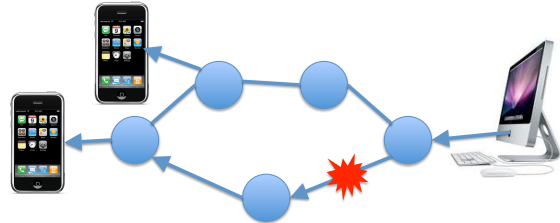


Routing Convergence

Mike Freedman
COS 461: Computer Networks

<http://www.cs.princeton.edu/courses/archive/spr14/cos461/>

Routing Changes



- **Topology changes:** new route to the same place
- **Host mobility:** route to a different place

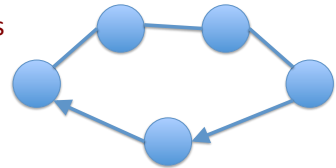
2

Topology Changes

3

Two Types of Topology Changes

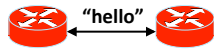
- **Planned**
 - Maintenance: shut down a node or link
 - Energy savings: shut down a node or link
 - Traffic engineering: change routing configuration
- **Unplanned Failures**
 - Fiber cut, faulty equipment, power outage, software bugs, ...



4

Detecting Topology Changes

- **Beaconing**
 - Periodic “hello” messages in both directions
 - Detect a failure after a few missed “hellos”



- **Performance trade-offs**
 - Detection delay
 - Overhead on link bandwidth and CPU
 - Likelihood of false detection

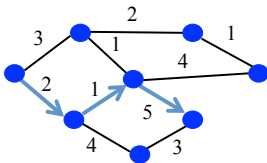
5

Routing Convergence: Link-State Routing

6

Convergence

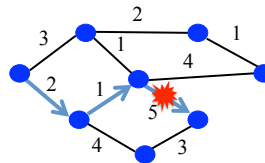
- **Control plane**
 - All nodes have consistent information
- **Data plane**
 - All nodes forward packets in a consistent way



7

Transient Disruptions

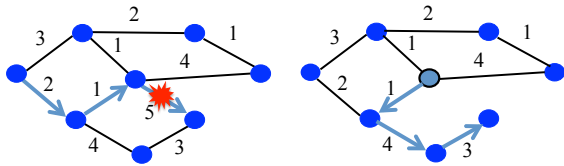
- **Detection delay**
 - A node does not detect a failed link immediately
 - ... and forwards data packets into a “blackhole”
 - Depends on timeout for detecting lost hellos



8

Transient Disruptions

- **Inconsistent link-state database**
 - Some routers know about failure before others
 - Inconsistent paths cause transient forwarding loops



9

Convergence Delay

- **Sources of convergence delay**
 - Detection latency
 - Updating control-plane information
 - Computing and install new forwarding tables
- **Performance during convergence period**
 - Lost packets due to blackholes and TTL expiry
 - Looping packets consuming resources
 - Out-of-order packets reaching the destination
- **Very bad for VoIP, online gaming, and video**

10

Reducing Convergence Delay

- **Faster detection**
 - Smaller hello timers, better link-layer technologies
- **Faster control plane**
 - Flooding immediately
 - Sending routing messages with high-priority
- **Faster computation**
 - Faster processors, and incremental computation
- **Faster forwarding-table update**
 - Data structures supporting incremental updates

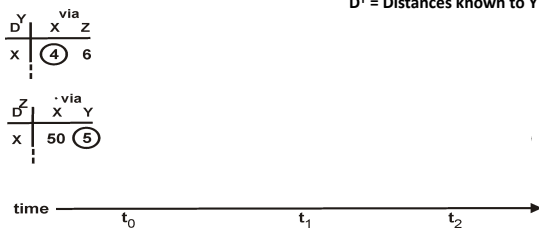
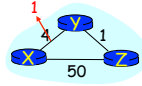
11

