# Routing Convergence

# Routing Changes



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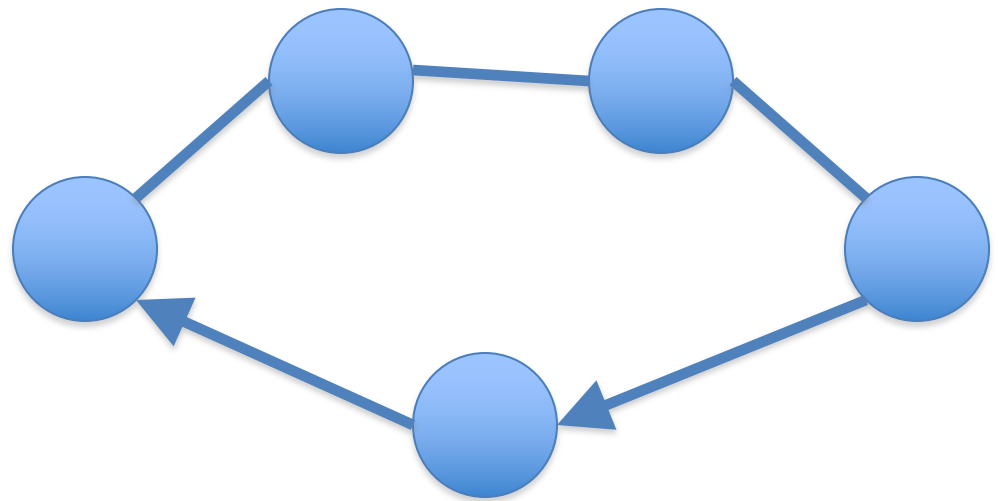- Host mobility: route to a different place

# Topology Changes

# 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
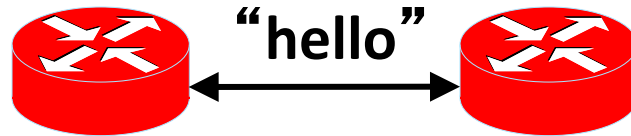  - Fiber cut,
    faulty equipment,
    power outage,
    software bugs, …

# 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

# Routing Convergence: Link-State Routing

# Convergence

- ## Control plane
  - All nodes have consistent information

- ## Data plane
  - All nodes forward packets in a consistent way

# 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

# Transient Disruptions

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

# 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

# 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

# Slow Convergence in Distance-Vector Routing

# Distance Vector: Link Cost Changes

- ## Link cost decreases and recovery

  - Node updates the distance table
  - If cost change in least cost path, notify neighbors

**$D^Y$ = Distances known to Y**

| $D^Y$ | via X | via Z |
|---|---|---|
| X | (4) | 6 |

| $D^Z$ | via X | via Y |
|---|---|---|
| X | 50 | (5) |

time ————————————————————————————————→

$t_0$            $t_1$            $t_2$

# Distance Vector: Link Cost Changes

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

**$D^Y$ = Distances known to Y**

| $D^Y$ | via X | Z |
|---|---|---|
| X | (4) | 6 |

| $D^Y$ | via X | Z |
|---|---|---|
| X | (1) | 6 |

| $D^Z$ | via X | Y |
|---|---|---|
| X | 50 | (5) |

| $D^Z$ | via X | Y |
|---|---|---|
| X | 50 | (5) |

**c(X,Y) change**

time ──────────────────────────────────────►
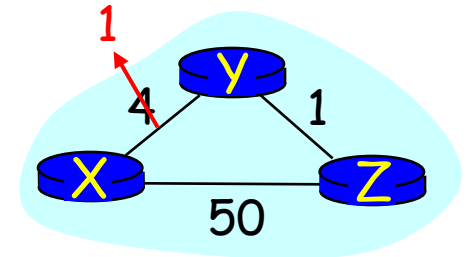
$t_0$          $t_1$          $t_2$

# Distance Vector: Link Cost Changes

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



**$D^Y$ = Distances known to Y**

**"good news travels fast"**

| $D^Y$ | via X | Z |
|---|---|---|
| X | ④ | 6 |

| $D^Z$ | via X | Y |
|---|---|---|
| X | 50 | ⑤ |

| $D^Y$ | X | Z |
|---|---|---|
| X | ① | 6 |

| $D^Z$ | X | Y |
|---|---|---|
| X | 50 | ⑤ |

| $D^Y$ | X | Z |
|---|---|---|
| X | ① | 6 |

| $D^Z$ | X | Y |
|---|---|---|
| X | 50 | ② |

| $D^Y$ | X | Z |
|---|---|---|
| X | ① | 3 |

| $D^Z$ | X | Y |
|---|---|---|
| X | 50 | ② |

**c(X,Y) change**

time → $t_0$      $t_1$      $t_2$

15

# Distance Vector: Link Cost Changes

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



```
      via                        via
Y                          
D | X   Z                  D | X   Z
------------               ------------
X | (4)  6                 X | 60  (6)
  :                          :
```

```
      via                        via
Z                          
D | X   Y                  D | X   Y
------------               ------------
X | 50  (5)                X | 50  (5)
  :                          :
```

