

# Lecture 12 Routing

Peter Steenkiste  
Departments of Computer Science and  
Electrical and Computer Engineering  
Carnegie Mellon University

15-441 Networking, Spring 2006  
<http://www.cs.cmu.edu/~prs/15-441>

Peter A. Steenkiste, SCS, CMU

## Packet Forwarding

- Use IP address in the header to determine where to forward the packet
  - » Longest prefix match in the forwarding table
  - » Send packet out appropriate port
- Today's lecture: how to create and manage the forwarding table
  - » Focus on intra-domain routing
  - » Route selection is based on the optimization of a routing metric

Peter A. Steenkiste, SCS, CMU

## Outline

- Distance vector routing
- RIP
- Link state routing
- OSPF

Peter A. Steenkiste, SCS, CMU

## Abstraction: Represent Network as a Graph

- Represent each router as node
- Direct link between routers represented by edge
  - » Symmetric links  $\Rightarrow$  undirected graph
- Edge "cost"  $c(x,y)$  denotes measure of "cost" of using link
  - » delay, \$ cost, or congestion level
- Must determine least cost path for every node pair
  - Path cost  $d(x,y)$  = sum of link costs

Peter A. Steenkiste, SCS, CMU

## Routes from Node A

Dest	Cost	Next Hop
A	0	A
B	4	B
C	6	E
D	7	B
E	2	E
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

Peter A. Steenkiste, SCS, CMU

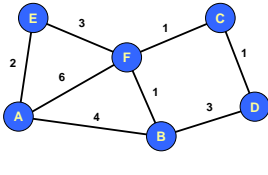
## Classes of Routing Solutions

- Centralized
  - » Collect graph structure in one place
  - » Use standard graph algorithm Does this work?
  - » Disseminate routing tables
- Partially Distributed
  - » Every node collects complete graph structure
  - » Each locally computes shortest paths from it
  - » Each generates own routing table
  - » "Link-state" algorithm
- Fully Distributed
  - » No one has copy of graph
  - » Nodes construct their own tables iteratively
  - » Each sends information about its table to neighbors
  - » "Distance-Vector" algorithm

Peter A. Steenkiste, SCS, CMU

## Distance-Vector Method

Initial Table for A		
Dest	Cost	Next Hop
A	0	A
B	4	B
C	$\infty$	-
D	$\infty$	-
E	2	E
F	6	F

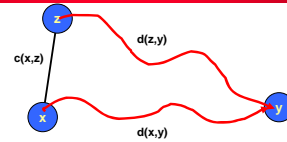


- Idea is to iteratively improve least-cost path to any destination based on information learned from neighbors
  - Table lists cost/next hop of best known path to destination
  - Initially table only has entries for directly connected nodes

Peter A. Steenkiste, SCS, CMU

7

## Distance-Vector Update



- Update( $x,y,z$ )
  - $d \leftarrow c(x,z) + d(z,y)$  # 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  $d(x,y), \text{nextHop}(x,y)$  # Existing cost / next hop

Peter A. Steenkiste, SCS, CMU

8

## Algorithm

- Bellman-Ford algorithm
- Repeat
  - For every node x
    - For every neighbor z
      - For every destination y
        - $d(x,y) \leftarrow \text{Update}(x,y,z)$
- Until Converge
- Nodes x can run the algorithm in fully distributed fashion

Peter A. Steenkiste, SCS, CMU

9

## Start

### Optimum 1-hop paths

Table for A			Table for B		
Dst	Cst	Hop	Dst	Cst	Hop
A	0	A	A	4	A
B	4	B	B	0	B
C	$\infty$	-	C	$\infty$	-
D	$\infty$	-	D	3	D
E	2	E	E	$\infty$	-
F	6	F	F	1	F

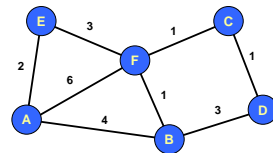


Table for C			Table for D			Table for E			Table for F		
Dst	Cst	Hop	Dst	Cst	Hop	Dst	Cst	Hop	Dst	Cst	Hop
A	$\infty$	-	A	$\infty$	-	A	2	A	A	6	A
B	$\infty$	-	B	3	B	B	$\infty$	-	B	1	B
C	0	C	C	1	C	C	$\infty$	-	C	1	C
D	1	D	D	0	D	D	$\infty$	-	D	$\infty$	-
E	$\infty$	-	E	$\infty$	-	E	0	E	E	3	E
F	1	F	F	$\infty$	-	F	3	F	F	0	F

