

Chapter 5

Network Layer:

The Control Plane

A note on the use of these Powerpoint slides:

We're making these slides freely available to all (faculty, students, readers). They're in PowerPoint form so you see the animations; and can add, modify, and delete slides (including this one) and slide content to suit your needs. They obviously represent a *lot* of work on our part. In return for use, we only ask the following:

- If you use these slides (e.g., in a class) that you mention their source (after all, we'd like people to use our book!)
- If you post any slides on a www site, that you note that they are adapted from (or perhaps identical to) our slides, and note our copyright of this material.

Thanks and enjoy! JFK/KWR

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



Computer Networking: A Top Down Approach

7th edition

Jim Kurose, Keith Ross

Pearson/Addison Wesley

April 2016

Network Layer (Routing)

Network-layer functions

Recall: two network-layer functions:

- *forwarding*: move packets from router's input to appropriate router output

data plane

- *routing*: determine route taken by packets from source to destination

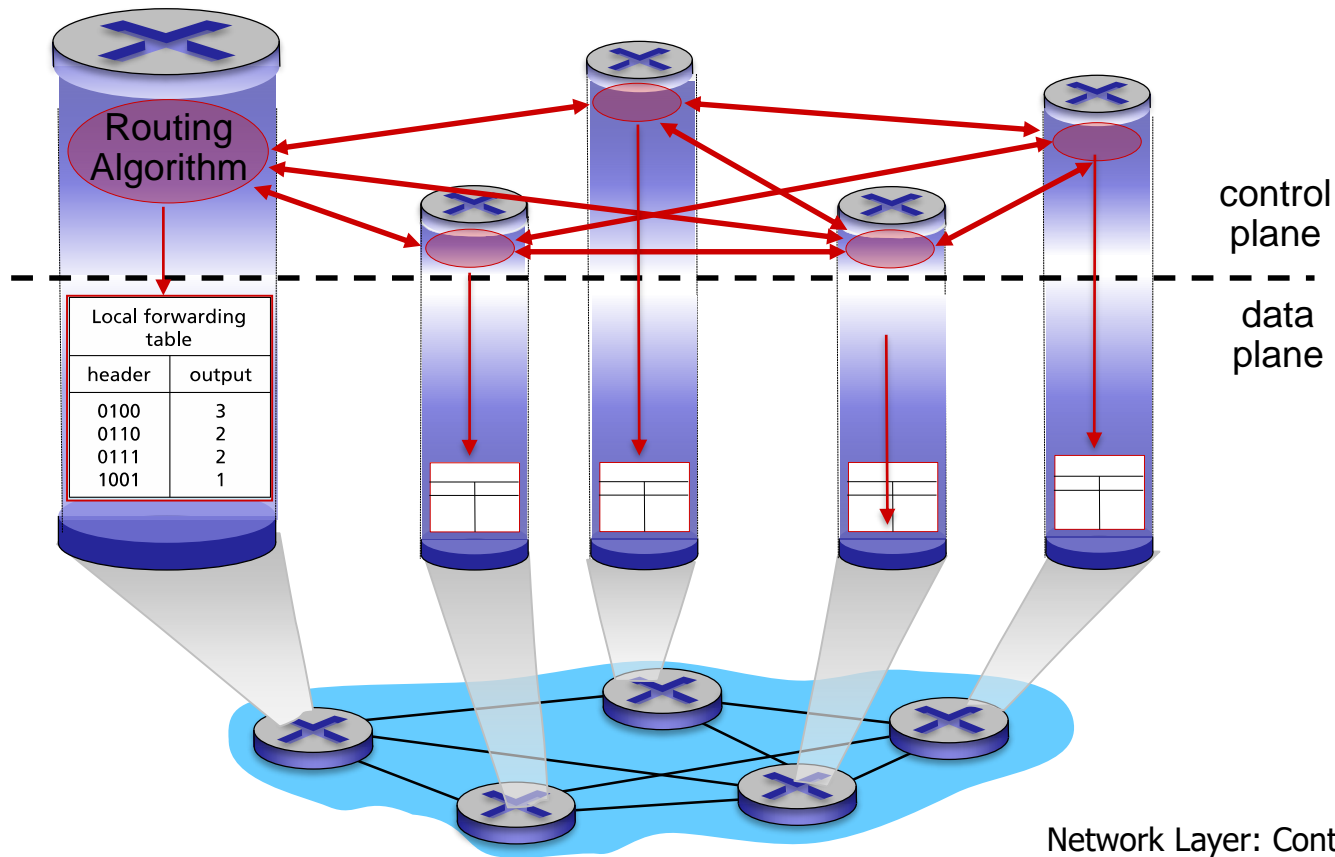
control plane

Two approaches to structuring network control plane:

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

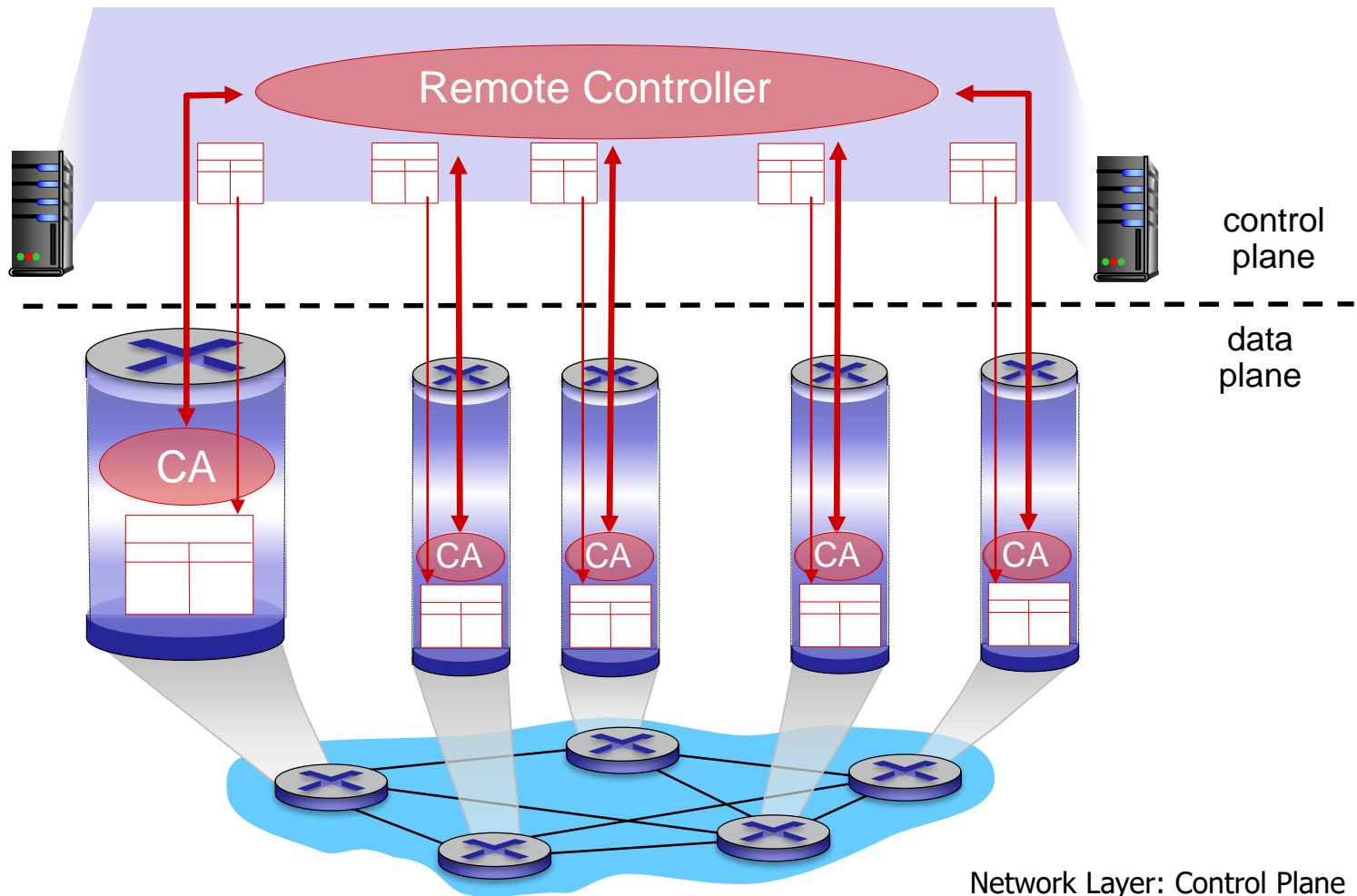
Per-router control plane

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



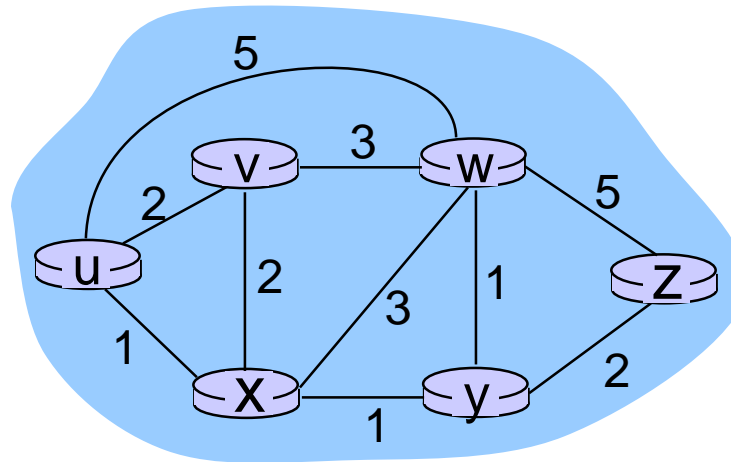
Routing Protocols

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!

Graph abstraction of the network

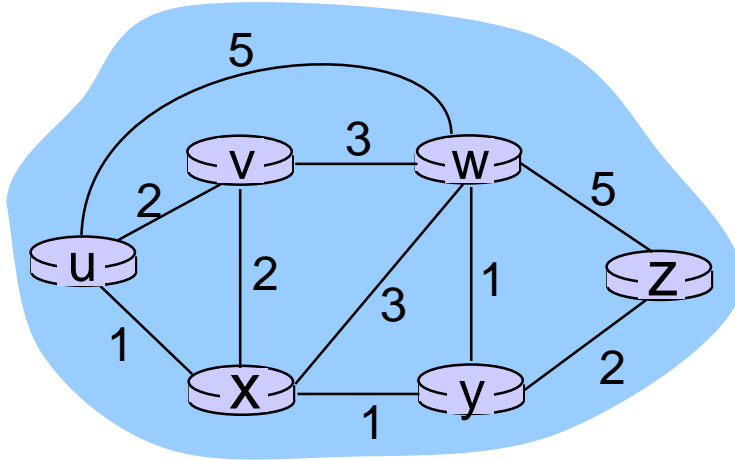


graph: $G = (N, E)$

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

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

Graph abstraction: costs



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

cost could always be 1, 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

Routing algorithm classification

Q: global or decentralized information?

global:

- all routers have complete topology, link cost info
- “link state” algorithms

decentralized:

- router knows physically-connected neighbors, link costs to neighbors
- iterative process of computation, exchange of info with neighbors
- “distance vector” algorithms

Q: static or dynamic?

static:

- routes change slowly over time

dynamic:

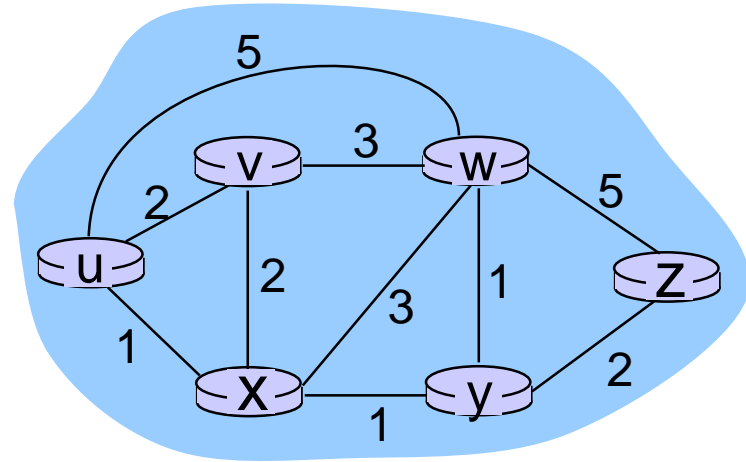
- routes change more quickly
 - periodic update
 - in response to link cost changes

