

ELEC / COMP 177 – Fall 2010

Computer Networking

→ Routing Protocols (2)

Some slides from Kurose and Ross, *Computer Networking*, 5th Edition

Schedule

- **Project #2** – Due Thursday, Nov 10th
- **Homework #5** – Due Thursday, Nov 17th
- *Later this semester:*
 - *Homework #6 - Presentation on security/privacy*
 - *Topic selection* – Due Tuesday, Nov 22nd
 - *Slides* – Due Monday, Nov 28th
 - *Present!* – Tuesday, Nov 29th (and Thursday?)
 - *Project #3* – Due Tue, Dec 6th

Recap – Forwarding versus Routing

- Forwarding
 - Move packets from router's input to appropriate router output
 - Router does a *longest prefix match* (LPM) on entries in the forwarding table to determine output port
- Routing
 - Determine path (route) taken by packets from source to destination
 - Routing algorithms

Recap – Routing Algorithm Classification

- **Global Information**

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

- **Decentralized**

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

Recap – Link State – Dijkstra's Algorithm

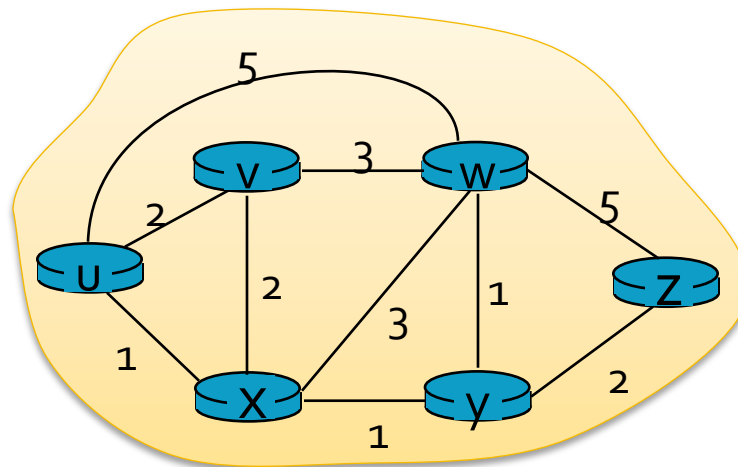
- Network topology and link costs are known to all nodes
 - Accomplished via “link state broadcast”
 - All nodes have same info
- Computes least cost paths from one node (source) to all other nodes
 - Produces **forwarding table** for that node
- Iterative: after k iterations, know least cost path to k destinations

Notation:

- u : the source (“you”)
- $c(x,y)$: link cost from node x to y ; $= \infty$ if not direct neighbors
- $D(v)$: current value of cost of path from source to dest. v
- $p(v)$: predecessor node along path from source to v
- N' : set of nodes whose least cost path definitively known

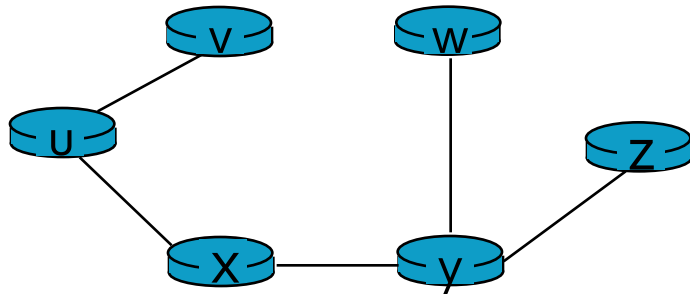
Recap – Dijkstra's Algorithm

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					



Recap – Dijkstra's Algorithm

Resulting shortest-path tree from u:



Resulting forwarding table in u:

<u>destination</u>	<u>link</u>
v	(u,v)
x	(u,x)
y	(u,x)
w	(u,x)
z	(u,x)

Recap – Distance Vector Algorithm

Iterative, asynchronous:

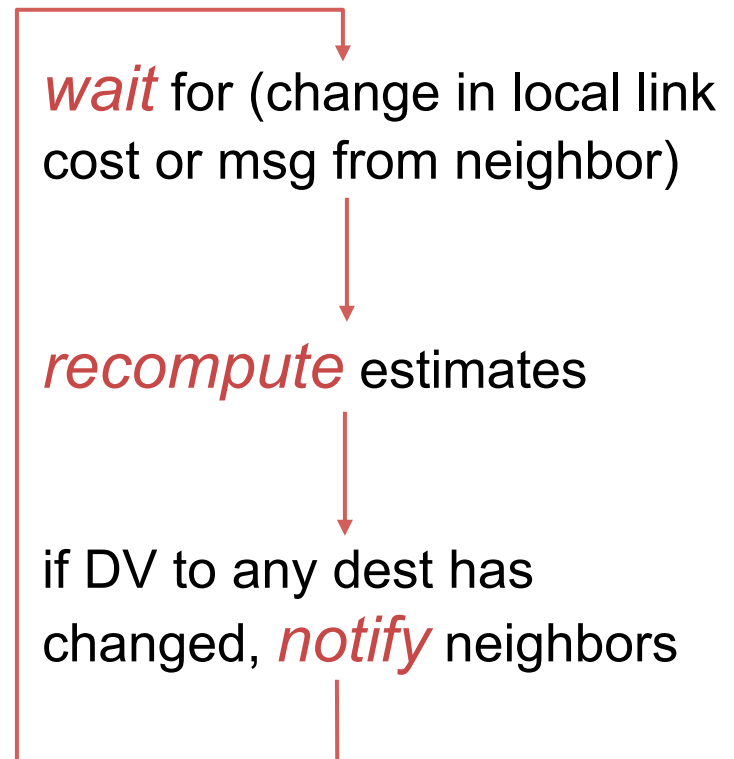
each local iteration caused by:

- 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:



Recap – Distance Vector – Bellman-Ford Equation

Define:

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

Then:

Something I know...