Peter A. Steenkiste, SCS, CMU

10

## Iteration #1

### Optimum 2-hop paths

Table for A			Table for B		
Dst	Cst	Hop	Dst	Cst	Hop
A	0	A	A	4	A
B	4	B	B	0	B
C	7	F	C	2	F
D	7	B	D	3	D
E	2	E	E	4	F
F	5	E	F	1	F

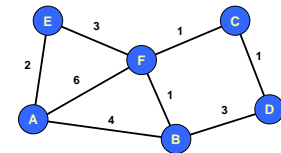


Table for C			Table for D			Table for E			Table for F		
Dst	Cst	Hop	Dst	Cst	Hop	Dst	Cst	Hop	Dst	Cst	Hop
A	7	F	A	7	B	A	2	A	A	5	B
B	2	F	B	3	B	B	4	F	B	1	B
C	0	C	C	1	C	C	4	F	C	1	C
D	1	D	D	0	D	D	$\infty$	-	D	2	C
E	4	F	E	$\infty$	-	E	0	E	E	3	E
F	1	F	F	2	C	F	3	F	F	0	F

Peter A. Steenkiste, SCS, CMU

11

## Iteration #2

### Optimum 3-hop paths

Table for A			Table for B		
Dst	Cst	Hop	Dst	Cst	Hop
A	0	A	A	4	A
B	4	B	B	0	B
C	6	E	C	2	F
D	7	B	D	3	D
E	2	E	E	4	F
F	5	E	F	1	F

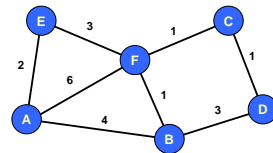


Table for C			Table for D			Table for E			Table for F		
Dst	Cst	Hop	Dst	Cst	Hop	Dst	Cst	Hop	Dst	Cst	Hop
A	6	F	A	7	B	A	2	A	A	5	B
B	2	F	B	3	B	B	4	F	B	1	B
C	0	C	C	1	C	C	4	F	C	1	C
D	1	D	D	0	D	D	5	F	D	2	C
E	4	F	E	5	C	E	0	E	E	3	E
F	1	F	F	2	C	F	3	F	F	0	F

Peter A. Steenkiste, SCS, CMU

12

### 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

**"good news travels fast"**

$\begin{matrix} D & X & Z \\ X & 4 & 6 \end{matrix}$	$\begin{matrix} D & X & Z \\ X & 1 & 6 \end{matrix}$	$\begin{matrix} D & X & Z \\ X & 1 & 6 \end{matrix}$	$\begin{matrix} D & X & Z \\ X & 1 & 3 \end{matrix}$
$\begin{matrix} D & X & Y \\ X & 50 & 5 \end{matrix}$	$\begin{matrix} D & X & Y \\ X & 50 & 5 \end{matrix}$	$\begin{matrix} D & X & Y \\ X & 50 & 2 \end{matrix}$	$\begin{matrix} D & X & Y \\ X & 50 & 2 \end{matrix}$

time  $t_0 \quad t_1 \quad t_2$

algorithm terminates

13

### Distance Vector: Link Cost Changes

**Link cost changes:**

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

$\begin{matrix} D & X & Z \\ X & 4 & 6 \end{matrix}$	$\begin{matrix} D & X & Z \\ X & 60 & 6 \end{matrix}$	$\begin{matrix} D & X & Z \\ X & 60 & 6 \end{matrix}$	$\begin{matrix} D & X & Z \\ X & 60 & 8 \end{matrix}$	$\begin{matrix} D & X & Z \\ X & 60 & 8 \end{matrix}$
$\begin{matrix} D & X & Y \\ X & 50 & 5 \end{matrix}$	$\begin{matrix} D & X & Y \\ X & 50 & 5 \end{matrix}$	$\begin{matrix} D & X & Y \\ X & 50 & 7 \end{matrix}$	$\begin{matrix} D & X & Y \\ X & 50 & 7 \end{matrix}$	$\begin{matrix} D & X & Y \\ X & 50 & 9 \end{matrix}$

time  $t_0 \quad t_1 \quad t_2 \quad t_3 \quad t_4$

algorithm continues on!

