - store most recent LSP from each node
 - forward LSP to all nodes but one that sent it
 - generate new LSP periodically; increment SEQNO
 - start SEQNO at 0 when reboot
 - decrement TTL of each stored LSP; discard when TTL=0

Link State



Flooding of link-state packets. (a) LSP arrives at node X; (b) X floods LSP to A and C; (c) A and C flood LSP to B (but not X); (d) flooding is complete

Route Calculation

- Dijkstra's shortest path algorithm
- Let
 - N denotes set of nodes in the graph
 - $l(i, j)$ denotes non-negative cost (weight) for edge (i, j)
 - s denotes this node
 - M denotes the set of nodes incorporated so far
 - $C(n)$ denotes cost of the path from s to node n

```
 $M = \{s\}$ 
```

```
for each  $n$  in  $N - \{s\}$ 
```

```
     $C(n) = l(s, n)$ 
```

```
while ( $N \neq M$ )
```

```
     $M = M$  union  $\{w\}$  such that  $C(w)$  is the minimum for  
    all  $w$  in  $(N - M)$ 
```

```
    for each  $n$  in  $(N - M)$ 
```

```
         $C(n) = \text{MIN}(C(n), C(w) + l(w, n))$ 
```

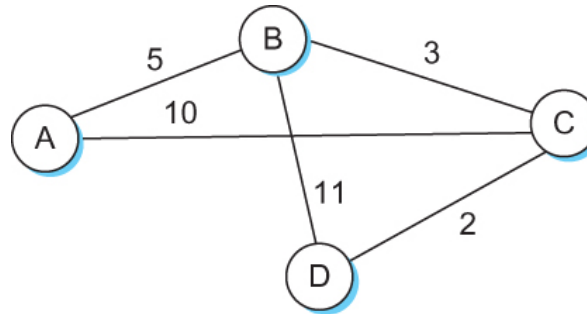

Shortest Path Routing

- In practice, each switch computes its routing table directly from the LSP's it has collected using a realization of Dijkstra's algorithm called the forward search algorithm
- Specifically each switch maintains two lists, known as Tentative and Confirmed
- Each of these lists contains a set of entries of the form (Destination, Cost, NextHop)

Shortest Path Routing

- The algorithm
 - Initialize the **Confirmed** list with an entry for myself; this entry has a cost of 0
 - For the node just added to the **Confirmed** list in the previous step, call it node **Next**, select its LSP
 - For each neighbor (**Neighbor**) of **Next**, calculate the cost (**Cost**) to reach this **Neighbor** as the sum of the cost from myself to **Next** and from **Next** to **Neighbor**
 - If **Neighbor** is currently on neither the **Confirmed** nor the **Tentative** list, then add (**Neighbor**, **Cost**, **Nexthop**) to the **Tentative** list, where **Nexthop** is the direction I go to reach **Next**
 - If **Neighbor** is currently on the **Tentative** list, and the **Cost** is less than the currently listed cost for the **Neighbor**, then replace the current entry with (**Neighbor**, **Cost**, **Nexthop**) where **Nexthop** is the direction I go to reach **Next**
 - If the **Tentative** list is empty, stop. Otherwise, pick the entry from the **Tentative** list with the lowest cost, move it to the **Confirmed** list, and return to Step 2.

Shortest Path Routing



Step	Confirmed	Tentative	Comments
1	(D,0,-)		Since D is the only new member of the confirmed list, look at its LSP.
2	(D,0,-)	(B,11,B) (C,2,C)	D's LSP says we can reach B through B at cost 11, which is better than anything else on either list, so put it on Tentative list; same for C.
3	(D,0,-) (C,2,C)	(B,11,B)	Put lowest-cost member of Tentative (C) onto Confirmed list. Next, examine LSP of newly confirmed member (C).
4	(D,0,-) (C,2,C)	(B,5,C) (A,12,C)	Cost to reach B through C is 5, so replace (B,11,B). C's LSP tells us that we can reach A at cost 12.
5	(D,0,-) (C,2,C) (B,5,C)	(A,12,C)	Move lowest-cost member of Tentative (B) to Confirmed, then look at its LSP.
6	(D,0,-) (C,2,C) (B,5,C)	(A,10,C)	Since we can reach A at cost 5 through B, replace the Tentative entry.
7	(D,0,-) (C,2,C) (B,5,C) (A,10,C)		Move lowest-cost member of Tentative (A) to Confirmed, and we are all done.

