

Network Layer:

Distance Vector Routing, Link State Routing
Global Internet Routing (Interdomain, BGP)

Qiao Xiang, Congming Gao, **Qiang Su**

<https://sngroup.org.cn/courses/cnns-xmuf25/index.shtml>

11/25/2025

Outline

- ❑ Admin and recap
- ❑ Network control plane
 - Routing
 - Link weights assignment
 - Routing computation
 - Distance vector protocols (distributed computing)
 - synchronous Bellman-Ford (SBF)
 - asynchronous Bellman-Ford (ABF)
 - properties of DV
 - DV w/ loop prevention
 - reverse poison
 - *destination-sequenced DV (DSDV)*

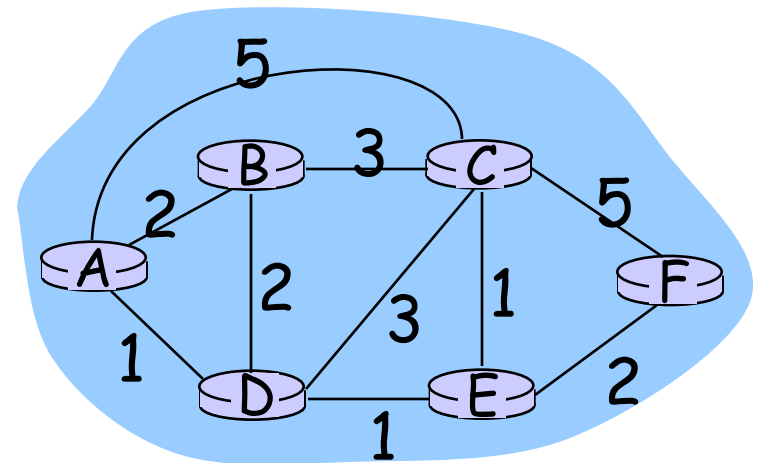
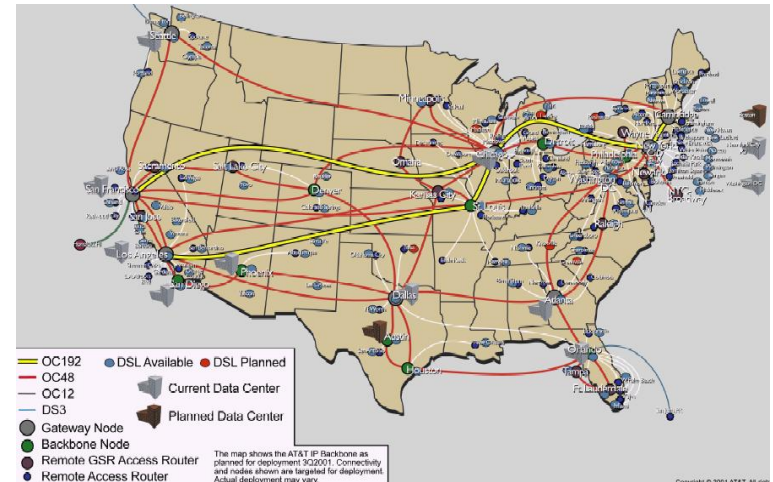
Recap: Routing Context

Routing

Goal: determine "good" paths (sequences of routers) thru networks from source to dest.

Often depends on a graph abstraction:

- graph nodes are routers
- graph edges are physical links
 - links have properties: delay, capacity, \$ cost, **policy**

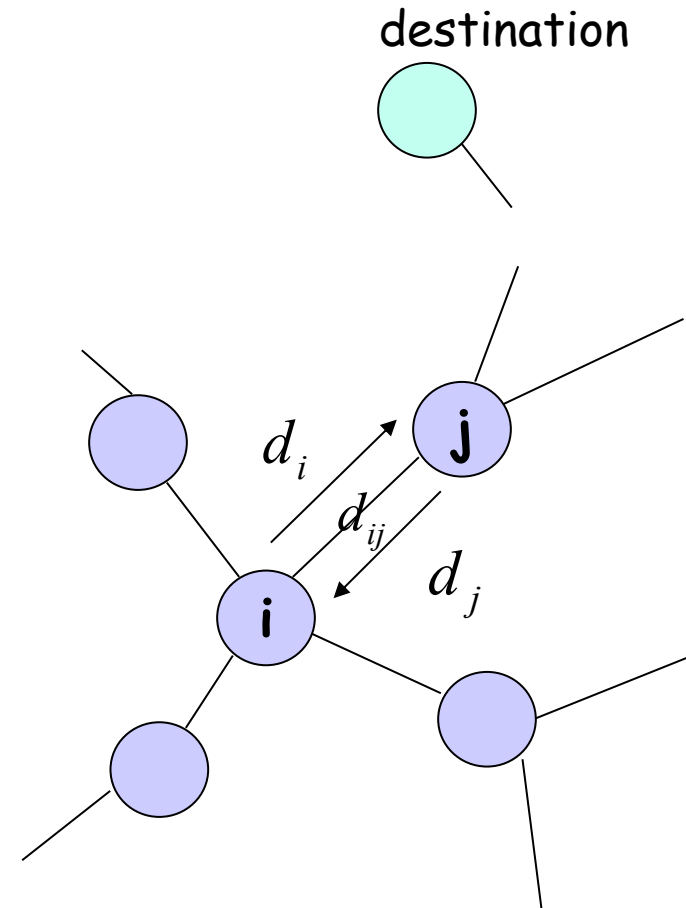


Recap: Routing Computation using Distance Vector/Bellman-Ford Routing

- Distributed computation:
At node i , computes

$$d_i = \min_{j \in N(i)} (d_{ij} + d_j)$$

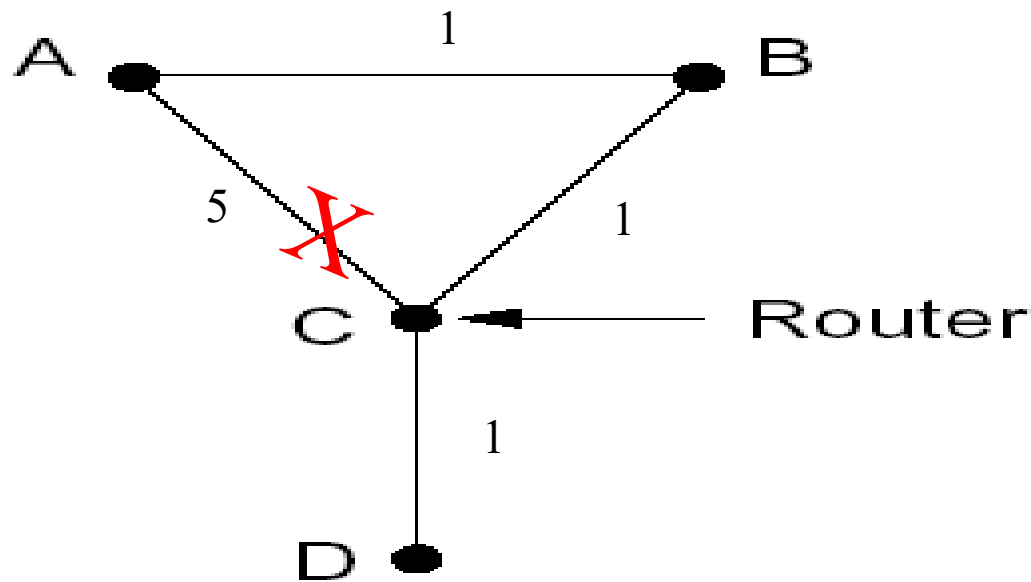
- One way to understand BFA is to consider it as a dynamic programming alg, propagating from dest to other nodes



Recap: Destination-Sequenced Distance Vector protocol (DSDV)

