

CS 3700

Networks and Distributed Systems

Lecture 8: Inter Domain Routing

Revised 2/4/2014

Network Layer, Control Plane

2

Data Plane

Application

Presentation

Session

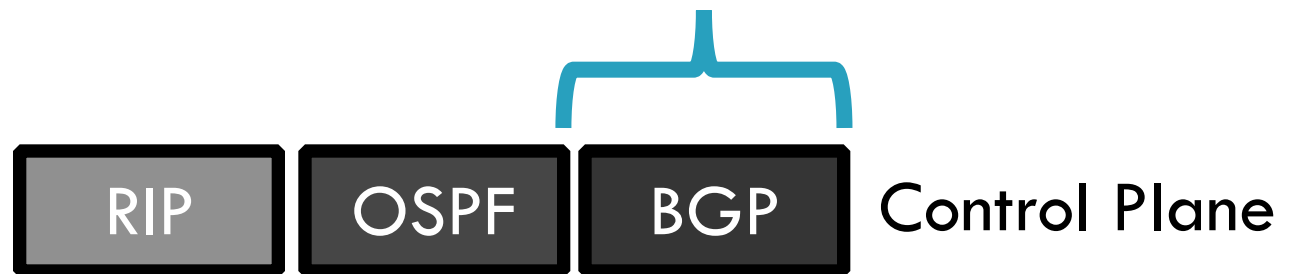
Transport

Network

Data Link

Physical

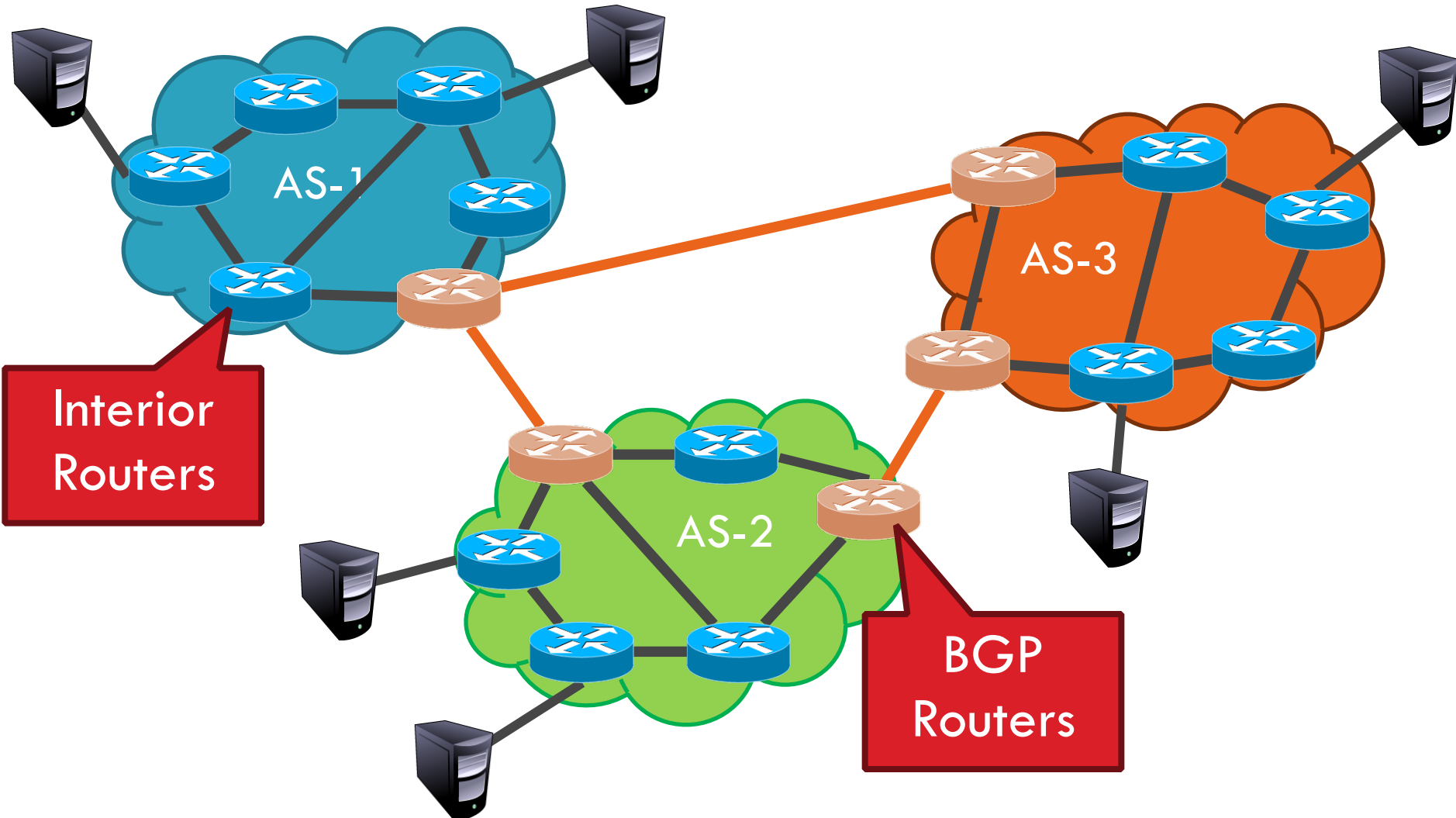
- Function:
 - ▣ Set up routes between networks
- Key challenges:
 - ▣ Implementing provider policies
 - ▣ Creating stable paths