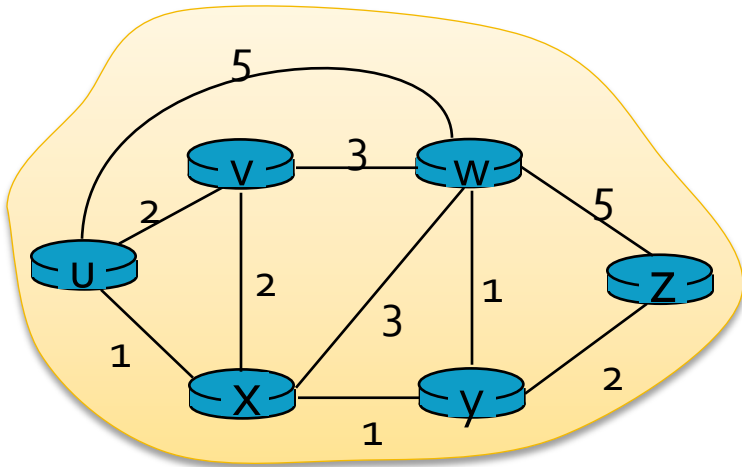
Something my neighbor told me...

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

where min is taken over all neighbors v of x

Recap – Distance Vector – Bellman-Ford

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



B-F equation says:

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

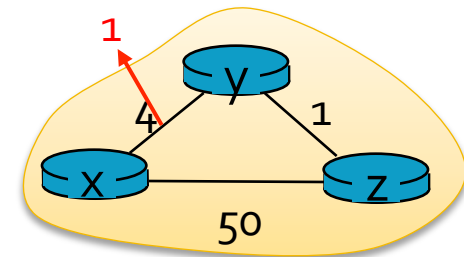
$$= \min \{ 2 + 5, \\ 1 + 3, \\ 5 + 3 \} = 4 \quad (\text{by way of } x!)$$

The node that provides the minimum cost is entered in the router forwarding table as the next hop

Distance Vector – Link Cost Changes

When a link cost changes:

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



At time t_0 , y detects the link-cost change, updates its DV, and informs its neighbors.

At time t_1 , z receives the update from y and updates its table. It computes a new least cost to x and sends its neighbors its DV.

At time t_2 , y receives z 's update and updates its distance table. y 's least costs do not change and hence y does *not* send any message to z .

“Good
news
travels
fast”

Distance Vector – Link Cost Changes

“Bad news travels slow”

Previously: $D_y(x)=4$, $D_y(z)=1$, $D_z(y)=1$, $D_z(x)=5$

At time t_0 , y detects the link-cost change. What is $D_y(x)$?

$$D_y(x) = \min\{c(y,x)+D_x(x), c(y,z)+D_z(x)\} = \{60+0, 1+5\} = 6$$

This “best route” is wrong, but y doesn’t know that!

y updates its DV, and informs its neighbors.

At time t_1 , z receives the update from y and updates its table.

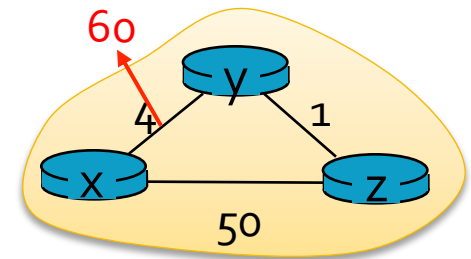
$$D_z(x) = \min\{c(z,x)+D_x(x), c(z,y)+D_y(x)\} = \{50+0, 1+6\} = 7$$

This “best route” is also wrong, but z doesn’t know that!

Now we have an infinite loop!

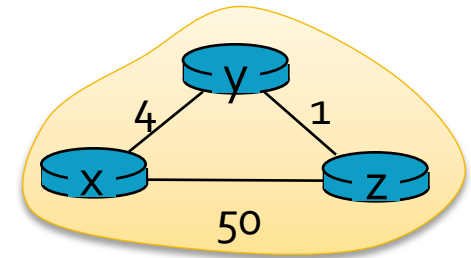
z computes a new least cost to x and sends its neighbors its DV.

Count to Infinity Problem – we escape this once z learns that its least-cost path to x is the direct link (cost 50), as the cost through y keeps incrementing (“counting”) upwards by 1 each iteration



Distance Vector – Link Cost Changes

- Poisoned-reverse “solution”
 - If z routes to x through y, then z advertises to y that its distance to x is infinity
 - Prevents y from creating a routing loop through z
- Doesn't work for more complicated networks (loops involving 3 or more routers)



Today

- Continue discussing network layer
 - **Routing algorithms used in the Internet**
 - **Routing Information Protocol (RIP)**
 - **Open Shortest Path First (OSPF)**
 - **Border Gateway Protocol (BGP)**

Recap – Hierarchical Routing

- Our routing discussion thus far has been idealized
 - All routers are identical
 - The network is “flat”
- This is not true in practice!
- **Problem 1 – Scale**
 - Hundreds of millions of destinations:
 - Can’t store all destinations in routing tables!
 - Routing table exchange would swamp links!
 - Distance-vector would never converge
- **Problem 2 - Administrative autonomy**
 - Internet = network of networks
 - Each network admin wants to control routing in his/her own network

