

Copyright © William C. Cheng

Forwarding V.S. Routing

- **Forwarding:** the process of moving packets from input to output based on:
 - = the forwarding table
 - = information in the packet
- **Routing:** process by which the forwarding table is built and maintained:
 - = one or more routing protocols
 - = procedures (algorithms) to convert routing info to forwarding table

Computer Communications - CSC1 551

Copyright © William C. Cheng

A Router And its Components

Cisco 7xxx Router

Computer Communications - CSC1 551

Copyright © William C. Cheng

Two Main Approaches

- **DV:** Distance-vector protocols
 - = you tell your neighbors what you know about everyone
- **LS:** Link state protocols
 - = you tell everyone about your neighbors

Computer Communications - CSC1 551

Copyright © William C. Cheng

CS551

Unicast Routing

Bill Cheng

<http://merlot.usc.edu/cs551-f12>

Computer Communications - CSC1 551

Copyright © William C. Cheng

Forwarding Examples

- **To forward unicast packets a router uses:**
 - = destination IP address
 - = longest matching prefix in forwarding table
- **To forward multicast packets:**
 - = source + destination IP address and incoming interface
 - = longest and exact match algorithms

Computer Communications - CSC1 551

Copyright © William C. Cheng

Factors Affecting Routing

Routing algorithms view the network as a graph

Problem: find lowest cost path between two nodes

➤ Factors

- = static: topology
- = dynamic: load
- = policy

Computer Communications - CSC1 551

Copyright © William C. Cheng

8

Distance Vector

- ↳ $D^X(Y,Z)$: distance to Y via Z in node X's distance table (Z is X's direct neighbor)
- ↳ $c(X,Z)$: cost from X to X's direct neighbor Z
- ↳ $D^X(Y,Z) = c(X,Z) + \min_w \{D^Z(Y,w)\}$ where w is a direct neighbor of Z

Computer Communications - CSCI 551

Copyright © William C. Cheng

10

Example - Initial Distances

Info at node	A	B	C	D	E
A	0	7	~	~	1
B	7	0	1	~	8
C	~	1	0	2	~
D	~	~	2	0	2
E	1	8	~	2	0

Computer Communications - CSCI 551

Copyright © William C. Cheng

12

E Receives D's Routes

Info at node	A	B	C	D	E
A	0	7	~	~	1
B	7	0	1	~	8
C	~	1	0	2	~
D	~	~	2	0	2
E	1	8	~	2	0

Computer Communications - CSCI 551

Copyright © William C. Cheng

7

Distance Vector Protocols

- ↳ Employed in the early Arpanet
- ↳ Distributed next hop computation
- ↳ adaptive
- ↳ Asynchronous, iterative
- ↳ Unit of information exchange
- ↳ vector of distances to destinations
- ↳ Distributed Bellman-Ford Algorithm

Computer Communications - CSCI 551

Copyright © William C. Cheng

9

Distributed Bellman-Ford

- ↳ Start Conditions:
 - ↳ each router starts with a vector of distances to all directly attached networks
- ↳ Send step:
 - ↳ each router advertises its current vector to all neighbouring routers
- ↳ Receive step:
 - ↳ upon receiving vectors from each of its neighbors, router computes its own **distance** to each neighbor
 - ↳ then, for every network X, router finds that neighbor who is closer to X than to any other neighbor
 - ↳ router updates its cost to X. After doing this for all X, router goes to send step if routing information has changed

Computer Communications - CSCI 551

Copyright © William C. Cheng

11

Example - Initial Distances

Info at node	A	B	C	D	E
A	0	7	~	~	1
B	7	0	1	~	8
C	~	1	0	2	~
D	~	~	2	0	2
E	1	8	~	2	0

Computer Communications - CSCI 551

Copyright © William C. Cheng

14

A Receives B's Routes

Info at node	A	B	C	D	E
A	0	7	~	~	1
B	7	0	1	~	8
C	~	1	0	2	~
D	~	~	2	0	2
E	1	8	4	2	0

