

Network layer: “control plane” roadmap

- introduction
- routing protocols
 - link state
 - distance vector
- intra-ISP routing: OSPF
- routing among ISPs: BGP
- SDN control plane
- Internet Control Message Protocol
- network management, configuration
 - SNMP
 - NETCONF/YANG

Network-layer functions

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

data plane

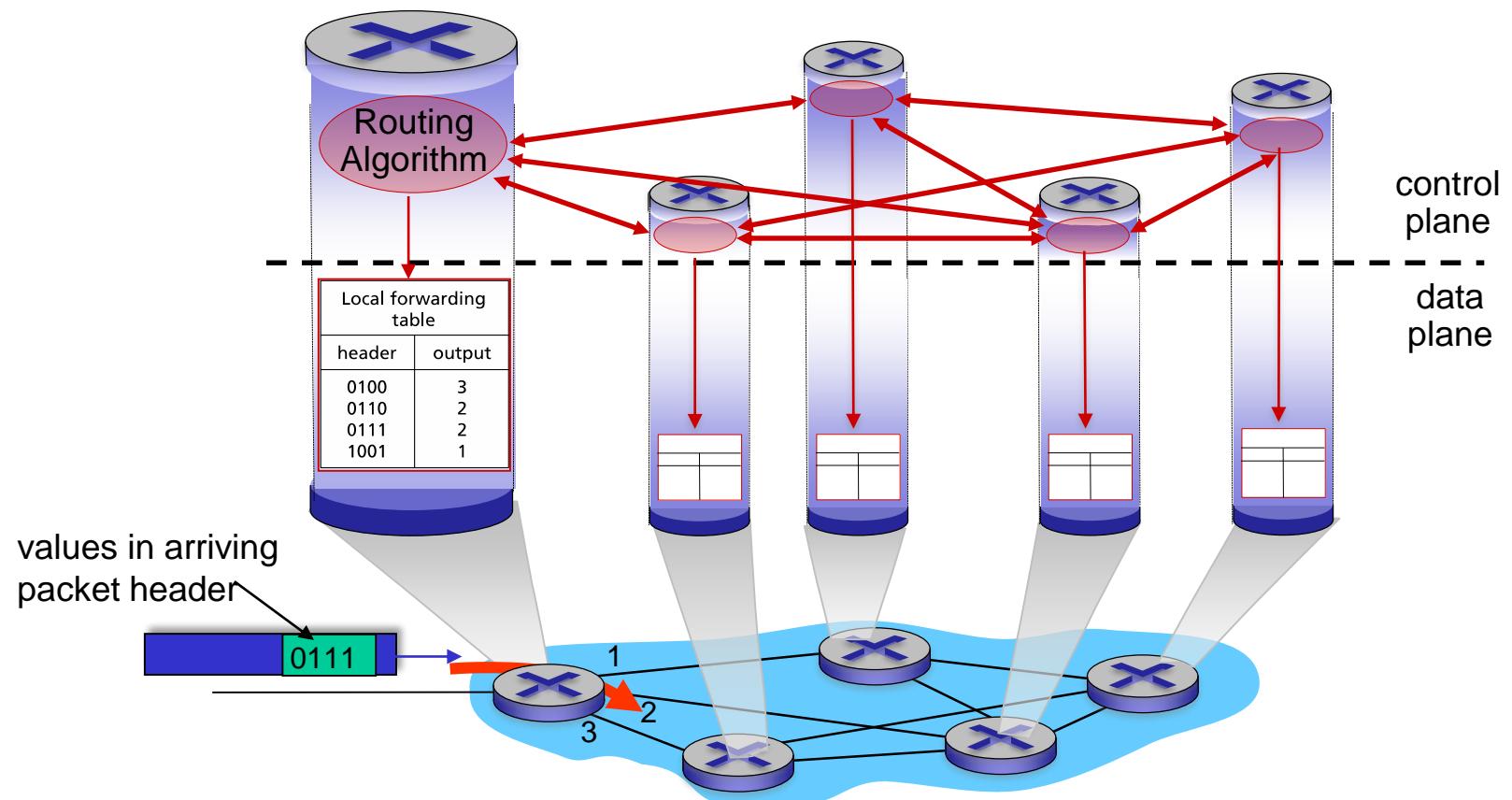
control plane

Two approaches to structuring network control plane:

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

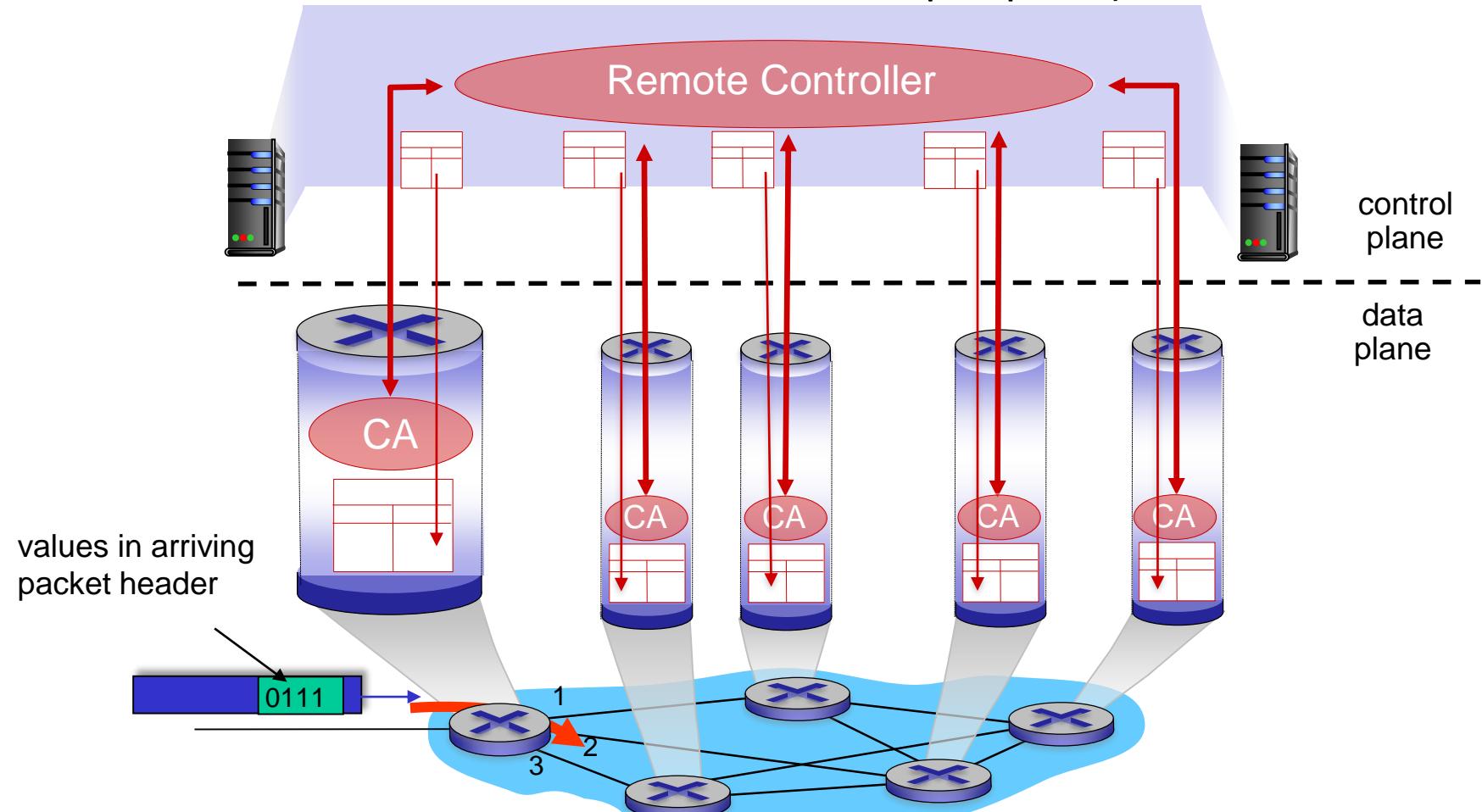
Per-router control plane

Individual routing algorithm components *in each and every router* interact in the control plane (routing processor maintains routing tables and computes the forwarding table for that router, each router has a routing component that communicates with the routing components in other routers to compute the values for its forwarding table)



Software-Defined Networking (SDN) control plane

Remote controller computes, installs forwarding tables in routers (routing processor communicates with the remote controller via the control agent, receive forwarding table entries computed by the remote controller, and install these entries in the router's input ports)



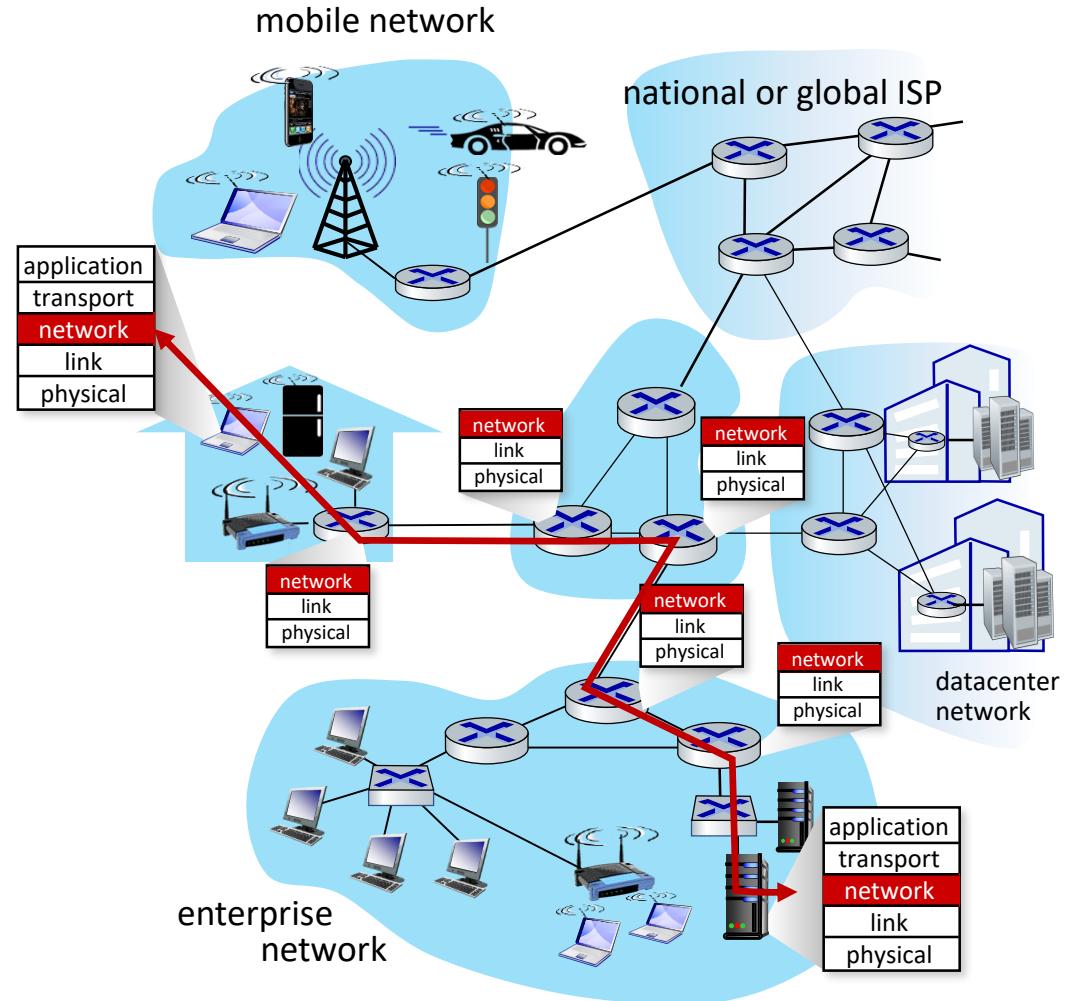
Network layer: “control plane” roadmap

- introduction
- **routing protocols**
 - link state
 - distance vector
- intra-ISP routing: OSPF
- routing among ISPs: BGP
- SDN control plane
- Internet Control Message Protocol
- network management, configuration
 - SNMP
 - NETCONF/YANG

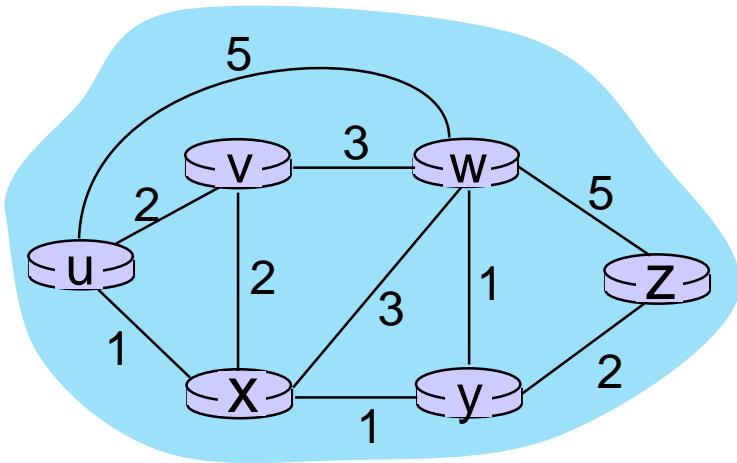
Routing protocols

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

- **path:** sequence of routers where packets traverse from given initial source host to final destination host
- **“good”:** least “cost”, “fastest”, “least congested”



Graph abstraction: link costs



$c_{a,b}$: cost of *direct* link connecting a and b
e.g., $c_{w,z} = 5$, $c_{u,z} = \infty$

may reflect the physical length of the corresponding link, the link speed, or the monetary cost associated with a link

graph: $G = (N, E) \rightarrow$ A graph is used to formulate routing problems

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

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

Routing algorithm classification (2 classifications)

global (centralized) : all routers have *complete* topology, link cost info (complete, global knowledge about the network is used)

- “link state” algorithms

decentralized: iterative process of computation, exchange of info with neighbors

- routers initially only know link costs to attached neighbors (No node has complete information about the costs of all network)
- “distance vector” algorithms

static: routes change slowly over time (as a result of human intervention)

dynamic: routes change more quickly (change the routing paths as the network traffic loads or topology change)

