# Chapter 4
# Network Layer

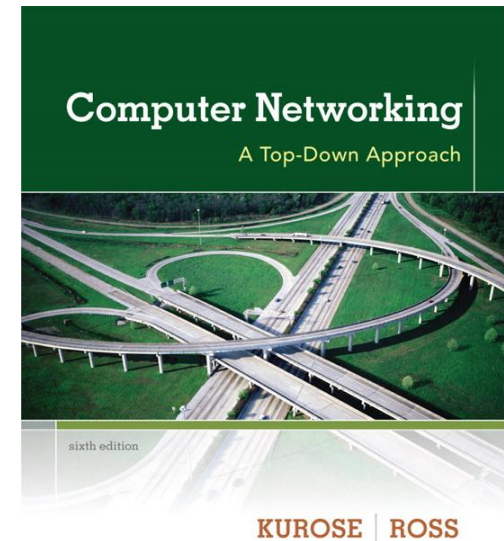## Part 3 (of 3):
## BGP & Broadcast/multicast

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

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

# Chapter 4: outline

# Hierarchical OSPF

boundary router

backbone router

backbone

area border routers

internal routers

area 1

area 2
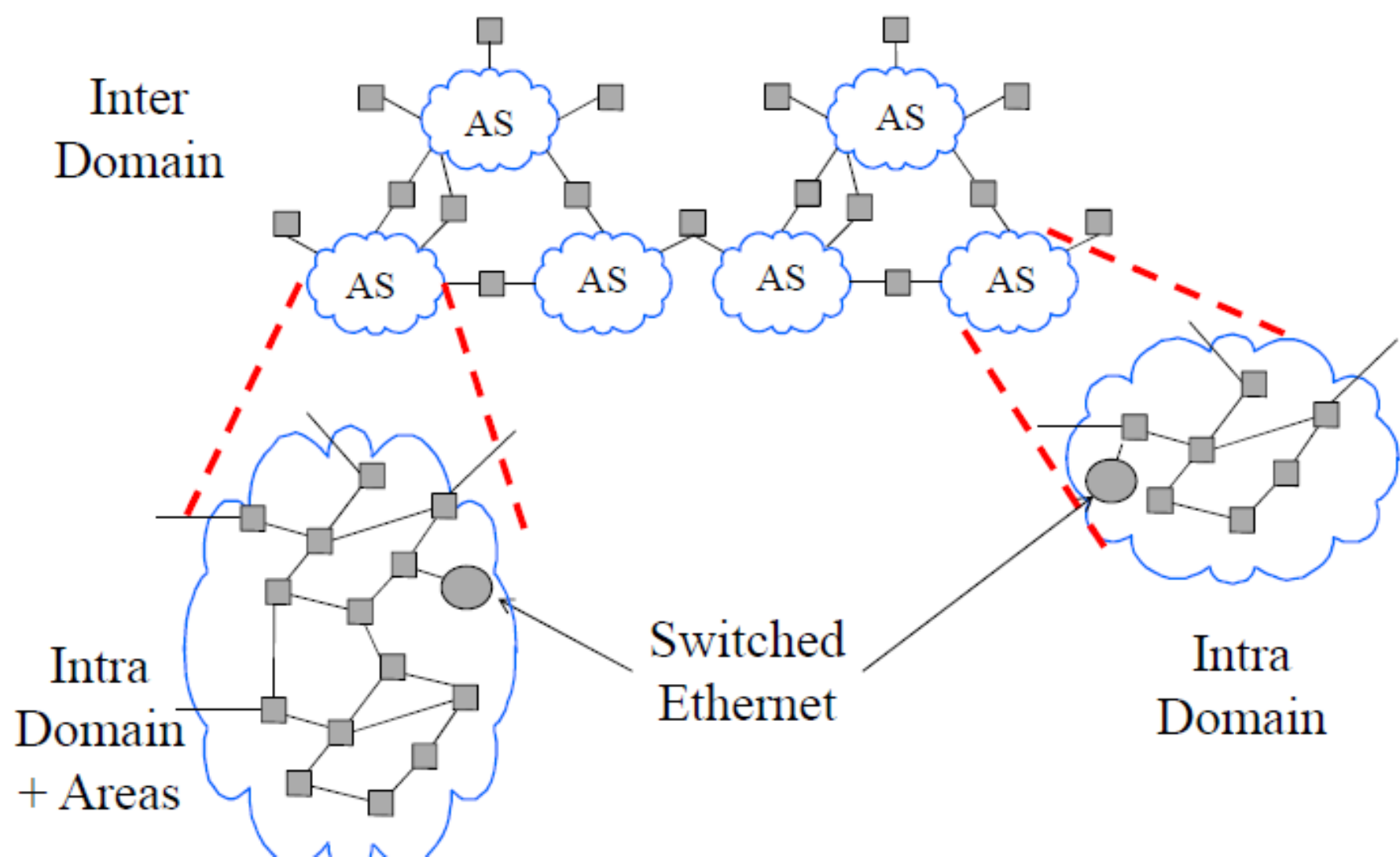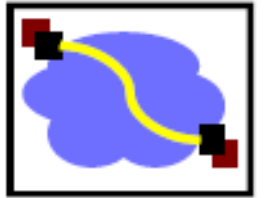
area 3

# 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.

# Inter and Intra-Domain Routing

Inter Domain

AS

AS

AS  AS  AS  AS

Intra Domain + Areas

Switched Ethernet

Intra Domain

# Internet's Area Hierarchy

- ## What is an Autonomous System (AS)?
  - A set of routers under a single technical administration, using an *interior gateway protocol (IGP)* and common metrics to route packets within the AS and using an *exterior gateway protocol (EGP)* to route packets to other AS's

- ## Each AS assigned unique ID
  - Only transit domains really need it

- ## ASes peer with other ASes at network exchanges
  - "Gateway routers" forward packets across ASes

# AS Numbers (ASNs)

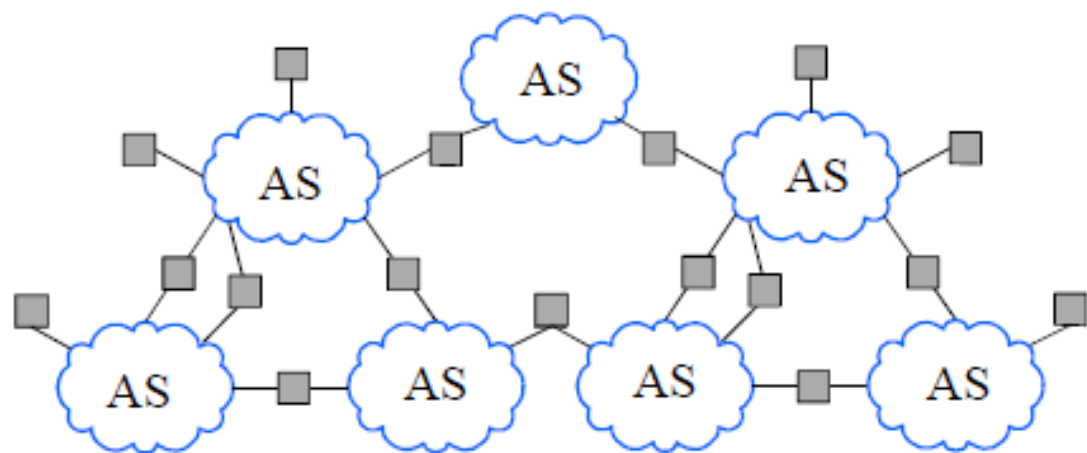ASNs are 16 bit values     64512 through 65535 are "private"

- Genuity: 1
- MIT: 3
- CMU: 9
- UC San Diego: 7377
- AT&T: 7018, 6341, 5074, …
- UUNET: 701, 702, 284, 12199, …
- Sprint: 1239, 1240, 6211, 6242, …
- …

