# Chapter 4: Network Layer

# Hierarchical Routing 階層式繞徑

Our routing study thus far - idealization

- □ all routers identical
- □ network "flat"

... *not* true in practice

scale: with 200 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 網路自治
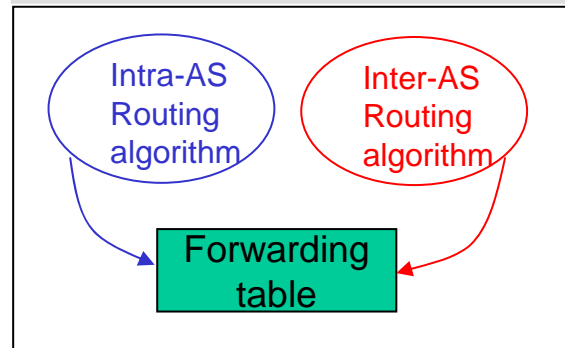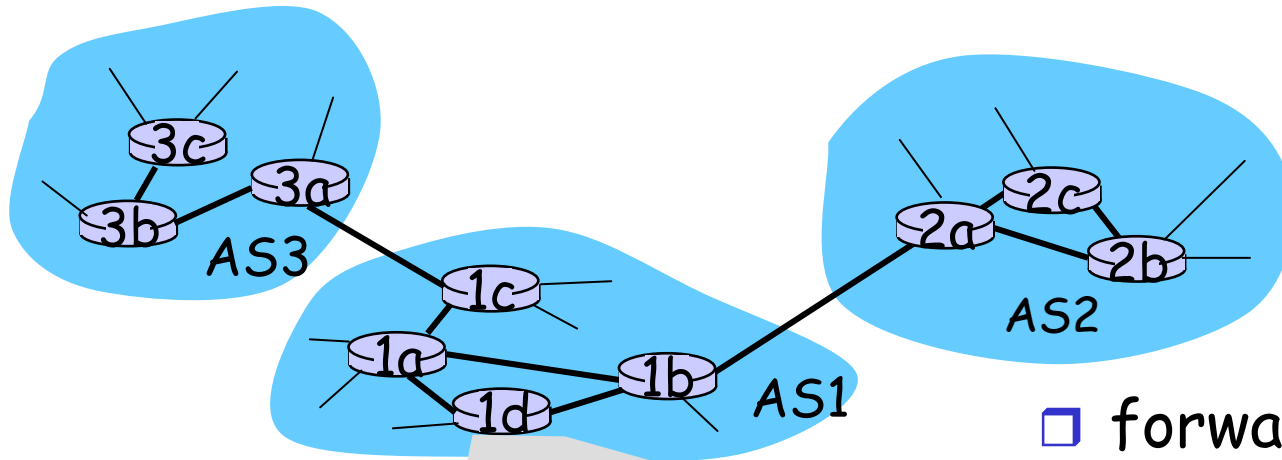
# Hierarchical Routing 階層式繞徑

□ aggregate routers into regions, "autonomous systems" (AS)
自治系統

□ 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

匣道路由器

□ Direct link to router in another AS

# Interconnected ASes



- □ forwarding table configured by both intra- and inter-AS routing algorithm
  - ○ intra-AS sets entries for internal dests
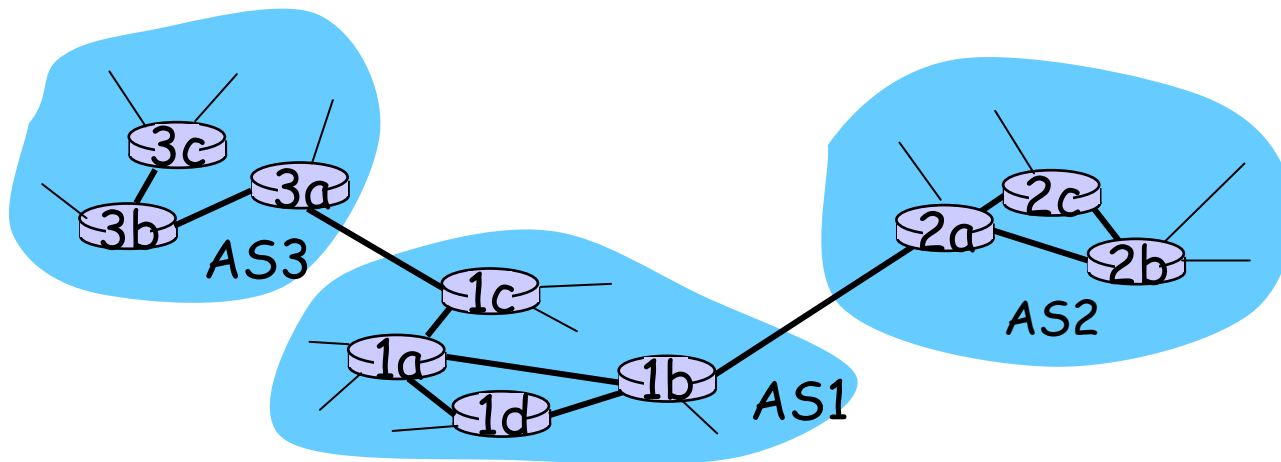  - ○ inter-AS & Intra-As sets entries for external dests

# Inter-AS tasks AS間的繞徑

□ suppose router in AS1 receives datagram dest outside of AS1

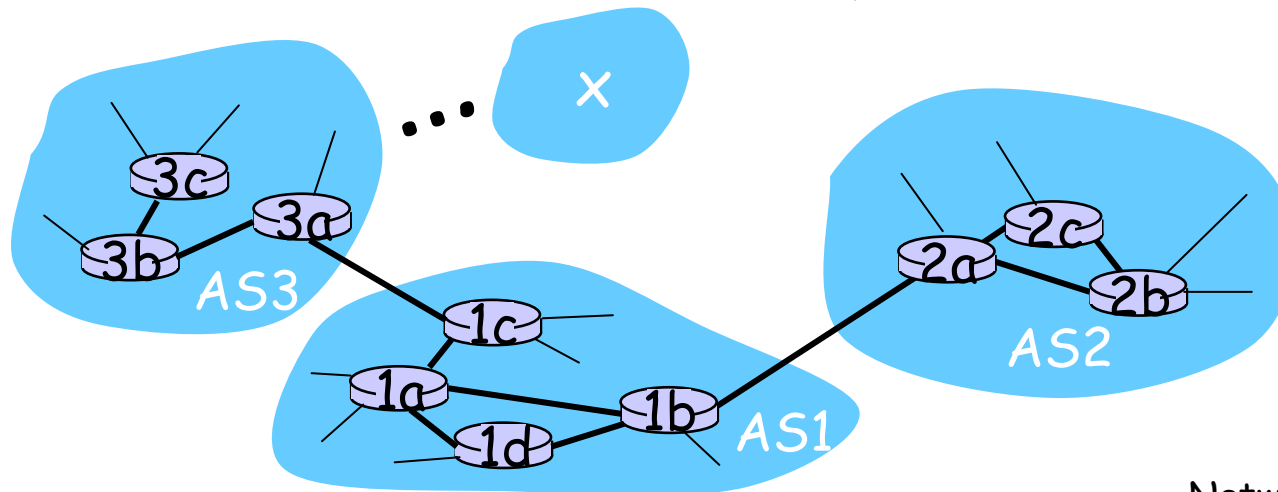○ router should forward packet to gateway router, but which one?

