# Last time

☐ Transitioning to IPv6
- ◆ Tunneling

- ◆ Gateways

☐ Routing
- ◆ Graph abstraction

- ◆ Link-state routing
  - Dijkstra's Algorithm

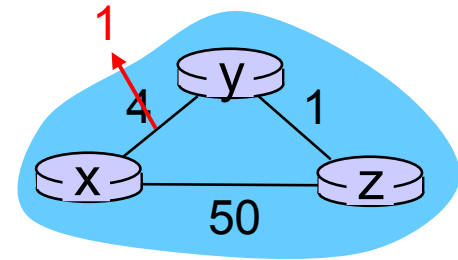- ◆ Distance-vector routing
  - Bellman-Ford Equation

# This time

- Distance vector link cost changes

- Hierarchical routing

- Routing protocols

# Distance Vector: link cost changes

**Link cost changes:**

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

*"good news travels fast"*

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

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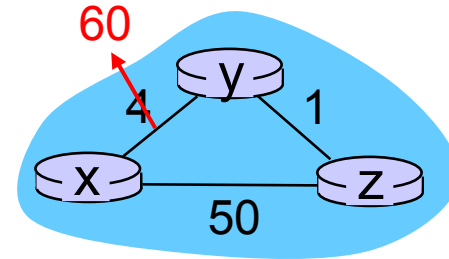
At time $t_2$, $y$ receives **z**'s update and updates its distance table. $y$'s least costs do not change and hence $y$ does *not* send any message to **z**.

# Distance Vector: link cost changes

## Link cost changes:

☐  good news travels fast

☐  bad news travels slow - "count to infinity" problem!

☐  44 iterations before algorithm stabilizes: see text

## Poisoned reverse:

☐  If Z routes through Y to get to X :

♦  Z tells Y its (Z's) distance to X is infinite (so Y won't route to X via Z)

☐  Will this completely solve count to infinity problem?

# Comparison of LS and DV algorithms

## Message complexity

- **LS:** with n nodes, E links, O(nE) msgs sent
- **DV:** exchange between neighbours only
    - ◆ convergence time varies

## Speed of Convergence

- **LS:** O(n²) algorithm requires O(nE) msgs
    - ◆ may have oscillations
- **DV:** convergence time varies
    - ◆ may be routing loops
    - ◆ count-to-infinity problem

## Robustness: what happens if router malfunctions?

### LS:

- ◆ node can advertise incorrect *link* cost
- ◆ each node computes only its *own* table

### DV:

- ◆ DV node can advertise incorrect *path* cost
- ◆ each node's table used by others
    - • error propagates through network

# Chapter 4: Network Layer

# Hierarchical Routing

Our routing study thus far - idealization
- all routers identical
- network "flat"
… *not* true in practice

scale: with 200 million destinations:
- can't store all destinations in routing tables!

- routing table exchange would swamp links!