ASNs represent units of routing policy
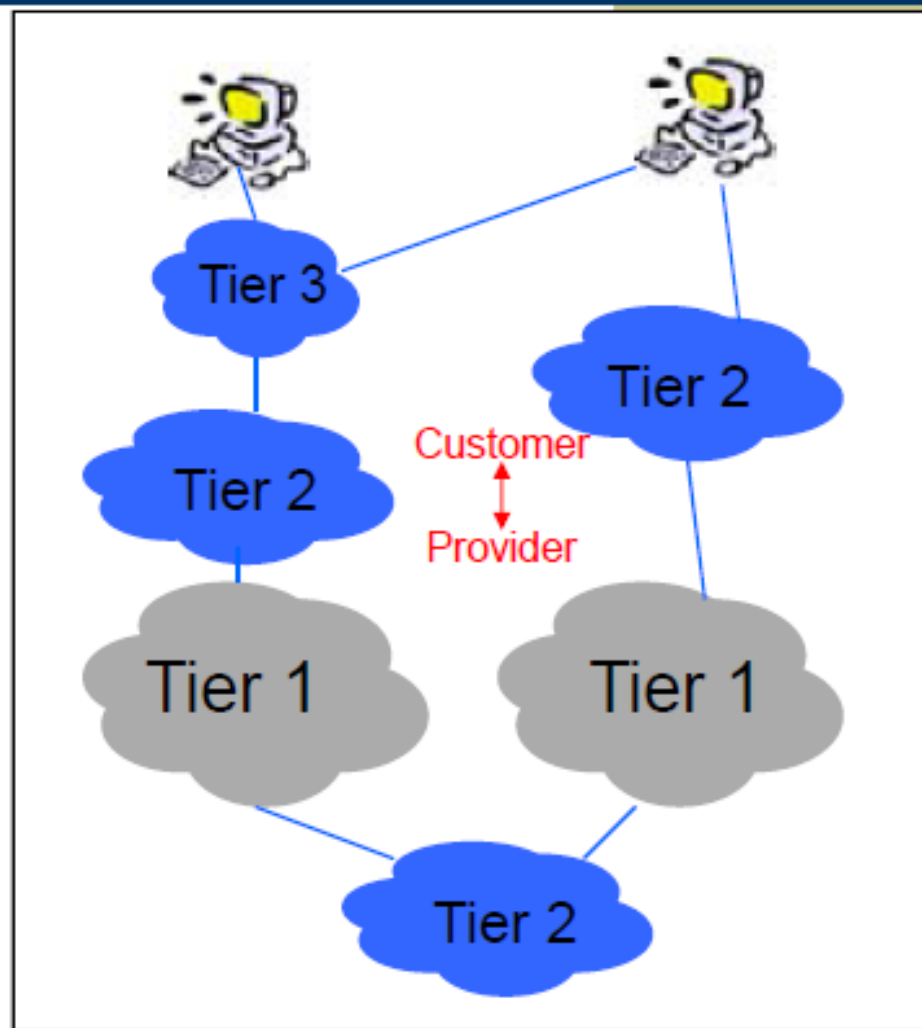
# A Logical View of the Internet?



- Logical consequence of hierarchy: repeat the intra-domain solutions at inter-net level
  - Based on IP and OSPF style routing protocols
- NOT TRUE!
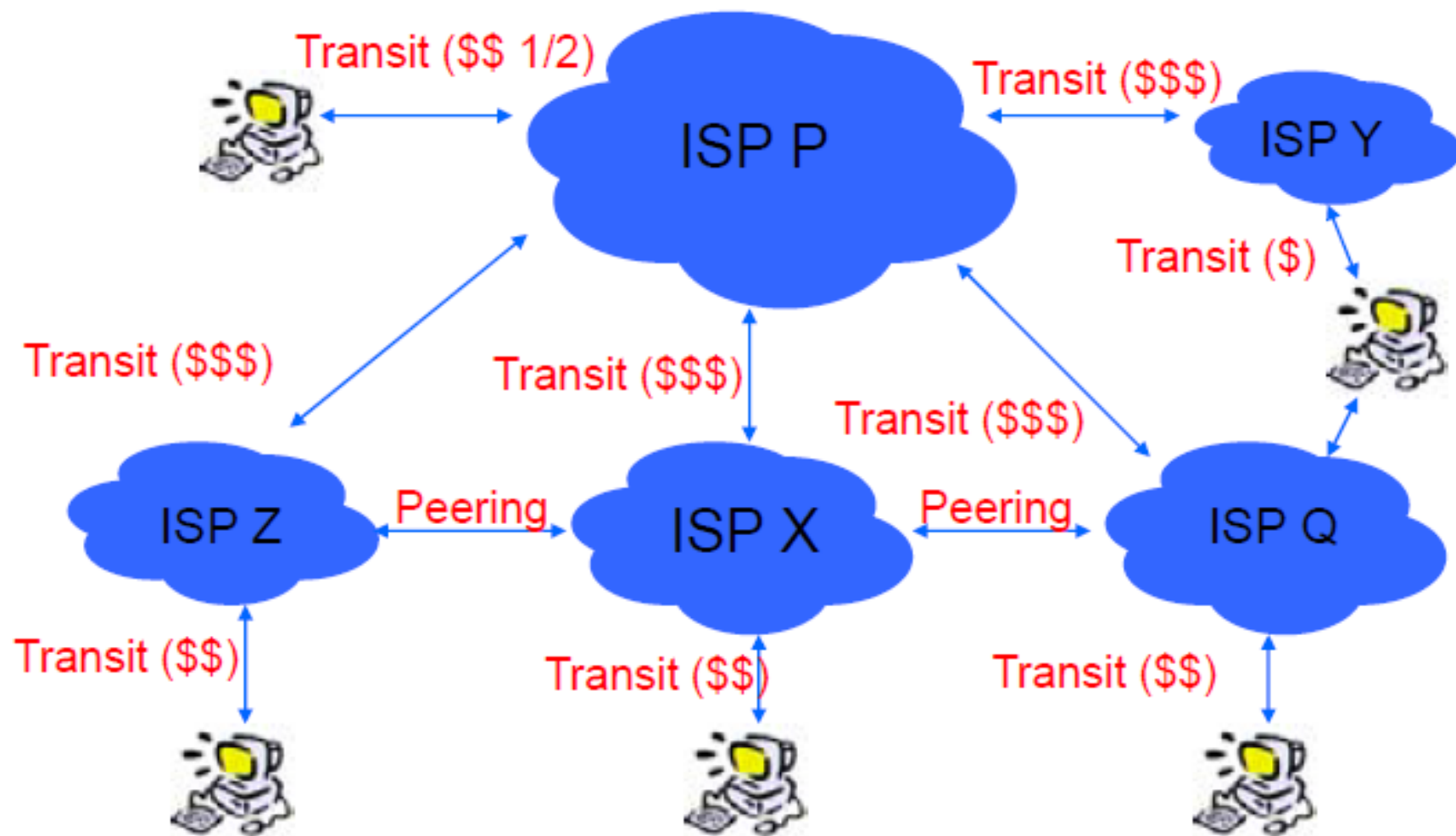  - Lots of problems with this picture

# A Logical View of the Internet

- ASes play different roles in the Internet
- Tier 1 ISP: gobal, internet wide connectivity
- Tier 2 ISP: regional or country-wide
- Tier 3 ISP: local
- Emergent property:
  - Businesses specialize
  - Business relationships
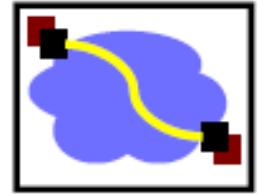
# A More Interesting Example

Transit ($$ 1/2)

Transit ($$$)

ISP P

Transit ($$$)

ISP Y

Transit ($)

Transit ($$$)

Transit ($$$)

Transit ($$$)

ISP Z

Peering

ISP X

Peering

ISP Q

Transit ($$)

Transit ($$)

Transit ($$)

# Policy and Economics Rules

- ## WHY?
  - Consider the economics of the Internet
  - Why does an ISP forward packets?

