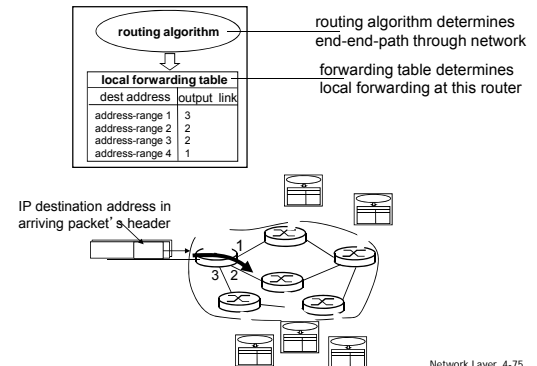


Chapter 4: outline

- 4.1 introduction
- 4.2 virtual circuit and datagram networks
- 4.3 what's inside a router
- 4.4 IP: Internet Protocol
 - datagram format
 - IPv4 addressing
 - ICMP
 - IPv6
- 4.5 routing algorithms
 - link state
 - distance vector
 - hierarchical routing
- 4.6 routing in the Internet
 - RIP
 - OSPF
 - BGP
- 4.7 broadcast and multicast routing

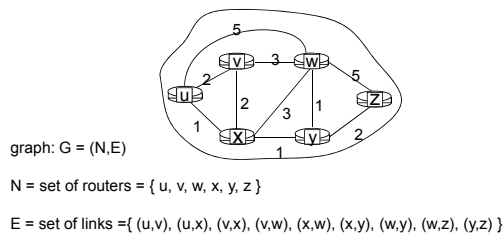
Network Layer 4-74

Interplay between routing, forwarding



Network Layer 4-75

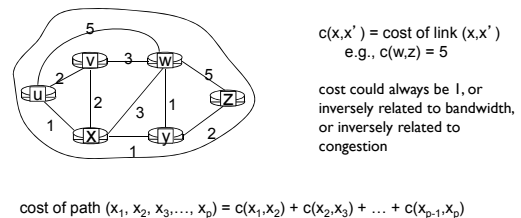
Graph abstraction



aside: graph abstraction is useful in other network contexts, e.g., P2P, where N is set of peers and E is set of TCP connections

Network Layer 4-76

Graph abstraction: costs



key question: what is the least-cost path between u and z ?
routing algorithm: algorithm that finds that least cost path

Network Layer 4-77

Routing algorithm classification

Q: 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

Q: static or dynamic?

static:

- ❖ routes change slowly over time

dynamic:

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

Network Layer 4-78

Chapter 4: outline

4.1 introduction

4.2 virtual circuit and datagram networks

4.3 what's inside a router

4.4 IP: Internet Protocol

- datagram format
- IPv4 addressing
- ICMP
- IPv6

4.5 routing algorithms

- link state
- distance vector
- hierarchical routing

4.6 routing in the Internet

- RIP
- OSPF
- BGP

4.7 broadcast and multicast routing

Network Layer 4-79

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 4-80

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 $D(v) = \infty$

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'**

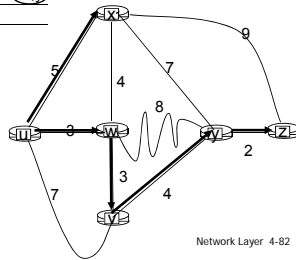
Network Layer 4-81

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	7,u	(3,u)	5,u	∞	∞
1	uw	6,w	(5,u)	11,w	∞	∞
2	uwx	(6,w)		11,w	14,x	∞
3	uwxv			(10,y)	14,x	∞
4	uwxvy				(12,y)	∞
5	uwxvyz					∞

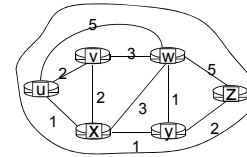
notes:

- construct shortest path tree by tracing predecessor nodes
- ties can exist (can be broken arbitrarily)



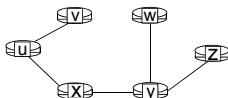
Dijkstra's algorithm: another example

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



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)

Network Layer 4-84

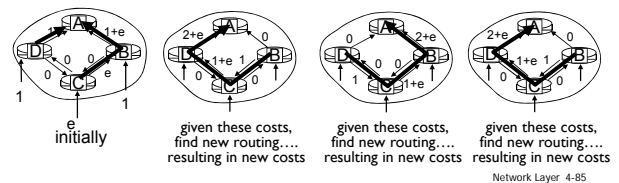
Dijkstra's algorithm, discussion

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., support link cost equals amount of carried traffic:



Chapter 4: outline

- 4.1 introduction
- 4.2 virtual circuit and datagram networks
- 4.3 what's inside a router
- 4.4 IP: Internet Protocol
 - datagram format
 - IPv4 addressing
 - ICMP
 - IPv6
- 4.5 routing algorithms
 - link state
 - distance vector
 - hierarchical routing
- 4.6 routing in the Internet
 - RIP
 - OSPF
 - BGP
- 4.7 broadcast and multicast routing

Network Layer 4-86

Distance vector algorithm

Bellman-Ford equation (dynamic programming)

let

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

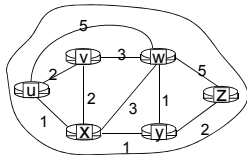
then

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

\downarrow \downarrow \downarrow
 \min taken over all neighbors v of x cost to neighbor v cost from neighbor v to destination y

Network Layer 4-87

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 achieving minimum is next hop in shortest path, used in forwarding table

Network Layer 4-88

Distance vector algorithm

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

Network Layer 4-89

Distance vector algorithm

key idea:

- ❖ from time-to-time, each node sends its own distance vector estimate to neighbors
- ❖ when 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 to the actual least cost $d_x(y)$

Network Layer 4-90

Distance vector algorithm

iterative, asynchronous:

each local iteration

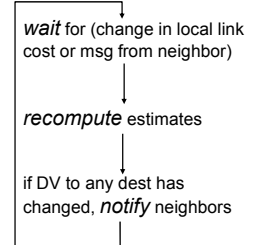
caused by:

- ❖ local link cost change
- ❖ DV update message from neighbor

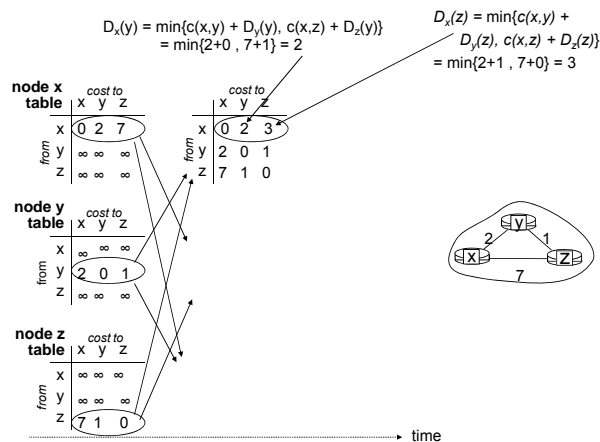
distributed:

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

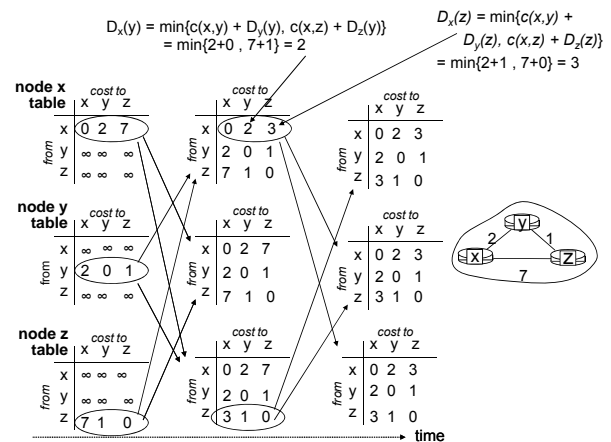
each node:



Network Layer 4-91



Network Layer 4-92

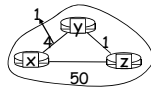


Network Layer 4-93

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”

t_0 : y detects link-cost change, updates its DV, informs its neighbors.

t_1 : z receives update from y, updates its table, computes new least cost to x, sends its neighbors its DV.

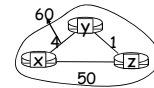
t_2 : y receives z's update, updates its distance table. y's least costs do not change, so y does not send a message to z.