- ❑ BGP Basics
- ❑ Stable Paths Problem
- ❑ BGP in the Real World

ASs, Revisited

4



AS Numbers

5

- Each AS identified by an ASN number
 - 16-bit values (latest protocol supports 32-bit ones)
 - 64512 – 65535 are reserved
- Currently, there are > 20000 ASNs
 - AT&T: 5074, 6341, 7018, ...
 - Sprint: 1239, 1240, 6211, 6242, ...
 - Northeastern: 156
 - North America ASs → <ftp://ftp.arin.net/info/asn.txt>

Inter-Domain Routing

6

- Global connectivity is at stake!
 - Thus, all ASs must use the same protocol
 - Contrast with intra-domain routing

Inter-Domain Routing

6

- Global connectivity is at stake!
 - ▣ Thus, all ASs must use the same protocol
 - ▣ Contrast with intra-domain routing
- What are the requirements?
 - ▣ Scalability
 - ▣ Flexibility in choosing routes
 - Cost
 - Routing around failures

Inter-Domain Routing

6

- Global connectivity is at stake!
 - ▣ Thus, all ASs must use the same protocol
 - ▣ Contrast with intra-domain routing
- What are the requirements?
 - ▣ Scalability
 - ▣ Flexibility in choosing routes
 - Cost
 - Routing around failures
- Question: link state or distance vector?

Inter-Domain Routing

6

- Global connectivity is at stake!
 - ▣ Thus, all ASs must use the same protocol
 - ▣ Contrast with intra-domain routing
- What are the requirements?
 - ▣ Scalability
 - ▣ Flexibility in choosing routes
 - Cost
 - Routing around failures
- Question: link state or distance vector?
 - ▣ Trick question: BGP is a **path vector** protocol

BGP

7

- Border Gateway Protocol
 - De facto inter-domain protocol of the Internet
 - Policy based routing protocol
 - Uses a Bellman-Ford path vector protocol

BGP

7

- Border Gateway Protocol
 - ▣ De facto inter-domain protocol of the Internet
 - ▣ Policy based routing protocol
 - ▣ Uses a Bellman-Ford path vector protocol
- Relatively simple protocol, but...
 - ▣ Complex, manual configuration

BGP

7

- Border Gateway Protocol
 - ▣ De facto inter-domain protocol of the Internet
 - ▣ Policy based routing protocol
 - ▣ Uses a Bellman-Ford path vector protocol
- Relatively simple protocol, but...
 - ▣ Complex, manual configuration
 - ▣ Entire world sees advertisements
 - Errors can screw up traffic globally

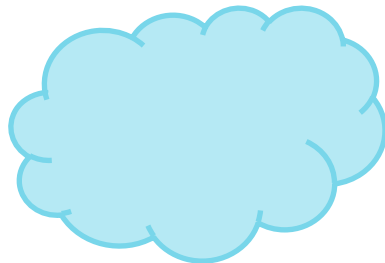
BGP

7

- Border Gateway Protocol
 - ▣ De facto inter-domain protocol of the Internet
 - ▣ Policy based routing protocol
 - ▣ Uses a Bellman-Ford path vector protocol
- Relatively simple protocol, but...
 - ▣ Complex, manual configuration
 - ▣ Entire world sees advertisements
 - Errors can screw up traffic globally
 - ▣ Policies driven by economics
 - How much \$\$\$ does it cost to route along a given path?
 - Not by performance (e.g. shortest paths)