- ## Emergent property: "Valley-free" routing
  - Number links as (+1, 0, -1) for provider, peer and customer
  - In any path should only see sequence of +1, followed by at most one 0, followed by sequence of -1
  - $-1 \rightarrow 0 \rightarrow +1$ corresponds to a valley and means an ISP is forwarding packets for free
    - Worse: it is paying its providers for it
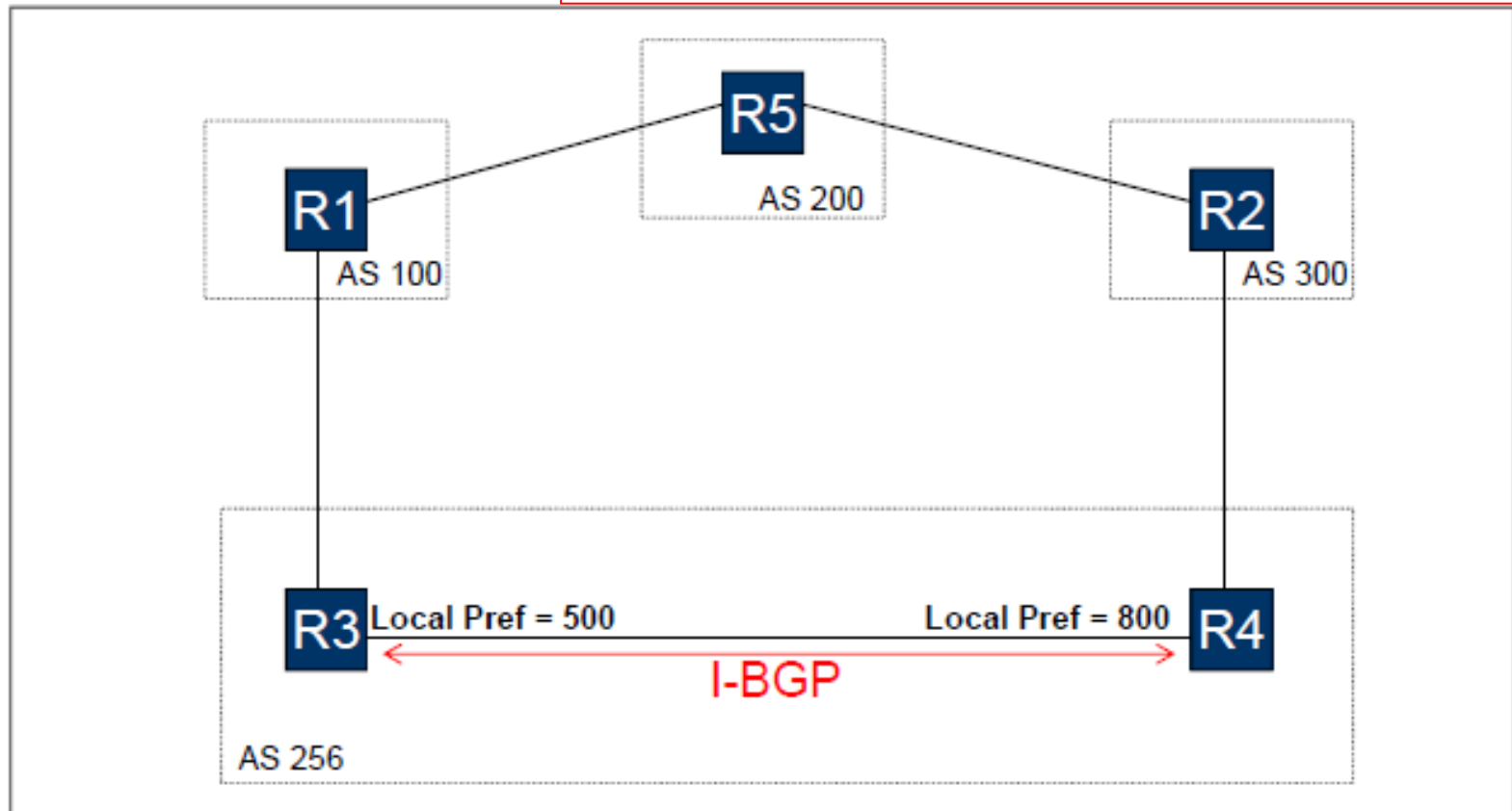
# Internet inter-AS routing: BGP

❖ 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"*

# LOCAL PREF

- Local (within an AS) mechanism to provide relative priority among BGP routers

The route with the highest local preference value is preferred

# LOCAL PREF – Common Uses

- Routers have a default LOCAL PREF
  - Can be changed for specific ASes

- Peering vs. transit
  - Prefer to use peering connection, why?

# LOCAL PREF – Common Uses

- Routers have a default LOCAL PREF
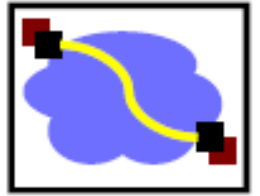  - Can be changed for specific ASes

Large ISPs have no advantage to peer with potential customers

- Peering vs. transit
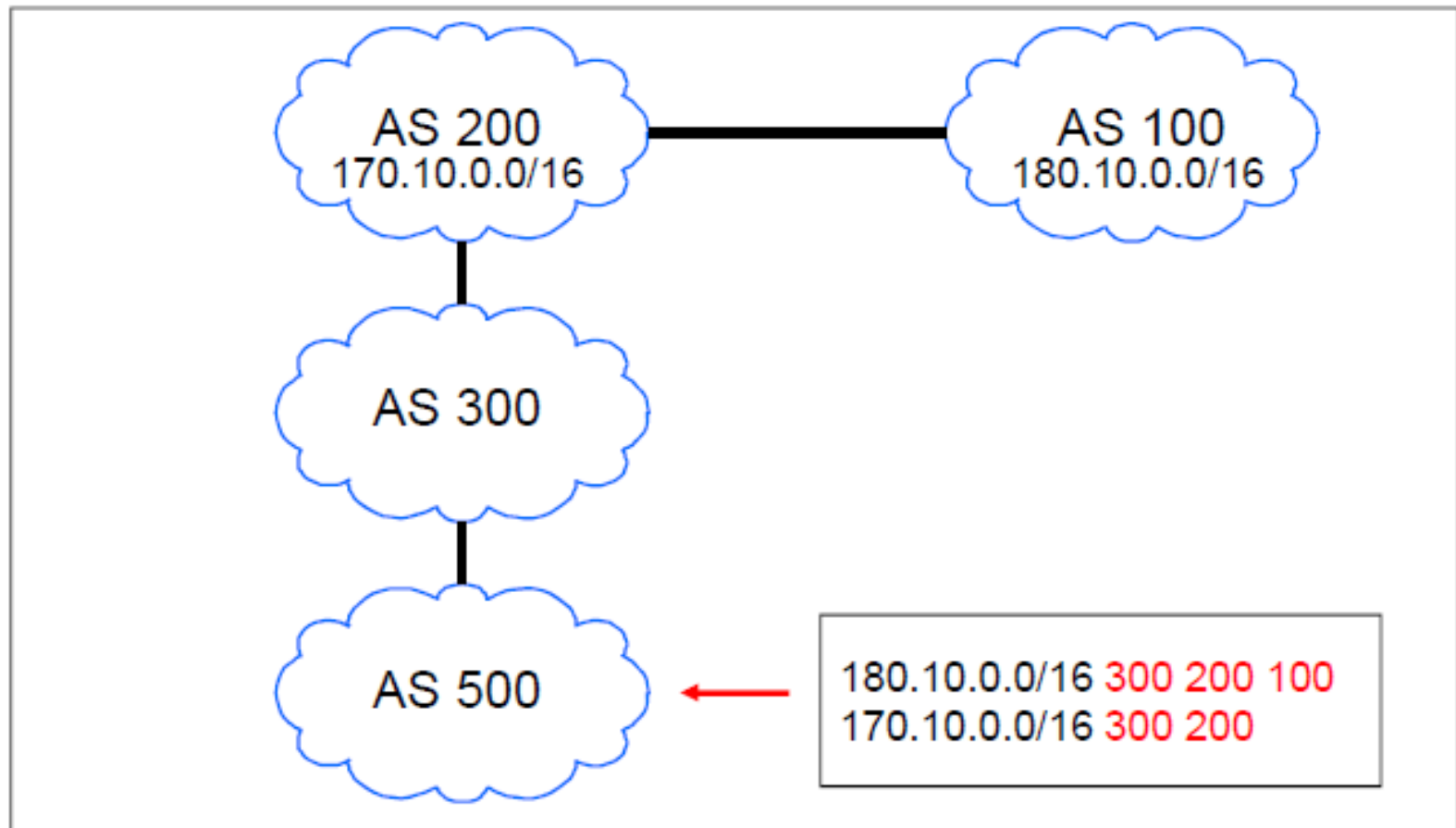  - Prefer to use peering connection, why?