c(X,Y) change

time ————————————————————————————————————→
        t₀        t₁        t₂        t₃        t₄
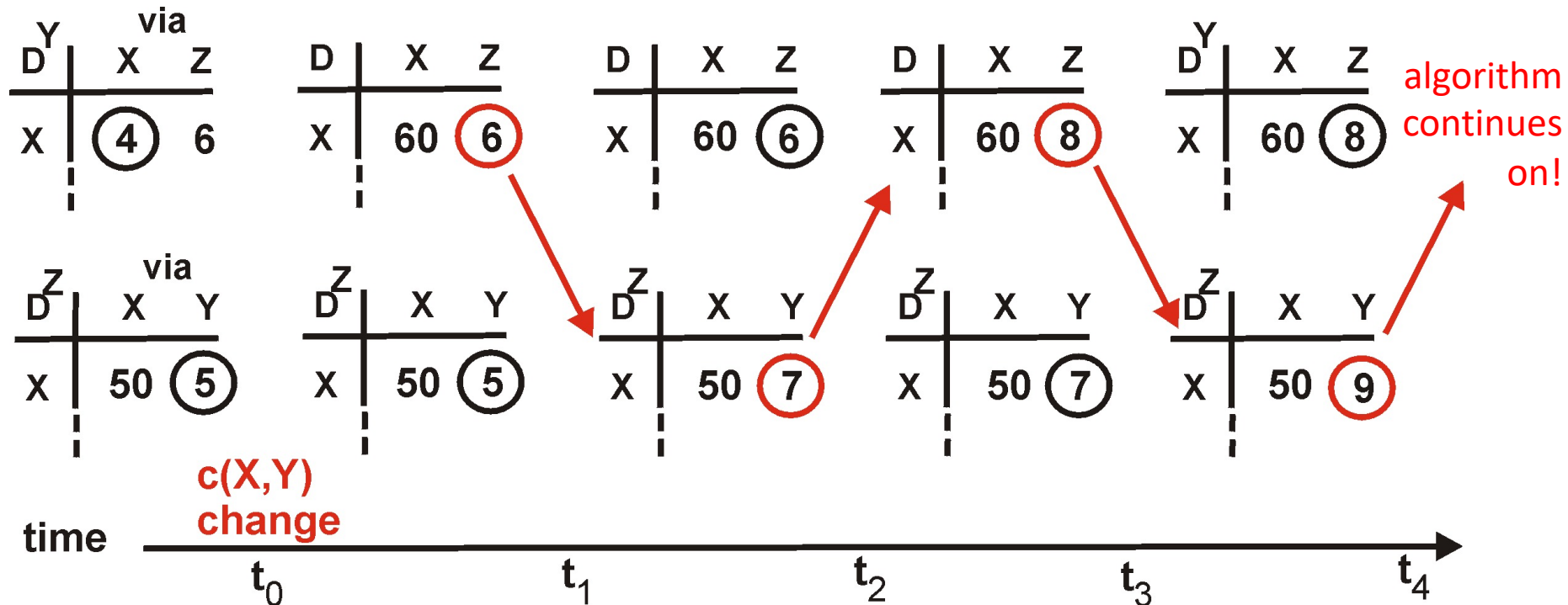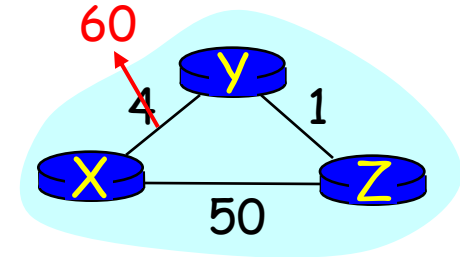
# Distance Vector: Link Cost Changes

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



| $D^Y$ | via X | Z |
|---|---|---|
| X | (4) | 6 |

| D | X | Z |
|---|---|---|
| X | 60 | (6) |

| D | X | Z |
|---|---|---|
| X | 60 | (6) |

| D | X | Z |
|---|---|---|
| X | 60 | (8) |

| $D^Y$ | X | Z |
|---|---|---|
| X | 60 | (8) |

algorithm continues on!

| $D^Z$ | via X | Y |
|---|---|---|
| X | 50 | (5) |

| $D^Z$ | X | Y |
|---|---|---|
| X | 50 | (5) |

| $D^Z$ | X | Y |
|---|---|---|
| X | 50 | (7) |

| $D^Z$ | X | Y |
|---|---|---|
| X | 50 | (7) |

| $D^Z$ | X | Y |
|---|---|---|
| X | 50 | (9) |

time

$c(X,Y)$ change

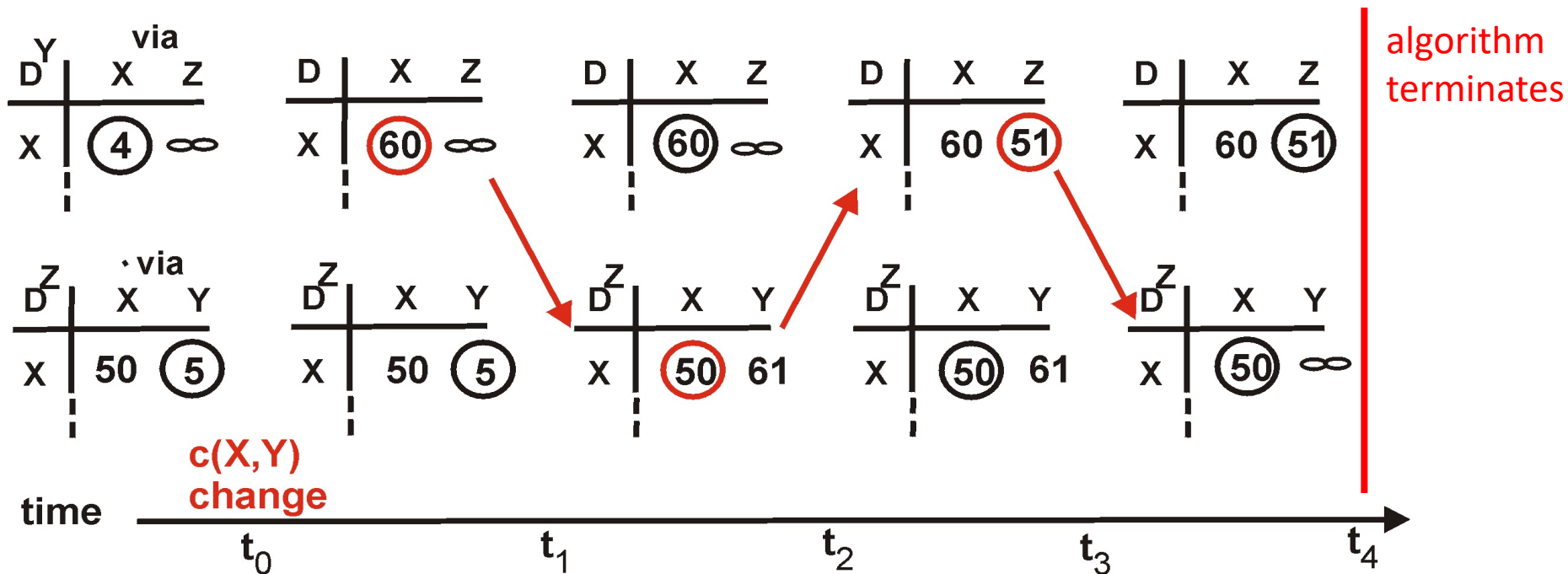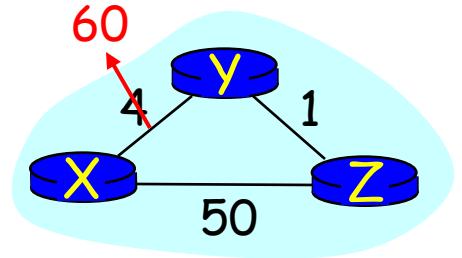$t_0$   $t_1$   $t_2$   $t_3$   $t_4$

