

# National University of Computer & Emerging Sciences

## CS 3001 - COMPUTER NETWORKS

### Lecture 21 Chapter 4

3<sup>rd</sup> November, 2022

Nauman Moazzam Hayat  
[nauman.moazzam@lhr.nu.edu.pk](mailto:nauman.moazzam@lhr.nu.edu.pk)

Office Hours: 02:30 pm till 06:00 pm (Every Tuesday & Thursday)

# Chapter 4

## Network Layer

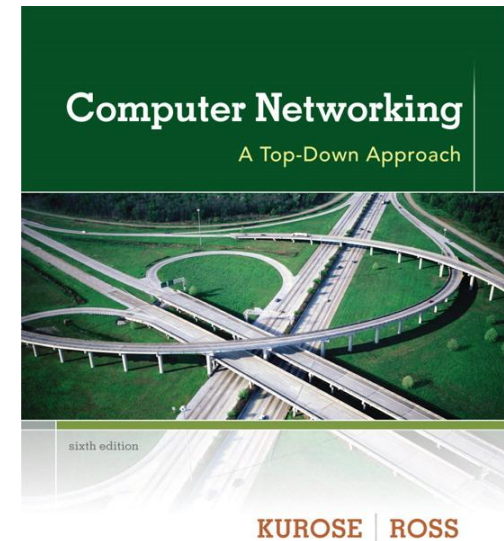
### A note on the use of these ppt slides:

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 1996-2013  
J.F Kurose and K.W. Ross, All Rights Reserved



**Computer  
Networking: A Top  
Down Approach**  
6<sup>th</sup> edition  
Jim Kurose, Keith Ross  
Addison-Wesley  
March 2012

# Chapter 4: outline

## 4.1 introduction

## 4.2 virtual circuit and datagram networks

## 4.3 what's inside a router

## 4.4 IP: Internet Protocol

- datagram format
- IPv4 addressing
- ICMP
- IPv6

## 4.5 routing algorithms

- link state
- distance vector
- hierarchical routing

## 4.6 routing in the Internet

- RIP
- OSPF
- BGP

## 4.7 broadcast and multicast routing

# Hierarchical routing

our routing study thus far - idealization

- ❖ all routers identical
- ❖ network “flat”

... *not* true in practice

*scale:* with 600 million destinations:

- ❖ can't store all dest's 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

# Hierarchical routing

- ❖ aggregate routers into regions, “**autonomous systems**” (AS)
- ❖ routers in same AS run same routing protocol
  - “**intra-AS**” routing protocol (or Interior Gateway Protocols i.e. IGP)
  - routers in different AS can run different intra-AS routing protocol

## *gateway router:*

- ❖ at “edge” of its own AS
- ❖ has link to router in another AS

## *(First Hop / Default Router:*

- ❖ Typically the host is connected to one router, called the first hop / default router for that host)

# Why different Intra-, Inter-AS routing ?

## *policy:*

- ❖ inter-AS: admin wants control over how its traffic routes, who routes through its net.
- ❖ intra-AS: single admin, so no policy decisions needed

