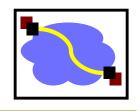


14-736 Distributed Systems

Lecture 26:

Internet Routing as a DS Problem

Outline

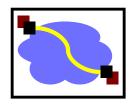


Intro

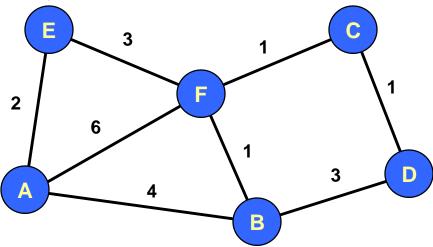
Distance Vector

Link State

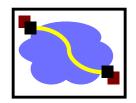
Graph Model



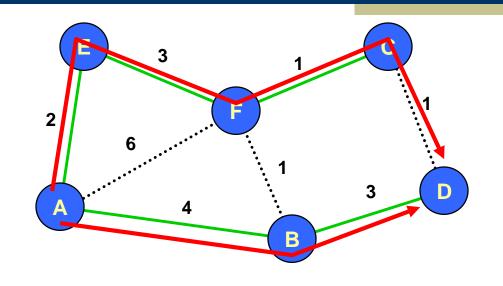
- Represent each router as node
- Direct link between routers represented by edge
 - Symmetric links ⇒ undirected graph
- Edge "cost" c(x,y) denotes measure of difficulty of using link
 - delay, \$ cost, or congestion level
- Task
 - Determine least cost path from every node to every other node
 - Path cost d(x,y) = sum of link costs



Routes from Node A



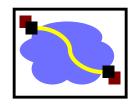
Forwarding Table for A						
Dest	Cost	Next Hop				
Α	0	А				
В	4	В				
С	6	Е				
D	7	В				
Е	2	Е				
F	5	E				



Properties

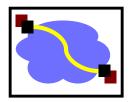
- Some set of shortest paths forms tree
 - Shortest path spanning tree
- Solution not unique
 - E.g., A-E-F-C-D also has cost 7

Key Challenges



- This would be easy if the graph didn't change
- But, what happens as links are removed and are added?
 - Links fail
 - Links heal
- Internet routing is basically a big distributed agreement problem
 - If routers agree on the graph, they can make consistent decisions
 - If routers disagree on the graph, badness can happen
 - Cycles are a notable problem
 - Lost packets are another (time out)

Ways to Compute Shortest Paths



Centralized

- Collect graph structure in one place
- Use standard graph algorithm
- Disseminate routing tables

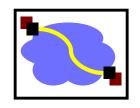
Link-state

- Every node collects complete graph structure
- Each computes shortest paths from it
- Each generates own routing table

Distance-vector

- No one has copy of graph
- Nodes construct their own tables iteratively
- Each sends information about its table to neighbors

Outline

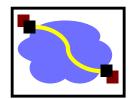


Intro

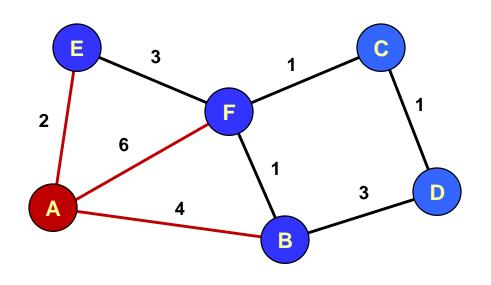
Distance Vector

Link State

Distance-Vector Method



Initial Table for A							
Dest	Cost	Next Hop					
Α	0	Α					
В	4	В					
С	8	_					
D	8	_					
Е	2	Е					
F	6	F					



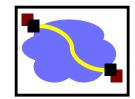
Idea

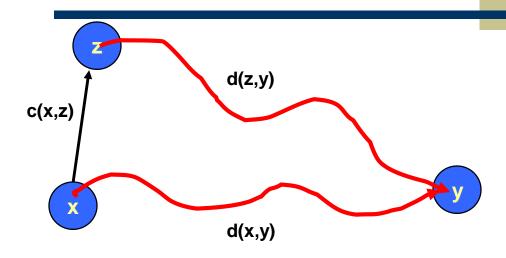
- At any time, have cost/next hop of best known path to destination
- Use cost ∞ when no path known