Link State Routing

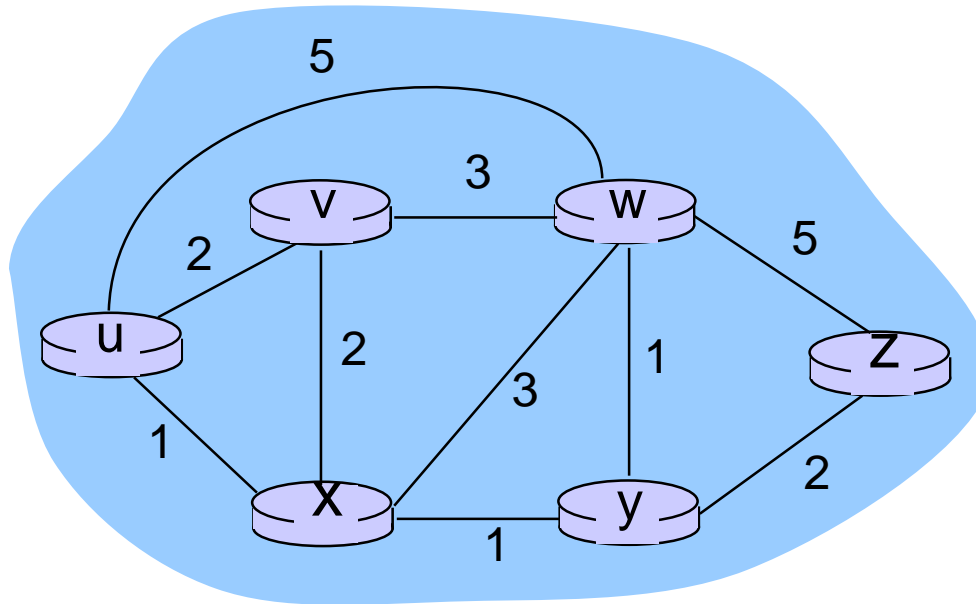
A link-state routing algorithm

Dijkstra's algorithm

- net topology, link costs known to all nodes
 - accomplished via “link state broadcast”
 - all nodes have same info
- computes least cost paths from one node (“source”) to all other nodes
 - gives *forwarding table* for that node
- iterative: after k iterations, know least cost path to k dest.'s

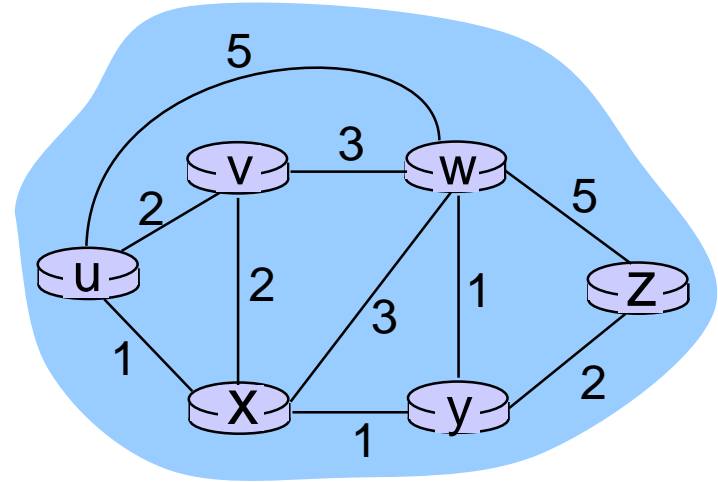
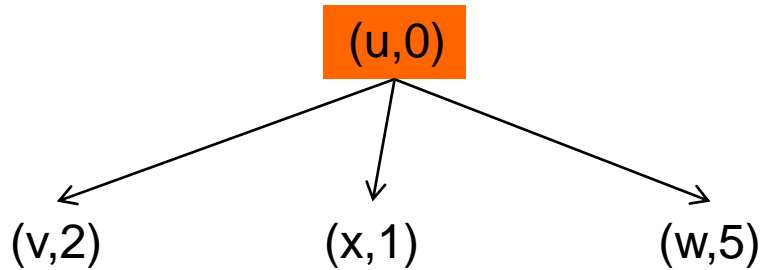


Dijkstra's algorithm



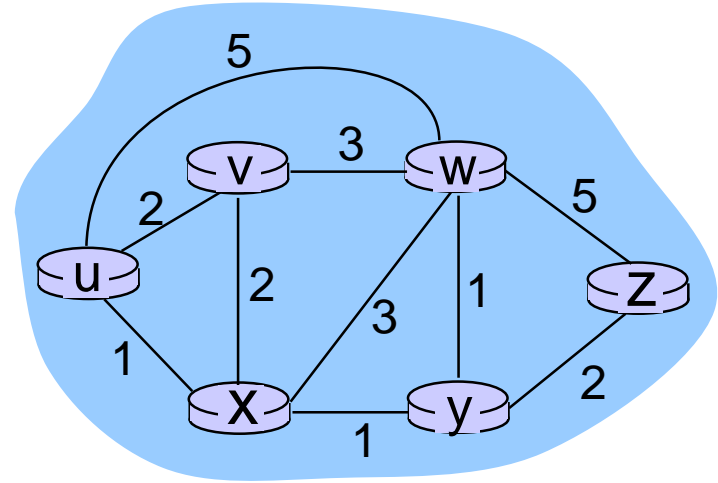
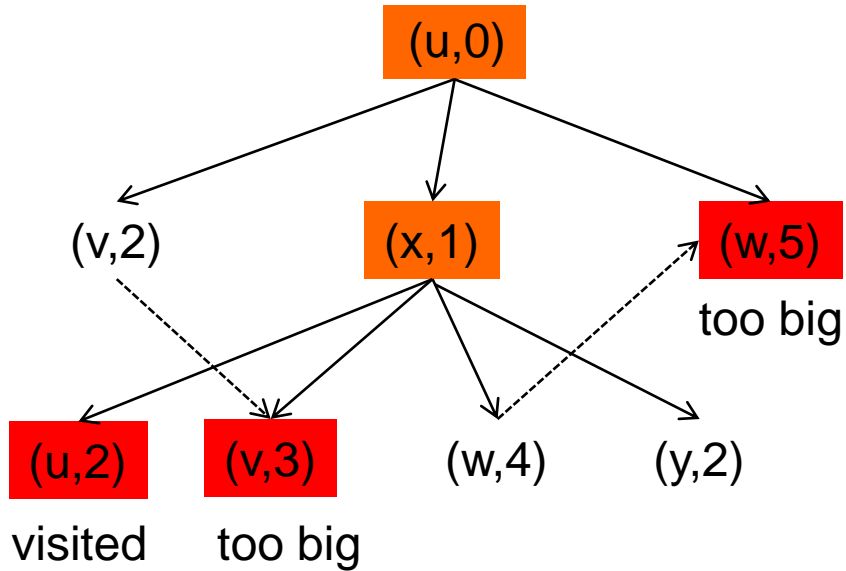
Let us try this starting from u!

Dijkstra's algorithm



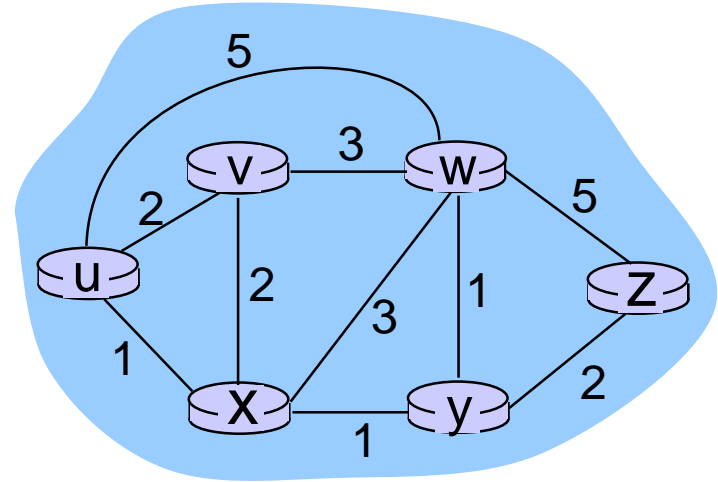
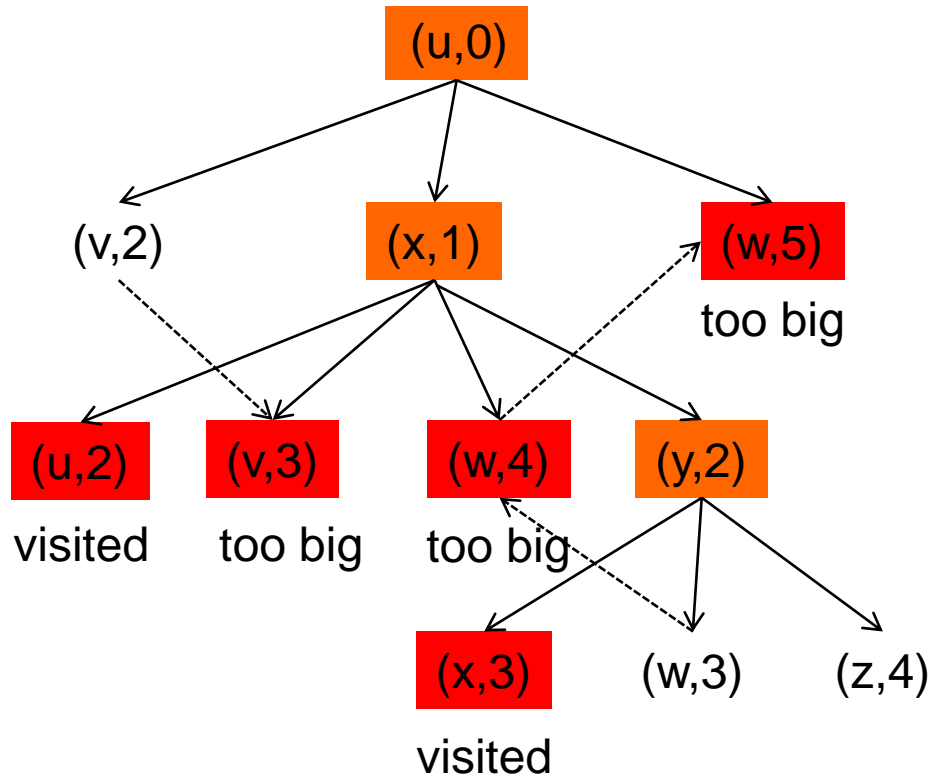
Q. Which neighbors of u has the “least cost” path?