- periodic updates or in response to link cost changes

Dijkstra's link-state routing algorithm

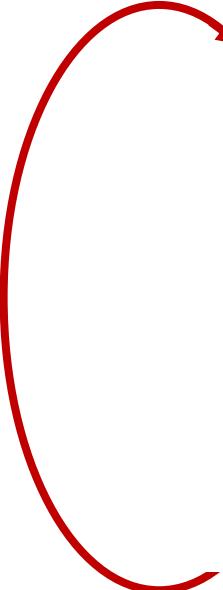
- **centralized:** network topology and link costs known to *all* nodes
 - accomplished via “link state broadcast” (each node broadcast link-state packets to all other nodes in the network)
 - all nodes have same info
- computes least cost paths from one node (“source”) to all other nodes
 - gives *forwarding table* for that node
- **iterative:** after k iterations, know least cost path to k destinations

notation

- $c_{x,y}$: direct link cost from node x to y ; $= \infty$ if not direct neighbors
- $D(v)$: *current* estimate of cost of least-cost-path from source to destination v
- $p(v)$: predecessor node along path from source to v
- N' : set of nodes whose least-cost-path *definitively* known

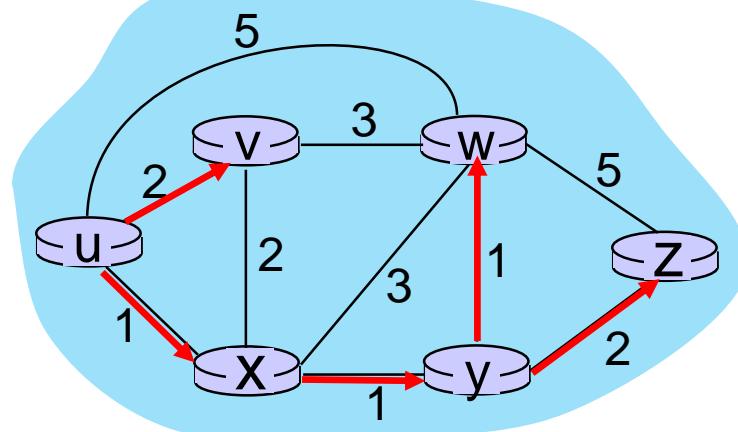
Dijkstra's link-state routing algorithm

```
1 Initialization:
2    $N' = \{u\}$                                 /* compute least cost path from u to all other nodes */
3   for all nodes  $v$ 
4     if  $v$  adjacent to  $u$                       /*  $u$  initially knows direct-path-cost only to direct neighbors */
5       then  $D(v) = c_{u,v}$                       /* but may not be minimum cost!
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 least-path-cost to  $v$  is either old least-cost-path to  $v$  or known
14    least-cost-path to  $w$  plus direct-cost from  $w$  to  $v$  */
15 until all nodes in  $N'$ 
```



Dijkstra's algorithm: an 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	u, x	2, u	4, x	2, x	∞	∞
2	u, x, y	2, u	3, y	∞	4, y	∞
3	u, x, y, v	∞	3, y	∞	4, y	∞
4	u, x, y, v, w	∞	∞	∞	4, y	∞
5	u, x, y, v, w, z	∞	∞	∞	∞	∞

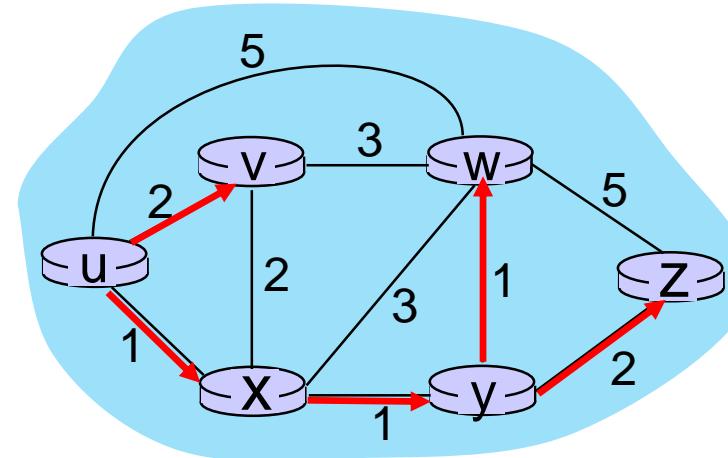


Initialization (step 0): For all a : if a adjacent to then $D(a) = c_{u,a}$

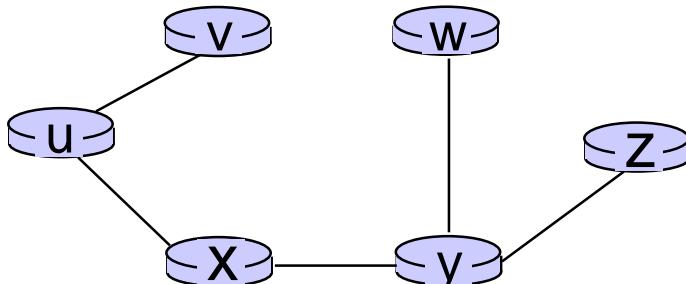
find a not in N' such that $D(a)$ is a minimum
 add a to N'
 update $D(b)$ for all b adjacent to a and not in N' :

$$D(b) = \min(D(b), D(a) + c_{a,b})$$

Dijkstra's algorithm: an example



resulting least-cost-path tree from u:



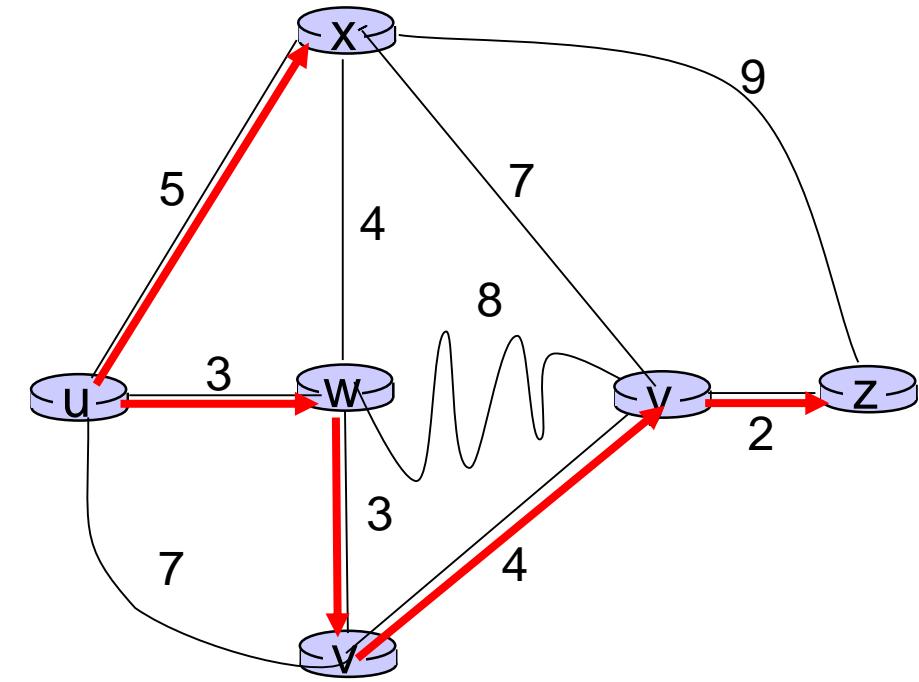
resulting forwarding table in u:

destination	outgoing link
v	(u,v)
x	(u,x)
y	(u,x)
w	(u,x)
z	(u,x)

route from u to v directly
route from u to all other destinations via x

Dijkstra's algorithm: another example

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



Network layer: “control plane” roadmap

- introduction
- **routing protocols**
 - link state
 - **distance vector**
- intra-ISP routing: OSPF
- routing among ISPs: BGP
- SDN control plane
- Internet Control Message Protocol
- network management, configuration
 - SNMP
 - NETCONF/YANG

Distance vector algorithm

Based on *Bellman-Ford* (BF) equation:

Bellman-Ford equation

Let $D_x(y)$: cost of least-cost path from x to y .

Then update distances based on neighbors:

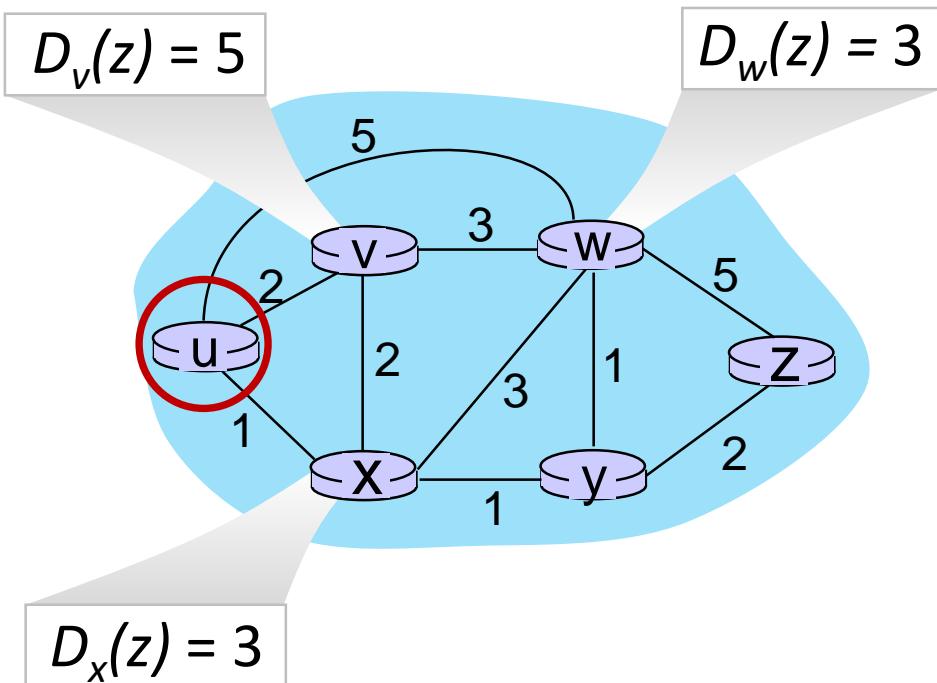
$$D_x(y) = \min_v \{ c_{x,v} + D_v(y) \}$$

\min taken over all neighbors v of x

v 's estimated least-cost-path cost to y
direct cost of link from x to v

Bellman-Ford Example

Suppose that u 's neighboring nodes (*having direct connection*), x, v, w know that for destination z :



Bellman-Ford 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 -x- is next hop on estimated least-cost path to destination -z-

Distance vector algorithm

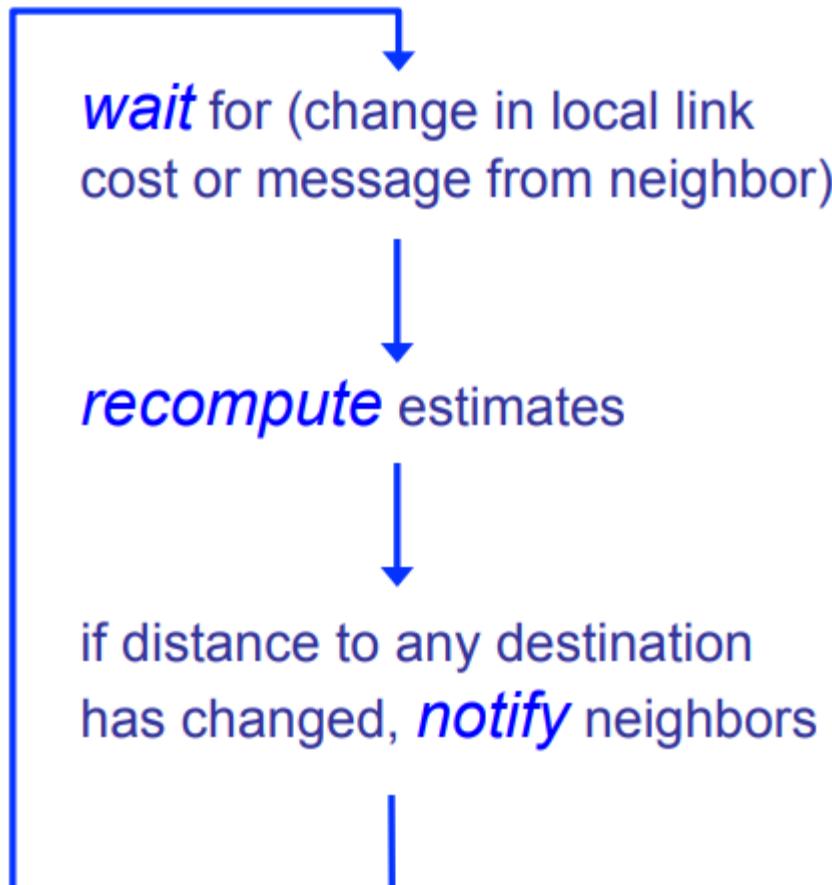
key idea:

- from time-to-time, each node sends its own distance vector estimate to neighbors
- when x receives new DV estimate from any neighbor, it updates its own DV using Bellman-Ford equation:

$$D_x(y) \leftarrow \min_v \{c_{x,v} + D_v(y)\} \text{ for each node } y \in N$$

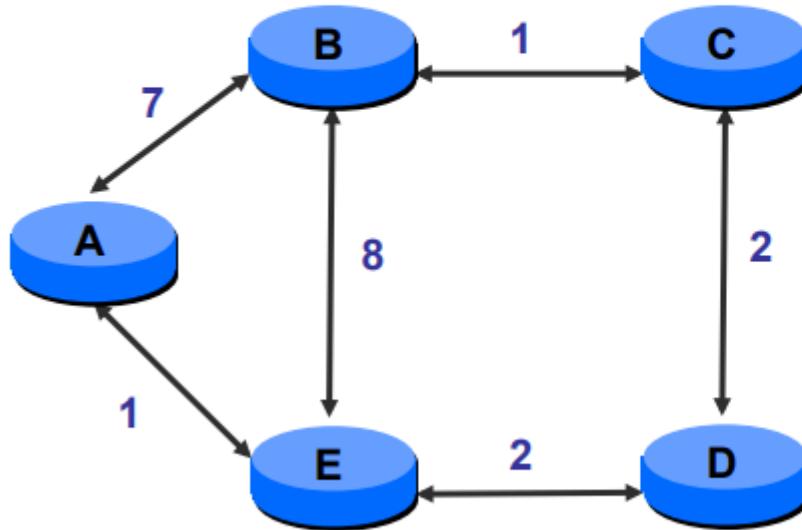
Distance vector algorithm

Each node:



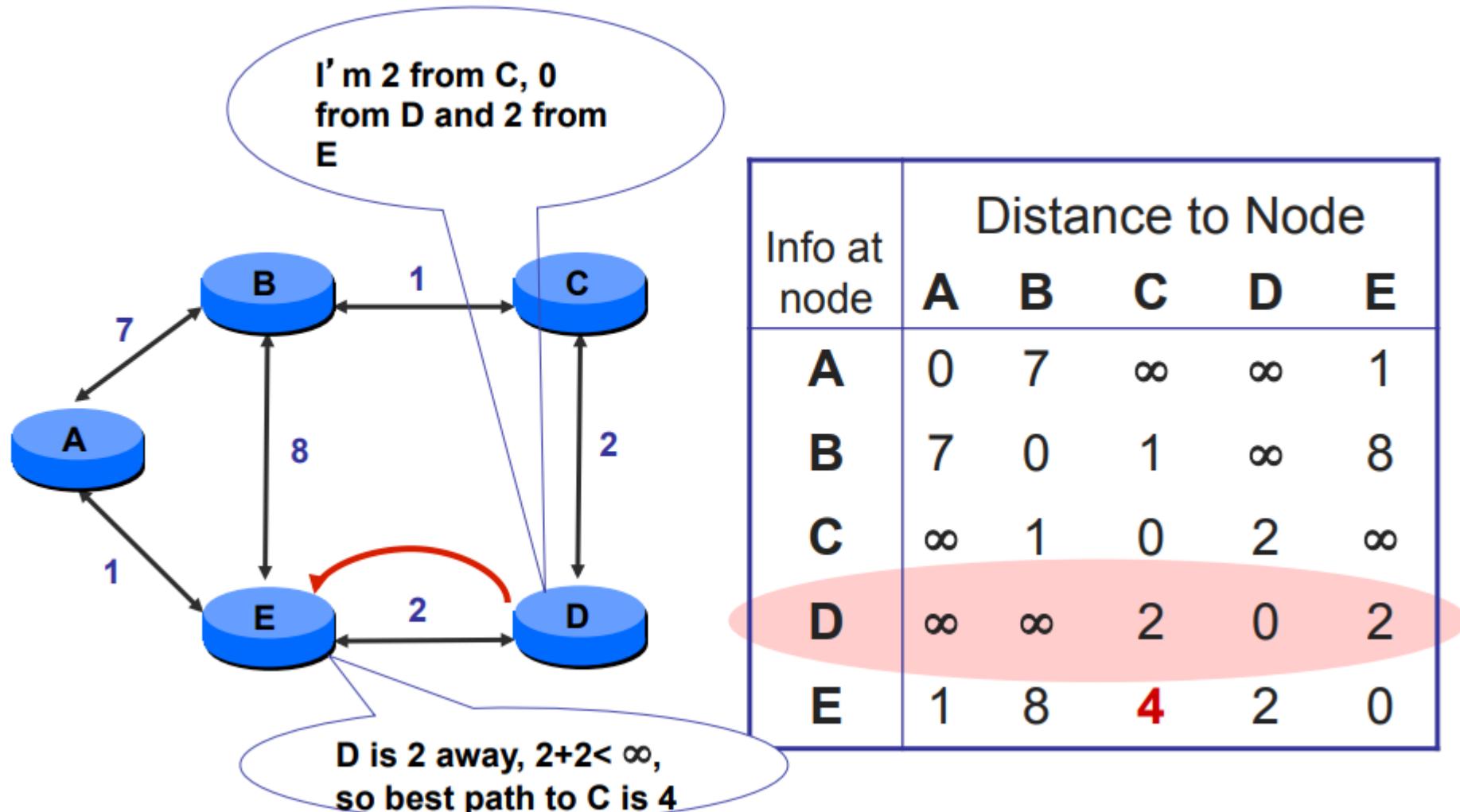
- $c(x,v)$ = cost for direct link from x to v
 - Node x maintains costs of direct links $c(x,v)$
- $D_x(y)$ = estimate of least cost from x to y
 - Node x maintains distance vector $\mathbf{D}_x = [D_x(y): y \in N]$
- Node x maintains its neighbors' distance vectors
 - For each neighbor v , x maintains $\mathbf{D}_v = [D_v(y): y \in N]$
- Each node v periodically sends D_v to its neighbors
 - And neighbors update their own distance vectors
 - $D_x(y) \leftarrow \min_v \{c(x,v) + D_v(y)\}$ for each node $y \in N$

Distance vector: example (initial state)

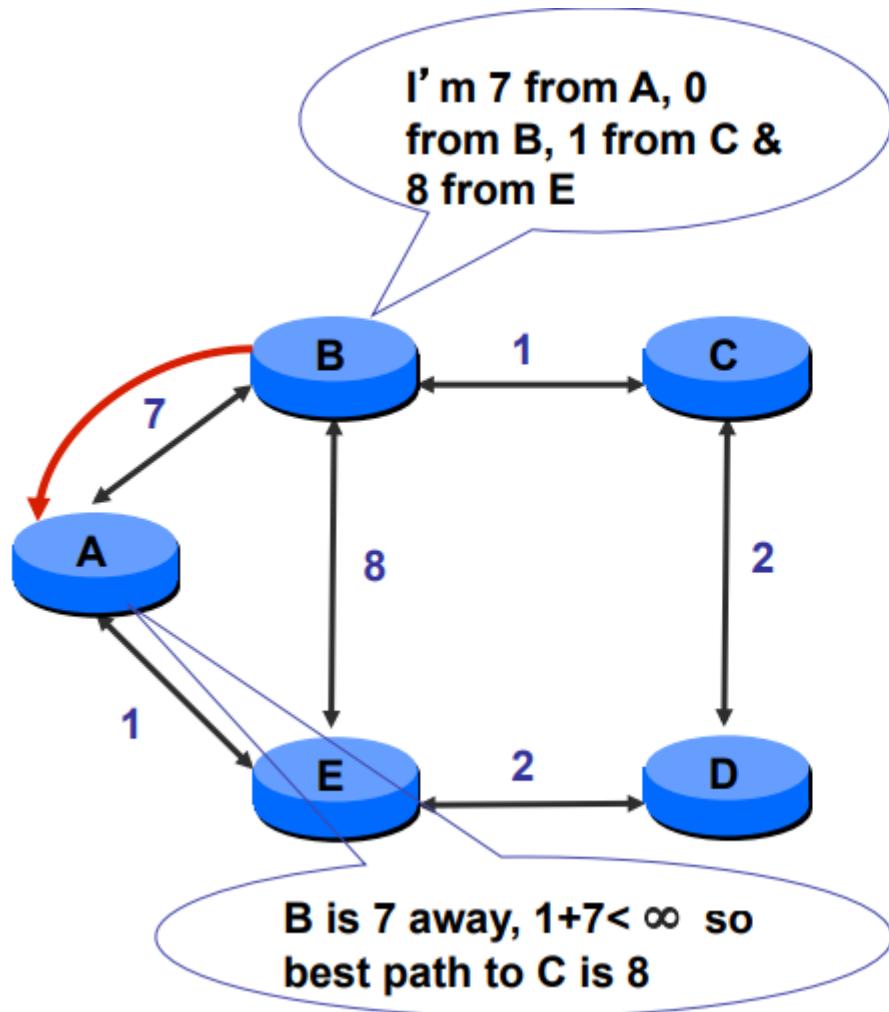


Info at node	Distance to Node				
	A	B	C	D	E
A	0	7	∞	∞	1
B	7	0	1	∞	8
C	∞	1	0	2	∞
D	∞	∞	2	0	2
E	1	8	∞	2	0

Distance vector: example (D sends vector to E)

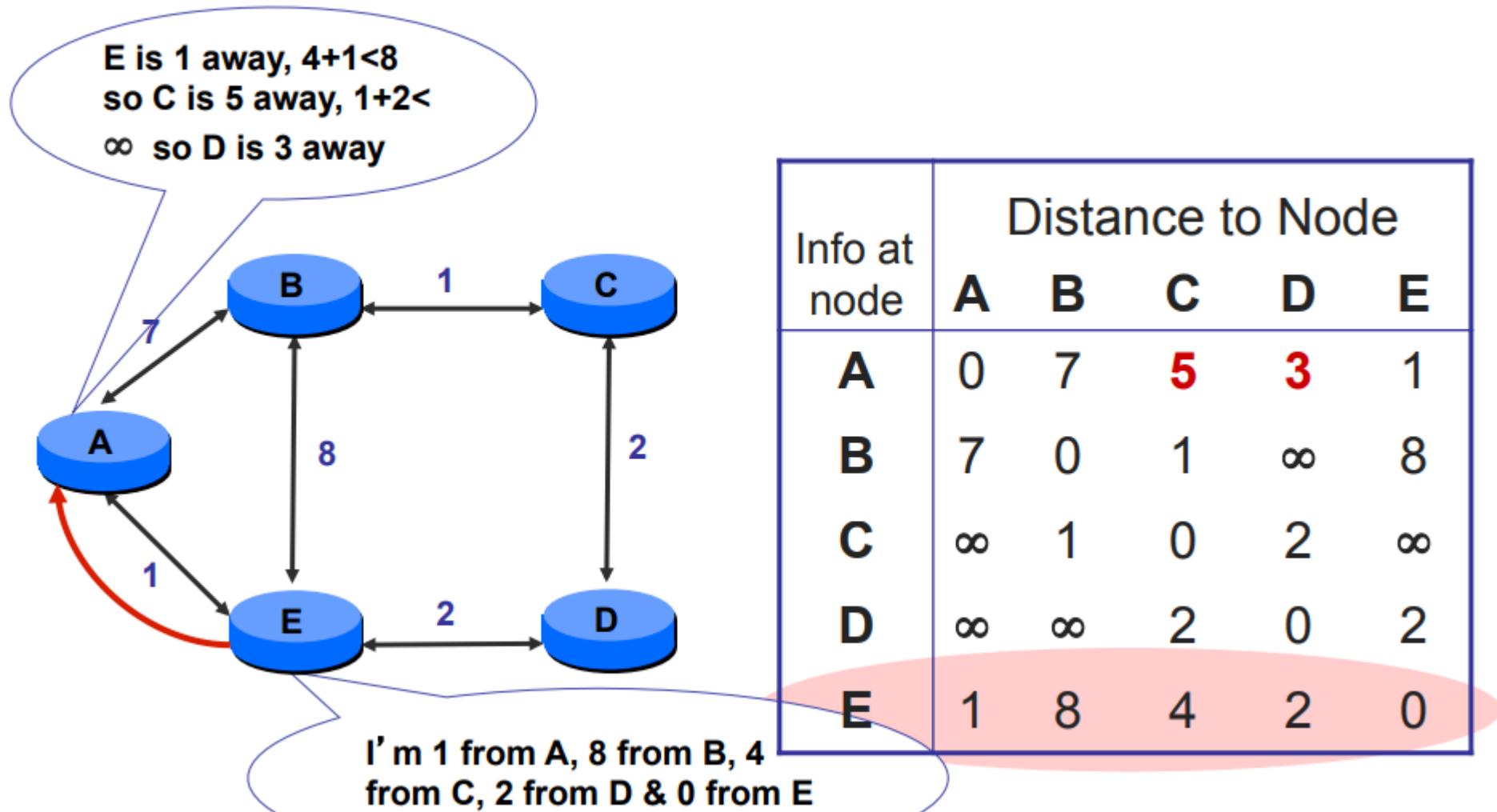


Distance vector: example (B sends vector to A)

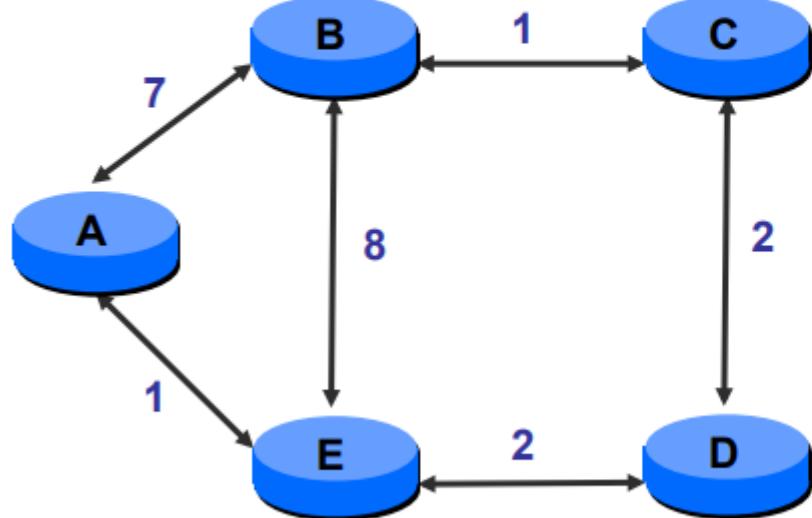


Info at node	Distance to Node				
	A	B	C	D	E
A	0	7	8	∞	1
B	7	0	1	∞	8
C	∞	1	0	2	∞
D	∞	∞	2	0	2
E	1	8	4	2	0

Distance vector: example (E sends vector to A)