1. learn which dests 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
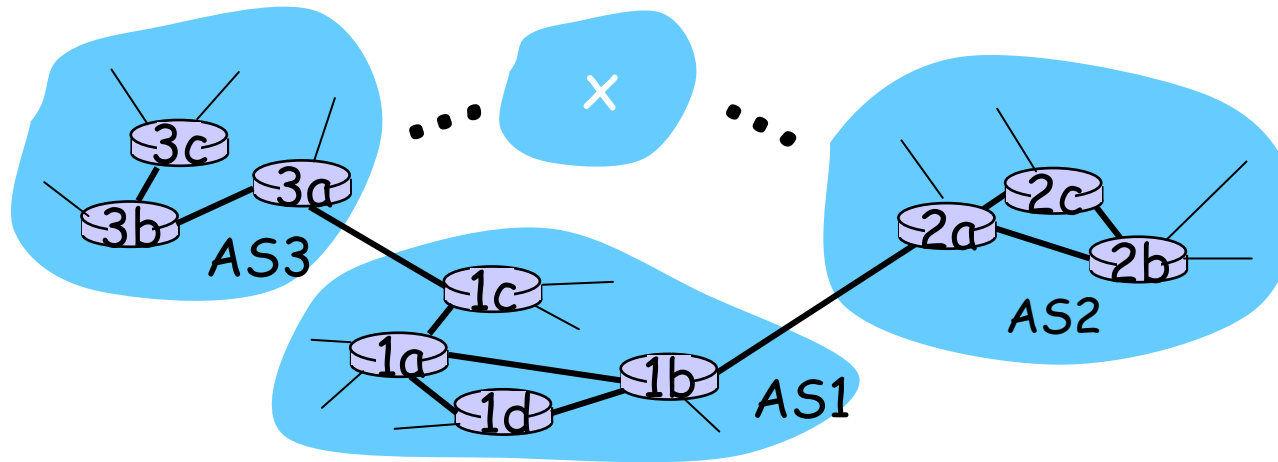## 設定Forwarding Table

- 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 $I$ is on the least cost path to 1c.

    - installs forwarding table entry *(x,I)*
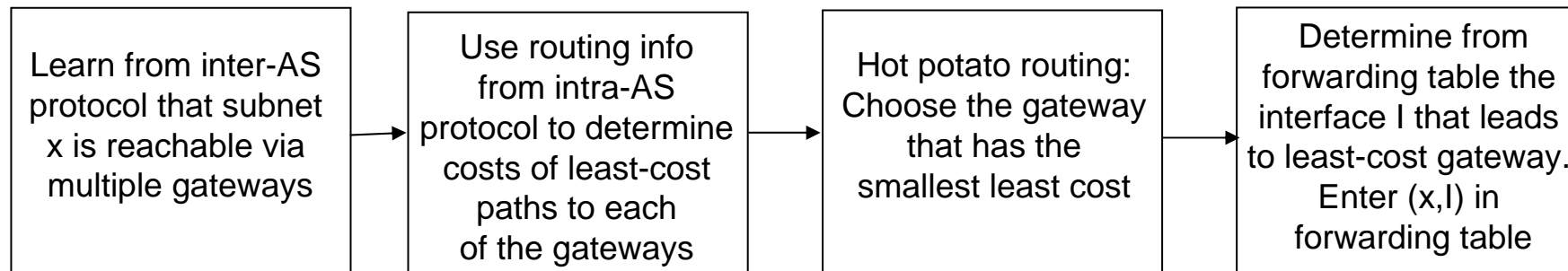
# Example: Choosing among multiple ASes

## 選擇路徑

- 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!

# 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 燙手山芋繞徑演算法: send packet towards closest of two routers. 選擇在本身AS中，較靠近的router

| Learn from inter-AS protocol that subnet x is reachable via multiple gateways | → | Use routing info from intra-AS protocol to determine costs of least-cost paths to each of the gateways | → | Hot potato routing: Choose the gateway that has the smallest least cost | → | Determine from forwarding table the interface I that leads to least-cost gateway. Enter (x,I) in forwarding table |

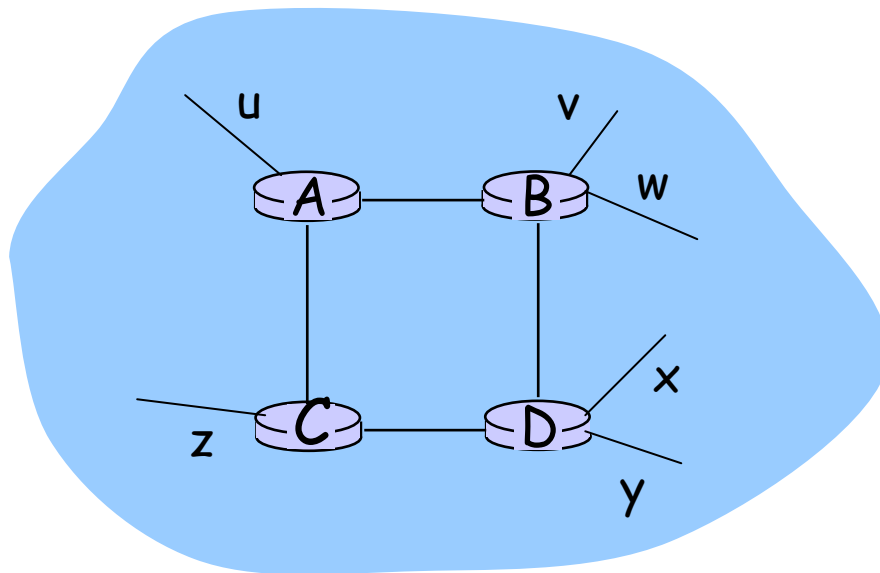# Chapter 4: Network Layer

# Intra-AS Routing 在自治系統內部繞送

□ also known as Interior Gateway Protocols (IGP)
  內部匣道協定

□ most common Intra-AS routing protocols:

  ○ RIP: Routing Information Protocol
    繞送資訊協定

  ○ OSPF: Open Shortest Path First
    最短開放路徑優先協定

  ○ IGRP: Interior Gateway Routing Protocol (Cisco proprietary)

# Chapter 4: Network Layer

# RIP ( Routing Information Protocol)

- distance vector algorithm 距離向量演算法
- included in BSD-UNIX Distribution in 1982
- distance metric: # of hops (max = 15 hops)



From router A to subsets:

| destination | hops |
|---|---|
| u | 1 |
| v | 2 |
| w | 2 |
| x | 3 |
| y | 3 |
| z | 2 |

# RIP advertisements (RIP 通告)

- *distance vectors:* exchanged among neighbors every 30 sec via Response Message (also called advertisement) 每三十秒和相鄰節點更新訊息
- each advertisement: list of up to 25 destination nets within AS

# RIP: Example



| Destination Network | Next Router | Num. of hops to dest. |
|---|---|---|
| w | A | 2 |
| y | B | 2 |
| z | B | 7 |
| x | -- | 1 |
| …. | …. | …. |

Routing table in D

# RIP: Example

Dest | Next | hops
--- | --- | ---
w | - | 1
x | - | 1
z | C | 4
.... | ... | ...