14

### Distance Vector: Split Horizon

**If Z routes through Y to get to X:**

- Z does not advertise its route to X back to Y

$\begin{matrix} D & X & Z \\ X & 4 & ? \end{matrix}$	$\begin{matrix} D & X & Z \\ X & 60 & ? \end{matrix}$	$\begin{matrix} D & X & Z \\ X & 60 & ? \end{matrix}$	$\begin{matrix} D & X & Z \\ X & 60 & 51 \end{matrix}$
$\begin{matrix} D & X & Y \\ X & 50 & 5 \end{matrix}$	$\begin{matrix} D & X & Y \\ X & 50 & 5 \end{matrix}$	$\begin{matrix} D & X & Y \\ X & 50 & 61 \end{matrix}$	$\begin{matrix} D & X & Y \\ X & 50 & 61 \end{matrix}$

time  $t_0 \quad t_1 \quad t_2 \quad t_3$

algorithm terminates

15

### 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)
- Will this completely solve count to infinity problem?

$\begin{matrix} D & X & Z \\ X & 4 & \infty \end{matrix}$	$\begin{matrix} D & X & Z \\ X & 60 & \infty \end{matrix}$	$\begin{matrix} D & X & Z \\ X & 60 & \infty \end{matrix}$	$\begin{matrix} D & X & Z \\ X & 60 & 51 \end{matrix}$	$\begin{matrix} D & X & Z \\ X & 60 & 51 \end{matrix}$
$\begin{matrix} D & X & Y \\ X & 50 & 5 \end{matrix}$	$\begin{matrix} D & X & Y \\ X & 50 & 5 \end{matrix}$	$\begin{matrix} D & X & Y \\ X & 50 & 61 \end{matrix}$	$\begin{matrix} D & X & Y \\ X & 50 & 61 \end{matrix}$	$\begin{matrix} D & X & Y \\ X & 50 & \infty \end{matrix}$

time  $t_0 \quad t_1 \quad t_2 \quad t_3 \quad t_4$

algorithm terminates

16

### Poison Reverse Failures

<b>Table for A</b>	<b>Table for B</b>	<b>Table for D</b>	<b>Table for F</b>
Dist Cost Hop	Dist Cost Hop	Dist Cost Hop	Dist Cost Hop
C 7 F C	C 8 A	C 9 B	C 1 C

**Forced Update**

<b>Table for A</b>	<b>Table for F</b>
Dist Cost Hop	Dist Cost Hop
C 13 D	C 1 C

**Better Route**

**Forced Update**

<b>Table for B</b>
Dist Cost Hop
C 14 A

**Forced Update**

<b>Table for D</b>
Dist Cost Hop
C 15 B

**Forced Update**

<b>Table for A</b>
Dist Cost Hop
C 19 D

**Forced Update**

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

17

### 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

18

## 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 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 to update their tables

Peter A. Steenkiste, SCS, CMU

19

## RIP Staleness / Oscillation Control

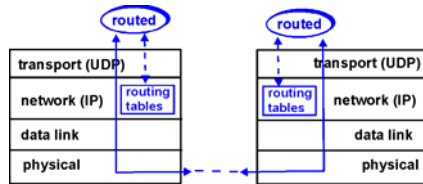
- **Small Infinity**
  - » Count to infinity doesn't take very long
- **Route Timer**
  - » Every route has timeout limit of 180 seconds
    - Reached when haven't received update from next hop for 6 periods
  - » If not updated, set to infinity
  - » Soft-state refresh → important concept!!!
- **Behavior**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

Peter A. Steenkiste, SCS, CMU

20

## RIP Table Processing

- RIP routing tables managed by application-level process called route-d (daemon)
- advertisements sent in UDP packets, periodically repeated



Peter A. Steenkiste, SCS, CMU

21

## Outline

- Distance vector routing
- RIP
- Link state routing
- OSPF

Peter A. Steenkiste, SCS, CMU

22

## Link State Protocol Concept

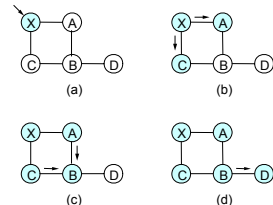
- **Every node gets complete copy of graph**
  - » Every node "floods" network with data about its outgoing links
- **Every node computes routes to every other node**
  - » Using single-source, shortest-path algorithm
- **Process performed whenever needed**
  - » When connections die / reappear

