

CSC358 Intro. to Computer Networks

Lecture 9: DV, Routing in the Internet

Amir H. Chinatei, Winter 2016

ahchinaei@cs.toronto.edu
<http://www.cs.toronto.edu/~ahchinaei/>

Many slides are (inspired/adapted) from the above source
 © all material copyright; all rights reserved for the authors



Office Hours: T 17:00–18:00 R 9:00–10:00 BA4222

TA Office Hours: W 16:00-17:00 BA3201 R 10:00-11:00 BA7172

csc358ta@cdf.toronto.edu
<http://www.cs.toronto.edu/~ahchinaei/teaching/2016jan/csc358/>

Distance vector algorithm

Bellman-Ford equation (dynamic programming)

let

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

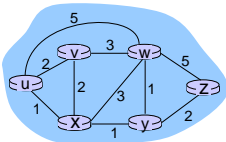
then

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

\min taken over all neighbors v of x
 $c(x,v)$ cost to neighbor v
 $d_v(y)$ cost from neighbor v to destination y

Network Layer 4-2

Bellman-Ford example



clearly, $d_v(z) = 5$, $d_x(z) = 3$, $d_w(z) = 3$

B-F equation says:

$$d_u(z) = \min \{ c(u,v) + d_v(z), c(u,x) + d_x(z), c(u,w) + d_w(z) \}$$

$$= \min \{ 2 + 5, 1 + 3, 5 + 3 \} = 4$$

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

Network Layer 4-3

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 4-4

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 4-5

Distance vector algorithm

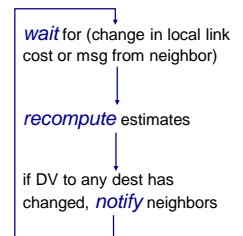
iterative, asynchronous:

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

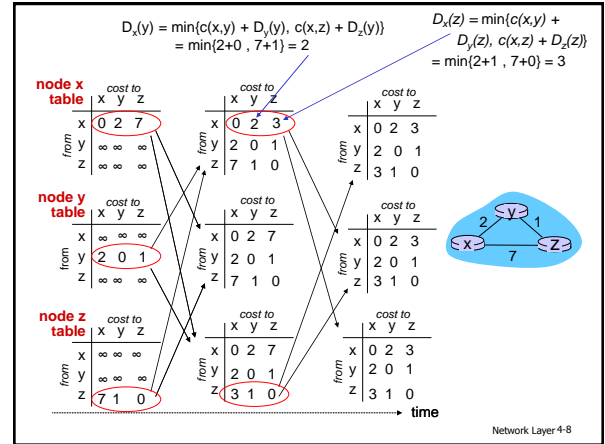
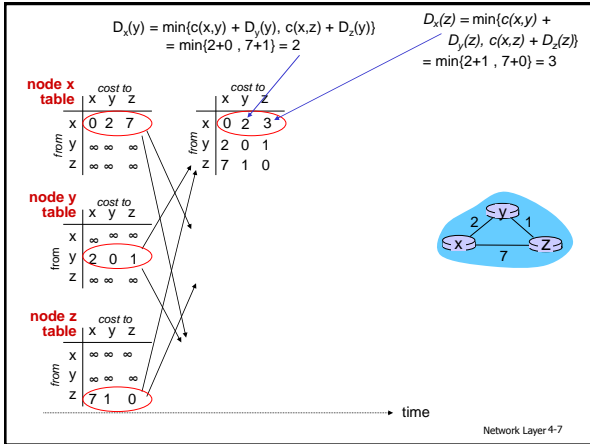
distributed:

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

each node:



Network Layer 4-6



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 4-9

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 4-10

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 4-11

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 4-12

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 4-13

Hierarchical routing

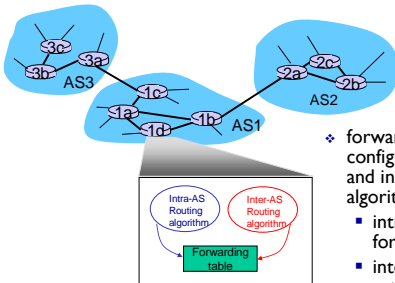
- ❖ aggregate routers into regions, "autonomous systems" (AS)
- ❖ routers in same AS run same routing protocol
 - "intra-AS" routing protocol
 - routers in different AS can run different intra-AS routing protocol

gateway router:

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

Network Layer 4-14

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 4-15

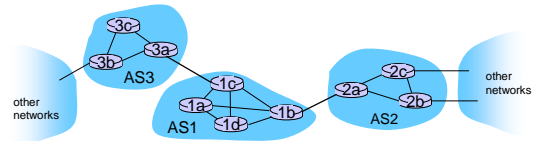
Inter-AS tasks

- ❖ 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

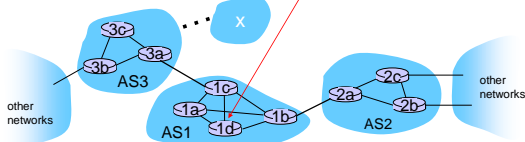
job of inter-AS routing!



Network Layer 4-16

Example: setting forwarding table in router 1d

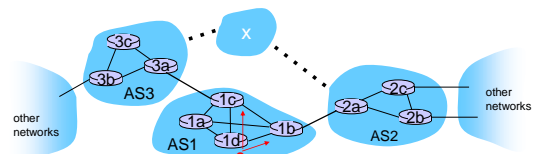
- ❖ suppose AS1 learns (via inter-AS protocol) that subnet x reachable via AS3 (gateway 1c), but not via AS2
 - inter-AS protocol propagates reachability info to all internal routers
- ❖ router 1d determines from intra-AS routing info that its interface l is on the least cost path to 1c
 - installs forwarding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



Network Layer 4-17

Example: choosing among multiple ASes

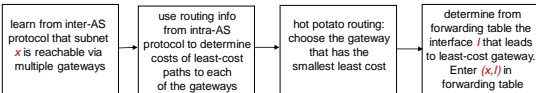
- ❖ now suppose AS1 learns from inter-AS protocol that subnet x is reachable from AS3 *and* from AS2
- ❖ to configure forwarding table, router 1d must determine which gateway it should forward packets towards for dest x
 - this is also job of inter-AS routing protocol!



Network Layer 4-18

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.



Network Layer 4-19

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 4-20

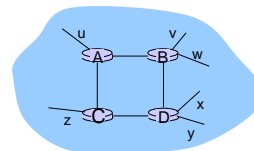
Intra-AS Routing

- ❖ 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 4-21

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)

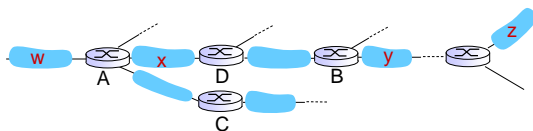


from router A to destination **subnets**:

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

Network Layer 4-22

RIP: example

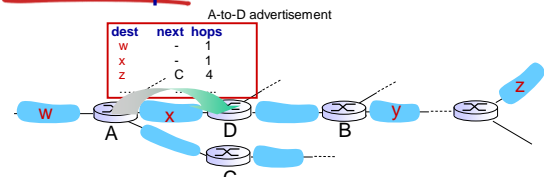


routing table in router D

destination subnet	next router	# hops to dest
w	A	2
y	B	2
z	B	7
x	--	1
....

Network Layer 4-23

RIP: example



routing table in router D

destination subnet	next router	# hops to dest
w	A	2
y	B	2
z	B → A	7 → 5
x	--	1
....

Network Layer 4-24

RIP: link failure, recovery

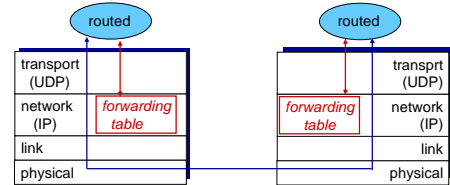
if no advertisement heard after 180 sec --> neighbor/link declared dead

- 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 4-25

RIP table processing

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



Network Layer 4-26

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

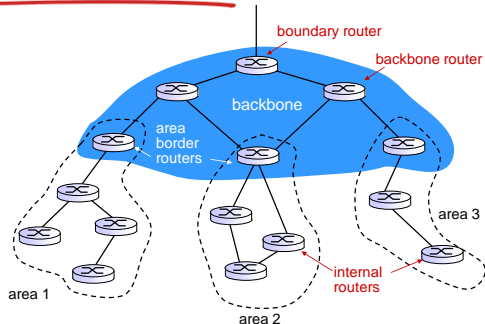
Network Layer 4-27

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 4-28

Hierarchical OSPF



Network Layer 4-29

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 4-30

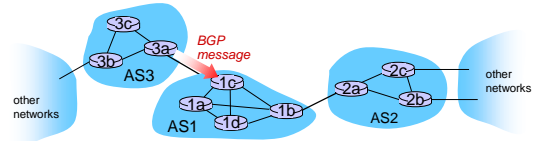
Internet inter-AS routing: BGP

- ❖ **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: “*am here*”