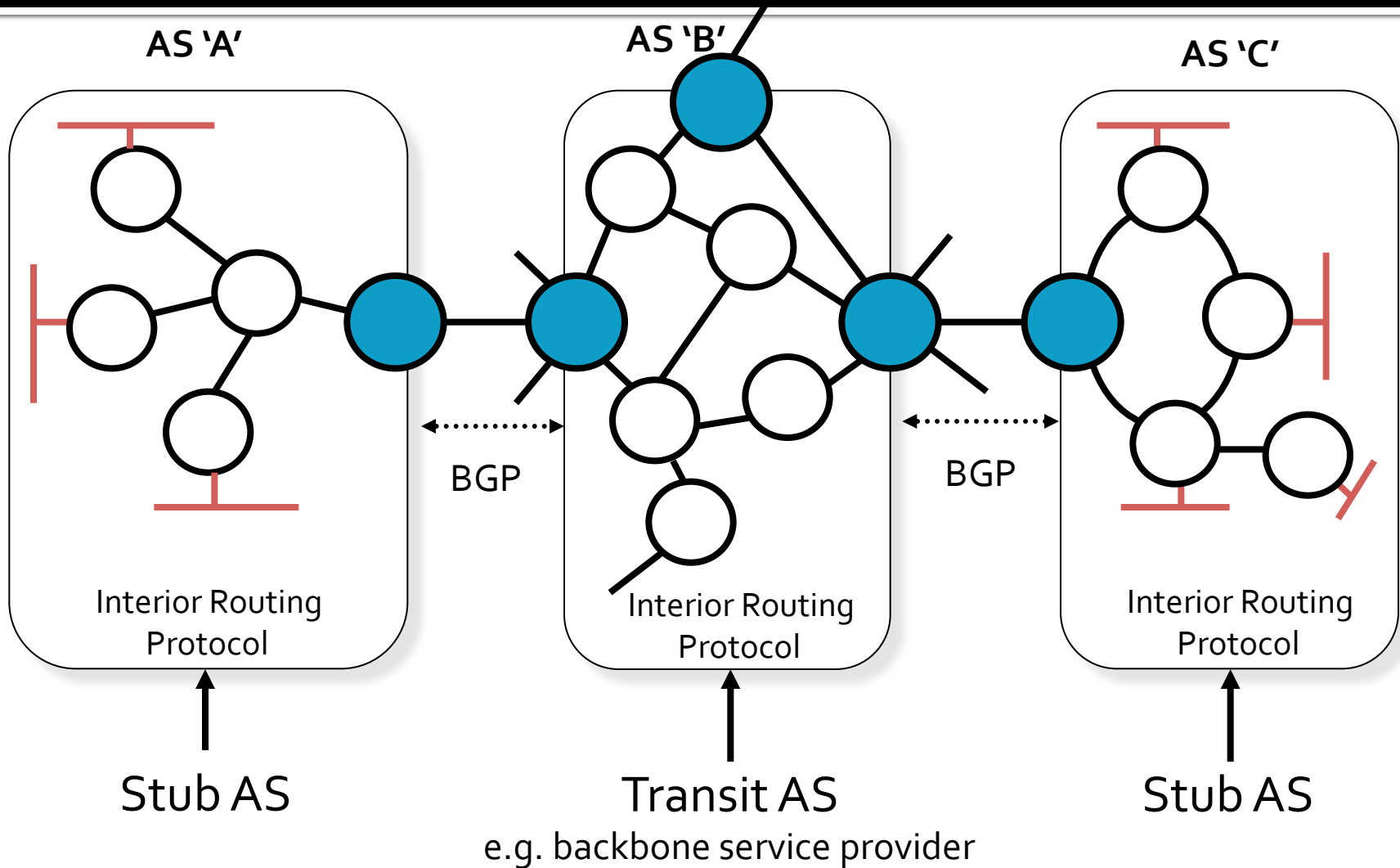
Recap – Hierarchical Routing

- Aggregate routers into regions (aka “**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
 - Direct link to router in another AS

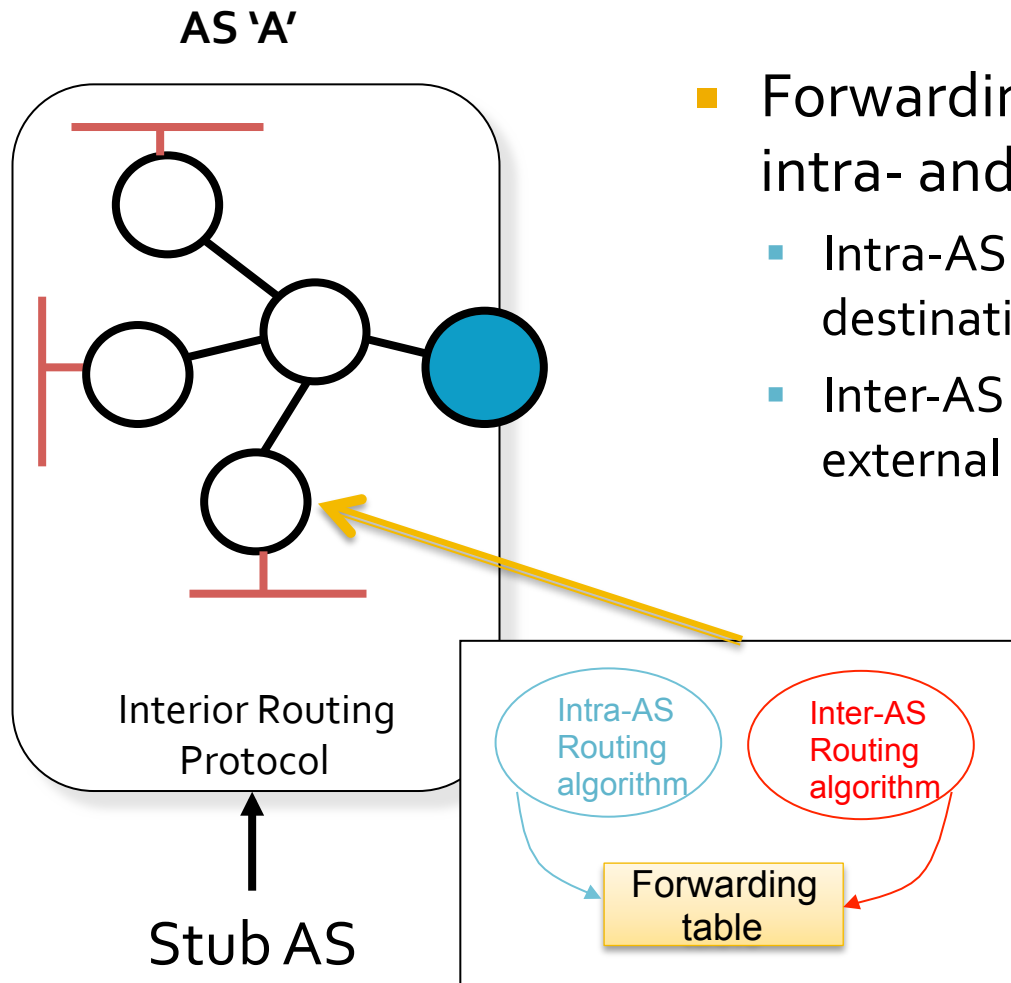
Routing in the Internet

- The Internet uses hierarchical routing
- The Internet is split into Autonomous Systems
 - “Independent” networks on the Internet
 - Typically owned/controlled by a single entity
 - Share a common routing policy
- Example autonomous systems
 - Pacific (18663), Exxon (1766), IBM (16807), Level3 (3356)
- Different routing protocols within and between autonomous systems
 - Interior gateway/routing protocol (e.g. OSPF)
 - Border gateway protocol (e.g. BGP)

