

Network Layer: Control Plane

EECS3214

18-03-05

© All material copyright 1996-2016
J.F. Kurose and K.W. Ross, All Rights Reserved

4-1

Chapter 5: network layer control plane

chapter goals: understand principles behind network control plane

- traditional routing algorithms
- SDN controllers
- Internet Control Message Protocol
- network management

and their instantiation, implementation in the Internet:

- OSPF, BGP, OpenFlow, ODL and ONOS controllers, ICMP, SNMP

Network Layer: Control Plane 5-2

Chapter 5: outline

5.1 introduction

5.2 routing protocols

- link state
- distance vector

5.3 intra-AS routing in the Internet: OSPF

5.4 routing among the ISPs: BGP

5.5 The SDN control plane

5.6 ICMP: The Internet Control Message Protocol

5.7 Network management and SNMP

Network Layer: Control Plane 5-3

Network-layer functions

Recall: two network-layer functions:

- *forwarding*: move packets from router's input to appropriate router's output *data plane*
- *routing*: determine route taken by packets from source to destination *control plane*

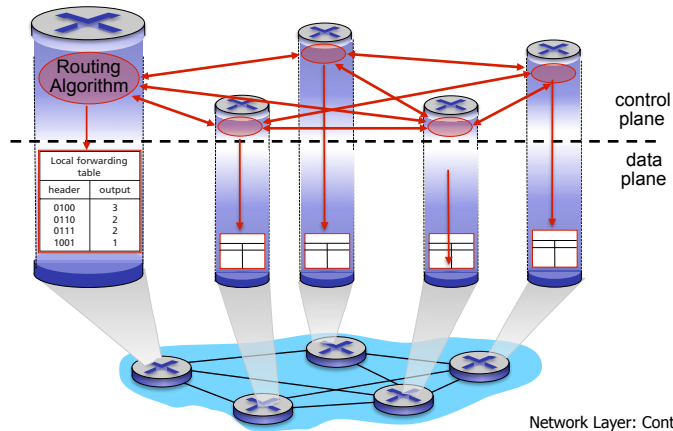
Two approaches to structuring network control plane:

- per-router control (traditional)
- logically centralized control (software defined networking)

Network Layer: Control Plane 5-4

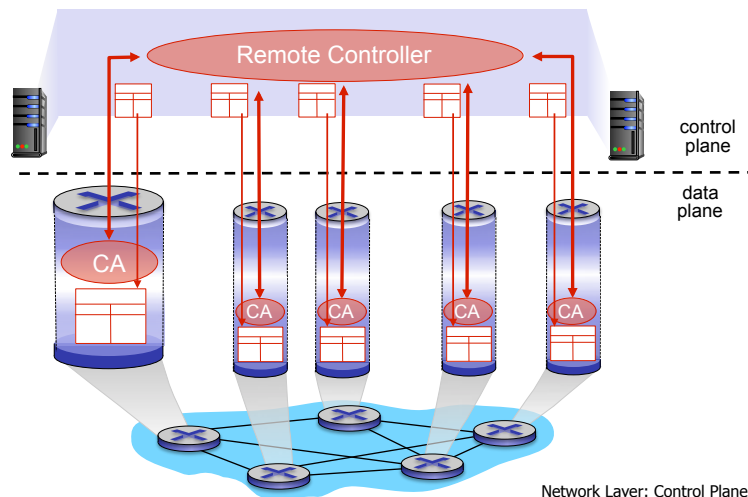
Per-router control plane

Individual routing algorithm components *in each and every router* interact with each other in control plane to compute forwarding tables



Logically centralized control plane

A distinct (typically remote) controller interacts with local control agents (CAs) in routers to compute forwarding tables



Chapter 5: outline

5.1 introduction

5.2 routing protocols

- link state
- distance vector

5.3 intra-AS routing in the Internet: OSPF

5.4 routing among the ISPs: BGP

5.5 The SDN control plane

5.6 ICMP: The Internet Control Message Protocol

5.7 Network management and SNMP

Network Layer: Control Plane 5-7

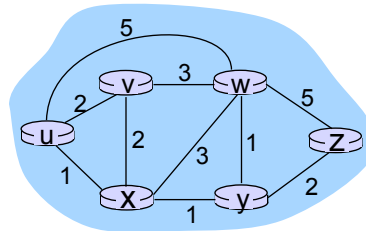
Routing protocols

Routing protocol goal: determine “good” paths (equivalently, routes), from sending hosts to receiving host, through network of routers

- path: sequence of routers packets will traverse in going from given initial source host to given final destination host
- “good”: least “cost”, “fastest”, “least congested”
- routing: a “top-10” networking challenge!

