

Lecture 23

- ❖ Sections 5.2.2, 5.3, 5.4
- ❖ Routing protocols
 - Distance vector routing protocol
- ❖ Intra-Autonomous Systems routing (OSPF)
- ❖ Routing among the ISPs: BGP

Network Layer Control Plane 5-25

25

Distance Vector Algorithm (5)

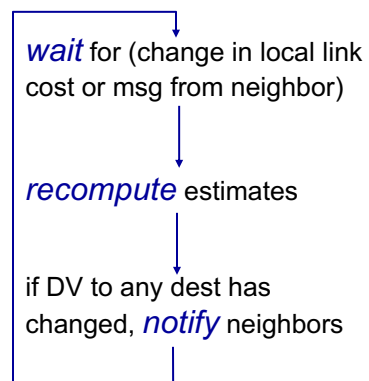
Distributed:

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

Iterative, asynchronous:

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

Each node:

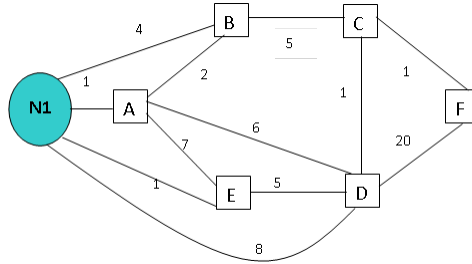


Network Layer Control Plane 4-26

26

Distance Vector Routing: Example

(X,i):
X: next hop to N1
i: distance to N1



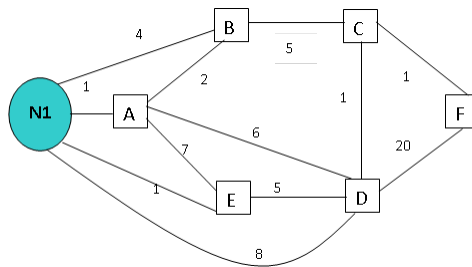
Iter.	A	B	C	D	E	F
0	(-,1)	(-,4)	(•,∞)	(-,8)	(-,1)	(•,∞)

Network Layer Control Plane 4-27

27

Distance Vector Routing: Example

(X,i):
X: next hop to N1
i: distance to N1



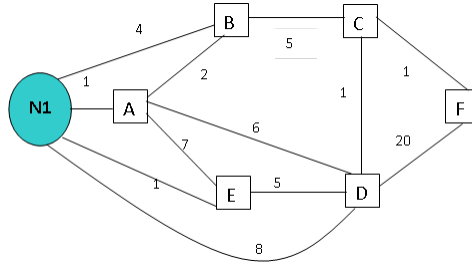
Iter.	A	B	C	D	E	F
0	(-,1)	(-,4)	(•,∞)	(-,8)	(-,1)	(•,∞)
1	(-,1)	(A,3)	(B,9)	(E,6)	(-,1)	(D,28)

Network Layer Control Plane 4-28

28

Distance Vector Routing: Example

(X,i):
X: next hop to N1
i: distance to N1



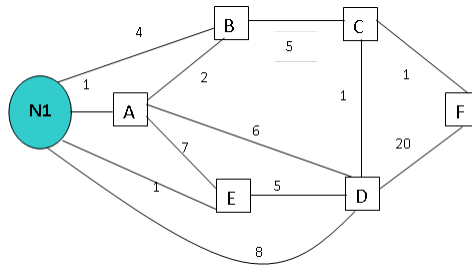
Iter.	A	B	C	D	E	F
0	(-,1)	(-,4)	(•,∞)	(-,8)	(-,1)	(•,∞)
1	(-,1)	(A,3)	(B,9)	(E,6)	(-,1)	(D,28)

Network Layer Control Plane 4-29

29

Distance Vector Routing: Example

(X,i):
X: next hop to N1
i: distance to N1



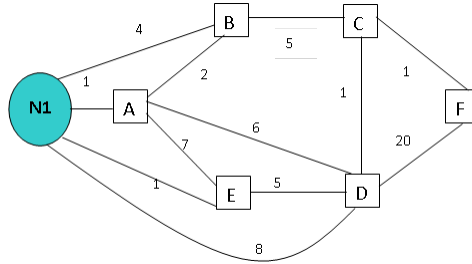
Iter.	A	B	C	D	E	F
0	(-,1)	(-,4)	(•,∞)	(-,8)	(-,1)	(•,∞)
1	(-,1)	(A,3)	(B,9)	(E,6)	(-,1)	(D,28)