Network Layer 4-94

Distance vector: link cost changes

link cost changes:

- ❖ node detects local link cost change
- ❖ 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 4-95

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 4-96

Chapter 4: outline

4.1 introduction

4.2 virtual circuit and datagram networks

4.3 what's inside a router

4.4 IP: Internet Protocol

- datagram format
- IPv4 addressing
- ICMP
- IPv6

4.5 routing algorithms

- link state
- distance vector
- hierarchical routing

4.6 routing in the Internet

- RIP
- OSPF
- BGP

4.7 broadcast and multicast routing

Network Layer 4-97

Hierarchical routing

our routing study thus far - idealization

- ❖ all routers identical
- ❖ network “flat”
- ... *not* true in practice

scale: with 600 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 4-98

Hierarchical routing

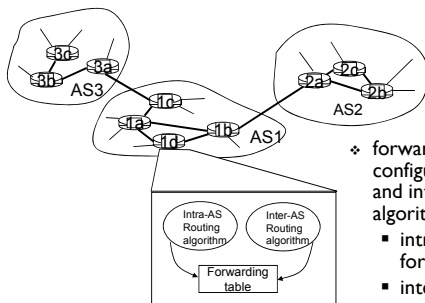
- ❖ aggregate routers into regions, “autonomous systems” (AS)
- ❖ routers in same AS run same routing protocol
 - “intra-AS” routing protocol
 - routers in different AS can run different intra-AS routing protocol

gateway router:

- ❖ at “edge” of its own AS
- ❖ has link to router in another AS

Network Layer 4-99

Interconnected ASes



- ❖ forwarding table 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 4-100

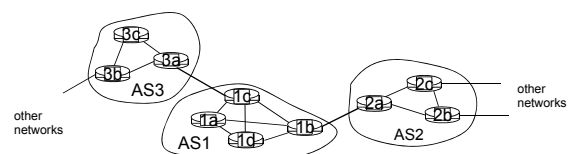
Inter-AS tasks

- ❖ suppose router in AS1 receives datagram destined outside of AS1:
 - router should forward packet to gateway router, but which one?

AS1 must:

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

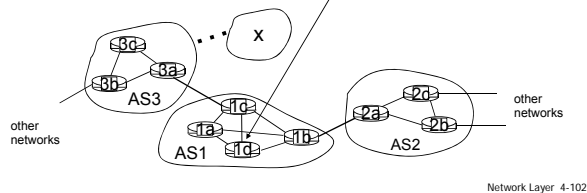
job of inter-AS routing!



Network Layer 4-101

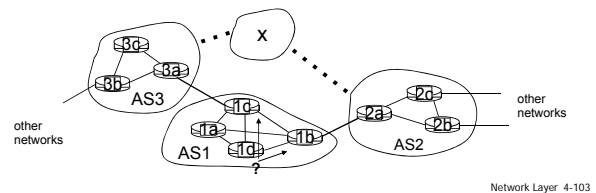
Example: setting forwarding table in router 1d

- ❖ suppose AS1 learns (via inter-AS protocol) that subnet x is 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)



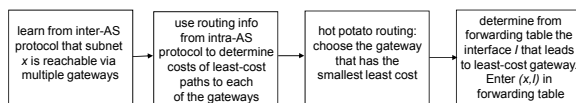
Example: choosing among multiple ASes

- ❖ 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 which gateway it should forward packets towards for dest x
 - this is also job of inter-AS routing protocol!



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.



Chapter 4: outline

- 4.1 introduction
- 4.2 virtual circuit and datagram networks
- 4.3 what's inside a router
- 4.4 IP: Internet Protocol
 - datagram format
 - IPv4 addressing
 - ICMP
 - IPv6
- 4.5 routing algorithms
 - link state
 - distance vector
 - hierarchical routing
- 4.6 routing in the Internet
 - RIP
 - OSPF
 - BGP
- 4.7 broadcast and multicast routing

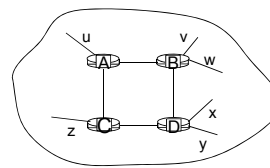
Intra-AS Routing

- ❖ 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 4-106

RIP (Routing Information Protocol)

- ❖ included in BSD-UNIX distribution in 1982
- ❖ distance vector algorithm
 - distance metric: # hops (max = 15 hops), each link has cost 1
 - DVs exchanged with neighbors every 30 sec in response message (aka advertisement)
 - each advertisement: list of up to 25 destination *subnets* (in IP addressing sense)

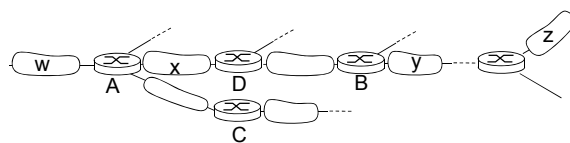


from router A to destination *subnets*:

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

Network Layer 4-107

RIP: example

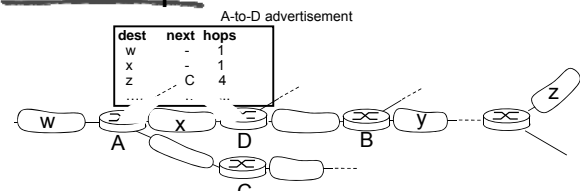


routing table in router D

destination subnet	next router	# hops to dest
w	A	2
y	B	2
x	--	1
....

Network Layer 4-108

RIP: example



A-to-D advertisement

dest	next	hops
w	-	1
x	-	1
z	C	4

routing table in router D

destination subnet	next router	# hops to dest
w	A	2
y	B	2
z	B A	2 5
x	--	1
....

Network Layer 4-109

RIP: link failure, recovery

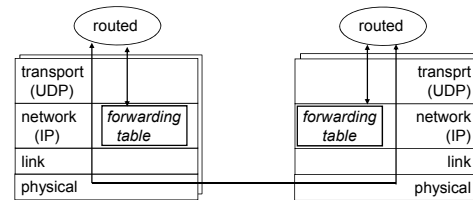
if no advertisement heard after 180 sec -->
neighbor/link declared dead

- 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 4-110

RIP table processing

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



Network Layer 4-111

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
- ❖ advertisements flooded to *entire* AS
 - carried in OSPF messages directly over IP (rather than TCP or UDP)
- ❖ *IS-IS* routing protocol: nearly identical to OSPF

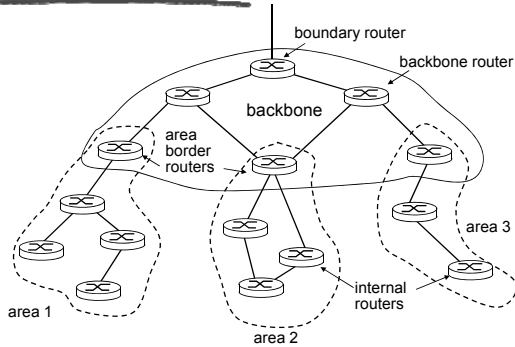
Network Layer 4-112

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 ToS; high for real time ToS)
- ❖ integrated uni- and multicast support:
 - Multicast OSPF (MOSPF) uses same topology data base as OSPF
- ❖ hierarchical OSPF in large domains.

Network Layer 4-113

Hierarchical OSPF



Network Layer 4-114

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 4-115

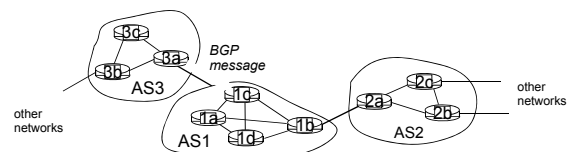
Internet inter-AS routing: BGP

- ❖ BGP (Border Gateway Protocol): *the de facto* inter-domain routing protocol
 - “glue that holds the Internet together”
- ❖ BGP provides each AS a means to:
 - eBGP: obtain subnet reachability information from neighboring ASs.
 - iBGP: propagate reachability information to all AS-internal routers.
 - determine “good” routes to other networks based on reachability information and policy.
- ❖ allows subnet to advertise its existence to rest of Internet: “*i am here*”

Network Layer 4-116

BGP basics

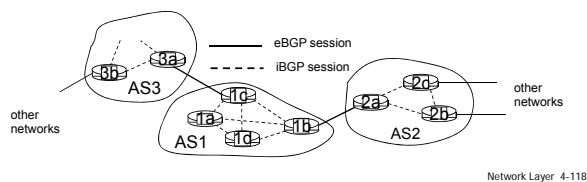
- ❖ BGP session: two BGP routers (“peers”) exchange BGP messages:
 - advertising *paths* to different destination network prefixes (“path vector” protocol)
 - exchanged over semi-permanent TCP connections
- ❖ when AS3 advertises a prefix to AS1:
 - AS3 *promises* it will forward datagrams towards that prefix
 - AS3 can aggregate prefixes in its advertisement



Network Layer 4-117

BGP basics: distributing path information

- ❖ 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 prefix in its forwarding table.