Network Layer: Control Plane 5-8

Graph abstraction of the network



graph: $G = (V, E)$

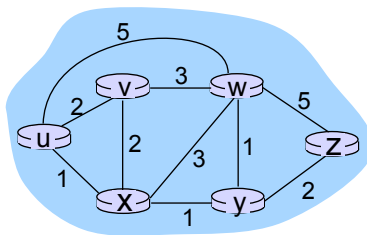
V = set of routers = $\{u, v, w, x, y, z\}$

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

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

Network Layer: Control Plane 5-9

Graph abstraction: costs



$c(x, x') = \text{cost of link } (x, x')$
e.g., $c(w, z) = 5$

cost could always be one (hop),
or inversely related to bandwidth,
or inversely related to congestion

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

key question: what is the least-cost path between u and z ?

routing algorithm: algorithm that finds that least cost path

Network Layer: Control Plane 5-10

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 updates
 - in response to link cost changes
 - more responsive but
 - more route oscillation; routing loops

Network Layer: Control Plane 5-11

Chapter 5: outline

5.1 introduction

5.2 routing protocols

- link state
- distance vector

5.3 intra-AS routing in the Internet: OSPF

5.4 routing among the ISPs: BGP

5.5 The SDN control plane

5.6 ICMP: The Internet Control Message Protocol

5.7 Network management and SNMP

Network Layer: Control Plane 5-12

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 paths to k destinations

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

Network Layer: Control Plane 5-13

Dijkstra's algorithm

1 **Initialization:**

- $N' = \{u\}$
- for all nodes v
- if v adjacent to u
- then $D(v) = c(u,v)$
- else $D(v) = \infty$

7

8 **Loop**

- find w not in N' such that $D(w)$ is a minimum
- add w to N'
- update $D(v)$ for all v adjacent to w and not in N' :
 $D(v) = \min(D(v), D(w) + c(w,v))$
- /* new cost to v is either old cost to v or known shortest path cost to w plus cost from w to v */
- until all nodes in N'**

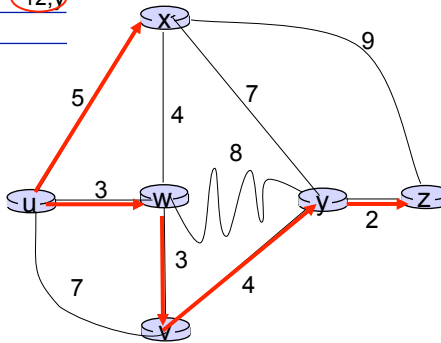
Network Layer: Control Plane 5-14

Dijkstra's algorithm: example

Step	N'	D(v)	D(w)	D(x)	D(y)	D(z)
		p(v)	p(w)	p(x)	p(y)	p(z)
0	u	7,u	3,u	5,u	∞	∞
1	uw	6,w		5,u	11,w	∞
2	uw x	6,w			11,w	14,x
3	uw x v				10,y	14,x
4	uw x v y					12,y
5	uw x v y z					

notes:

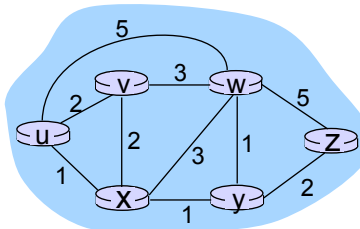
- ❖ construct shortest path tree by tracing predecessor nodes
- ❖ ties can exist (can be broken arbitrarily)



Network Layer: Control Plane 5-15

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

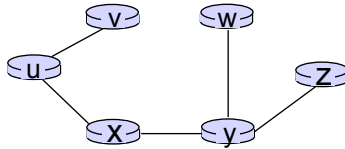


* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

Network Layer: Control Plane 5-16

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: Control Plane 5-17

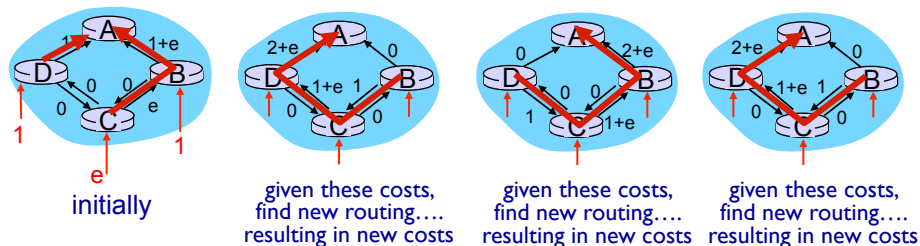
Dijkstra's algorithm, discussion

algorithm complexity: n nodes

- each iteration: need to check all nodes, w , not in V
- $n(n+1)/2$ comparisons: $O(n^2)$
- more efficient implementation using a heap: $O(n \log n)$

oscillations possible:

- e.g., support link cost equals amount of carried traffic



Network Layer: Control Plane 5-18

Chapter 5: outline

5.1 introduction

5.2 routing protocols

- link state
- distance vector

5.3 intra-AS routing in the Internet: OSPF

5.4 routing among the ISPs: BGP

5.5 The SDN control plane

5.6 ICMP: The Internet Control Message Protocol

5.7 Network management and SNMP

Network Layer: Control Plane 5-19

Distance vector algorithm

Bellman-Ford equation (dynamic programming)