$c(A,B) = 7$

Computer Communications - CSCI 551

Copyright © William C. Cheng

16

A Receives E's Routes

Info at node	A	B	C	D	E
A	0	7	8	~	1
B	7	0	1	~	8
C	~	1	0	2	~
D	~	~	2	0	2
E	1	8	4	2	0

$c(A,E) = 1$

Computer Communications - CSCI 551

Copyright © William C. Cheng

18

Final Distances

Info at node	A	B	C	D	E
A	0	6	5	3	1
B	6	0	1	3	5
C	5	1	0	2	4
D	3	3	2	0	2
E	1	5	4	2	0

Computer Communications - CSCI 551

Copyright © William C. Cheng

13

E Updates Cost to C

Info at node	A	B	C	D	E
A	0	7	~	~	1
B	7	0	1	~	8
C	~	1	0	2	~
D	~	~	2	0	2
E	1	8	4	2	0

Computer Communications - CSCI 551

Copyright © William C. Cheng

15

A Updates Cost to C

Info at node	A	B	C	D	E
A	0	7	8	~	1
B	7	0	1	~	8
C	~	1	0	2	~
D	~	~	2	0	2
E	1	8	4	2	0

Computer Communications - CSCI 551

Copyright © William C. Cheng

17

A Updates Cost to C & D

Info at node	A	B	C	D	E
A	0	7	5	3	1
B	7	0	1	~	8
C	~	1	0	2	~
D	~	~	2	0	2
E	1	8	4	2	0

Computer Communications - CSCI 551

Copyright © William C. Cheng

24

The Bouncing Effect (a.k.a. Count to Infinity Problem)

dest	np	A	B
A	1	51	2
B	1	51	1

dest	np	A	B	C
A	1	3	1	1
B	1	3	1	1
C	1	3	1	1

dest	np	A	B	C
A	1	51	2	1
B	1	51	1	1

Computer Communications - CSC1 551

Copyright © William C. Cheng

22

Final Distances After Link Failure (Cont..)

Info at node	A	B	C	D	E
A	0	7	8	10	1
B	7	0	1	3	8
C	8	1	0	2	9
D	10	3	2	0	11
E	1	8	9	11	0

dest	A	B	D
A	1	15	5
B	8	5	5
C	9	4	4
D	11	11	2

E's routing table:

Next hop	A	B	D
A	1	15	5
B	8	5	5
C	9	4	4
D	11	11	2

Computer Communications - CSC1 551

Copyright © William C. Cheng

20

Final Distances After Link Failure

Info at node	A	B	C	D	E
A	0	6	5	3	1
B	6	0	1	3	5
C	5	1	0	2	4
D	3	3	2	0	2
E	1	5	4	2	0

dest	A	B	D
A	1	14	5
B	7	8	5
C	6	9	4
D	4	11	2

E's routing table:

Next hop	A	B	D
A	1	14	5
B	7	8	5
C	6	9	4
D	4	11	2

Computer Communications - CSC1 551

Copyright © William C. Cheng

23

The Bouncing Effect (a.k.a. Count to Infinity Problem)

dest	np	A	B
A	1	51	2
B	1	51	1

dest	np	A	B	C
A	1	3	1	1
B	1	3	1	1
C	1	3	1	1

dest	np	A	B	C
A	1	51	2	1
B	1	51	1	1

Computer Communications - CSC1 551

Copyright © William C. Cheng

21

Final Distances After Link Failure (Cont..)