Initially

Only have entries for directly connected nodes

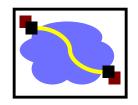
Distance-Vector Update





Update(x,y,z)
 d ← c(x,z) + d(z,y) # Cost of path from x to y with first hop z
 if d < d(x,y)
 # Found better path
 return d,z # Updated cost / next hop
 else
 return d(x,y), nexthop(x,y) # Existing cost / next hop

Algorithm



- Bellman-Ford algorithm
- Repeat

```
For every node x

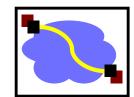
For every neighbor z

For every destination y

d(x,y) ← Update(x,y,z)
```

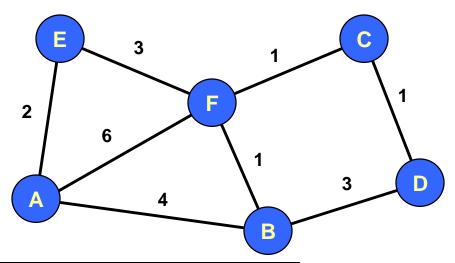
Until converge

Start



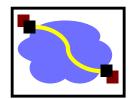
Optimum 1-hop paths

Ta	able for	Α	Table for B			
Dst	Cst	Нор	Dst	Cst	Нор	
Α	0	Α	Α	4	Α	
В	4	В	В	0	В	
С	∞	_	С	8	-	
D	∞	-	D	3	D	
E	2	Е	Е	8	_	
F	6	F	F	1	F	
				_	-	



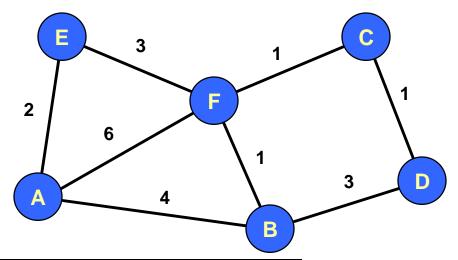
T	able fo	or C	Ta	Table for D		Table for E		Table for F			
Dst	Cst	Нор	Dst	Cst	Нор	Dst	Cst	Нор	Dst	Cst	Нор
Α	8	1	Α	8	_	Α	2	Α	Α	6	Α
В	∞	1	В	3	В	В	∞	_	В	1	В
С	0	С	C	1	С	С	∞	_	С	1	С
D	1	D	D	0	D	D	∞	_	D	∞	_
Е	∞	_	Е	8	_	Е	0	Е	Е	3	Е
F	1	F	F	8	_	F	3	F	F	0	F

Iteration #1



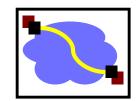
Optimum 2-hop paths

Ta	able for	А	Table for B			
Dst	Cst	Нор	Dst	Cst	Нор	
Α	0	Α	Α	4	Α	
В	4	В	В	0	В	
С	7	F	С	2	F	
D	7	В	D	3	D	
Е	2	Е	Е	4	F	
F	5	Е	F	1	F	



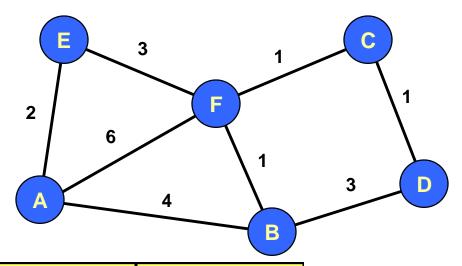
Та	able for	С	Table for D		Table for E			Table for F			
Dst	Cst	Нор	Dst	Cst	Нор	Dst	Cst	Нор	Dst	Cst	Нор
Α	7	F	Α	7	В	Α	2	Α	Α	5	В
В	2	F	В	3	В	В	4	F	В	1	В
С	0	С	С	1	С	С	4	F	С	1	С
D	1	D	D	0	D	D	∞	_	D	2	С
Е	4	F	Е	∞	_	Е	0	Е	Е	3	Е
F	1	F	F	2	С	F	3	F	F	0	F