Dijkstra's algorithm



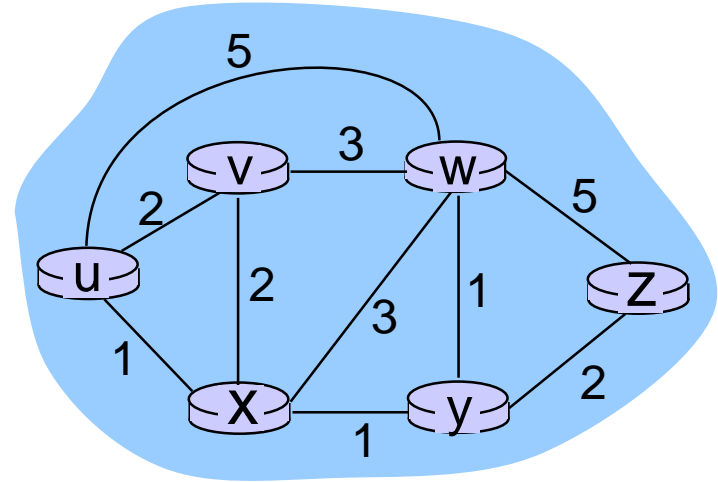
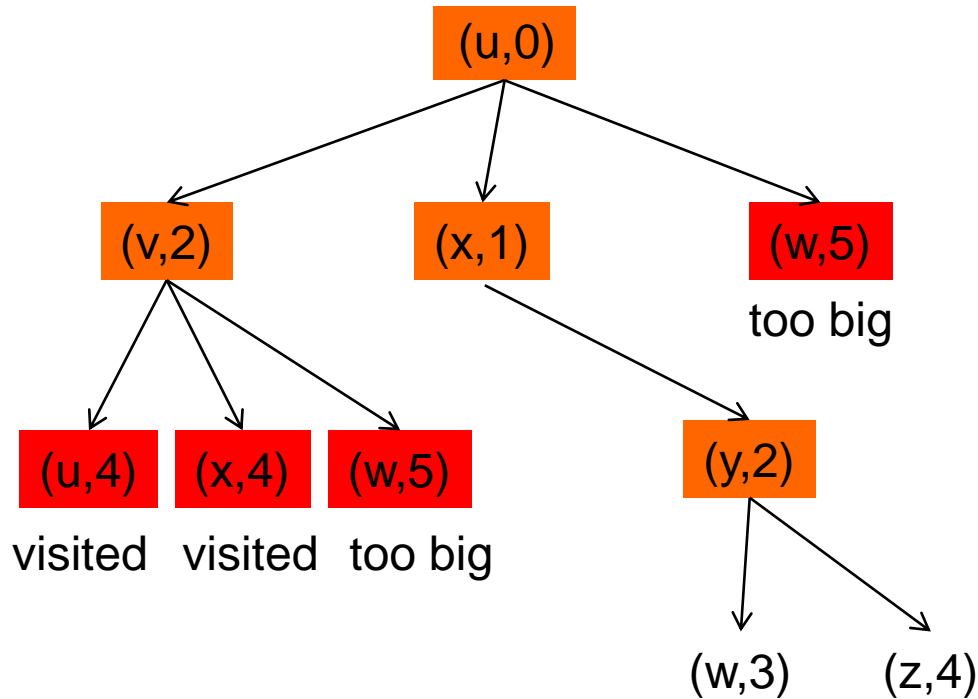
Q. Which unexplored node has the “least cost” path?

Dijkstra's algorithm



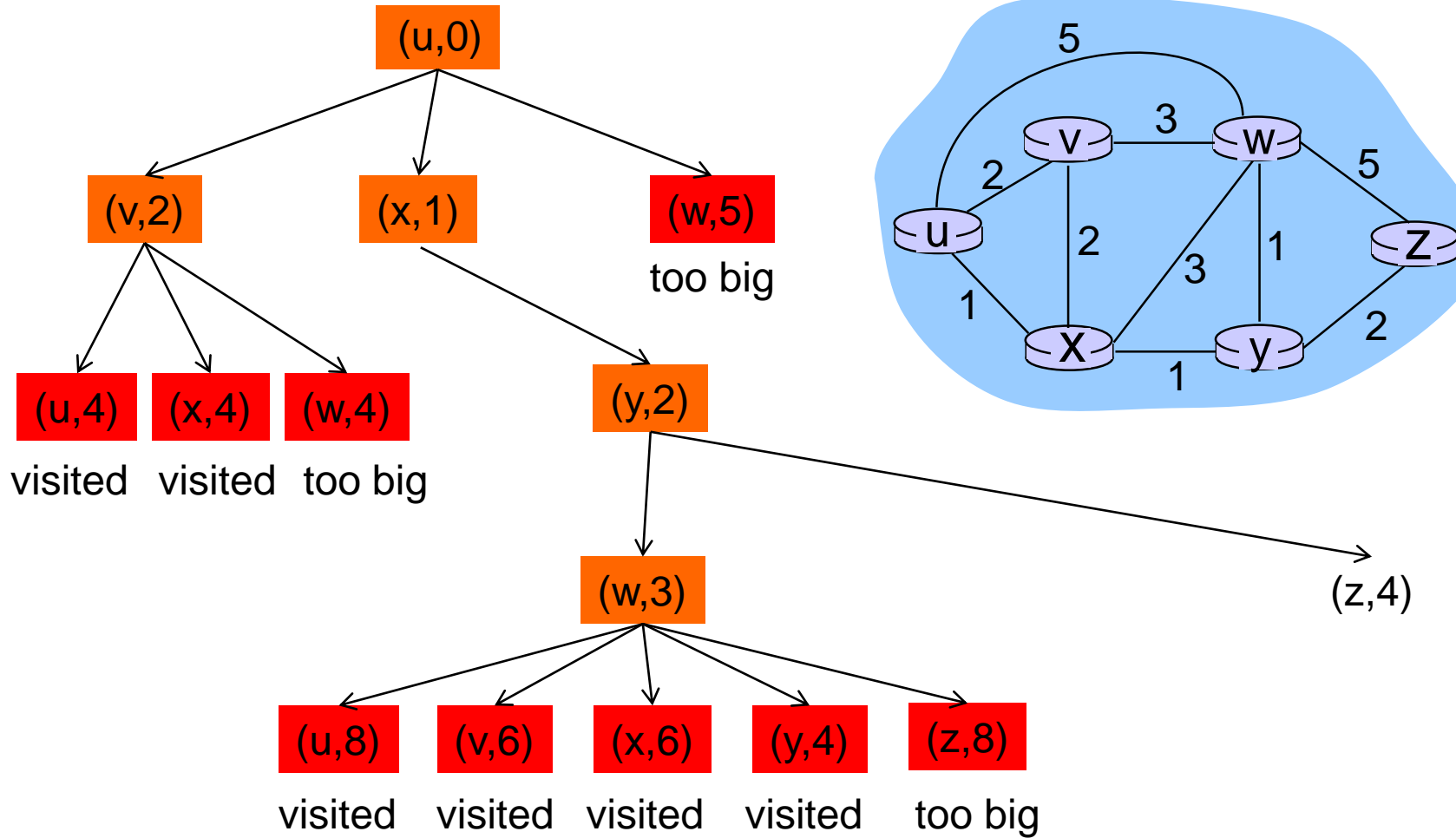
Q. Which unexplored node has the “least cost” path?

Dijkstra's algorithm

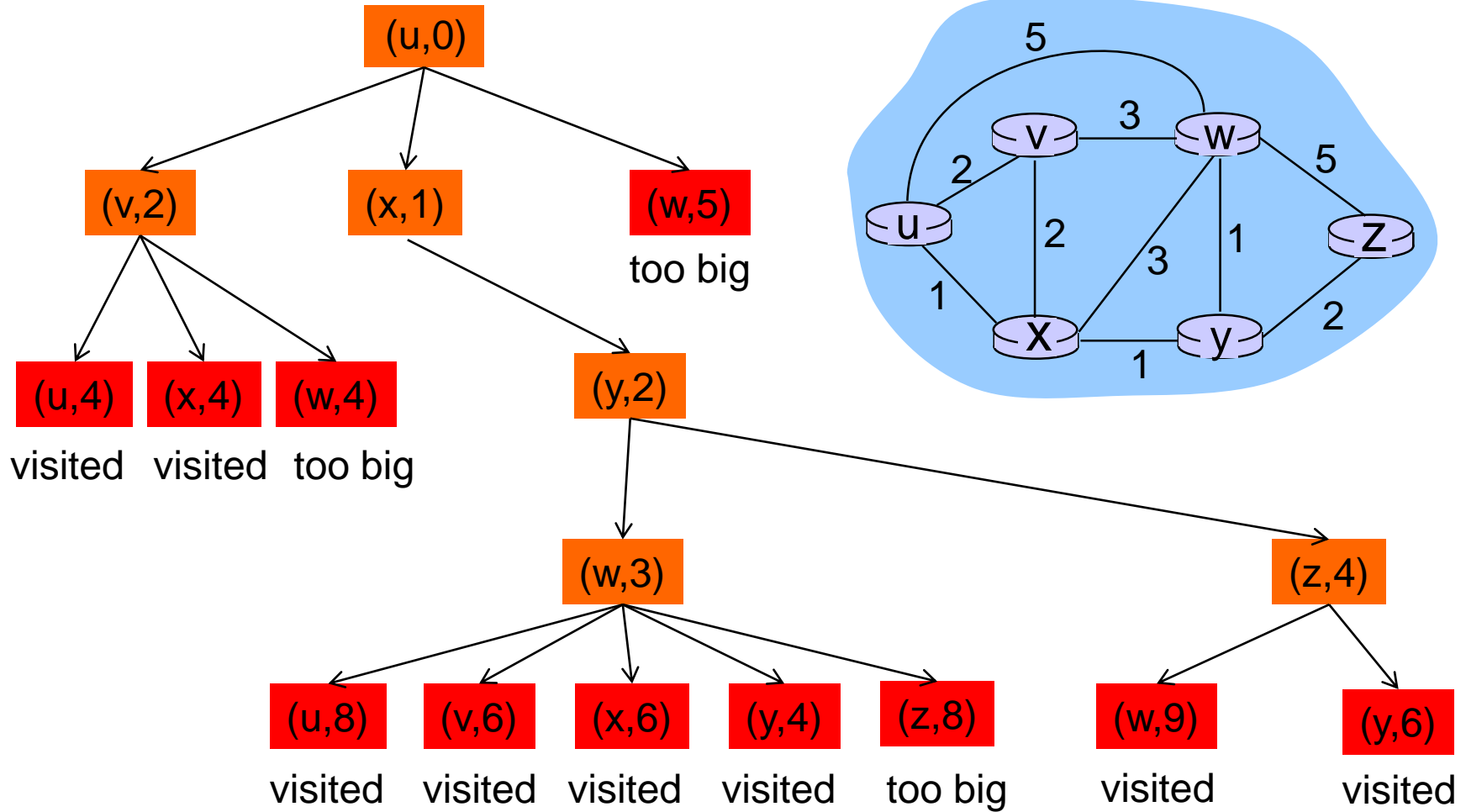


Q. Which unexplored node has the “least cost” path?

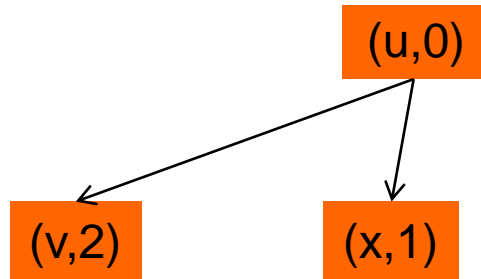
Dijkstra's algorithm



Dijkstra's algorithm

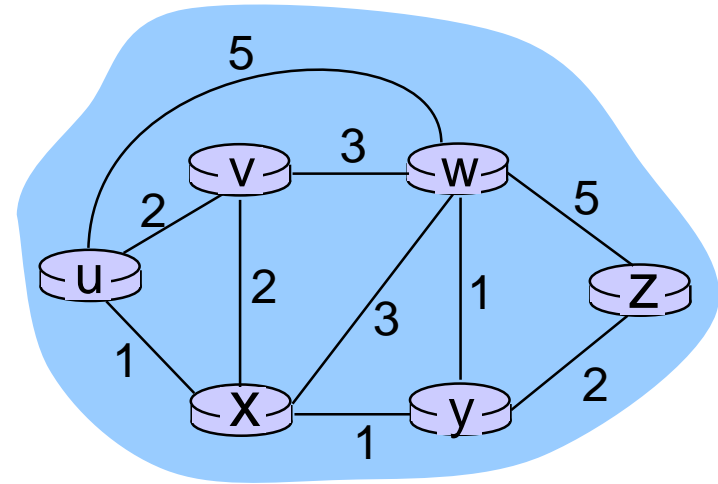
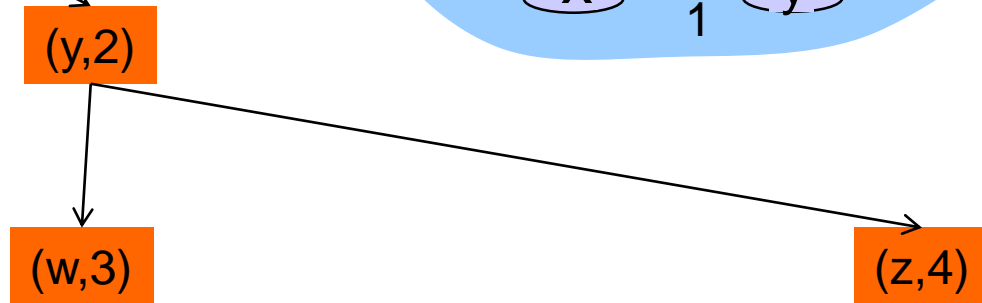