Info at node	A	B	C	D	E
A	0	6	5	3	1
B	6	0	1	3	5
C	5	1	0	2	4
D	3	3	2	0	2
E	1	5	4	2	0

dest	A	B	D
A	1	14	5
B	7	8	5
C	6	9	4
D	4	11	2

E's routing table:

Next hop	A	B	D
A	1	14	5
B	7	8	5
C	6	9	4
D	4	11	2

Computer Communications - CSC1 551

Copyright © William C. Cheng

19

View From a Node

Info at node	A	B	C	D	E
A	0	6	5	3	1
B	6	0	1	3	5
C	5	1	0	2	4
D	3	3	2	0	2
E	1	5	4	2	0

dest	A	B	D
A	1	14	5
B	7	8	5
C	6	9	4
D	4	11	2

E's routing table:

Next hop	A	B	D
A	1	14	5
B	7	8	5
C	6	9	4
D	4	11	2

Computer Communications - CSC1 551

Copyright © William C. Cheng

30

Solution 1: Holdowns

- ↳ If metric increases, delay propagating information
- ↳ in our example, B delays advertising route
- ↳ C eventually thinks B's route is gone, picks its own route
- ↳ B then selects C as next hop
- ↳ Adversely affects convergence

Computer Communications - CSCI 551

Copyright © William C. Cheng

28

**C Sends Routes to B
B Updates Distance to A**

This is known as the *count to infinity* problem

Computer Communications - CSCI 551

Copyright © William C. Cheng

26

**C Sends Routes to B
B Updates Distance to A**

Computer Communications - CSCI 551

Copyright © William C. Cheng

29

How Are These Loops Caused?

- ↳ Observation 1: B's metric *increases*
- ↳ Observation 2: C picks B as next hop to A
- ↳ But, the *implicit path* from C to A includes itself!

Computer Communications - CSCI 551

Copyright © William C. Cheng

27

**B Sends Routes to C
C Updates Distance to A**

Computer Communications - CSCI 551

Copyright © William C. Cheng

25

**B Sends Routes to C
C Updates Distance to A**

Computer Communications - CSCI 551

Copyright © William C. Cheng

32

Example Where Split Horizon Fails

- When link breaks, C marks D as unreachable and reports that to A and B.
- Suppose A learns it first. A now thinks best path to D is through B. A reports D unreachable to B and a route of cost=3 to C.
- C thinks D is reachable through A at cost 4 and reports that to B.
- B reports a cost 5 to A who reports new cost to C.
- etc...

Computer Communications - CSC1 551

Copyright © William C. Cheng

34

Computing Implicit Paths

v	u
w	u
z	w
y	z
x	y

- To reduce the space requirements
- propagate for each destination not only the cost but also its predecessor
- can recursively compute the path
- space requirements independent of diameter

Computer Communications - CSC1 551

Copyright © William C. Cheng

36

Distance Vector in Practice

- RIP and RIP2
- uses split-horizon/poison reverse
- BGP/DRP
- propagates entire path
- path also used for effecting policies

Computer Communications - CSC1 551

Copyright © William C. Cheng

31

Other "Solutions"

- Split horizon
- B does not advertise route to C
- Poisoned reverse
- B advertises route to C with infinite distance
- works for two node loops
- does not work for loops with more nodes

Computer Communications - CSC1 551

Copyright © William C. Cheng

33

Avoiding the Bouncing Effect

- Select loop-free paths
- One way of doing this:
 - each route advertisement carries entire path
 - if a router sees itself in path, it rejects the route
- BGP does it this way
- Space proportional to diameter

[Cheng, Riley et al]

Computer Communications - CSC1 551

Copyright © William C. Cheng

35

Loop Freedom at Every Instant

- Does bouncing effect avoid loops?
 - No! *Transient* loops are still possible
 - Why? Because implicit path information may be stale
- Only way to fix this
 - ensure that you have up-to-date information by explicitly querying

Computer Communications - CSC1 551

Copyright © William C. Cheng

37

Link State Algorithms

Computer Communications - CSC1 551

Copyright © William C. Cheng

38

Basic Steps

- Each node assumed to know state of links to its neighbors
- Step 1:** Each node broadcasts its state to all other nodes
- Step 2:** Each node locally computes shortest paths to all other nodes from global state

Computer Communications - CSC1 551

Copyright © William C. Cheng

39

Building Blocks

- Reliable broadcast mechanism
- flooding
- sequence number issues
- Shortest path tree (SPT) algorithm
- Dijkstra's SPT algorithm

Computer Communications - CSC1 551

Copyright © William C. Cheng

40

Link State Packets (LSPs)

- Periodically, each node creates a Link state packet containing:
 - Node ID
 - List of neighbors and link cost
 - Sequence number
 - Time to live (TTL)
- Node outputs LSP on *all* its links

Computer Communications - CSC1 551

Copyright © William C. Cheng

41

Reliable Flooding

- When node *i* receives LSP from node *j*:
 - If LSP is the most recent LSP from *j* that *i* has seen so far, *i* saves it in database and forwards a copy on all links except link LSP was received on.
 - Otherwise, discard LSP.