administrative autonomy
- internet = network of networks

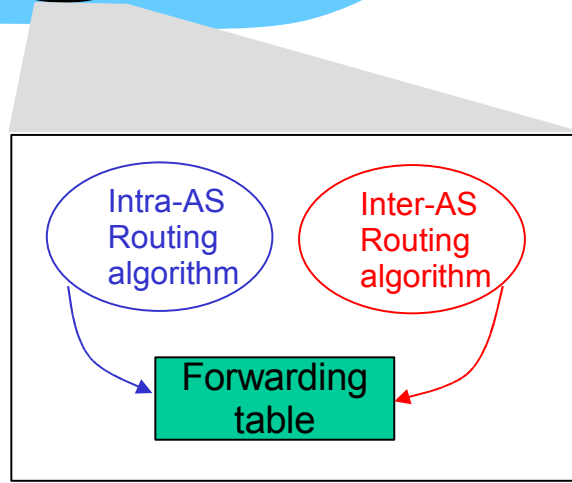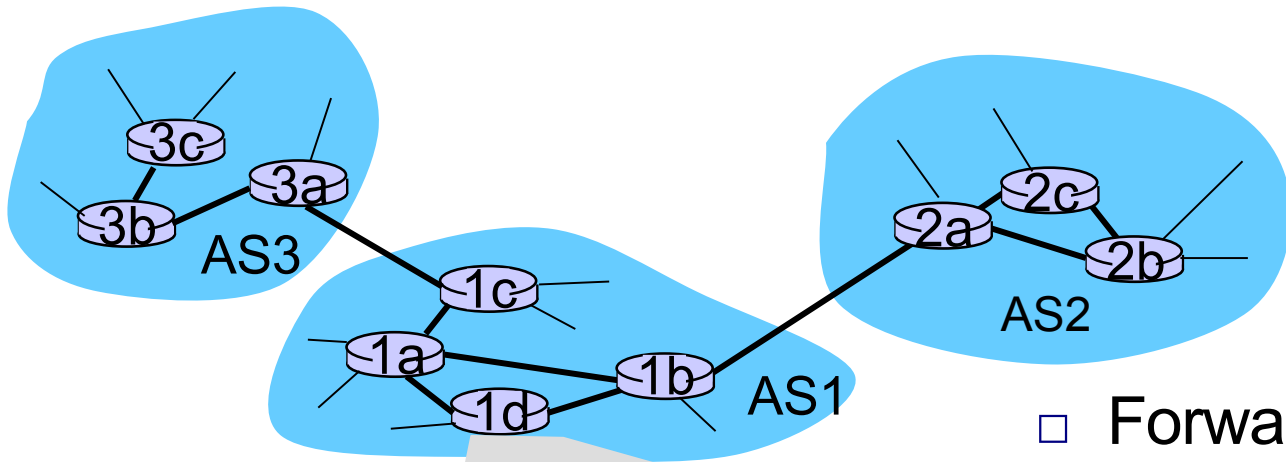- each network admin may want to control routing in his or her own network

# Hierarchical Routing

- Aggregate routers into regions, "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

# Interconnected ASes



Intra-AS
Routing
algorithm

Inter-AS
Routing
algorithm

Forwarding
table

□ Forwarding table is configured by both intra- and inter-AS routing algorithm

♦ Intra-AS sets entries for internal destinations

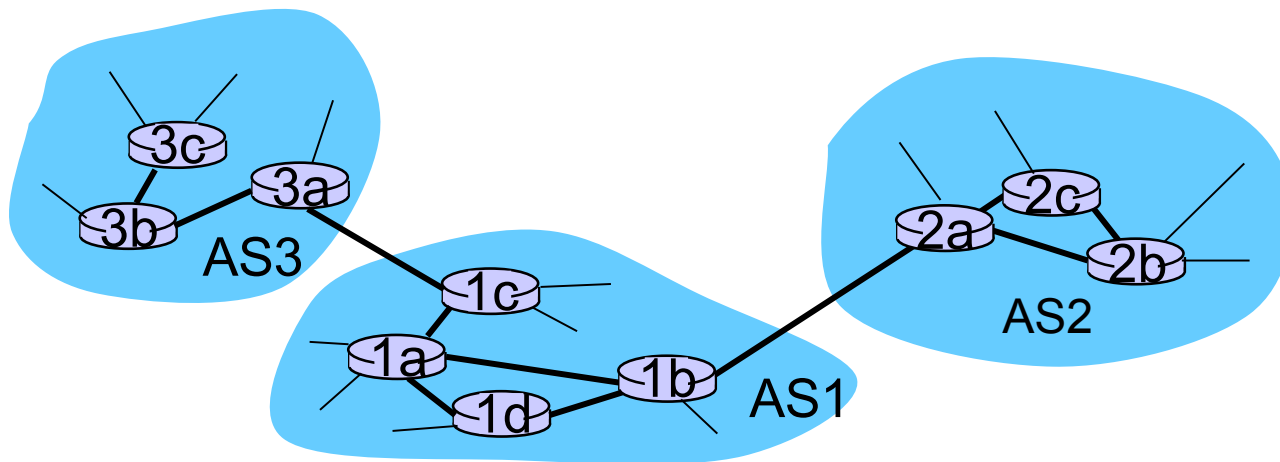♦ Inter-AS & Intra-AS sets entries for external destinations