Iteration #2



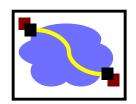
Optimum 3-hop paths

Ta	able for	Α	Table for B			
Dst	Cst	Нор	Dst	Cst	Нор	
Α	0	Α	Α	4	Α	
В	4	В	В	0	В	
С	6	Е	С	2	F	
D	7	В	D	3	D	
Е	2	Е	Е	4	F	
F	5	Е	F	1	F	



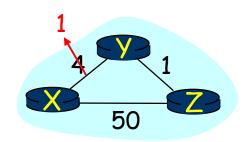
Ta	able for	С	Ta	able for	D	Table for E		Table for F			
Dst	Cst	Нор	Dst	Cst	Нор	Dst	Cst	Нор	Dst	Cst	Нор
Α	6	F	Α	7	В	Α	2	Α	Α	5	В
В	2	F	В	3	В	В	4	F	В	1	В
С	0	С	O	1	С	C	4	F	C	1	С
D	1	D	D	0	D	D	5	F	D	2	С
Е	4	F	Ш	5	С	Е	0	Е	Е	3	Е
F	1	F	F	2	С	F	3	F	F	0	F

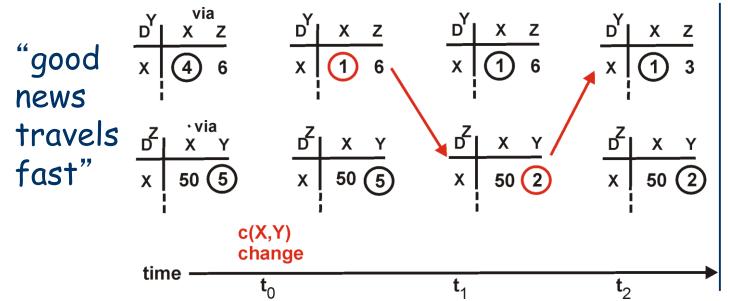
Distance Vector: Link Cost Changes



Link cost changes:

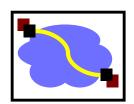
- Node detects local link cost change
- Updates distance table
- If cost change in least cost path, notify neighbors





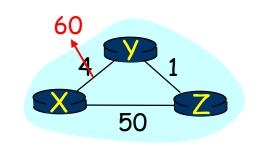
algorithm terminates

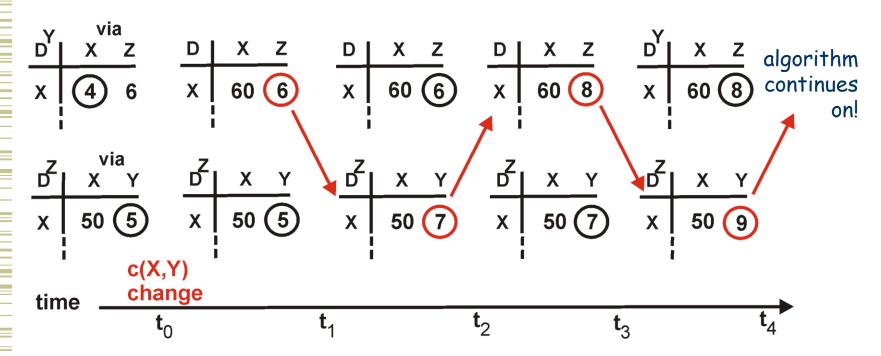
Distance Vector: Link Cost Changes



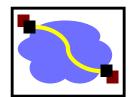
Link cost changes:

- Good news travels fast
- Bad news travels slow -"count to infinity" problem!



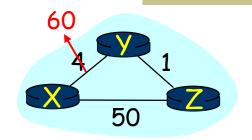


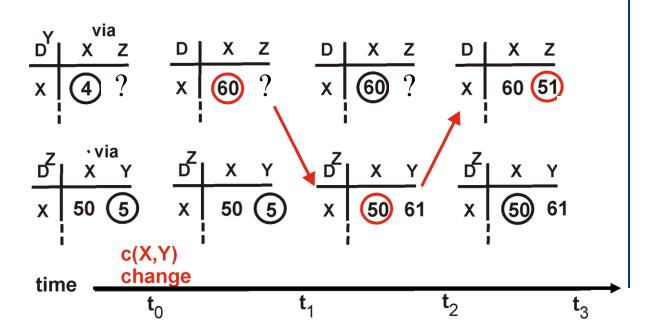
Distance Vector: Split Horizon



If Z routes through Y to get to X:

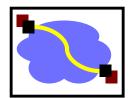
Z does not advertise its route to X back to Y





algorithm terminates

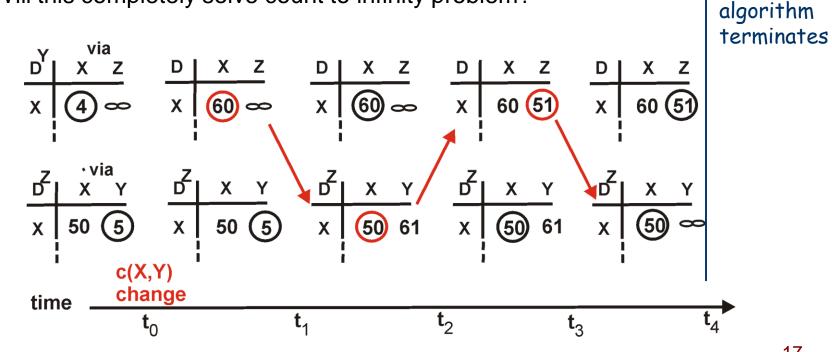
Distance Vector: Poison Reverse



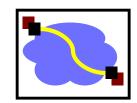
If Z routes through Y to get to X:

- Z tells Y its (Z's) distance to X is infinite (so Y won't route to X via Z)
- Immediate notification of unreachability, rather than split horizon timeout waiting for advertisement
- 60 50

Will this completely solve count to infinity problem?



Poison Reverse Failures



Ta	able for	Α	Table for B		Table for D			Table for F			
Dst	Cst	Нор	Dst	Cst	Нор	Dst	Cst	Нор	Dst	Cst	Нор
С	7	F	С	8	Α	С	9	В	С	1	С

Table for A Cst Hop **Dst** С

Table for A Hop Dst Cst С 13

Better Route

Forced

Update

Forced Update

Table for B							
Dst	Dst Cst Ho						
С	14	Α					

Forced Update

Table for D					
Dst	Cst	Нор			
С	15	В			

Table for A Cst Hop Dst С 19 D

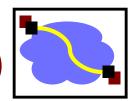
Forced Update **Forced Update**

Table for F							
Dst	Dst Cst						
С	8	1					

A В

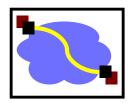
- Iterations don't converge
- "Count to infinity"
- Solution
 - Make "infinity" smaller
 - What is upper bound on maximum path length?

Routing Information Protocol (RIP)



- Earliest IP routing protocol (1982 BSD)
 - Current standard is version 2 (RFC 1723)
- Features
 - Every link has cost 1
 - "Infinity" = 16
 - Limits to networks where everything reachable within 15 hops
- Sending Updates
 - Every router listens for updates on UDP port 520
 - RIP message can contain entries for up to 25 table entries

RIP Updates



Initial

- When router first starts, asks for copy of table for every neighbor
- Uses it to iteratively generate own table

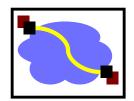
Periodic

- Every 30 seconds, router sends copy of its table to each neighbor
- Neighbors use it to iteratively update their tables

Triggered

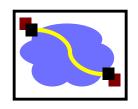
- When every entry changes, send copy of entry to neighbors
 - Except for one causing update (split horizon rule)
- Neighbors use it to update their tables

RIP Staleness / Oscillation Control



- Small Infinity
 - Count to infinity doesnt take very long
- Route Timer
 - Every route has timeout limit of 180 seconds
 - Reached when haven't received update from next hop for 6 periods
 - If not updated, set to infinity
 - Soft-state refresh → important concept!
- Behavior
 - When router or link fails, can take minutes to stabilize