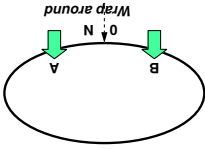
Computer Communications - CSC1 551

Copyright © William C. Cheng

42

Sequence Number Space Issues

- Problem: sequence number may wrap around
- Solution: treat space as circular, continue after wrap around:
 - A is less than B if
 - $A < B$ and $B - A < N/2$, or
 - $A > B$ and $A - B > N/2$



Computer Communications - CSC1 551

Copyright © William C. Cheng

44

One Solution: LSP Aging

- Nodes periodically decrement age (TTL) of stored LSPs
- LSPs expire when TTL reaches 0
- LSP is re-flooded once TTL = 0
- (haven't heard from you for a while, how are you doing?)
- Rebooted router waits until all LSPs have expired
- Trade-off between frequency of LSPs and router wait after reboot

Computer Communications - CSC1 551

Copyright © William C. Cheng

45

Lollipop Operation

Router comes up and starts with $-N/2$, then $-N/2 + 1$, $-N/2 + 2$, etc.

- When seq number becomes positive, wrap around as before
- a is older than b if:
 - $a > 0$ and $a < b$, or
 - $a > 0$, $a < b$ and $b - a < N/4$
 - $a < 0$, $b < 0$, $a > b$, and $a - b < N/4$

Computer Communications - CSC1 551

Copyright © William C. Cheng

46

Is Aging Still Needed?

- Yes! Stale LSPs are still possible
- suppose a router is down but not detected
- = net partitions and then heals
- Aging ensures that old state is eventually flushed out of the network

Computer Communications - CSC1 551

Copyright © William C. Cheng

43

Problem: Router Failure

- A failed router and comes up but does not remember the last sequence number it used before it crashed
- New LSPs may be ignored if they have lower sequence number

Computer Communications - CSC1 551

Copyright © William C. Cheng

45

A Better Solution

- Lollipop Sequence space [Perman83]
 - Divide sequence space N into 3 spaces:
 - Negative space: $-N/2 - 0$
 - The number 0
 - Positive space: 0 to $N/2 - 1$

Computer Communications - CSC1 551

Copyright © William C. Cheng

47

Lollipop Operation (Cont..)

- Newly booted router always starts with oldest seq num ($-N/2$)
- New rule:
 - if router R1 gets older LSP from router R2, R1 informs R2 of the sequence number in R1's LSP
- Newly booted router discovers its seq num before it crashed and resumes

Computer Communications - CSC1 551

Copyright © William C. Cheng

50

Computer Communications - CSC1 551

Example

step	SPT	B	C	D	E	F
0	A	2,A	5,A	1,A	~	~
1	A	2,A	3,C	2,D	~	~
2	A	2,A	3,C	2,D	2,E	~
3	A	2,A	3,C	2,D	2,E	4,F

Copyright © William C. Cheng

51

Computer Communications - CSC1 551

Example

step	SPT	B	C	D	E	F
0	A	2,A	5,A	1,A	~	~
1	A	2,A	3,C	2,D	~	~
2	A	2,A	3,C	2,D	2,E	~
3	A	2,A	3,C	2,D	2,E	4,F

Copyright © William C. Cheng

52

Computer Communications - CSC1 551

Example

step	SPT	B	C	D	E	F
0	A	2,A	5,A	1,A	~	~
1	A	2,A	3,C	2,D	~	~
2	A	2,A	3,C	2,D	2,E	~
3	A	2,A	3,C	2,D	2,E	4,F

Copyright © William C. Cheng

49

Computer Communications - CSC1 551

SPT Algorithm (Dijkstra)