# Inter-AS tasks

- Suppose router in AS1 receives datagram whose destination is outside of AS1
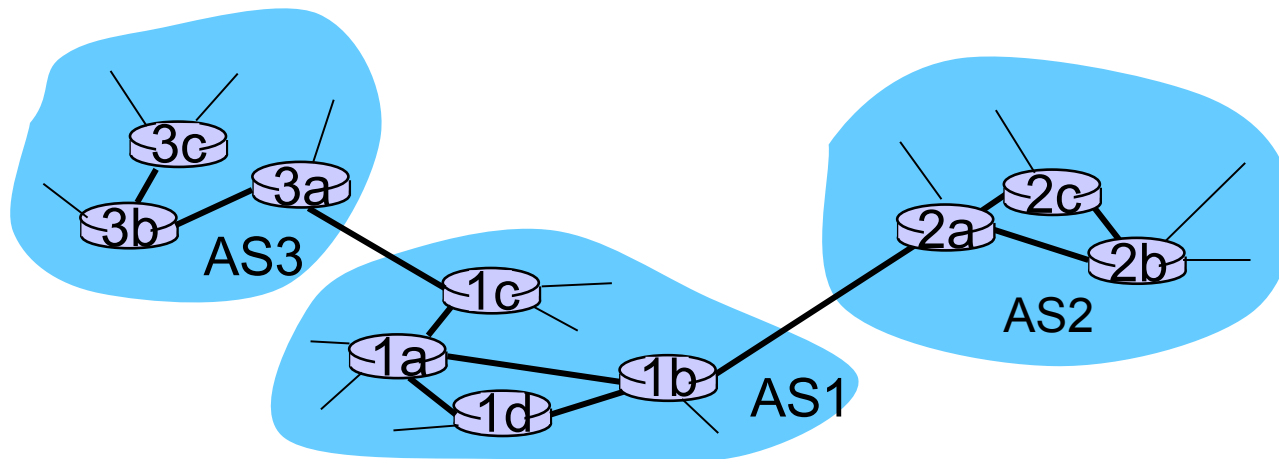  - ♦ Router should forward packet towards one of the gateway routers, but which one?

- to learn which destinations are reachable through AS2 and which through AS3
- to propagate this reachability info to all routers in AS1
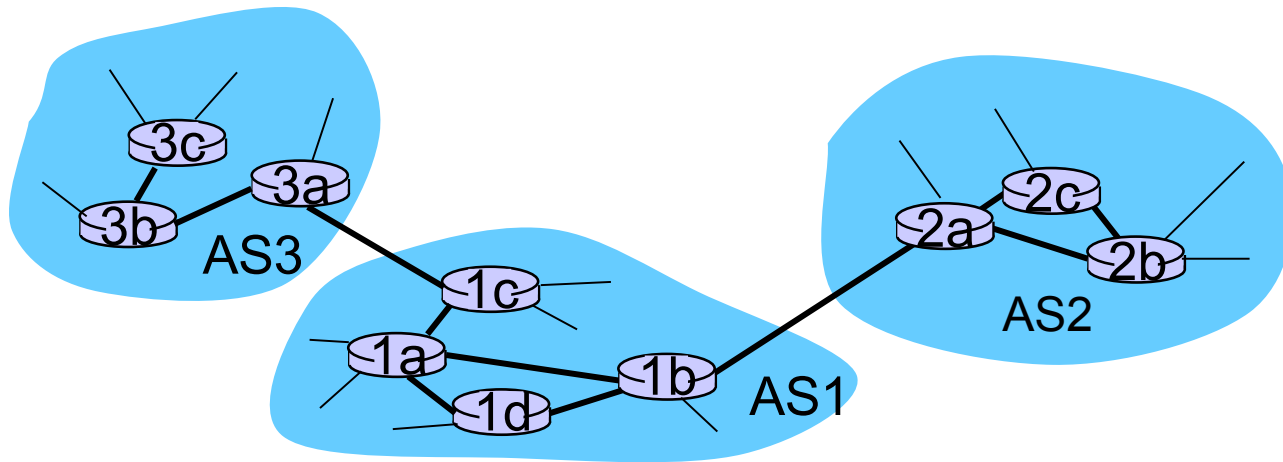
Job of inter-AS routing!