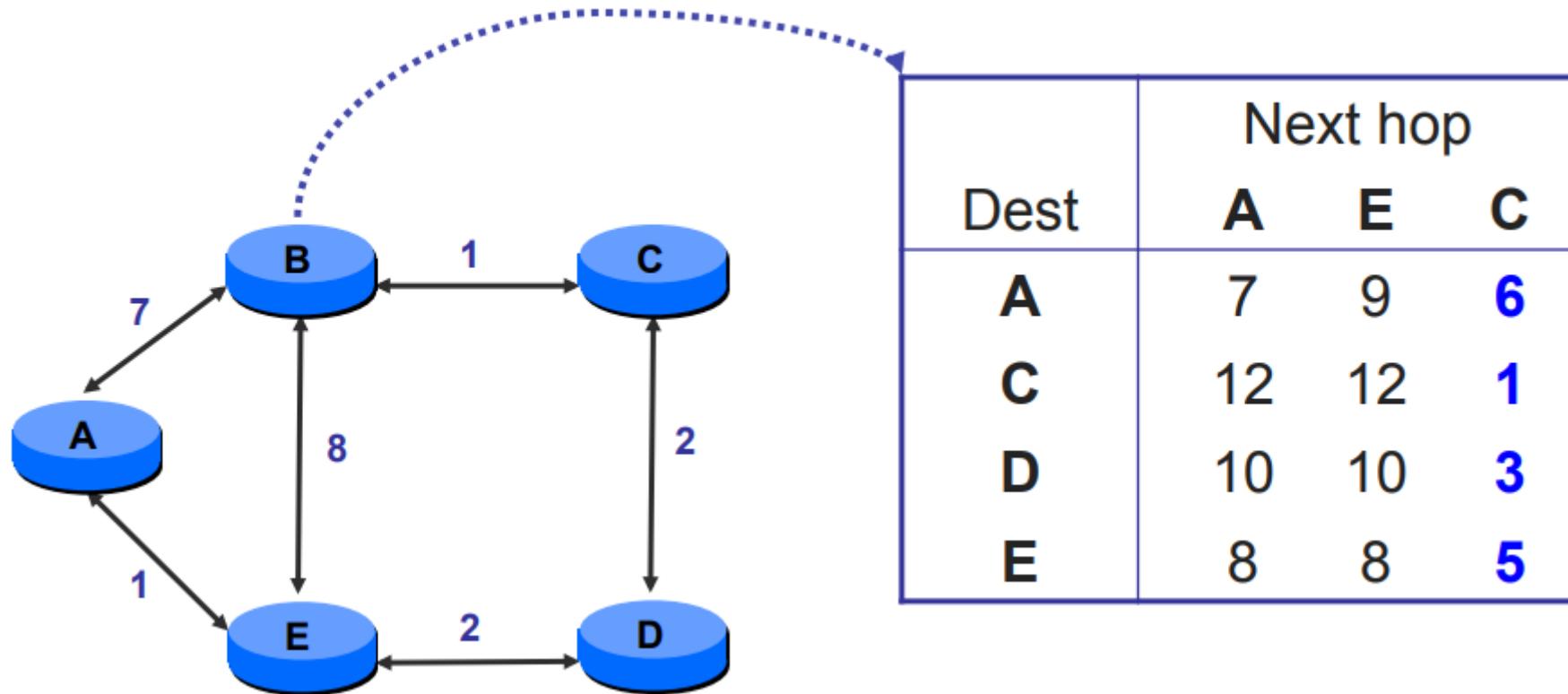


Distance vector: example (until convergence)



Info at node	Distance to Node				
	A	B	C	D	E
A	0	6	5	3	1
B	6	0	1	3	5
C	5	1	0	2	4
D	3	3	2	0	2
E	1	5	4	2	0

Distance vector: example (B's distance vectors)



Network layer: “control plane” roadmap

- introduction
- routing protocols
- **intra-ISP routing: OSPF**
- routing among ISPs: BGP
- SDN control plane
- Internet Control Message Protocol
- network management, configuration
 - SNMP
 - NETCONF/YANG

Internet approach to scalable routing

- * organizing routers into autonomous systems (ASs)
- * **AS**: consisting of a group of routers that are under the same administrative control
- * the routers in an ISP, and the links that interconnect them, constitute a single AS
- * routers within the same AS all run the same routing algorithm (intra-autonomous system routing protocol)

Internet approach to scalable routing

most common intra-AS routing protocols:

- **RIP: Routing Information Protocol [RFC 1723]**
 - classic DV: DVs exchanged every 30 secs
 - no longer widely used
- **EIGRP: Enhanced Interior Gateway Routing Protocol**
 - DV based
 - formerly Cisco-proprietary for decades (became open in 2013 [RFC 7868])
- **OSPF: Open Shortest Path First [RFC 2328]**
 - link-state routing
 - ISO standard, not RFC standard

OSPF (Open Shortest Path First) routing

- “open”: publicly available
- classic link-state
 - each router floods OSPF link-state advertisements (directly over IP rather than using TCP/UDP) to all other routers in entire AS
 - multiple link costs metrics possible: bandwidth, delay
 - each router has full topology, uses Dijkstra’s algorithm to compute forwarding table
- *security*: all OSPF messages authenticated (to prevent malicious intrusion)

Hierarchical OSPF

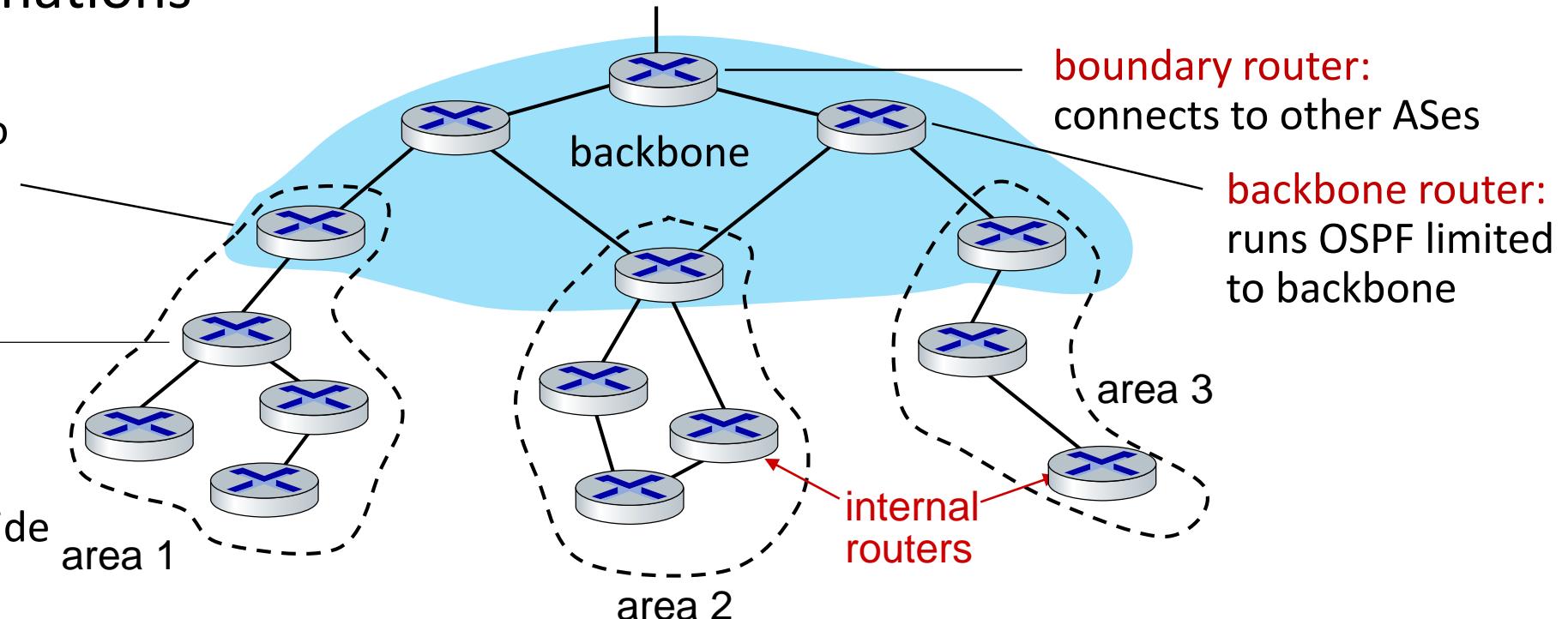
- an OSPF autonomous system can be configured hierarchically into areas
- **two-level hierarchy:** local area, backbone.
- link-state advertisements flooded only in area, or backbone
- each node has detailed area topology; only knows direction to reach other destinations

area border routers:

“summarize” distances to destinations in own area, advertise in backbone

local routers:

- flood LS in area only
- compute routing within area
- forward packets to outside via area border router



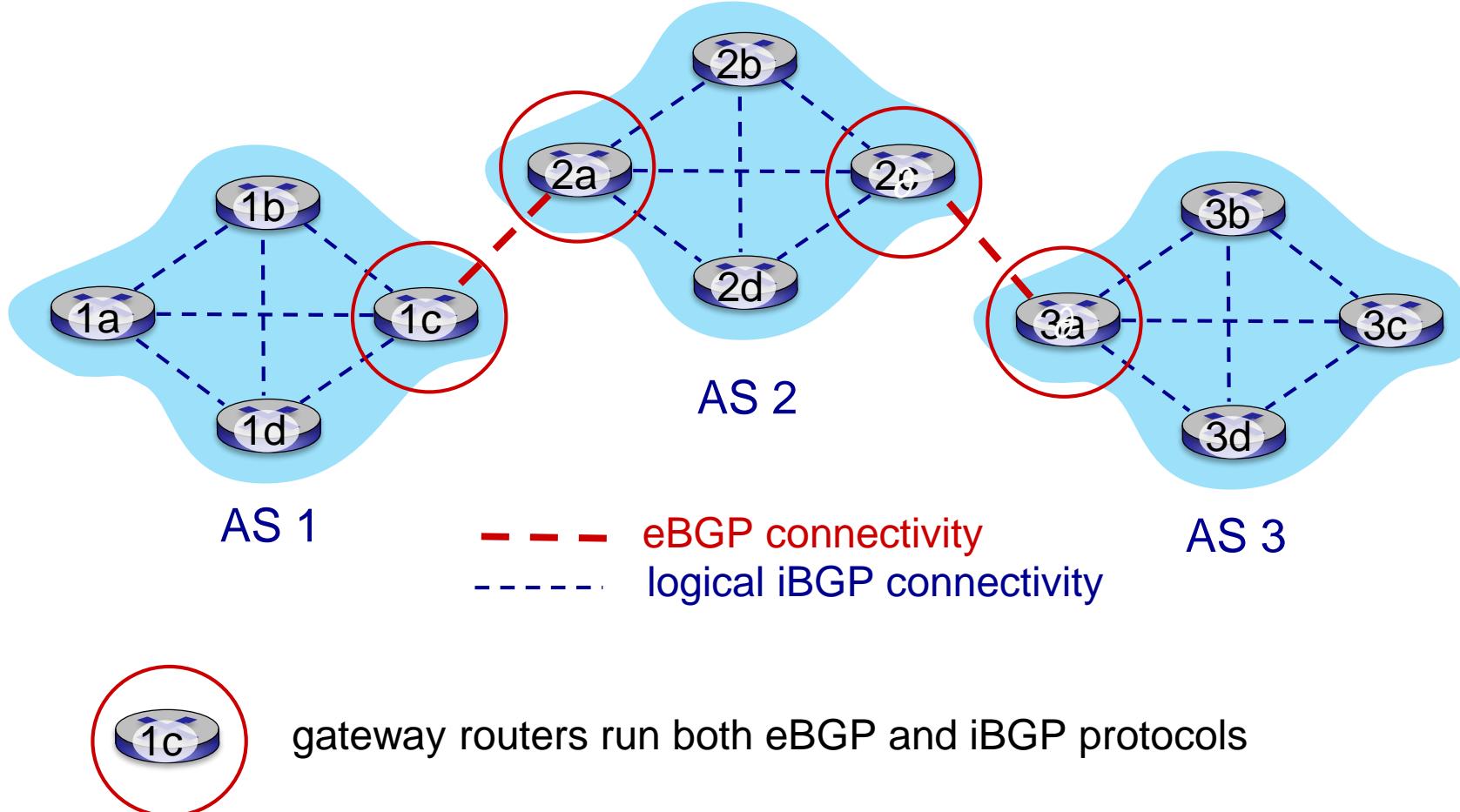
Network layer: “control plane” roadmap

- introduction
- routing protocols
- intra-ISP routing: OSPF
- **routing among ISPs: BGP**
- SDN control plane
- Internet Control Message Protocol
- network management, configuration
 - SNMP
 - NETCONF/YANG

Internet inter-AS routing: BGP

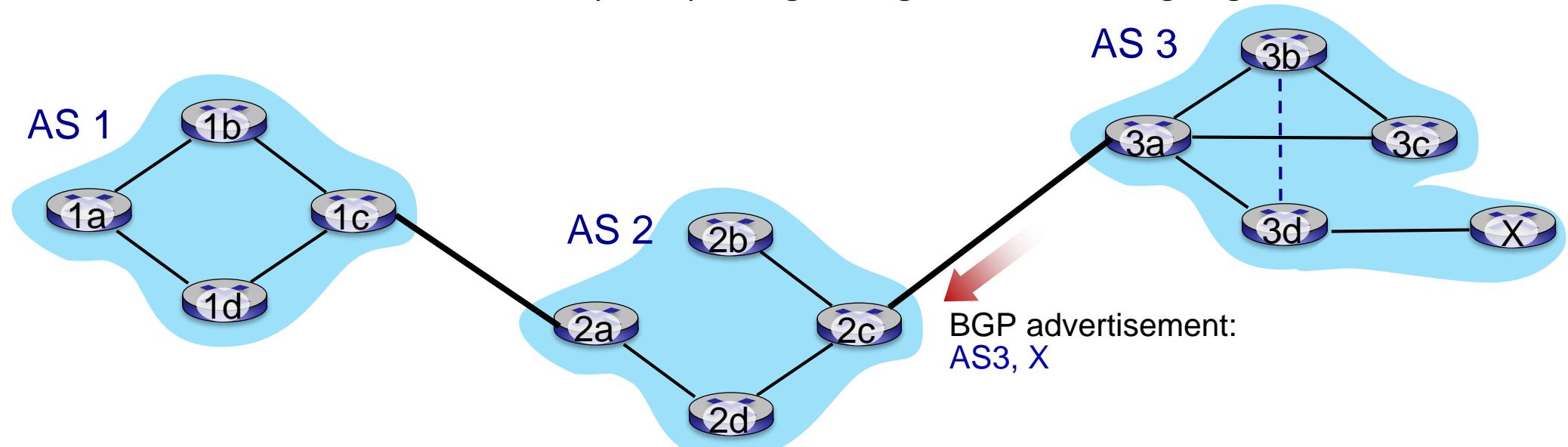
- BGP (Border Gateway Protocol): *the de facto* inter-domain routing protocol
 - “glue that holds the Internet together”
- allows subnet to advertise its existence, and the destinations it can reach, to rest of Internet: *“I am here, here is who I can reach, and how”*
- BGP provides each AS a means to: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