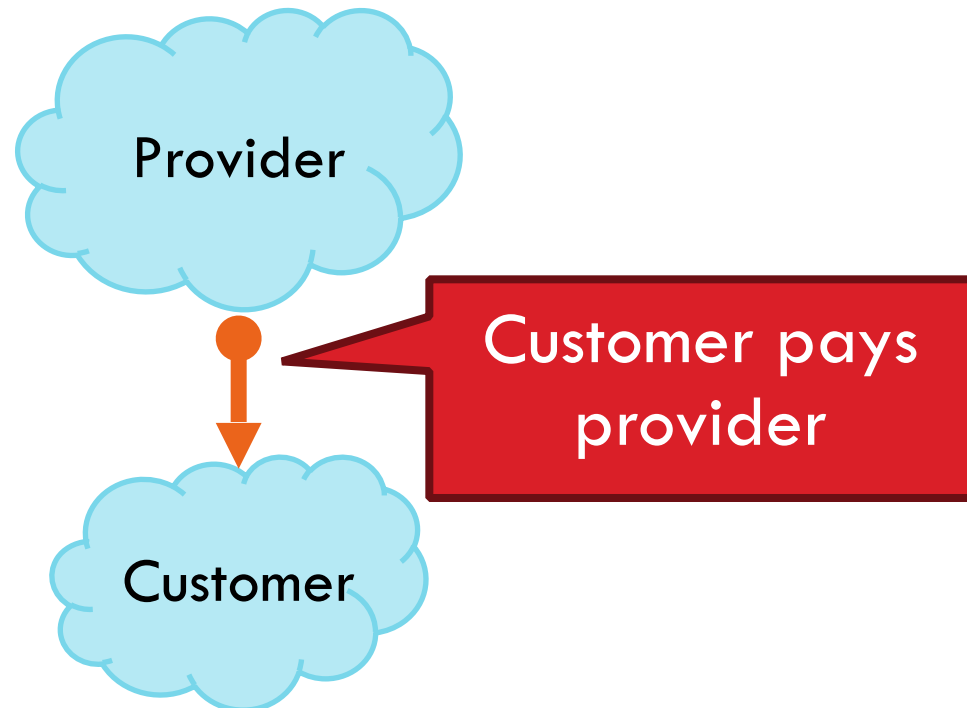
BGP Relationships

8



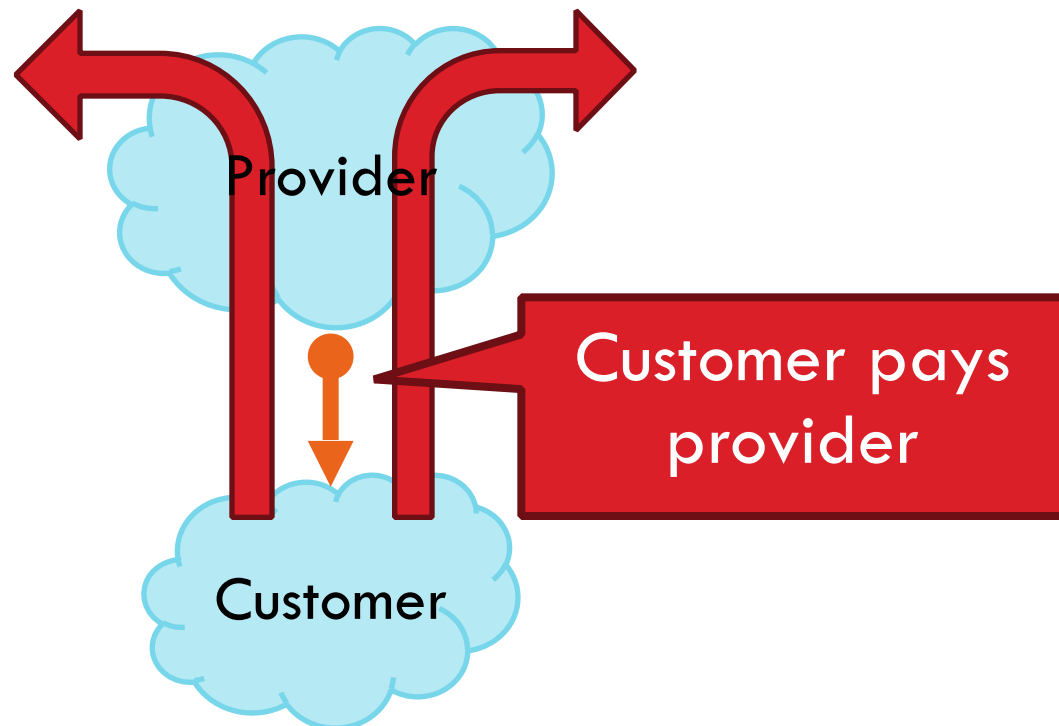
BGP Relationships

8