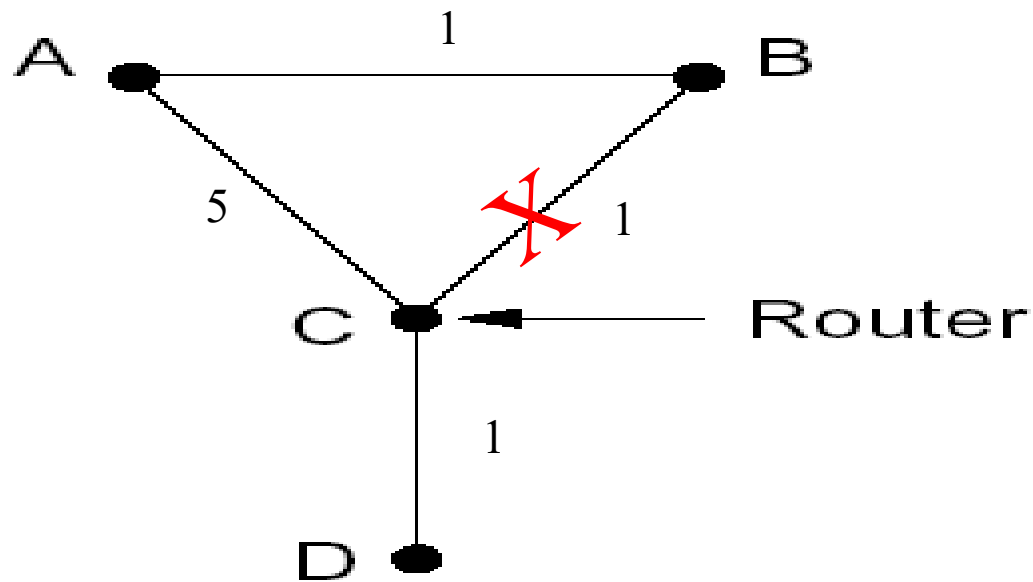
- ❑ Basic idea: use sequence numbers to partition computation
 - tags each route with a sequence number
 - each destination node D periodically advertises monotonically increasing even-numbered sequence numbers
 - when a node realizes that **the link it uses to reach destination D is broken**, it advertises an **infinite** metric and a **sequence number which is one greater** than the previous route (i.e., an odd seq. number)
 - the route is repaired by a later even-number advertisement from the destination

Recap: Example



Will this trigger an update?

Recap: Example



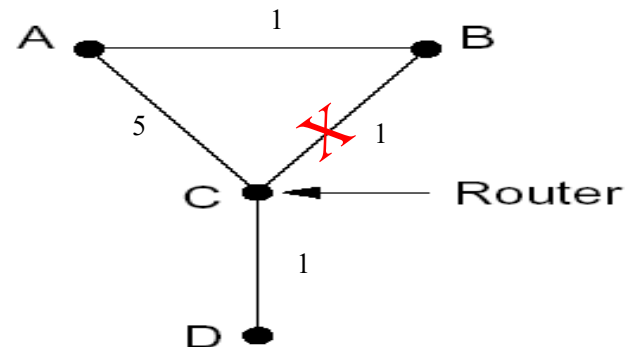
Will this trigger an update?

Outline

- ❑ Admin and recap
- ❑ Network control plane
 - Routing
 - Link weights assignment
 - Routing computation
 - Distance vector protocols (distributed computing)
 - synchronous Bellman-Ford (SBF)
 - asynchronous Bellman-Ford (ABF)
 - properties of DV
 - DV w/ loop prevention
 - reverse poison
 - destination-sequenced DV (DSDV)
 - *diffusive update algorithm (DUAL) and EIGRP*

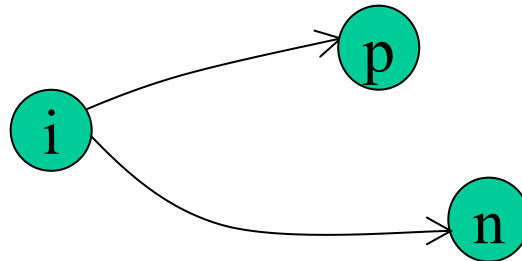
Basic Idea

- ❑ DSDV guarantees no loop, but at the price of not using any backup path before destination re-announces reachability.
- ❑ Basic idea: Sufficient condition to guarantee no loop using backup paths (called switching)?



Key Idea: Feasible Successors

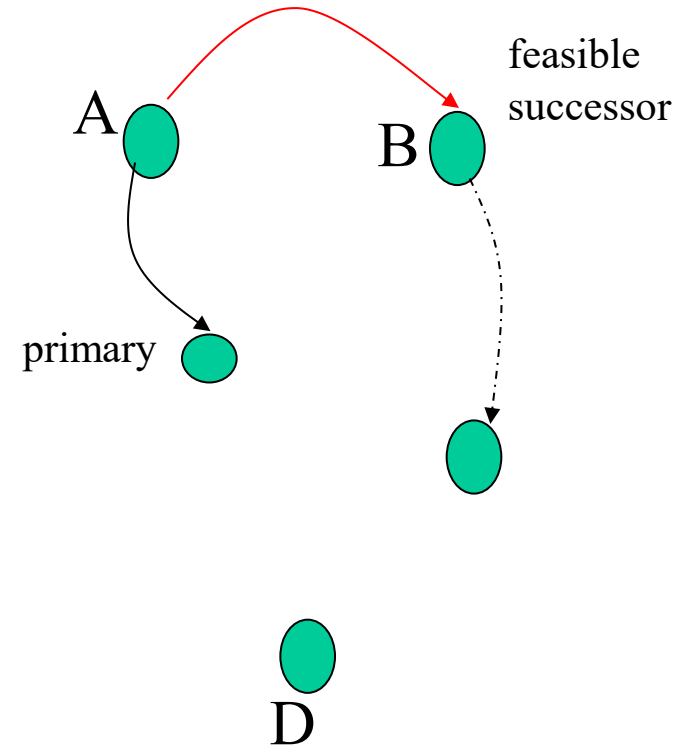
- If the reported distance of a neighbor n is lower than the total distance using primary (current shortest), the neighbor n is a feasible successor



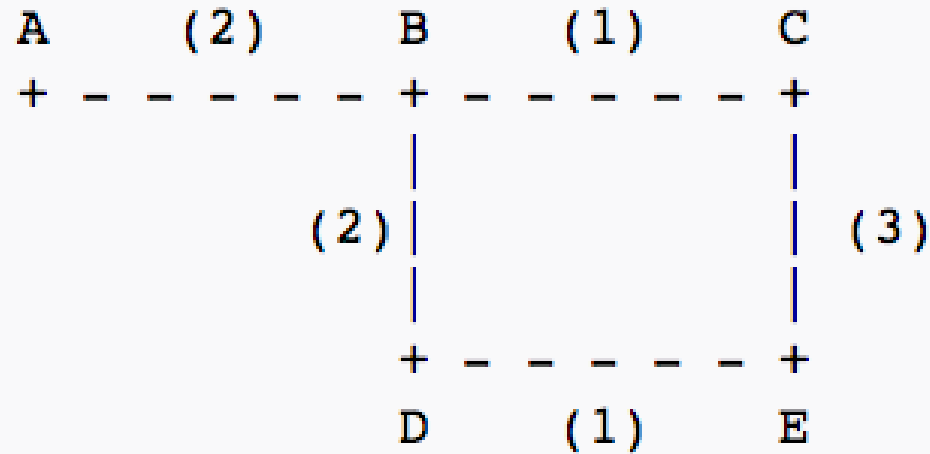
$$d_n + d_{i \rightarrow n} \geq d_{\text{primary}} + d_{i \rightarrow \text{primary}} > d_n$$

Intuition

- Since the reported distance of B is lower than my total distance, B cannot be using me (along a path) to reach the destination



Example



□ Assume A is destination, consider E

	Reported Dist.	Total Dist.
Neighbor C	3	6
Neighbor D	4	5

Summary: Distance Vector Routing

□ Basic DV protocol

- take away: use monotonicity as a technique to understand liveness/convergence
 - highly recommended reading of Bersekas/Gallager chapter