Autonomous Systems



Forwarding Table



- Forwarding table configured by both intra- and inter-AS routing algorithm
 - Intra-AS sets entries for internal destinations
 - Inter-AS & intra-As sets entries for external destinations

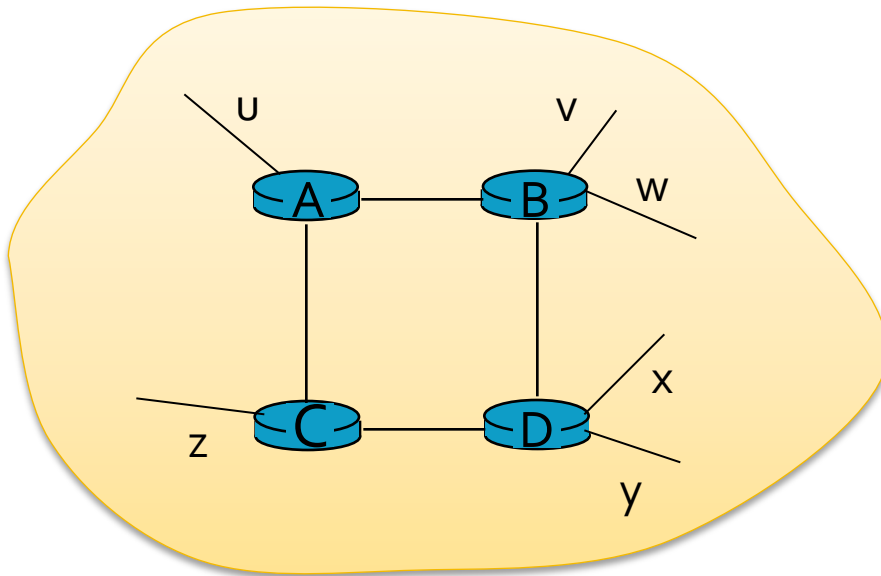
Intra-AS Routing

- Routing *inside* the autonomous system
- 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)

Routing Information Protocol (RIP)

Routing Information Protocol (RIP)

- **Distance vector** algorithm
- Included in BSD-UNIX Distribution in 1982
- Distance metric: # of hops (max = 15 hops)



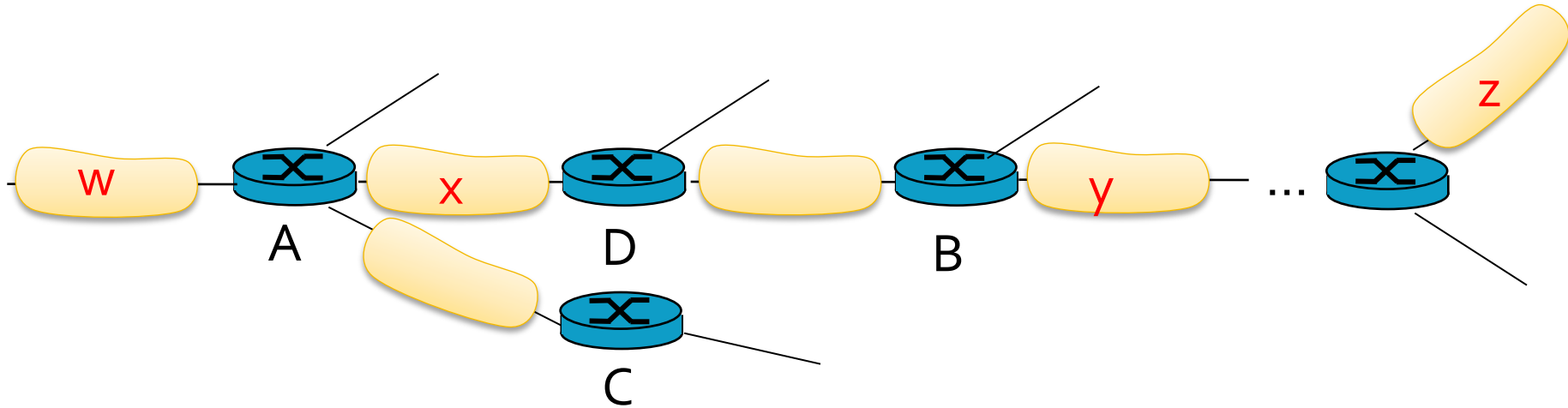
From router A to subnets:

<u>destination</u>	<u>hops</u>
U	1
V	2
W	2
X	3
Y	3
Z	2

RIP advertisements

- Distance vectors
 - Exchanged among neighbors every 30 seconds via Response Message (also called **advertisement**)
- Each advertisement lists up to 25 destination subnets within AS

RIP: Example



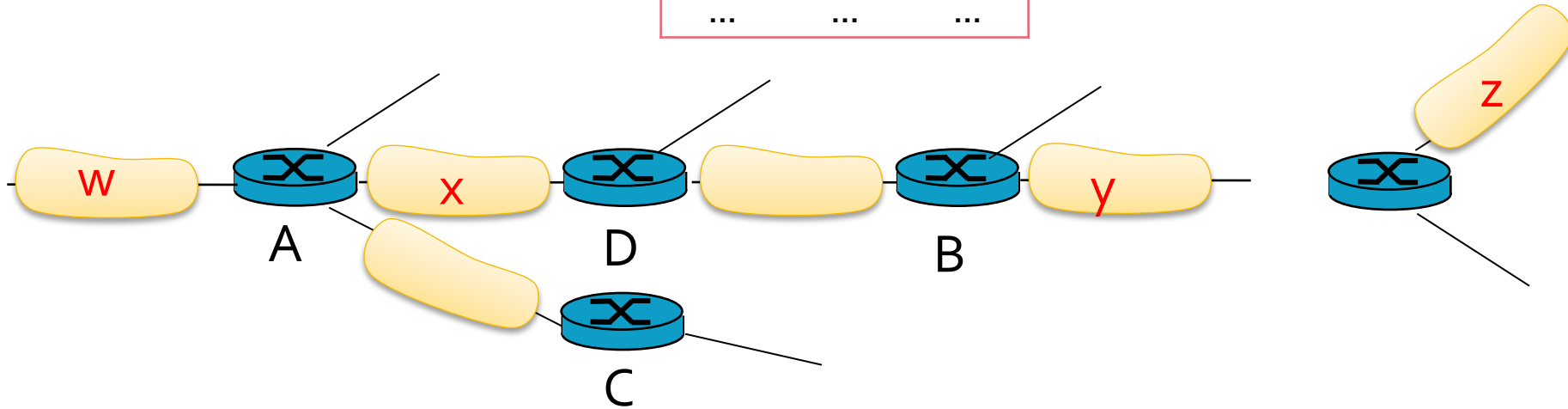
Routing/Forwarding table in D:

Destination Network	Next Router	# of Hops to Destination
w	A	2
y	B	2
z	B	7
x	--	1
...