Let

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

then

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

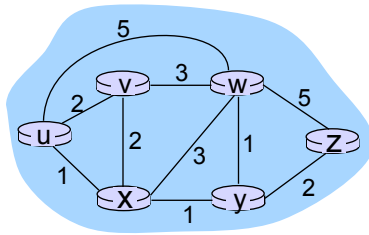
\min taken over all neighbors v of x

$c(x,v)$ cost to neighbor v

$d_v(y)$ cost from neighbor v to destination y

Network Layer: Control Plane 5-20

Bellman-Ford example



Compute cost from u to z:
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: Control Plane 5-21

Distance vector algorithm: variables

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

Network Layer: Control Plane 5-22

Distance vector algorithm: key idea

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: Control Plane 5-23

Distance vector algorithm: properties

iterative, asynchronous:

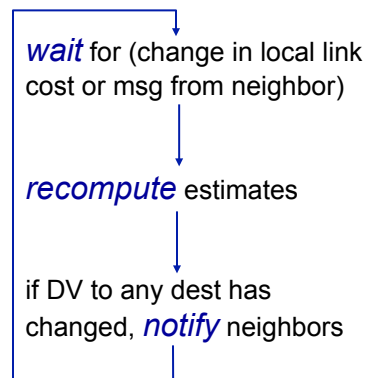
each local iteration caused by:

- local link cost change
- DV update message from neighbor

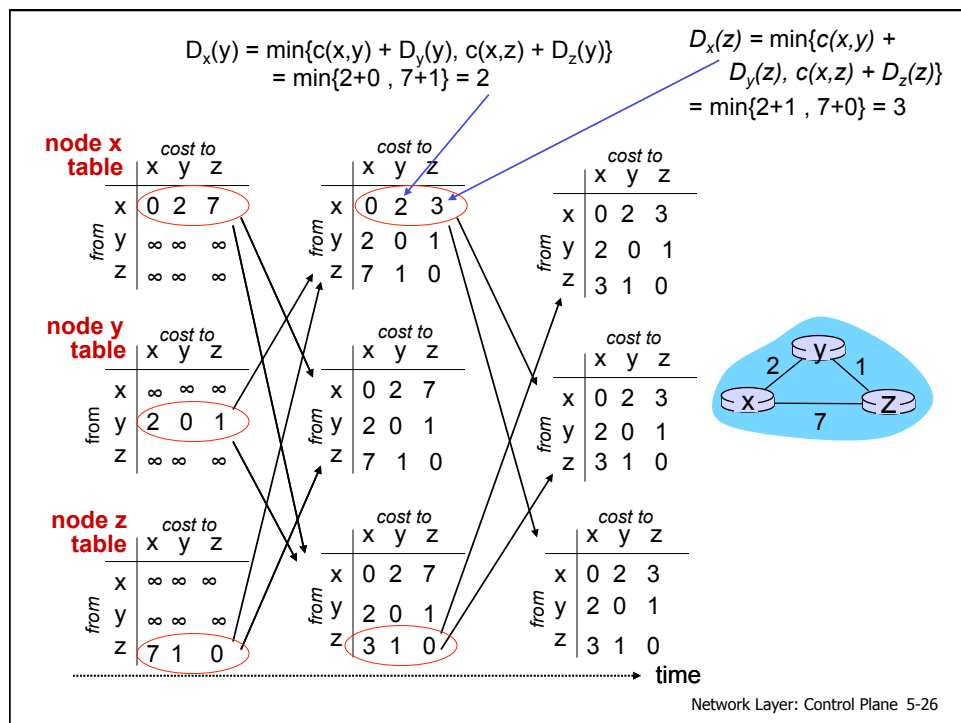
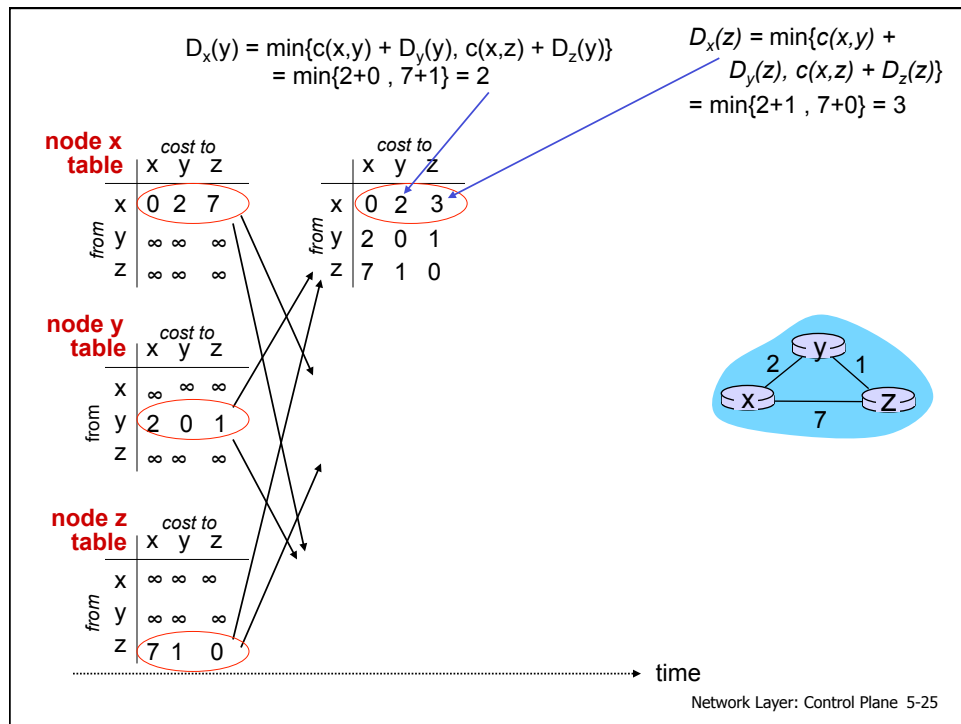
distributed:

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

each node:



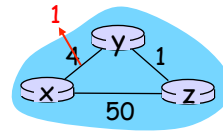
Network Layer: Control Plane 5-24



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.

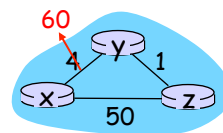
* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

Network Layer: Control Plane 5-27

Distance vector: link cost changes (2)

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: Control Plane 5-28

Comparison of LS and DV algorithms

message complexity

- **LS:** with n nodes, E links, $O(nE)$ messages sent
- **DV:** exchange between neighbors only
 - convergence time varies

speed of convergence

- **LS:** $O(n^2)$ algorithm requires $O(nE)$ messages
 - may have oscillations
- **DV:** convergence time varies
 - may cause 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
 - errors propagate through network

Network Layer: Control Plane 5-29

Chapter 5: outline

5.1 introduction

5.2 routing protocols

- link state
- distance vector

5.3 intra-AS routing in the Internet: OSPF

5.4 routing among the ISPs: BGP

5.5 The SDN control plane

5.6 ICMP: The Internet Control Message Protocol

5.7 Network management and SNMP

Network Layer: Control Plane 5-30

Making routing scalable

our routing study thus far - idealized

- all routers identical
- network “flat”

... *not* true in practice

scale: with billions of destinations:

- can't store all destinations in routing tables!
- routing table exchange would swamp links!

administrative autonomy

- internet = network of networks
- each network administrator may want to control routing in their own network

Network Layer: Control Plane 5-31

Internet approach to scalable routing

Aggregate routers into regions known as “*autonomous systems*” (AS) (a.k.a. “domains”)

intra-AS routing

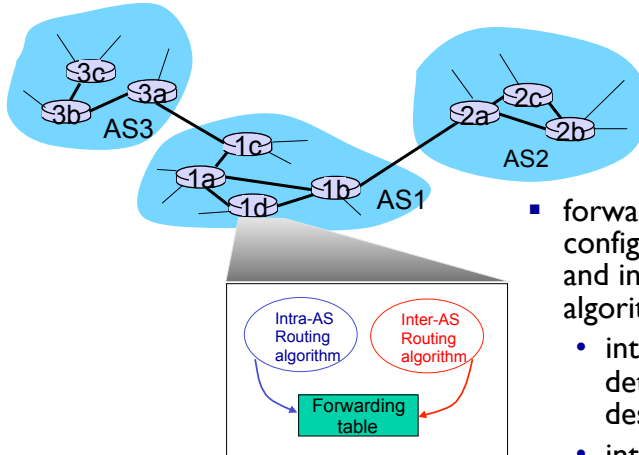
- routing among hosts, routers in same AS (“network”)
- all routers in AS must run *same* intra-domain protocol
- routers in *different* AS can run *different* intra-domain routing protocols
- gateway router: at “edge” of its own AS, has link(s) to router(s) in other AS'es

inter-AS routing

- routing among AS'es
- gateways perform inter-domain routing (as well as intra-domain routing)

Network Layer: Control Plane 5-32

Interconnected ASes



- forwarding table configured by both intra- and inter-AS routing algorithms
 - intra-AS routing determine entries for destinations within AS
 - inter-AS and intra-AS determine entries for external destinations

Network Layer: Control Plane 5-33

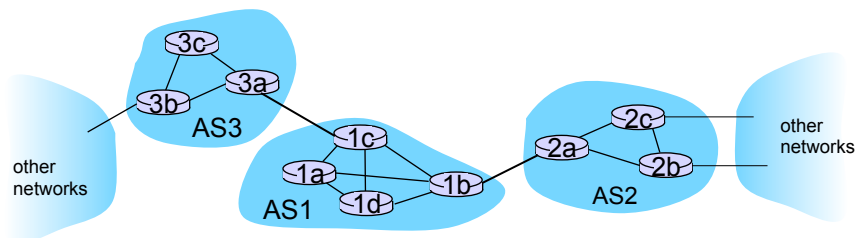
Inter-AS tasks

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

job of inter-AS routing!