Dijkstra's algorithm



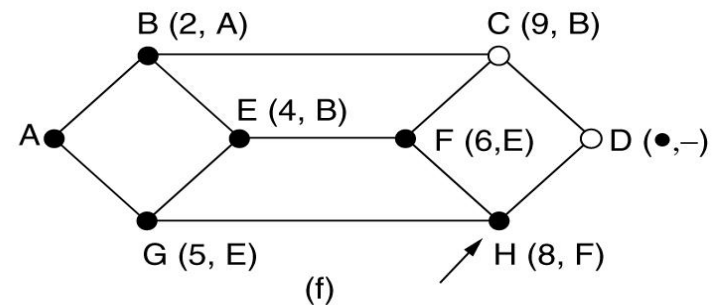
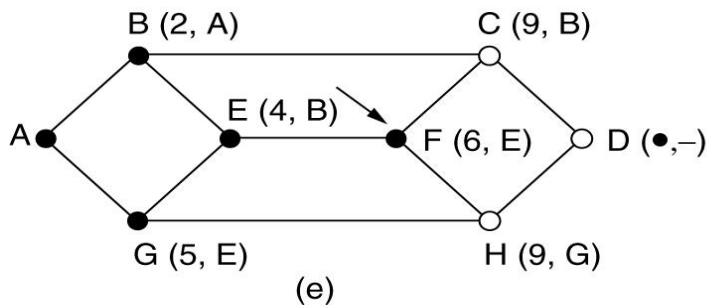
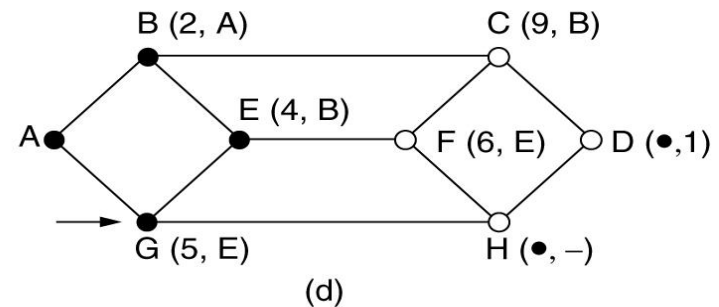
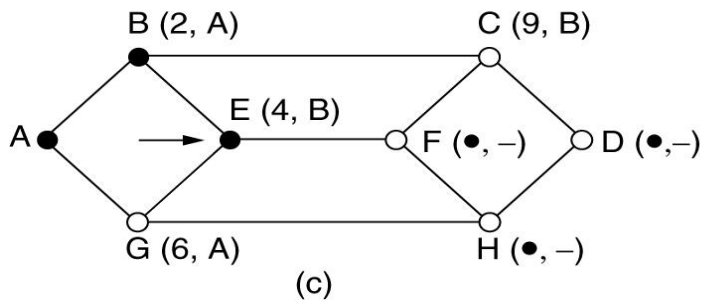
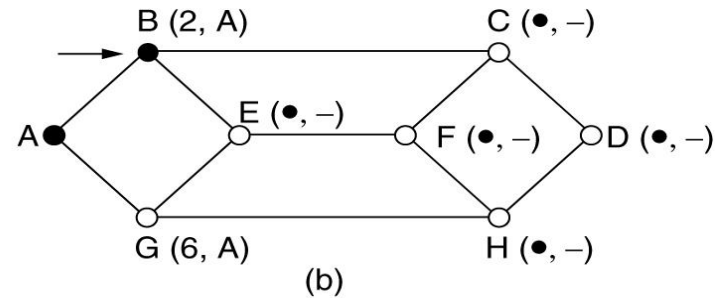
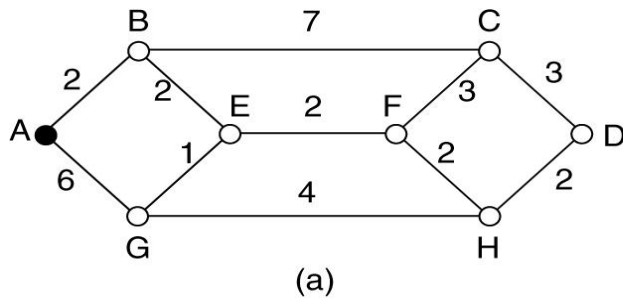
Forwarding table for u:

Node	Cost	Next Hop
u	0	-
v	2	v
w	3	x
x	1	x
y	2	x
z	4	x



Try this starting from z!

Dijkstra's algorithm (2nd example)



The arrows indicate the working node. (#,n) where # is the cost and n is the parent.

Distance Vector Routing

Distance vector algorithm

Bellman-Ford equation (dynamic programming)

let

$D(x,y)$ = cost of **least-cost path** from x to y

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

then

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

cost from neighbor v to destination y

cost to neighbor v

min taken over all neighbors v of x

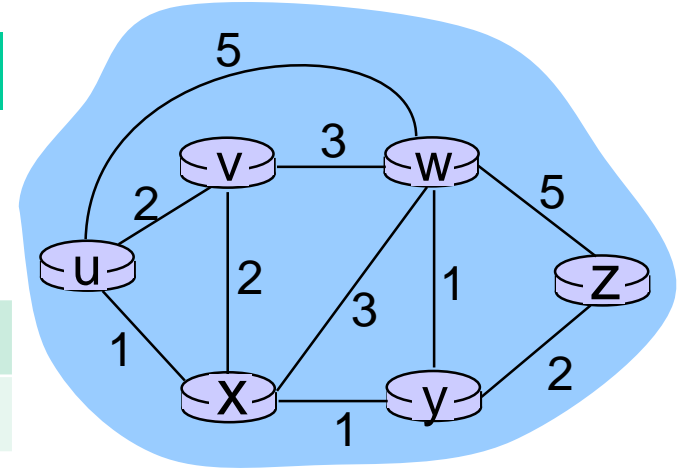
Distance vector algorithm

$D(n)$ means least cost paths from "n" to all nodes.

$D(u)$		$D(v)$					
		u	v	w	x	y	z
u	0	2	0	3	2	?	?
v	2						
w	5						
x	1						
y	?						
z	?						

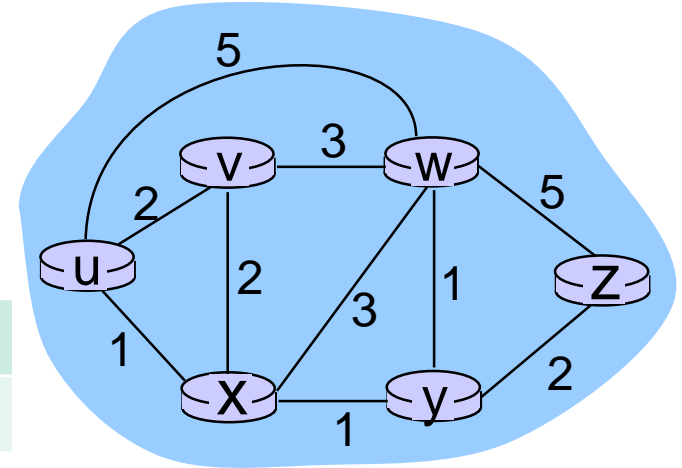
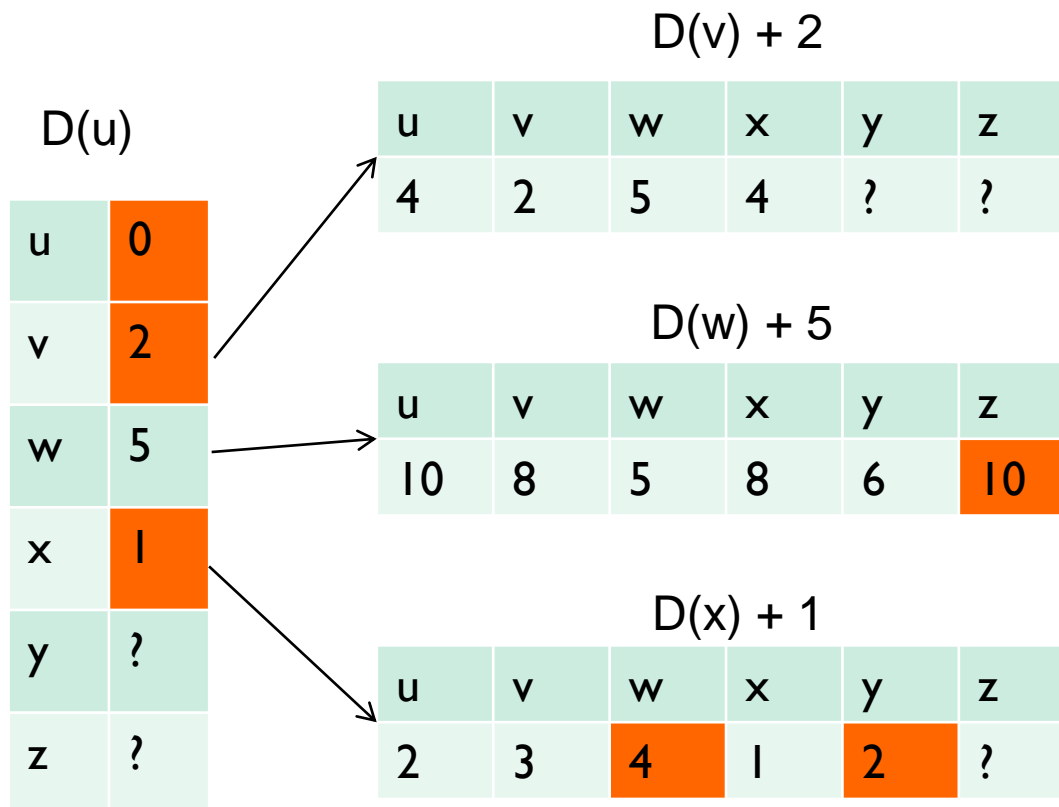
$D(w)$		u	v	w	x	y	z
u	5	5	3	0	3	1	5
v							
w							
x							
y							
z							