$SPT = \{a\}$
 for all nodes v
 if v adjacent to a then $D(v) = cost(a, v)$
 else $D(v) = infinity$
 Loop
 find w not in SPT , where $D(w)$ is min
 add w in SPT
 for all v adjacent to w and not in SPT
 $D(v) = \min(D(v), D(w) + C(w, v))$
 until all nodes are in SPT

Copyright © William C. Cheng

53

Computer Communications - CSC1 551

Example

step	SPT	B	C	D	E	F
0	A	2,A	5,A	1,A	~	~
1	A	2,A	3,C	2,D	~	~
2	A	2,A	3,C	2,D	2,E	~
3	A	2,A	3,C	2,D	2,E	4,F

Copyright © William C. Cheng

54

Computer Communications - CSC1 551

Example

step	SPT	B	C	D	E	F
0	A	2,A	5,A	1,A	~	~
1	A	2,A	3,C	2,D	~	~
2	A	2,A	3,C	2,D	2,E	~
3	A	2,A	3,C	2,D	2,E	4,F

Copyright © William C. Cheng

60

LS v.s. DV (Cont..)

Robustness:

- LS can broadcast incorrect/corrupted LSP
- localized problem
- DV can advertise incorrect paths to all destinations
- incorrect calculation can spread to entire network

- In LS, nodes must compute consistent routes independently
- must protect against LSDB corruption

- In DV, routes are computed relative to other nodes

Computer Communications - CSC1 551

Copyright © William C. Cheng

59

LS v.s. DV (Cont..)

- Msg size:
 - LS: small
 - DV: potentially large
- Msg exchange:
 - LS: O(nE)
 - DV: only to neighbors
- Convergence speed:
 - LS: fast
 - DV: fast with triggered updates
- Space requirements:
 - LS maintains entire topology
 - DV maintains only neighbor state

Computer Communications - CSC1 551

Copyright © William C. Cheng

58

LS v.s. DV

- In DV send everything you know to your neighbors
- In LS send info about your neighbors to everyone

Computer Communications - CSC1 551

Copyright © William C. Cheng

57

Link State Characteristics

- With consistent LSDBs, all nodes compute consistent loop-free paths
- Limited by Dijkstra computation overhead, space requirements
- Can still have transient loops

Computer Communications - CSC1 551

Copyright © William C. Cheng

56

Link State Algorithm

Flooding:

- Periodically distribute link-state advertisement (LSA) to neighbors
- LSA contains delays to each neighbor
- Install received LSA in LS database
- Re-distribute LSA to all neighbors

Path Computation

- Use Dijkstra's shortest path algorithm to compute distances to all destinations
- Install <destination, next-hop> pair in forwarding table

Computer Communications - CSC1 551

Copyright © William C. Cheng

55

Example

step	SPT	D(b),P(b)	D(c),P(c)	D(d),P(d)	D(e),P(e)	D(f),P(f)
0	A	2,A	5,A	1,A	~	~
1	AD	4,D	2,D	~	~	~
2	ADE	2,A	3,E	~	~	~
3	ADEB	3,E	~	~	~	~
4	ADEBC	~	~	~	~	~
5	ADEBCF	~	~	~	~	~

Computer Communications - CSC1 551

Copyright © William C. Cheng

Computer Communications - CSC1 551

What Makes Routing Hard?

- Scalability to many hosts
- Reliability and robustness
- Dealing with changes
 - = some changes (link goes down) should be dealt with ASAP,
 - some (link goes up and down) should be suppressed
- Congestion
 - = why not route around congestion?
 - routing algorithm takes too long to react to congestion
- Distributed computation (and debugging)
- Routing and business/policy issues

Copyright © William C. Cheng

Computer Communications - CSC1 551

Example Area Hierarchy

The diagram shows three sub-ASes: 11, 12, and 13. Sub-AS 11 contains networks 11.1, 11.1.1, 11.2, 11.2.1, 11.2.2, 11.2.3, 11.2.4, and 11.2.5. Sub-AS 12 contains networks 12.1, 12.1.1, 12.2, 12.2.1, and 12.2.2. Sub-AS 13 contains networks 13.1, 13.2, 13.1.1, and 13.1.2. A yellow box highlights network 11.2.1. Connections exist between sub-ASes and between networks within sub-ASes.