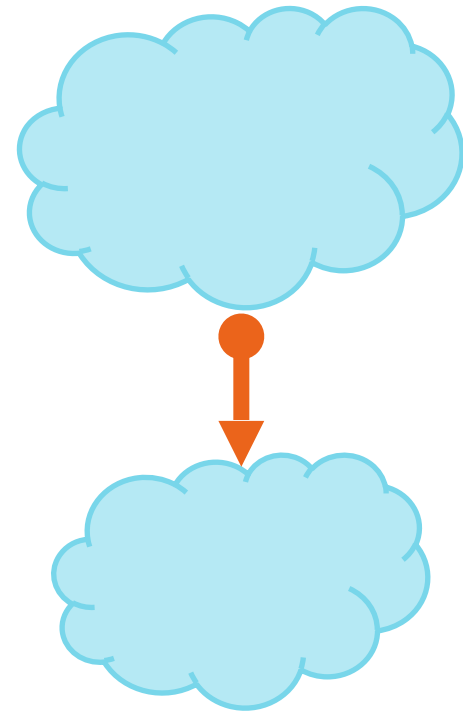
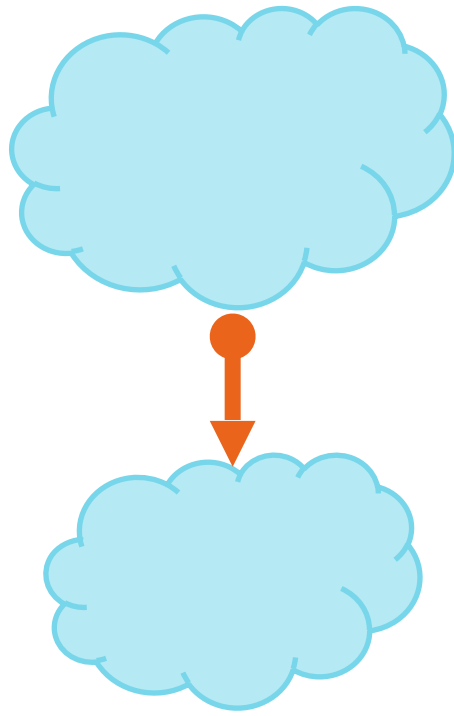
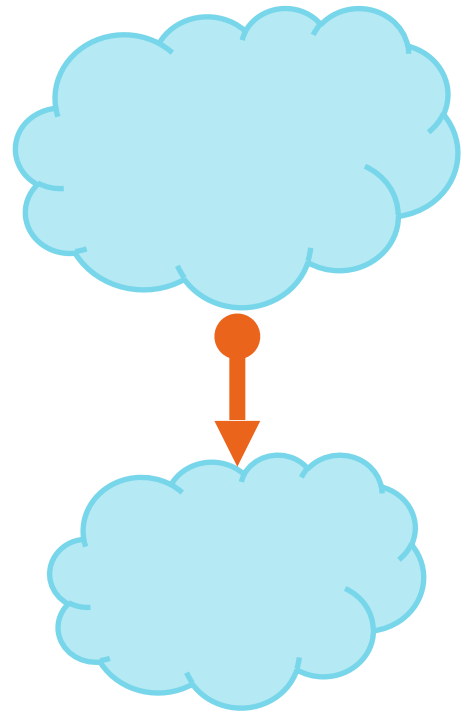
BGP Relationships