Advertisement from A to D



| Destination Network | Next Router | Num. of hops to dest. |
| --- | --- | --- |
| w | A | 2 |
| y | B | 2 |
| z | ~~B~~ A | ~~7~~ 5 |
| x | -- | 1 |
| .... | .... | .... |

Routing table in D

# RIP: Link Failure and Recovery

If no advertisement heard after 180 sec --> neighbor/link declared dead
180秒後沒有RIP 通告，表示連結失效

- ○ routes via neighbor invalidated
- ○ new advertisements sent to neighbors
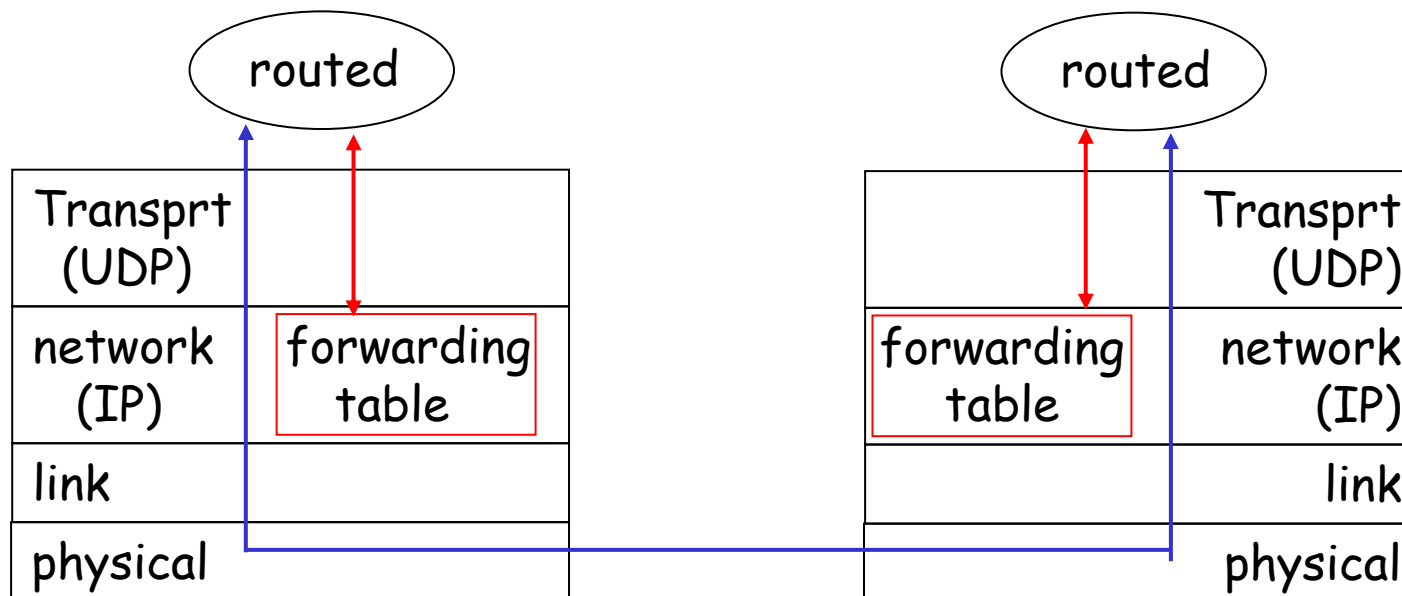- ○ neighbors in turn send out new advertisements (if tables changed)
- ○ link failure info quickly (?) propagates to entire net
- ○ *poison reverse* used to prevent ping-pong loops (infinite distance = 16 hops)

# RIP Table processing

□ RIP routing tables managed by **application-level** process called route-d (daemon) 由應用層處理

□ advertisements sent in UDP packets, periodically repeated

routed                                    routed

| Transprt (UDP) |                 |        |                 | Transprt (UDP) |
|---|---|---|---|---|
| network (IP) | forwarding table |        | forwarding table | network (IP) |
| link |                 |        |                 | link |
| physical |                 |        |                 | physical |

# Chapter 4: Network Layer

# OSPF (Open Shortest Path First)

□ "open": publicly available  公開可使用的

□ uses Link State algorithm  使用連結狀態演算法
  ○ LS packet dissemination  散布LS封包
  ○ topology map at each node  每個節點都知道網路狀態
  ○ route computation using Dijkstra's algorithm 計算最短路徑

□ OSPF advertisement carries one entry per neighbor router

□ advertisements disseminated to entire AS (via flooding)
  ○ carried in OSPF messages directly over IP (rather than TCP or UDP

# OSPF "advanced" features (not in RIP)

- Security 安全性: all OSPF messages authenticated (to prevent malicious intrusion) 必需自己加密

- multiple same-cost paths allowed (only one path in RIP) 多條同成本路徑

- For each link, multiple cost metrics for different TOS (e.g., satellite link cost set "low" for best effort; high for real time)

- integrated uni- and multicast support: 單播與群播的支援
  - Multicast OSPF (MOSPF) uses same topology data base as OSPF

- hierarchical OSPF in large domains. 支援階層架構

# Hierarchical OSPF 階層式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:* 網域邊境路由器 "summarize" distances to nets in own area, advertise to other Area Border routers.

□ *backbone routers:* 主幹路由器 run OSPF routing limited to backbone.

□ *boundary routers:* 邊境路由器 connect to other AS's.