Network Layer Control Plane 4-30

30

Distance Vector Routing: Example

(X,i):
X: next hop to N1
i: distance to N1



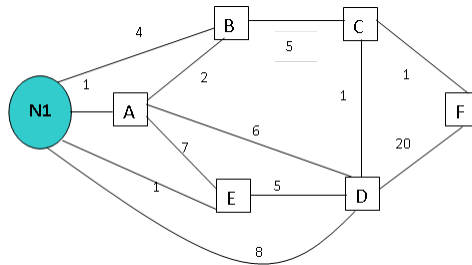
Iter.	A	B	C	D	E	F
0	(-,1)	(-,4)	(•,∞)	(-,8)	(-,1)	(•,∞)
1	(-,1)	(A,3)	(B,9)	(E,6)	(-,1)	(D,28)
2	(-,1)	(A,3)	(D,7)	(E,6)	(-,1)	(C,10)

Network Layer Control Plane 4-31

31

Distance Vector Routing: Example

(X,i):
X: next hop to N1
i: distance to N1



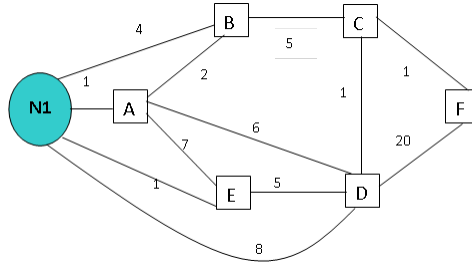
Iter.	A	B	C	D	E	F
0	(-,1)	(-,4)	(•,∞)	(-,8)	(-,1)	(•,∞)
1	(-,1)	(A,3)	(B,9)	(E,6)	(-,1)	(D,28)
2	(-,1)	(A,3)	(D,7)	(E,6)	(-,1)	(C,10)

Network Layer Control Plane 4-32

32

Distance Vector Routing: Example

(X,i):
X: next hop to N1
i: distance to N1



Iter.	A	B	C	D	E	F
0	(-,1)	(-,4)	(·,∞)	(-,8)	(-,1)	(·,∞)
1	(-,1)	(A,3)	(B,9)	(E,6)	(-,1)	(D,28)
2	(-,1)	(A,3)	(D,7)	(E,6)	(-,1)	(C,10)
3	(-,1)	(A,3)	(D,7)	(E,6)	(-,1)	(C,8)

Network Layer Control Plane 4-33

33

$$D_x(y) = \min\{c(x,y) + D_y(y), c(x,z) + D_z(y)\}$$

$$= \min\{2+0, 7+1\} = 2$$

$$D_x(z) = \min\{c(x,y) + D_y(z), c(x,z) + D_z(z)\}$$

$$= \min\{2+1, 7+0\} = 3$$

node x table

cost to			
	x	y	z
from x	0	2	7
from y	∞	∞	∞
from z	∞	∞	∞

node y table

cost to			
	x	y	z
from x	∞	∞	∞
from y	2	0	1
from z	∞	∞	∞

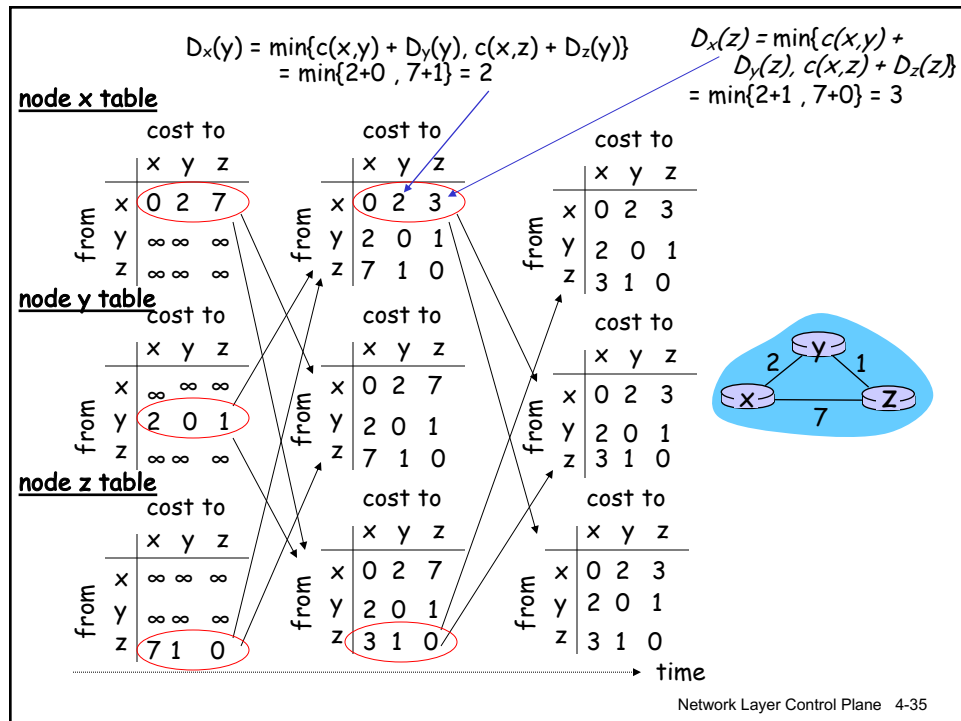
node z table

cost to			
	x	y	z
from x	∞	∞	∞
from y	∞	∞	∞
from z	7	1	0

time →

Network Layer Control Plane 4-34

34

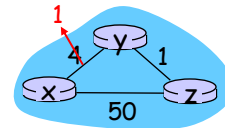


35

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.