Hierarchical Routing

- 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

How to Make Routing Scale

- Flat versus Hierarchical Addresses
- Inefficient use of Hierarchical Address Space
 - class C with 2 hosts ($2/255 = 0.78\%$ efficient)
 - class B with 256 hosts ($256/65535 = 0.39\%$ efficient)
- Still Too Many Networks
 - routing tables do not scale
 - route propagation protocols do not scale

Supernetting

- Assign block of contiguous network numbers to nearby networks
- Called CIDR: Classless Inter-Domain Routing
- Represent blocks with a single pair
`(first_network_address, count)`
- Restrict block sizes to powers of 2
- Use a bit mask (CIDR mask) to identify block size
- All routers must understand CIDR addressing

Routing in the Internet

- The Global Internet consists of Autonomous Systems (AS) interconnected with each other:
 - Stub AS: small corporation
 - Multihomed AS: large corp. (no transit)
 - Transit AS: provider
- Two level routing:
 - Intra-AS: administrator is responsible for choice
 - Inter-AS: unique standard

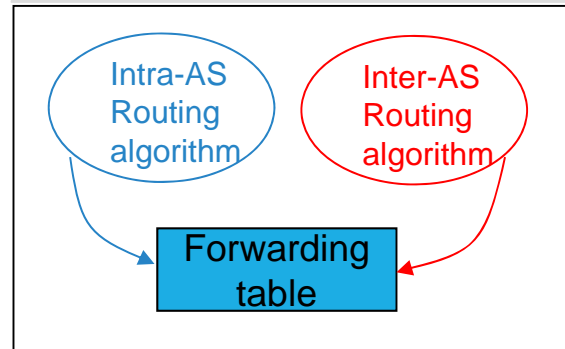
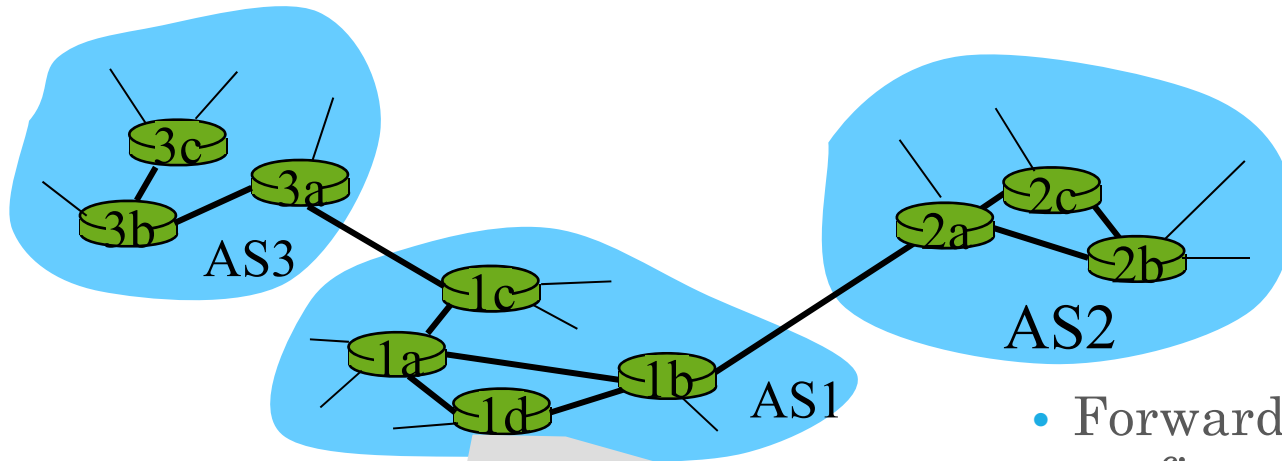
Hierarchical Routing

- 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

Interconnected ASes



- Forwarding table is configured by both intra- and inter-AS routing algorithm
 - Intra-AS sets entries for internal dests
 - Inter-AS & Intra-As sets entries for external dests