$D(x)$		u	v	w	x	y	z
u	1	1	2	3	0	1	?
v							
w							
x							
y							
z							



Q. What is the new $D(u)'$ if u receives $D(v)$, $D(w)$ and $D(x)$?

Distance vector algorithm

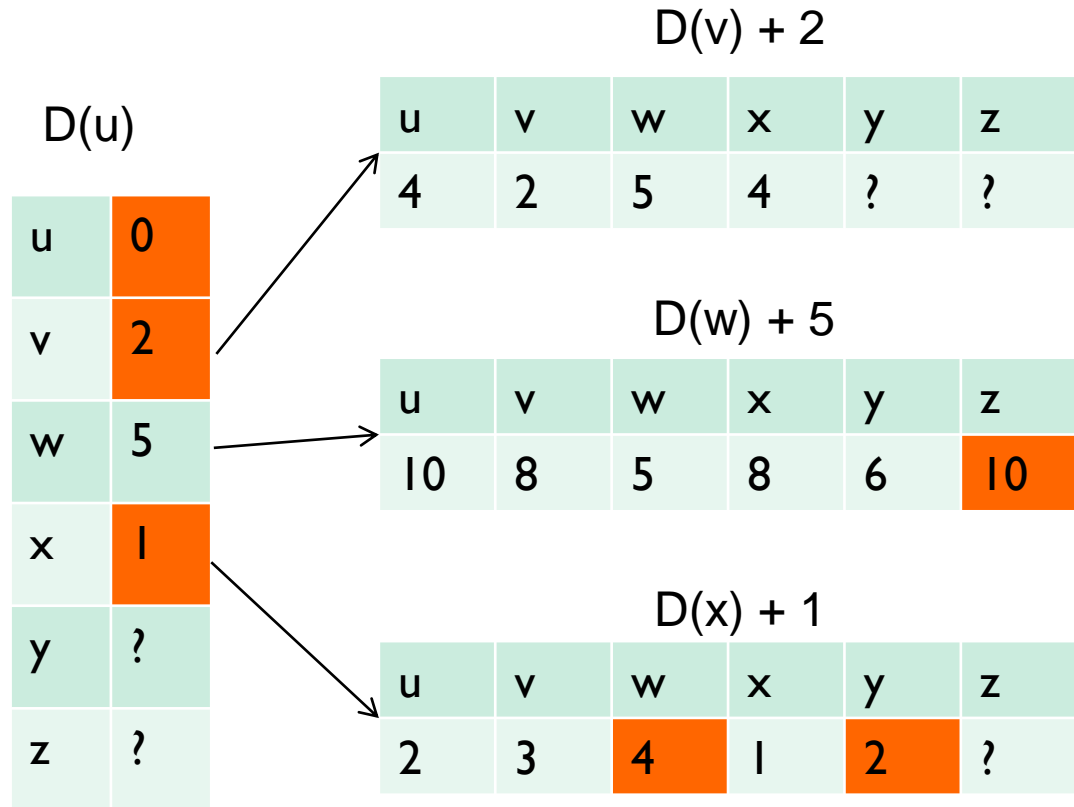


Distance is the smallest among all DVs

Distance vector algorithm

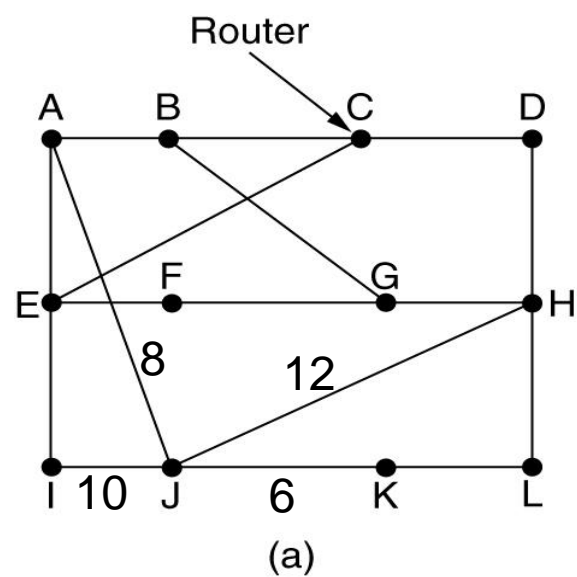
Forwarding table for u
after 1st round

To	DV	Next Hop
u	0	-
v	2	v
w	4	x
x	1	x
y	2	x
z	10	w



“u” sends its new DV(u)’ to all its neighbours “v”, “w”, and “x”, which then update their DVs accordingly.

Distance vector algorithm 2nd example



J receives 4 updates from A,I,H and K. It updates its DV table with next hop. It then transmits its new estimated to all its neighbors.

					New estimated delay from J	
To	A	I	H	K	Line	
A	0	24	20	21	8	A
B	12	36	31	28	20	A
C	25	18	19	36	28	I
D	40	27	8	24	20	H
E	14	7	30	22	17	I
F	23	20	19	40	30	I
G	18	31	6	31	18	H
H	17	20	0	19	12	H
I	21	0	14	22	10	I
J	9	11	7	10	0	–
K	24	22	22	0	6	K
L	29	33	9	9	15	K

JA delay is 8	JI delay is 10	JH delay is 12	JK delay is 6
---------------	----------------	----------------	---------------

Vectors received from J's four neighbors

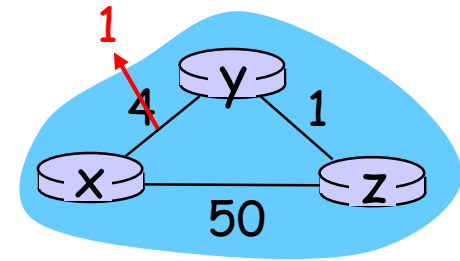
New routing table for J	
8	A
20	A
28	I
20	H
17	I
30	I
18	H
12	H
10	I
0	–
6	K
15	K

(b)

Distance vector: link cost changes

link cost changes (decrease):

- ❖ node detects local link cost change
- ❖ updates routing info, recalculates distance vector
- ❖ if DV changes, notify neighbors



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.

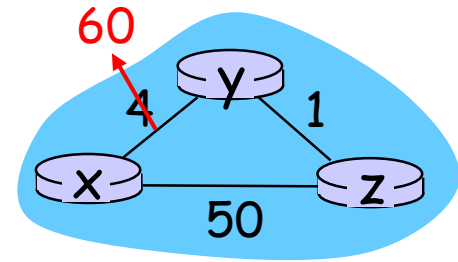
t_2 : y receives z's update, updates its distance table. y's least costs do *not* change, so y does *not* send a message to z.

“good news travels fast”

Distance vector: link cost changes

link cost changes (increase):

- ❖ node detects local link cost change
- ❖ “count to infinity” problem!
- ❖ 44 iterations before algorithm stabilizes



“bad news travels slow”

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?
NO!

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

Internet Routing

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 admin may want to control routing in its own network

Internet approach to scalable routing

aggregate routers into regions known as
“Autonomous Systems” (**AS**) (a.k.a. “domains”)

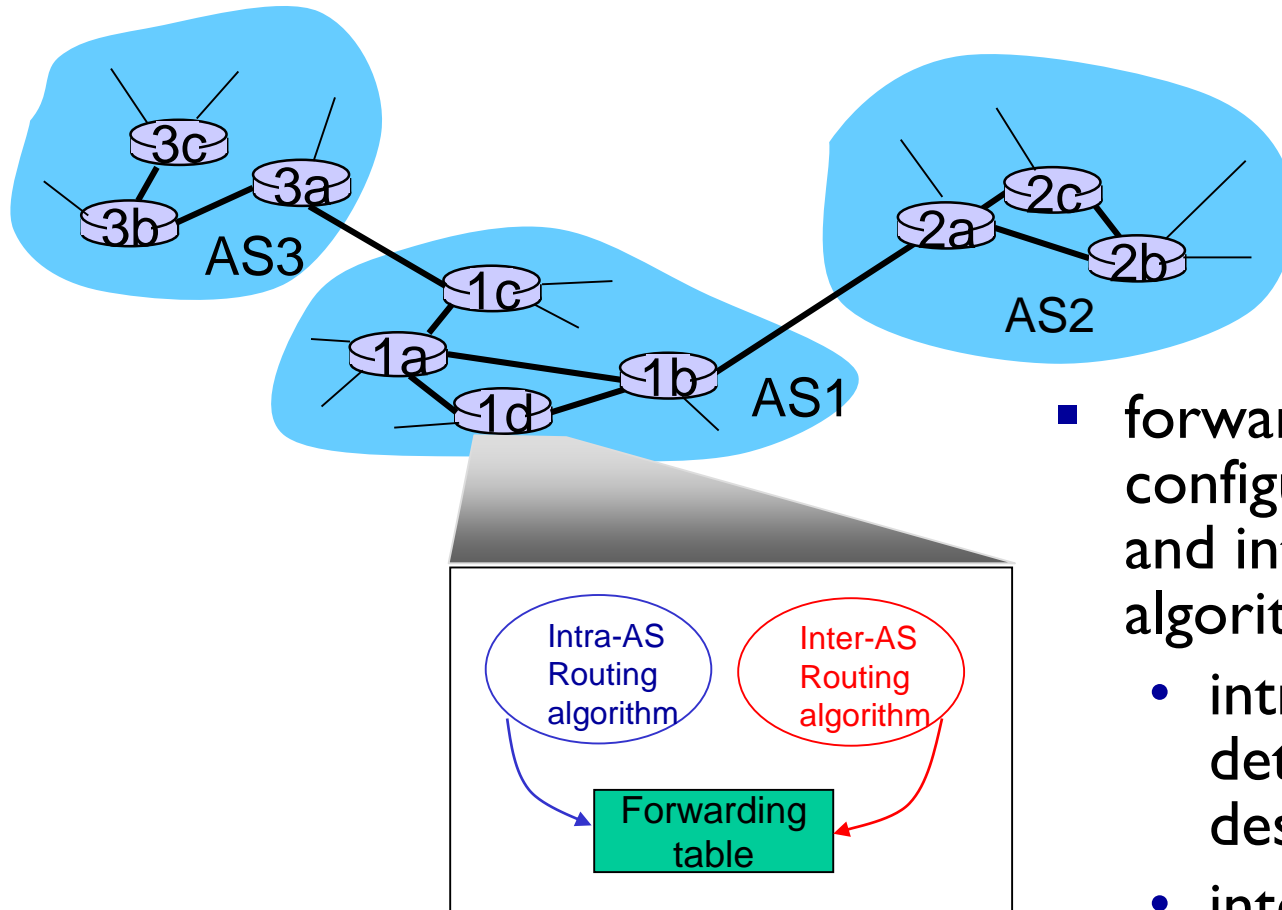
intra-AS routing

- 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 protocol
- 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)

Interconnected ASes



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

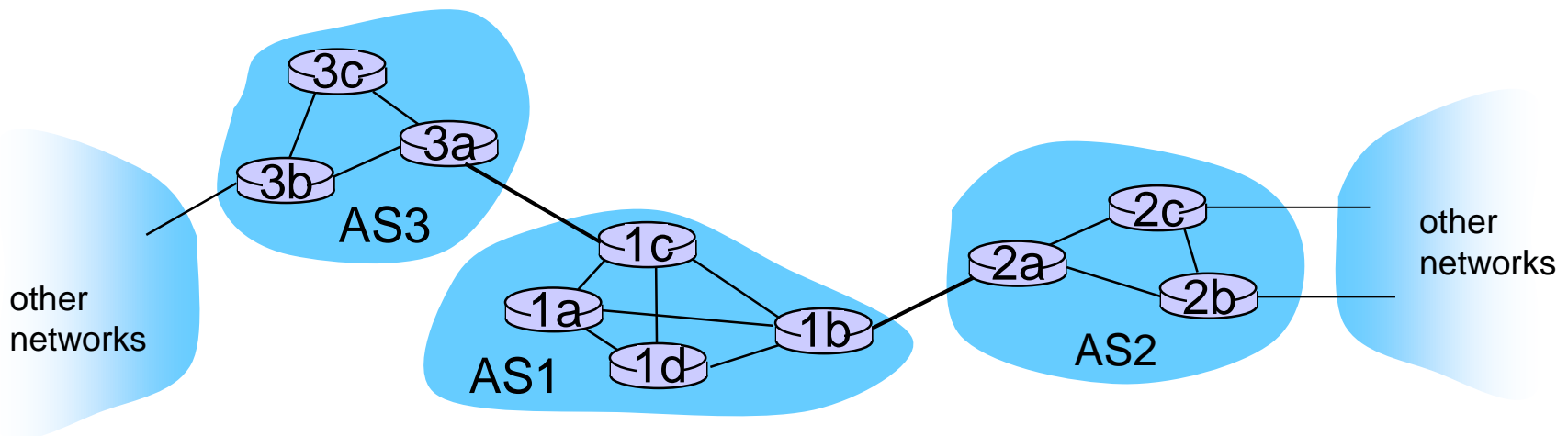
Inter-AS tasks

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

job of inter-AS routing!