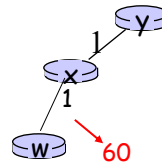
Network Layer Control Plane 4-36

36

Distance Vector: link cost changes

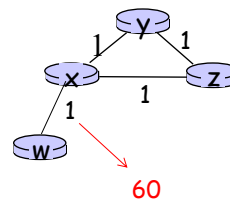
Link cost changes:

- ❖ **bad news travels slow** - "count to infinity" problem!
- ❖ Many iterations before algorithm stabilizes: how many?



Poisoned reverse:

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



Network Layer Control Plane 4-37

37

Comparison of LS and DV algorithms

Message complexity

- ❖ **LS:** with n nodes, E links, $O(nE)$ msgs sent
- ❖ **DV:** exchange between neighbors only

Speed of Convergence

- ❖ **LS:** exactly n iterations
- ❖ **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 through network

Network Layer Control Plane 4-38

38

Chapter 5: Network Layer

5.1 Introduction

5.2 Routing algorithms

5.2.1 link-state routing algorithm

5.2.2 distance vector routing

5.3 Intra-Autonomous System (AS) routing: OSPF

Network Layer Control Plane 4-39

39

Internet approach to scalable routing

aggregate routers into regions known as "autonomous systems" (AS) (a.k.a. "domains")

intra-AS (aka "intra-domain"): routing among *within same AS* ("network")

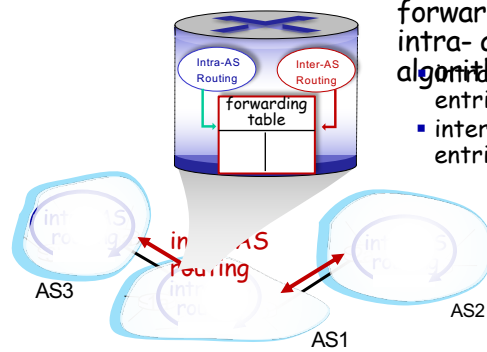
- all routers in AS must run same intra-domain protocol
- routers in different AS can run different intra-domain routing protocols
- **gateway router**: at "edge" of its own AS, has link(s) to router(s) in other AS'es

inter-AS (aka "inter-domain"): routing *among AS'es*

- gateways perform inter-domain routing (as well as intra-domain routing)

40

Interconnected ASes



forwarding table configured by intra- and inter-AS routing algorithms

- inter-AS routing determine entries for destinations within AS
- inter-AS & intra-AS determine entries for external destinations

41

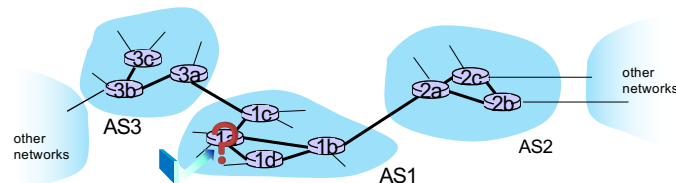
Inter-AS routing: a role in intradomain forwarding

- suppose router in AS1 receives datagram destined outside of AS1:

? • router should forward packet to gateway router in AS1, but which one?

AS1 inter-domain routing must:

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



42

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 4-43

43

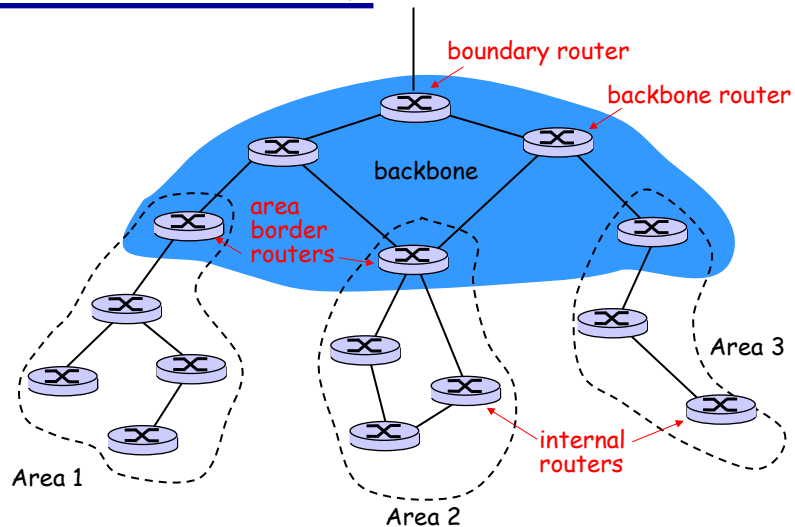
OSPF "advanced" features (not in RIP)