□ Fix counting-to-infinity problem

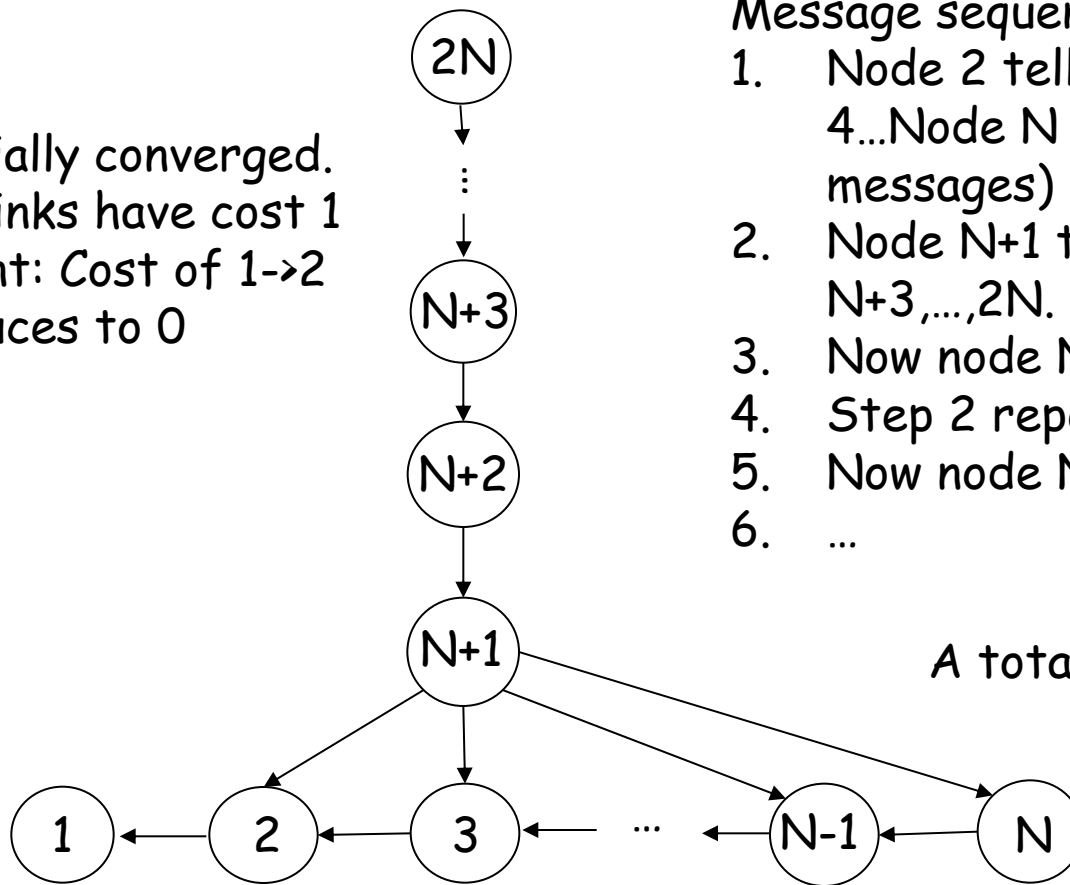
- DSDV
 - Idea: uses sequence number to avoid routing loops
 - seq# partitions routing updates from different outside events
 - within same event, no loop so long each node only decreases its distance
 - Analysis: use global invariants to understand/design safety/no routing loops
- EIRGP (DUAL)
 - Idea: introduces a sufficient condition for local recovery

Discussion: Distance Vector Routing

- ❑ What do you like about distributed, distance vector routing?
- ❑ What do you **not** like about distributed, distance vector routing?

Churns of DV: One Example

Initially converged.
All links have cost 1
Event: Cost of 1→2
reduces to 0



Message sequences

1. Node 2 tells 3. Node 3 tells 4...Node N tells $N+1$. ($N-1$ messages)
2. Node $N+1$ tells $N+2$, $N+2$ tells $N+3$, ..., $2N$. ($N-1$ messages)
3. Now node $N-1$ tells node $N+1$
4. Step 2 repeats
5. Now node $N-2$ tells node $N+1$
6. ...

A total of $O(N^2)$ messages

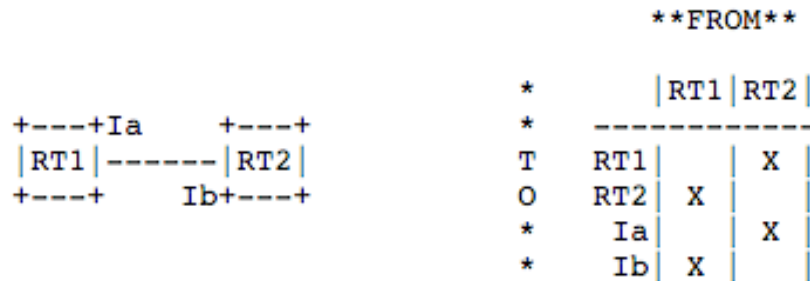
Outline

- ❑ Admin and recap
- ❑ Network control plane
 - Routing
 - Link weights assignment
 - Routing computation
 - Distance vector protocols (distributed computing)
 - *Link state protocols (distributed state synchronization)*

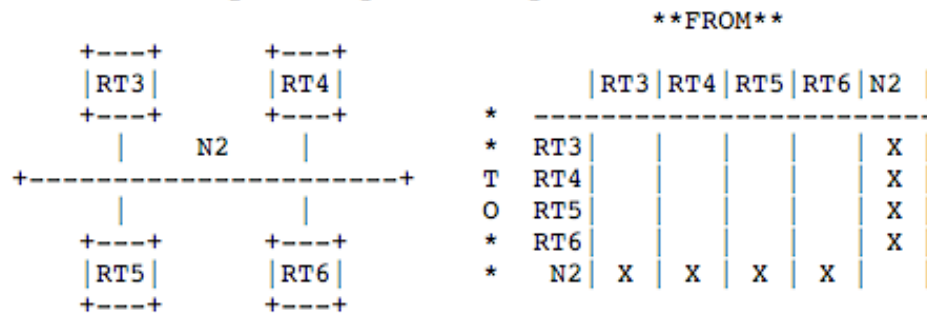
Link-State Routing

- ❑ Basic idea: Not distributed computing, only distributed state distribution
- ❑ Net topology, link costs are distributed to all nodes
 - all nodes have same info
 - Each node computes its shortest paths from itself to all other nodes
 - standard Dijkstra's algorithm as path compute alg
 - Allows multiple same-cost paths
 - Multiple cost metrics per link (for type of service routing)
- ❑ Most commonly used routing protocol (e.g., OSPF/ISIS) by most networks in Internet

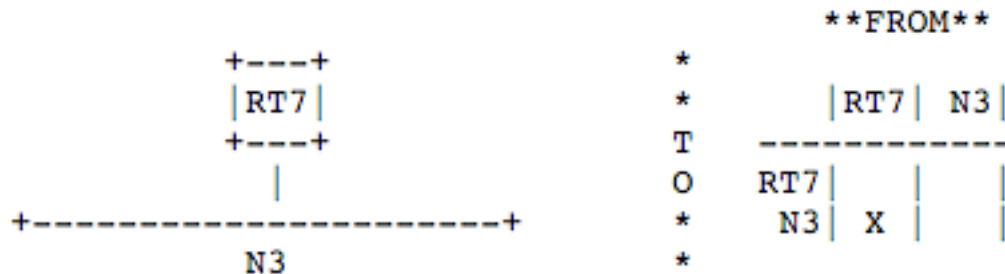
Example: Link State and Directed Graph (OSPFv2)



Physical point-to-point networks



Multi-access networks



Stub multi-access networks

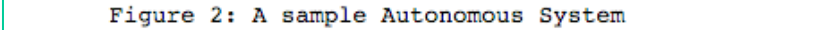


Figure 3: The resulting directed graph

Figure 3: The resulting directed graph

Outline

- ❑ Admin and recap
- ❑ Network control plane
 - Routing
 - Link weights assignment
 - Routing computation
 - Distance vector protocols (distributed computing)
 - *Link state protocols (distributed state synchronization)*
 - data structure to be distributed
 - *state distribution protocol*