Peter A. Steenkiste, SCS, CMU

23

## 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
- **Sequence number used to suppress duplicates**



Peter A. Steenkiste, SCS, CMU

24

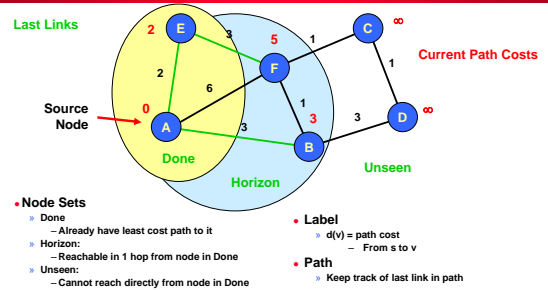
## Dijkstra's Algorithm

- **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

Peter A. Steenkiste, SCS, CMU

25

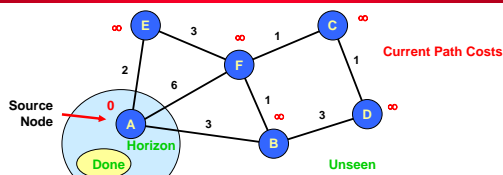
## Dijkstra's Algorithm: Concept



Peter A. Steenkiste, SCS, CMU

26

## Dijkstra's Algorithm: Initially

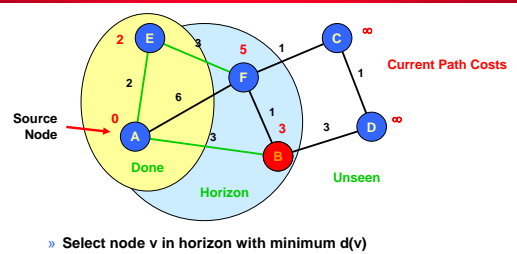


- **No nodes done**
- **Source in horizon**

Peter A. Steenkiste, SCS, CMU

27

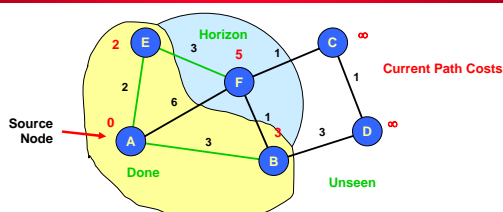
## Dijkstra's Algorithm: Selection



Peter A. Steenkiste, SCS, CMU

28

## Dijkstra's Algorithm: Selection

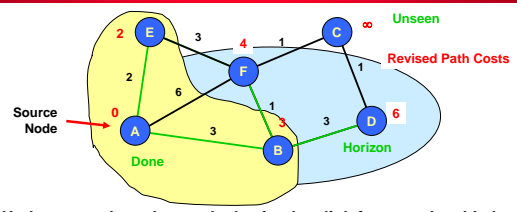


- » Add selected node to Done

Peter A. Steenkiste, SCS, CMU

29

## Dijkstra's Algorithm: Update



- **Update costs based on paths having last link from newly added Done node**
  - » Could change values of nodes in horizon
  - » Could add new nodes to horizon
- **Update link information**

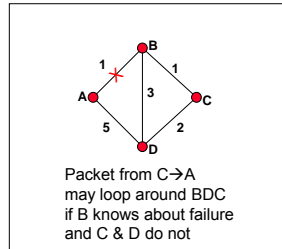
Peter A. Steenkiste, SCS, CMU

30

## Link State Characteristics

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

\*Link State Data Base



Peter A. Steenkiste, SCS, CMU

31

## OSPF Routing Protocol

- Open
  - » Open standard created by IETF
- Shortest-path first
  - » Another name for Dijkstra's algorithm
- RIP viewed as outmoded
  - » Good when networks small and routers had limited memory & computational power
- OSPF Advantages
  - » Fast convergence when configuration changes

Peter A. Steenkiste, SCS, CMU

32

## 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

Peter A. Steenkiste, SCS, CMU

33

## 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

Peter A. Steenkiste, SCS, CMU

34

## 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

Peter A. Steenkiste, SCS, CMU

35

## Comparison of LS and DV Algorithms

### Message complexity

- **LS:** with  $n$  nodes,  $E$  links,  $O(nE)$  messages
- **DV:** exchange between neighbors only

### Space requirements:

- » LS maintains entire topology
- » DV maintains only neighbor state

### Speed of Convergence

- **LS:** Complex computation
  - » But...can forward before computation
  - » may have oscillations
- **DV:** convergence time varies
  - » may be routing loops
  - » count-to-infinity problem
  - » (faster with triggered updates)

Peter A. Steenkiste, SCS, CMU

36

## Comparison of LS and DV Algorithms

**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
  - errors propagate thru network
- Other tradeoffs
  - Making LSP flood reliable