eBGP, iBGP connections



BGP basics

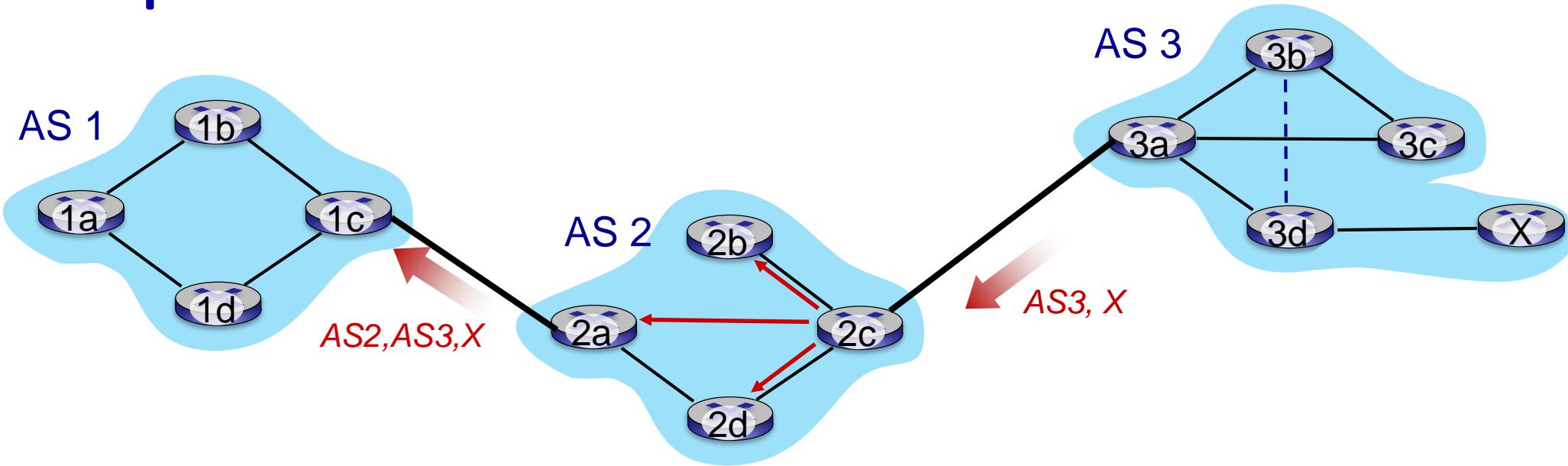
- **BGP session:** two BGP routers (“peers”) exchange BGP messages over TCP connection:
 - advertising *paths* to different destination network prefixes (BGP is a “path vector” protocol)
- when AS3 gateway 3a advertises **path AS3,X** to AS2 gateway 2c:
 - AS3 *promises* to AS2 it will forward datagrams towards X
- when AS2 gateway 2a advertises **path AS3,AS2,X** to AS1 gateway 1c:
 - It means that AS1 can arrive at X by first passing through AS2 and then going to AS3



Path attributes and BGP routes

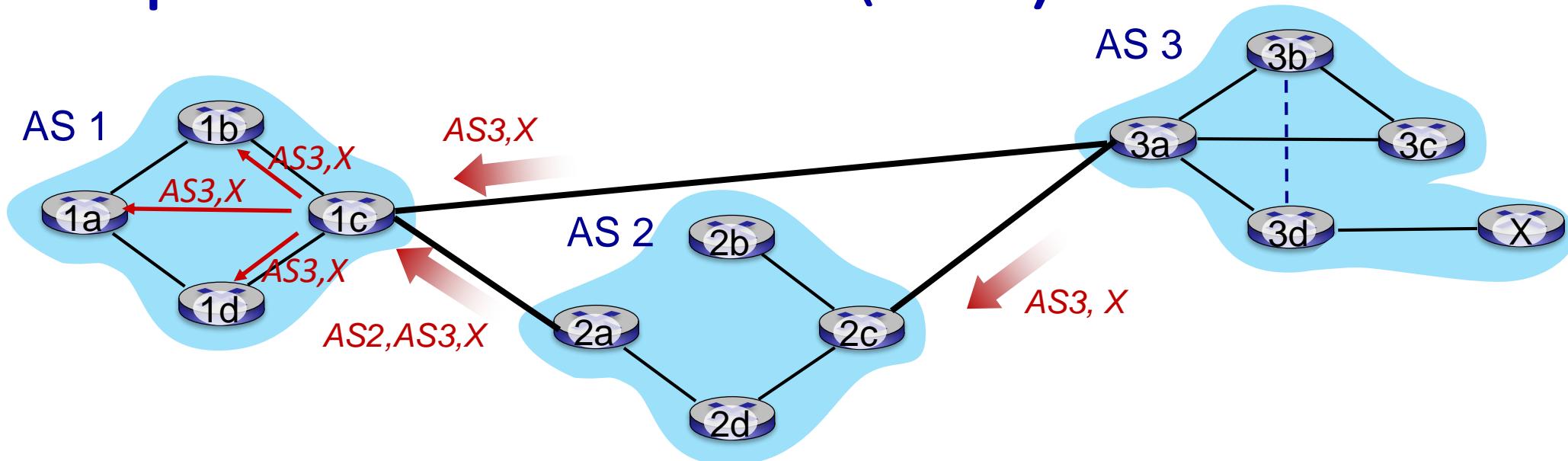
- policy-based routing:
 - gateway receiving route advertisement uses *import policy* to accept/decline path
 - AS policy also determines whether to *advertise* path to other neighboring ASes

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

BGP path advertisement (more)



gateway router may learn about multiple paths to destination:

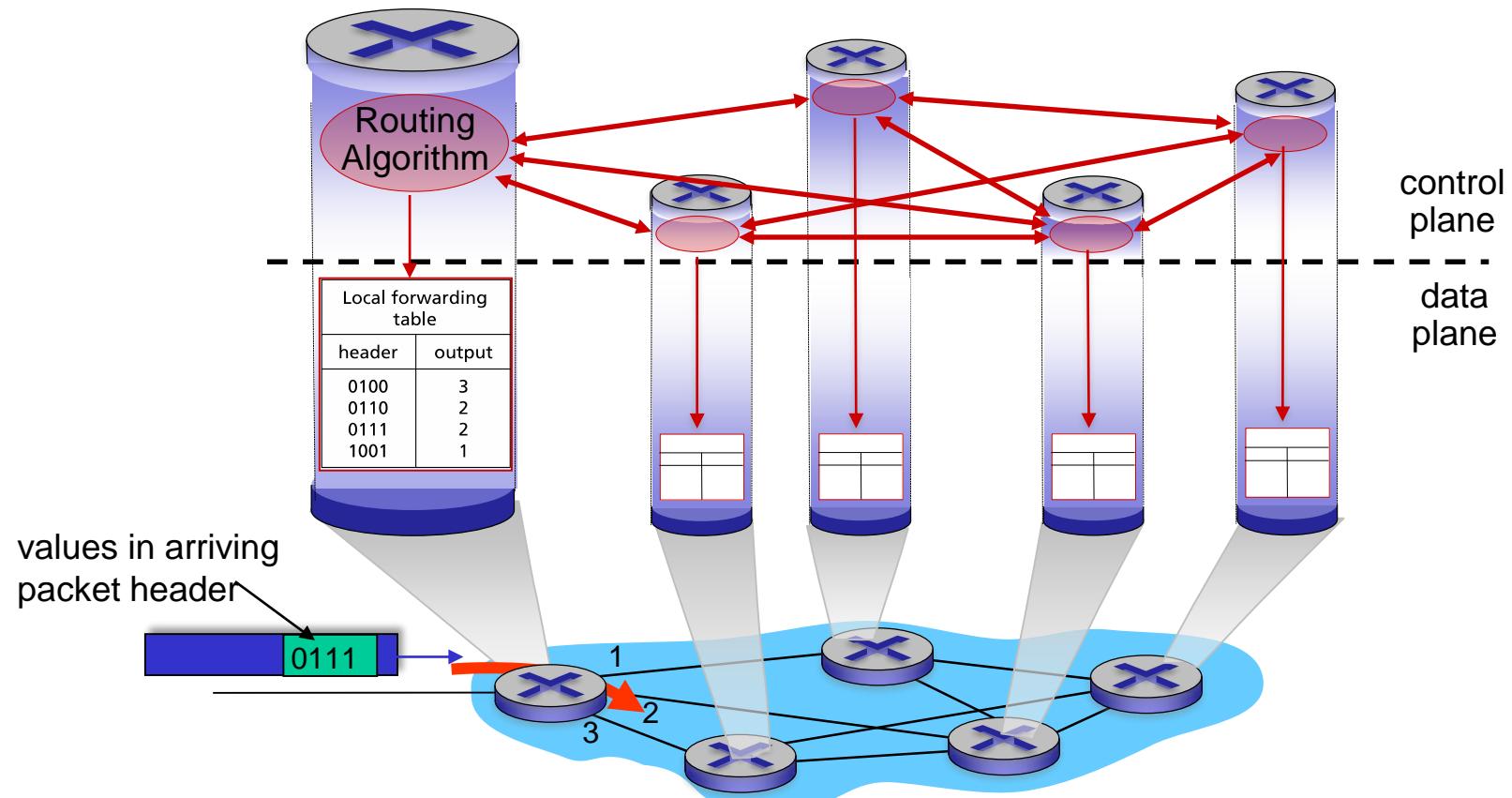
- AS1 gateway router 1c learns path $AS2, AS3, X$ from 2a
- AS1 gateway router 1c learns path $AS3, X$ from 3a
- based on *policy*, AS1 gateway router 1c chooses path $AS3, X$ and advertises path within AS1 via iBGP

Network layer: “control plane” roadmap

- introduction
- routing protocols
- intra-ISP routing: OSPF
- routing among ISPs: BGP
- **SDN control plane**
- Internet Control Message Protocol
- network management, configuration
 - SNMP
 - NETCONF/YANG

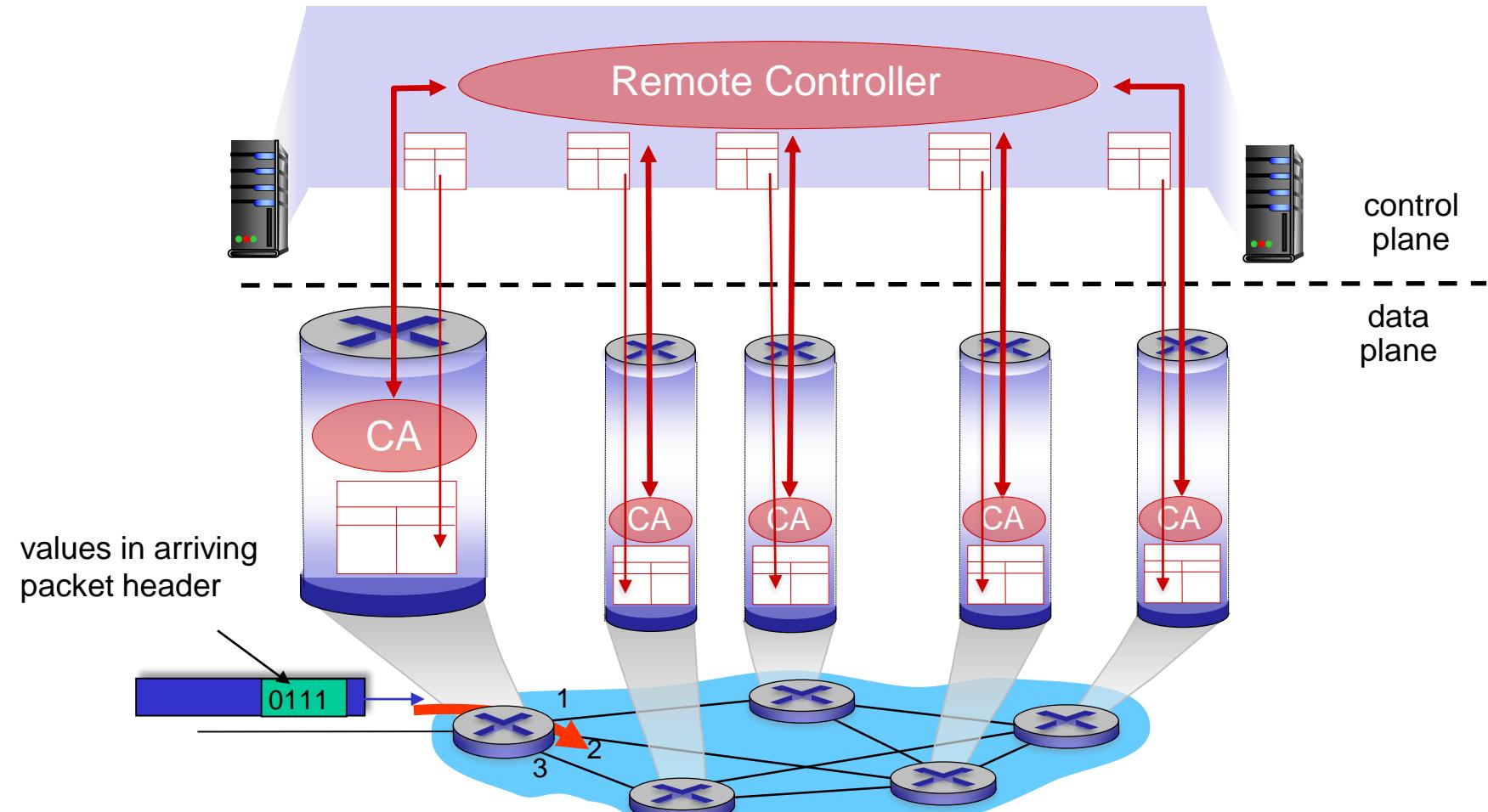
Per-router control plane

Individual routing algorithm components *in each and every router* interact in the control plane to computer forwarding tables



Software-Defined Networking (SDN) control plane

Remote controller computes, installs forwarding tables in routers

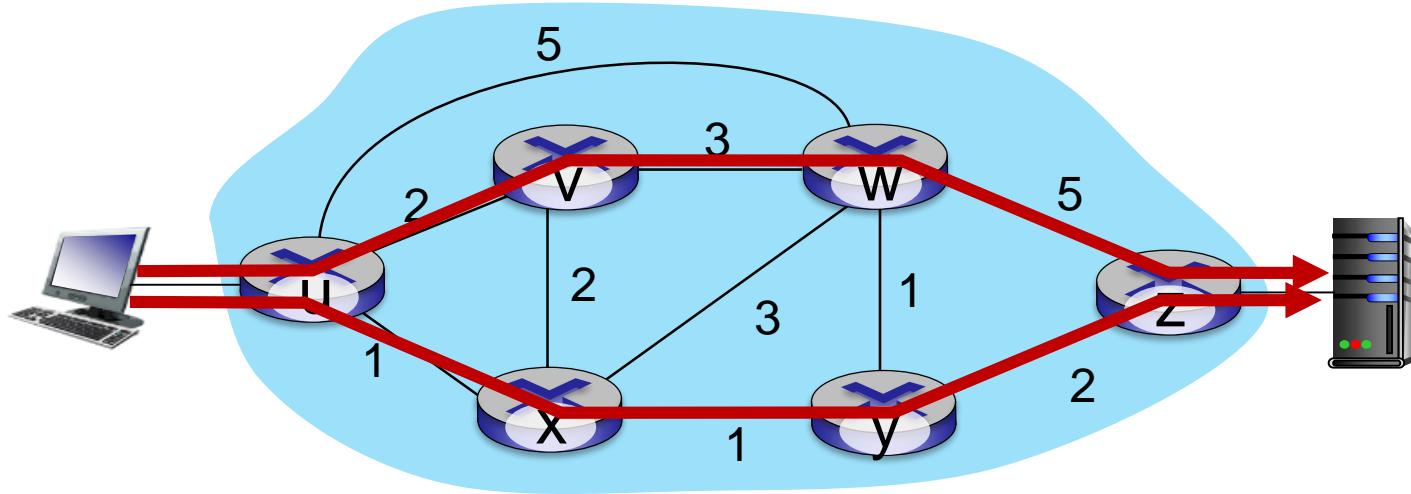


Software defined networking (SDN)

Why a logically centralized control plane?