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)
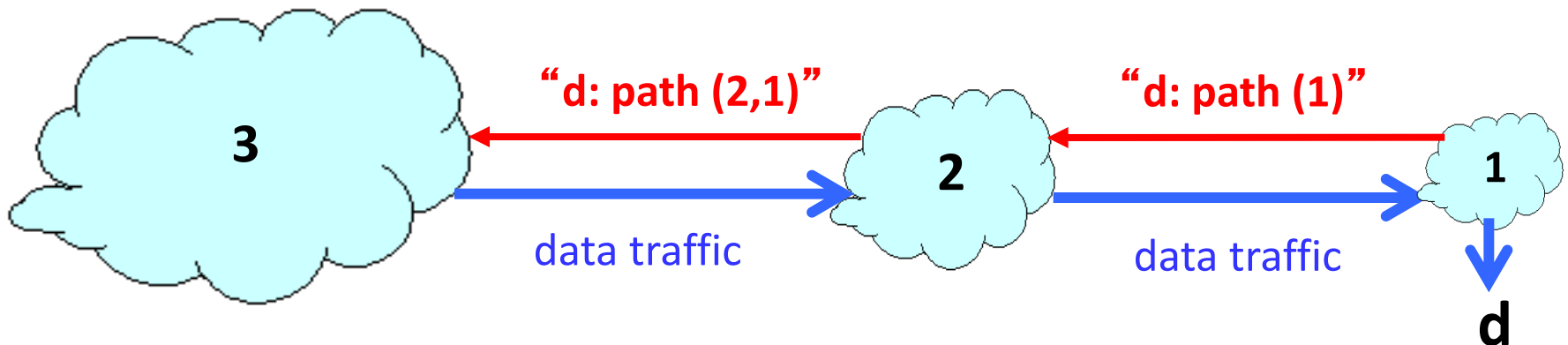    - Still, can have problems in larger networks



algorithm terminates

| $D^Y$ | X | Z |
|---|---|---|
| X | ④ | ∞ |

| D | X | Z |
|---|---|---|
| X | ⑥⓪ | ∞ |

| D | X | Z |
|---|---|---|
| X | ⑥⓪ | ∞ |

| D | X | Z |
|---|---|---|
| X | 60 | ㊴ |

| D | X | Z |
|---|---|---|
| X | 60 | ㊴ |

| $D^Z$ | X | Y |
|---|---|---|
| X | 50 | ⑤ |

| $D^Z$ | X | Y |
|---|---|---|
| X | 50 | ⑤ |

| $D^Z$ | X | Y |
|---|---|---|
| X | ㊵ | 61 |

| $D^Z$ | X | Y |
|---|---|---|
| X | ㊵ | 61 |

| $D^Z$ | X | Y |
|---|---|---|
| X | ㊵ | ∞ |

time

c(X,Y) change

$t_0$    $t_1$    $t_2$    $t_3$    $t_4$

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

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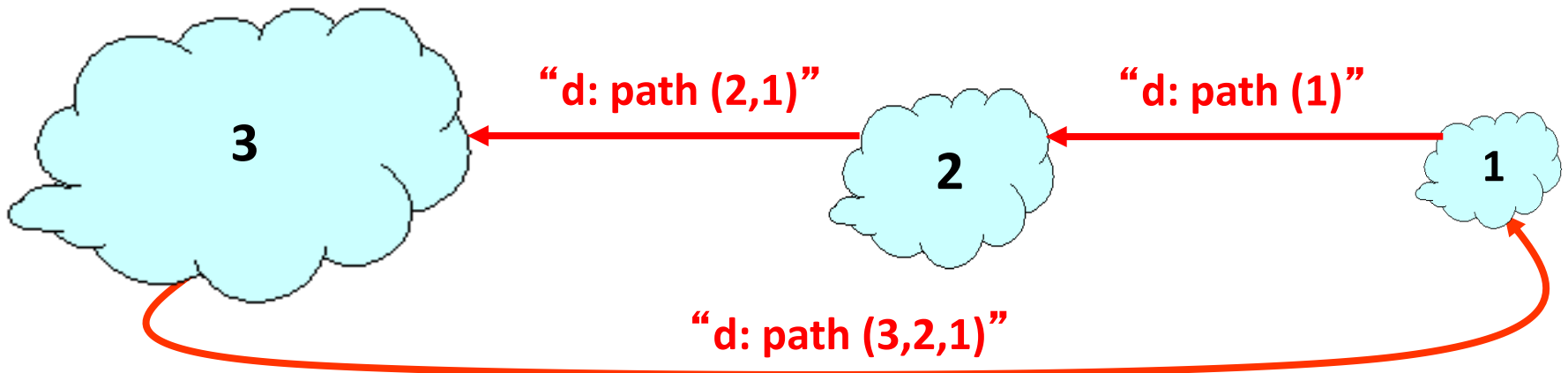
# Path-Vector Routing