Slow Convergence in Distance-Vector Routing

12

Distance Vector: Link Cost Changes

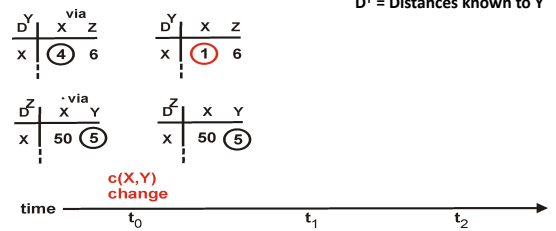
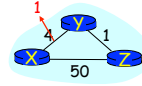
- Link cost decreases and recovery
 - Node updates the distance table
 - If cost change in least cost path, notify neighbors



13

Distance Vector: Link Cost Changes

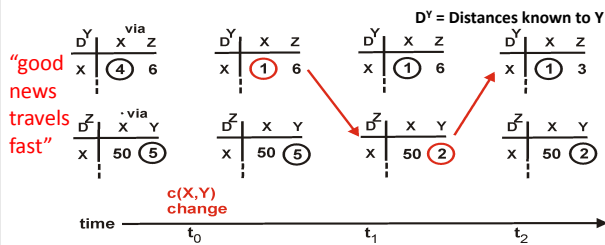
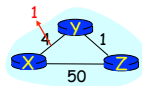
- Link cost decreases and recovery
 - Node updates the distance table
 - If cost change in least cost path, notify neighbors



14

Distance Vector: Link Cost Changes

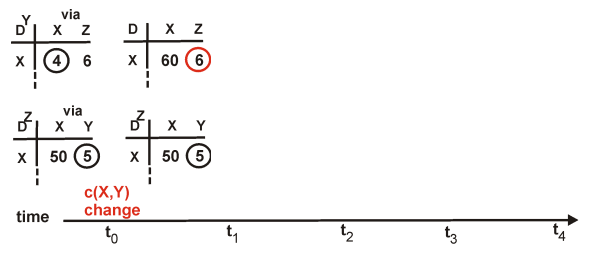
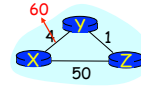
- Link cost decreases and recovery
 - Node updates the distance table
 - If cost change in least cost path, notify neighbors



15

Distance Vector: Link Cost Changes

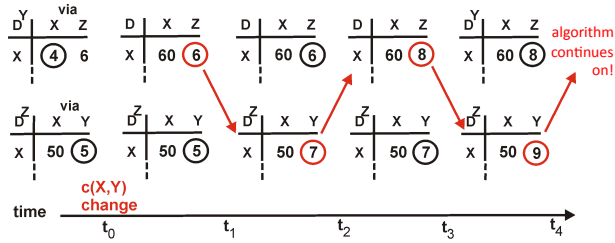
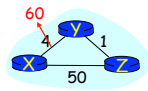
- Link cost increases and failures
 - Bad news travels slowly
 - “Count to infinity” problem!



16

Distance Vector: Link Cost Changes

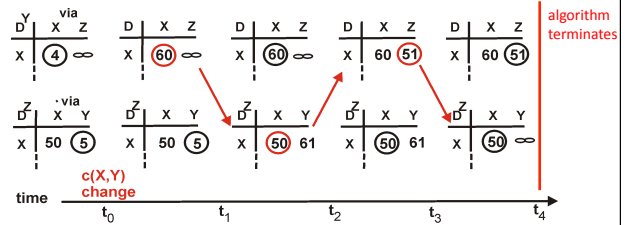
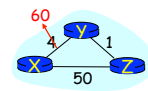
- Link cost increases and failures
 - Bad news travels slowly
 - “Count to infinity” problem!