Intra-AS Routing

- also known as *interior gateway protocols (IGP)*
- most common intra-AS routing protocols:
 - RIP: Routing Information Protocol
 - OSPF: Open Shortest Path First
 - IGRP: Interior Gateway Routing Protocol

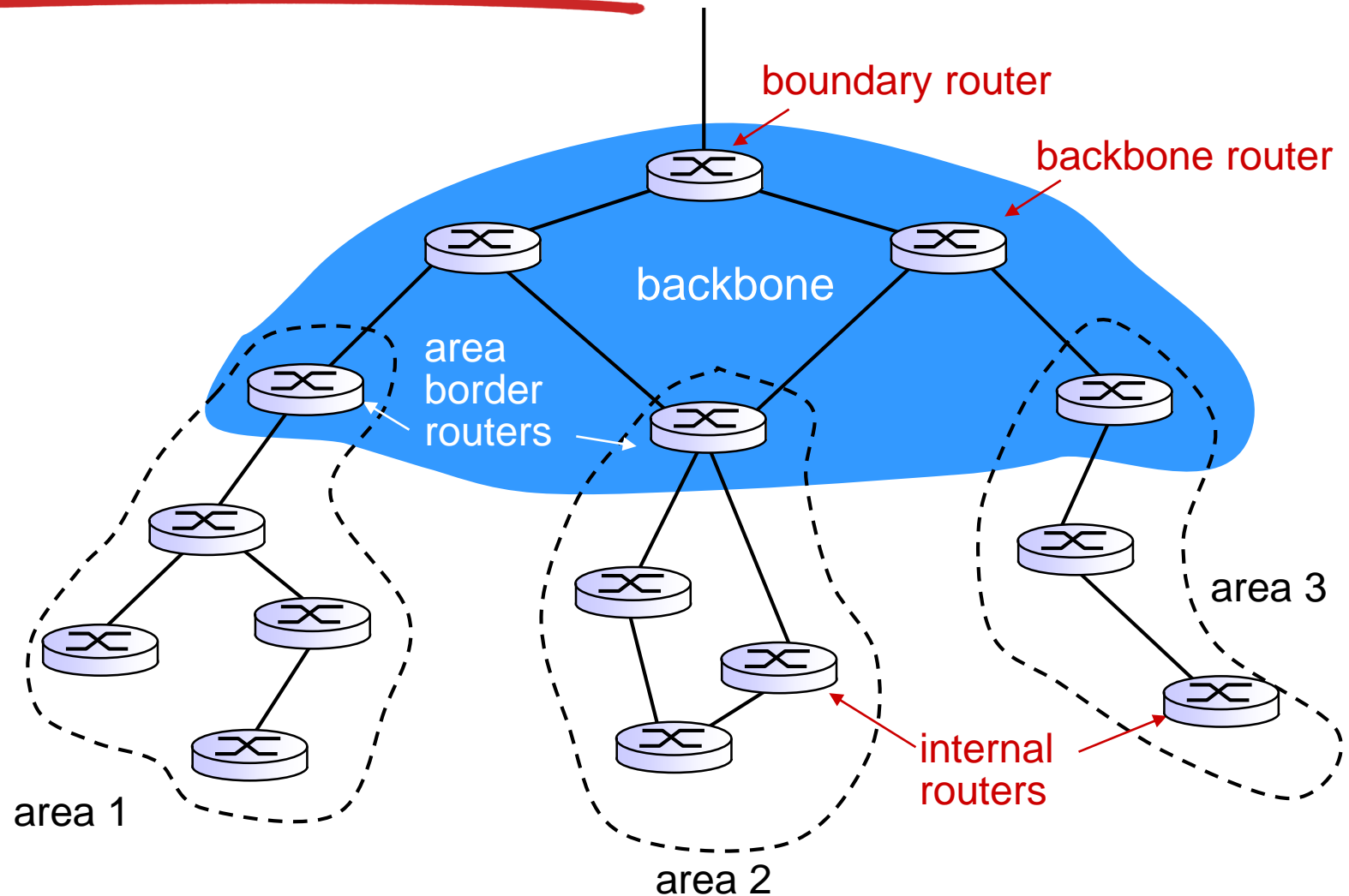
OSPF (Open Shortest Path First)

- “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
 - carried in OSPF messages **directly over IP** (rather than TCP or UDP)
 - link state: for each attached link

OSPF “advanced” features

- **security**: all OSPF messages **authenticated** (to prevent malicious intrusion)
- **multiple** same-cost **paths** allowed (only one path in RIP)
- integrated uni- and **multi-cast** support:
 - Multicast OSPF (MOSPF) uses same topology data base as OSPF
- **hierarchical** OSPF in large domains.

Hierarchical OSPF



Hierarchical OSPF

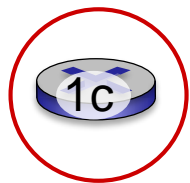
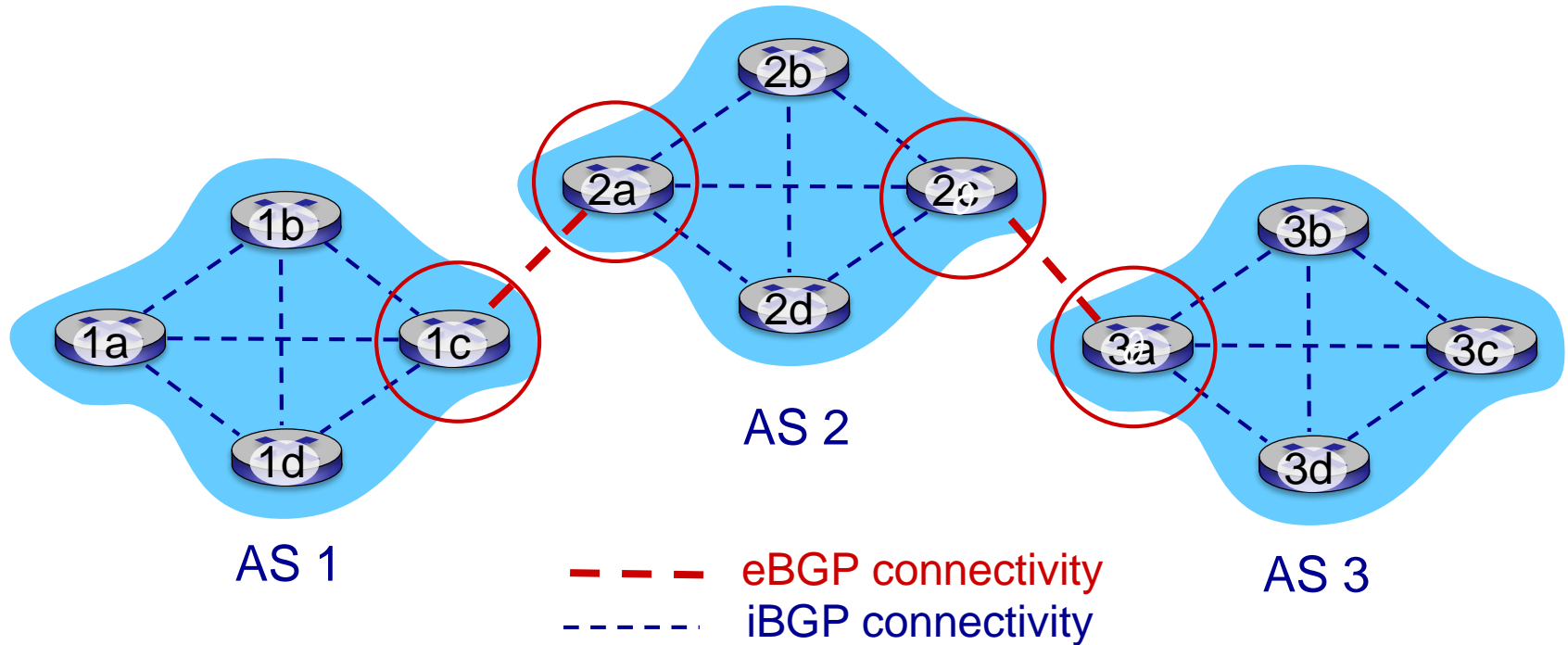
- *two-level hierarchy*: local area, backbone.
 - link-state advertisements only in **local** area
 - each nodes has detailed area topology; only know direction to other areas via backbone network.
- *area border routers*: “summarize” distances to nets in own area, advertise to other **Area Border** routers.
- *backbone routers*: run OSPF routing limited to backbone.
- *boundary routers*: connect to other AS'es.

BGP

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** ASes
 - **iBGP:** propagate reachability information to all AS-**internal** routers.
 - determine “good” routes to other networks based on reachability information and *policy*
- allows subnet to advertise its existence to rest of Internet: *“I am here”*

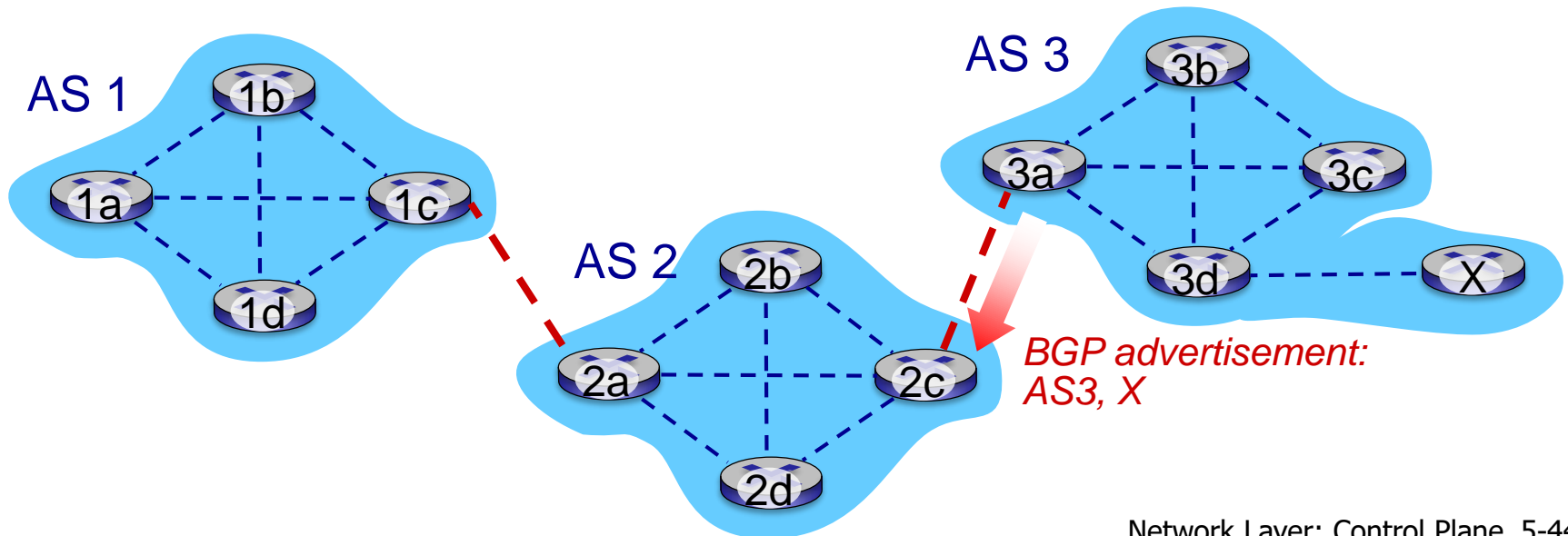
eBGP, iBGP connections



gateway routers run both eBGP and iBGP protocols

BGP basics

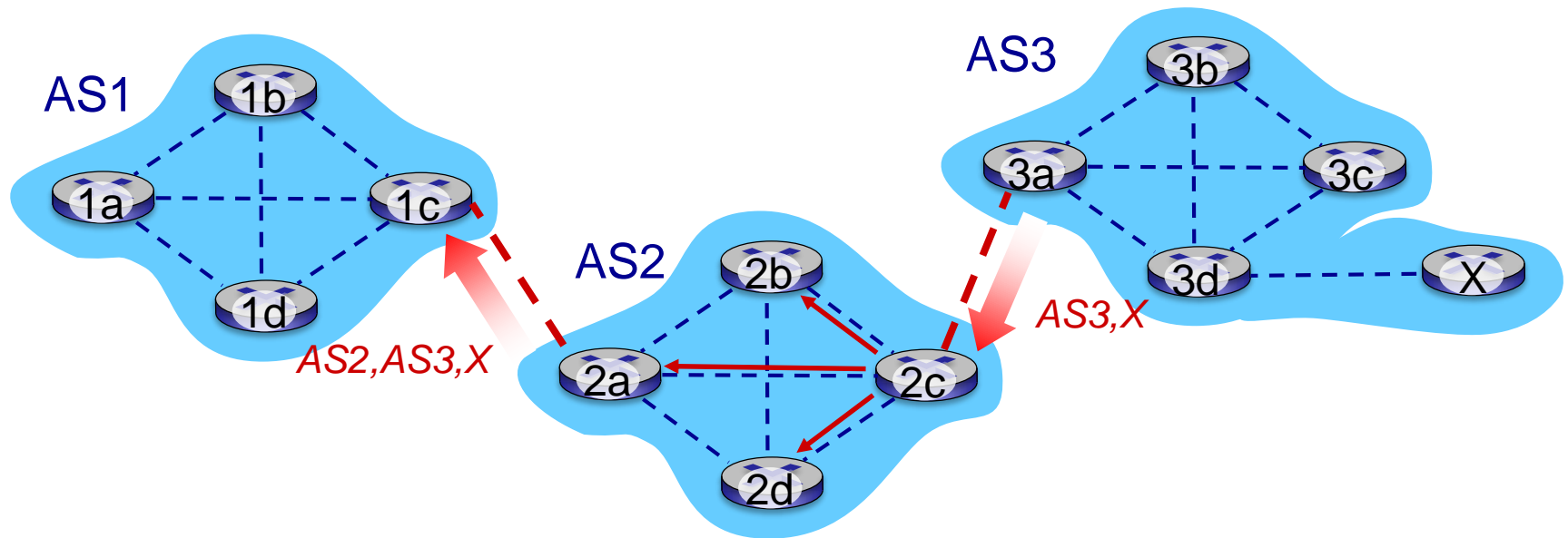
- **BGP session:** two BGP routers (“peers”) exchange BGP messages over semi-permanent TCP connection:
 - advertising *paths* to different destination network prefixes (BGP is a “**path vector**” protocol)
- when AS3 gateway router 3a advertises path **AS3,X** to AS2 gateway router 2c:
 - AS3 *promises*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



