# Logistics
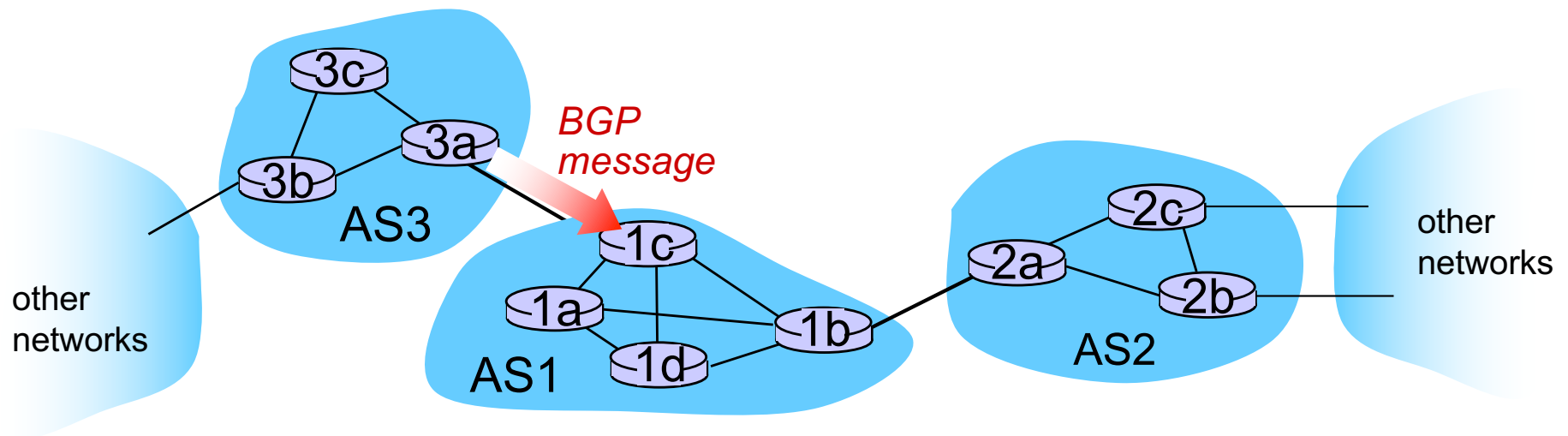
❖ Reading and HW5
  ▪ Chapter 5 (seventh edition)
  ▪ Tues 4/4 at 10pm

❖ Wireshark Labs
  ▪ 2 of them! IP and ICMP
  ▪ Both due Sunday 4/9
  ▪ <u>May work with a partner</u>, write both names

❖ Today:
  ▪ Network layer, moving towards link layer

# How does entry get in forwarding table?
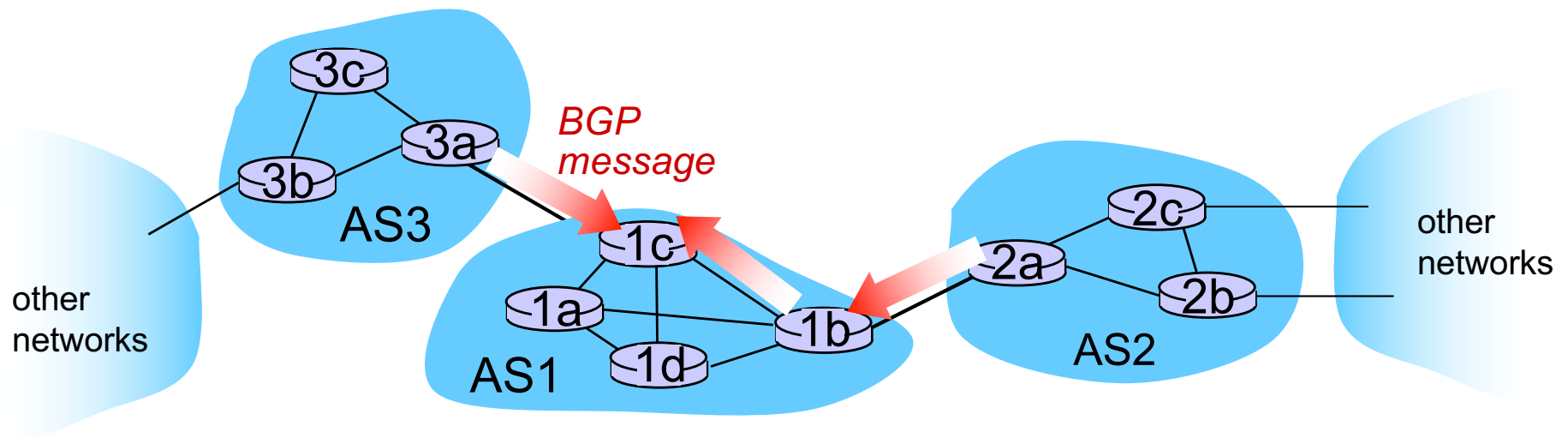
## High-level overview

1. Router becomes aware of prefix
2. Router determines output port for prefix
3. Router enters prefix-port in forwarding table

# Router becomes aware of prefix



- ❖ BGP message contains "routes"
- ❖ "route" is a prefix and attributes: AS-PATH, NEXT-HOP,...
- ❖ Example: route:
  - ❖ Prefix:138.16.64/22 ;  AS-PATH:  AS3  AS131 ;  NEXT-HOP:  201.44.13.125