- 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

# 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
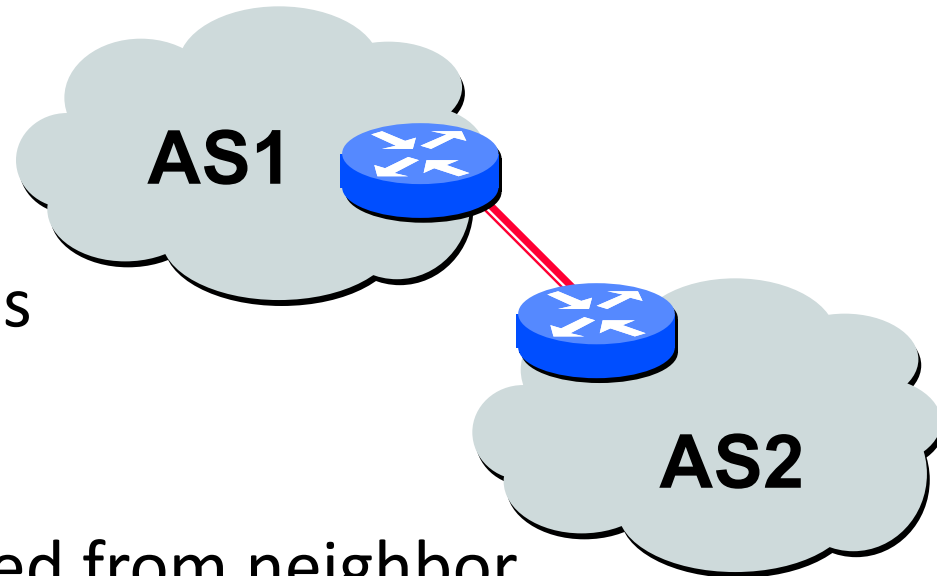  - E.g., node 1 simply discards the advertisement

"d: path (2,1)"    "d: path (1)"

3        2        1
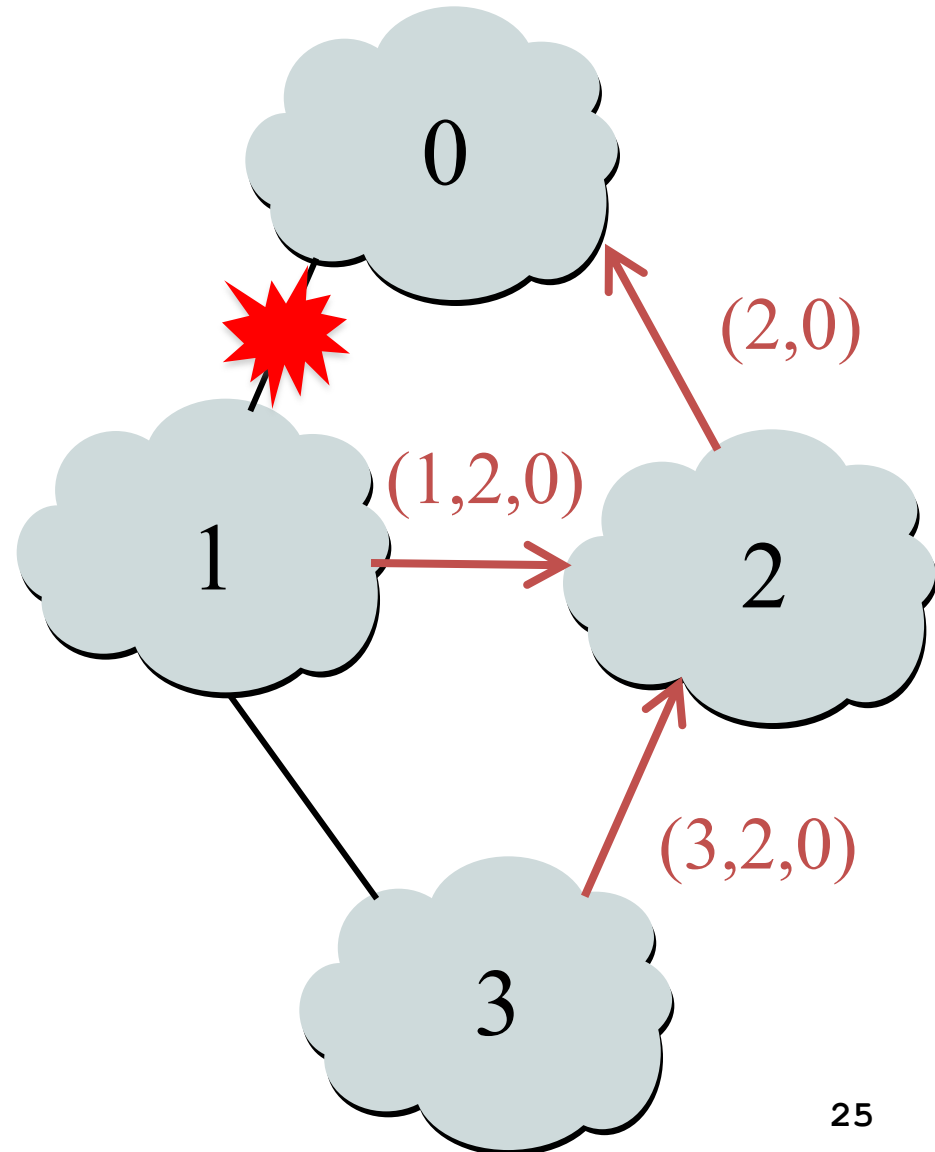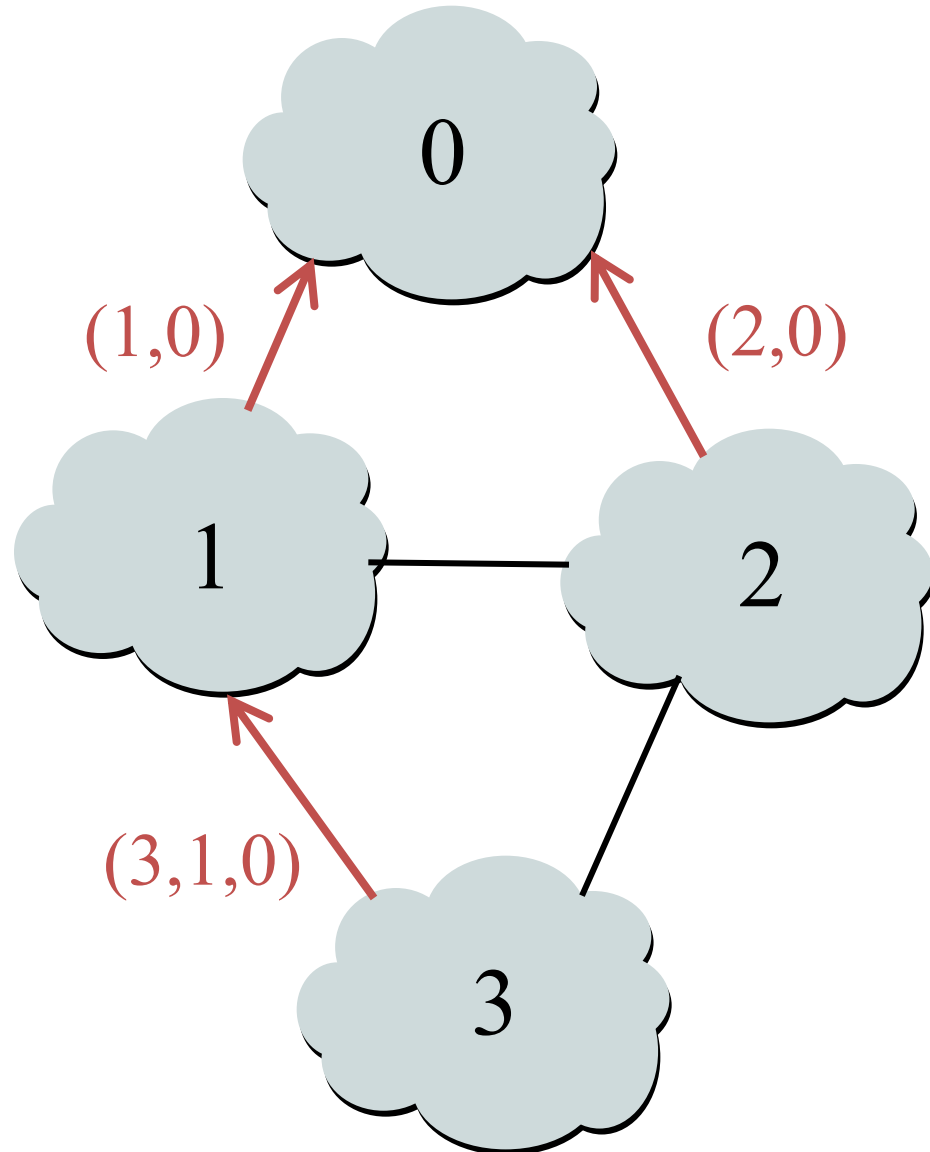
"d: path (3,2,1)"