- easier network management: avoid router misconfigurations, greater flexibility of traffic flows
- table-based forwarding allows “programming” routers
 - centralized “programming” easier: compute tables centrally and distribute
 - distributed “programming” more difficult: compute tables as result of distributed algorithm (protocol) implemented in each-and-every router
- open (non-proprietary) implementation of control plane

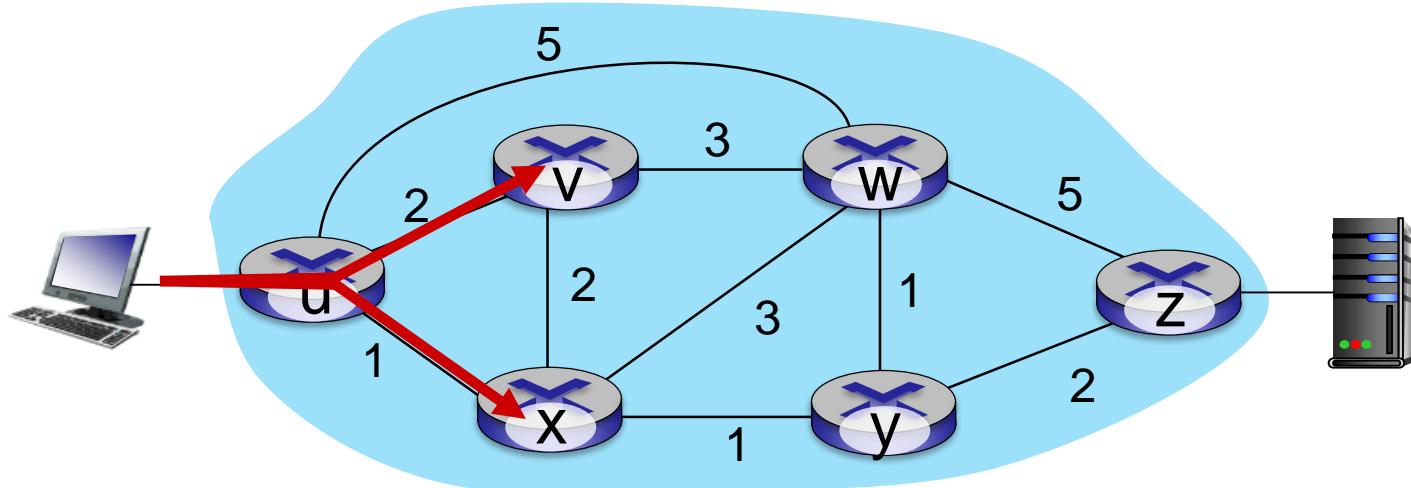
Traffic engineering: difficult with traditional routing



Q: what if network operator wants u-to-z traffic to flow along $uvwz$, rather than $uxyz$?

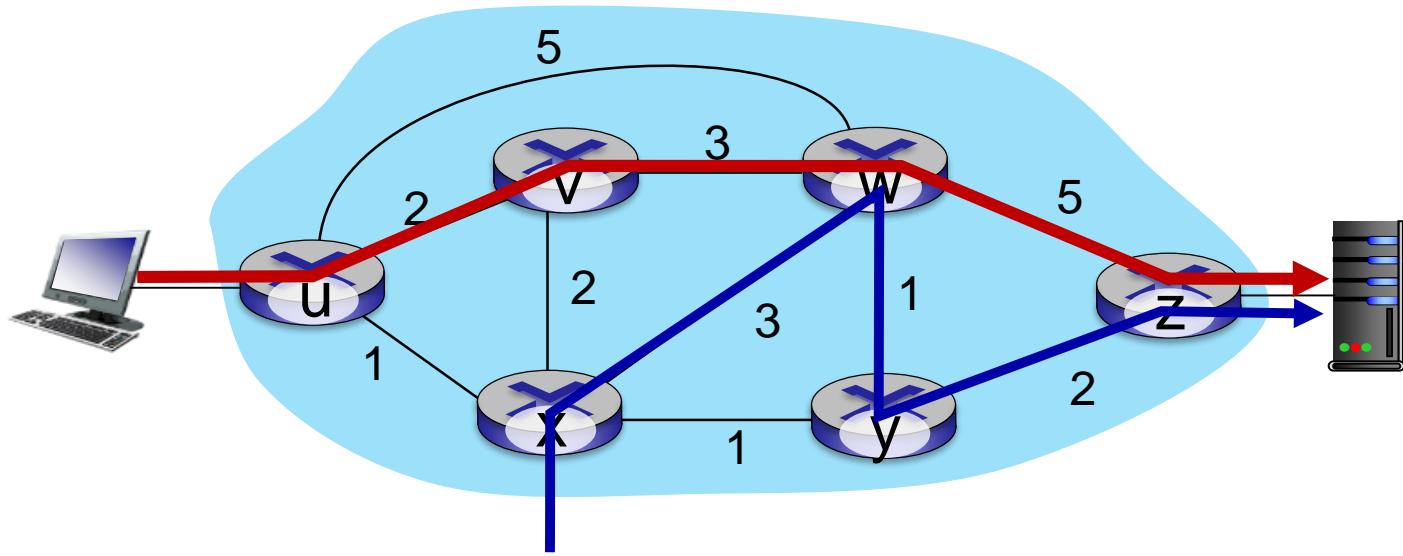
A: need to re-define link weights so traffic routing algorithm computes routes accordingly (or need a new routing algorithm)!

Traffic engineering: difficult with traditional routing



Q: what if network operator wants to split u-to-z traffic along uvwz *and* uxyz (load balancing)?
A: can't do it (or need a new routing algorithm)

Traffic engineering: difficult with traditional routing



Q: what if w wants to route blue and red traffic differently from w to z?

A: can't do it (with destination-based forwarding, and LS, DV routing)

Software defined networking (SDN)

4. programmable
control
applications

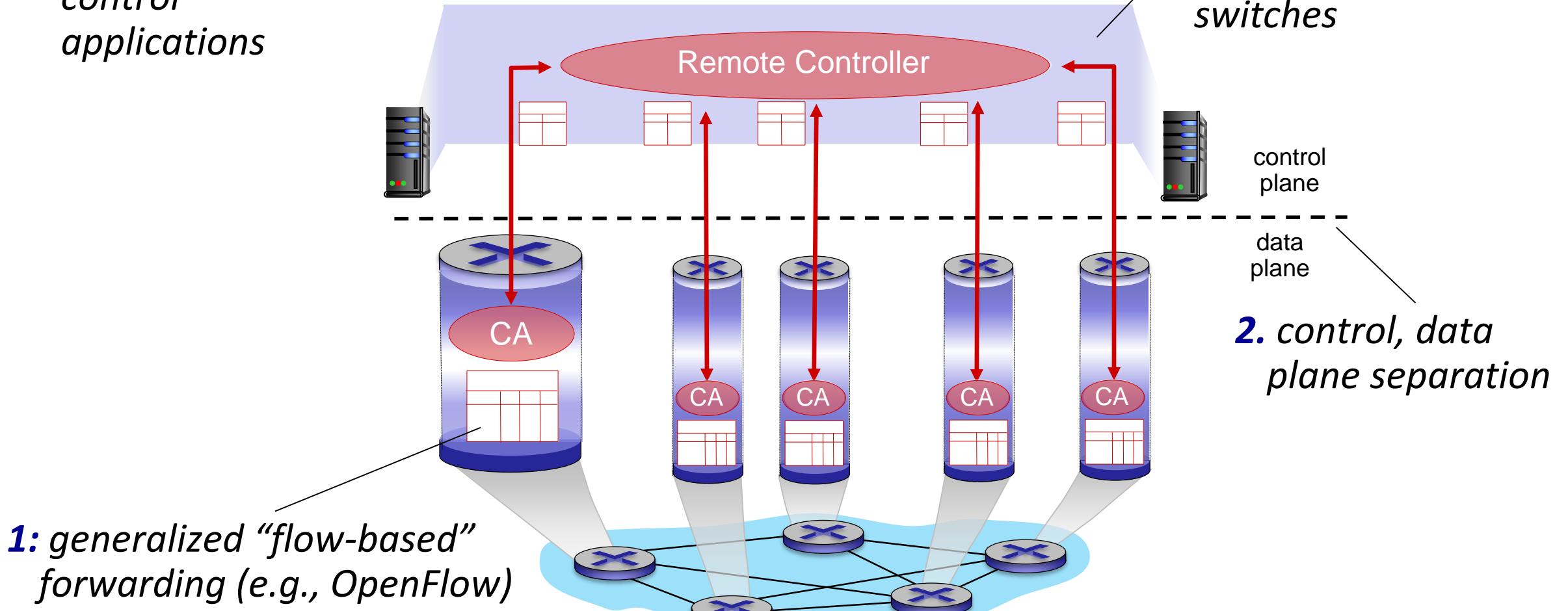
routing

access
control

...

load
balance

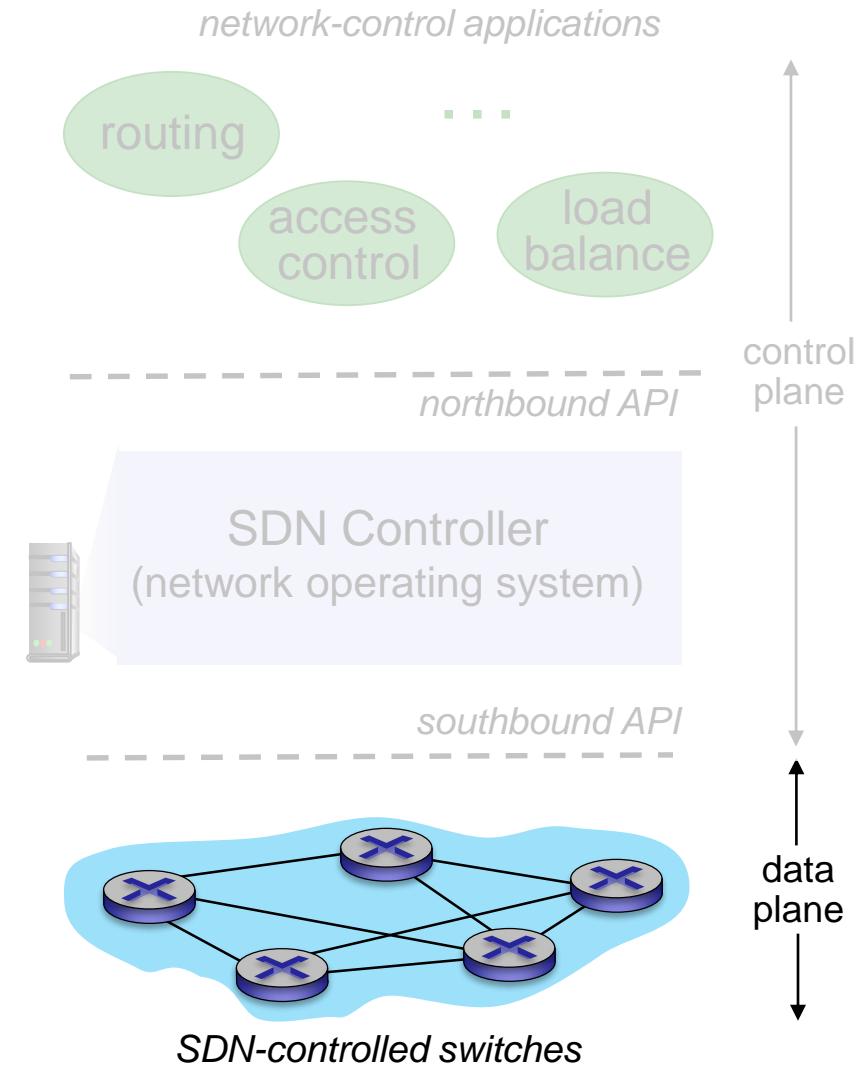
3. control plane functions
external to data-plane
switches



Software defined networking (SDN)

Data-plane switches:

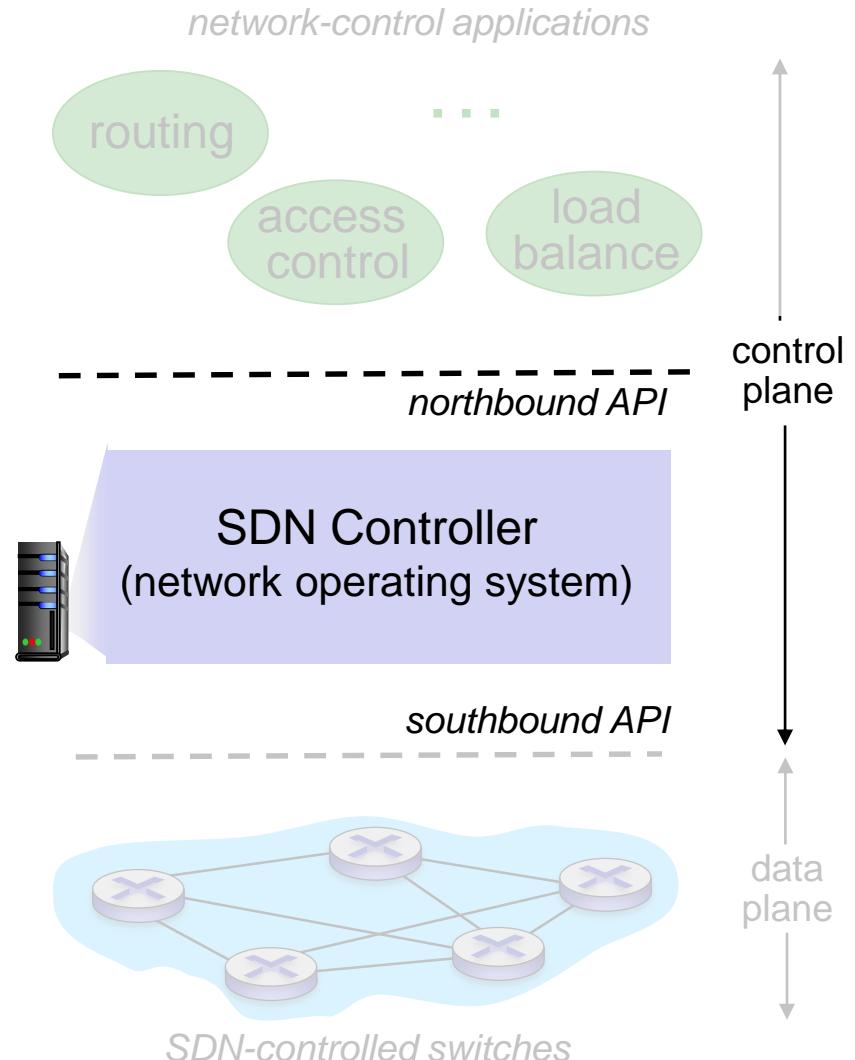
- fast, simple, commodity switches implementing generalized data-plane forwarding in hardware
- flow (forwarding) table computed, installed under controller supervision
- API for table-based switch control (e.g., OpenFlow)
 - defines what is controllable, what is not
- protocol for communicating with controller (e.g., OpenFlow)



Software defined networking (SDN)

SDN controller (network OS):

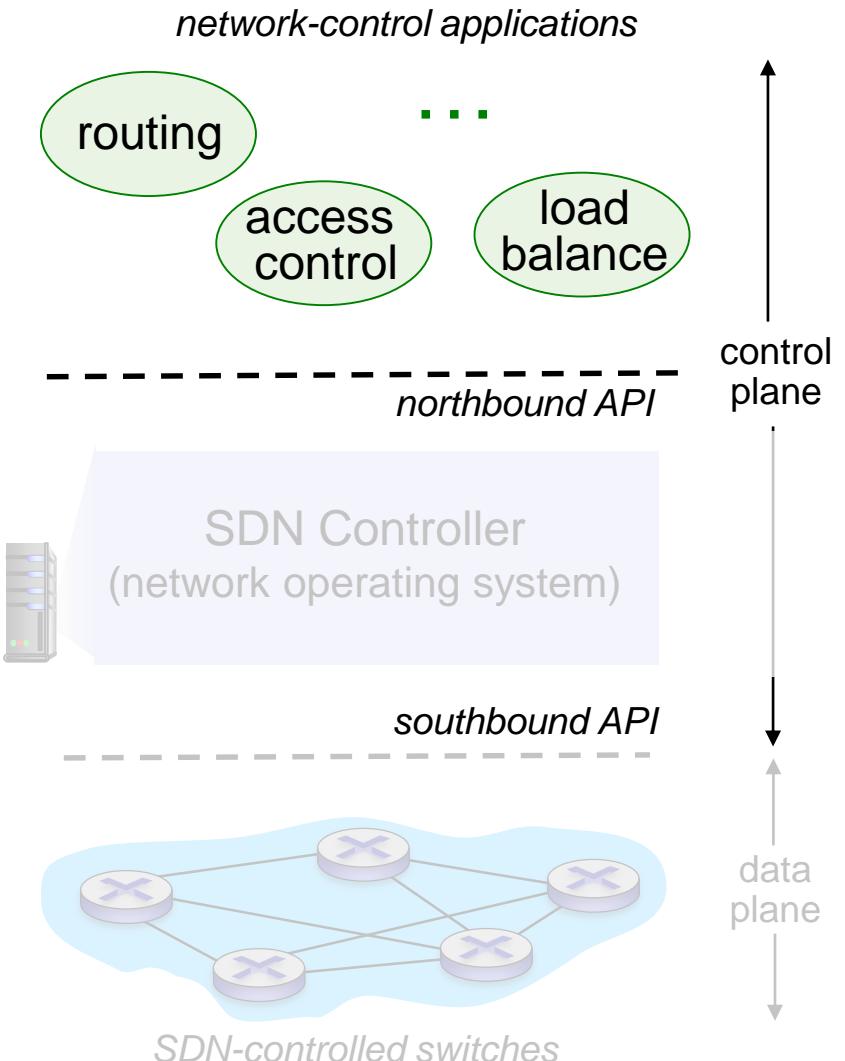
- maintain network state information
- interacts with network control applications “above” via northbound API
- interacts with network switches “below” via southbound API
- implemented as distributed system for performance, scalability, fault-tolerance, robustness



Software defined networking (SDN)

network-control apps:

- “brains” of control:
implement control functions
using lower-level services, API
provided by SDN controller
- *unbundled*: can be provided by
3rd party: distinct from routing
vendor, or SDN controller



Network layer: “control plane” roadmap

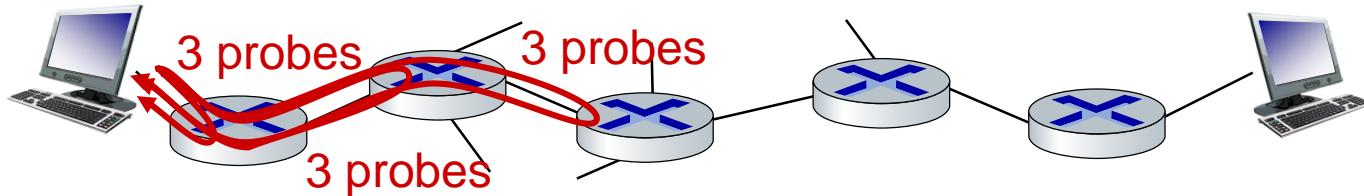
- introduction
- routing protocols
- intra-ISP routing: OSPF
- routing among ISPs: BGP
- SDN control plane
- **Internet Control Message Protocol**
- network management, configuration
 - SNMP
 - NETCONF/YANG

ICMP: internet control message protocol

- used by hosts and routers to communicate network-level information
 - error reporting: unreachable host, network, port, protocol
 - echo request/reply (used by ping)
- lies “above” IP:
 - ICMP messages carried in IP datagrams (as IP payloads)
- *ICMP message*: type, code plus first 8 bytes of IP datagram causing error

Type	Code	description
0	0	echo reply (ping)
3	0	dest. network unreachable
3	1	dest host unreachable
3	2	dest protocol unreachable
3	3	dest port unreachable
3	6	dest network unknown
3	7	dest host unknown
4	0	source quench (congestion control - not used)
8	0	echo request (ping)
9	0	route advertisement
10	0	router discovery
11	0	TTL expired
12	0	bad IP header

Traceroute and ICMP



- Traceroute program traces a route from a host to any other host, is implemented with ICMP messages
- source sends sets of UDP segments to destination
 - 1st set has TTL =1, 2nd set has TTL=2, etc.
- datagram in n th set arrives to n th router:
 - router discards datagram and sends source ICMP message (type 11, code 0)
 - ICMP message possibly includes name of router & IP address
- when ICMP message arrives at source: record RTTs

stopping criteria:

- UDP segment eventually arrives at destination host
- destination returns ICMP “port unreachable” message (type 3, code 3)
- source stops

Network layer: “control plane” roadmap

- introduction
- routing protocols
- intra-ISP routing: OSPF
- routing among ISPs: BGP
- SDN control plane
- Internet Control Message Protocol
- network management, configuration
 - SNMP
 - NETCONF/YANG

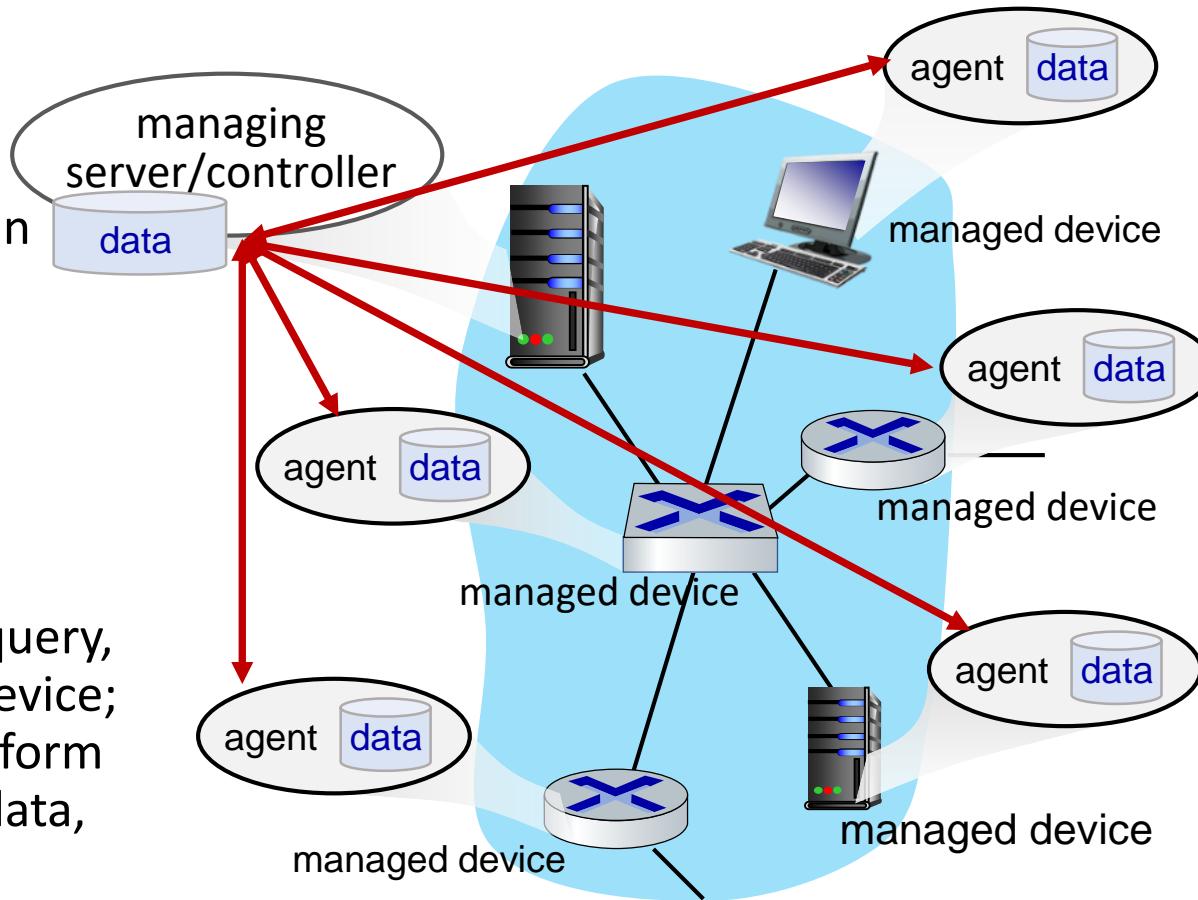
What is network management?

"**Network management** includes the deployment, integration and coordination of the hardware, software, and human elements to monitor, test, poll, configure, analyze, evaluate, and control the network and element resources to meet the real-time, operational performance, and Quality of Service requirements at a reasonable cost."

Components of network management

Managing server:
application, typically
with network
managers (humans) in
the loop

**Network
management
protocol:** used by
managing server to query,
configure, manage device;
used by devices to inform
managing server of data,
events.



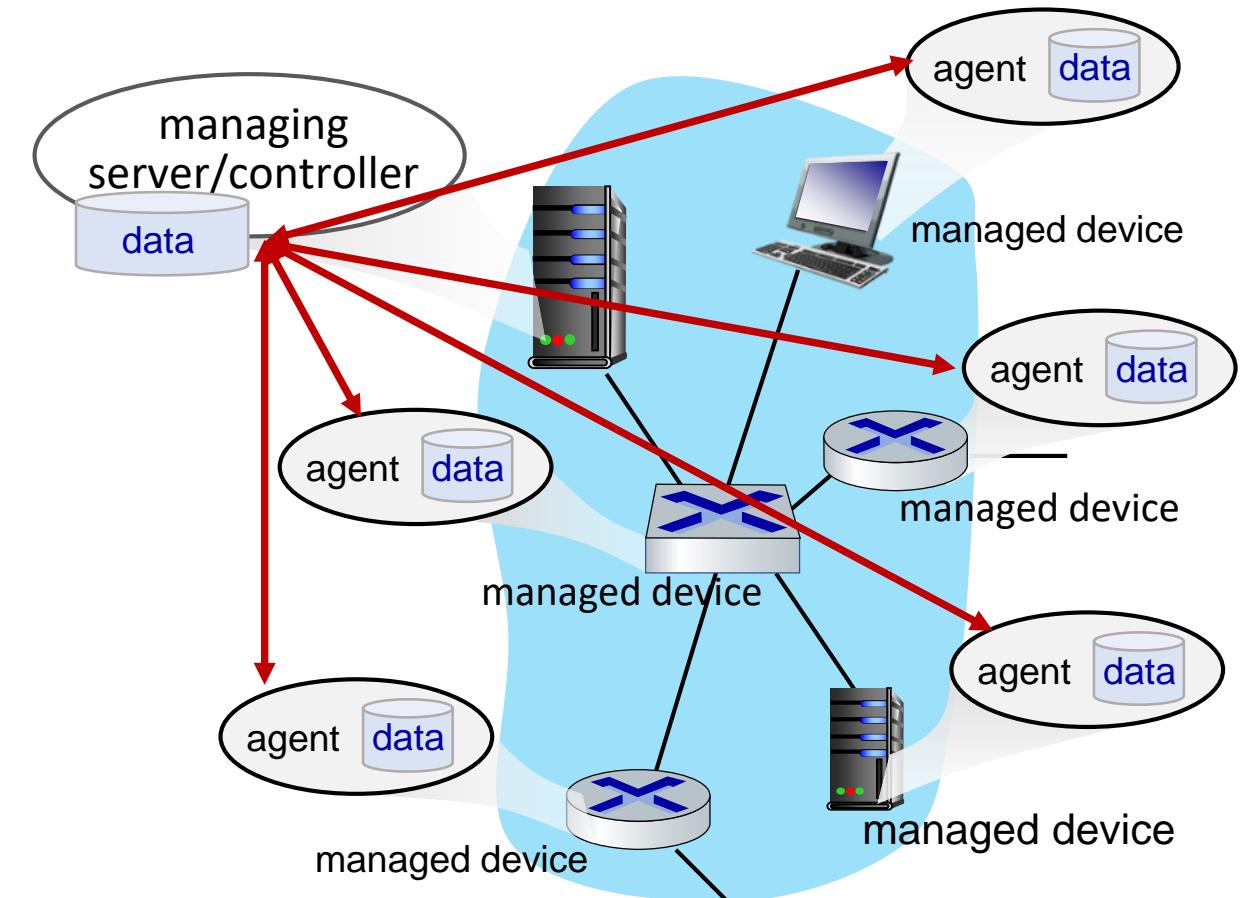
Managed device:
equipment with manageable,
configurable hardware,
software components