- AS's (11, 12, 13)
- sub-AS's (11.1, etc.)
- networks (11.2.1/2/4, etc.)
- routing table at 11.2.1:
 - 11.2.3/2/4: 11.2.1x
 - 11.2.4/2/4: 11.2.1y
 - 11.1/1/6: ?

Copyright © William C. Cheng

Computer Communications - CSC1 551

Example Area Hierarchy

The diagram is identical to the previous one, showing sub-ASes 11, 12, and 13 with their internal networks and connections. A yellow box highlights network 11.2.1.

- AS's (11, 12, 13)
- sub-AS's (11.1, etc.)
- networks (11.2.1/2/4, etc.)
- routing table at 11.2.1:
 - 11.2.3/2/4: 11.2.1x
 - 11.2.4/2/4: 11.2.1y
 - 11.1/1/6: ?

Copyright © William C. Cheng

Computer Communications - CSC1 551

LS v.s. DV (Cont...)

- DV risks:
 - = looping, convergence time, corrupted host can get all routes
 - = solutions are split horizon, poison reverse, path vectors
- LS risks:
 - = flooding of information, must know whole topology (hierarchy and aggregation are forces against this)
 - = Bottom line: no clear winner, but we see more frequent use of LS in the Internet

Copyright © William C. Cheng

Computer Communications - CSC1 551

Scaling to Big Networks

- Internet today is on the order of 1 million networks
- Approaches
 - = hierarchy and aggregation
- Key idea: want to know about you and not about some networks far away
- Two approaches:
 - = area hierarchy
 - approach used in the Internet
 - semi-manual aggregation
 - = landmark hierarchy
 - not directly used

Copyright © William C. Cheng

Computer Communications - CSC1 551

Example Area Hierarchy

The diagram is identical to the previous ones, showing sub-ASes 11, 12, and 13 with their internal networks and connections. A yellow box highlights network 11.2.1.

- AS's (11, 12, 13)
- sub-AS's (11.1, etc.)
- networks (11.2.1/2/4, etc.)
- routing table at 11.2.1:
 - 11.2.3/2/4: 11.2.1x
 - 11.2.4/2/4: 11.2.1y
 - 11.1/1/6: ?

Computer Communications - CSCI 551

Example Area Hierarchy

- AS's (11, 12, 13)
- sub-AS's (11.1, etc.)
- networks (11.2.1/24, etc.)
- routing table at 11.2.1:
 - 11.2.3/24: 11.2.1x
 - 11.2.4/24: 11.2.1x
 - 11.1/16: 11.2.1z
 - 12/8: 11.2.1w
 - 13/8: ?

Copyright © William C. Cheng

Computer Communications - CSCI 551

Example Area Hierarchy

- AS's (11, 12, 13)
- sub-AS's (11.1, etc.)
- networks (11.2.1/24, etc.)
- routing table at 11.2.1:
 - 11.2.3/24: 11.2.1x
 - 11.2.4/24: 11.2.1x
 - 11.1/16: 11.2.1z
 - 12/8: ?

Copyright © William C. Cheng

Computer Communications - CSCI 551

Example Area Hierarchy

- AS's (11, 12, 13)
- sub-AS's (11.1, etc.)
- networks (11.2.1/24, etc.)
- routing table at 11.2.1:
 - 11.2.3/24: 11.2.1x
 - 11.2.4/24: 11.2.1x
 - 11.1/16: 11.2.1z
 - 12/8: 11.2.1w
 - 13/8: 11.2.1y

is this the only way?

Copyright © William C. Cheng