Network Layer 4-31

BGP basics

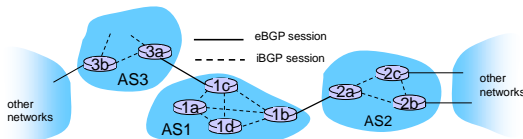
- ❖ **BGP session:** two BGP routers (“peers”) exchange BGP messages:
 - advertising *paths* to different destination network prefixes (“path vector” protocol)
 - exchanged over semi-permanent TCP connections
- ❖ when AS3 advertises a prefix to AS1:
 - AS3 *promises* it will forward datagrams towards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



Network Layer 4-32

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.



Network Layer 4-33

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 4-34

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 4-35

BGP messages

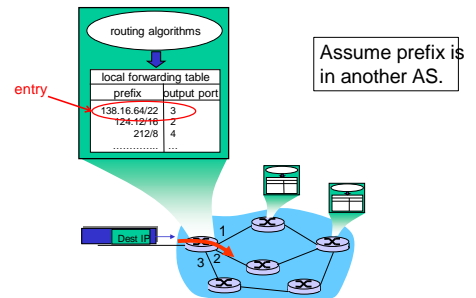
- ❖ BGP messages exchanged between peers over TCP connection
- ❖ BGP messages:
 - **OPEN:** opens TCP connection to peer and authenticates sender
 - **UPDATE:** advertises new path (or withdraws old)
 - **KEEPALIVE:** keeps connection alive in absence of UPDATES; also ACKs OPEN request
 - **NOTIFICATION:**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

Network Layer 4-36

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?

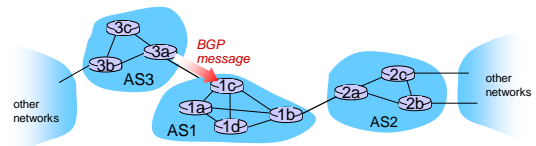


How does entry get in forwarding table?

High-level overview

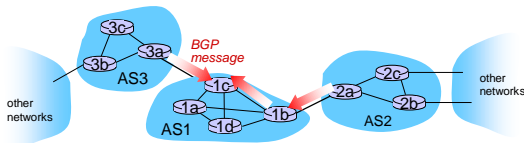
1. Router becomes aware of prefix
2. Router determines output port for prefix
3. Router enters prefix-port in forwarding table

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

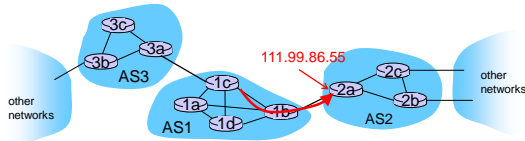
Example:

- ❖ 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!

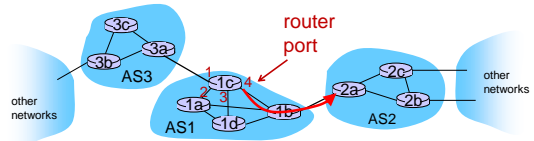
Find best intra-route to BGP route

- ❖ 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



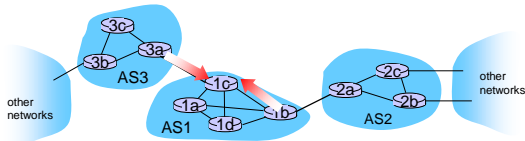
Router identifies port for route

- ❖ 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 AS131 since it is closer

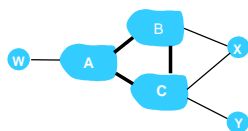


How does entry get in forwarding table?

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

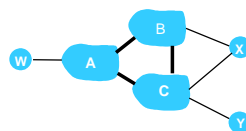


legend: provider network
 customer network:

- ❖ 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 4-47

BGP routing policy (2)



legend: provider network
 customer network:

- ❖ 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!