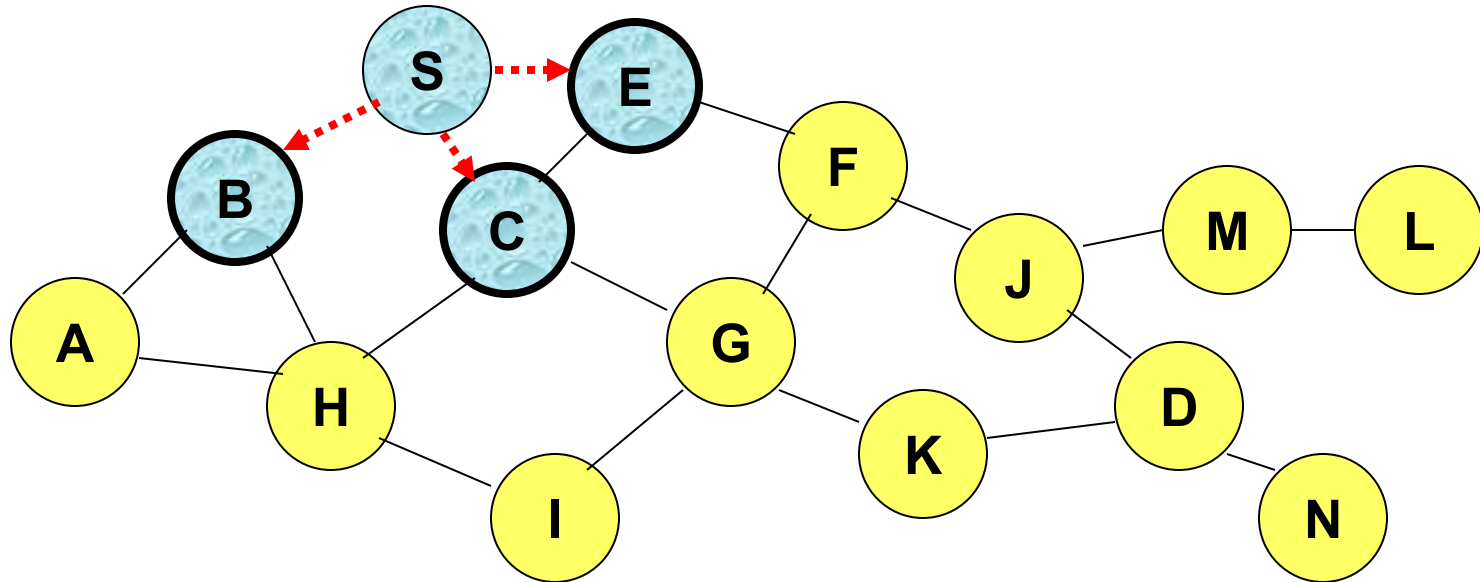
Basic Link State Broadcast Protocol

Basic event structure at node n

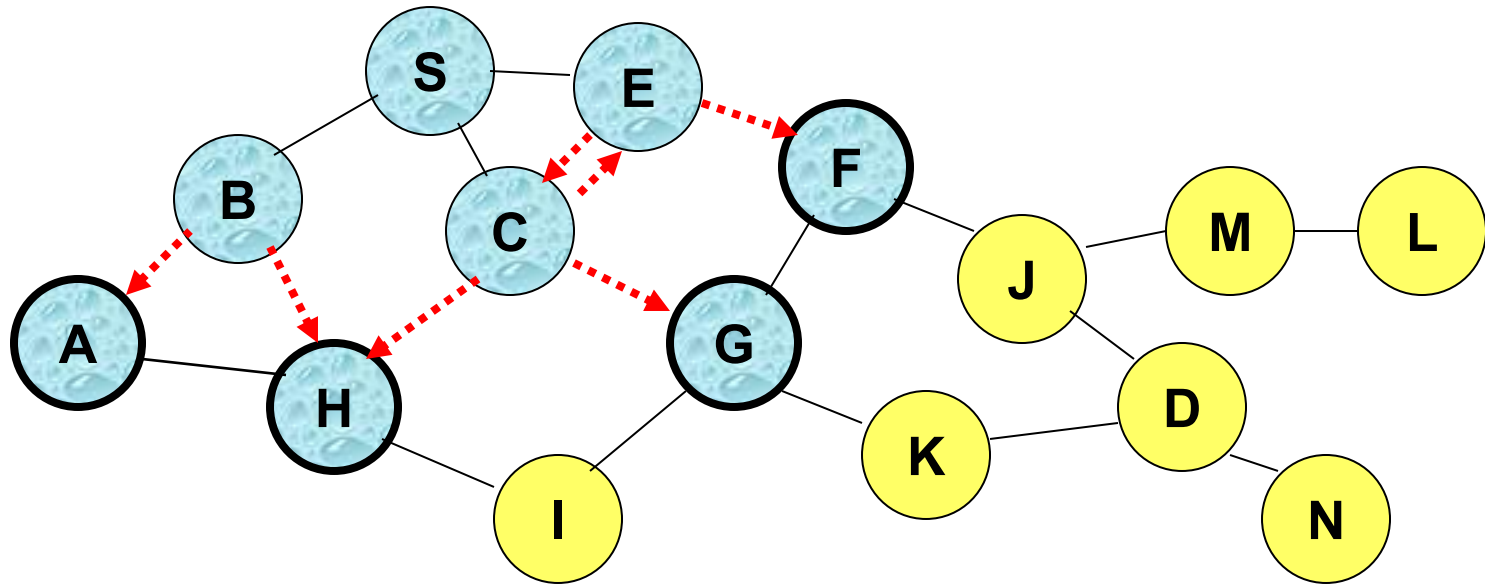
- on initialization:
 - broadcast $LSA[e]$ for each link e connected to n
- on state change to a link e connected to n :
 - broadcast $LSA[e] = \text{new status}$
- on receiving an $LSA[e]$:
 - if (does not have $LSA[e]$)
forwards $LSA[e]$ to all links except the incoming link

Link State Broadcast

Node S updates link states connected to it.



Link State Broadcast



To avoid forwarding the same link state announcement (LSA) multiple times (forming a loop), each node remembers the received LSAs.

- Second LSA[S] received by E from C is discarded
- Second LSA[S] received by C from E is discarded as well
- Node H receives LSA[S] from two neighbors, and will discard one of them

Discussion

- ❑ Issues of the basic link state protocol?
 - Recall: goal is to efficiently distribute to each node to a correct, complete link state map

Link State Broadcast: Issues

❑ Problem: Out of order delivery

- link down and then up
- A node may receive up first and then down

❑ Solution

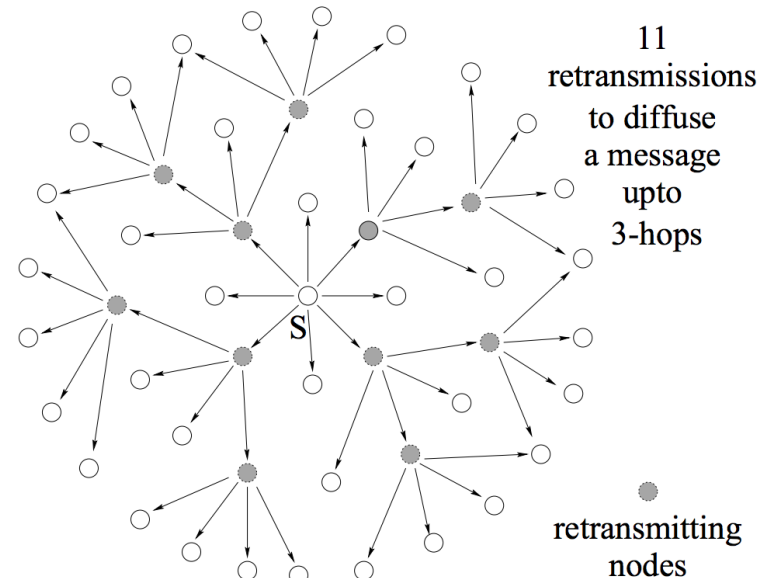
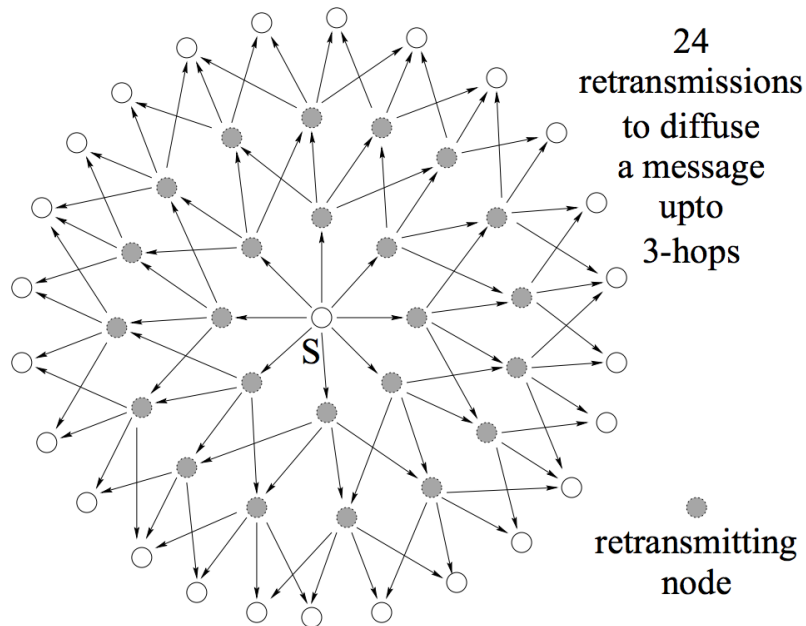
- Each link update is given a sequence number: (initiator, seq#, link, status)
 - the initiator should increase the seq# for each new update
- If the seq# of an update of a link is not higher than the highest seq# a router has seen, drop the update
- Otherwise, forward it to all links except the incoming link (real implementation using packet buffer)
- Problem of solution: seq# corruption
- Solution: age field (e.g., <https://tools.ietf.org/html/rfc1583#page-102>)