Outline

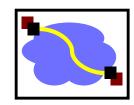


Intro

Distance Vector

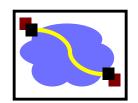
Link State

Link State Protocol Concept

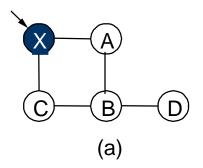


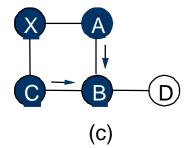
- Every node gets complete copy of graph
 - Every node "floods" network with data about its outgoing links
- Every node computes routes to every other node
 - Using single-source, shortest-path algorithm
- Process performed whenever needed
 - When connections die / reappear

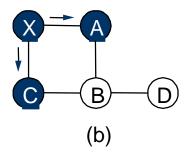
Sending Link States by Flooding

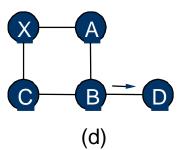


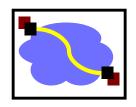
- X Wants to Send Information
 - Sends on all outgoing links
- When Node Y Receives Information from Z
 - Send on all links other than Z





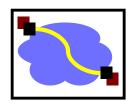


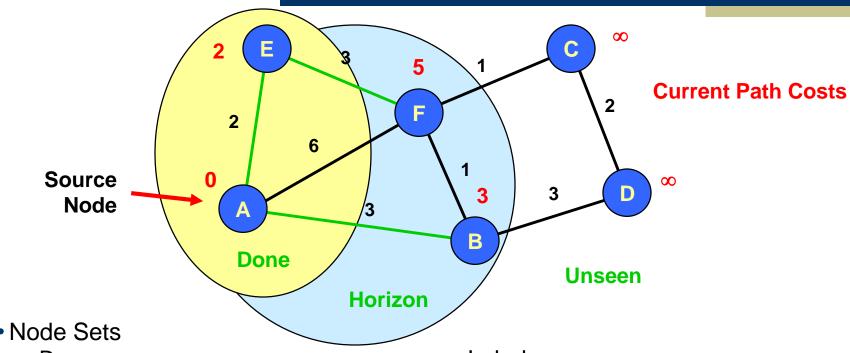




- Given
 - Graph with source node s and edge costs c(u,v)
 - Determine least cost path from s to every node v
- Shortest Path First Algorithm
 - Traverse graph in order of least cost from source

Dijkstra's Algorithm: Concept

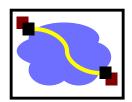


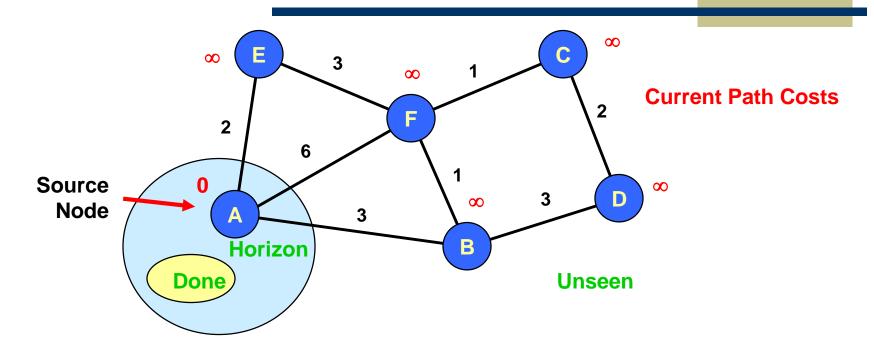


- Done
 - Already have least cost path to it
- Horizon:
 - Reachable in 1 hop from node in Done
- Unseen:
 - Cannot reach directly from node in Done

- Label
 - d(v) = path cost from s to v
- Path
 - Keep track of last link in path