# Chapter 4: Network Layer

□ 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
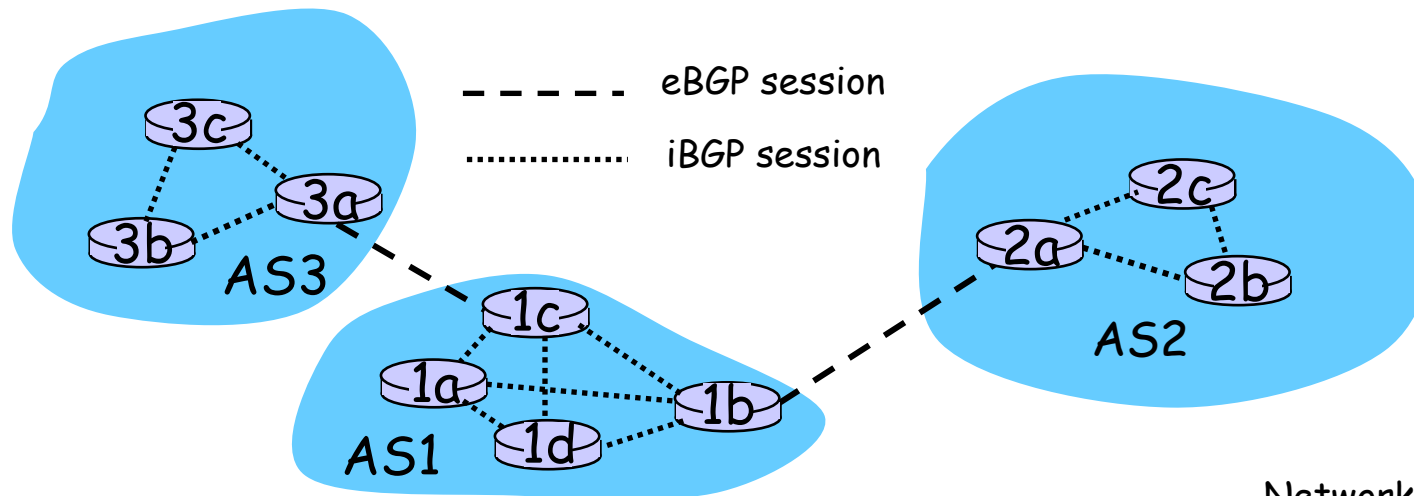
# Internet inter-AS routing: BGP
## 邊境匣道協定

□ **BGP (Border Gateway Protocol):** *the* de facto standard

□ BGP provides each AS a means to:

1. Obtain subnet reachability information from neighboring ASs. 從相鄰AS取得子網路的連通資訊

2. Propagate reachability information to all AS-internal routers. 傳播連通資訊給所有AS內部的路由器

3. Determine "good" routes to subnets based on reachability information and policy. 根據上列資訊及策略，判斷前往各子網路的"好"路徑

□ allows subnet to advertise its existence to rest of Internet: *"I am here"*
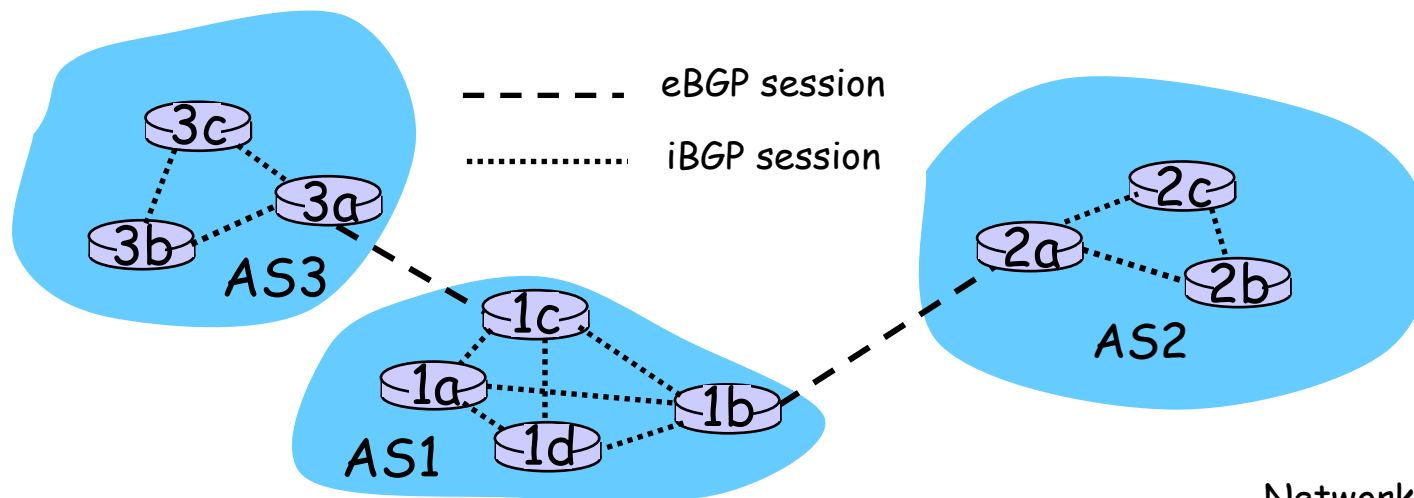
# BGP basics

- 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 prefix to AS1:
  - AS2 *promises* it will forward any addresses datagrams towards that prefix.
  - AS2 can aggregate prefixes 前置碼 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 do 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, creates entry for prefix in its forwarding table.

# Path attributes & BGP routes
## 路徑屬性與BGP繞送

□ 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)

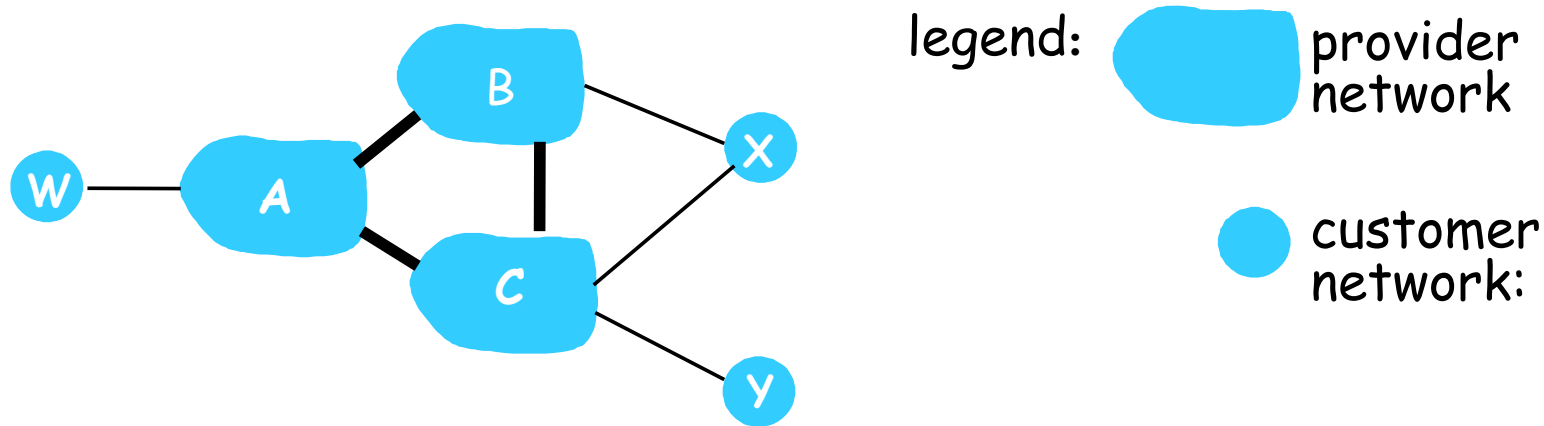□ when gateway router receives route advertisement, uses import policy to accept/decline.