Network Layer 4-48

Why different Intra-, Inter-AS routing ?

policy:

- ❖ inter-AS: admins want 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 4-49

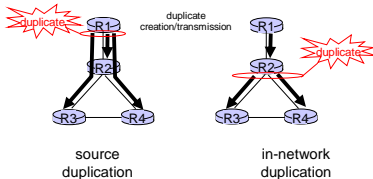
Chapter 4: outline

- | | |
|--|---|
| 4.1 introduction | 4.5 routing algorithms |
| 4.2 virtual circuit and datagram networks | <ul style="list-style-type: none"> ▪ link state ▪ distance vector ▪ hierarchical routing |
| 4.3 what's inside a router | 4.6 routing in the Internet |
| 4.4 IP: Internet Protocol | <ul style="list-style-type: none"> ▪ RIP ▪ OSPF ▪ BGP |
| <ul style="list-style-type: none"> ▪ datagram format ▪ IPv4 addressing ▪ ICMP ▪ IPv6 | 4.7 broadcast and multicast routing |

Network Layer 4-50

Broadcast routing

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



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

Network Layer 4-51

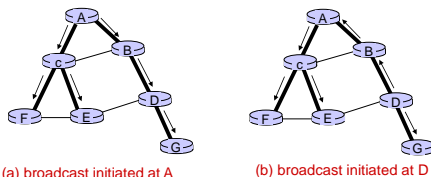
In-network duplication

- ❖ **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 4-52

Spanning tree

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



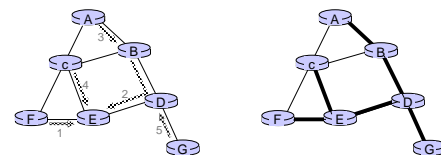
(a) broadcast initiated at A

(b) broadcast initiated at D

Network Layer 4-53

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



(a) stepwise construction of spanning tree (center: E)

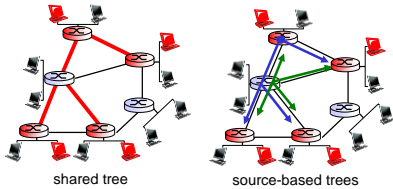
(b) constructed spanning tree

Network Layer 4-54

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 4-55

Approaches for building mcast trees

approaches:

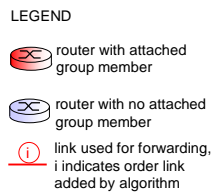
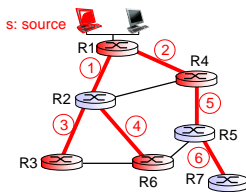
- ❖ **source-based tree:** one tree per source
 - shortest path trees
 - reverse path forwarding
- ❖ **group-shared tree:** group uses one tree
 - minimal spanning (Steiner)
 - center-based trees

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

Network Layer 4-56

Shortest path tree

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



Network Layer 4-57

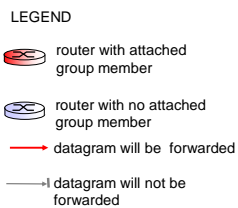
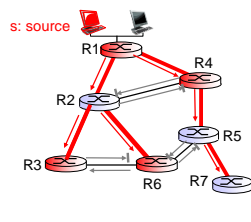
Reverse path forwarding

- ❖ 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 4-58

Reverse path forwarding: example

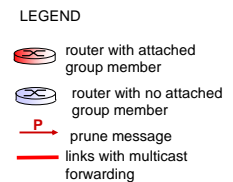
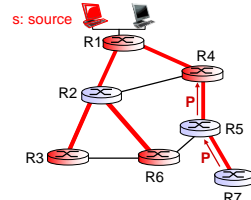


- ❖ result is a source-specific reverse SPT
 - may be a bad choice with asymmetric links

Network Layer 4-59

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 4-60

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 4-61

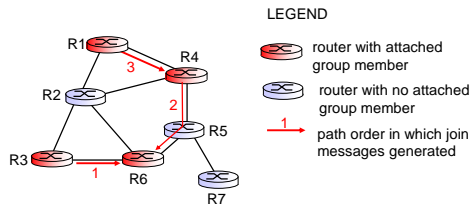
Center-based trees

- ❖ 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 4-62

Center-based trees: example

suppose R6 chosen as center:



Network Layer 4-63