neither ISP pays the other for the exchange of traffic

allowing traffic from another network to transit the provider's network. Smaller ISPs pay to the transit provider to connect to the rest of Internet

# AS_PATH

- List of traversed AS's



AS 200
170.10.0.0/16

AS 100
180.10.0.0/16

AS 300

AS 500

180.10.0.0/16 300 200 100
170.10.0.0/16 300 200

# Multi-Exit Discriminator (MED)

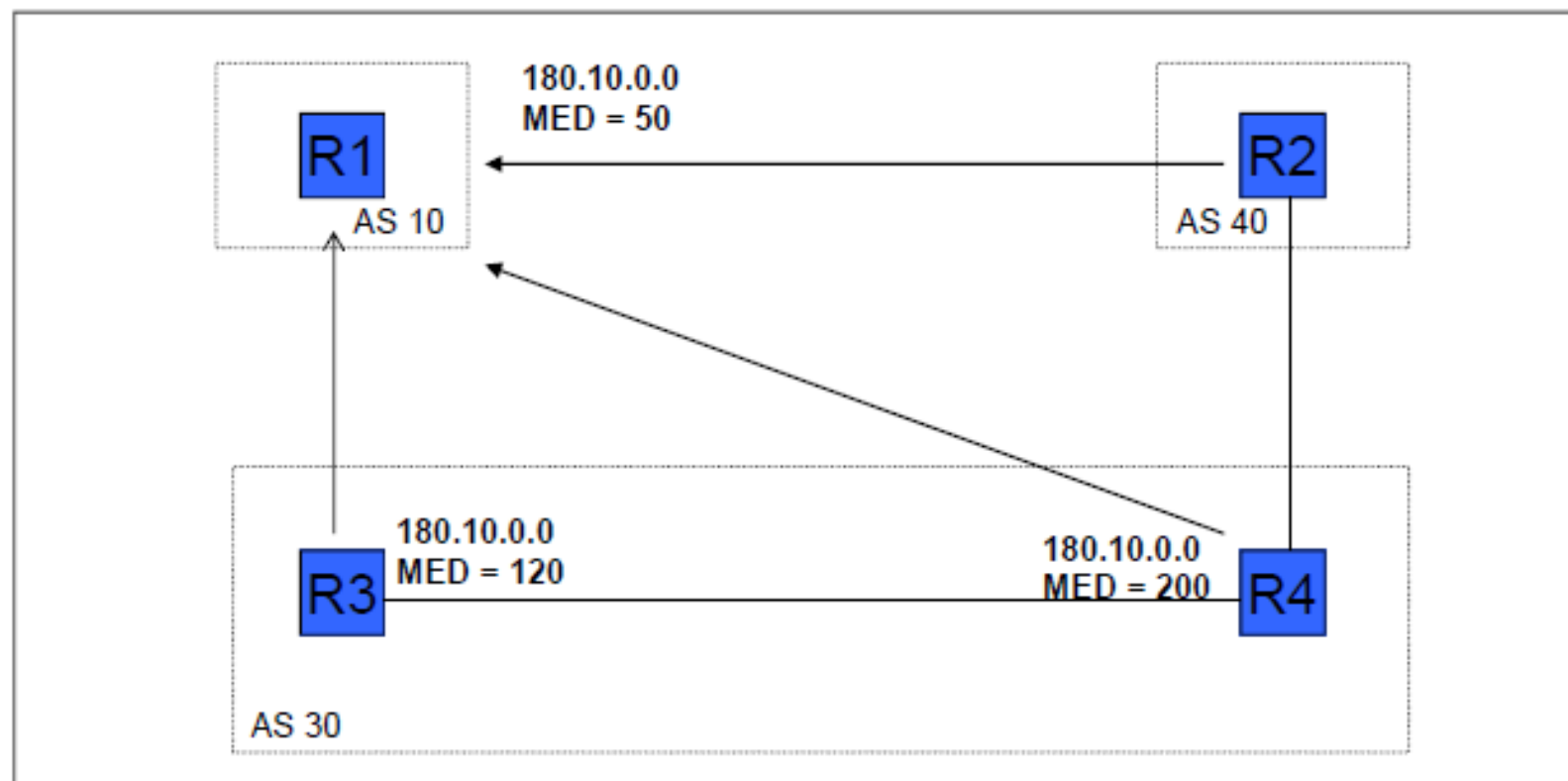- Hint to external neighbors about the preferred path into an AS

  **A lower MED value is preferred over a higher value**

- Used when two AS's connect to each other in more than one place
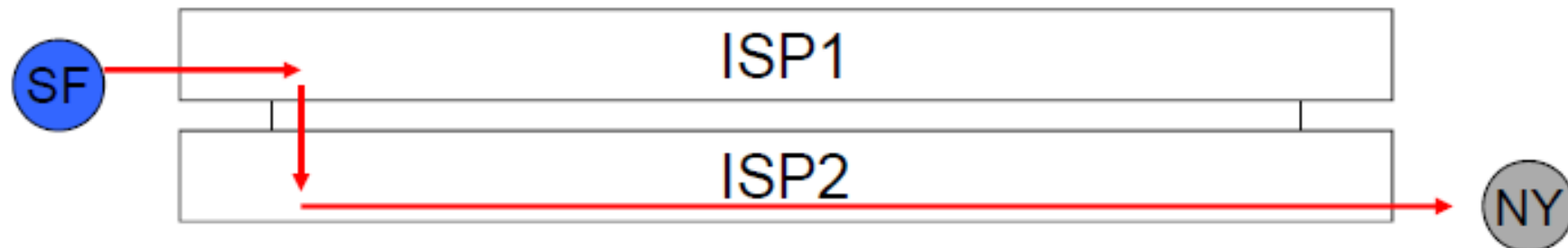
# MED

- Hint to R1 to use R3 over R4 link
- Cannot compare AS40's values to AS30's

# MED

- MED is typically used in provider/subscriber scenarios
- It can lead to unfairness if used between ISP because it may force one ISP to carry more traffic:
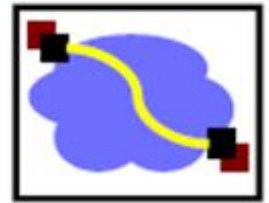