RIP: Example

Dest	Next	Hops
w	--	1
x	--	1
z	C	4
...

Advertisement from A to D



Routing/Forwarding table in D:

Destination Network	Next Router	# of Hops to Destination
w	A	2
y	B	2
z	B A	7 5
x	--	1
...

RIP: Link Failure and Recovery

- If no advertisement heard after 180 sec, the neighbor/link declared dead
- Failure recovery
 - 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

Open Shortest Path First (OSPF)

Open Shortest Path First Routing

- Networks are partitioned into “areas”
 - OSPF only runs within a specific area
 - Other protocols (i.e., BGP) used to route outside an area
- Link-state algorithm
 - Each node has full topology map
 - Route computation using **Dijkstra’s algorithm**

Open Shortest Path First Routing

- Routers periodically send “hello” and “link state” packets to their neighbors
 - Learn who your neighbors are dynamically
 - Decide link/router down if no more hellos
 - Announce changes to the topology
 - Broadcast throughout the area
 - Carried in OSPF messages directly over IP (rather than TCP or UDP)

Link State Advertisements

- Router link advertisements
 - Sent by all routers
 - State and cost of the router's links to the area
- Network link advertisements
 - Sent only by designated routers
 - Describes all routers attached to the network
- Summary link advertisements
 - Sent only by area border routers
 - Describes routes to other areas
- AS external link advertisements
 - Sent only by AS boundary routers
 - Describes routes to other autonomous systems

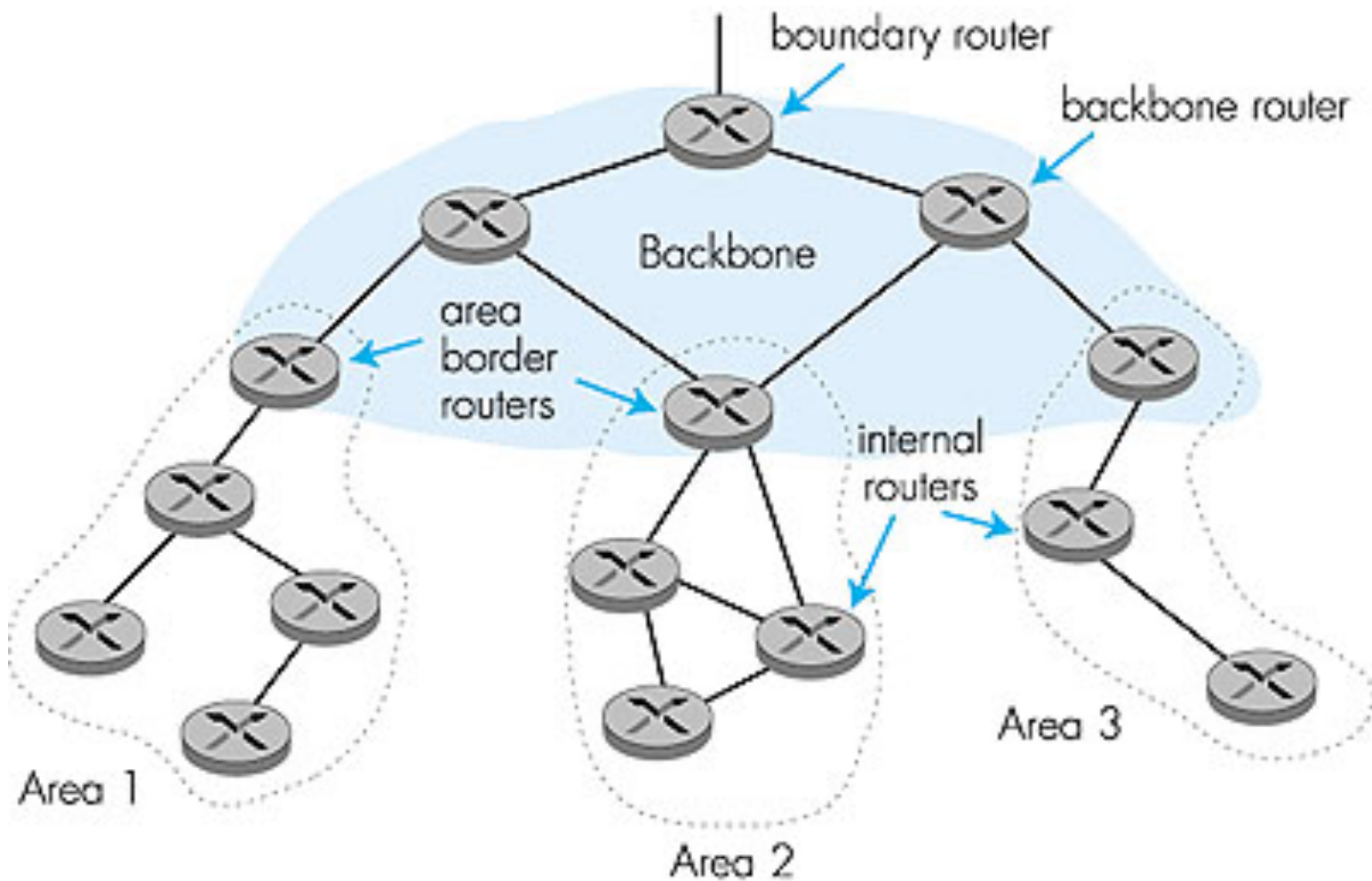
Reliable Flooding of LSPs

- Link state packets (LSP) delivered throughout the area
 - Flooded throughout the area
 - Sequence numbers and TTLs
- Reliable Flooding
 - If newer sequence number, then forward packet over all links other than the ingress link, otherwise drop packet
 - Resend unacknowledged packets
- Link State Detection
 - If no hello packets during dead interval, assume link is down