# Router may receive multiple routes



- ❖ Router may receive multiple routes for <u>same</u> prefix
- ❖ Has to select one route

# Select best BGP route to prefix

- ❖ Router selects route based on shortest AS-PATH
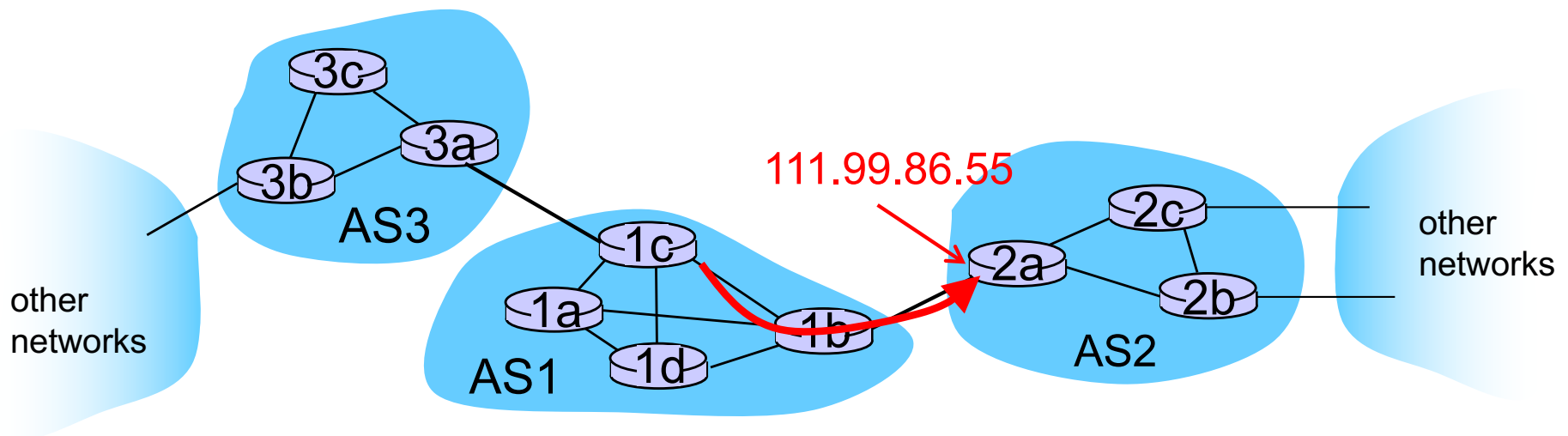
- ❖ Example:

  select

  - ❖ AS2 AS17  to 138.16.64/22
  - ❖ AS3 AS131 AS201 to 138.16.64/22
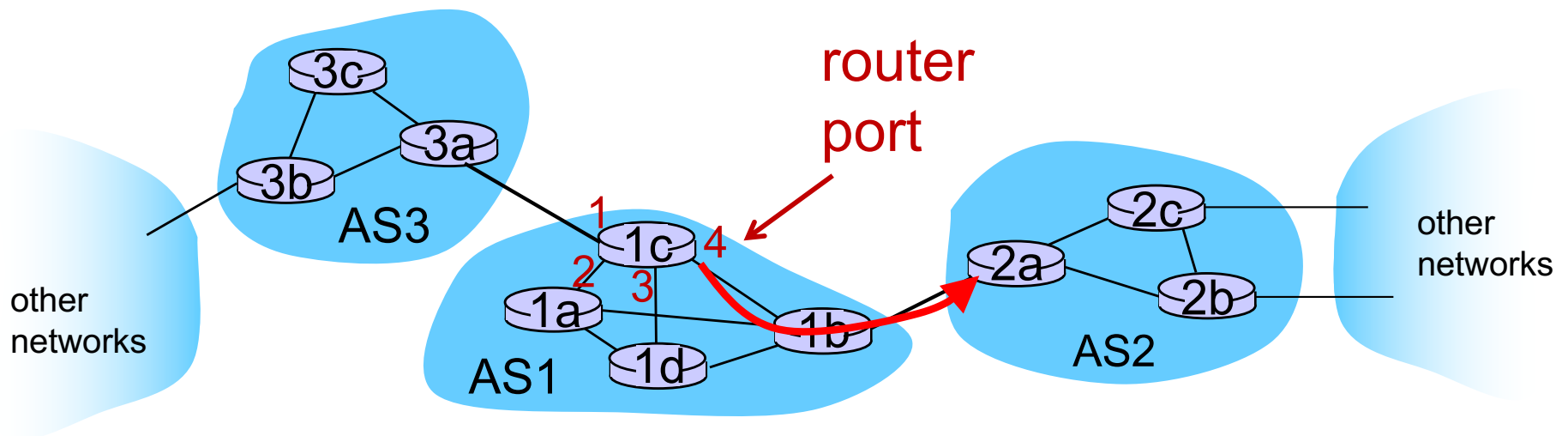
- ❖ What if there is a tie?