Data: device “state”
configuration data,
operational data,
device statistics

Network operator approaches to management

CLI (Command Line Interface)

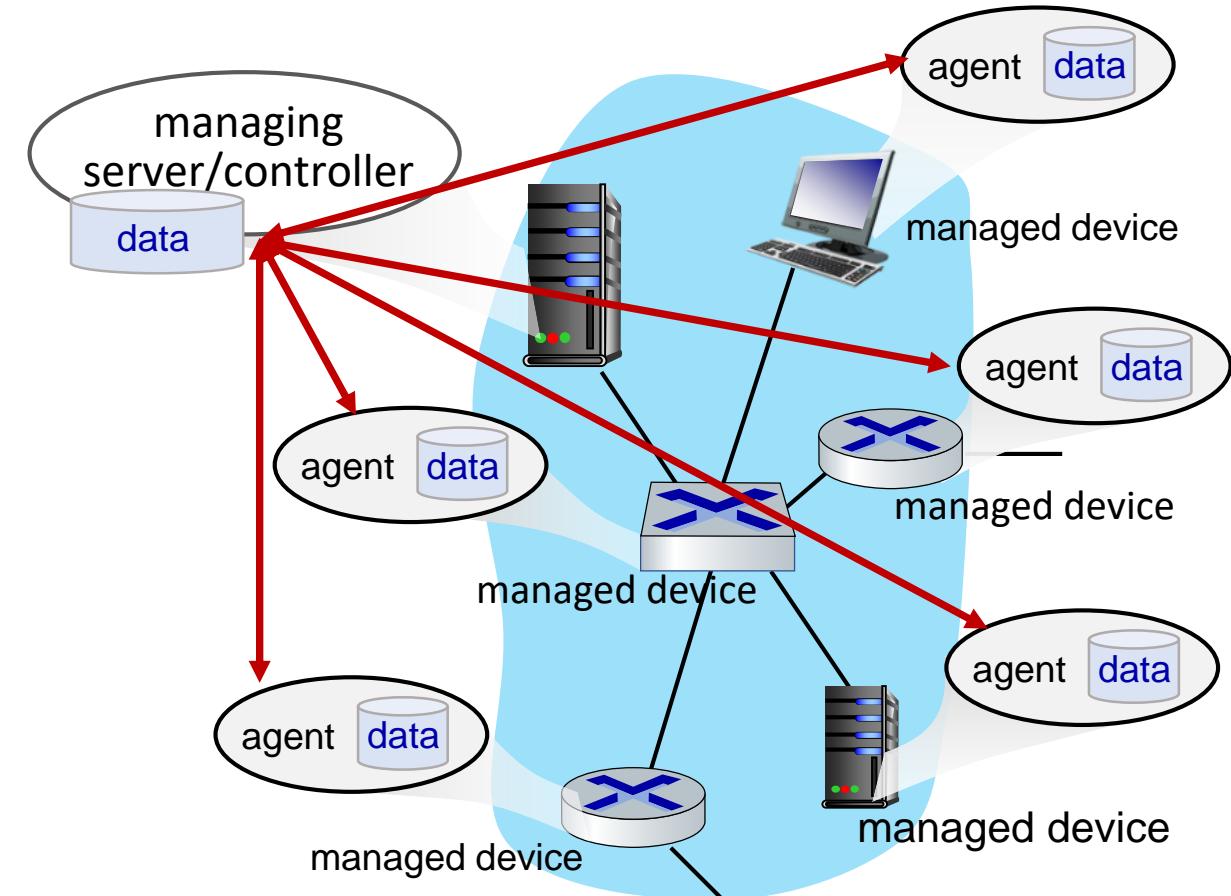
- operator issues (types, scripts) direct to individual devices (e.g. ssh)



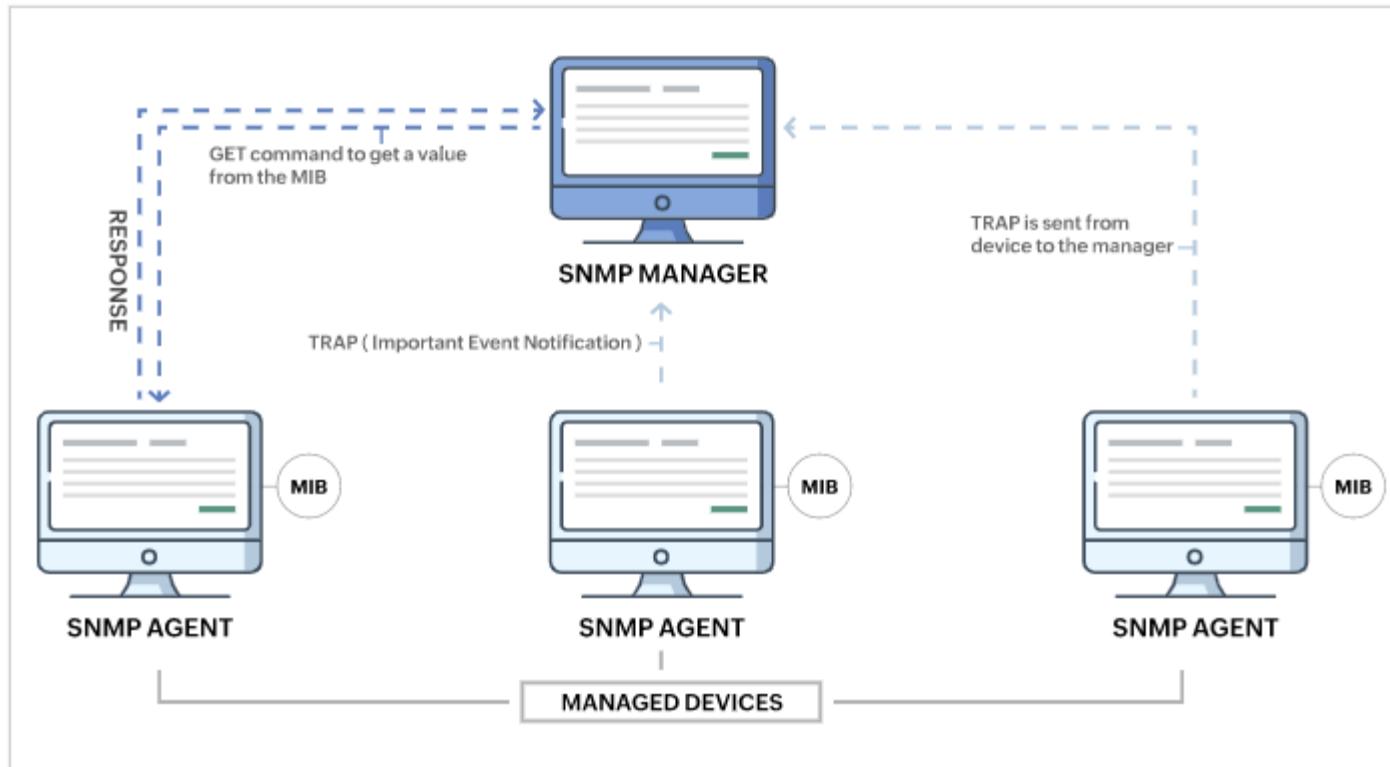
Network operator approaches to management

SNMP/MIB

- Simple Network Management Protocol (SNMP): an application–layer protocol, a network operator can query/set the data contained in a device’s Management Information Base (MIB) objects
- Management Information Base (MIB): a hierarchical virtual database of network objects, describes a device being monitored by a network management system
- MIB examples: on a printer; the different cartridge states and the number of printed files, on a switch; incoming and outgoing traffic, rate of package loss, number of packets addressed to a broadcast address ...



SNMP: Management Information Base (MIB)

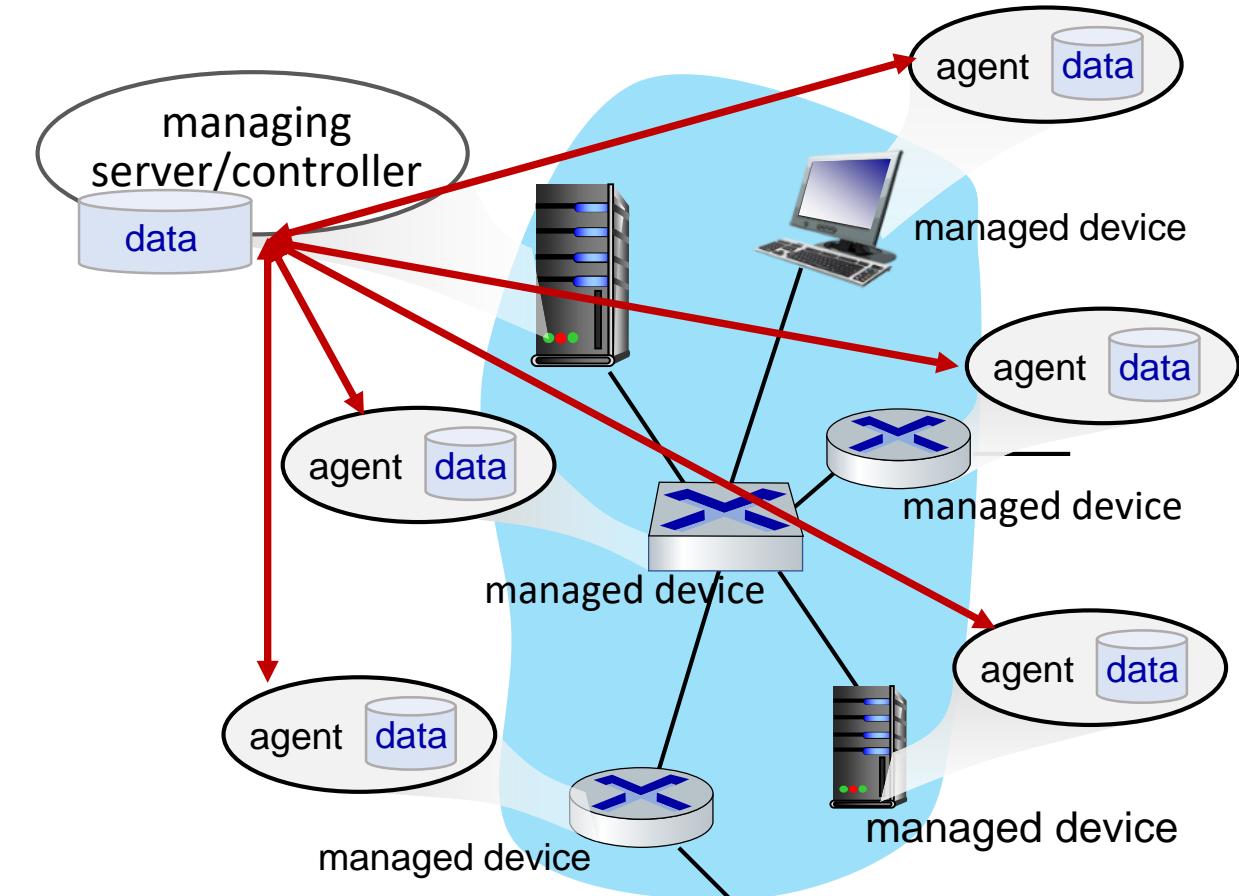


- SNMP tools (NMS and agents) like SolarWinds, Paessler (PRTG)
- SNMP uses UDP in TL, IP in NL
- SNMP commands help retrieve, manage, modify, and parse the data
- SNMP manager initiates the commands
- SNMP agents initiate the TRAPS command. (TRAPS is a signal sent to the manager by the agent when events occur)

Network operator approaches to management

NETCONF/YANG

- NETCONF is the standard for installing, manipulating and deleting configuration of network devices
- YANG is used to model both configuration and state data of network elements, modelling language representing data structures in an XML tree format, can be converted into any encoding format, e.g. XML or JSON



NETCONF overview

- **goal:** actively manage/configure devices network-wide
- operates between managing server and managed network devices
 - actions: retrieve, set, modify, activate configurations
 - query operational data and statistics
 - subscribe to notifications from devices
- uses a simple remote procedure call (RPC) mechanism (to implement communication between a client and a server)
 - NETCONF protocol messages encoded in XML
 - exchanged over secure, reliable transport (e.g., TLS) protocol
- standard application programming interfaces (APIs) are available on network devices for the NMS to manage the devices using NETCONF
- runs primarily over Secure Shell (SSH) transport

Sample NETCONF RPC message

```
01 <?xml version="1.0" encoding="UTF-8"?>
02 <rpc message-id="101" note message id
03   xmlns="urn:ietf:params:xml:ns:netconf:base:1.0">
04     <edit-config> change a configuration
05       <target>
06         <running/> change the running configuration
07       </target>
08     <config>
09       <top xmlns="http://example.com/schema/
1.2/config">
10         <interface>
11           <name>Ethernet0/0</name> change MTU of Ethernet 0/0 interface to 1500
12           <mtu>1500</mtu>
13         </interface>
14       </top>
15     </config>
16   </edit-config>
17 </rpc>
```