17

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)
 - Still, can have problems in larger networks



18

Redefining Infinity

- Avoid “counting to infinity”
 - By making “infinity” smaller!
- Routing Information Protocol (RIP)
 - All links have cost 1
 - Valid path distances of 1 through 15
 - ... with 16 representing infinity
- Used mainly in small networks

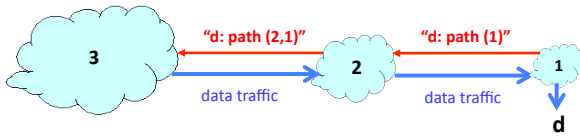
19

Reducing Convergence Time With Path-Vector Routing (e.g., Border Gateway Protocol)

20

Path-Vector Routing

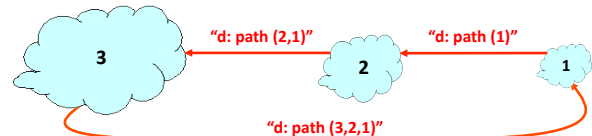
- Extension of distance-vector routing
 - Support flexible routing policies
 - Avoid count-to-infinity problem
- Key idea: advertise the entire path
 - Distance vector: send distance metric per dest d
 - Path vector: send the entire path for each dest d



21

Faster Loop Detection

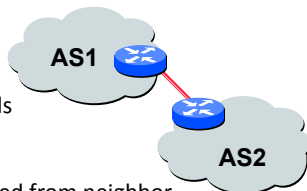
- Node can easily detect a loop
 - Look for its own node identifier in the path
 - E.g., node 1 sees itself in the path “3, 2, 1”
- Node can simply discard paths with loops
 - E.g., node 1 simply discards the advertisement



22

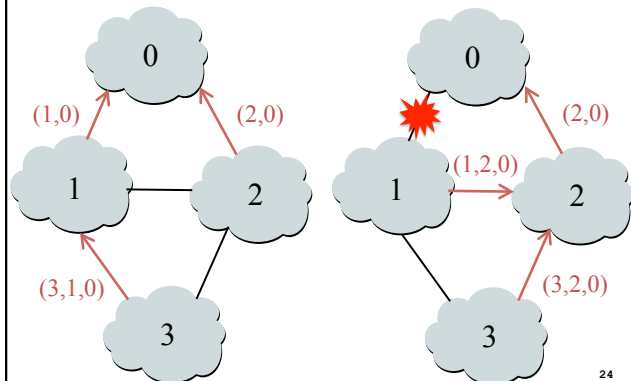
BGP Session Failure

- BGP runs over TCP
 - BGP only sends updates when changes occur
 - TCP doesn't detect lost connectivity on its own
- Detecting a failure
 - Keep-alive: 60 seconds
 - Hold timer: 180 seconds
- Reacting to a failure
 - Discard all routes learned from neighbor
 - Send new updates for any routes that change



23

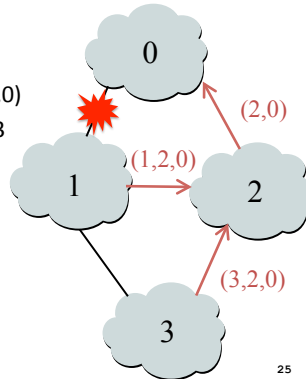
Routing Change: Before and After



24

Routing Change: Path Exploration

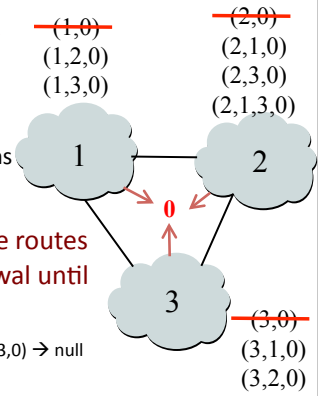
- **AS 1**
 - Delete the route (1,0)
 - Switch to next route (1,2,0)
 - Send route (1,2,0) to AS 3
- **AS 3**
 - Sees (1,2,0) replace (1,0)
 - Compares to route (2,0)
 - Switches to using AS 2



25

Routing Change: Path Exploration

- **Initial: All AS use direct**
 - ~~(1,0)~~
 - (1,2,0)
 - (1,3,0)
- **Then destination 0 dies**
 - All ASes lose direct path
 - All switch to longer paths
 - Eventually withdrawn



- **How many intermediate routes following (2,0) withdrawal until no route known to 2?**