# Find best intra-route to BGP route

- ❖ Use selected route's NEXT-HOP attribute
  - ■ Route's NEXT-HOP attribute is the IP address of the router interface that begins the AS PATH.
- ❖ Example:
  - ❖ AS-PATH: AS2 AS17 ; NEXT-HOP: 111.99.86.55
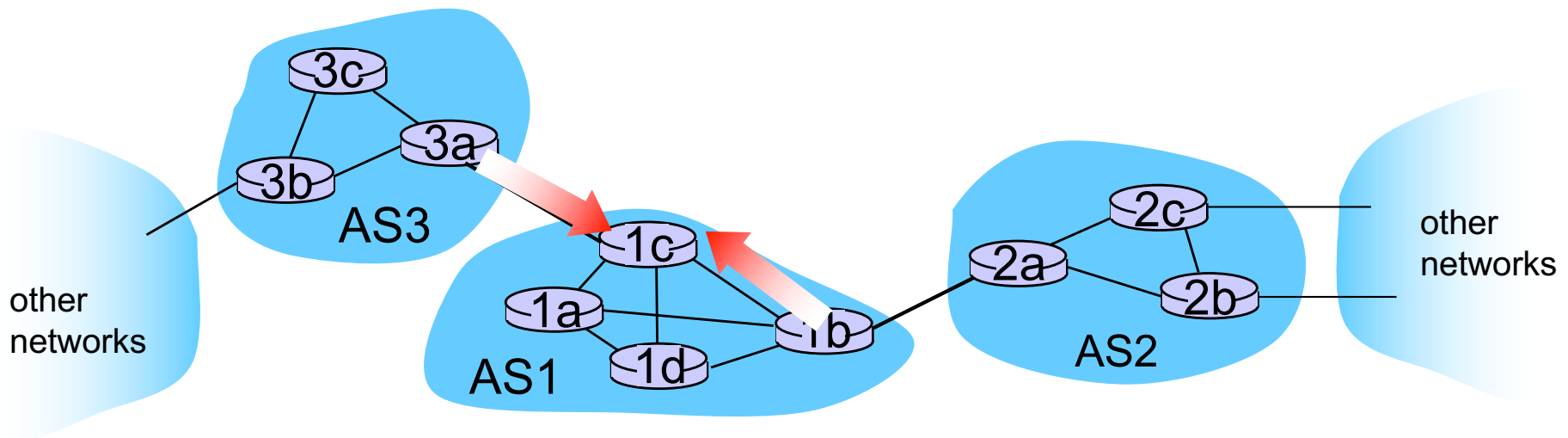- ❖ Router uses OSPF to find shortest path from 1c to 111.99.86.55

# Router identifies port for route

❖ Identifies port along the OSPF shortest path
❖ Adds prefix-port entry to its forwarding table:
  ▪ (138.16.64/22 , port 4)

# 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?
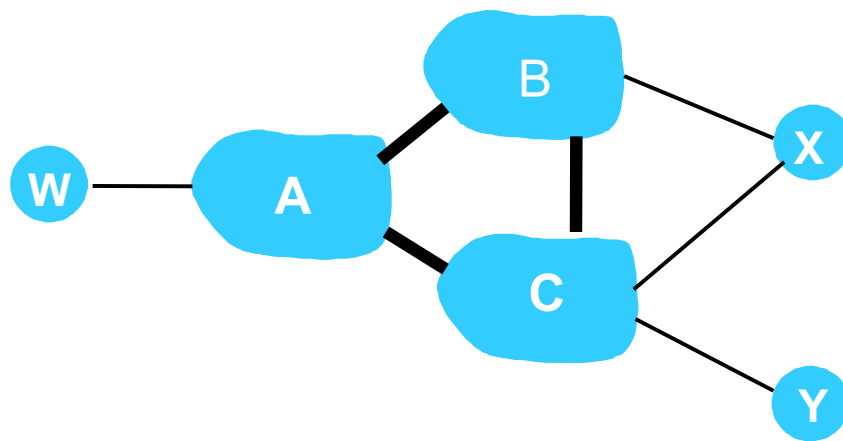- A: route AS3 AS201 since it is closer

# 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 to find best intra-AS route  leading to best inter-AS route
   - Router identifies router port for that best route
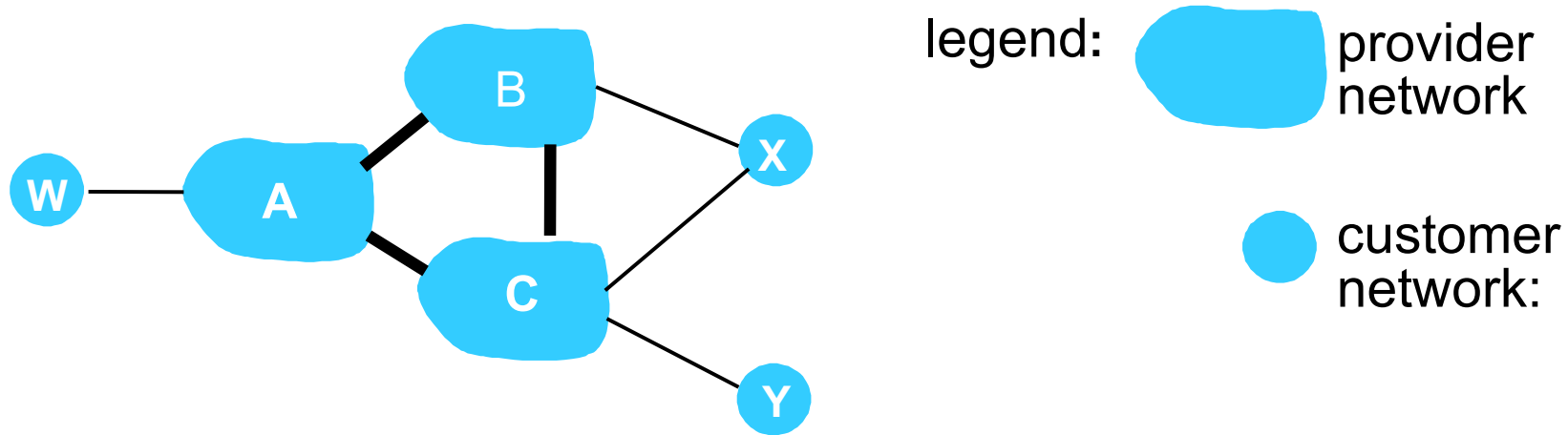3. Enter prefix-port entry in forwarding table