- ISP1 ignores MED from ISP2
- ISP2 obeys MED from ISP1
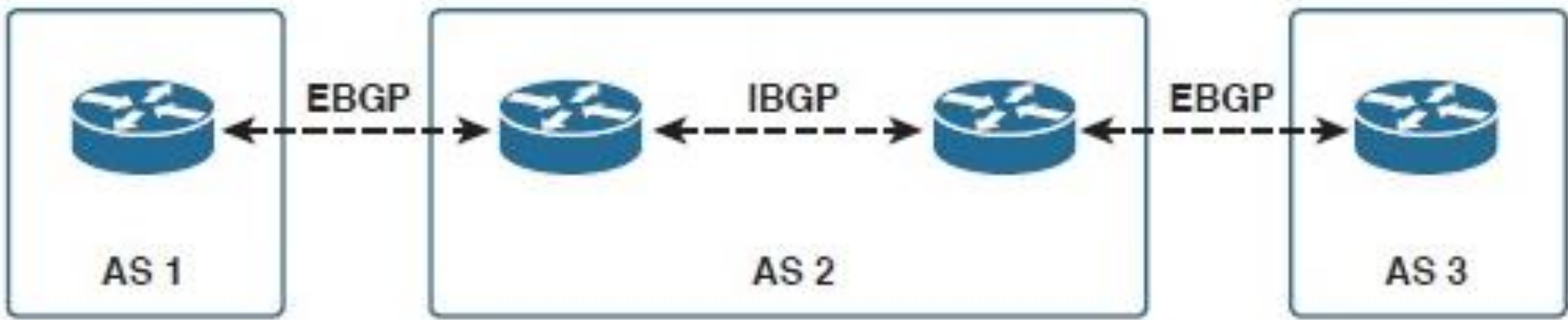- ISP2 ends up carrying traffic most of the way

# Path Selection Criteria

- Attributes + external (policy) information
- Rough ordering for path selection
  - Highest LOCAL-PREF
    - Captures business relationships and other factors
  - Shortest AS-PATH
  - Lowest MED (if routes learned from same neighbor)
  - eBGP over iBGP-learned
  - Lowest internal routing cost to border router
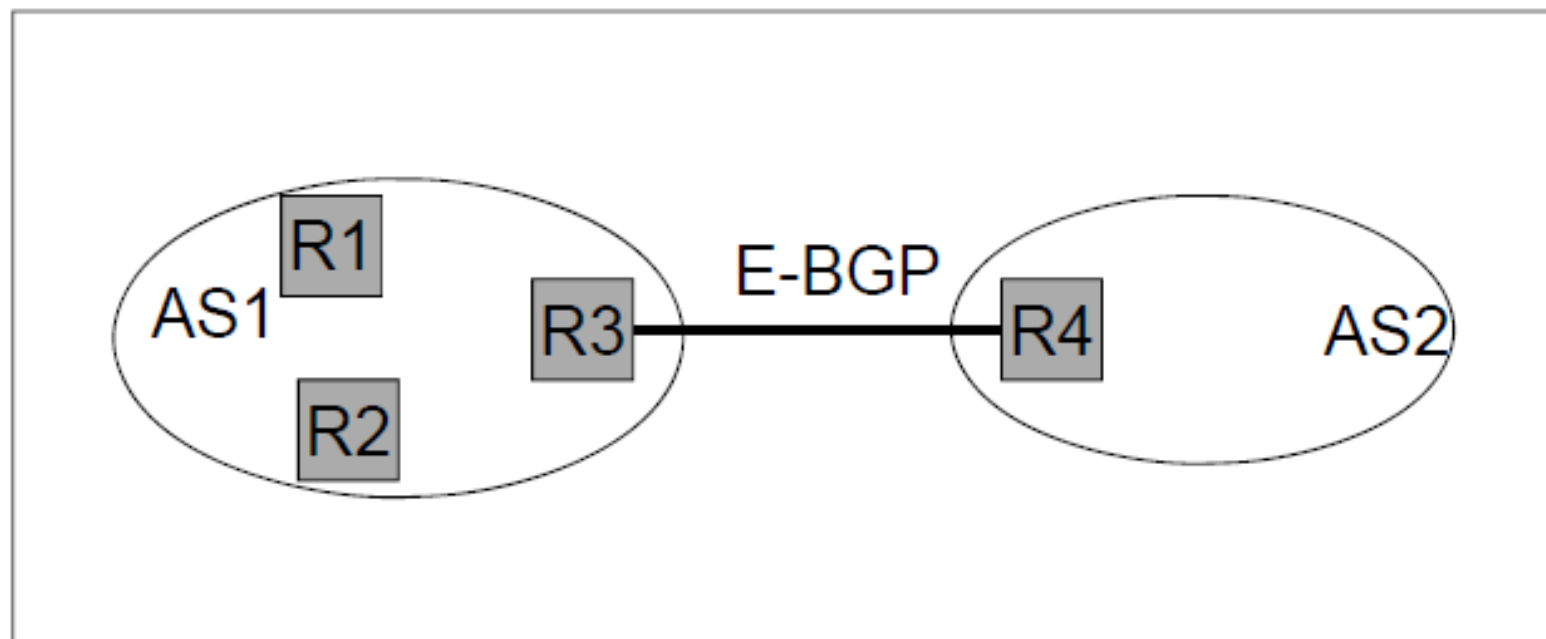  - Tie breaker, e.g., lowest router ID

**IBGP is necessary if BGP-advertised information must be passed within a given AS**

# Internal vs. External BGP

- BGP can be used by R3 and R4 to learn routes
- How do R1 and R2 learn routes?
- Border gateways also need to run an internal routing protcol
  - Establish connectivity between routers inside AS
- I-BGP: uses same messages as E-BGP
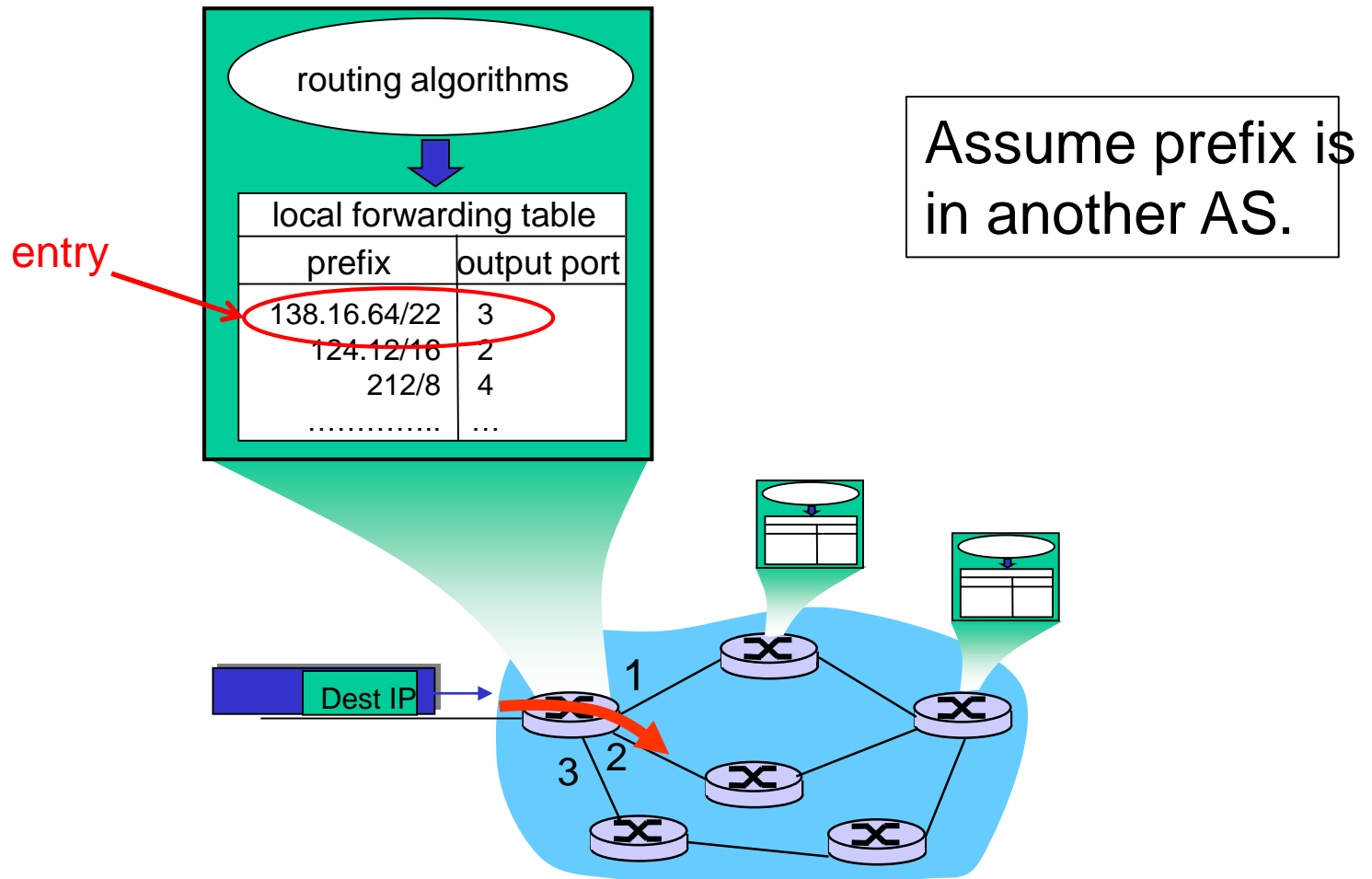
# Putting it Altogether:
## *How Does an Entry Get Into a Router's Forwarding Table?*

❖ Answer is complicated!