Link State Broadcast: Issues

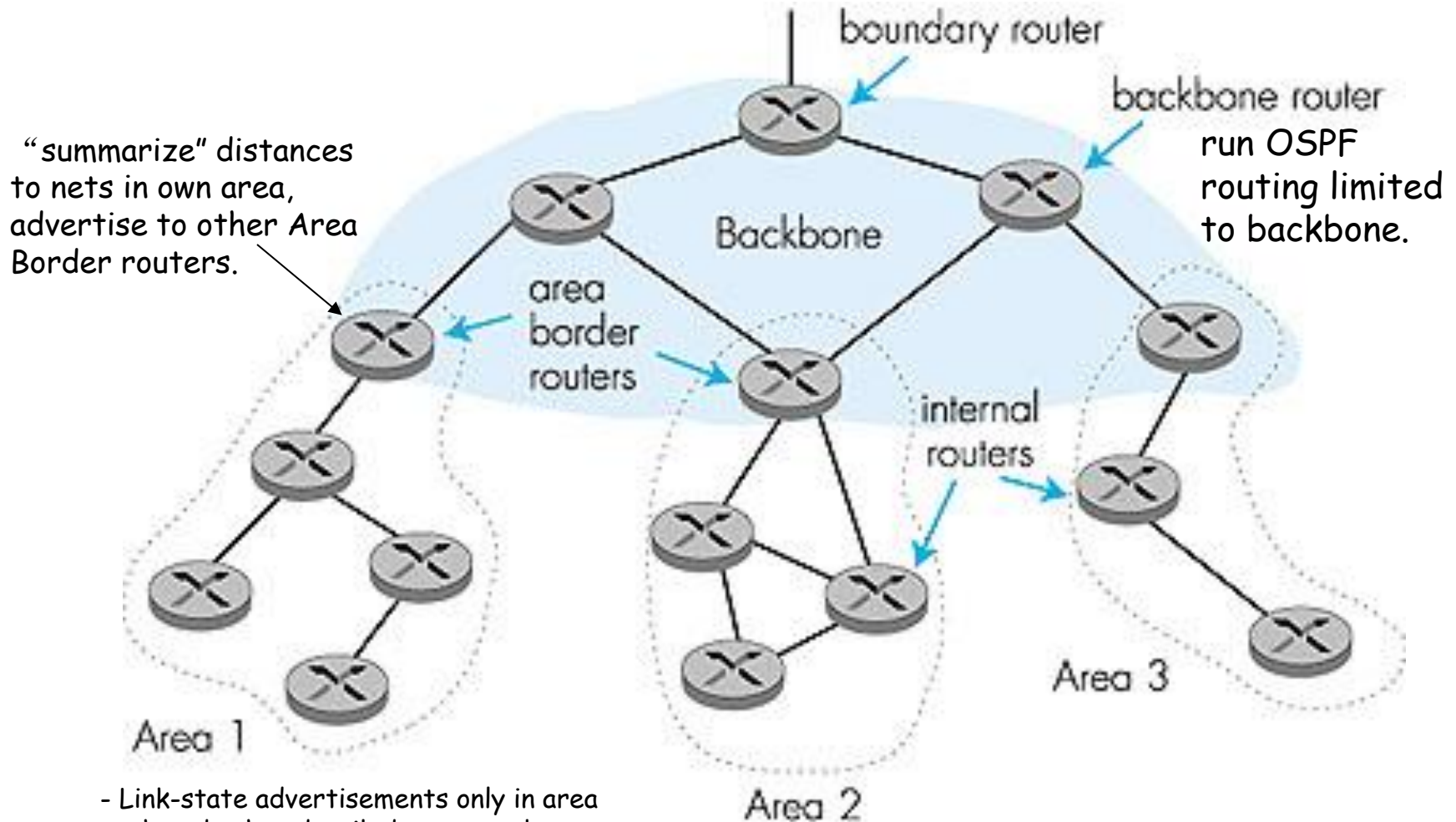
- ❑ Problem: network partition and then reconnect, how to sync across the reconnected components
- ❑ Solution: updates are sent periodically

Link State Broadcast: Issues

□ Problem: Broadcast redundancy



Hierarchical OSPF



- Link-state advertisements only in area
- each node has detailed area topology;
- only know direction (shortest path) to nets in other areas.

Two-level hierarchy: local area, backbone.

Summary: Link State

Basic LS protocol

- take away: instead of computing routing results using distributed computing, distributed computing is for only link state distribution (synchronization)

- ❑ Link state distribution can still have much complexity, e.g., out of order delivery, partition and reconnect, scalability

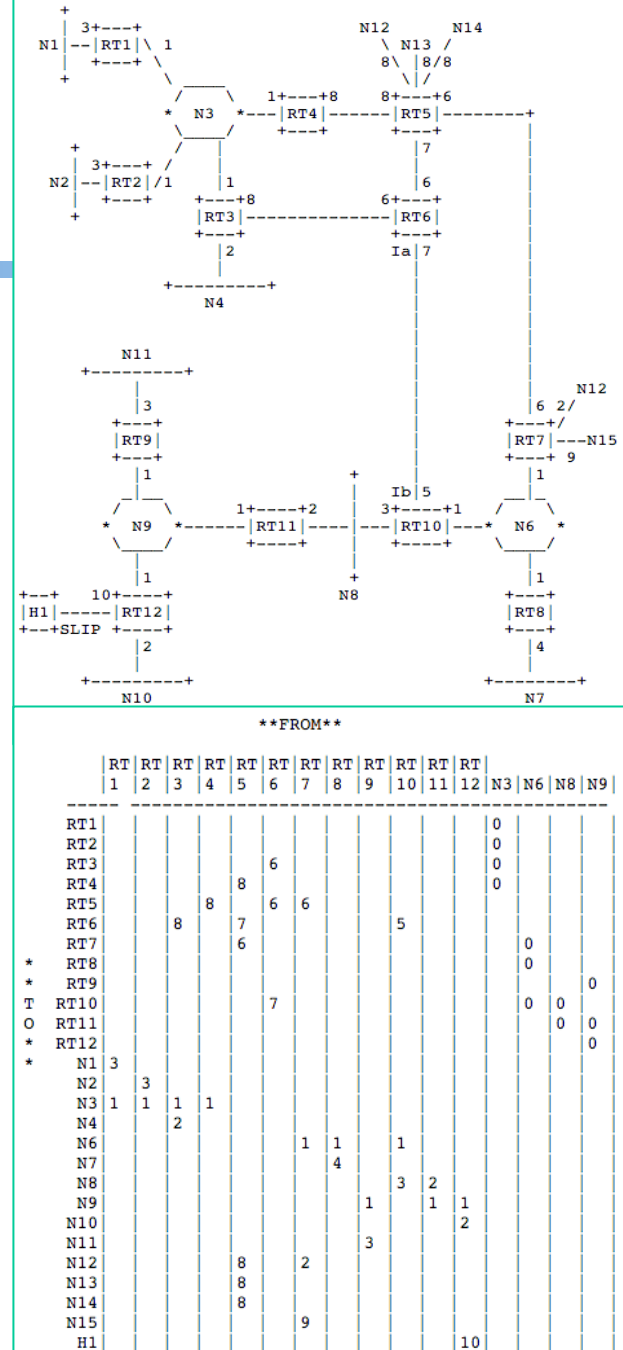
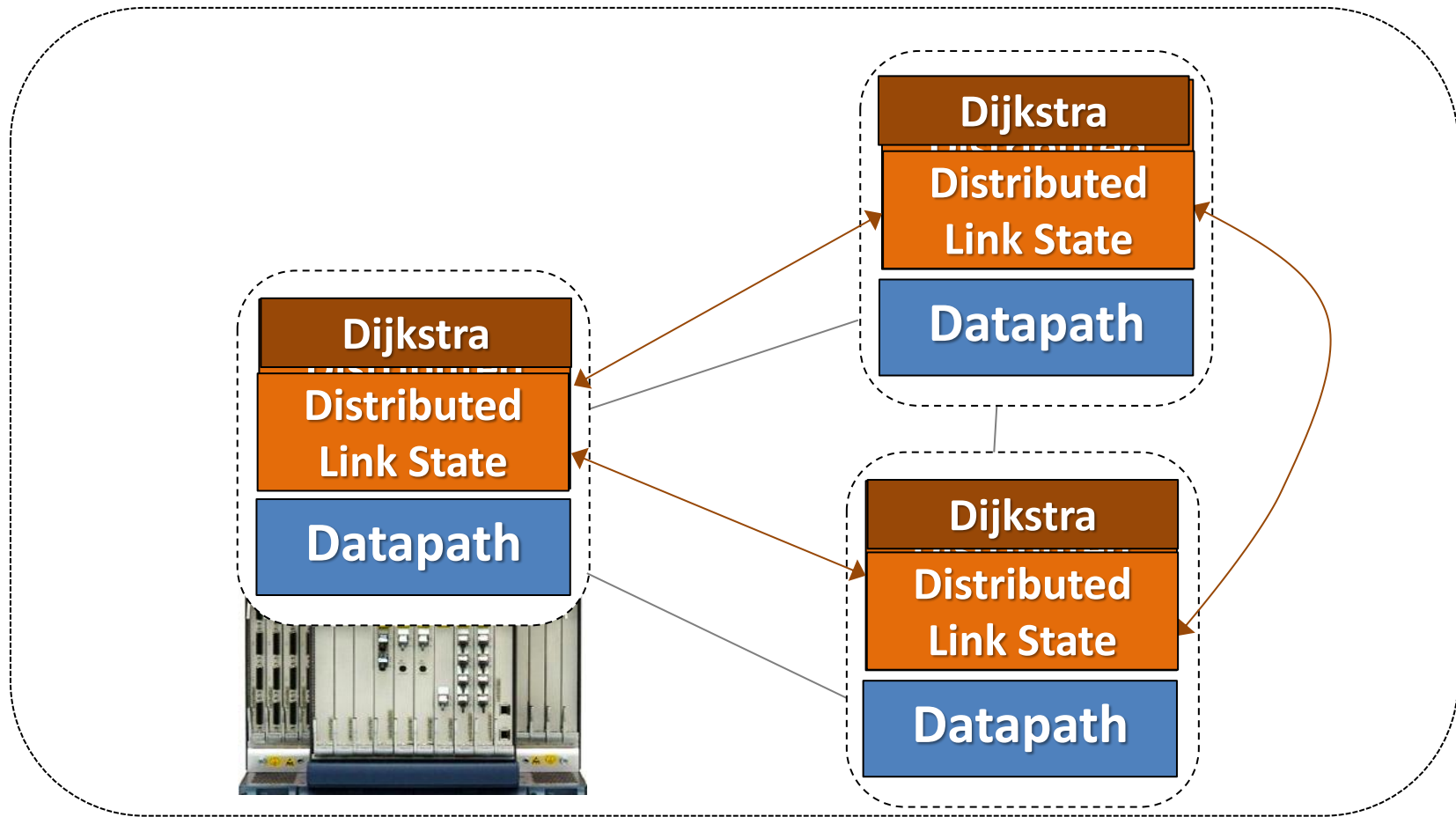