Path attributes and 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
 - NEXT-HOP: indicates specific internal-AS router to next-hop AS. (may be multiple links from current AS to next-hop-AS)
- ❖ gateway router receiving route advertisement uses import policy to accept/decline
 - e.g., never route through AS x
 - *policy-based routing*

Network Layer 4-119

BGP route selection

- ❖ router may learn about more than 1 route to destination AS, selects route based on:
 1. local preference value attribute: policy decision
 2. shortest AS-PATH
 3. closest NEXT-HOP router: hot potato routing
 4. additional criteria

Network Layer 4-120

BGP messages

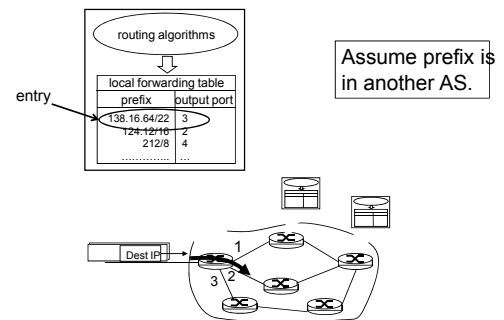
- ❖ BGP messages exchanged between peers over TCP connection
- ❖ BGP messages:
 - OPEN: opens TCP connection to peer and authenticates sender
 - UPDATE: advertises new path (or withdraws old)
 - KEEPALIVE: keeps connection alive in absence of UPDATES; also ACKs OPEN request
 - NOTIFICATION: reports errors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

Network Layer 4-121

Putting it Altogether: How Does an Entry Get Into a Router's Forwarding Table?

- ❖ Answer is complicated!
- ❖ Ties together hierarchical routing (Section 4.5.3) with BGP (4.6.3) and OSPF (4.6.2).
- ❖ Provides nice overview of BGP!

How does entry get in forwarding table?

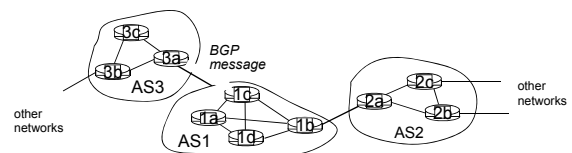


How does entry get in forwarding table?

High-level overview

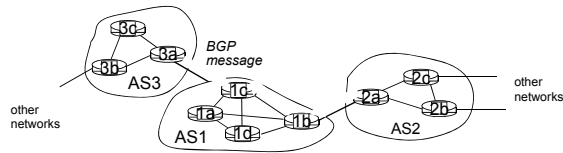
1. Router becomes aware of prefix
2. Router determines output port for prefix
3. Router enters prefix-port in forwarding table

Router becomes aware of prefix



- ❖ BGP message contains "routes"
- ❖ "route" is a prefix and attributes: AS-PATH, NEXT-HOP,...
- ❖ Example: route:
 - ❖ Prefix: 138.16.64/22 ; AS-PATH: AS3 AS131 ;
 - ❖ NEXT-HOP: 201.44.13.125

Router may receive multiple routes



- ❖ Router may receive multiple routes for same prefix
- ❖ Has to select one route

Select best BGP route to prefix

- ❖ Router selects route based on shortest AS-PATH

- ❖ Example:

❖ AS2 AS17 to 138.16.64/22

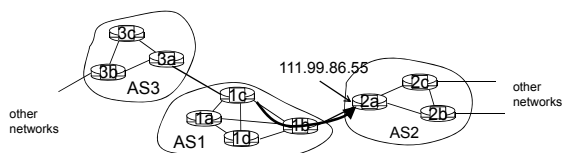
❖ AS3 AS131 AS201 to 138.16.64/22

- ❖ What if there is a tie? We'll come back to that!

select

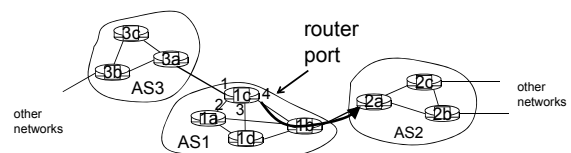
Find best intra-route to BGP route

- ❖ Use selected route's NEXT-HOP attribute
 - Route's NEXT-HOP attribute is the IP address of the router interface that begins the AS PATH.
- ❖ Example:
 - ❖ AS-PATH: AS2 AS17 ; NEXT-HOP: 111.99.86.55
- ❖ Router uses OSPF to find shortest path from 1c to 111.99.86.55