# Causes of BGP Routing Changes

- Topology changes
  - Equipment going up or down
  - Deployment of new routers or sessions

- BGP session failures
  - Due to equipment failures, maintenance, etc.
  - Or, due to congestion on the physical path

- Changes in routing policy
  - Changes in preferences in the routes
  - Changes in whether the route is exported

- Persistent protocol oscillation
  - Conflicts between policies in different ASes

23

# 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

**AS1**

**AS2**

- ## Reacting to a failure
  - Discard all routes learned from neighbor
  - Send new updates for any routes that change
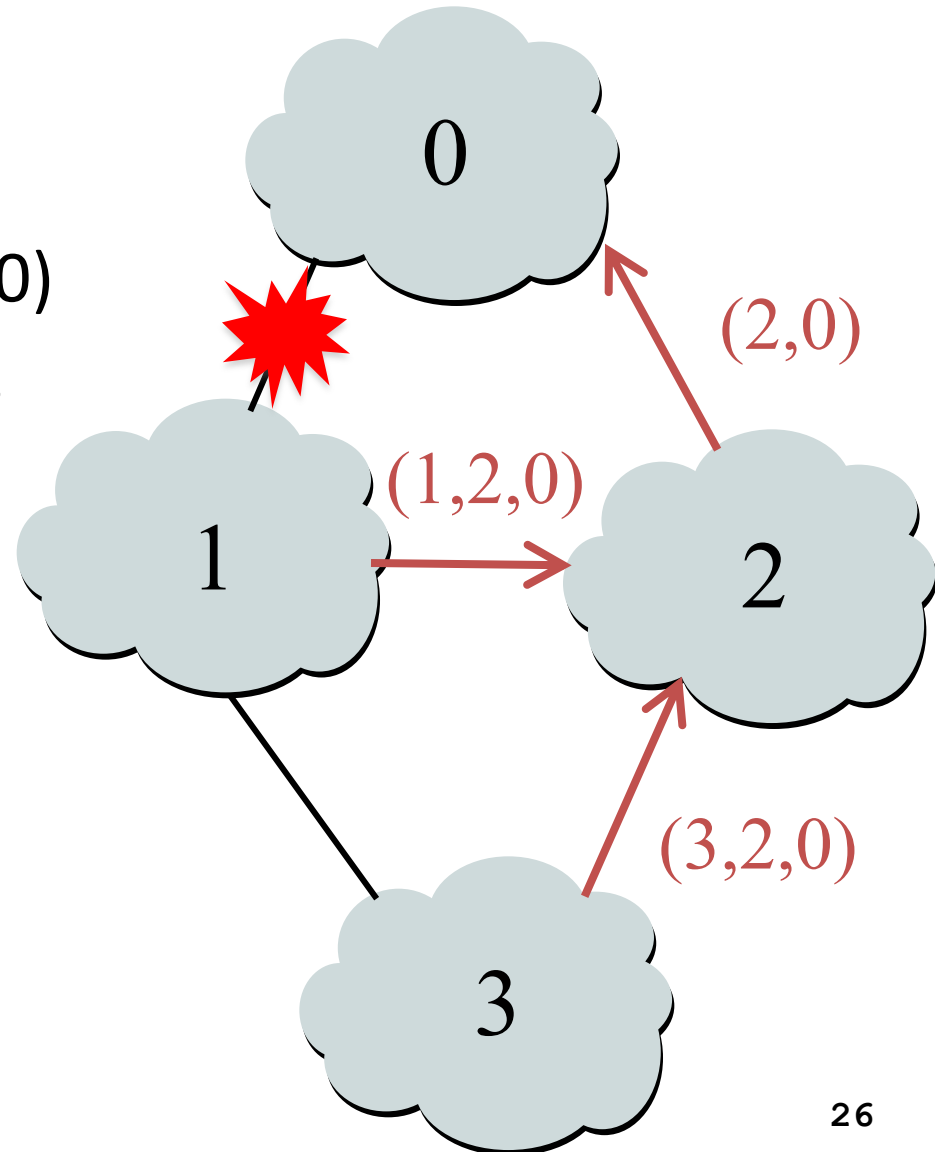
# Routing Change: Before and After

