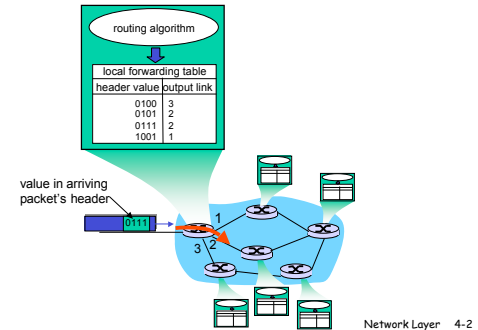


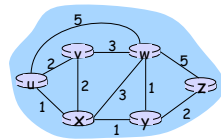
Routing Algorithms and Routing in the Internet

Network Layer 4-1

Interplay between routing and forwarding



Graph abstraction



Graph: $G = (N, E)$

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

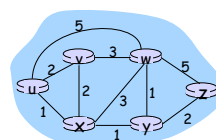
E = set of links = $\{(u,v), (u,x), (v,x), (v,w), (x,w), (x,y), (w,y), (w,z), (y,z)\}$

Remark: Graph abstraction is useful in other network contexts

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

Network Layer 4-3

Graph abstraction: costs



$c(x, x')$ = cost of link (x, x')

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

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

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

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

Routing algorithm: algorithm that finds least-cost path

Network Layer 4-4

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 4-5

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 dest's

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 dest. v
- $p(v)$: predecessor node along path from source to v
- N' : set of nodes whose least cost path definitively known

Network Layer 4-6

Dijkstra's Algorithm

```

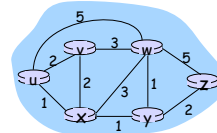
1 Initialization:
2 N' = {u}
3 for all nodes v
4   if v adjacent to u
5     then D(v) = c(u,v)
6   else D(v) = ∞
7
8 Loop
9   find w not in N' such that D(w) is a minimum
10  add w to N'
11  update D(v) for all v adjacent to w and not in N':
12    D(v) = min( D(v), D(w) + c(w,v) )
13  /* new cost to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

```

Network Layer 4-7

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	∞	∞
1	ux	2,u	4,x	2,x	∞	∞
2	uxy	2,u	3,y	4,y	4,y	∞
3	uxyv	2,u	3,y	4,y	4,y	4,y
4	uxyvw	2,u	3,y	4,y	4,y	4,y
5	uxyvwz	2,u	3,y	4,y	4,y	4,y



Network Layer 4-8

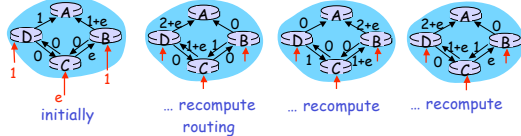
Dijkstra's algorithm, discussion

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 4-9

Distance Vector Algorithm (1)

Bellman-Ford Equation (dynamic programming)

Define

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

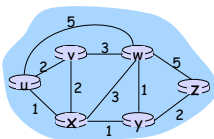
Then

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

where min is taken over all neighbors of x

Network Layer 4-10

Bellman-Ford example (2)



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

Network Layer 4-11

Distance Vector Algorithm (3)

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

Network Layer 4-12

Distance vector algorithm (4)

Basic idea:

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

$$D_x(y) \leftarrow \min_v \{c(x,v) + D_v(y)\} \text{ for each node } y \in N$$
- Under minor, natural conditions, the estimate $D_x(y)$ converge the actual least cost $d_x(y)$

Network Layer 4-13

Distance Vector Algorithm (5)

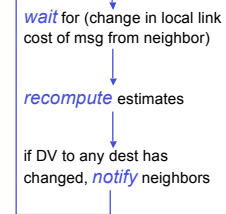
Iterative, asynchronous:

- each local iteration caused by:
 - local link cost change
 - DV update message from neighbor