(2,0) → (2,1,0) → (2,3,0) → (2,1,3,0) → null

~~(3,0)~~
(3,1,0)
(3,2,0)

BGP Converges Slowly

- **Path vector avoids count-to-infinity**
 - But, ASes still must explore many alternate paths to find highest-ranked available path
- **Fortunately, in practice**
 - Most popular destinations have stable BGP routes
 - Most instability lies in a few unpopular destinations
- **Still, lower BGP convergence delay is a goal**
 - Can be tens of seconds to tens of minutes

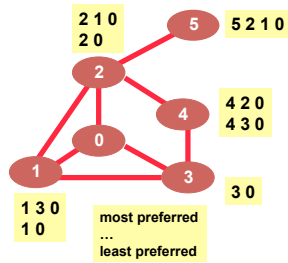
27

BGP Instability

28

Stable Paths Problem (SPP) Instance

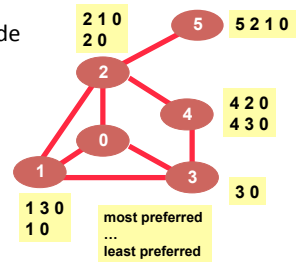
- **Node**
 - BGP-speaking router
 - Node 0 is destination
- **Edge**
 - BGP adjacency
- **Permitted paths**
 - Set of routes to 0 at each node
 - Ranking of the paths



29

Stable Paths Problem (SPP) Instance

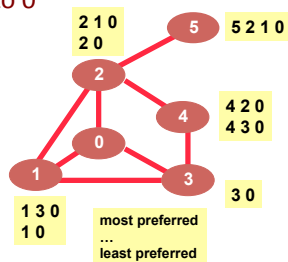
- **Solution**
 - Path assignment per node
 - Can be the “null” path
- **If node u has path uwP**
 - {u,w} is edge in graph
 - w is assigned path wP
- **Each node is assigned**
 - Highest ranked path consistent with its neighbors



30

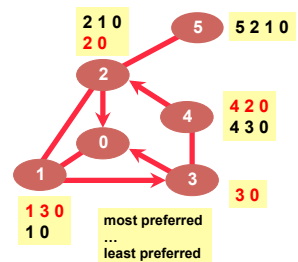
Stable Paths Problem (SPP) Instance

- **1 will use a direct path to 0**
(A) True (B) False
- **5 has a path to 0**
(A) True (B) False



31

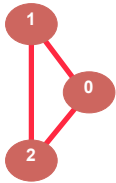
Stable Paths Problem (SPP) Instance



32

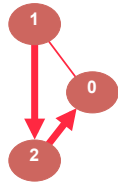
SPP May Have Multiple Solutions

1 2 0
1 0



2 1 0
2 0

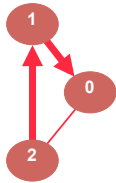
1 2 0
1 0



2 1 0
2 0

First solution

1 2 0
1 0

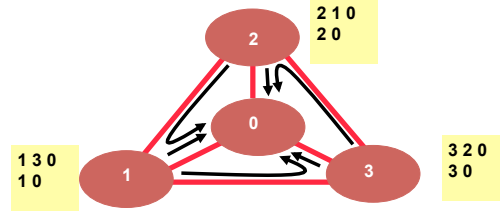


2 1 0
2 0

Second solution

33

An SPP May Have No Solution



1 3 0
1 0

2 1 0
2 0

3 2 0
3 0

34

Avoiding BGP Instability

- **Detecting conflicting policies**
 - Computationally expensive
 - Requires too much cooperation
- **Detecting oscillations**
 - Observing the repetitive BGP routing messages
- **Restricted routing policies and topologies**
 - Policies based on business relationships