# BGP routing policy



legend:
provider network

customer network:

- A,B,C are *provider networks*
- X,W,Y are customer (of provider networks)
- X is *dual-homed:* 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

# BGP routing policy (2)



legend:
provider network

customer network:

- ❖ A advertises path AW to B
- ❖ B advertises path BAW to X
- ❖ Should B advertise path BAW to C?
    - No way! B gets no "revenue" for routing CBAW since neither W nor C are B's customers
    - B wants to force C to route to w via A
    - B wants to route *only* to/from its customers!

# 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

# 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
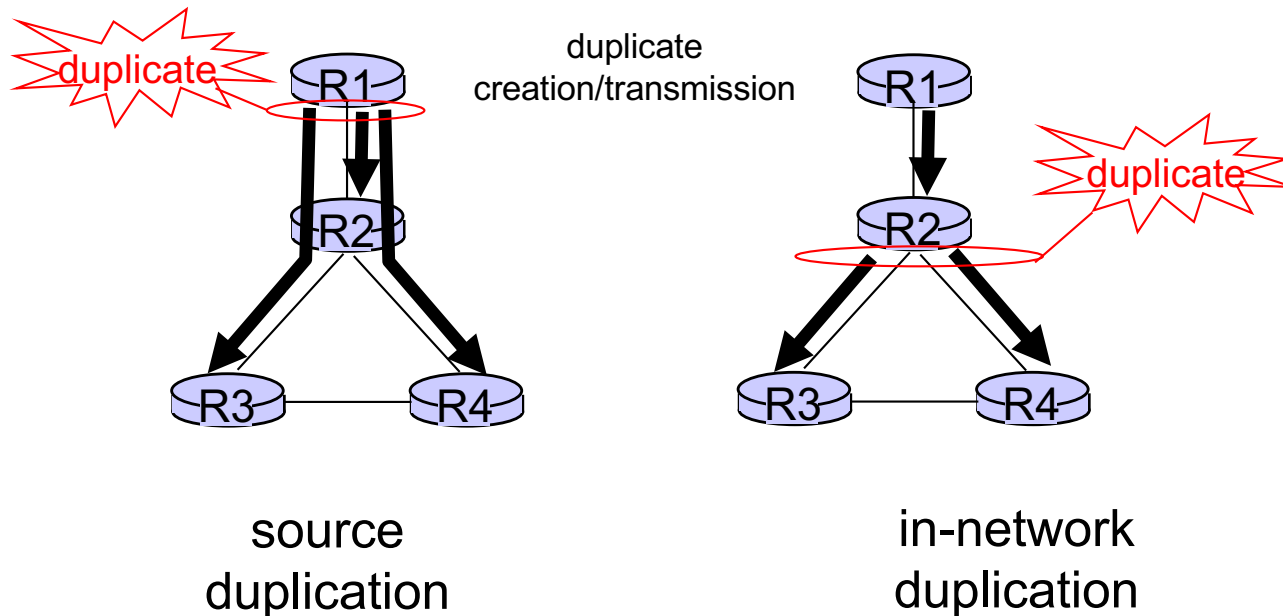- distance vector
- hierarchical routing

4.6 routing in the Internet
- RIP
- OSPF
- BGP

4.7 broadcast and multicast routing

# Broadcast routing

❖ deliver packets from source to all other nodes
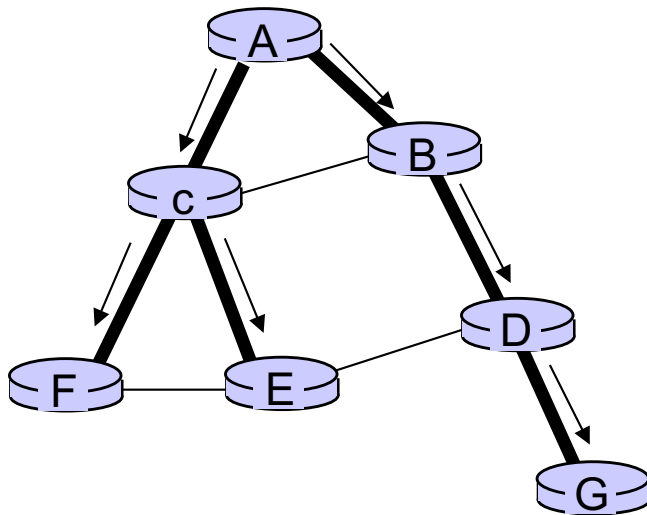
❖ source duplication is inefficient:



source
duplication

in-network
duplication

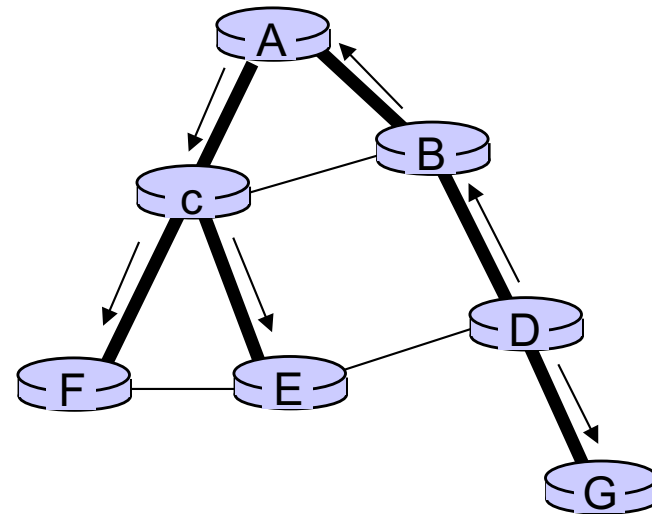❖ source duplication: how does source determine recipient addresses?

# In-network duplication

❖ *flooding:* when node receives broadcast packet, sends copy to all neighbors
  ▪ problems: cycles & broadcast storm
❖ *controlled flooding:* node only broadcasts pkt if it hasn't broadcast same packet before
  ▪ node keeps track of packet ids already broadacsted
  ▪ or reverse path forwarding (RPF): only forward packet if it arrived on shortest path between node and source
❖ *spanning tree:*
  ▪ no redundant packets received by any node

# Spanning tree

- ❖ first construct a spanning tree
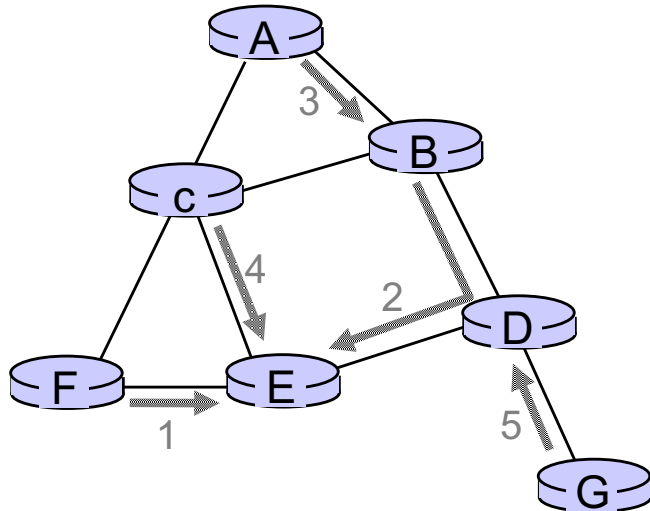- ❖ nodes then forward/make copies only along spanning tree
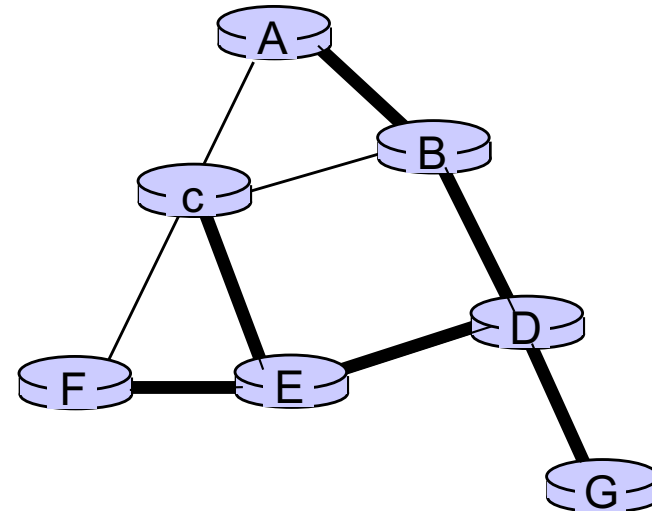


(a) broadcast initiated at A

(b) broadcast initiated at D

# Spanning tree: creation

❖ center node

❖ each node sends unicast join message to center node

  ■ message forwarded until it arrives at a node already belonging to spanning tree