8



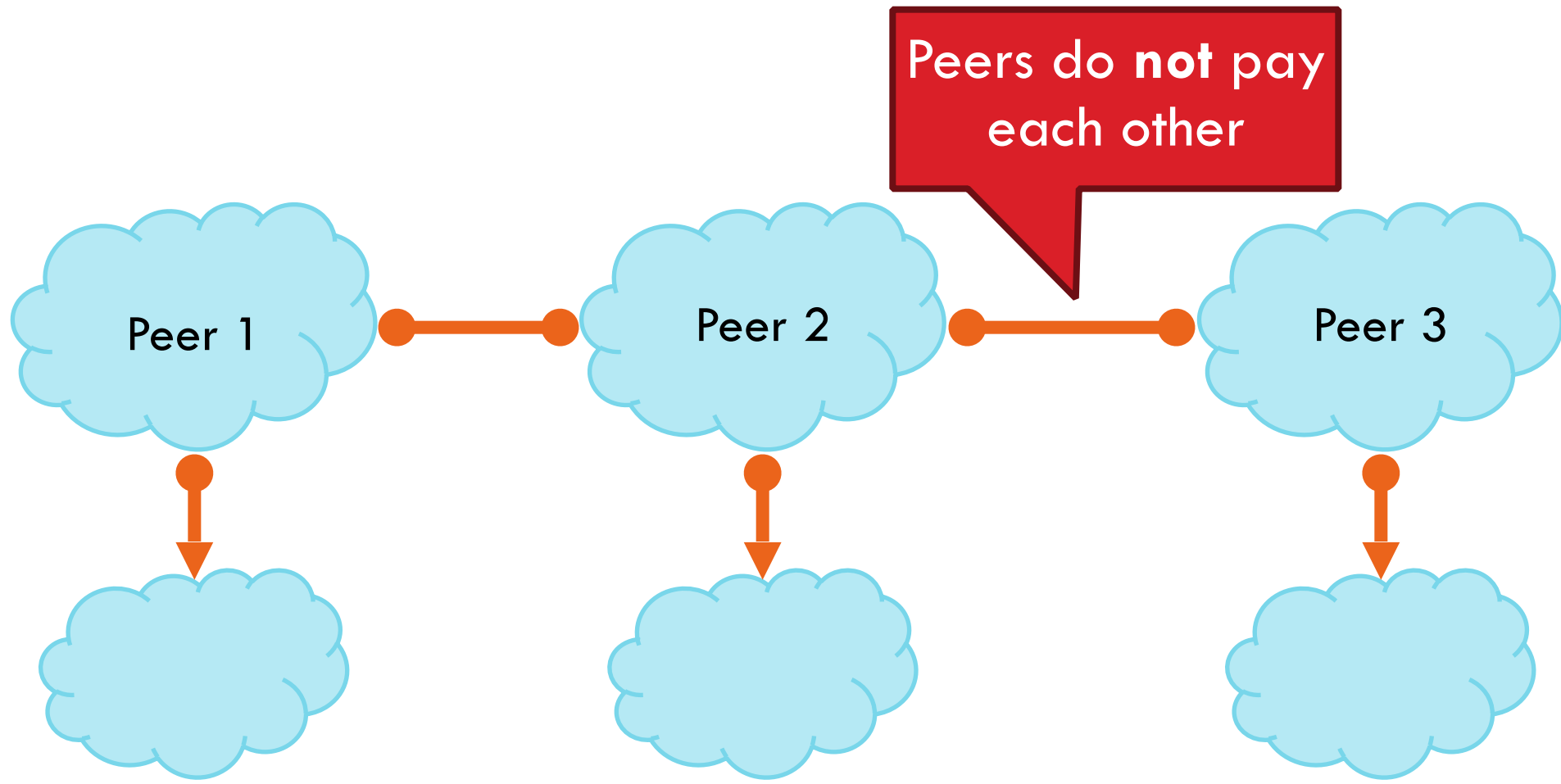
BGP Relationships

8



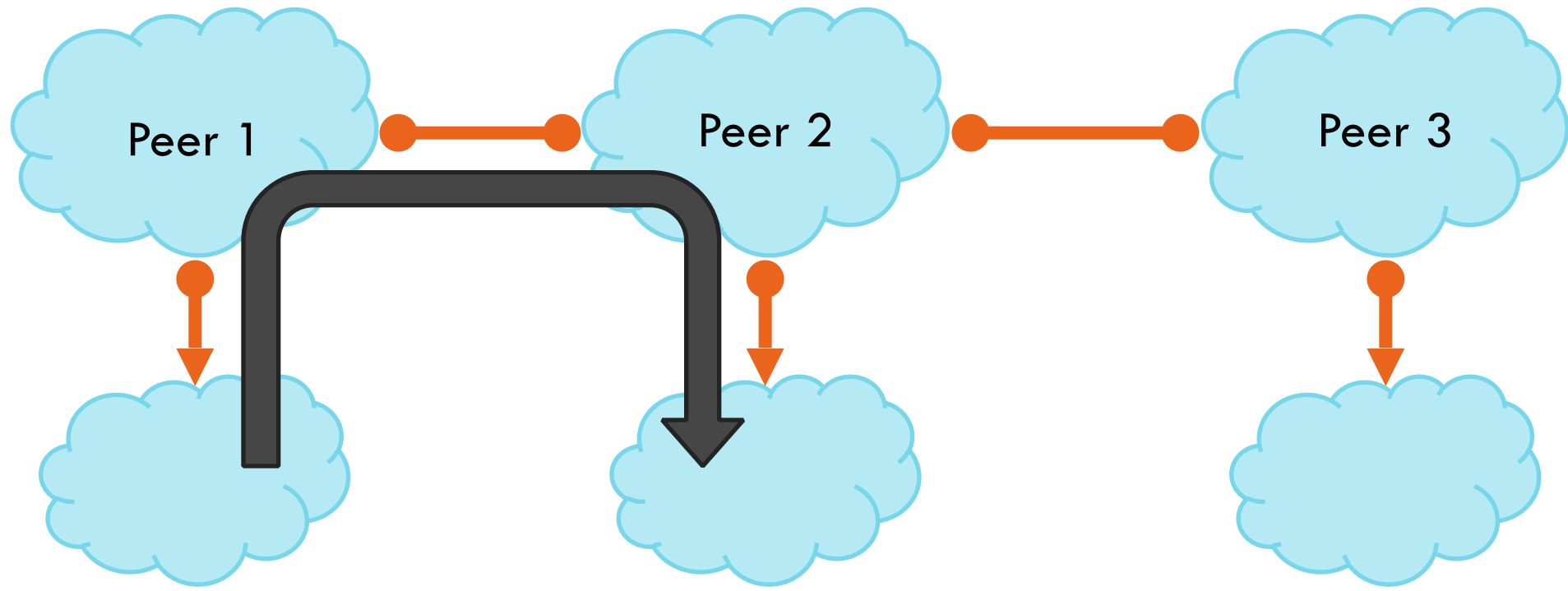