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; high for real time)
- **Hierarchical** OSPF in large domains

Hierarchical OSPF



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

Routing Across Borders

- OSPF doesn't scale
 - Broadcasts all link states to all routers
 - Calculates shortest path to all routers
- Autonomous systems are independent
 - Run by different organizations
 - May use different link cost metrics

Routing Across Borders

- Need a “border gateway protocol”
 - Global routing protocol across autonomous systems
- Global connectivity is at stake!
 - Must settle on one protocol
- What are the requirements?
 - Scalability
 - Flexibility in choosing routes

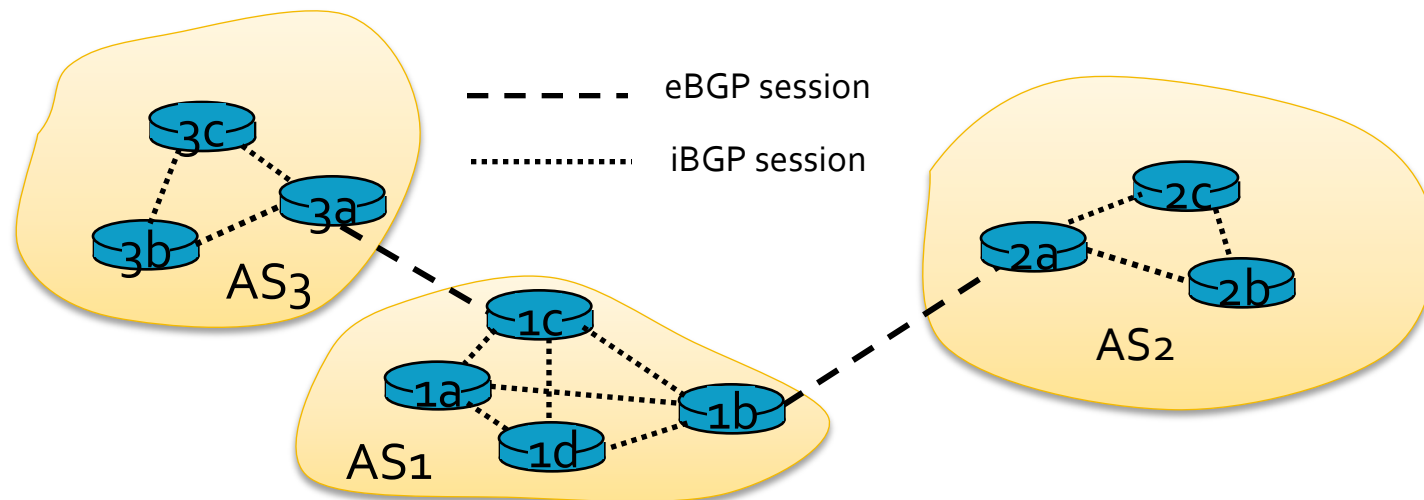
Border Gateway Protocol (BGP)

Internet Inter-AS routing: BGP

- BGP is the **de facto standard**
- BGP provides each AS a means to:
 - Obtain subnet reachability information from neighboring ASs
 - Propagate reachability information to all routers inside an AS
 - Determine “good” routes to subnets based on reachability information and policy
- Allows subnet to advertise its existence to rest of Internet: “**I am here**”

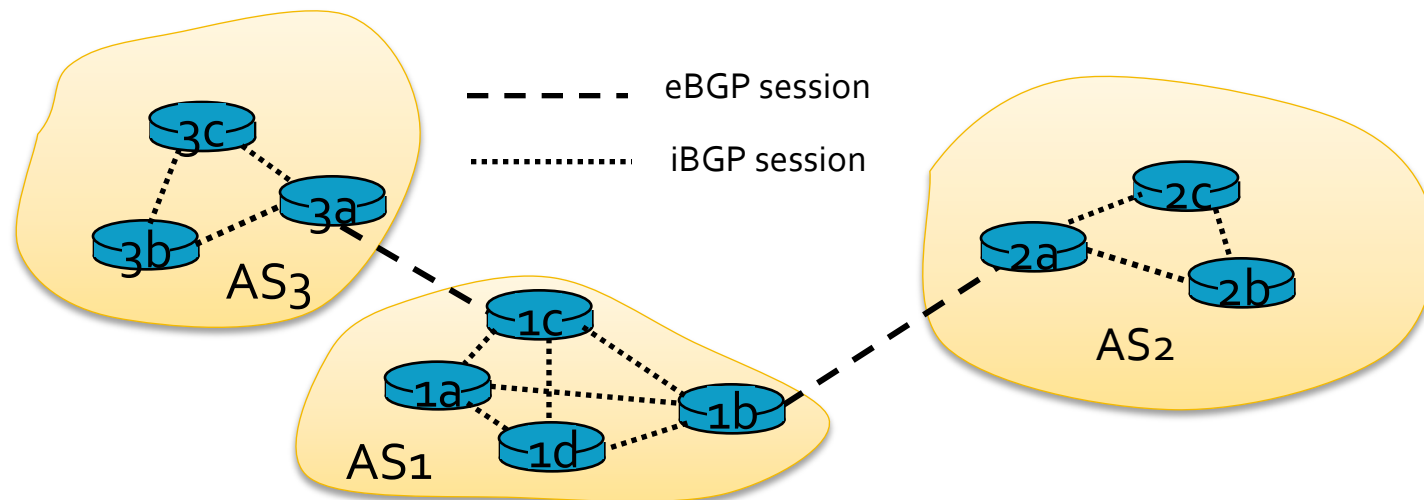
BGP Basics

- Pairs of routers (BGP peers) exchange routing info over semi-permanent TCP connections: **BGP sessions**
 - BGP sessions need not correspond to physical links.
- When AS2 advertises a prefix to AS1:
 - AS2 *promises* it will forward datagrams towards that prefix.
 - AS2 can aggregate prefixes in its advertisement



Distributing Reachability Info

- 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.