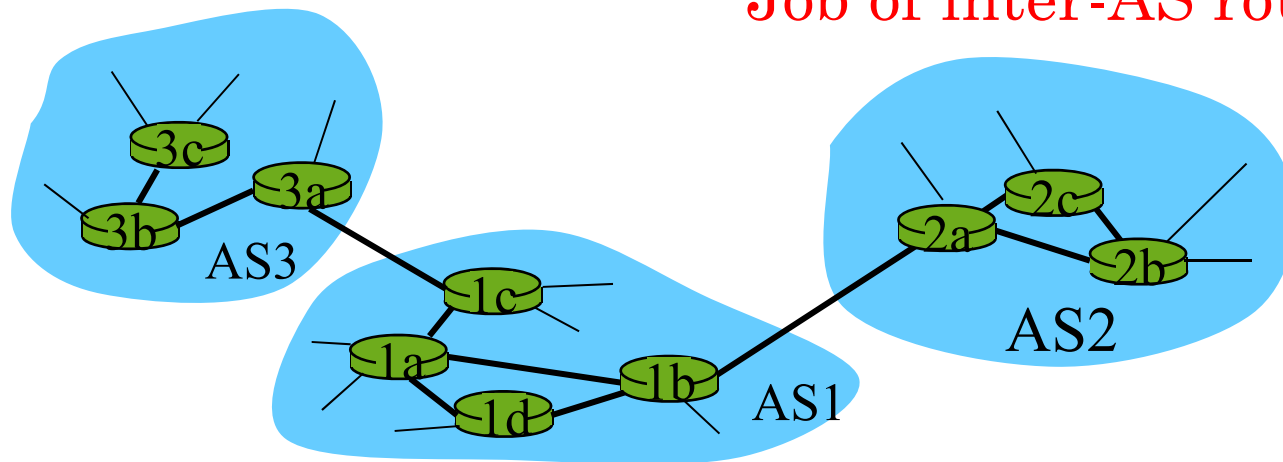
Inter-AS tasks

- Suppose router in AS1 receives datagram for which dest is outside of AS1
 - Router should forward packet towards one of the gateway routers, but which one?

AS1 needs:

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

Job of inter-AS routing!

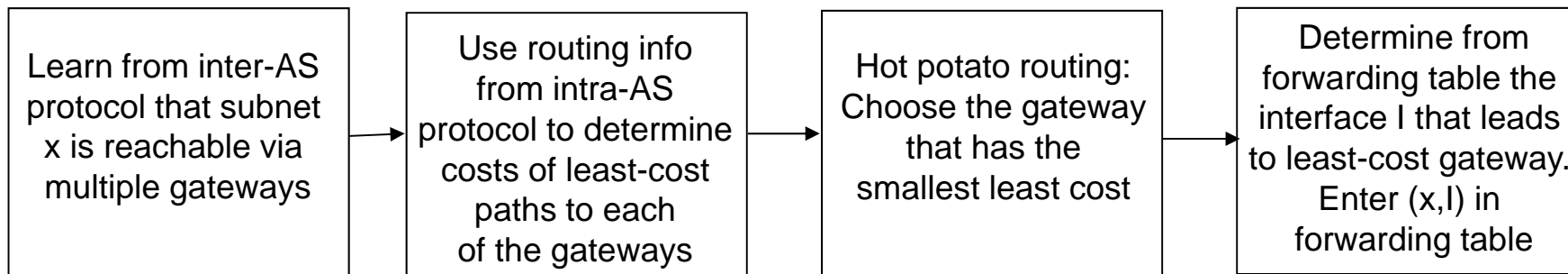


Example: Setting forwarding table in router 1d

- Suppose AS1 learns from the inter-AS protocol that subnet x is reachable from AS3 (gateway 1c) but not from 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.
- Puts in forwarding table entry (x, I) .

Example: Choosing among multiple ASes

- Now suppose AS1 learns from the 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 the job on inter-AS routing protocol!
- **Hot potato routing:** send packet towards closest of two routers.