BGP Relationships

8



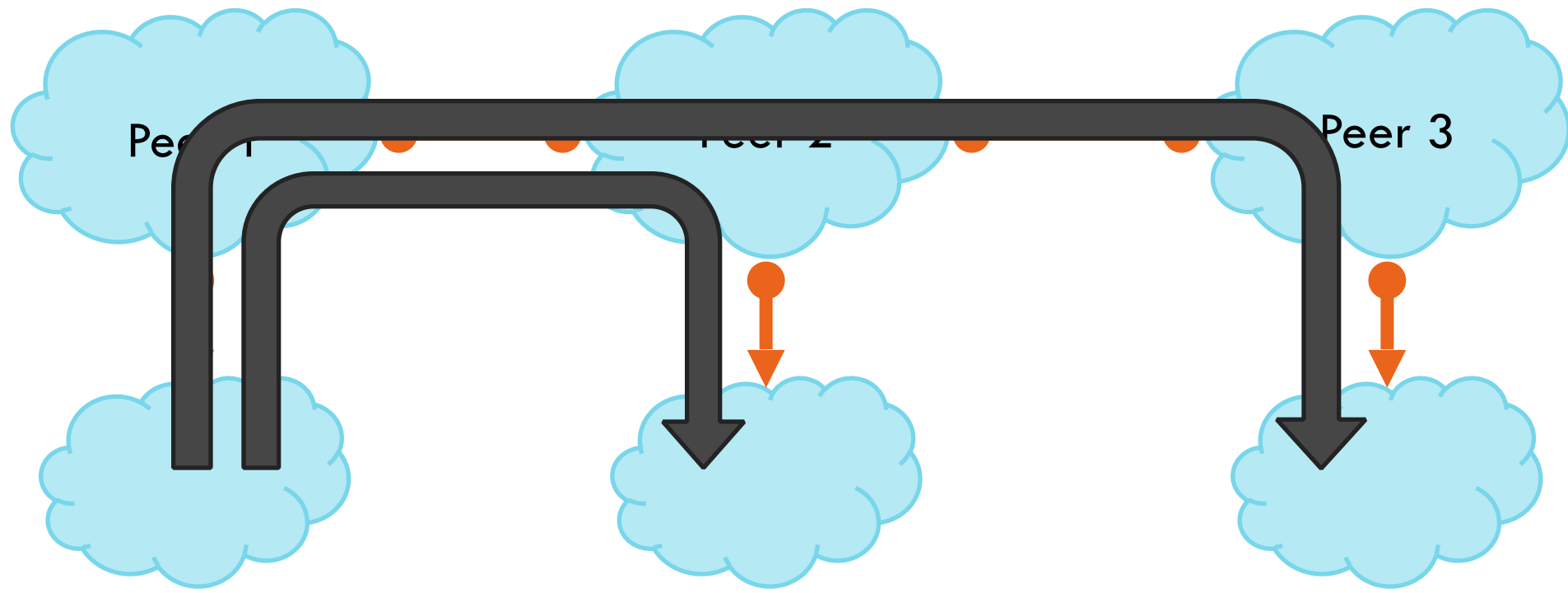
BGP Relationships