# BGP route selection (BGP繞送選擇)

□   router may learn about more than 1 route to some prefix. Router must select route.

□   elimination rules:

1.  local preference value attribute: policy decision

2.  shortest AS-PATH

3.  closest NEXT-HOP router: hot potato routing
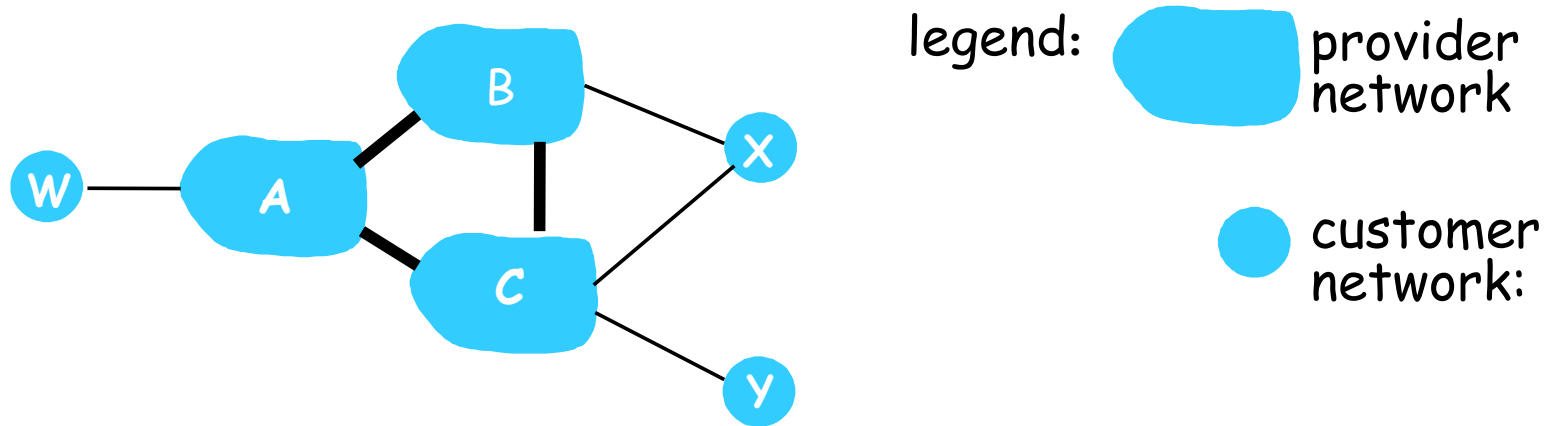
4.  additional criteria

# 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 繞送策略



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- and 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
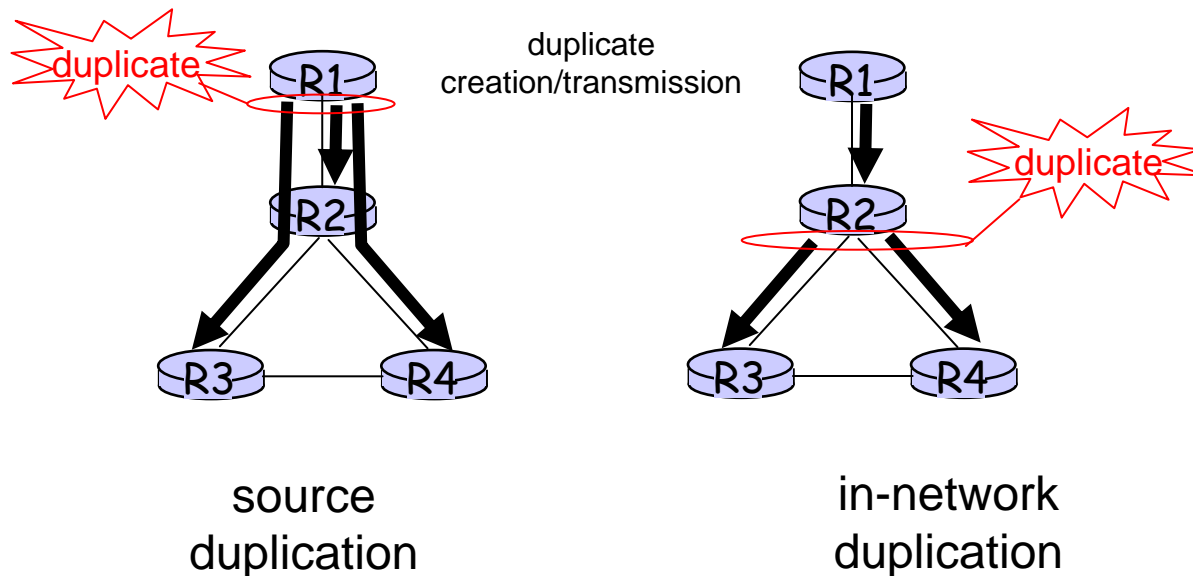- Inter-AS: policy may dominate over performance

# Chapter 4: Network Layer

# Broadcast Routing 廣播

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



source
duplication

in-network
duplication
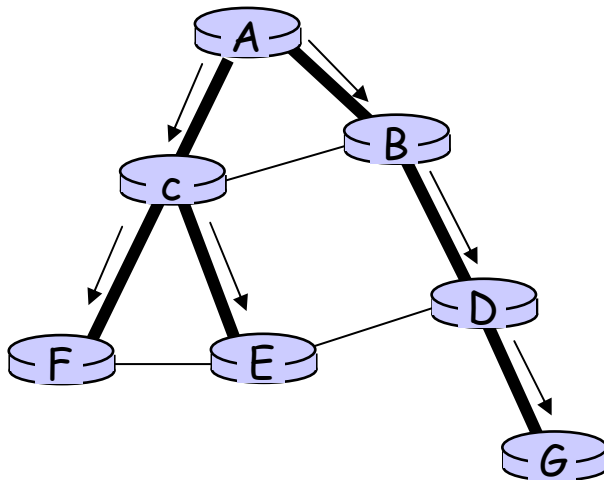
☐ source duplication: how does source
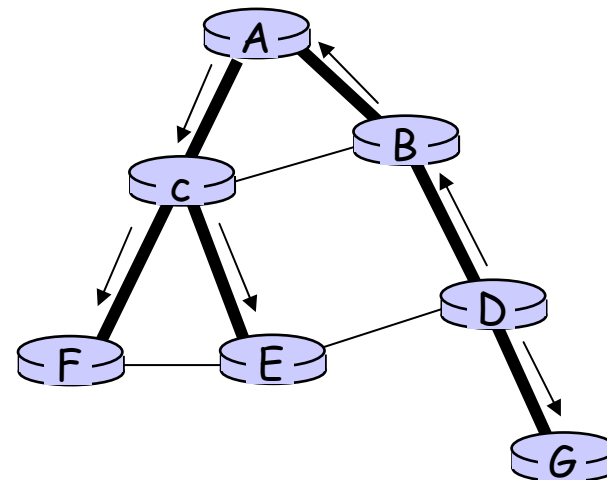determine recipient addresses?