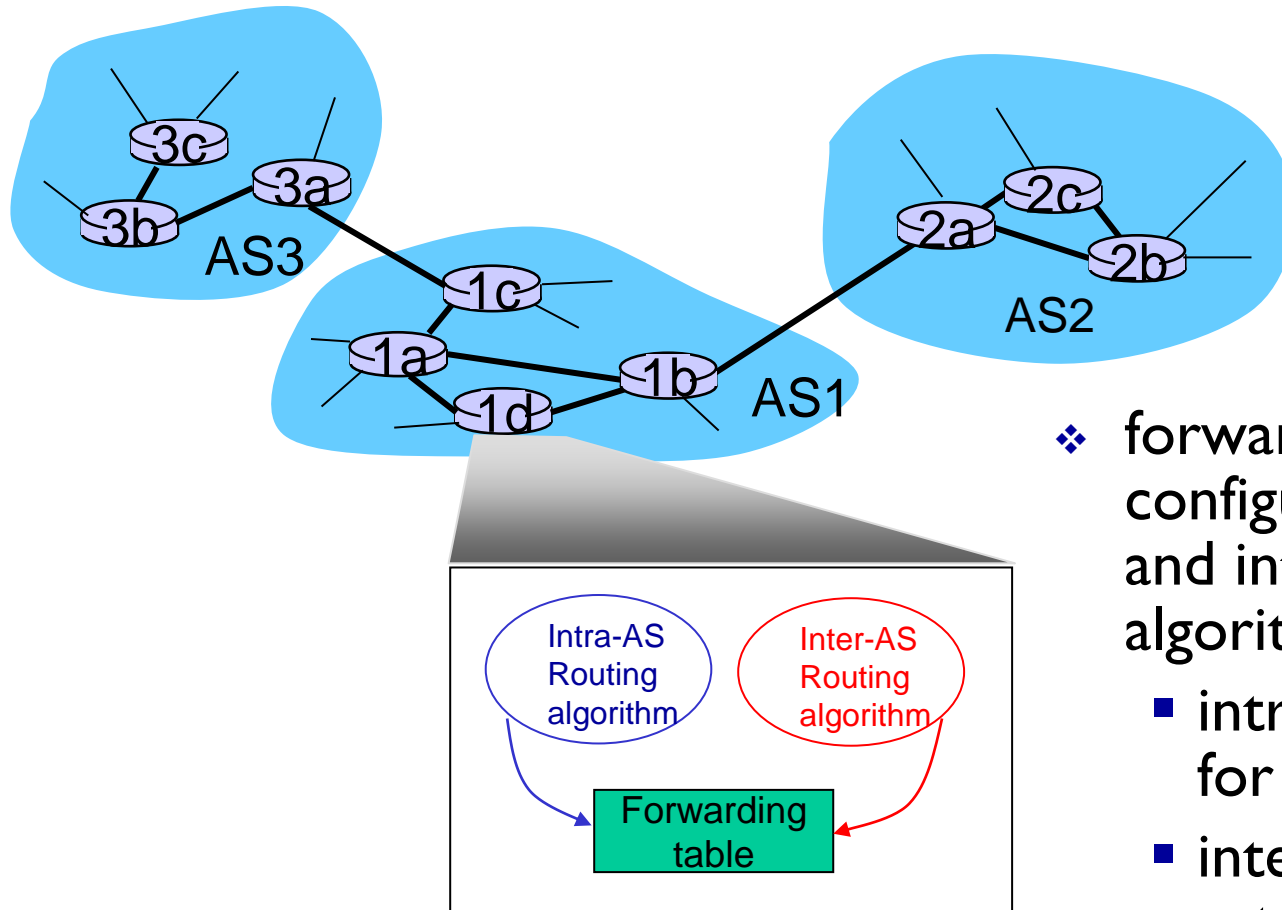
## *scale:*

- ❖ hierarchical routing saves table size, reduced update traffic

## *performance:*

- ❖ intra-AS: can focus on performance
- ❖ inter-AS: policy may dominate over performance

# Interconnected ASes



- ❖ forwarding table configured by both intra- and inter-AS routing algorithm
  - intra-AS sets entries for internal dests
  - inter-AS & intra-AS sets entries for external dests