# Routing Change: Path Exploration

- ## AS 1
  - Delete the route (1,0)
  - Switch to next route (1,2,0)
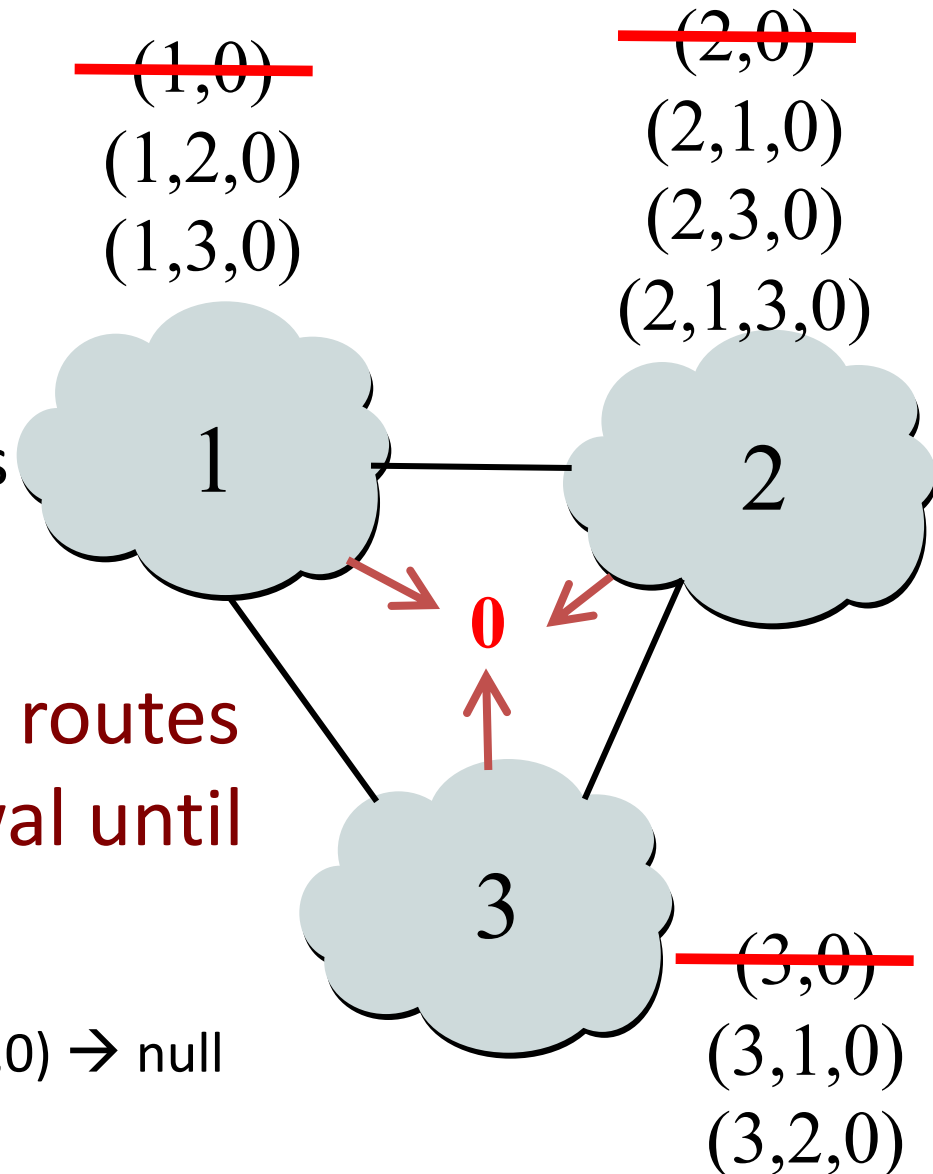  - Send route (1,2,0) to AS 3

- ## AS 3
  - Sees (1,2,0) replace (1,0)
  - Compares to route (2,0)
  - Switches to using AS 2

# Routing Change: Path Exploration

- Initial: All AS use direct

- 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?

$(2,0) \rightarrow (2,1,0) \rightarrow (2,3,0) \rightarrow (2,1,3,0) \rightarrow$ null