# In-network duplication

□ flooding: when node receives brdcst pckt, sends copy to all neighbors
  ○ Problems: cycles & broadcast storm

□ controlled flooding 受控制的溢出: node only brdcsts pkt if it hasn't brdcst same packet before
  ○ Node keeps track of pckt ids already brdcsted
  ○ Or reverse path forwarding (RPF): only forward pckt if it arrived on shortest path between node and source

□ spanning tree 展開樹
  ○ No redundant packets received by any node

# Spanning Tree 展開樹

- First construct a spanning tree
- Nodes forward copies only along spanning tree
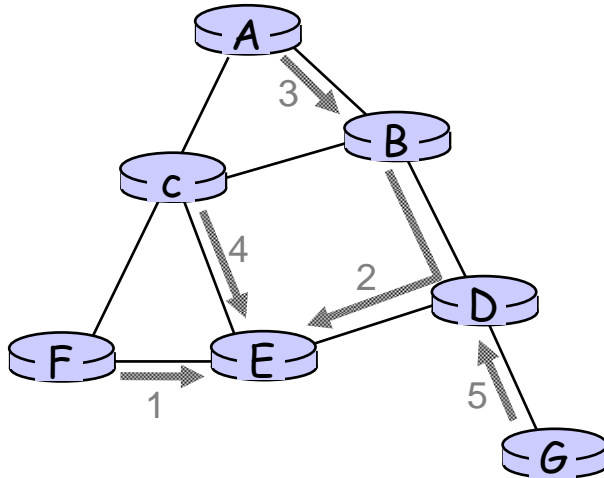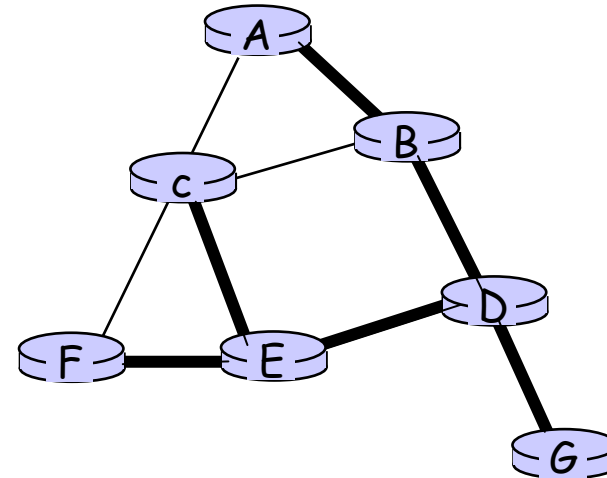


(a) Broadcast initiated at A          (b) Broadcast initiated at D

# 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**
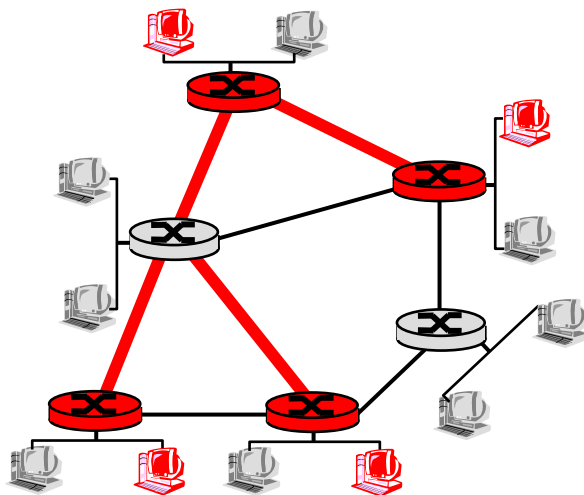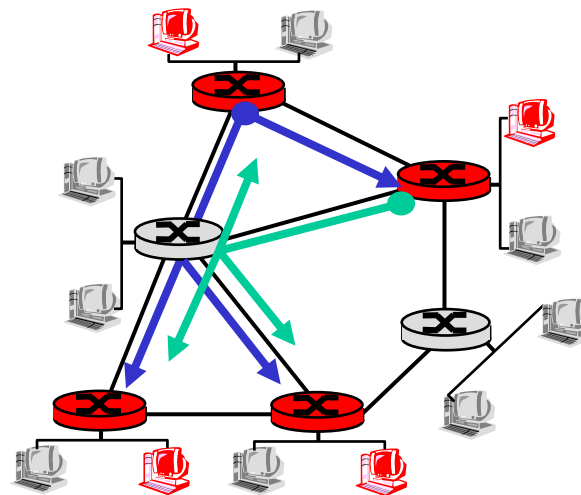
**(b) Constructed spanning tree**