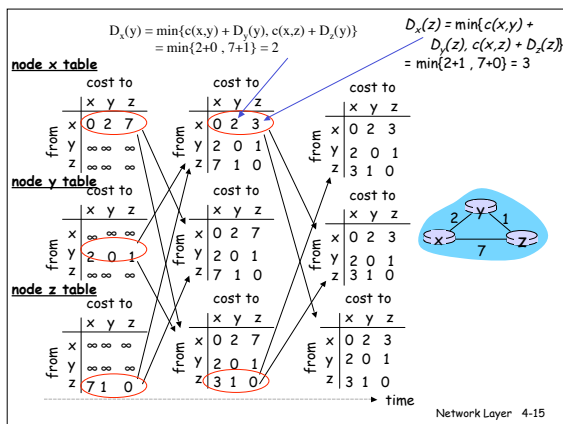
Distributed:

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

Each node:



Network Layer 4-14

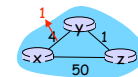


Network Layer 4-15

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.

Network Layer 4-16

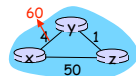
Distance Vector: link cost changes

Link cost changes:

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

Poisoned reverse:

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



Network Layer 4-17

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: $O(n^2)$ 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 4-18

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 4-19

Hierarchical Routing

- aggregate routers into regions, "autonomous systems" (AS)

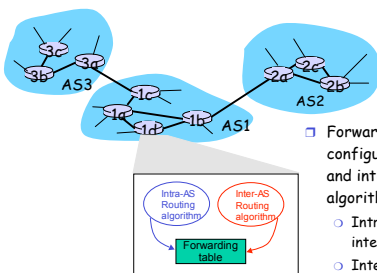
Gateway router

- Direct link to router in another AS

- routers in same AS run same routing protocol
 - "intra-AS" routing protocol
 - routers in different AS can run different intra-AS routing protocol

Network Layer 4-20

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

Network Layer 4-21

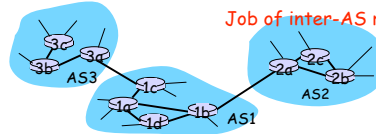
Inter-AS tasks

- 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!



Network Layer 4-22

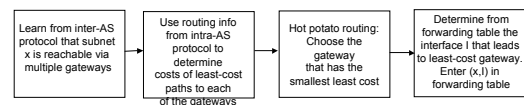
Example: Setting forwarding table in router 1d

- 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) .

Network Layer 4-23

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.



Network Layer 4-24

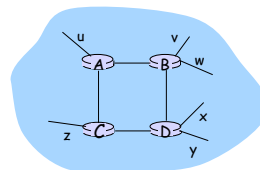
Intra-AS Routing

- ❑ 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 4-25

RIP (Routing Information Protocol)

- ❑ Distance vector algorithm
- ❑ Included in BSD-UNIX Distribution in 1982
- ❑ Distance metric: # of hops (max = 15 hops)



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

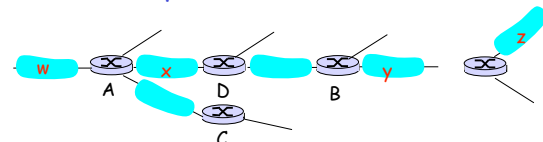
Network Layer 4-26

RIP advertisements

- ❑ 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 4-27

RIP: Example

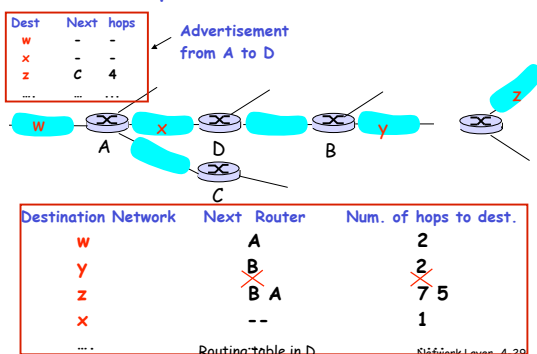


Destination Network	Next Router	Num. of hops to dest.
w	A	2
y	B	2
z	B	7
x	--	1
...