# Example: Setting forwarding table in router 1d

☐ Suppose AS1 learns (via inter-AS protocol) that subnet *x* is 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.

☐ Puts in forwarding table entry *(x,I)*.

# Example: Choosing among multiple ASes

- Now suppose AS1 learns from the inter-AS protocol that subnet *x* is reachable from AS3 *and* from AS2.

- To configure the forwarding table, router 1d must determine towards which gateway it should forward packets for destination x.

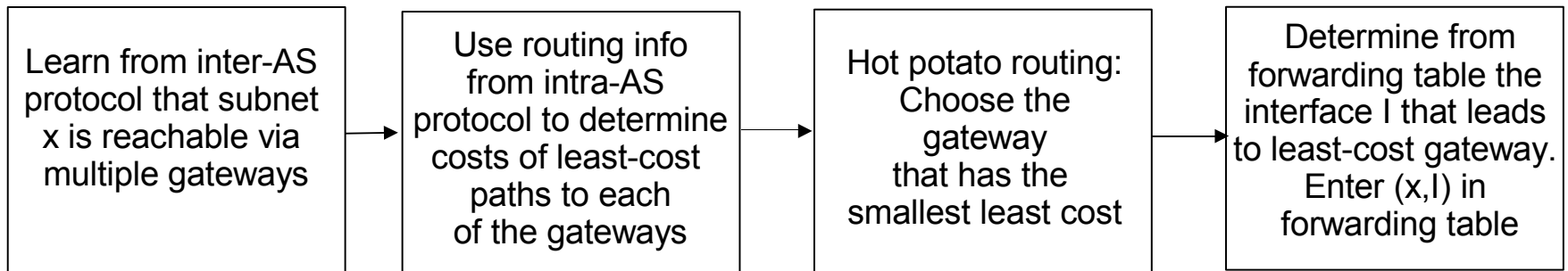- This is also the job of the inter-AS routing protocol!

# Example: Choosing among multiple ASes

- Now suppose AS1 learns from the inter-AS protocol that subnet *x* is reachable from AS3 *and* from AS2.
- To configure the forwarding table, router 1d must determine towards which gateway it should forward packets for destination x.
- This is also the job of the inter-AS routing protocol!
- Hot potato routing: send packet towards closest of two routers.

| 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

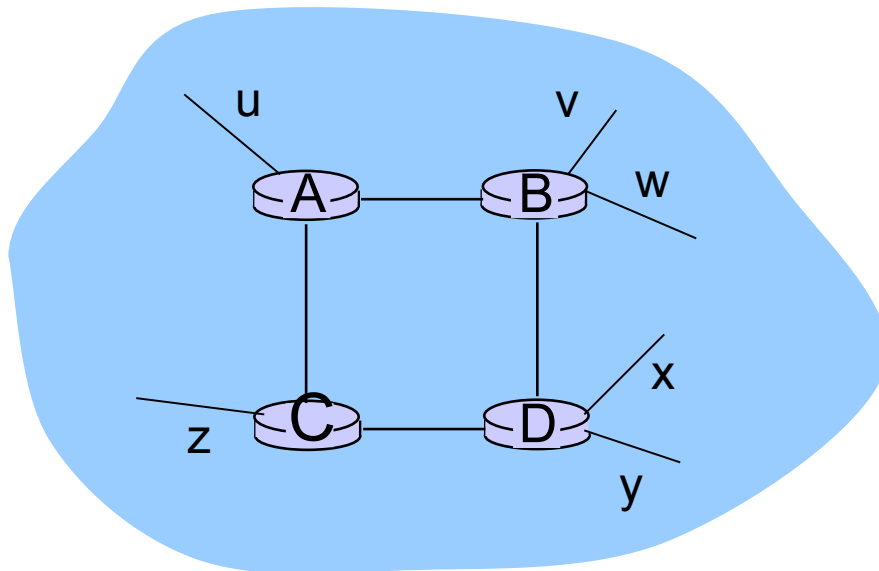☐ Also known as <span style="color:red">Interior Gateway Protocols (IGP)</span>

☐ Most common Intra-AS routing protocols:

♦ RIP: Routing Information Protocol

♦ OSPF: Open Shortest Path First

♦ IGRP: Interior Gateway Routing Protocol (Cisco proprietary)

# 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

# RIP ( Routing Information Protocol)

☐ Distance vector algorithm

☐ Included in BSD-UNIX Distribution in 1982

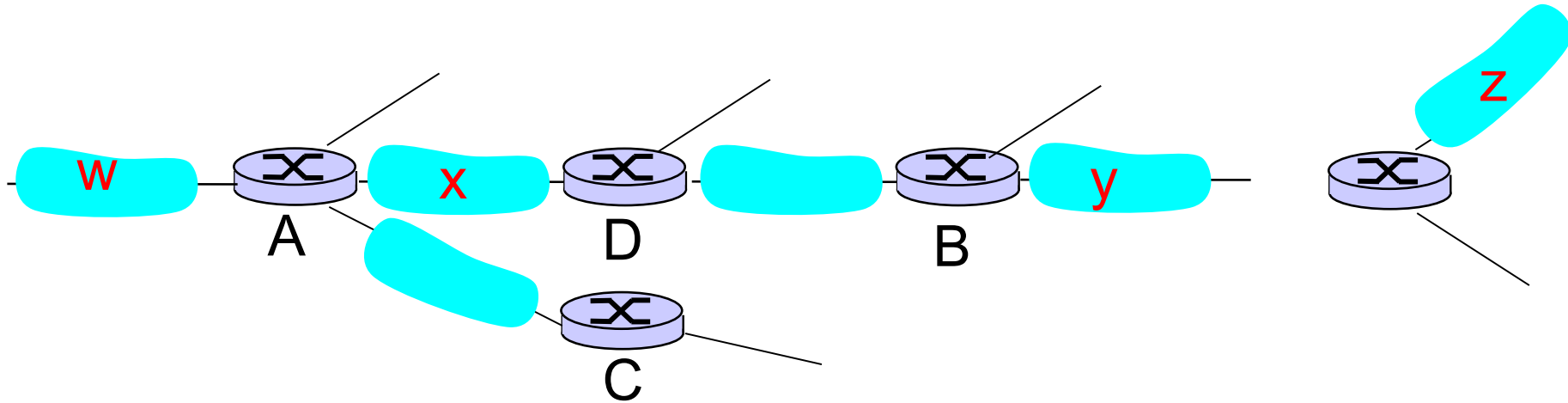☐ Distance metric: # of hops (max = 15 hops)

From router A to subsets:

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

# RIP advertisements

□ Distance vectors: exchanged among neighbours every 30 sec via Response Message (also called an advertisement)

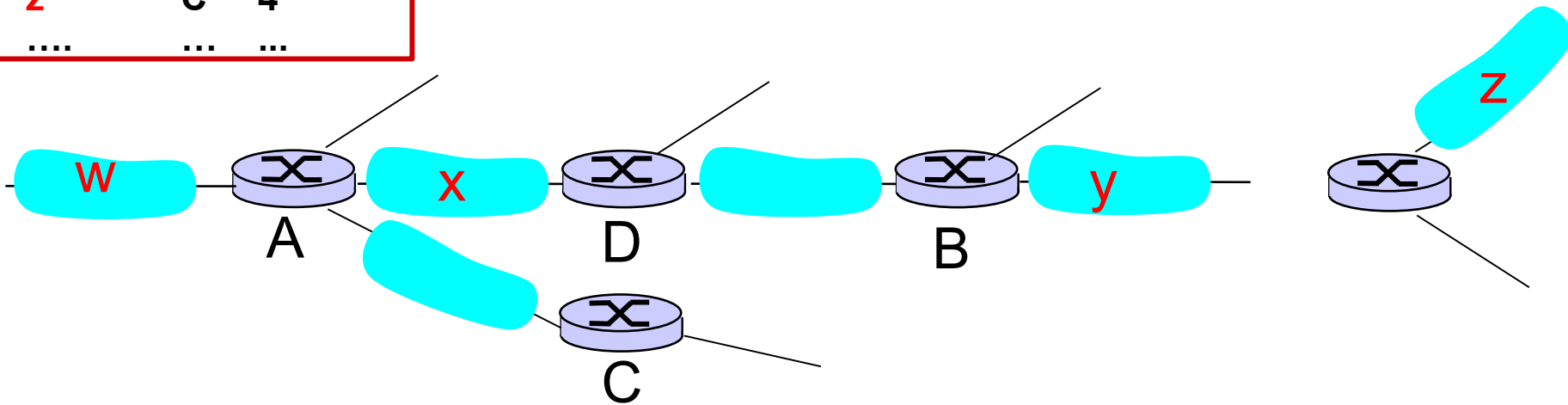□ Each advertisement lists up to 25 destination networks within AS