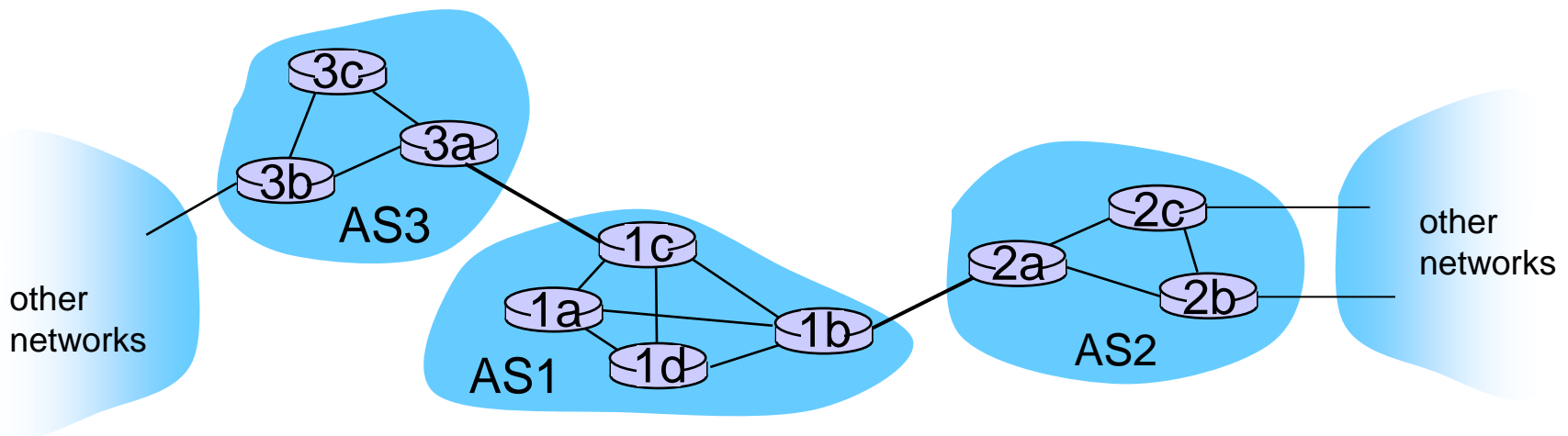
# Inter-AS tasks

- ❖ 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

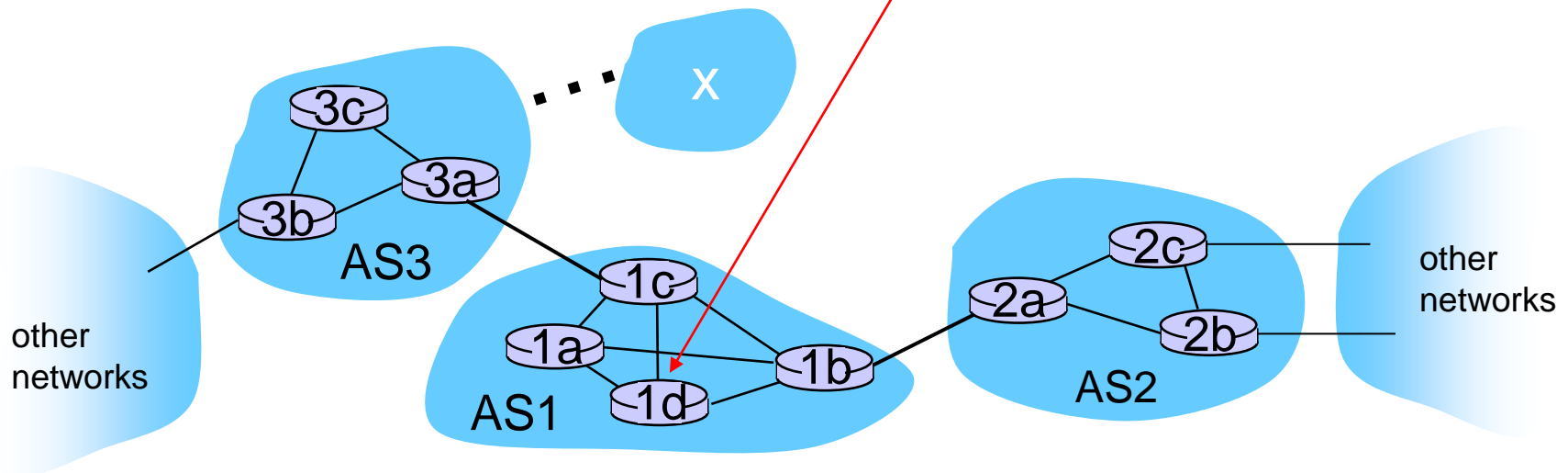
*job of inter-AS routing!*





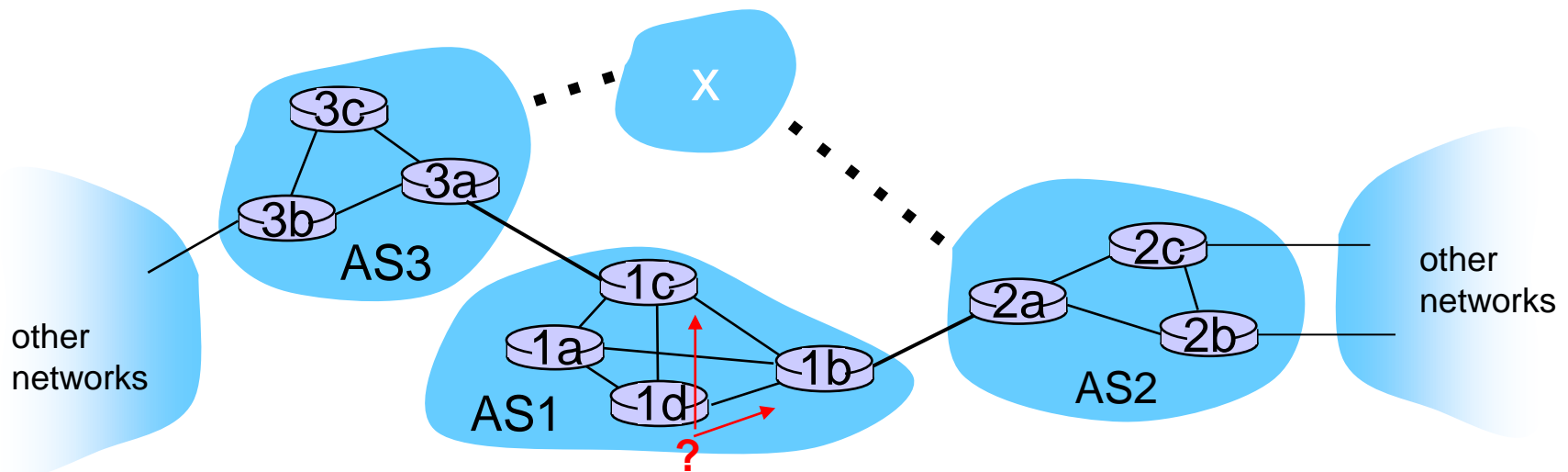
# Example: setting forwarding table in router 1d

- ❖ suppose AS1 learns (via inter-AS protocol) that subnet **x** reachable via AS3 (gateway 1c), but not via AS2
  - inter-AS protocol propagates reachability info to all internal routers
- ❖ router 1d determines from intra-AS routing info that its interface **1** is on the least cost path to 1c
  - installs forwarding table entry **(x,1)**



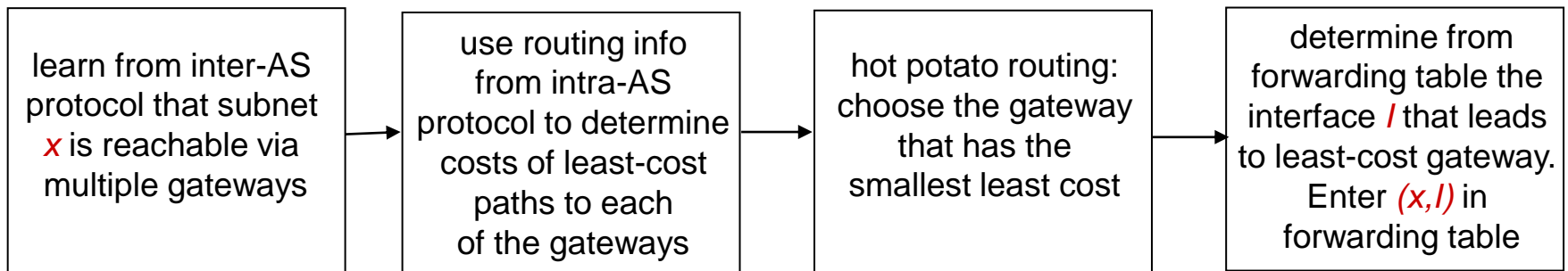
# Example: choosing among multiple ASes

- ❖ now suppose AS1 learns from inter-AS protocol that subnet **x** is reachable from AS3 *and* from AS2.
- ❖ to configure forwarding table, router 1d must determine which gateway it should forward packets towards for dest **x**
  - this is also job of inter-AS routing protocol!