# Multicast Routing: Problem Statement

## 群播繞徑

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

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



Shared tree            Source-based trees

# Approaches for building mcast trees
## 群播樹的建立

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 最短路徑樹

☐ 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

R1 · 1 · 2 · R4 · R2 · 3 · 4 · 5 · R5 · R3 · 6 · R6 · R7

# Reverse Path Forwarding
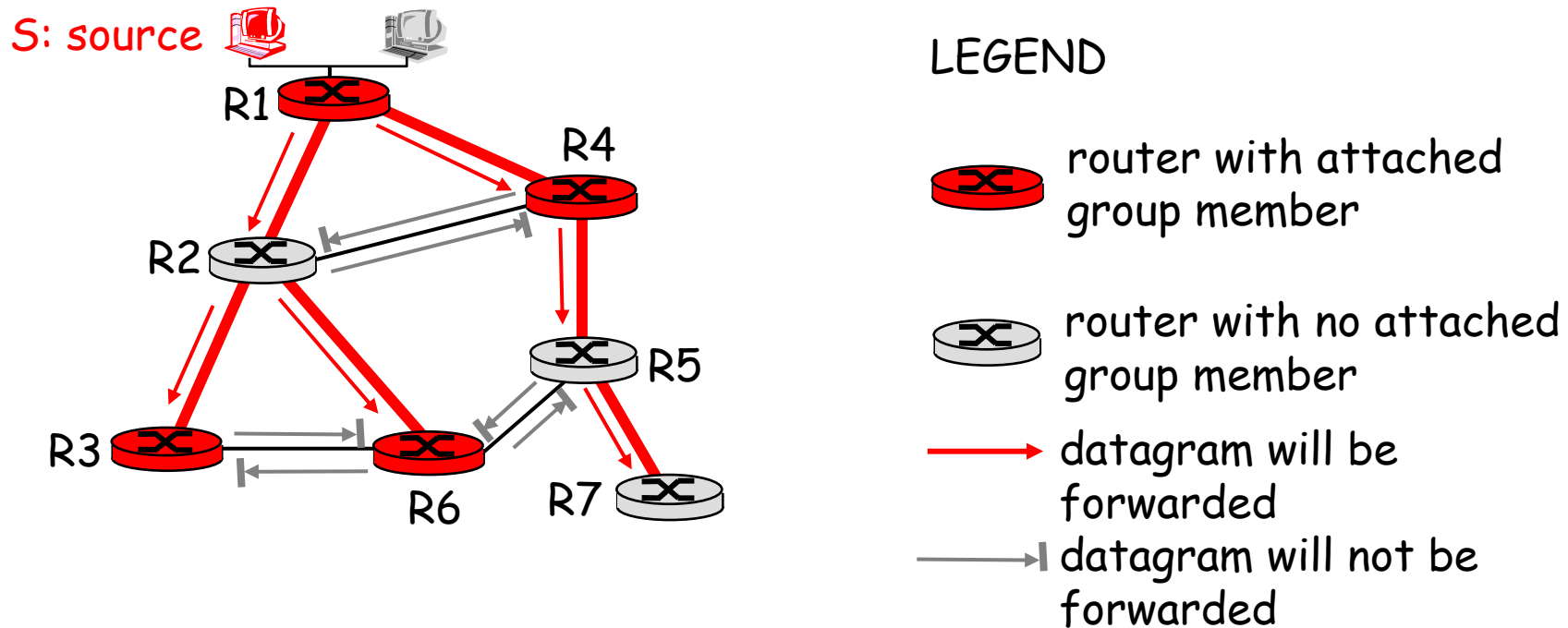## 反向路徑轉送

- ❑ 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 forwarded

datagram will not be forwarded
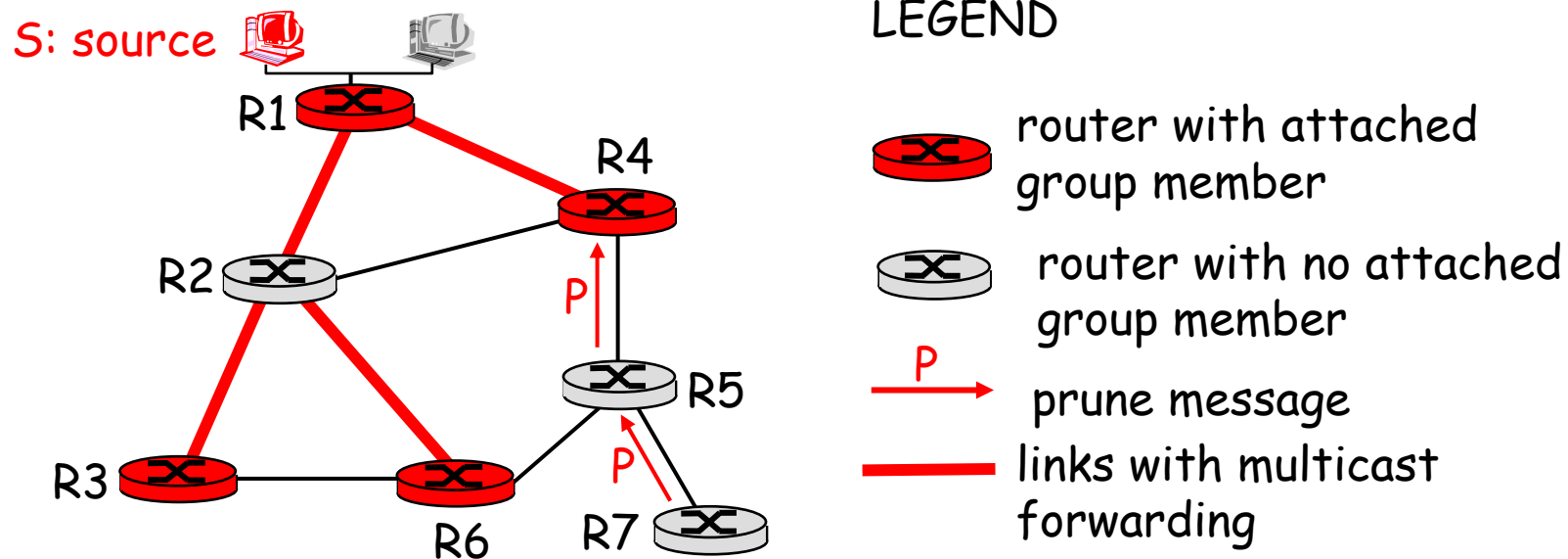
- result is a source-specific *reverse* SPT
  - may be a bad choice with asymmetric links

# 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

R1

R4

R2

P

R5

R3

P

R6    R7

LEGEND

router with attached group member

router with no attached group member

P →  prune message
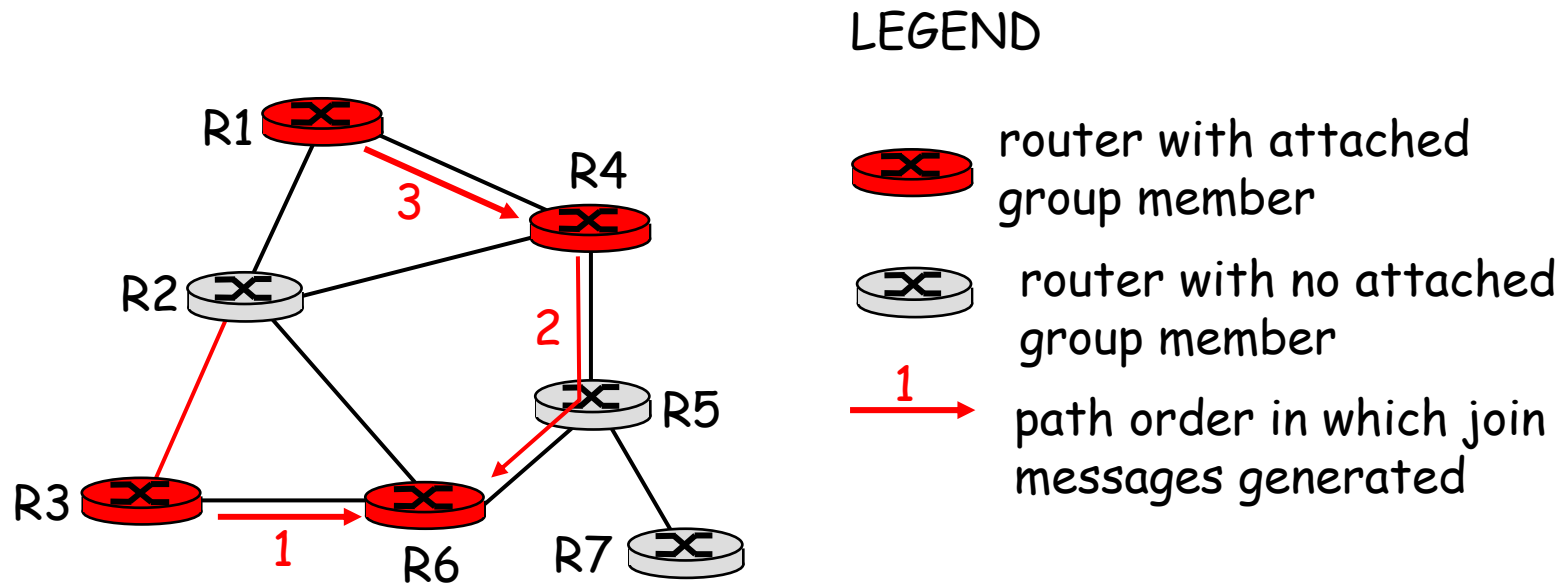
—— links with multicast forwarding