❖ Ties together hierarchical routing (Section 4.5.3) with BGP (4.6.3) and OSPF (4.6.2).

❖ Provides nice overview of BGP!

# How does entry get in forwarding table?



routing algorithms

local forwarding table

| prefix | output port |
|---|---|
| 138.16.64/22 | 3 |
| 124.12/16 | 2 |
| 212/8 | 4 |
| …………… | … |

entry

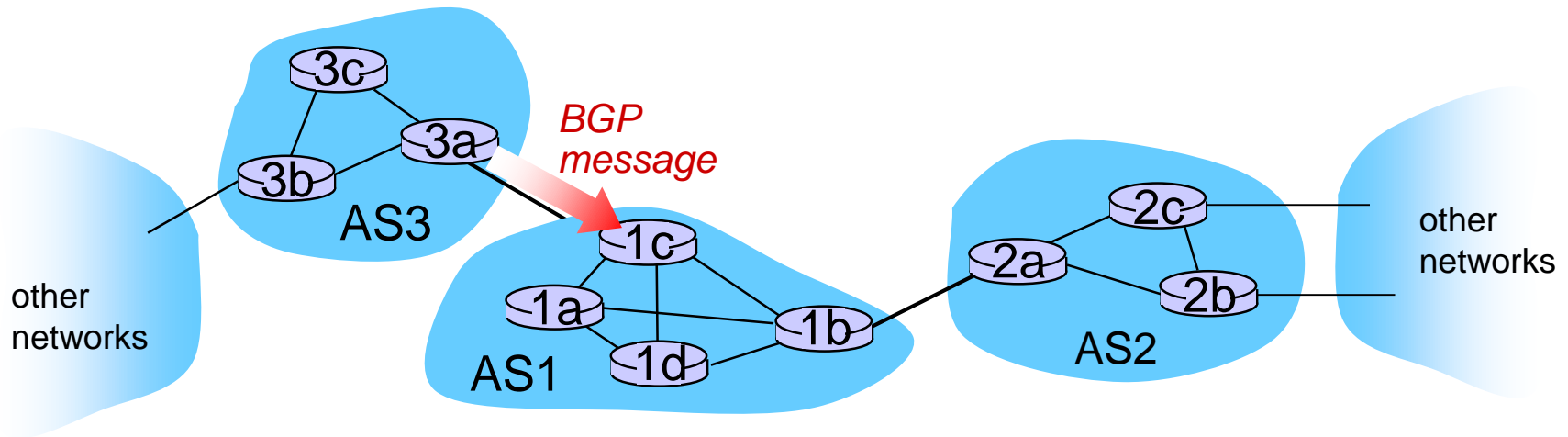Assume prefix is in another AS.

Dest IP

1
3 2

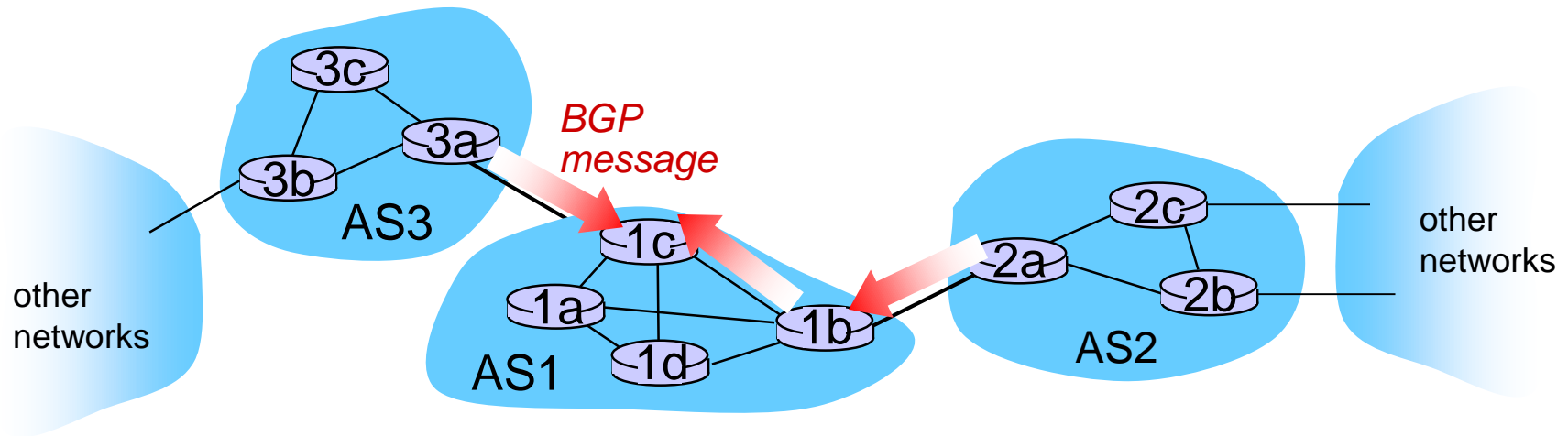# How does entry get in forwarding table?

## 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

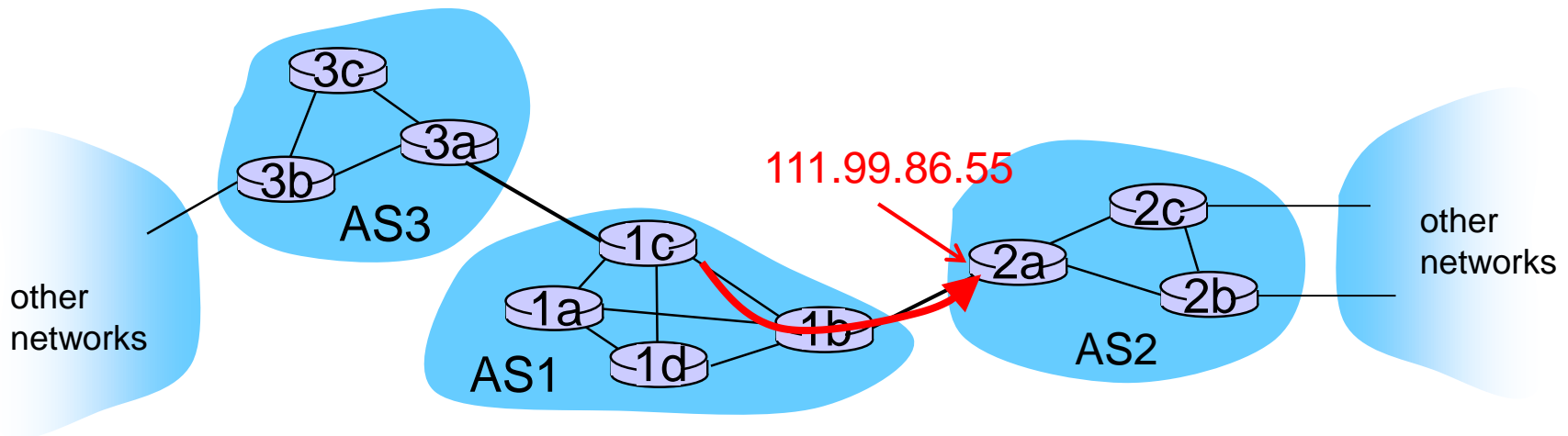❖ Router selects route based on shortest AS-PATH