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

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

**1**

**2**

**0**

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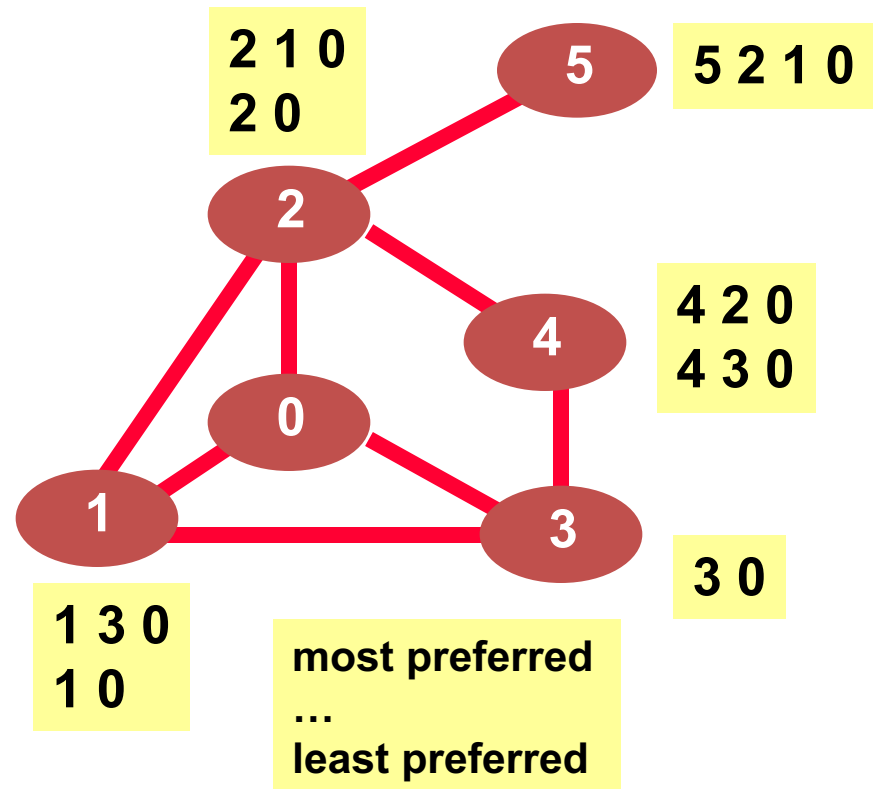
~~(3,0)~~
(3,1,0)
(3,2,0)

# 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

# BGP Instability

# Stable Paths Problem (SPP) Instance

- ## Node
  - BGP-speaking router
  - Node 0 is destination

- ## Edge
  - BGP adjacency

- ## Permitted paths
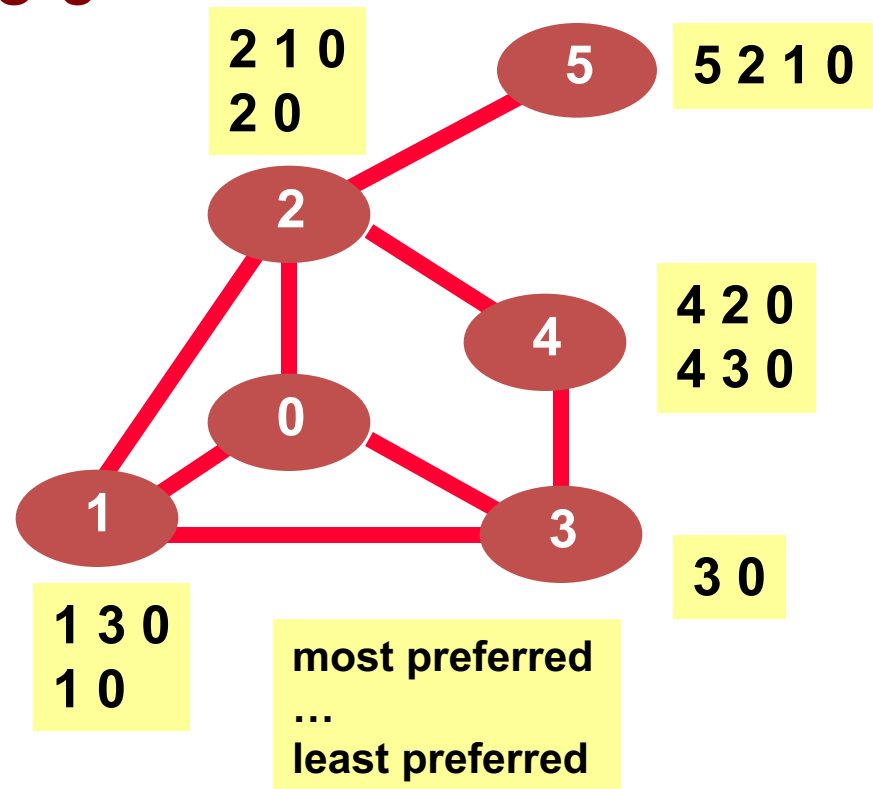  - Set of routes to 0 at each node
  - Ranking of the paths

2 1 0
2 0

5 2 1 0

4 2 0
4 3 0

3 0

1 3 0
1 0

**most preferred
...
least preferred**

# Stable Paths Problem (SPP) Instance

- ## 1 will use a direct path to 0
  (A) True    (B)  False

- ## 5 has a path to 0
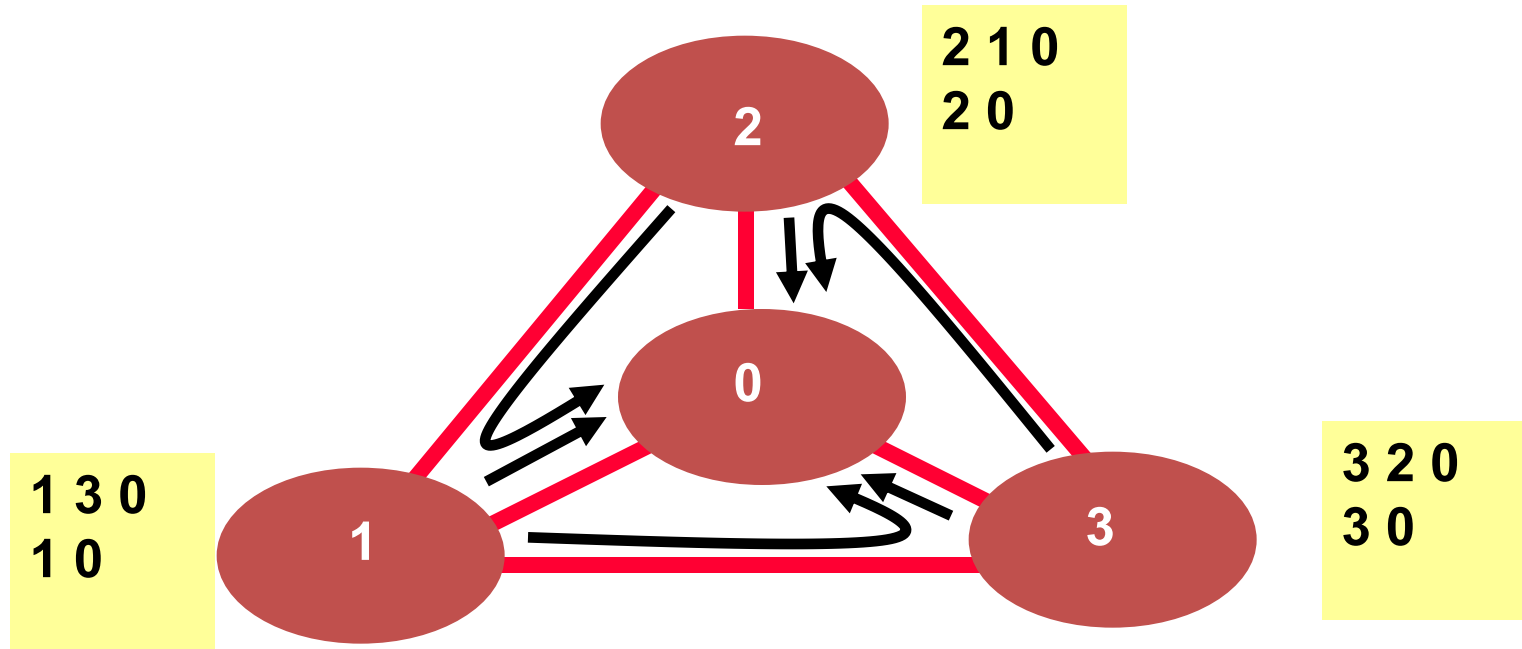  (A) True    (B)  False