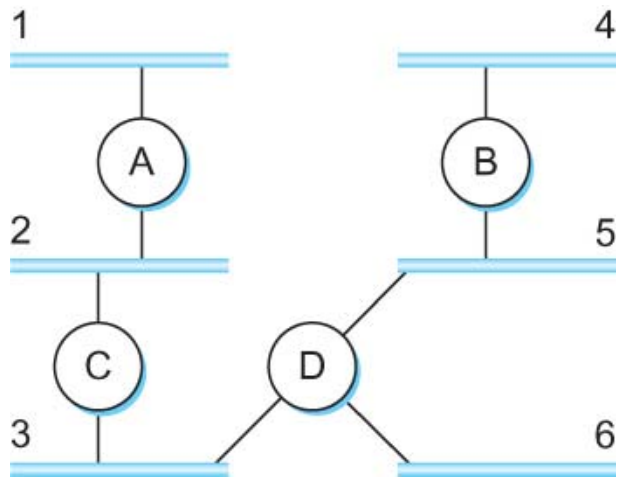
Intra-AS Routing

- 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)
- 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

Routing Information Protocol (RIP)



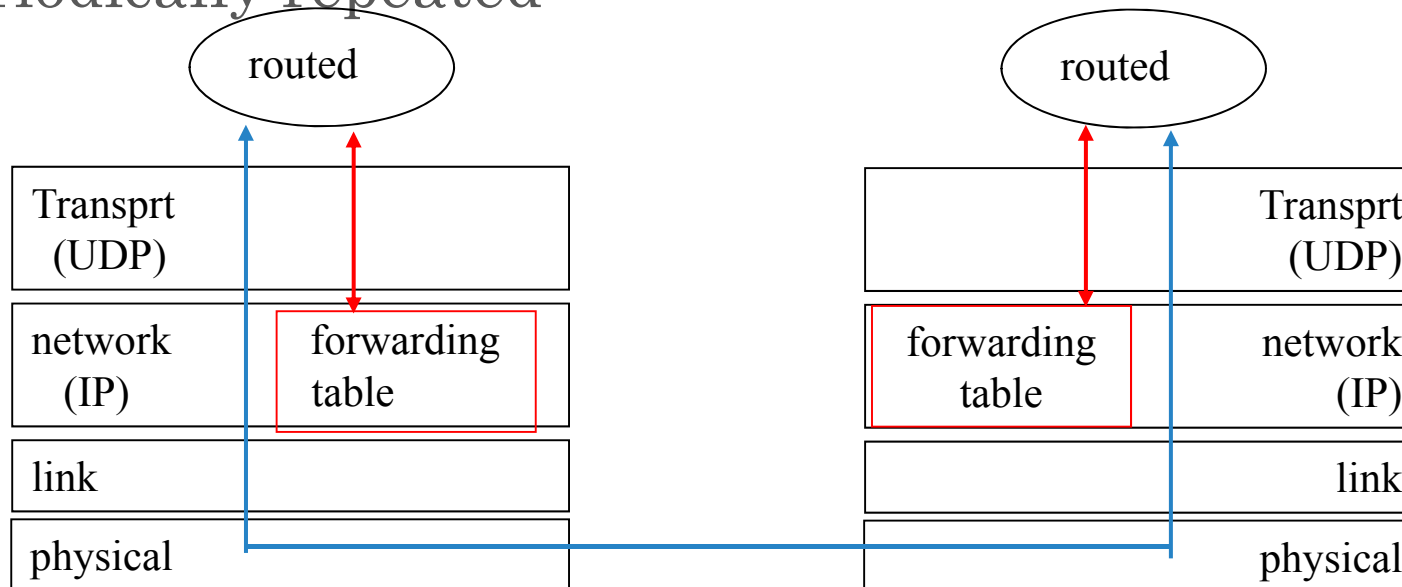
Example Network
running RIP
Format

0		8		16		31	
Command		Version		Must be zero			
Family of net 1				Route Tags			
Address prefix of net 1							
Mask of net 1							
Distance to net 1							
Family of net 2				Route Tags			
Address prefix of net 2							
Mask of net 2							
Distance to net 2							

RIPv2 Packet

RIP Table processing

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



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)

Open Shortest Path First (OSPF)

0	8	16	31
Version	Type	Message length	
SourceAddr			
AreaId			
Checksum		Authentication type	
Authentication			

OSPF Header Format Advertisement

LS Age		Options	Type = 1
Link-state ID			
Advertising router			
LS sequence number			
LS checksum		Length	
0	Flags	0	Number of links
Link ID			
Link data			
Link type	Num_TOS	Metric	
Optional TOS information			
More links			

OSPF Link State

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.