# 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 --> neighbour/link declared dead

- ♦ routes via neighbour invalidated
- ♦ new advertisements sent to neighbours
- ♦ neighbours 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)
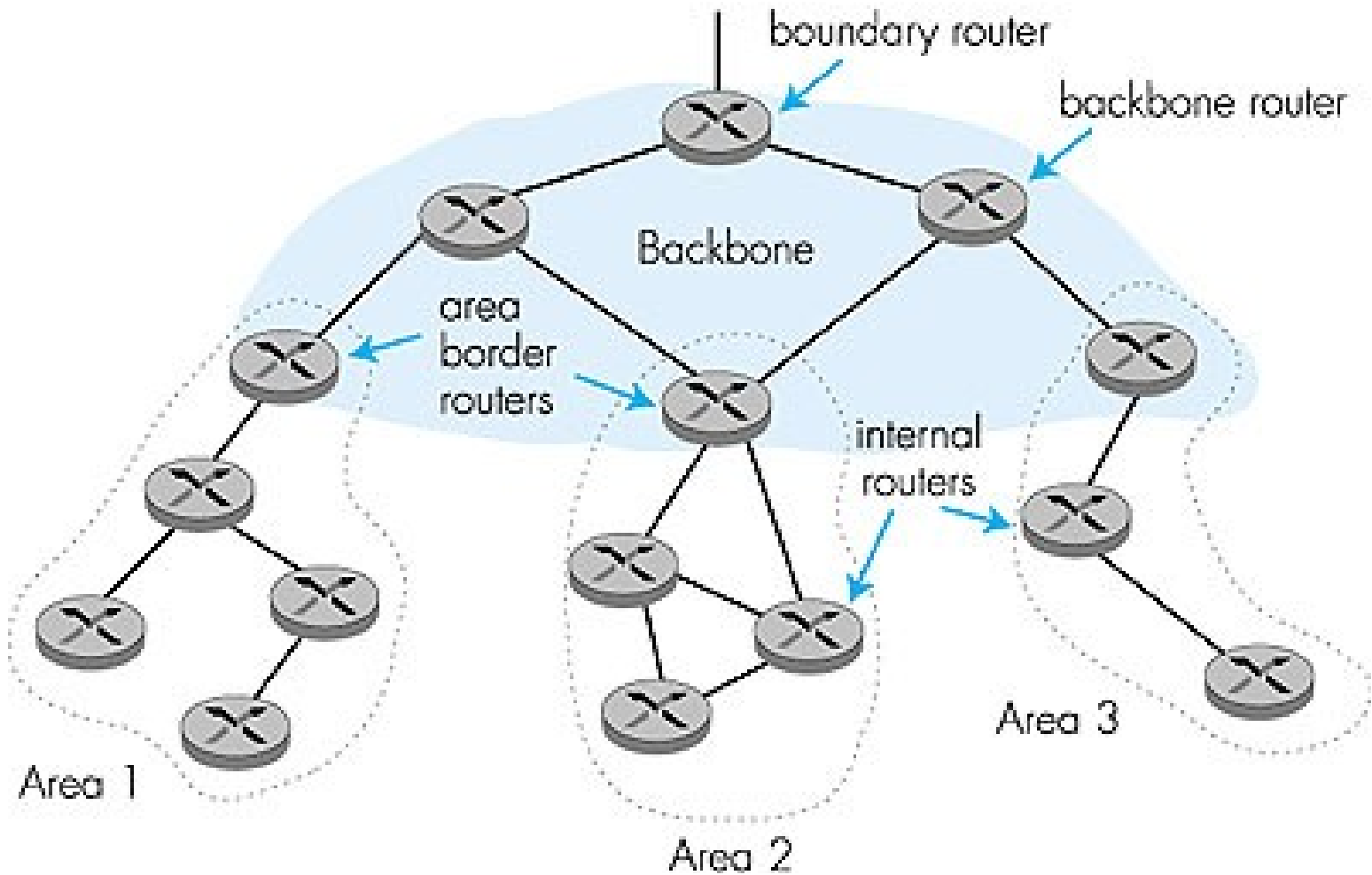
# Chapter 4: Network Layer