Network Layer: Control Plane 5-34

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 (IS-IS protocol essentially same as OSPF)
 - IS: Intermediate System
 - IGRP: Interior Gateway Routing Protocol (Cisco proprietary for decades, until 2016)

Network Layer: Control Plane 5-35

OSPF (Open Shortest Path First)

- “open”: publicly available
- uses link-state algorithm
 - link state packet dissemination
 - topology map at each node
 - route computation using Dijkstra’s algorithm
- router floods OSPF link-state advertisements to all other routers in *entire* AS
 - link state: for each attached link
 - carried in OSPF messages directly over IP (rather than TCP or UDP)
- *IS-IS routing* protocol: nearly identical to OSPF

Network Layer: Control Plane 5-36

OSPF: More Details

- Carries OSPF messages directly over IP
 - protocol 89
 - supports reliable message transfer
 - supports link state broadcast
- Checks if links are operational
 - sending HELLO messages
- Obtains link-state database from neighbor routers
- OSPF messages are sent
 - if link states change
 - periodically (30 minutes)

Network Layer

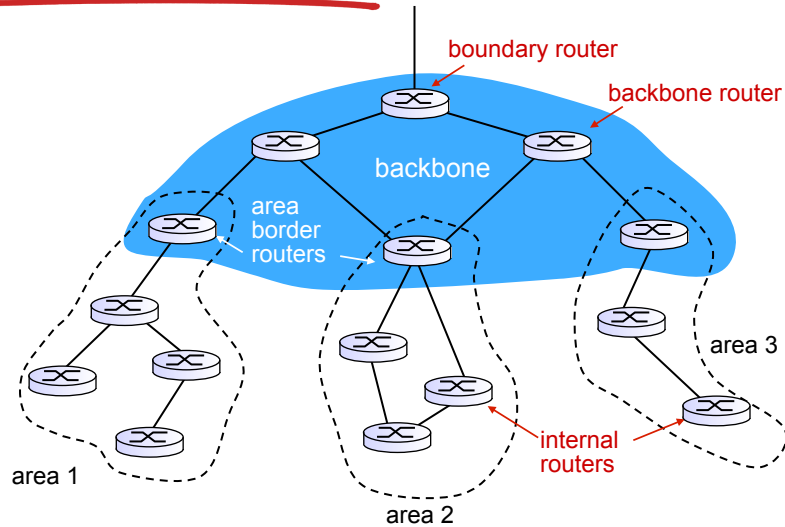
4-37

OSPF “advanced” features

- **security**: all OSPF messages authenticated (to prevent malicious intrusion)
- **multiple** same-cost **paths** allowed (only one path allowed in RIP)
- integrated unicast and **multicast** support:
 - Multicast OSPF (**MOSPF**, RFC 1584) uses same topology database as OSPF
- **hierarchical** OSPF in large domains

Network Layer: Control Plane 5-38

Hierarchical OSPF



Network Layer: Control Plane 5-39

Hierarchical OSPF

- **two-level hierarchy:** local area, backbone.
 - link-state advertisements only in area
 - each nodes has detailed area topology, but only knows direction (shortest path) to networks in other areas.
- **area border routers:** “summarize” distances to networks in its own area; advertise this info to other area border routers.
- **backbone routers:** run OSPF routing within the backbone.
- **boundary routers:** connect to other AS' es.

Network Layer: Control Plane 5-40

Hierarchical OSPF Routing

Packet routing within a domain:

- local router → area border router (intra-area)
- → backbone router(s)
- → area border router of destination network
- → destination local router

Inter-AS routing:

- local router → area border router (intra-area)
- → backbone router(s)
- → boundary router
- → another AS ...