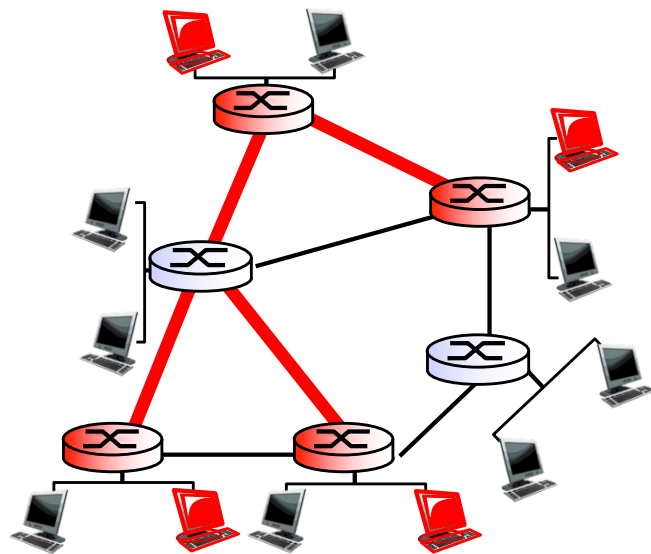(a) stepwise construction of spanning tree (center: E)
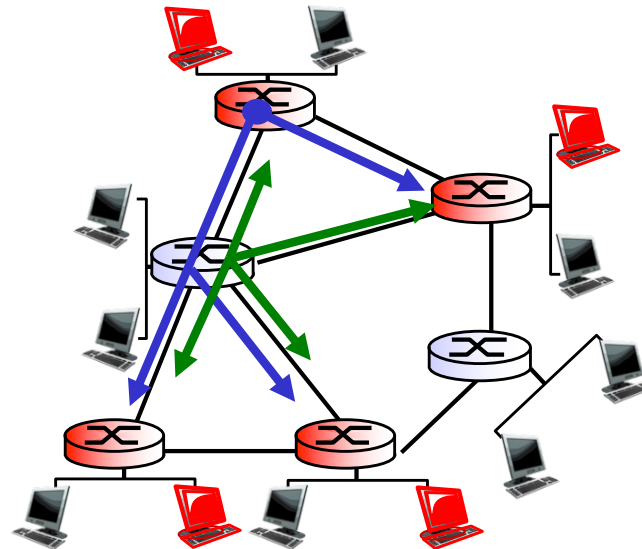
(b) constructed spanning tree

# Multicast routing: problem statement

*goal:* find a tree (or trees) connecting routers having local mcast group members

- ❖ *tree:* not all paths between routers used
- ❖ *shared-tree:* same tree used by all group members
- ❖ *source-based:* different tree from each sender to rcvrs



shared tree

source-based trees

**legend**

group member

not group member

router with a group member

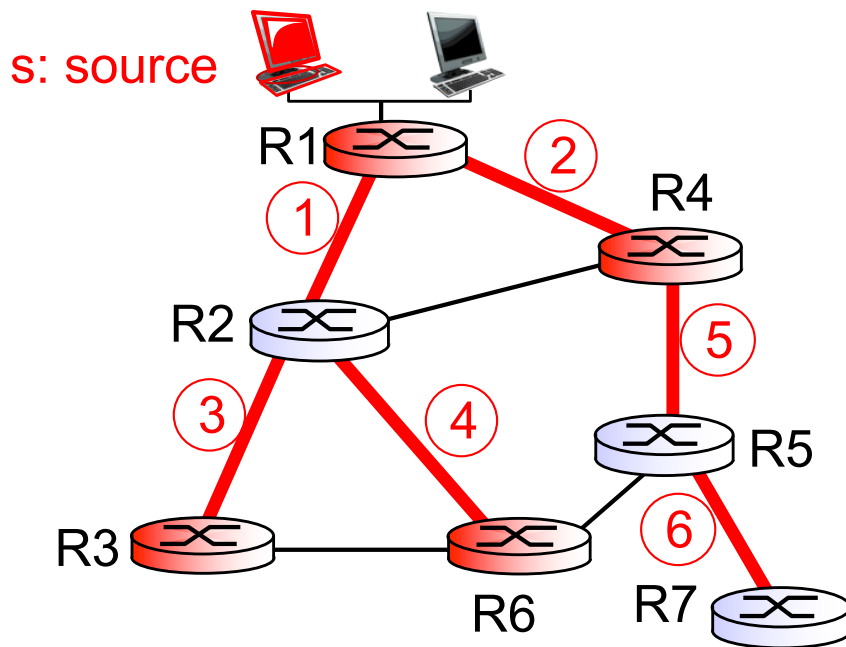router without group member

# Approaches for building mcast trees

approaches:

❖ *source-based tree:* one tree per source
  - shortest path trees
  - reverse path forwarding

❖ *group-shared tree:* group uses one tree
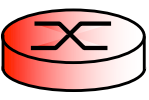  - minimal spanning (Steiner)
  - center-based trees

…we first look at basic approaches, then specific protocols adopting these approaches

# Shortest path tree

❖ mcast forwarding tree: tree of shortest path routes from source to all receivers
  ▪ Dijkstra's algorithm