Figure 3: The resulting directed graph

Roadmap: Routing Computation Architecture Spectrum

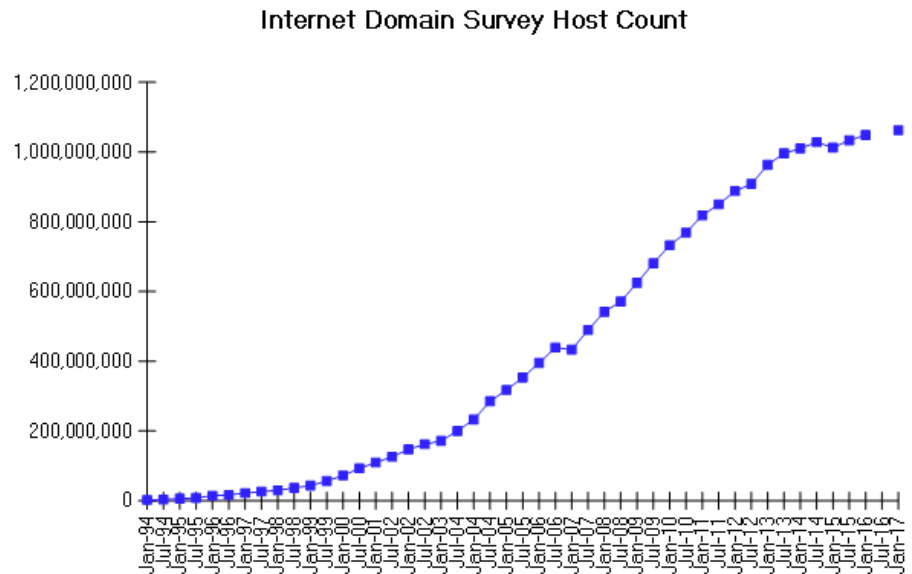
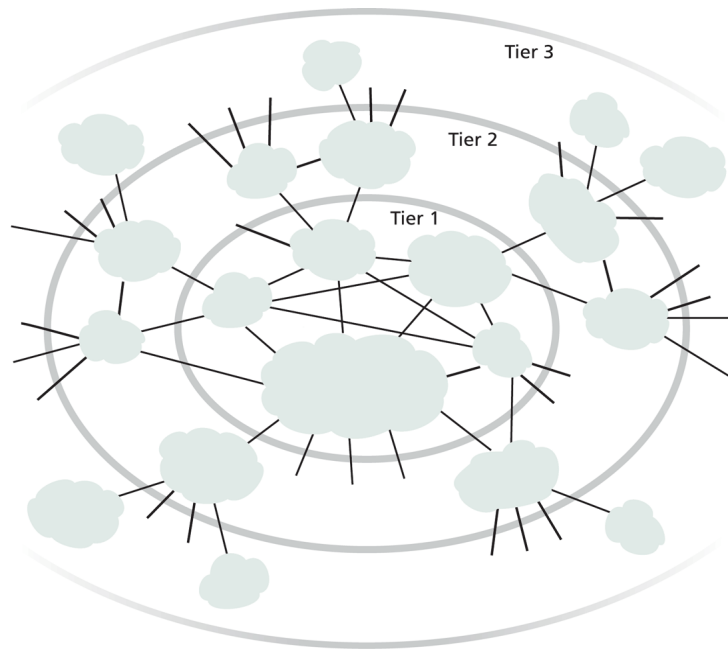


Outline

- ❑ Admin and recap
- ❑ Network control plane
 - Routing
 - Link weights assignment
 - Routing computation
 - Basic routing computation protocols
 - Distance vector protocols (distributed computing)
 - Link state protocols (distributed state synchronization)
 - Global Internet routing

Exercise

- ❑ Does it work to use DV or LS as we discussed for global Internet routing?



Source: Internet Systems Consortium (www.isc.org)

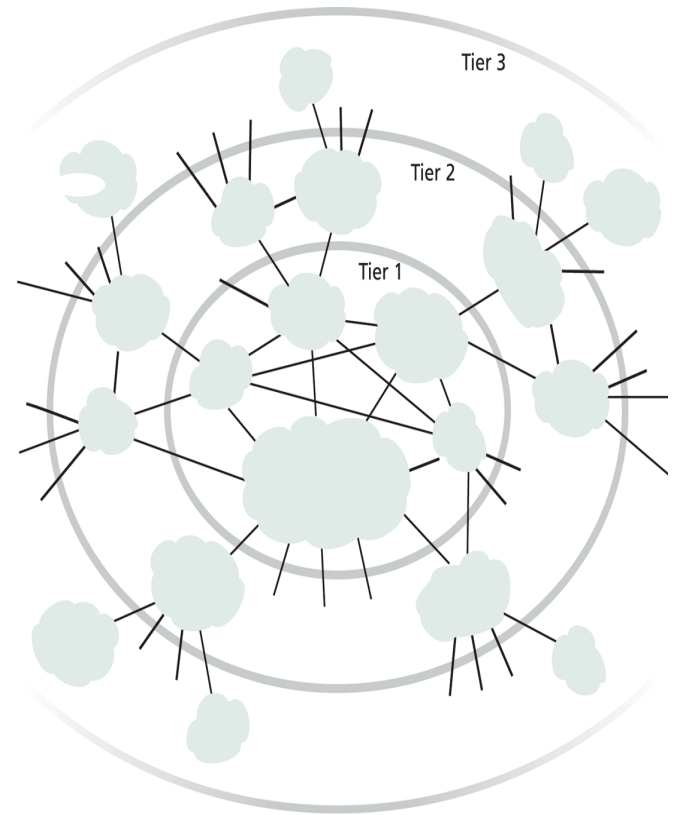
Requirements and Solution of Current Global Internet Routing