Network Layer

4-41

Chapter 5: outline

5.1 introduction

5.2 routing protocols

- link state
- distance vector

5.3 intra-AS routing in the Internet: OSPF

5.4 routing among the ISPs:
BGP

5.5 The SDN control plane

5.6 ICMP: The Internet
Control Message
Protocol

5.7 Network management
and SNMP

Network Layer: Control Plane 5-42

Internet inter-AS routing: BGP

- **BGP (Border Gateway Protocol)**(RFC 4271): *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 AS'es (gateway routers)
 - **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 the rest of Internet: *“I am here”*

Network Layer: Control Plane 5-43

BGP Forwarding Tables

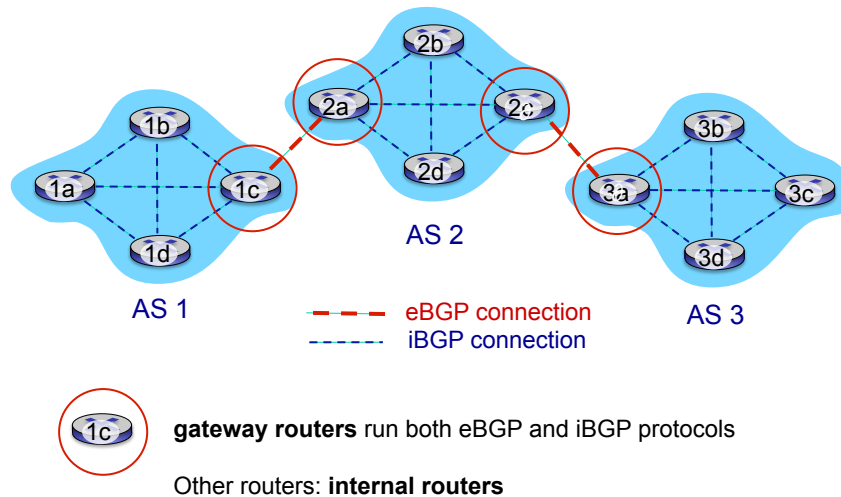
- Packets routed using CIDRized prefixes, e.g., 138.16.68/22
- Forwarding table entries have the form (*prefix, interface*)

Prefix Match	Interface
00	0
010	1
011	2
10	2
11	3

Network Layer

4-44

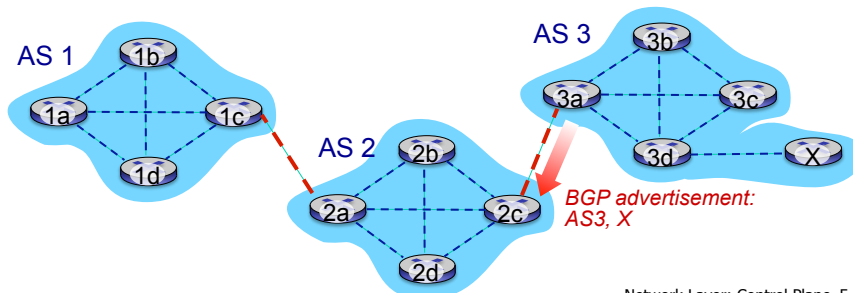
eBGP, iBGP connections



Network Layer: Control Plane 5-45

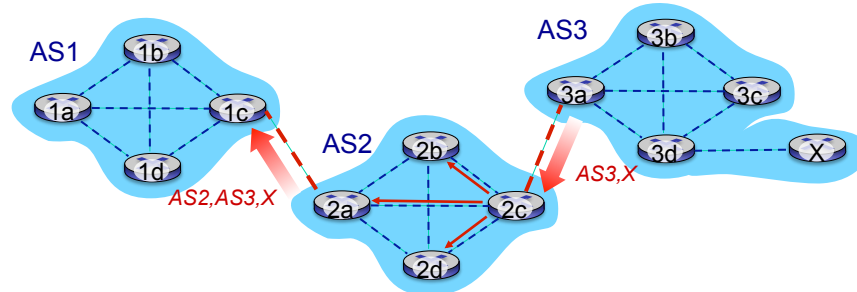
BGP basics

- **BGP connection:** two BGP routers (“peers”) exchange BGP messages over a semi-permanent TCP connection (port 179):
 - advertising *paths* to different destination network *prefixes* (BGP is a “path vector” protocol)
- when AS3 gateway router 3a advertises path **AS3,X** to AS2 gateway router 2c:
 - AS3 *promises* to AS2 it will forward datagrams towards X



Network Layer: Control Plane 5-46

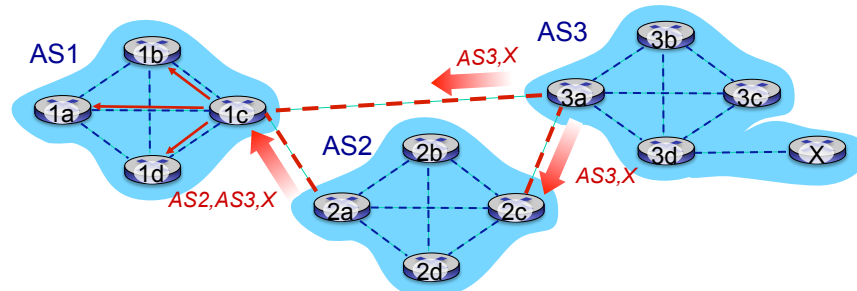
BGP path advertisement



- AS2 router 2c receives path advertisement **AS3,X** (via eBGP) from AS3 router 3a
- Based on AS2 policy, AS2 router 2c accepts path AS3,X, propagates (via iBGP) to all AS2 routers
- Based on AS2 policy, AS2 router 2a advertises (via eBGP) path **AS2,AS3,X** to AS1 router 1c