s: source

LEGEND

router with attached group member

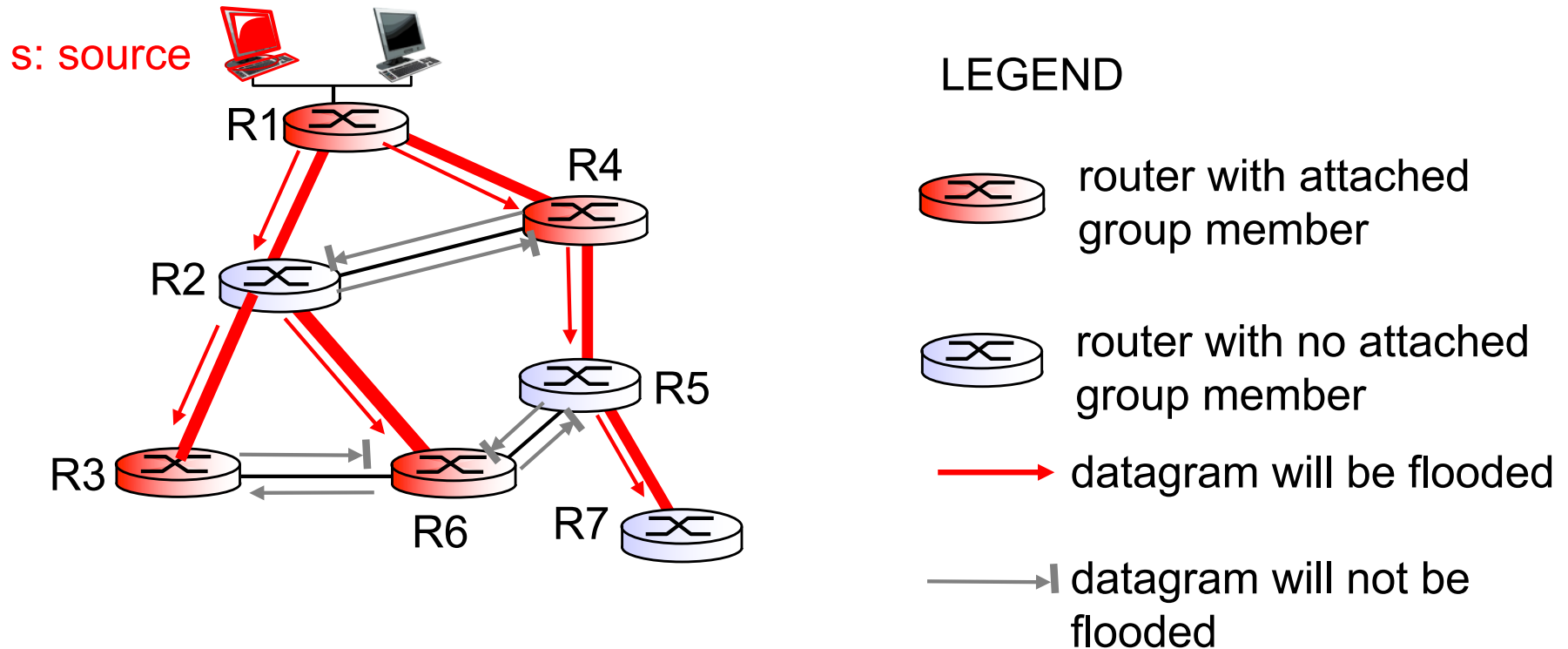router with no attached group member

(i) link used for forwarding, i indicates order link added by algorithm

# Reverse path forwarding

❖ rely on router's knowledge of unicast shortest path from it to sender

❖ each router has simple forwarding behavior:

*if* (mcast datagram received on incoming link on
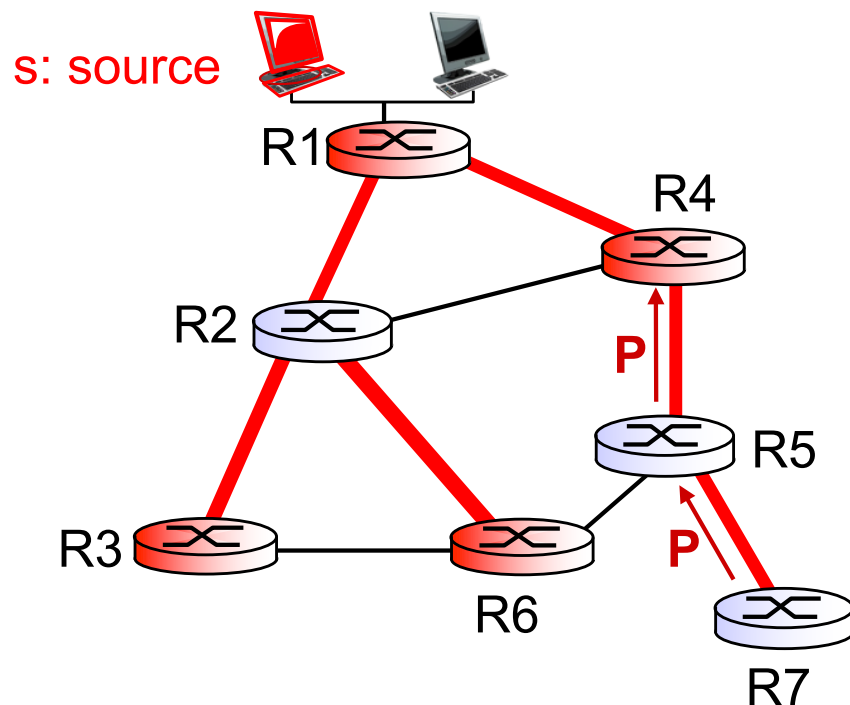   shortest path back to center)
  *then* flood datagram onto all outgoing links
  *else* ignore datagram

# Reverse path forwarding: example



s: source

LEGEND

router with attached
group member

router with no attached
group member

→ datagram will be flooded

→| datagram will not be
flooded

❖ result is a source-specific *reverse* SPT
  ▪ may be a bad choice with asymmetric links