- ❑ Scalability: handle network size (#devices) much higher than typical DV or LS can handle
 - Solution: Introduce **new abstraction** to reduce network (graph) size

- ❑ Autonomy: allow each network to have individual preference of routing (full control of its internal routing; control/preference of routing spanning multiple networks)
 - Solution: **hierarchical routing** and **policy routing**

New Abstraction: Autonomous Systems (AS)

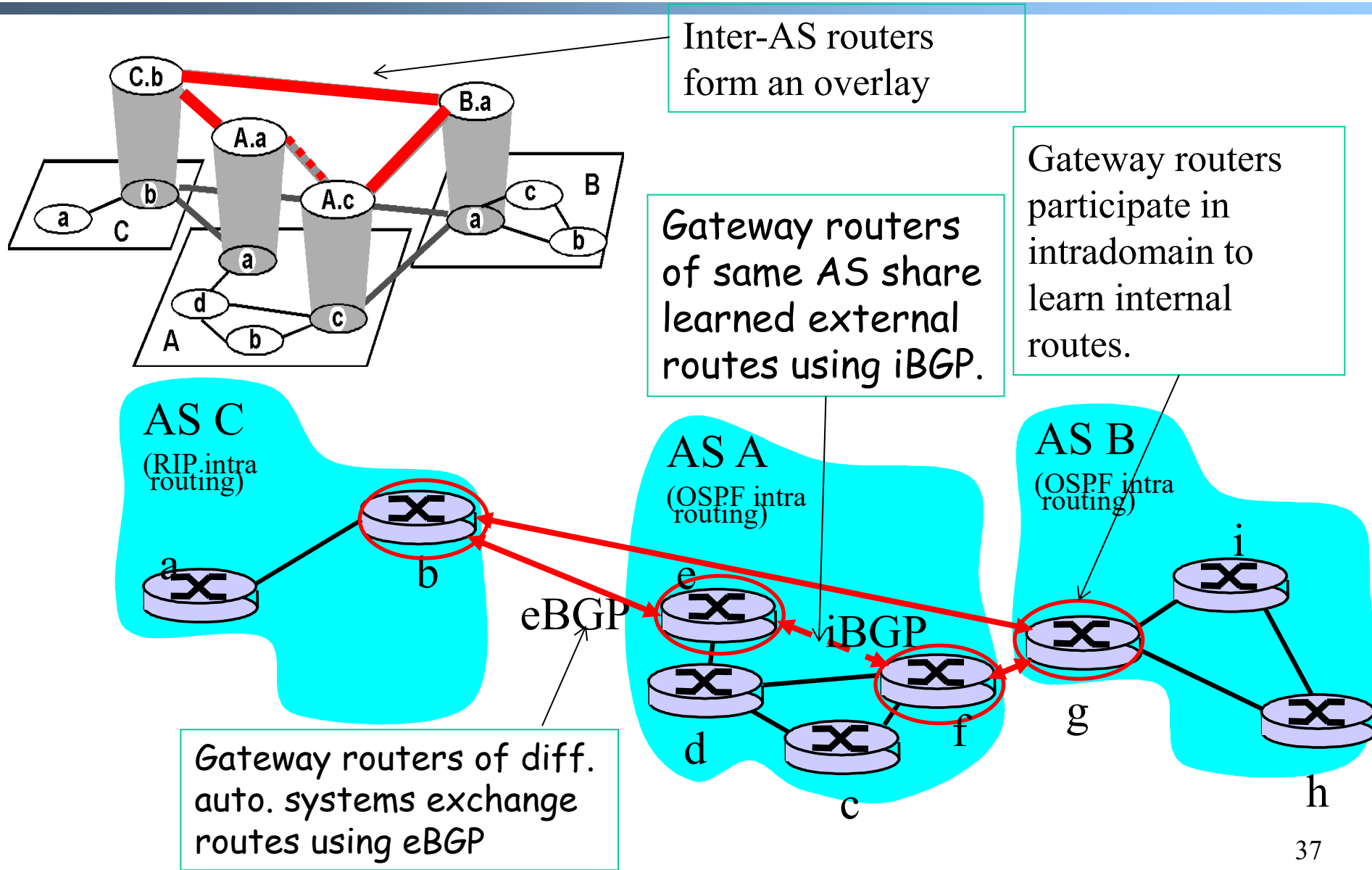
- Abstract each network as an autonomous system (AS), identified by an AS number (ASN)
- Conceptually the global routing graph consists of only autonomous systems as nodes



Routing with Autonomous Systems

- ❑ Internet routing is divided into intra-AS routing and inter-AS routing
 - Intra-AS routing (also called intradomain routing)
 - A protocol running inside an AS is called an Interior Gateway Protocol (IGP), each AS can choose its own protocol, such as RIP, E/IGRP, OSPF, IS-IS
 - Inter-AS routing (also called interdomain routing)
 - A protocol runs among autonomous systems is also called an Exterior Gateway Protocol (EGP)
 - The de facto EGP protocol is BGP

Routing with Autonomous Systems



Summary: Internet Routing Architecture

- ❑ Autonomous systems have flexibility to choose their own intradomain routing protocols
 - allows autonomy
- ❑ Only a small # of routers (gateways) from each AS in the interdomain level
 - improves scalability
- ❑ Interdomain routing using AS topology instead of detailed topology
 - improves scalability/privacy

Outline

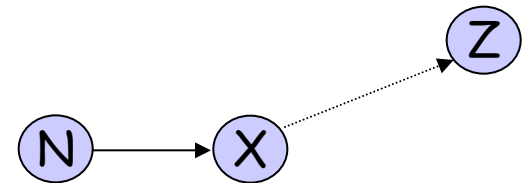
- ❑ Admin and recap
- ❑ Network control plane
 - Routing
 - Link weights assignment
 - Routing computation
 - Basic routing computation protocols
 - **Global Internet routing**
 - Basic architecture
 - *BGP (Border Gateway Protocol): The de facto Inter-domain routing standard*