Hierarchical OSPF

- Two-level hierarchy: local area, backbone.
 - Link-state advertisements only in area
 - each nodes has detailed area topology; only know 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.

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

EGP: Exterior Gateway Protocol

- Overview
 - designed for tree-structured Internet
 - concerned with *reachability*, not optimal routes
- Protocol messages
 - neighbor acquisition: one router requests that another be its peer; peers exchange reachability information
 - neighbor reachability: one router periodically tests if the another is still reachable; exchange HELLO/ACK messages; uses a k-out-of-n rule
 - routing updates: peers periodically exchange their routing tables (distance-vector)

Internet inter-AS routing: BGP

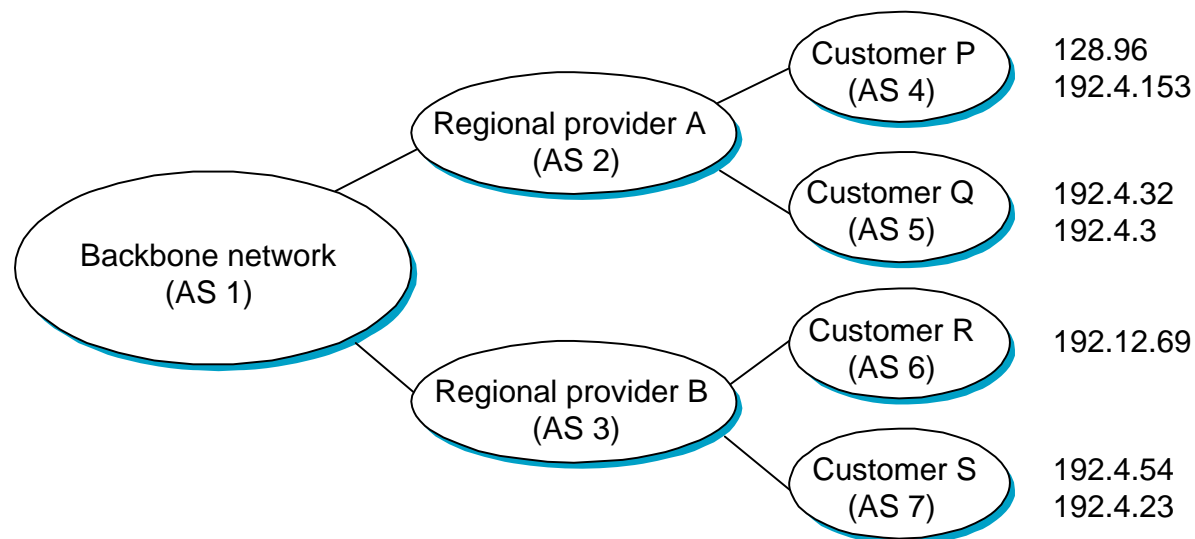
- **BGP (Border Gateway Protocol):** *the* de facto standard
- BGP provides each AS a means to:
 1. Obtain subnet reachability information from neighboring ASs.
 2. Propagate the reachability information to all routers internal to the AS.
 3. Determine “good” routes to subnets based on reachability information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

BGP-4: Border Gateway Protocol

- AS Types
 - stub AS: has a single connection to one other AS
 - carries local traffic only
 - multihomed AS: has connections to more than one AS
 - refuses to carry transit traffic
 - transit AS: has connections to more than one AS
 - carries both transit and local traffic
- Each AS has:
 - one or more border routers
 - one BGP *speaker* that advertises:
 - local networks
 - other reachable networks (transit AS only)
 - gives *path* information

BGP Example

- Speaker for AS2 advertises reachability to P and Q
 - network 128.96, 192.4.153, 192.4.32, and 192.4.3, can be reached directly from AS2



- Speaker for backbone advertises
 - networks 128.96, 192.4.153, 192.4.32, and 192.4.3 can be reached along the path (AS1, AS2).
- Speaker can cancel previously advertised paths