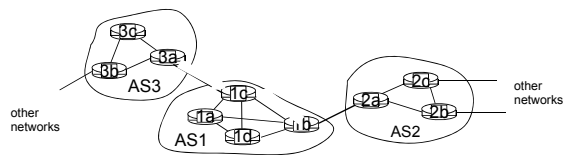
Router identifies port for route

- ❖ Identifies port along the OSPF shortest path
- ❖ Adds prefix-port entry to its forwarding table:
 - (138.16.64/22, port 4)



Hot Potato Routing

- ❖ Suppose there two or more best inter-routes.
- ❖ Then choose route with closest NEXT-HOP
 - Use OSPF to determine which gateway is closest
 - Q: From 1c, chose AS3 AS131 or AS2 AS17?
 - A: route AS3 AS201 since it is closer

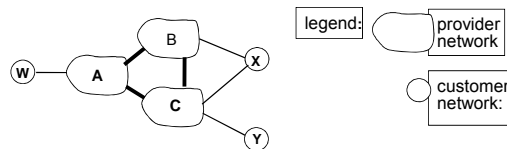


How does entry get in forwarding table?

Summary

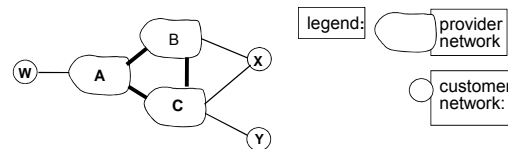
1. Router becomes aware of prefix
 - via BGP route advertisements from other routers
2. Determine router output port for prefix
 - Use BGP route selection to find best inter-AS route
 - Use OSPF to find best intra-AS route leading to best inter-AS route
 - Router identifies router port for that best route
3. Enter prefix-port entry in forwarding table

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

BGP routing policy (2)



- ❖ A advertises path AW to B
- ❖ B advertises path BAW to X
- ❖ Should B advertise path BAW to C?
 - 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!

Why different Intra-, 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 4-134

Chapter 4: outline

4.1 introduction

4.2 virtual circuit and datagram networks

4.3 what's inside a router

4.4 IP: Internet Protocol

- datagram format
- IPv4 addressing
- ICMP
- IPv6

4.5 routing algorithms

- link state
- distance vector
- hierarchical routing

4.6 routing in the Internet

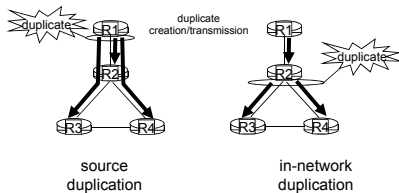
- RIP
- OSPF
- BGP

4.7 broadcast and multicast routing

Network Layer 4-135

Broadcast routing

- ❖ deliver packets from source to all other nodes
- ❖ source duplication is inefficient:



- ❖ source duplication: how does source determine recipient addresses?

Network Layer 4-136

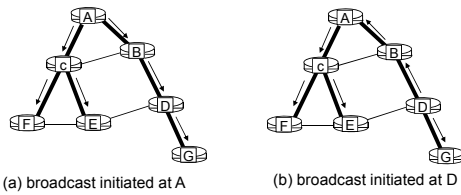
In-network duplication

- ❖ *flooding*: when node receives broadcast packet, sends copy to all neighbors
 - problems: cycles & broadcast storm
- ❖ *controlled flooding*: node only broadcasts pkt if it hasn't broadcast same packet before
 - node keeps track of packet ids already broadcasted
 - or reverse path forwarding (RPF): only forward packet if it arrived on shortest path between node and source
- ❖ *spanning tree*:
 - no redundant packets received by any node

Network Layer 4-137

Spanning tree

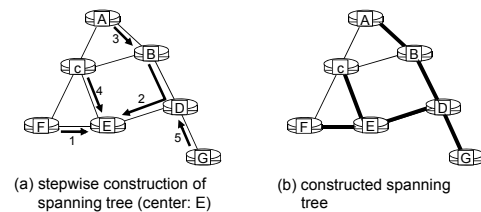
- ❖ first construct a spanning tree
- ❖ nodes then forward/make copies only along spanning tree



Network Layer 4-138

Spanning tree: creation

- ❖ center node
- ❖ each node sends unicast join message to center node
 - message forwarded until it arrives at a node already belonging to spanning tree