Path attributes and BGP routes

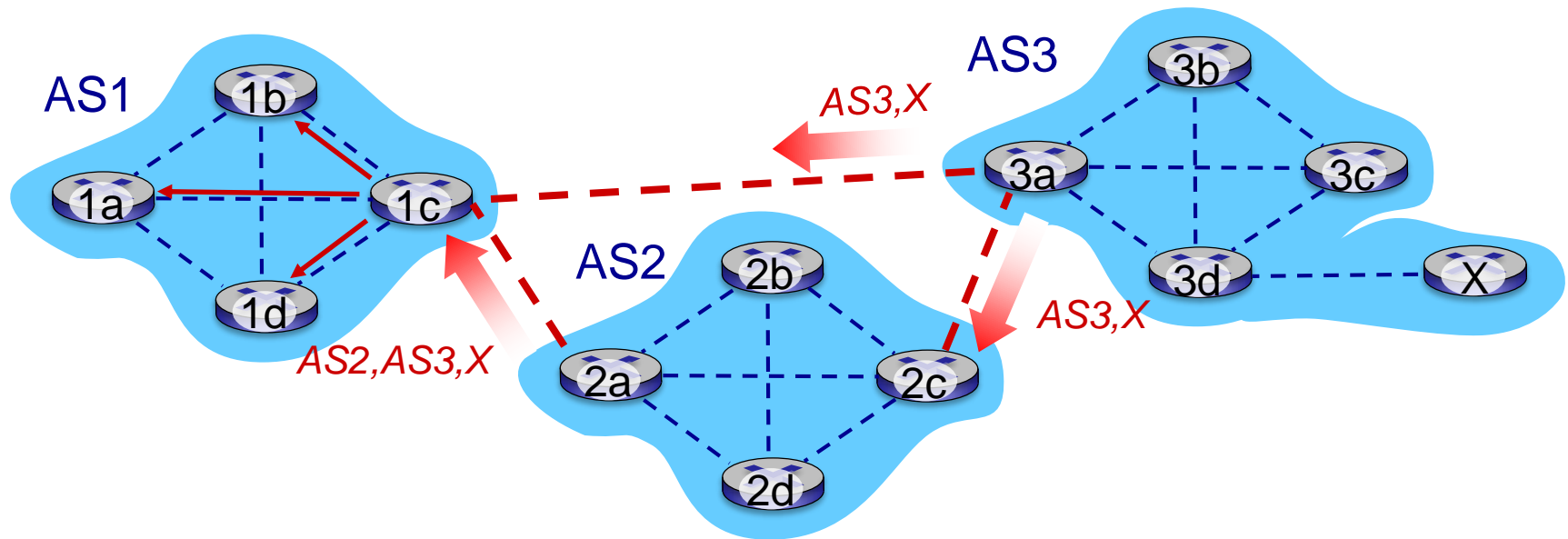
- advertised prefix includes BGP attributes
 - prefix + attributes = “route”
- two important attributes:
 - **AS-PATH**: list of **ASes** through which prefix advertisement has passed
 - **NEXT-HOP**: indicates **specific internal-AS** router to next-hop AS
- *Policy-based routing*:
 - gateway receiving 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 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



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

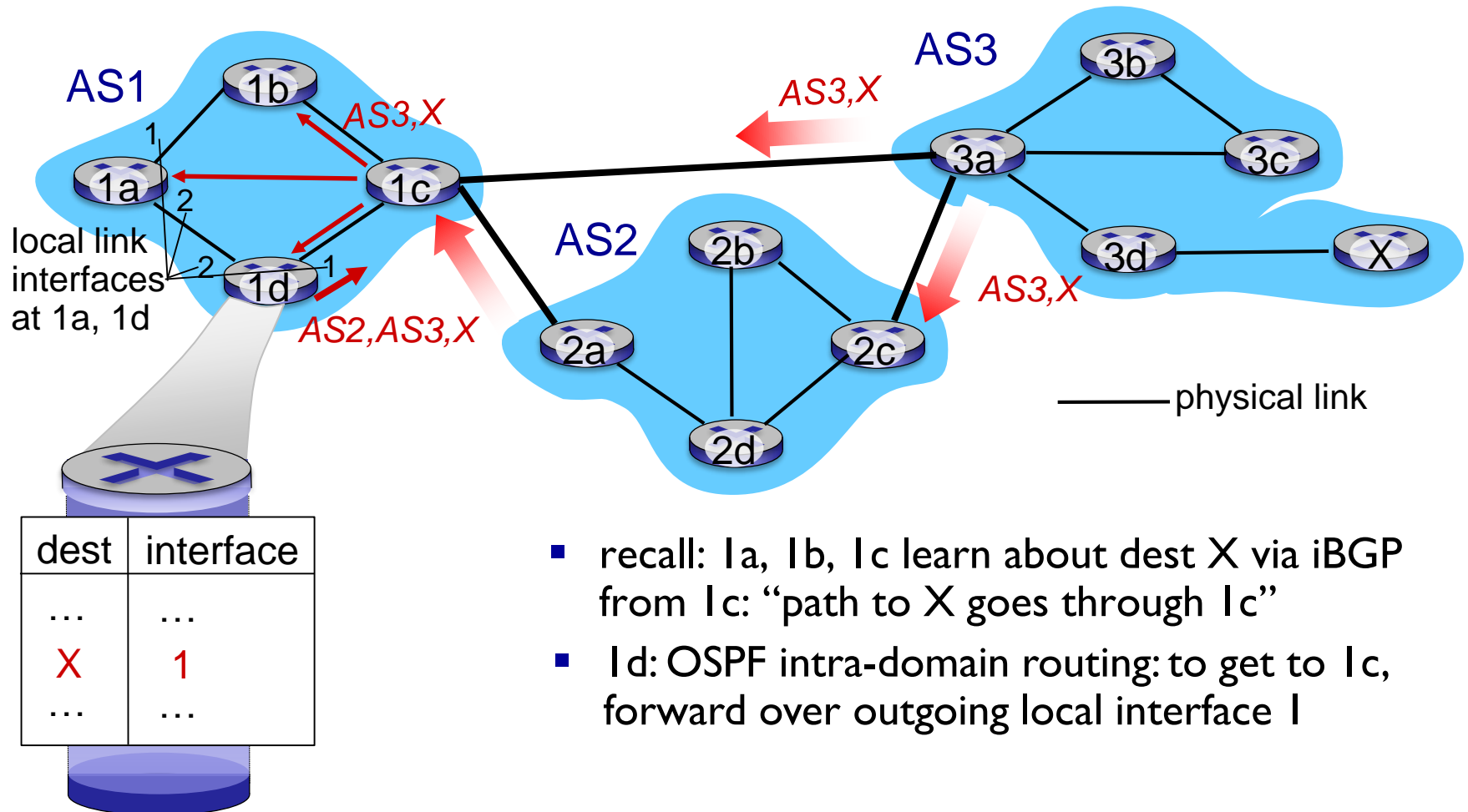
- **AS1** gateway router **1c** learns path **AS2,AS3,X** from **2a**
- **AS1** gateway router **1c** learns path **AS3,X** from **3a**
- Based on policy, **AS1** gateway router **1c** chooses path **AS3,X**, and advertises path within **AS1** via iBGP

BGP messages

- BGP messages exchanged between peers over TCP connection
- BGP messages:
 - **OPEN:** opens TCP connection to remote BGP peer and authenticates sending BGP peer
 - **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