❖ Example:

select

❖ AS2 AS17  to 138.16.64/22
❖ AS3 AS131 AS201 to 138.16.64/22

# 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

3c

3a

3b

AS3

111.99.86.55

2c

other networks

1c

2a

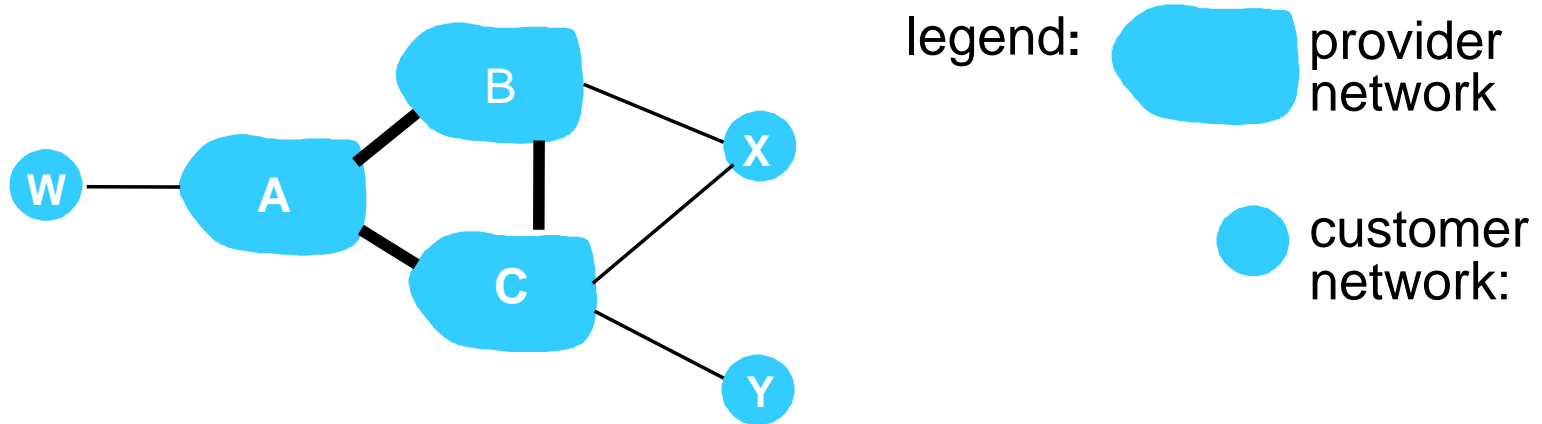other networks

1a

2b

1b

AS2

1d

AS1

# How does entry get in forwarding table?
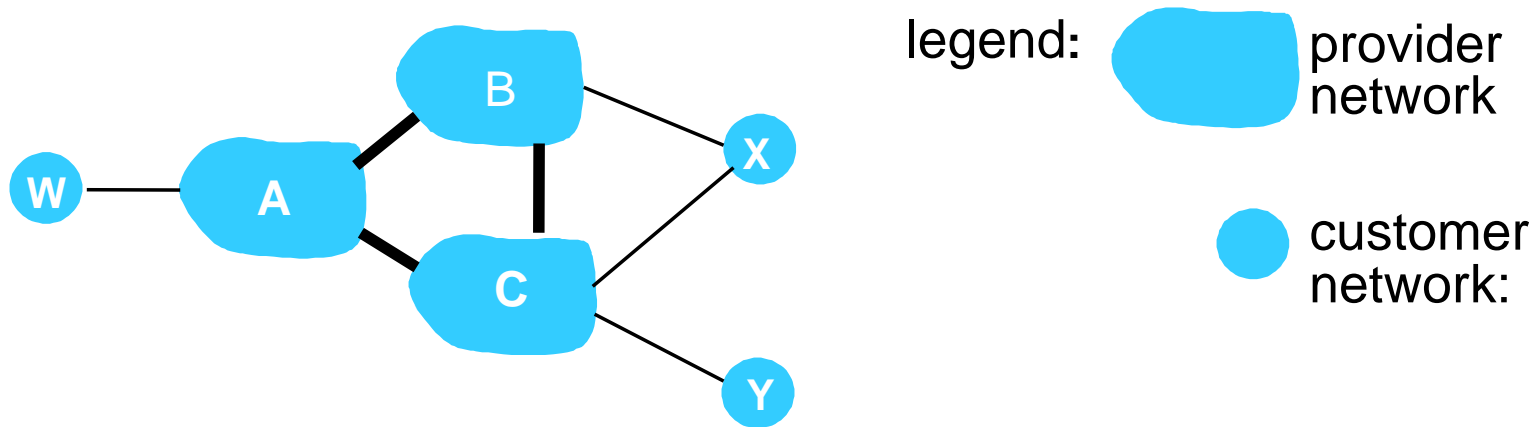
## Summary

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

# Broadcast routing

❖ deliver packets from source to all other nodes

❖ source duplication is inefficient:
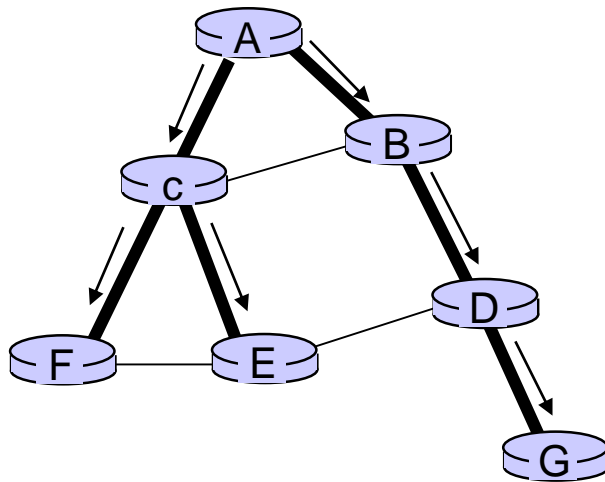


source duplication

in-network duplication

❖ source duplication: how does source determine recipient addresses?
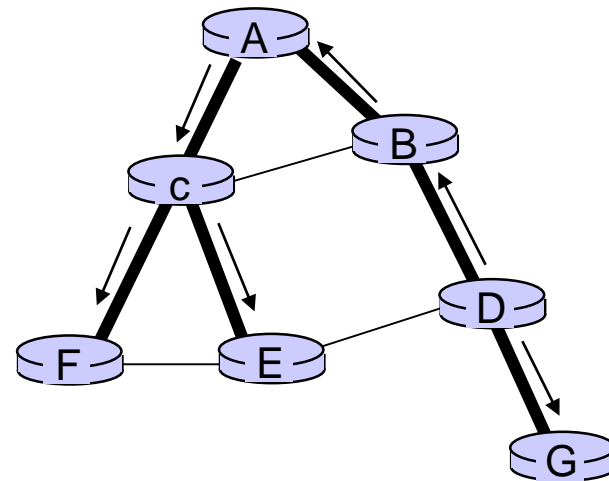
# In-network duplication