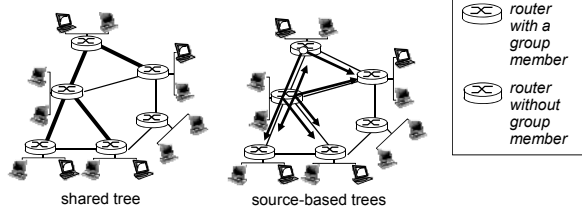


Network Layer 4-139

Multicast routing: problem statement

goal: find a tree (or trees) connecting routers having local mcast group members

- ❖ tree: not all paths between routers used
- ❖ shared-tree: same tree used by all group members
- ❖ source-based: different tree from each sender to rcvrs



Network Layer 4-140

Approaches for building mcast trees

approaches:

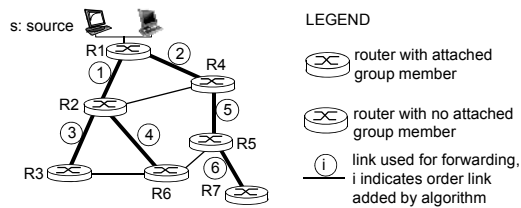
- ❖ source-based tree: one tree per source
 - shortest path trees
 - reverse path forwarding
- ❖ group-shared tree: group uses one tree
 - minimal spanning (Steiner)
 - center-based trees

...we first look at basic approaches, then specific protocols adopting these approaches

Network Layer 4-141

Shortest path tree

- ❖ mcast forwarding tree: tree of shortest path routes from source to all receivers
 - Dijkstra's algorithm



Network Layer 4-142

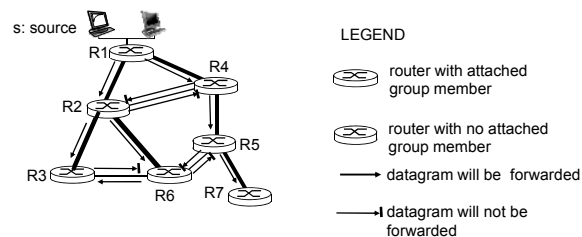
Reverse path forwarding

- ❖ rely on router's knowledge of unicast shortest path from it to sender
- ❖ each router has simple forwarding behavior:

if (mcast datagram received on incoming link on shortest path back to center)
then flood datagram onto all outgoing links
else ignore datagram

Network Layer 4-143

Reverse path forwarding: example

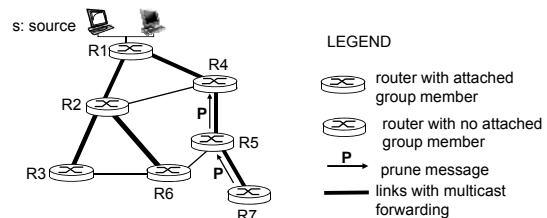


- ❖ result is a source-specific reverse SPT
 - may be a bad choice with asymmetric links

Network Layer 4-144

Reverse path forwarding: pruning

- ❖ forwarding tree contains subtrees with no mcast group members
 - no need to forward datagrams down subtree
 - "prune" msgs sent upstream by router with no downstream group members



Network Layer 4-145

Shared-tree: steiner tree

- ❖ *steiner tree*: minimum cost tree connecting all routers with attached group members
- ❖ problem is NP-complete
- ❖ excellent heuristics exists
- ❖ not used in practice:
 - computational complexity
 - information about entire network needed
 - monolithic: rerun whenever a router needs to join/leave

Network Layer 4-146

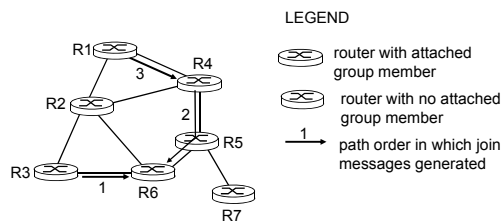
Center-based trees

- ❖ single delivery tree shared by all
- ❖ one router identified as “center” of tree
- ❖ to join:
 - edge router sends unicast *join-msg* addressed to center router
 - *join-msg* “processed” by intermediate routers and forwarded towards center
 - *join-msg* either hits existing tree branch for this center, or arrives at center
 - path taken by *join-msg* becomes new branch of tree for this router

Network Layer 4-147

Center-based trees: example

suppose R6 chosen as center:



Network Layer 4-148

Internet Multicasting Routing: DVMRP

- ❖ DVMRP: distance vector multicast routing protocol, RFC1075
- ❖ *flood and prune*: reverse path forwarding, source-based tree
 - RPF tree based on DVMRP’s own routing tables constructed by communicating DVMRP routers
 - no assumptions about underlying unicast
 - initial datagram to mcast group flooded everywhere via RPF
 - routers not wanting group: send upstream prune msgs

Network Layer 4-149

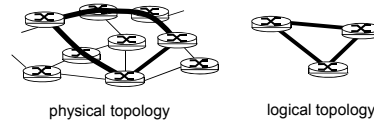
DVMRP: continued...