Network Layer: Control Plane 5-47

BGP path advertisement (2)



gateway router may learn about **multiple** paths to destination:

- AS1 gateway router 1c learns path **AS2,AS3,X** from 2a
- AS1 gateway router 1c learns path **AS3,X** from 3a
- Based on policy, AS1 gateway router 1c chooses path **AS3,X**, and *advertises path within AS1 via iBGP*

Network Layer: Control Plane 5-48

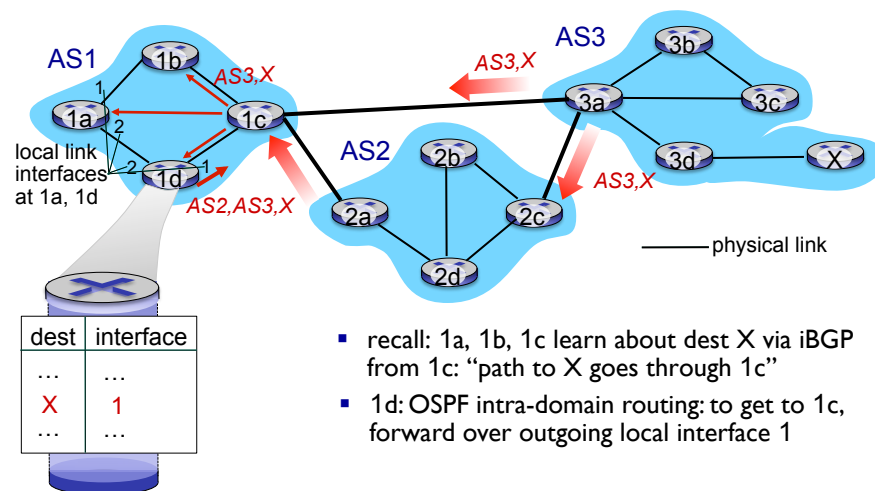
Path attributes and BGP routes

- advertised prefix includes BGP attributes
 - prefix + attributes = “route”
- two important attributes:
 - **AS-PATH**: list of ASes through which prefix advertisement has passed
 - **NEXT-HOP**: IP addresses of the router interface that begins the AS-PATH

Network Layer: Control Plane 5-49

BGP, OSPF, forwarding table entries

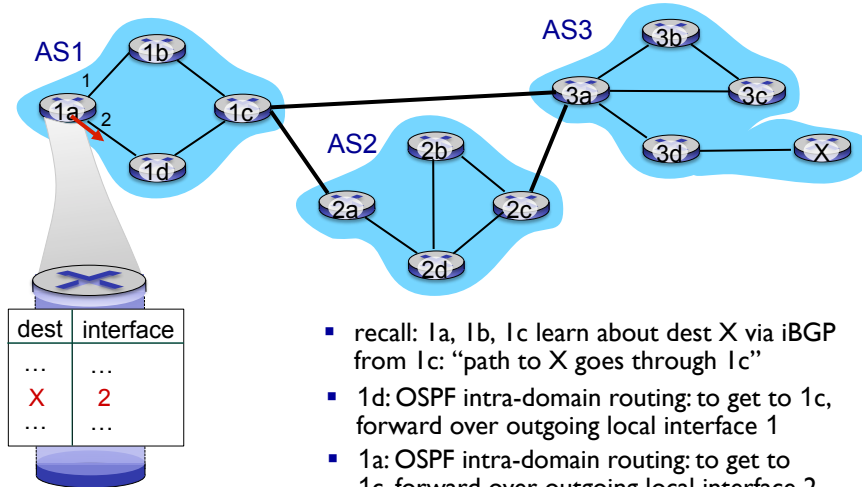
Q: how does router set forwarding table entry to distant prefix?



Network Layer: Control Plane 5-50

BGP, OSPF, forwarding table entries (2)

Q: how does router set forwarding table entry to distant prefix?