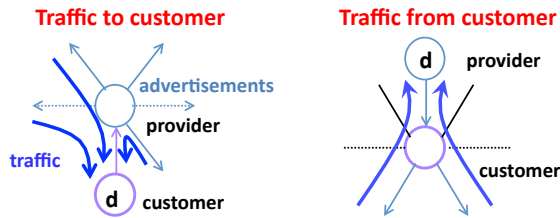
35

AS (Autonomous System) Business Relationships

36

Customer-Provider Relationship

- Customer pays provider for access to Internet
 - Provider exports its customer routes to everybody
 - Customer exports provider routes only to its customers

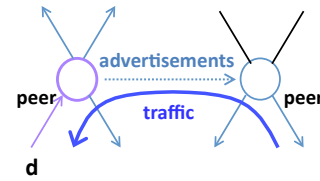


37

Peer-Peer Relationship

- Peers exchange traffic between their customers
 - AS exports only customer routes to a peer
 - AS exports a peer's routes only to its customers

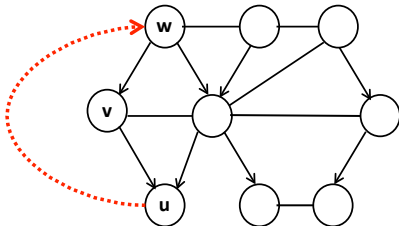
Traffic to/from the peer and its customers



38

Hierarchical AS Relationships

- Provider-customer graph is directed and acyclic
 - If u is a customer of v and v is a customer of w
 - ... then w is not a customer of u

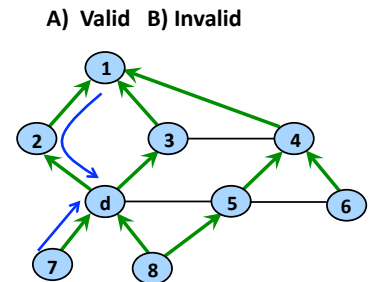


39

Valid and Invalid Paths

- Path 1 2 d
- Path 7 d
- Path 5 8 d
- Path 6 4 3 d
- Path 8 5 d
- Path 6 5 d
- Path 1 4 3 d

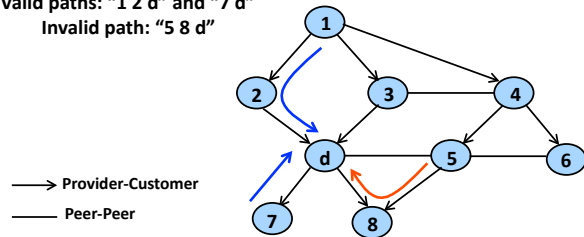
→ Provider-Customer
 — Peer-Peer



40

Valid and Invalid Paths

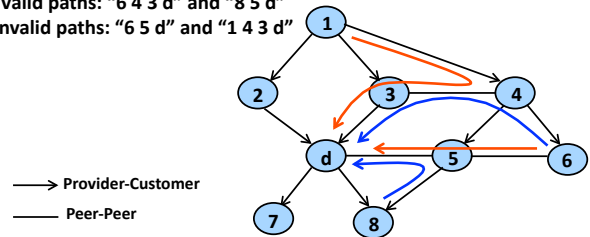
Valid paths: "1 2 d" and "7 d"
Invalid path: "5 8 d"



41

Valid and Invalid Paths

Valid paths: "6 4 3 d" and "8 5 d"
Invalid paths: "6 5 d" and "1 4 3 d"



42

Local Control, Global Stability: "Gao-Rexford Conditions"

1. **Route export**
 - Don't export routes learned from a peer or provider to another peer or provider
 2. **Global topology**
 - Provider-customer relationship graph is acyclic
 - E.g., my customer's customer is not my provider
 3. **Route selection**
 - Prefer routes through customers over routes through peers and providers
- **Guaranteed to converge to unique, stable solution**

43

Conclusion

- **The only constant is change**
 - Planned topology and configuration changes
 - Unplanned failure and recovery
- **Routing-protocol convergence**
 - Transient period of disagreement
 - Blackholes, loops, and out-of-order packets
- **Routing instability**
 - Permanent conflicts in routing policy
 - Leading to bi-stability or oscillation

44