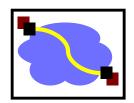
Dijkstra's Algorithm: Initially

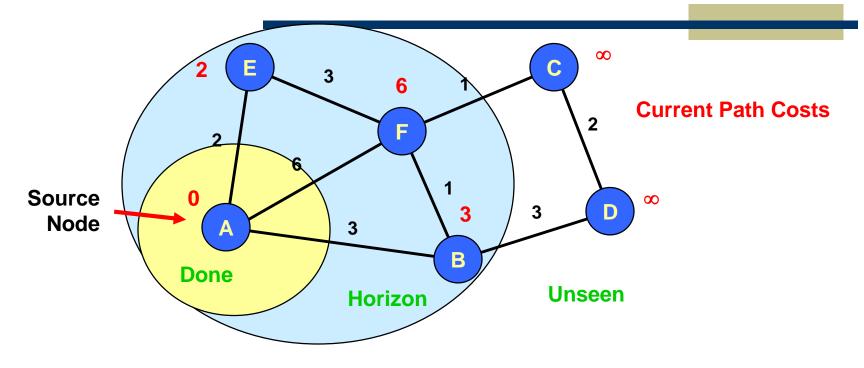




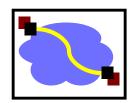
- No nodes done
- Source in horizon

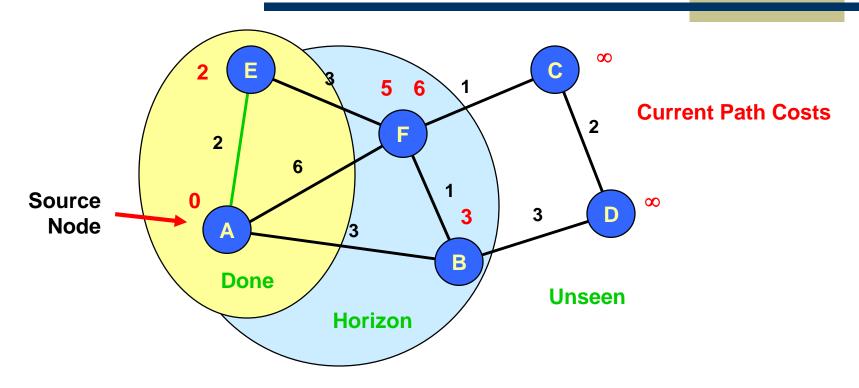
Dijkstra's Algorithm: Initially



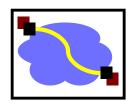


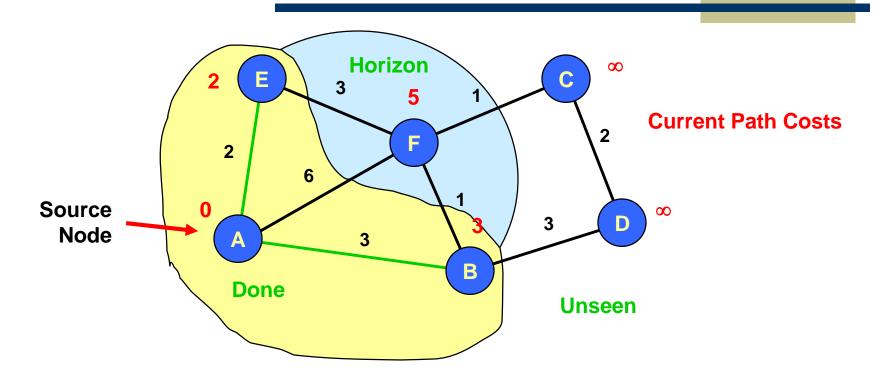
- d(v) to node A shown in red
 - Only consider links from done nodes



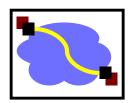


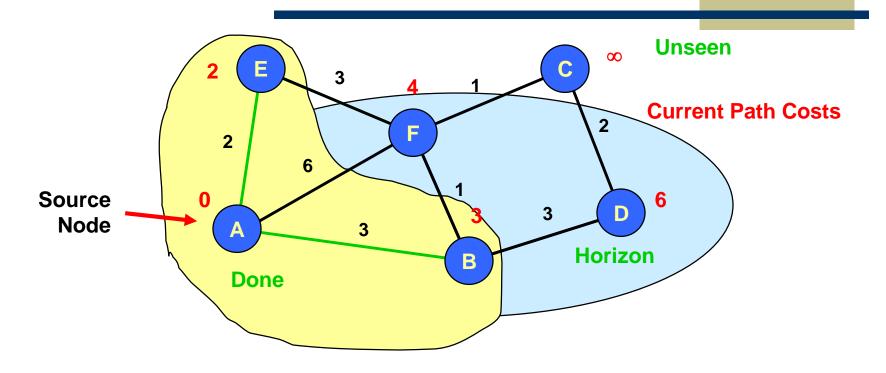
- Select node v in horizon with minimum d(v)
- Add link used to add node to shortest path tree
- Update d(v) information