❖ *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
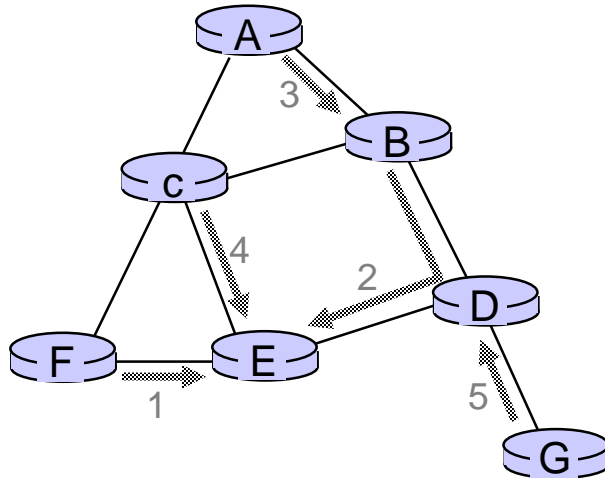❖ nodes then forward/make copies only along spanning tree
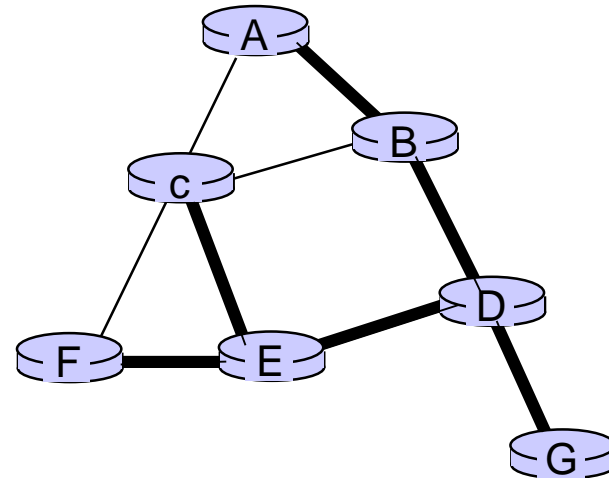


(a) broadcast initiated at A

(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
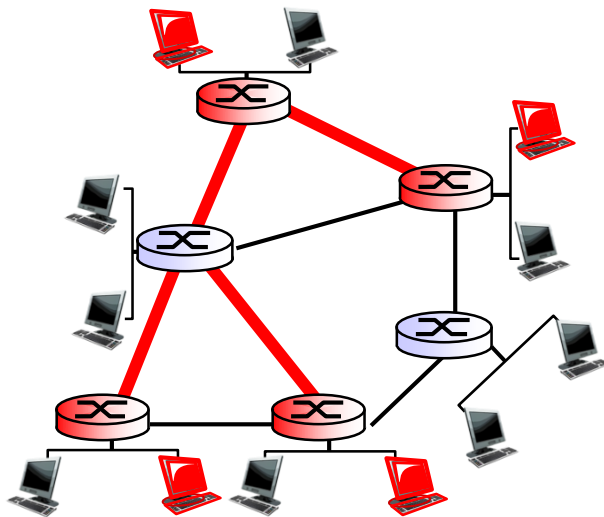- ❖ *source-based:* different tree from each sender to rcvrs
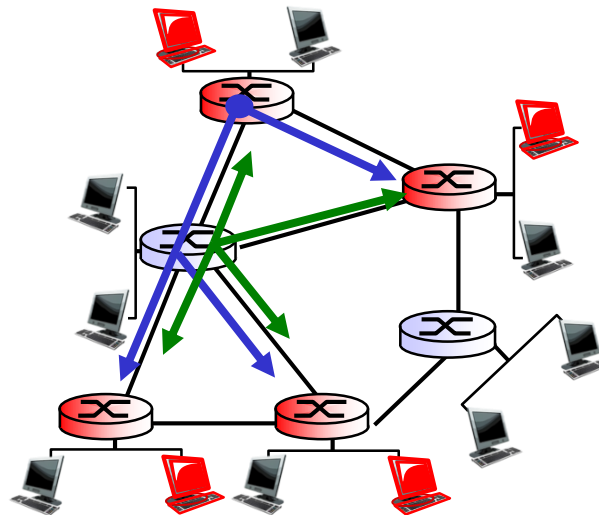
legend

group member

not group member

router with a group member

router without group member

shared tree

source-based trees
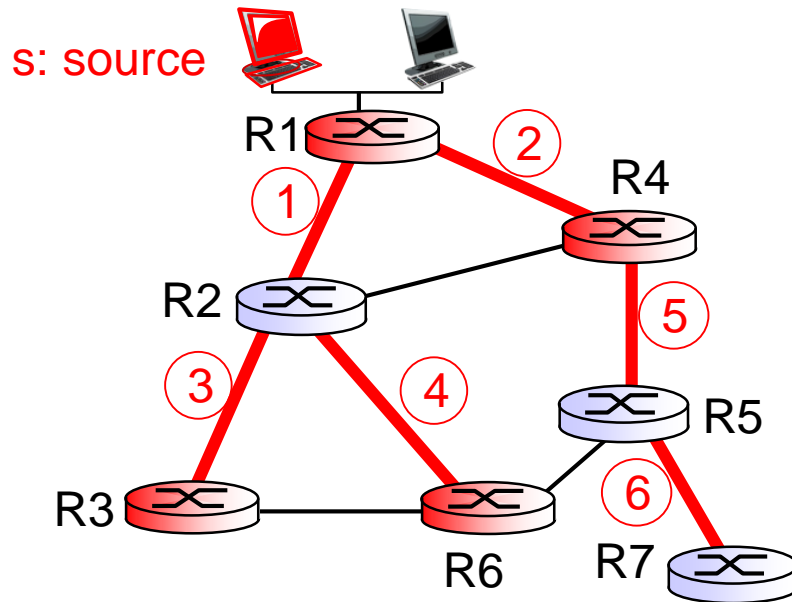
# Approaches for building mcast trees

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 (one) 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
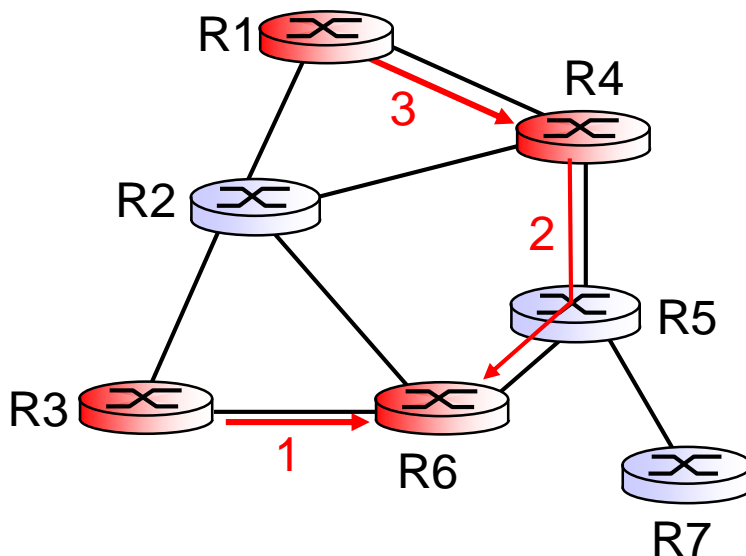
# Shared-tree: steiner tree

- ❖ *One (global) mcast tree*
- ❖ *steiner tree:* minimum cost tree connecting all routers with attached group members
- ❖ problem is NP-complete
- ❖ excellent heuristics exist
- ❖ 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
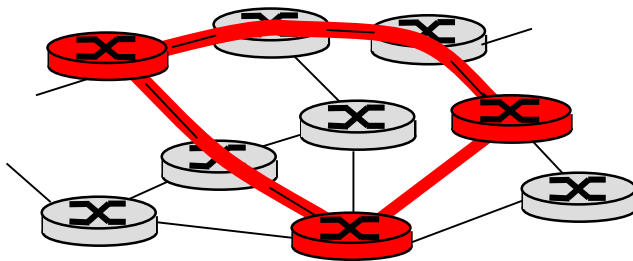
suppose R6 chosen as center:

LEGEND
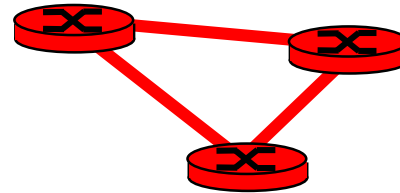


router with attached group member

router with no attached group member

1 → path order in which join messages generated

# 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: *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