- recall: 1a, 1b, 1c learn about dest X via iBGP from 1c: "path to X goes through 1c"
- 1d: OSPF intra-domain routing: to get to 1c, forward over outgoing local interface 1
- 1a: OSPF intra-domain routing: to get to 1c, forward over outgoing local interface 2

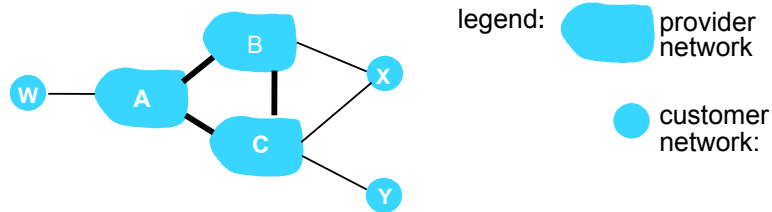
Network Layer: Control Plane 5-51

BGP route selection

- router may learn about more than one route to destination AS, selects route based on the following elimination rules:
 1. local preference value attribute: policy decision
 2. shortest AS-PATH
 3. closest NEXT-HOP router: hot potato routing
 4. additional criteria
- **Policy-based routing:**
 - gateway receiving route advertisement uses **import policy** to accept/decline path (e.g., never route through AS Y).
 - AS policy also determines whether to **advertise** path to other neighboring ASes

Network Layer: Control Plane 5-52

BGP: achieving policy via advertisements

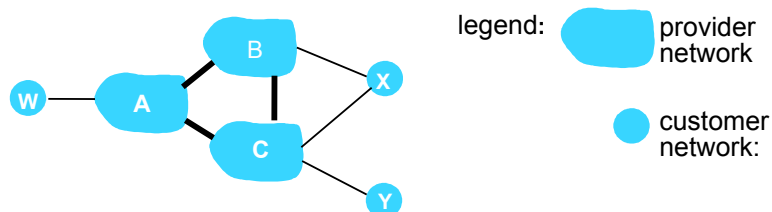


Suppose an ISP only wants to route traffic to/from its customer networks (does not want to carry transit traffic between other ISPs)

- A advertises path A_w to B and to C
- B *chooses not to advertise* B_{A_w} to C:
 - B gets no “revenue” for routing C_{B_{A_w}}, since none of C, A, w are B’s customers
 - C does not learn about C_{B_{A_w}} path
- C will route C_{A_w} (not using B) to get to w

Network Layer: Control Plane 5-53

BGP: achieving policy via advertisements (2)

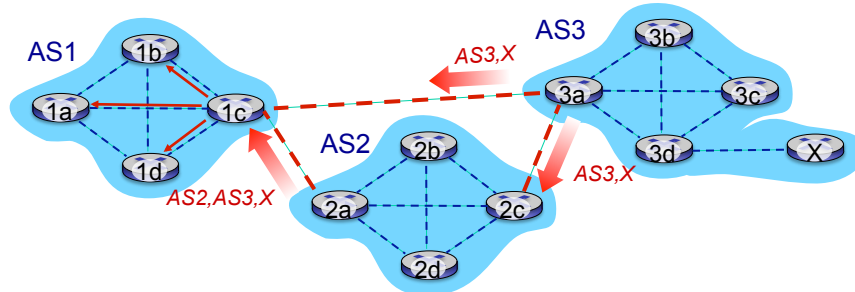


Suppose an ISP only wants to route traffic to/from its customer networks (does not want to carry transit traffic between other ISPs)

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

Network Layer: Control Plane 5-54

Shortest AS-PATH

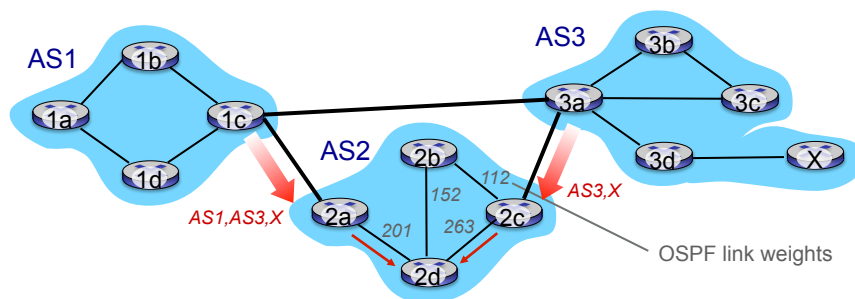


gateway router may learn about **multiple** paths to destination:

- AS1 gateway router 1c learns path **AS2,AS3,X** from 2a
- AS1 gateway router 1c learns path **AS3,X** from 3a
- Shortest AS-PATH is **AS3,X**

Network Layer: Control Plane 5-55

Hot Potato Routing



- 2d learns (via iBGP) it can route to X via 2a or 2c
- hot potato routing**: choose local gateway that has least intra-domain cost (e.g., 2d chooses 2a, even though more AS hops to X); does not consider inter-domain cost!
- selfish; may lead to longer end-to-end delay → apply shortest AS-PATH before hot potato.

Network Layer: Control Plane 5-56

BGP messages

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

Network Layer: Control Plane 5-57

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: Control Plane 5-58

Chapter 5: outline

5.1 introduction

5.2 routing protocols

- link state
- distance vector

5.3 intra-AS routing in the Internet: OSPF

5.4 routing among the ISPs: BGP

5.5 The SDN control plane

5.6 ICMP: The Internet Control Message Protocol

5.7 Network management and SNMP