2 1 0
2 0

5 2 1 0

4 2 0
4 3 0

3 0

1 3 0
1 0

most preferred
…
least preferred

# Stable Paths Problem (SPP) Instance



**2 1 0**
**2 0**

**5**

**5 2 1 0**

**2**

**4 2 0**
**4 3 0**

**4**

**0**

**1**

**3**

**3 0**

**1 3 0**
**1 0**

**most preferred**
**…**
**least preferred**

# SPP May Have Multiple Solutions

1 2 0
1 0

2 1 0
2 0

1 2 0
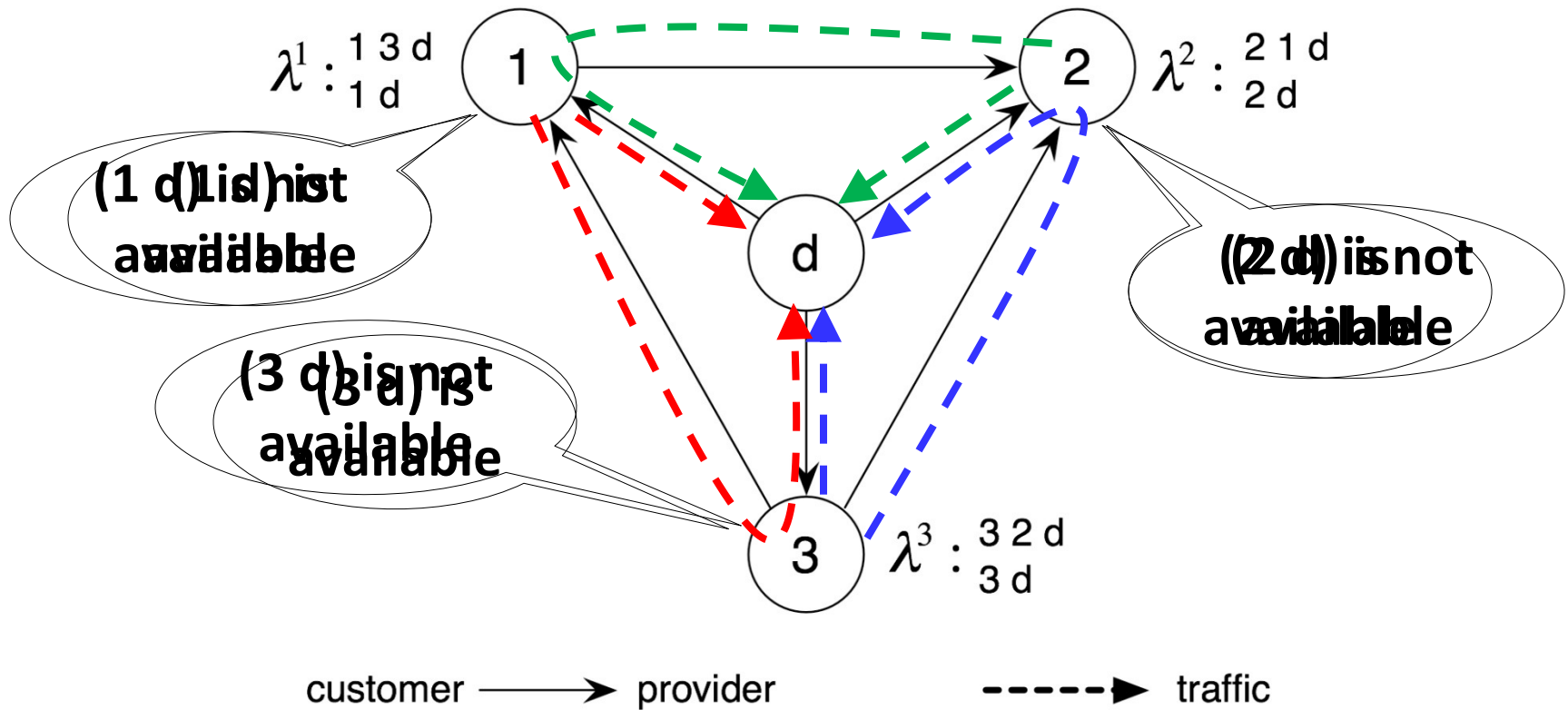1 0

2 1 0
2 0

**First solution**

1 2 0
1 0

2 1 0
2 0

**Second solution**

# An SPP May Have No Solution

# BGP Not Guaranteed to Converge



Example known as a "dispute wheel"

# 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

# 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