Routing Table in D

Network Layer 4-28

RIP: Example



Network Layer 4-29

RIP: Link Failure and Recovery

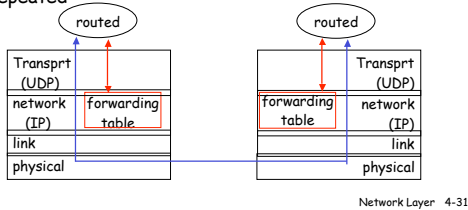
If no advertisement heard after 180 sec --> neighbor/link declared dead

- 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 4-30

RIP Table processing

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



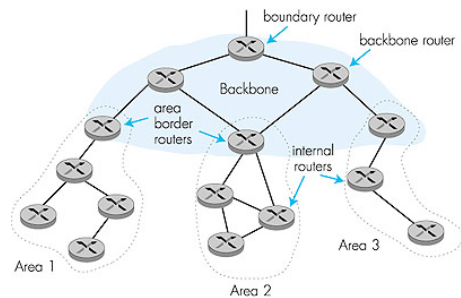
OSPF (Open Shortest Path First)

- ❑ "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)
- Network Layer 4-32

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.
- Network Layer 4-33

Hierarchical OSPF



Hierarchical OSPF

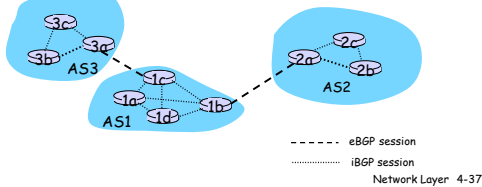
- ❑ **Two-level hierarchy**: local area, backbone.
 - Link-state advertisements only in area
 - each nodes has detailed area topology; only know direction (shortest path) to nets in other areas.
 - ❑ **Area border routers**: "summarize" distances to nets in own area, advertise to other Area Border routers.
 - ❑ **Backbone routers**: run OSPF routing limited to backbone.
 - ❑ **Boundary routers**: connect to other AS's.
- Network Layer 4-35

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"*
- Network Layer 4-36

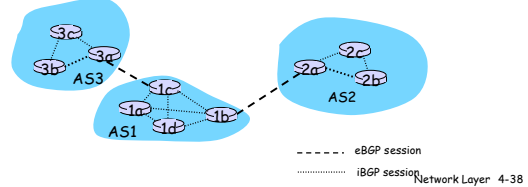
BGP basics

- Pairs of routers (BGP peers) exchange routing info over semi-permanent TCP conns: **BGP sessions**
- Note that BGP sessions do not correspond to physical links.
- When AS2 advertises a prefix to AS1, AS2 is **promising** it will forward any datagrams destined to that prefix towards the prefix.
 - AS2 can aggregate prefixes in its advertisement



Distributing reachability info

- With eBGP session between 3a and 1c, AS3 sends prefix reachability info to AS1.
- 1c can then use iBGP to distribute this new prefix reach info to all routers in AS1
- 1b can then re-advertise the new reach info to AS2 over the 1b-to-2a eBGP session
- When router learns about a new prefix, it creates an entry for the prefix in its forwarding table.



Path attributes & BGP routes

- When advertising a prefix, advert includes BGP attributes.
 - prefix + attributes = "route"
- Two important attributes:
 - **AS-PATH**: contains the ASs through which the advert for the prefix passed: AS 67 AS 17
 - **NEXT-HOP**: Indicates the specific internal-AS router to next-hop AS. (There may be multiple links from current AS to next-hop-AS.)
- When gateway router receives route advert, uses **import policy** to accept/decline.

Network Layer 4-39

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

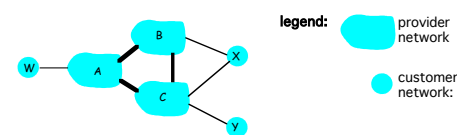
Network Layer 4-40

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

Network Layer 4-41

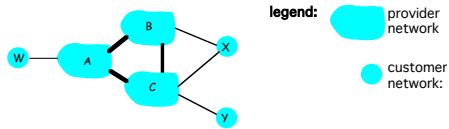
BGP routing policy



- A,B,C are **provider networks**
- X,W,Y are customer (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

Network Layer 4-42

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!

Network Layer 4-43

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 4-44