Border Gateway Protocol (BGP-4)

- BGP uses “path vectors” (AS_PATH)
 - Advertises complete “paths” – a list of autonomous systems
 - “The network 171.64/16 can be reached via the path {AS₁, AS₅, AS₁₃}”
 - Makes no use of distance vectors or link states
- Path selection
 - Supports classless inter-domain routing
 - Paths with loops are detected locally and ignored
 - Local policies pick the preferred path among options
 - When a link/router fails, the path is “withdrawn”

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)
- When gateway router receives route advertisement, uses **import policy** to accept/decline

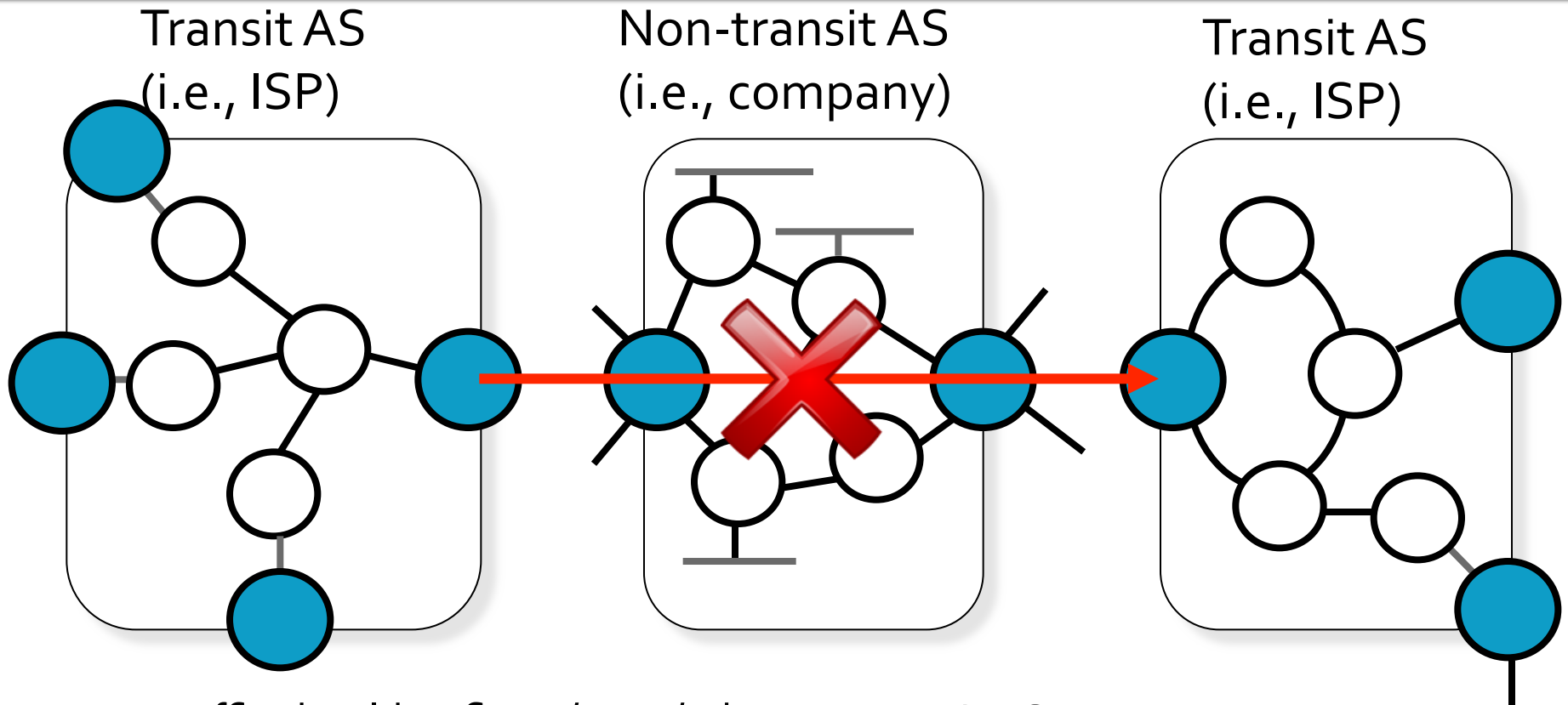
BGP route selection

- Router may learn about more than 1 route to some prefix
 - Must select best route
- Elimination rules:
 1. Local preference value attribute: policy decision
 2. Shortest AS-PATH
 3. Closest NEXT-HOP router: hot potato routing
 4. Additional criteria

BGP messages

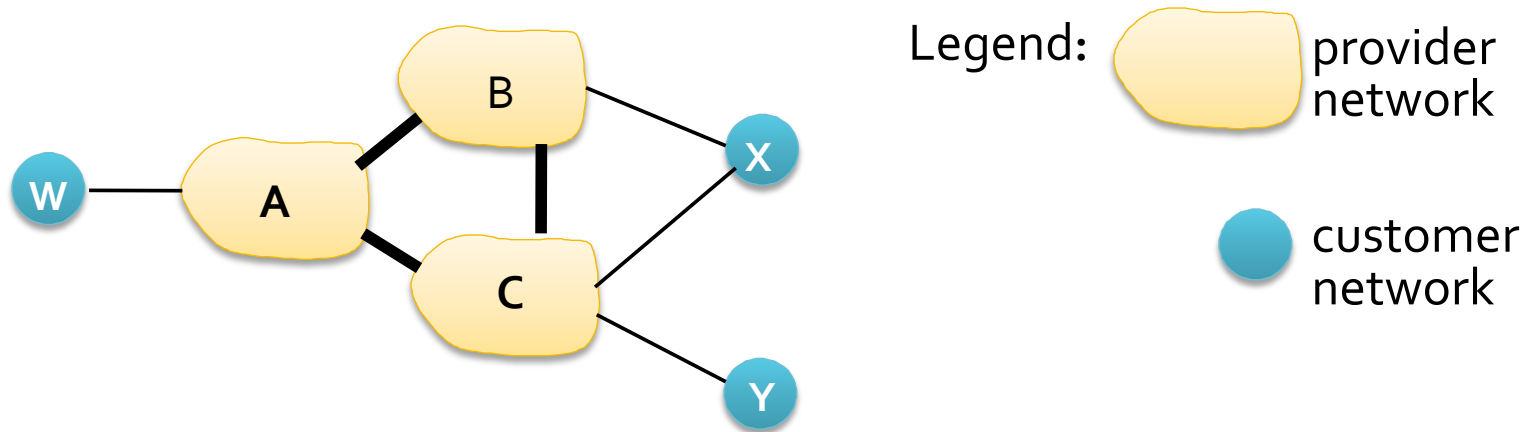
- BGP messages exchanged using TCP.
- BGP messages:
 - **OPEN**: opens TCP connection to peer and authenticates sender
 - **UPDATE**: advertises new path (or withdraws old)
 - **KEEPALIVE** keeps connection alive in absence of UPDATES; also ACKs OPEN request
 - **NOTIFICATION**: reports errors in previous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

BGP Routing Policy (1)



- Traffic shouldn't flow *through* the non-transit AS
 - Paying ISPs for connectivity, not to route traffic for them!
 - Don't advertise any BGP routes between transit AS's
 - Pacific is dual-homed to TCTC (Time Warner) and SWIS (AT&T)

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 \$\$\$ for routing CBAW since neither W nor C are customers of B
 - B wants to force C to route to w via A
 - B wants to route only to/from its customers!