# 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 neighbour 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: 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)

□ Integrated uni- and multicast support:

♦ Multicast OSPF (MOSPF) uses same topology data base as OSPF

□ Hierarchical OSPF in large domains.

# Hierarchical OSPF

# Hierarchical OSPF

- Two-level hierarchy: local area, backbone.
  - Link-state advertisements only in area
  - Each node has detailed area topology; only knows 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.

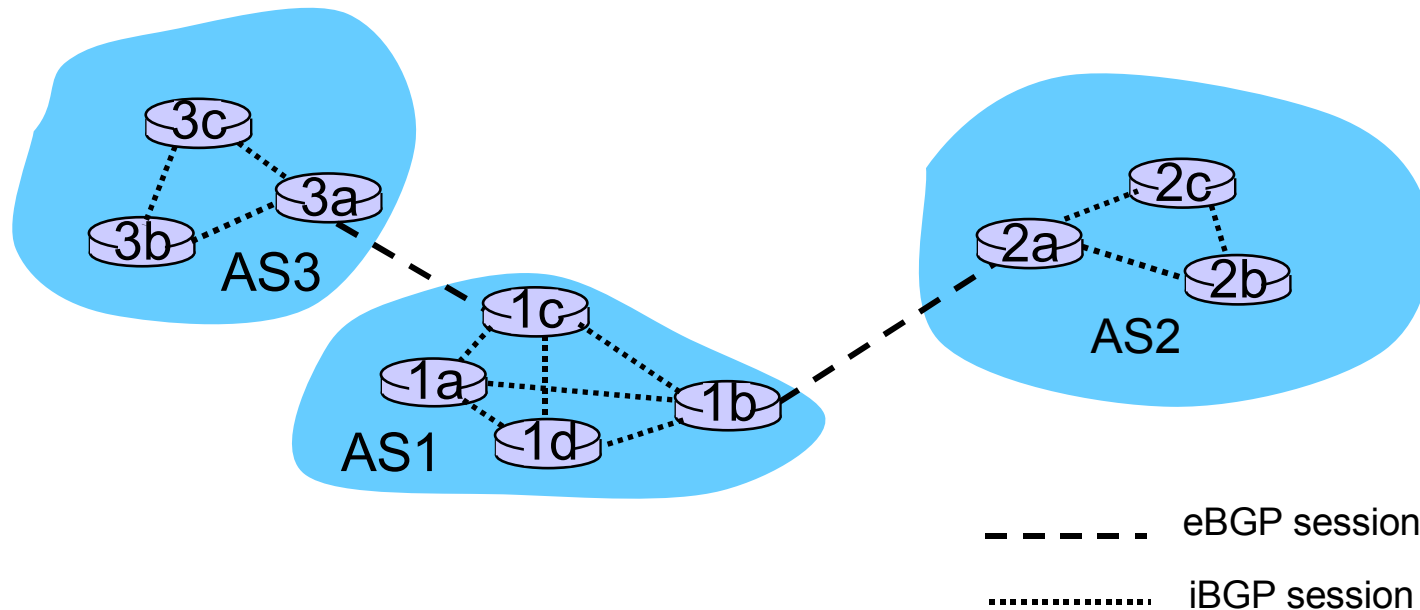# Chapter 4: Network Layer

# Internet inter-AS routing: BGP

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

□ BGP provides each AS a means to:
  ♦ Obtain subnet reachability information from neighbouring ASs.
  ♦ Propagate reachability information to all AS-internal routers.
  ♦ Determine "good" routes to subnets based on reachability information and policy.