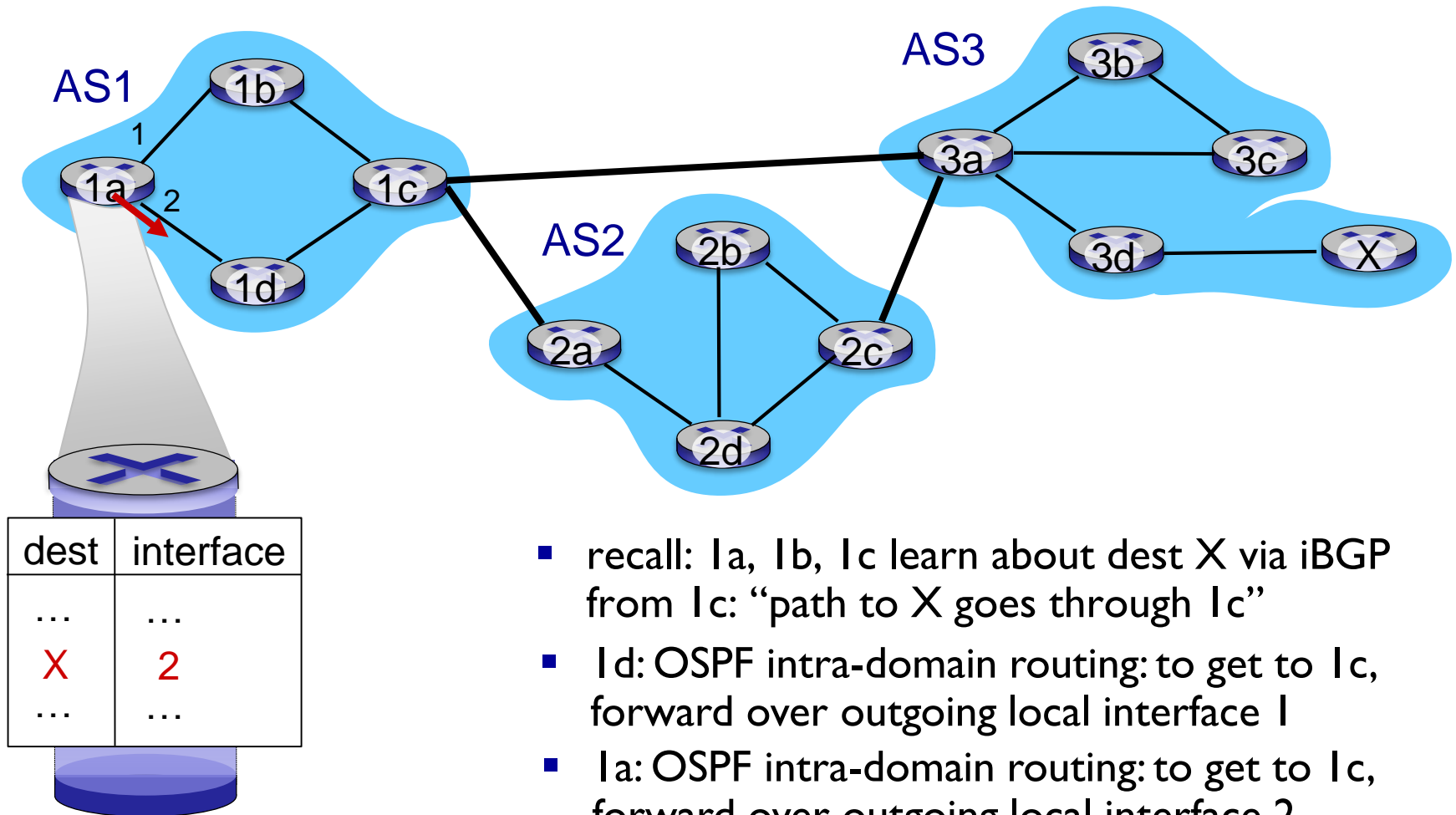
BGP, OSPF, forwarding table entries

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



BGP, OSPF, forwarding table entries

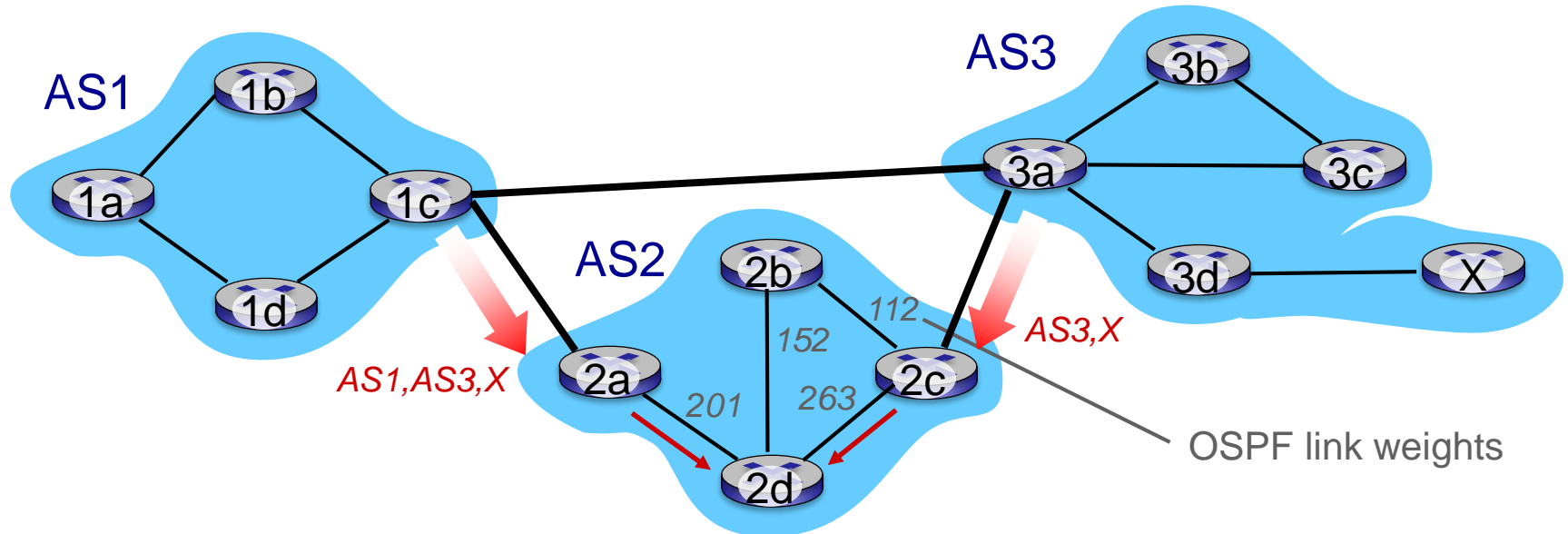
Q: how does router set forwarding table entry to distant prefix?



BGP route selection

- router may learn about more than one 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

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): don't worry about inter-domain cost!

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

ICMP

ICMP: internet control message protocol

- used by hosts & routers to communicate network-level information

- error reporting:
unreachable host, network, port, protocol
- echo request/reply (used by ping)

- network-layer “above” IP:

- ICMP msgs carried in IP datagrams

- **ICMP message:** type, code plus first 8 bytes of IP datagram causing error

<u>Type</u>	<u>Code</u>	<u>description</u>
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

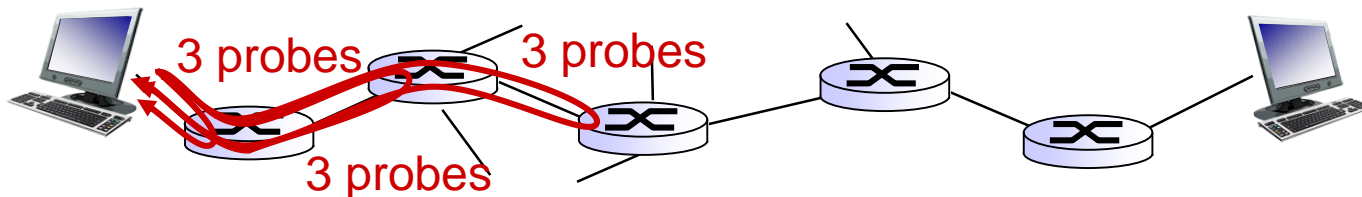
Traceroute and ICMP

- source sends series of UDP segments to destination
 - first set has TTL = 1
 - second set has TTL=2, etc.
 - unlikely port number
- when datagram in n th set arrives to n th router:
 - router discards datagram and sends source ICMP message (type 11, code 0)
 - ICMP message include name of router & IP address

- when ICMP message arrives, source records RTTs

stopping criteria:

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



SNMP

What is network management?

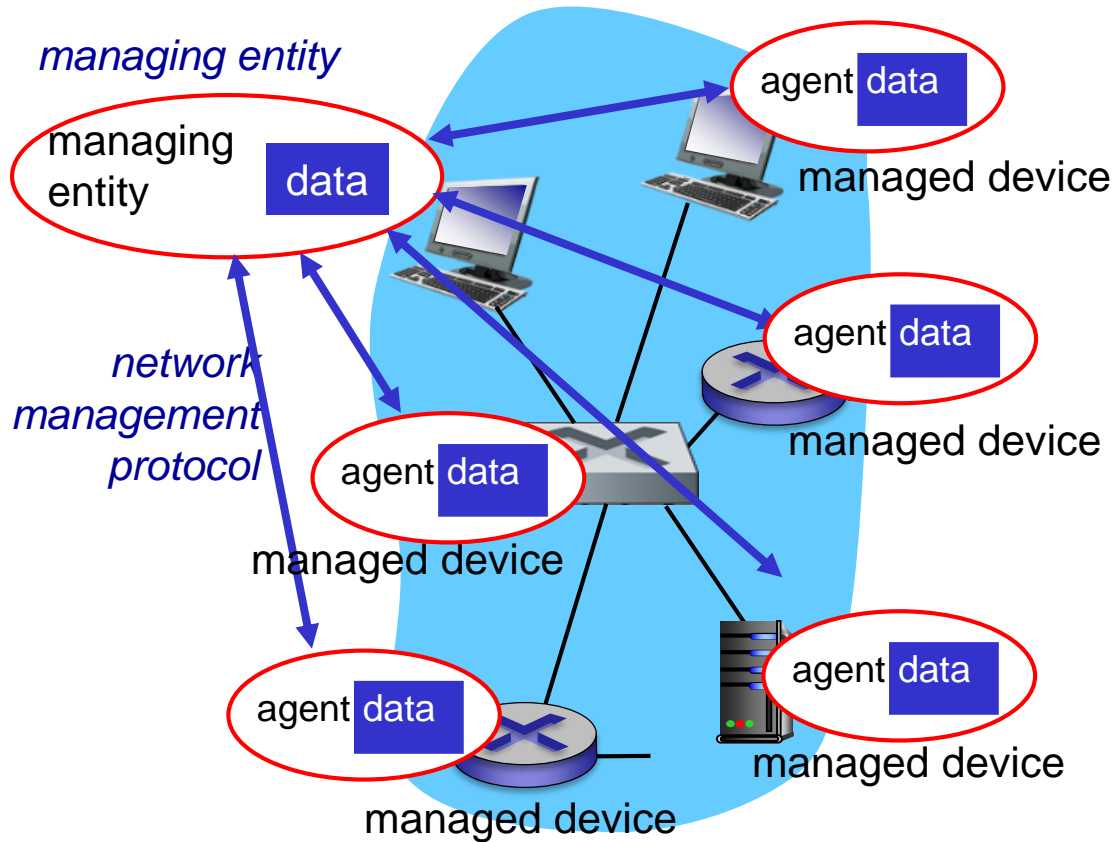
- **autonomous systems (aka “network”)**: 1000s of interacting hardware/software components
- other complex systems requiring monitoring, control:
 - jet airplane
 - nuclear power plant
 - others?



"**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."

Infrastructure for network management

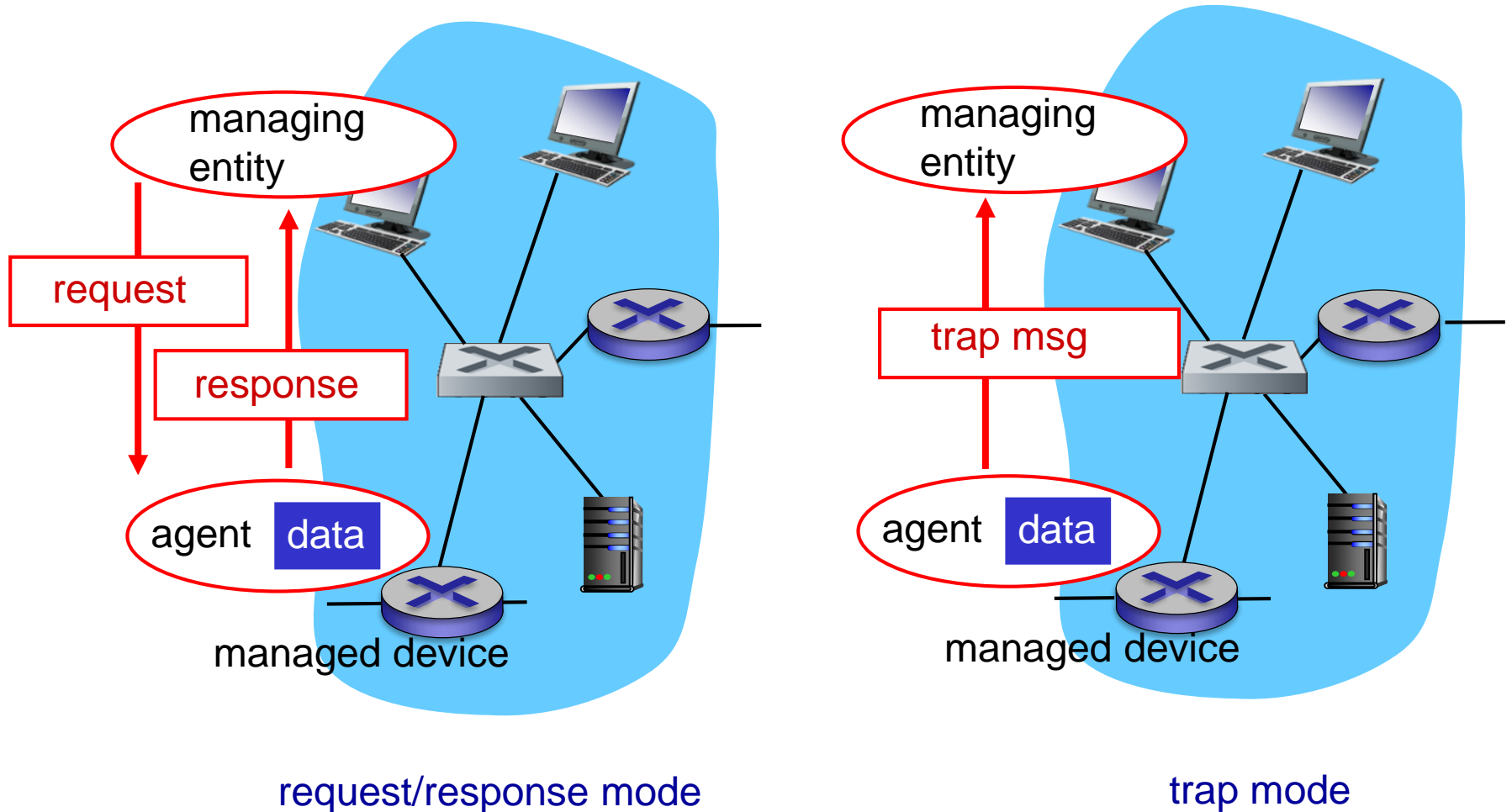
definitions:



managed devices
contain *managed objects* whose data is gathered into a **Management Information Base (MIB)**

SNMP protocol

Two ways to convey MIB info, commands:



SNMP protocol: message types

Message type

Function

GetRequest
GetNextRequest
GetBulkRequest

manager-to-agent: “get me data”
(data instance, next data in list, block of data)

InformRequest

manager-to-manager: here's MIB value

SetRequest

manager-to-agent: set MIB value

Response

Agent-to-manager: value, response to Request

Trap

Agent-to-manager: inform manager of exceptional event