Why Different Intra- and Inter-AS routing ?

■ Policy

- Inter-AS: admin wants control over how its traffic is routed and who routes through its net
- Intra-AS: single admin, so no policy decisions needed

■ Scale

- Hierarchical routing saves table size and reduces update traffic

■ Performance

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

Traceroute with AS numbers

```
dhcp-10-10-207-20:~ shafer$ traceroute -a www.msu.ru
traceroute to www.msu.ru (193.232.113.151), 64 hops max, 52 byte packets
 1  [AS0] 138.9.253.252 (138.9.253.252)  0.740 ms  0.741 ms  1.290 ms
 2  [AS0] 74.202.6.5 (74.202.6.5)  5.245 ms  15.006 ms  5.142 ms
 3  [AS4323] sjc1-pr1-xe-0-0-0-0.us.twtelecom.net (66.192.251.170)  6.414 ms  6.640 ms  17.283 ms
 4  [AS6453] if-10-0-0-56.core3.sqn-sanjose.as6453.net (209.58.116.50)  6.628 ms *
    [AS6453] if-13-0-0-55.core3.sqn-sanjose.as6453.net (66.198.97.9)  7.056 ms
 5  [AS6453] if-9-0-0.mcore4.pdi-paloalto.as6453.net (216.6.33.6)  68.184 ms
    [AS6453] if-6-0-0-1145.mcore4.pdi-paloalto.as6453.net (216.6.86.45)  8.120 ms
    [AS6453] if-9-0-0.mcore4.pdi-paloalto.as6453.net (216.6.33.6)  491.007 ms
 6  [AS11029] if-0-0-0-892.mcore3.njy-newark.as6453.net (209.58.124.25)  78.807 ms  109.426 ms
78.890 ms
 7  [AS15706] if-4-0-0.core1.fv0-frankfurt.as6453.net (195.219.69.29)  167.206 ms  167.461 ms
167.002 ms
 8  [AS15706] if-0-0-0.core1.frl-frankfurt.as6453.net (195.219.69.54)  171.256 ms  171.844 ms
174.118 ms
 9  [AS6453] if-7-1-0-1310.core1.stk-stockholm.as6453.net (195.219.131.45)  1180.587 ms  437.592
ms  586.125 ms
10 [AS6453] ix-4-0-1.core1.stk-stockholm.as6453.net (195.219.131.22)  200.475 ms  200.301 ms
201.106 ms
11 [AS3267] b57-1-gw.spb.runnet.ru (194.85.40.129)  216.199 ms  216.117 ms  214.311 ms
12 [AS3267] bl16-1-gw.spb.runnet.ru (194.85.40.78)  214.723 ms  214.463 ms  214.494 ms
13 [AS3267] bm18-1-gw.spb.runnet.ru (194.85.40.169)  214.608 ms  214.504 ms  214.493 ms
14 [AS3267] tv11-1-gw.msk.runnet.ru (194.85.40.137)  214.260 ms  214.360 ms  214.478 ms
15 [AS3267] m9-2-gw.msk.runnet.ru (194.85.40.53)  214.752 ms  214.496 ms  214.882 ms
16 [AS3267] msu.msk.runnet.ru (194.190.255.234)  214.197 ms  214.907 ms  214.656 ms
17 [AS2848] 193.232.127.12 (193.232.127.12)  214.501 ms  214.166 ms  214.531 ms
18 [AS2848] 193.232.113.151 (193.232.113.151)  214.864 ms !Z  214.666 ms !Z  214.522 ms !Z
```


AS Numbers in Traceroute

AS	Name
0	Reserved (local use) – Pacific is here...
4323	Time Warner Telecom
6453	Tata Communications A Tier-1 ISP headquartered in India This is their Canada-based AS number
11029	Tata Communications (again!) Strange registry entry – corporate buyout?
15706	Tata Communications (yet again!) Strange registry entry – corporate buyout?
3267	Runnet - State Institute of Information Technologies & Telecommunications (SIIT&T "Informika")
2848	Moscow State University

Problems

- BGP designed for policy, not performance
- Susceptible to misconfiguration
 - Intentionally / accidentally announce routes to networks you cannot reach
- Incompatible policies might render networks unreachable

BGP, Censorship, and You

1. February 2008 - Pakistan government orders Pakistan Telecom (AS 17557) to block access to YouTube
2. Pakistan Telecom advertises a route for 208.65.153/24 (YouTube) to its customers leading to a black hole
3. That route is accidentally advertised to its provider (PCCW)
 - This is more specific than YouTube's (AS 36561) real advertisements (208.65.152/22)
 - Multiple routes → More specific route preferred
4. PCCW failed to verify that Pakistan Telecom actually owned YouTube's netblock (very common)
 - BGP uses transitive trust – PCCW trusted P.T., and upstream providers trusted PCCW
5. Within about 3 minutes, a large fraction of the Internet had the bad route
 - YouTube traffic was routed to AS 17557 instead of AS 36561
 - AS 17557 can then just drop the received traffic

We Want Our Videos Back!

6. ~1 hour later, YouTube advertises that its addresses have been hijacked to its providers
 - YouTube verifiably owns that address space and its AS number
7. Autonomous systems stop using the bad route
 - YouTube also advertises its own /25 routes
8. ~1 hour later, Pakistan Telecom's provider (Hong Kong-based PCCW) withdraws bogus routes to AS 17557