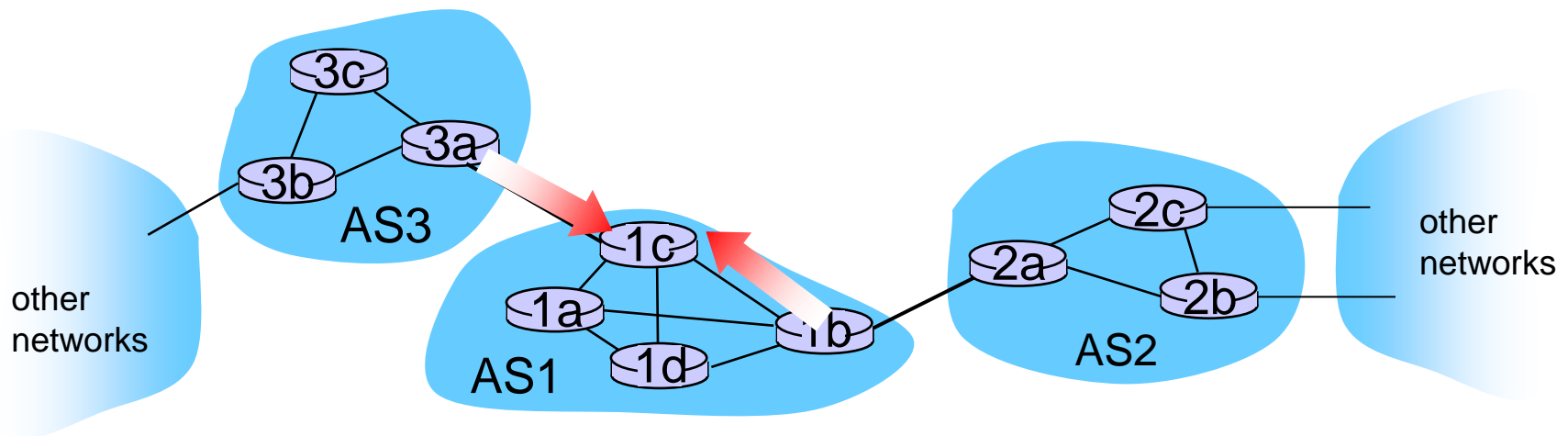
# Example: choosing among multiple ASes

- ❖ now suppose AS1 learns from inter-AS protocol that subnet **x** is reachable from AS3 *and* from AS2.
- ❖ to configure forwarding table, router 1d must determine towards which gateway it should forward packets for dest **x**
  - this is also job of inter-AS routing protocol!
- ❖ **hot potato routing: send** packet towards closest (least cost) of two routers.



# Hot Potato Routing

- ❖ Suppose there two or more best inter-routes.
- ❖ Then choose route with closest NEXT-HOP
  - Use OSPF to determine which gateway is closest
  - Q: From 1c, chose AS3 AS131 or AS2 AS17?
  - A: route AS3 AS131 since it is closer



# Hot-Potato vs Cold-Potato Routing

- **Hot-potato** routing is the practice of passing traffic off to another AS as quickly as possible (thus using their network for wide-area transit.)
  - normal behavior of most peering agreements. It has the effect that the network receiving the data bears the cost of carrying it between cities. When the traffic ratio (traffic in both directions between peers) is reasonably even, this is considered fair.
- **Cold-potato** routing is the opposite: where the source AS holds onto the packet until it is as near to the destination as possible.
  - This is more expensive to do, but keeps the traffic under the network administrator's control for longer, allowing operators of well-provisioned networks to offer a higher quality of service to their customers. It can also be preferred when connecting to content providers.

## Example

- Consider the case of two ISPs, A & B, who both have global networks. Additionally, they have peering agreements in both Europe and in Asia, which allows them to exchange data packets destined for the other's network at either location.
- Suppose a European customer of ISP A wants to transmit a data packet to an Asian customer of ISP B. ISP A will receive the packet in Europe and has to decide where to send the packet next.
- The first option is to hand off the packet to ISP B in Europe, and let ISP B carry the packet to Asia to be delivered to its destination. This is **hot-potato routing**, since ISP A hands off the packet at the earliest opportunity.
- The second option is for ISP A to carry the packet to Asia on its own internal network, and hand off to ISP B in Asia. This is called **cold-potato**, since ISP A keeps the packet in its internal network for as long as possible.

Source: Wikipedia

# Chapter 4: outline

## 4.1 introduction

## 4.2 virtual circuit and datagram networks

## 4.3 what's inside a router

## 4.4 IP: Internet Protocol

- datagram format
- IPv4 addressing
- ICMP
- IPv6

## 4.5 routing algorithms

- link state
- distance vector
- hierarchical routing

## 4.6 routing in the Internet

- RIP
- OSPF
- BGP

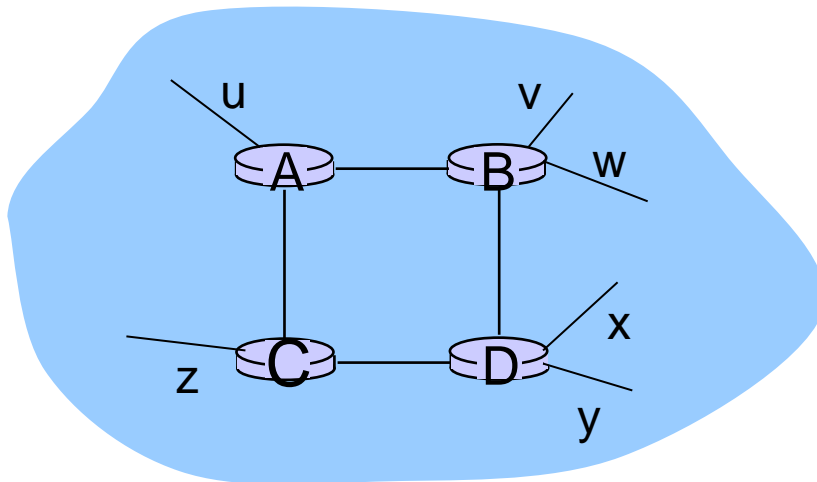
## 4.7 broadcast and multicast routing

# 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
  - *IS-IS (closely related to OSPF)*
  - IGRP: Interior Gateway Routing Protocol (Cisco proprietary)