- ❖ *soft state*: DVMRP router periodically (1 min.) “forgets” branches are pruned:
 - mcast data again flows down unpruned branch
 - downstream router: re prune or else continue to receive data
- ❖ routers can quickly regraft to tree
 - following IGMP join at leaf
- ❖ odds and ends
 - commonly implemented in commercial router

Network Layer 4-150

Tunneling

Q: how to connect “islands” of multicast routers in a “sea” of unicast routers?



- ❖ mcast datagram encapsulated inside “normal” (non-multicast-addressed) datagram
- ❖ normal IP datagram sent thru “tunnel” via regular IP unicast to receiving mcast router (recall IPv6 inside IPv4 tunneling)
- ❖ receiving mcast router unencapsulates to get mcast datagram

Network Layer 4-151

PIM: Protocol Independent Multicast

- ❖ not dependent on any specific underlying unicast routing algorithm (works with all)
- ❖ two different multicast distribution scenarios :

dense:

- ❖ group members densely packed, in “close” proximity.
- ❖ bandwidth more plentiful

sparse:

- ❖ # networks with group members small wrt # interconnected networks
- ❖ group members “widely dispersed”
- ❖ bandwidth not plentiful

Network Layer 4-152

Consequences of sparse-dense dichotomy:

dense

- ❖ group membership by routers *assumed* until routers explicitly prune
- ❖ *data-driven* construction on mcast tree (e.g., RPF)
- ❖ bandwidth and non-group-router processing *profligate*

sparse:

- ❖ no membership until routers explicitly join
- ❖ *receiver-driven* construction of mcast tree (e.g., center-based)
- ❖ bandwidth and non-group-router processing *conservative*

Network Layer 4-153

PIM- dense mode

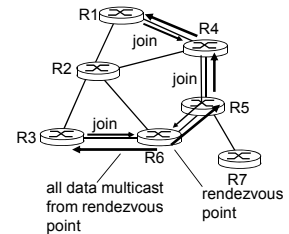
flood-and-prune RPF: similar to DVMRP but...

- ❖ underlying unicast protocol provides RPF info for incoming datagram
- ❖ less complicated (less efficient) downstream flood than DVMRP reduces reliance on underlying routing algorithm
- ❖ has protocol mechanism for router to detect it is a leaf-node router

Network Layer 4-154

PIM - sparse mode

- ❖ center-based approach
- ❖ router sends *join* msg to rendezvous point (RP)
 - intermediate routers update state and forward *join*
- ❖ after joining via RP, router can switch to source-specific tree
 - increased performance: less concentration, shorter paths

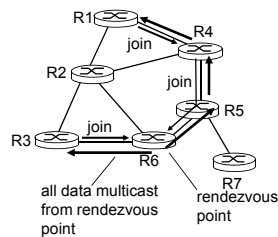


Network Layer 4-155

PIM - sparse mode

sender(s):

- ❖ unicast data to RP, which distributes down RP-rooted tree
- ❖ RP can extend mcast tree upstream to source
- ❖ RP can send *stop* msg if no attached receivers
 - "no one is listening!"



Network Layer 4-156

Chapter 4: done!

- 4.1 introduction
- 4.2 virtual circuit and datagram networks
- 4.3 what's inside a router
- 4.4 IP: Internet Protocol
 - datagram format, IPv4 addressing, ICMP, IPv6

- 4.5 routing algorithms
 - link state, distance vector, hierarchical routing
- 4.6 routing in the Internet
 - RIP, OSPF, BGP
- 4.7 broadcast and multicast routing

- ❖ understand principles behind network layer services:
 - network layer service models, forwarding versus routing
 - how a router works, routing (path selection), broadcast, multicast
- ❖ instantiation, implementation in the Internet

Network Layer 4-157