# 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: an 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

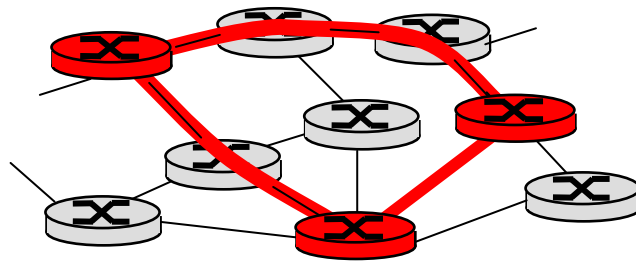  - routers not wanting group: send upstream prune msgs

# 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 routers
  ○ Mbone routing done using DVMRP
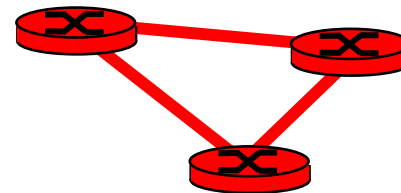
# Tunneling

**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
- ❑ receiving mcast router unencapsulates to get mcast datagram

# PIM: Protocol Independent Multicast

❑ not dependent on any specific underlying unicast routing algorithm (works with all)

❑ two different multicast distribution scenarios :

*Dense*:

❑ group members densely packed, in "close" proximity.

❑ bandwidth more plentiful

*Sparse:*

❑ # networks with group members small wrt # interconnected networks

❑ group members "widely dispersed"

❑ bandwidth not plentiful

# Consequences of Sparse-Dense Dichotomy:

## Dense

- group membership by routers *assumed* until routers explicitly prune
- *data-driven* construction on mcast tree (e.g., RPF)
- bandwidth and non-group-router processing *profligate*

## Sparse:

- no membership until routers explicitly join
- *receiver- driven* construction of mcast tree (e.g., center-based)
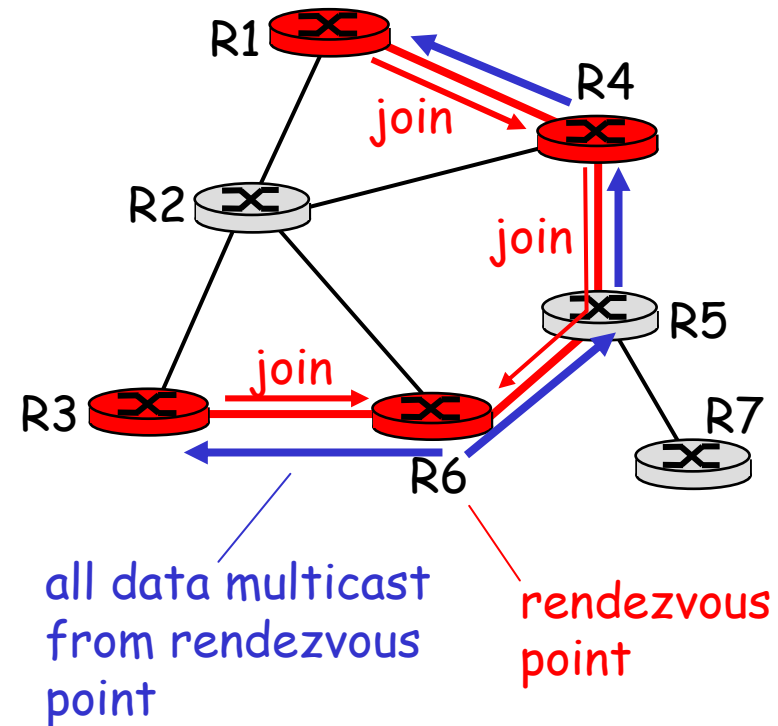- bandwidth and non-group-router processing *conservative*

# PIM- Dense Mode

flood-and-prune RPF, similar to DVMRP but

❑ underlying unicast protocol provides RPF info for incoming datagram

❑ less complicated (less efficient) downstream flood than DVMRP reduces reliance on underlying routing algorithm

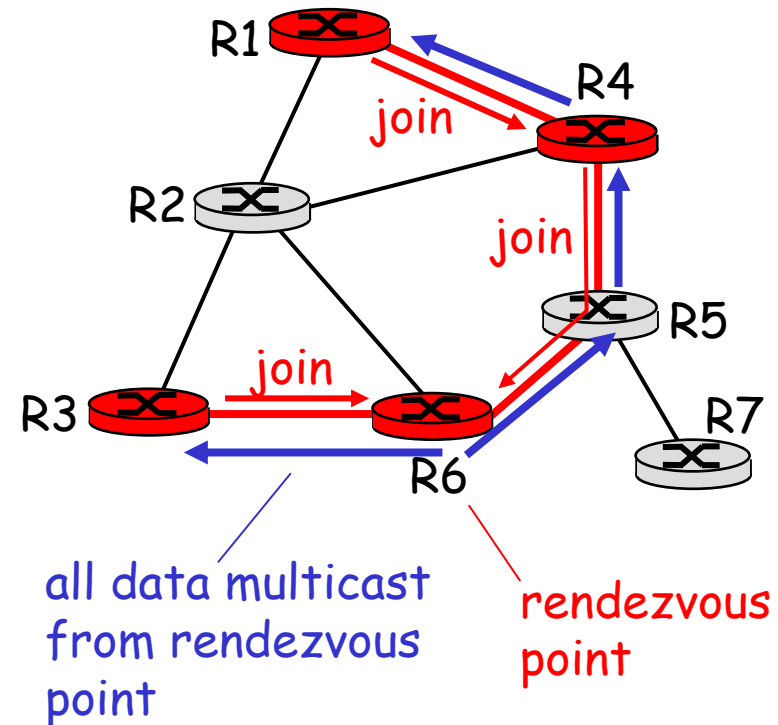❑ has protocol mechanism for router to detect it is a leaf-node router

# PIM - Sparse Mode

□ center-based approach

□ router sends *join* msg
  to rendezvous point
  (RP)

  ▫ intermediate routers
    update state and
    forward *join*

□ after joining via RP,
  router can switch to
  source-specific tree

  ▫ increased performance:
    less concentration,
    shorter paths



all data multicast
from rendezvous
point

rendezvous
point

# PIM - Sparse Mode

sender(s):

□ unicast data to RP, which distributes down RP-rooted tree

□ RP can extend mcast tree upstream to source

□ RP can send *stop* msg if no attached receivers

　○ "no one is listening!"



all data multicast from rendezvous point

rendezvous point

# Chapter 4: summary