# RIP ( Routing Information Protocol)

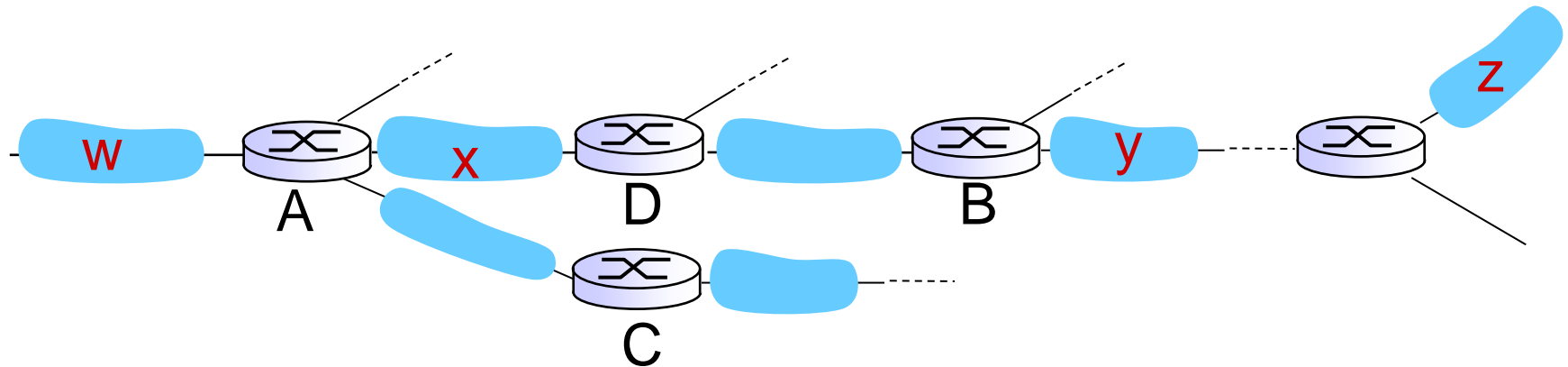
- ❖ included in BSD (Berkeley Software Division) -UNIX distribution in 1982
- ❖ distance vector algorithm
  - distance metric: # hops (max = 15 hops), each link has cost 1 (i.e. uses hop count as cost metric, The maximum cost of a path is limited to 15, thus limiting the use of RIP to autonomous systems that are fewer than 15 hops in diameter)
  - DVs exchanged with neighbors every 30 sec in response message (aka RIP advertisement)
  - each advertisement: list of up to 25 destination *subnets within the AS* (in IP addressing sense)  
from router A to destination subnets:



<u>subnet</u>	<u>hops</u>
u	1
v	2
w	2
x	3
y	3
z	2



# RIP: example



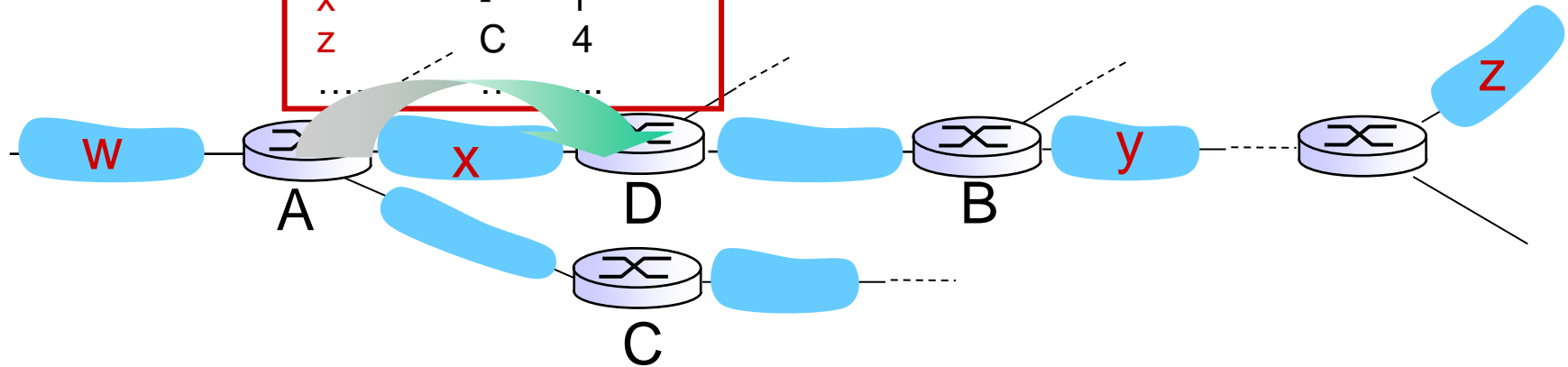
routing table in router D

destination subnet	next router	# hops to dest
W	A	2
y	B	2
z	B	7
x	--	1
....	....	....

# RIP: example

A-to-D advertisement

dest	next	hops
W	-	1
X	-	1
Z	C	4
...	...	...



routing table in router D

destination subnet	next router	# hops to dest
W	A	2
y	B	2
Z	<del>B</del> → A	<del>7</del> → 5
X	--	1
....	....	....