BGP Basic Operations

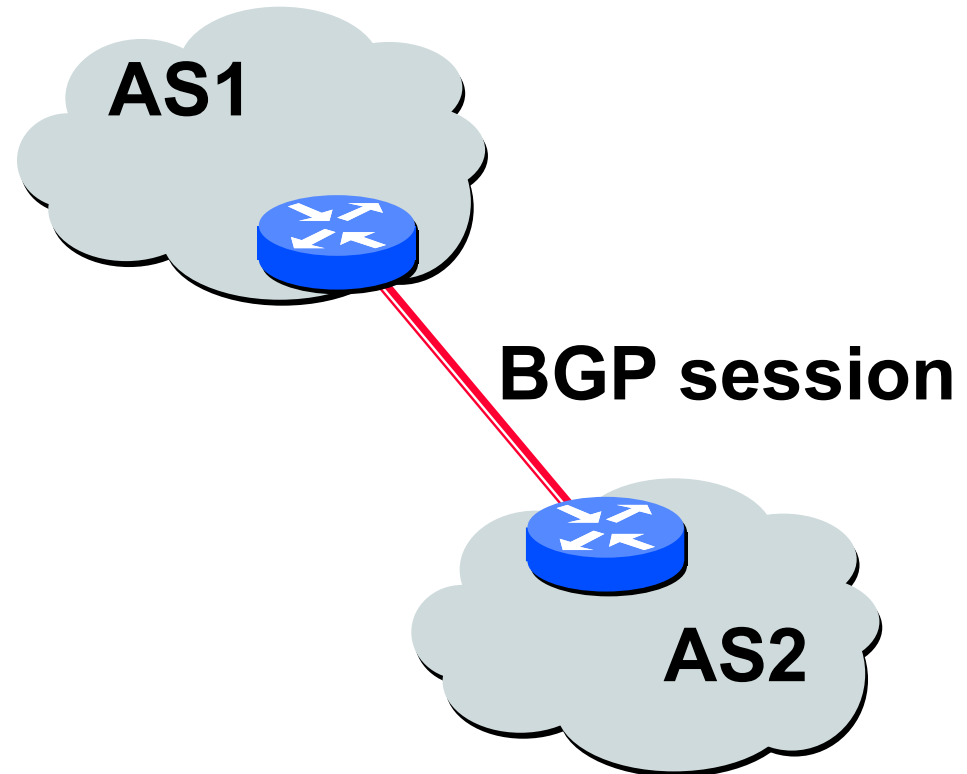
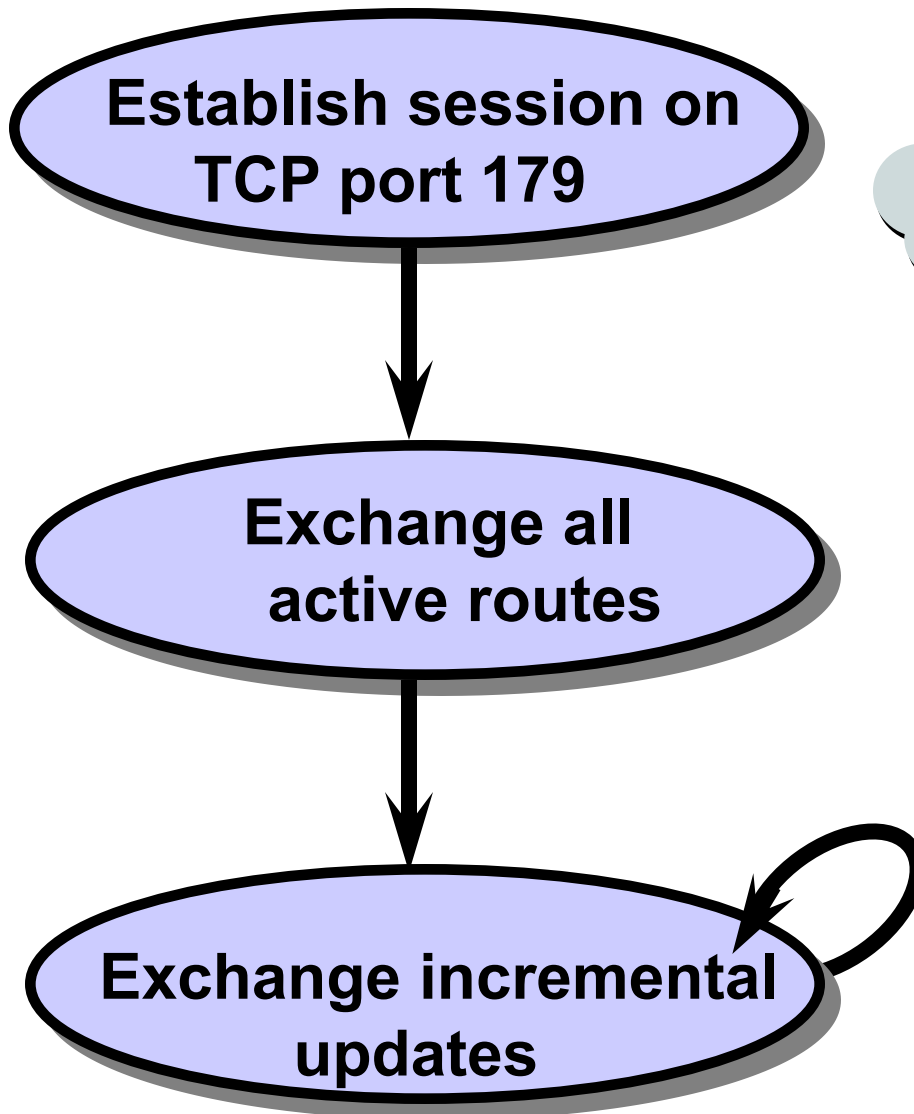
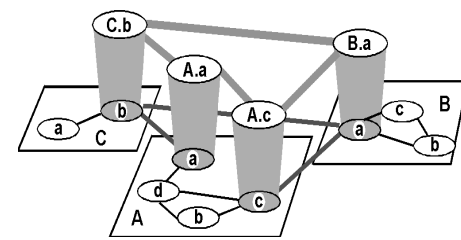
- BGP is a **Path Vector** protocol
 - similar to Distance Vector protocol
 - a border gateway sends to a neighbor *entire path* (i.e., *a sequence of ASNs*) to a destination, e.g.,
 - gateway X sends to neighbor N its path to dest. Z:

$$\text{path}(X,Z) = X, Y_1, Y_2, Y_3, \dots, Z$$

- if N selects $\text{path}(X, Z)$ advertised by X, then:
 $\text{path}(N,Z) = N, \text{path}(X,Z)$



BGP Basic Operations



while (connection is ALIVE)
exchange UPDATE message
select best available route
if route changes, export to neigh.

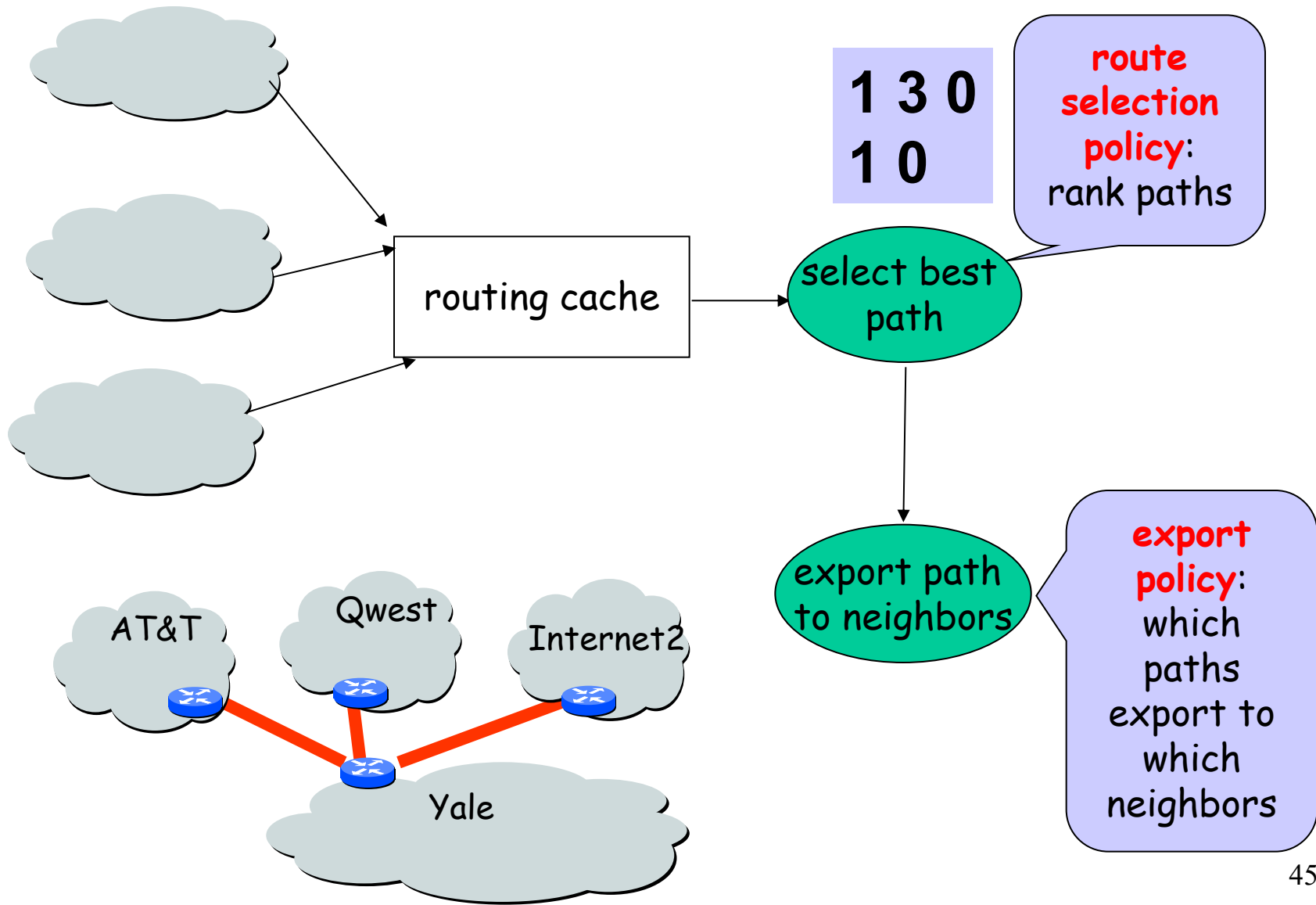
BGP Messages

- ❑ Four types of 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

Outline

- ❑ Admin and recap
- ❑ Network control plane
 - Routing
 - Link weights assignment
 - Routing computation
 - Basic routing computation protocols
 - Global Internet routing
 - Basic architecture
 - *BGP (Border Gateway Protocol): The de facto Inter-domain routing standard*
 - Basic operations
 - *BGP as a policy routing framework (control interdomain routes)*

BGP Policy Routing Framework: Decision Components



BGP Example (1)