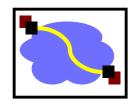


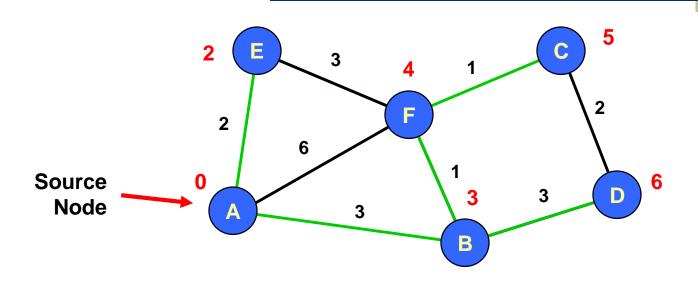
Repeat...





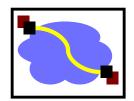
- Update d(v) values
 - Can cause addition of new nodes to horizon





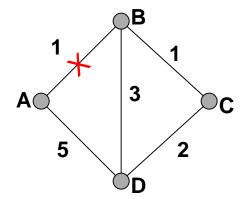
Final tree shown in green

Link State Characteristics



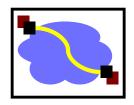
- With consistent LSDBs*, all nodes compute consistent loop-free paths
- Can still have transient loops

*Link State Data Base



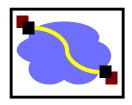
Packet from C→A
may loop around BDC
if B knows about failure
and C & D do not

OSPF Routing Protocol



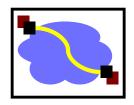
- Open
 - Open standard created by IETF
- Shortest-path first
 - Another name for Dijkstra's algorithm
- More prevalent than RIP

OSPF Reliable Flooding



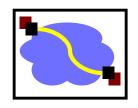
- Transmit link state advertisements
 - Originating router
 - Typically, minimum IP address for router
 - Link ID
 - ID of router at other end of link
 - Metric
 - Cost of link
 - Link-state age
 - Incremented each second
 - Packet expires when reaches 3600
 - Sequence number
 - Incremented each time sending new link information

OSPF Flooding Operation



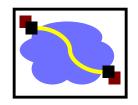
- Node X Receives LSA from Node Y
 - With Sequence Number q
 - Looks for entry with same origin/link ID
- Cases
 - No entry present
 - Add entry, propagate to all neighbors other than Y
 - Entry present with sequence number p < q
 - Update entry, propagate to all neighbors other than Y
 - Entry present with sequence number p > q
 - Send entry back to Y
 - To tell Y that it has out-of-date information
 - Entry present with sequence number p = q
 - Ignore it

Flooding Issues



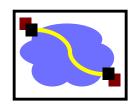
- When should it be performed
 - Periodically
 - When status of link changes
 - Detected by connected node
- What happens when router goes down & back up
 - Sequence number reset to 0
 - Other routers may have entries with higher sequence numbers
 - Router will send out LSAs with number 0
 - Will get back LSAs with last valid sequence number p
 - Router sets sequence number to p+1 & resends

Adoption of OSPF



- RIP viewed as outmoded
 - Good when networks small and routers had limited memory & computational power
- OSPF Advantages
 - Fast convergence when configuration changes

Comparison of LS and DV Algorithms



Message complexity

- LS: with n nodes, E links, O(nE) messages
- <u>DV:</u> exchange between neighbors only

Speed of Convergence

- <u>LS:</u> Relatively fast
 - Complex computation, but can forward before computation
 - may have transient loops
- <u>DV</u>: convergence time varies
 - may have routing loops
 - count-to-infinity problem
 - faster with triggered updates

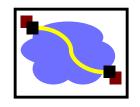
Space requirements:

- LS maintains entire topology
- DV maintains only neighbor state

Robustness: router malfunctions

- <u>LS</u>: Node can advertise incorrect link cost
 - Each node computes its own table
- DV: Node can advertise incorrect path cost
 - Each node's table used by others (error propagates)

Outline



Distance Vector

Link State

Routing Hierarchy