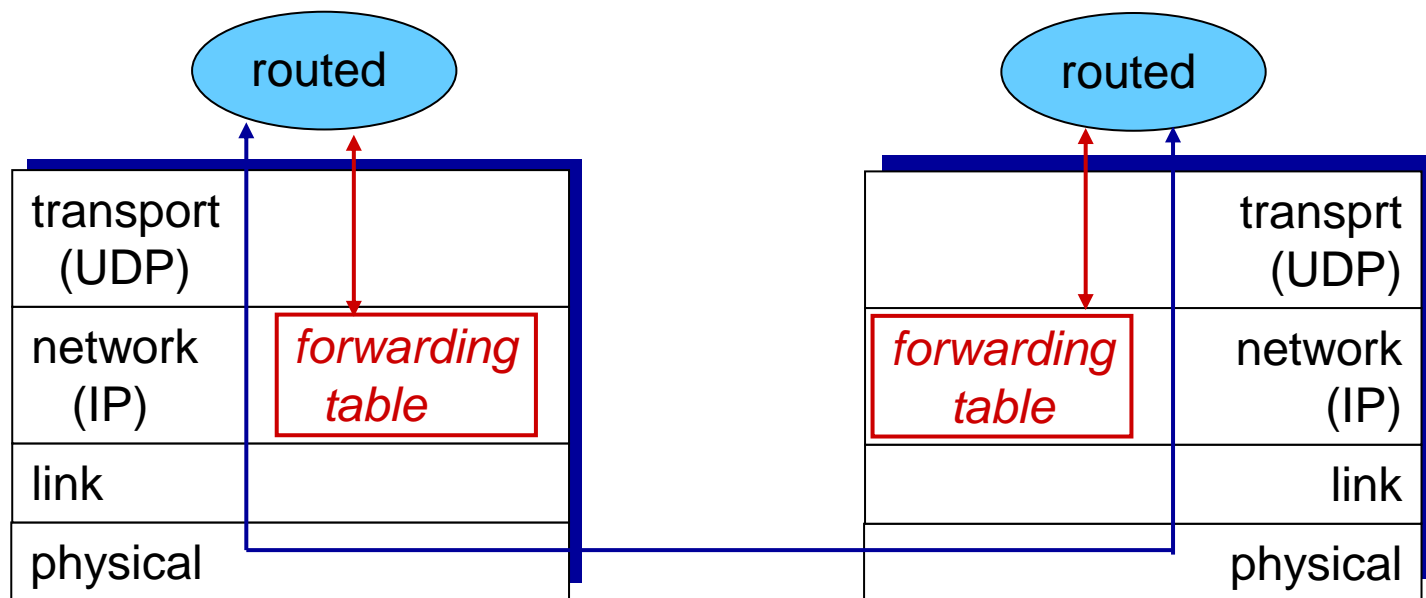
# RIP: link failure, recovery

if no advertisement heard after 180 sec -->  
neighbor/link declared dead

- 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
- *poison reverse* used to prevent ping-pong loops

# RIP table processing

- ❖ RIP routing tables managed by *application-level* process called route-d (daemon)
- ❖ advertisements sent in UDP packets (UDP port # 520), periodically repeated
- ❖ Since it uses UDP, thus it is an App Layer protocol



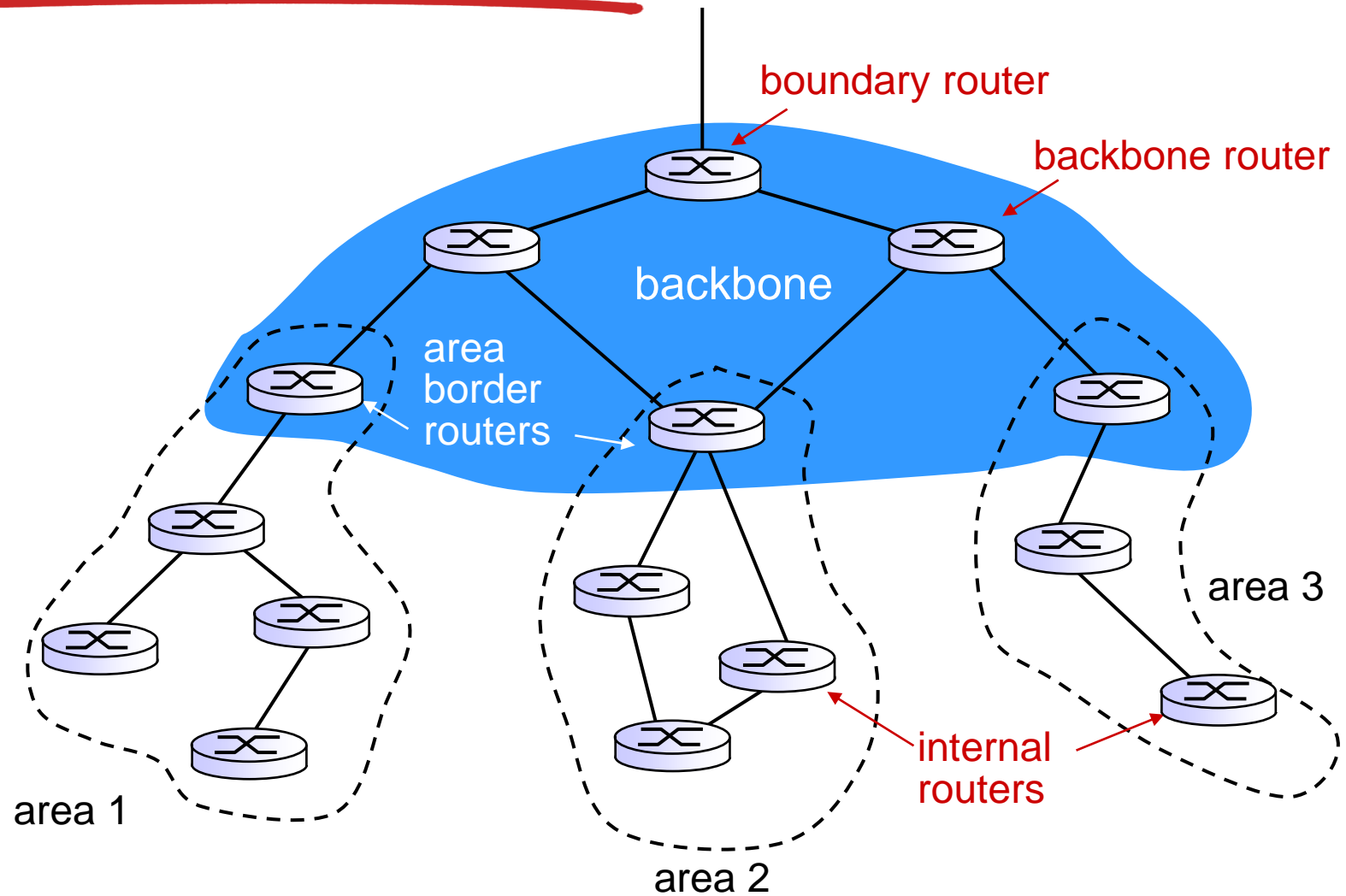
# OSPF (Open Shortest Path First)