- ❖ for each link, multiple cost metrics for different **TOS** (hop, delay, cost, throughput, and reliability) -> 5 routing trees can be constructed
- ❖ integrated uni- and **multicast** support:
 - Multicast OSPF (MOSPF) uses same topology data base as OSPF
- ❖ **hierarchical** OSPF in large domains.

Network Layer 4-44

44

Hierarchical OSPF



Network Layer 4-45

45

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 4-46

46

Chapter 5: Network Layer

5.1 Introduction

5.2 Routing algorithms

5.2.1 link-state routing algorithm

5.2.2 distance vector routing

5.3 Intra-Autonomous System (AS) routing: OSPF

5.4 Routing among the ISPs

Network Layer Control Plane 4-47

47

Internet inter-AS routing: BGP

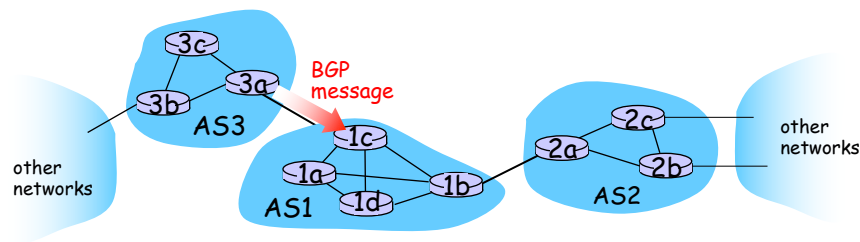
- ❖ **BGP (Border Gateway Protocol):** *the de facto* inter-domain routing protocol
 - “glue that holds the Internet together”
- ❖ BGP provides each AS a means to:
 - **eBGP:** obtain subnet reachability information from neighboring ASs.
 - **iBGP:** propagate reachability information to all AS-internal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

Network Layer 4-48

48

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

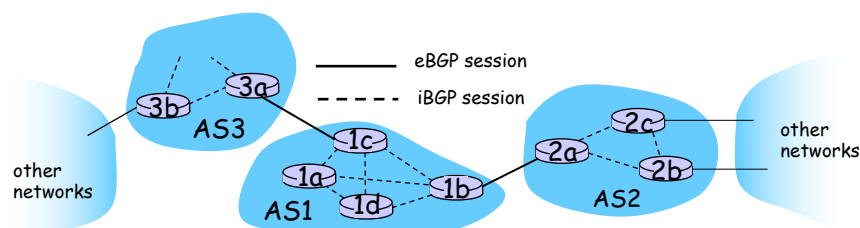


Network Layer 4-49

49

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 4-50

50

Path attributes & 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 4-51

51

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 4-52

52

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 4-53

53

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

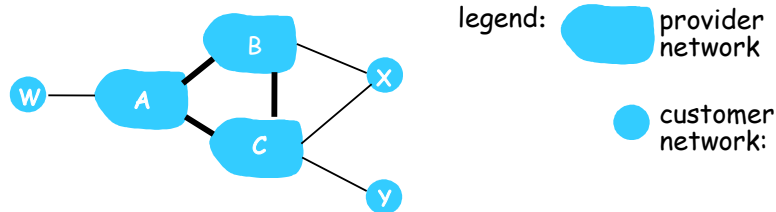
Performance:

- ❖ Intra-AS: can focus on performance
- ❖ Inter-AS: policy may dominate over performance

Network Layer 4-54

54

BGP routing policy

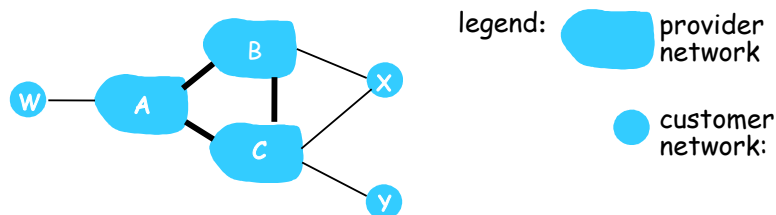


- ❖ A,B,C are **provider networks**
- ❖ X,W,Y are customer (of provider networks)
- ❖ X is **dual-homed**: attached to two networks
 - X does not want to route from B via X to C
 - .. so X will not advertise to B a route to C

Network Layer 4-55

55

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!

Network Layer 4-56

56