# Reverse path forwarding: pruning

- ❖ forwarding tree contains subtrees with no mcast group members
    - no need to forward datagrams down subtree
    - "prune" msgs sent upstream by router with no downstream group members



s: source

LEGEND

router with attached group member

router with no attached group member

P → prune message
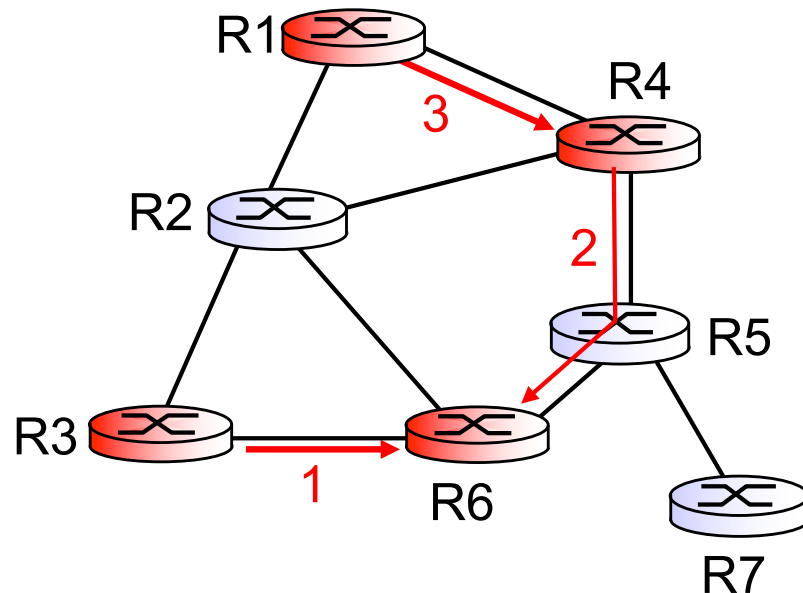
links with multicast forwarding

# Shared-tree: steiner tree

❖ *steiner tree:* minimum cost tree connecting all routers with attached group members
❖ problem is NP-complete
❖ excellent heuristics exists
❖ not used in practice:
  ▪ computational complexity
  ▪ information about entire network needed
  ▪ monolithic: rerun whenever a router needs to join/leave

# Center-based trees

- ❖ single delivery tree shared by all
- ❖ one router identified as *"center"* of tree
- ❖ to join:
  - ▪ edge router sends unicast *join-msg* addressed to center router
  - ▪ *join-msg* "processed" by intermediate routers and forwarded towards center
  - ▪ *join-msg* either hits existing tree branch for this center, or arrives at center
  - ▪ path taken by *join-msg* becomes new branch of tree for this router

# Center-based trees: example

suppose R6 chosen as center:

LEGEND

router with attached group member

router with no attached group member

1 → path order in which join messages generated

# Internet Multicasting Routing: DVMRP
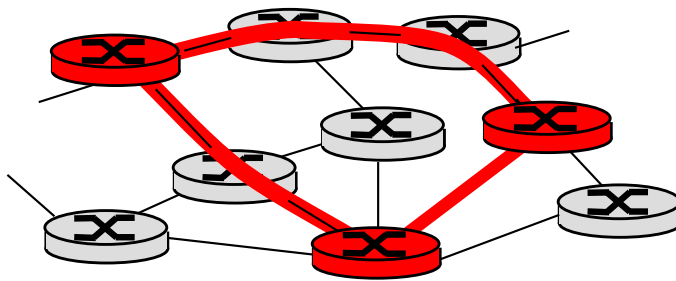
❖ **DVMRP:** distance vector multicast routing protocol, RFC1075

❖ *flood and prune:* reverse path forwarding, source-based tree

  ▪ RPF tree based on DVMRP's own routing tables constructed by communicating DVMRP routers

  ▪ no assumptions about underlying unicast

  ▪ initial datagram to mcast group flooded everywhere via RPF

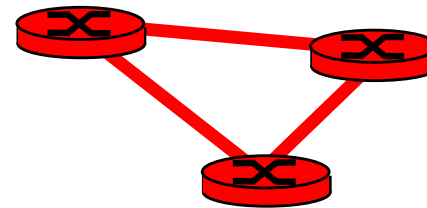  ▪ routers not wanting group: send upstream prune msgs

# DVMRP: continued...

❖ *soft state:* DVMRP router periodically (1 min.) "forgets" branches are pruned:
  - mcast data again flows down unpruned branch
  - downstream router: reprune or else continue to receive data

❖ routers can quickly regraft to tree
  - following IGMP join at leaf

❖ odds and ends
  - commonly implemented in commercial router

# Tunneling

*Q:* how to connect "islands" of multicast routers in a "sea" of unicast routers?



physical topology          logical topology

- ❖ mcast datagram encapsulated inside "normal" (non-multicast-addressed) datagram
- ❖ normal IP datagram sent thru "tunnel" via regular IP unicast to receiving mcast router (recall IPv6 inside IPv4 tunneling)
- ❖ receiving mcast router unencapsulates to get mcast datagram

# Chapter 4/5: *done!*

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

❖ understand principles behind network layer services:
- network layer service models, forwarding versus routing how a router works, routing (path selection), broadcast, multicast

❖ instantiation, implementation in the Internet