- ❖ “open”: publicly available
- ❖ uses link state algorithm
  - LS packet dissemination
  - topology map at each node
  - route computation using Dijkstra’s algorithm
- ❖ OSPF advertisement carries one entry per neighbor
- ❖ advertisements flooded to *entire* AS (every 30 minutes)
  - carried in OSPF messages directly over IP (rather than TCP or UDP) (upper layer field value = 89 for OSPF)
- ❖ *IS-IS routing* protocol: nearly identical to OSPF

## 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 ToS; high for real time ToS)
- ❖ integrated uni- and **multicast** support:
- ❖ **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.



# Chapter 4: outline

## 4.1 introduction

## 4.2 virtual circuit and datagram networks

## 4.3 what's inside a router

## 4.4 IP: Internet Protocol

- datagram format
- IPv4 addressing
- ICMP
- IPv6

## 4.5 routing algorithms

- link state
- distance vector
- hierarchical routing

## 4.6 routing in the Internet

- RIP
- OSPF
- BGP

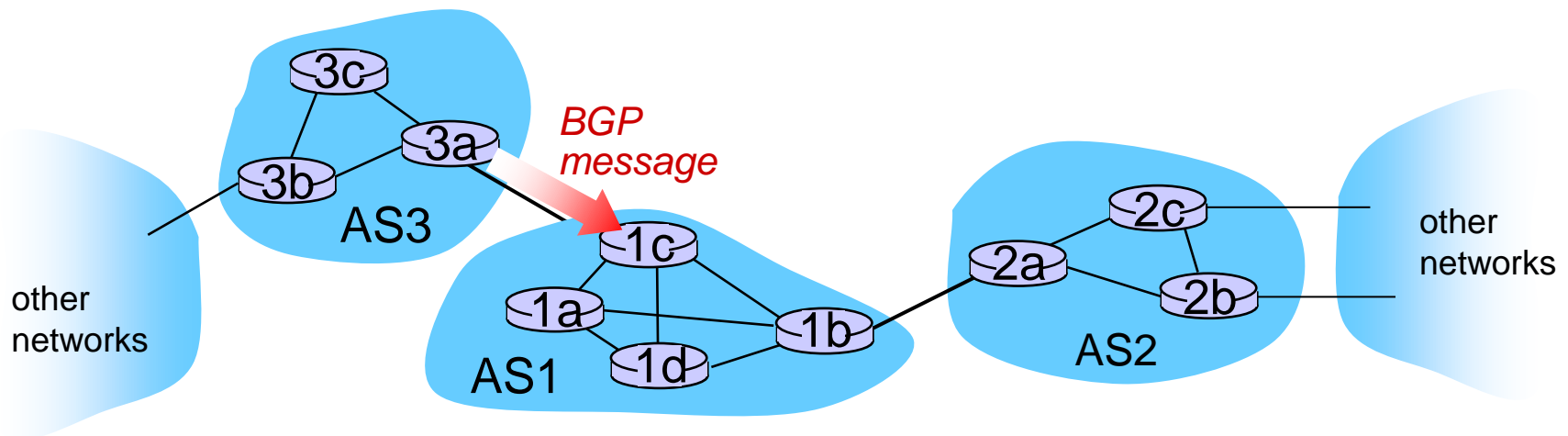
## 4.7 broadcast and multicast routing

# Internet inter-AS routing: BGP

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

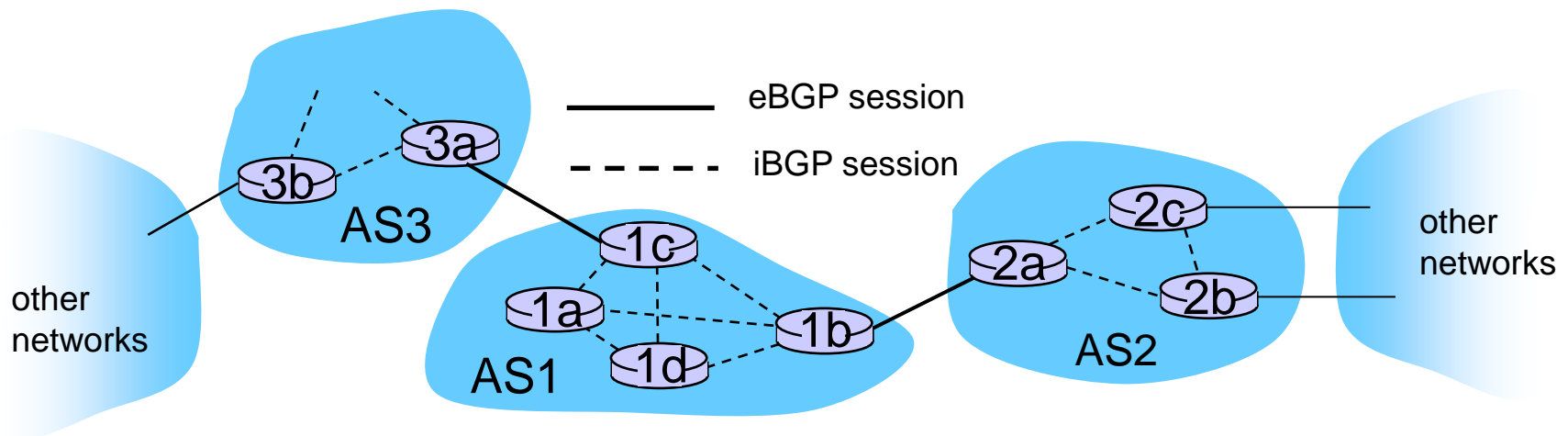
# BGP basics

- ❖ **BGP session:** two BGP routers (“peers”) exchange BGP messages:
  - advertising *paths* to different destination network prefixes (“path vector” protocol)
  - exchanged over semi-permanent TCP connections **using port 179**
- ❖ when AS3 advertises a prefix to AS1:
  - AS3 *promises* it will forward datagrams towards that prefix
  - AS3 can aggregate prefixes in its advertisement



# BGP basics: distributing path information

- ❖ using eBGP session between 3a and 1c, AS3 sends prefix reachability info to AS1.
  - 1c can then use iBGP to distribute new prefix info to all routers in AS1 **including gateway router 1b**