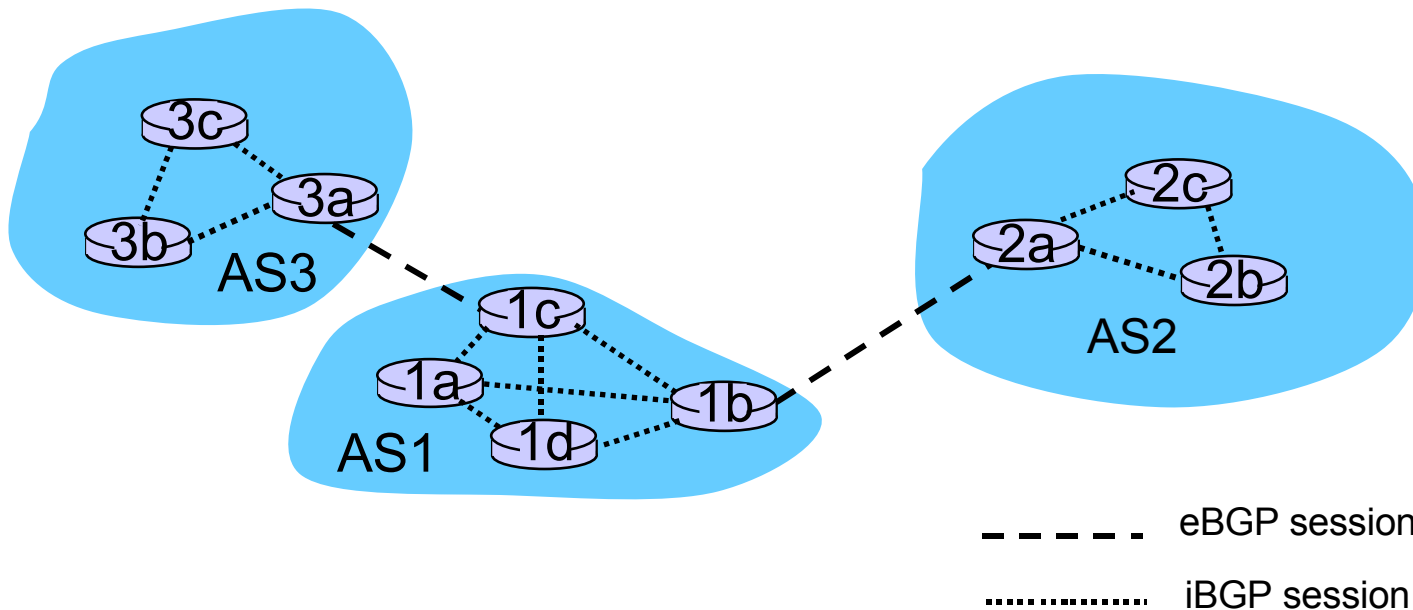□ allows subnet to advertise its existence to rest of Internet: *"I am here"*

# 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 is *promising* it will forward any datagrams destined to that prefix towards the prefix.
  - ♦ AS2 can aggregate prefixes in its advertisement



- - - - eBGP session

.............. iBGP session

# Distributing reachability info

- With eBGP session between 3a and 1c, AS3 sends prefix reachability info to AS1.
- 1c can then use iBGP do distribute this new prefix reach 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.



- - - -  eBGP session

............  iBGP session

# Path attributes & BGP routes

□ When advertising a prefix, advert includes BGP attributes.

♦ prefix + attributes = "route"

□ Two important attributes:

♦ AS-PATH: contains ASs through which prefix advertisement has passed: AS 67 AS 17

♦ NEXT-HOP: Indicates specific internal-AS router to next-hop AS. (There 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 one 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

# Recap

- Distance vector link cost changes
  - Count-to-infinity, poisoned reverse

- Hierarchical routing
  - Autonomous Systems
  - Inter-AS, Intra-AS routing

- Routing protocols
  - Intra-AS
    - RIP
    - OSPF
  - Inter-AS
    - BGP

# Next time

- BGP policy

- Broadcast / multicast routing

- ATM / MPLS

- Link virtualization

- Router internals