  - 1b can then re-advertise new reachability (**AS3's**) info to AS2 over 1b-to-2a eBGP session
- ❖ when router learns of new prefix, it creates entry for prefix in its forwarding table.



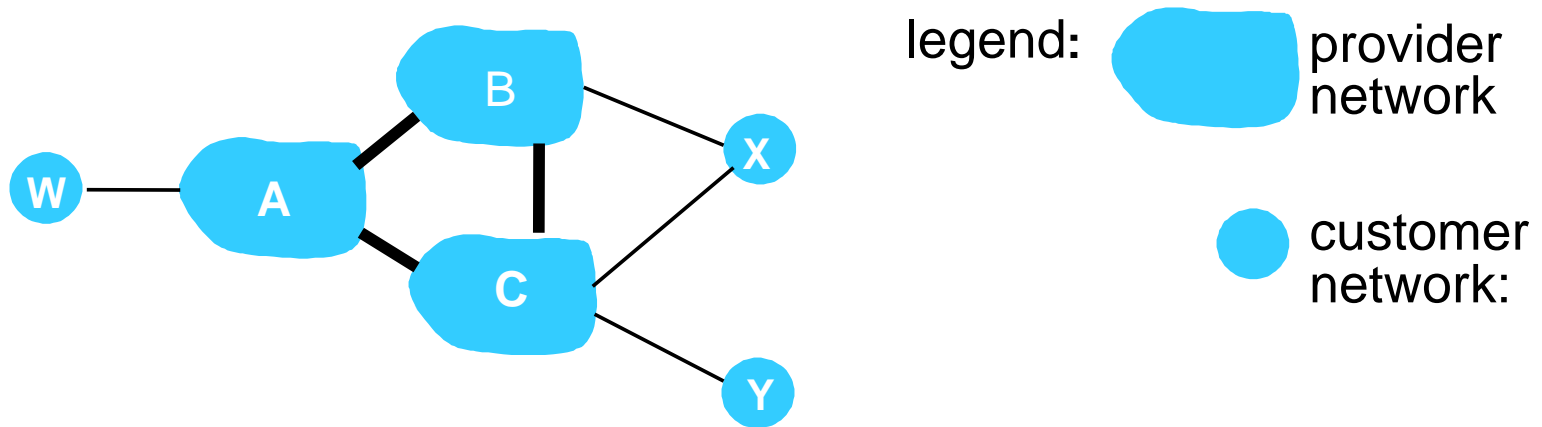
# Path attributes and BGP routes

- ❖ advertised prefix includes BGP attributes
  - prefix + attributes = “route”
- ❖ two important attributes:
  - **AS-PATH**: contains ASs through which prefix advertisement has passed: e.g., AS 67, AS 17
  - **NEXT-HOP**: indicates specific internal-AS router to next-hop AS. (may be multiple links from current AS to next-hop-AS)
- ❖ gateway router receiving route advertisement uses **import policy** to accept/decline
  - e.g., never route through AS x
  - *policy-based* routing

# BGP route selection

- ❖ router may learn about more than 1 route to destination AS, selects route based on:
  1. local preference value attribute: policy decision
  2. shortest AS-PATH
  3. closest NEXT-HOP router: hot potato routing
  4. additional criteria

# BGP routing policy



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

# How does entry get in forwarding table?

## Summary

1. Router becomes aware of prefix
  - via BGP route advertisements from other routers
2. Determine router output port for prefix
  - Use BGP route selection to find best inter-AS route
  - Use OSPF (typically) to find best intra-AS route leading to best inter-AS route
  - Router identifies router port for that best route
3. Enter prefix-port entry in forwarding table



# BGP messages

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

# Chapter 4: *done!*

---

# Quiz # 4 (Chapter - 4)

- *On: Tuesday 8<sup>th</sup> November, 2022 (During the lecture)*
- *Topics Included from Chapter 4 of the textbook:*
  - *4.1*
  - *4.4*
- *Quiz to be taken during own section class only*