8



BGP Relationships

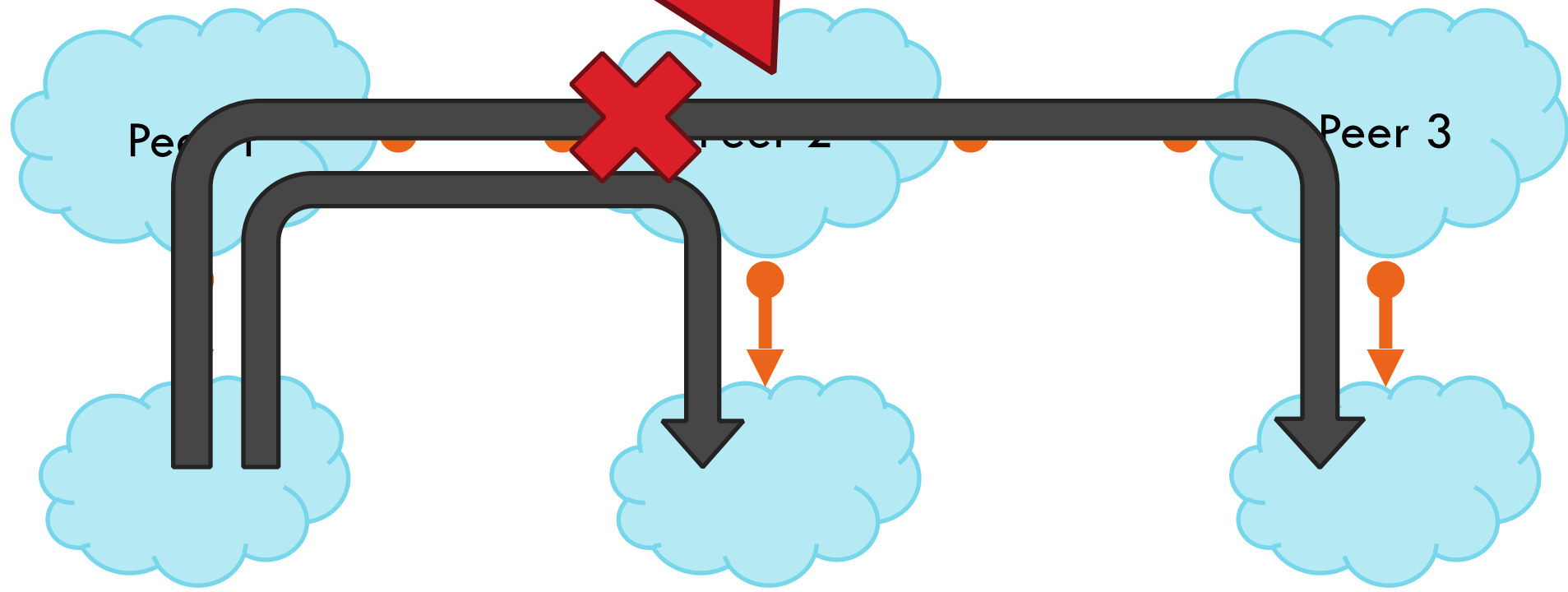
8



BGP Relationships

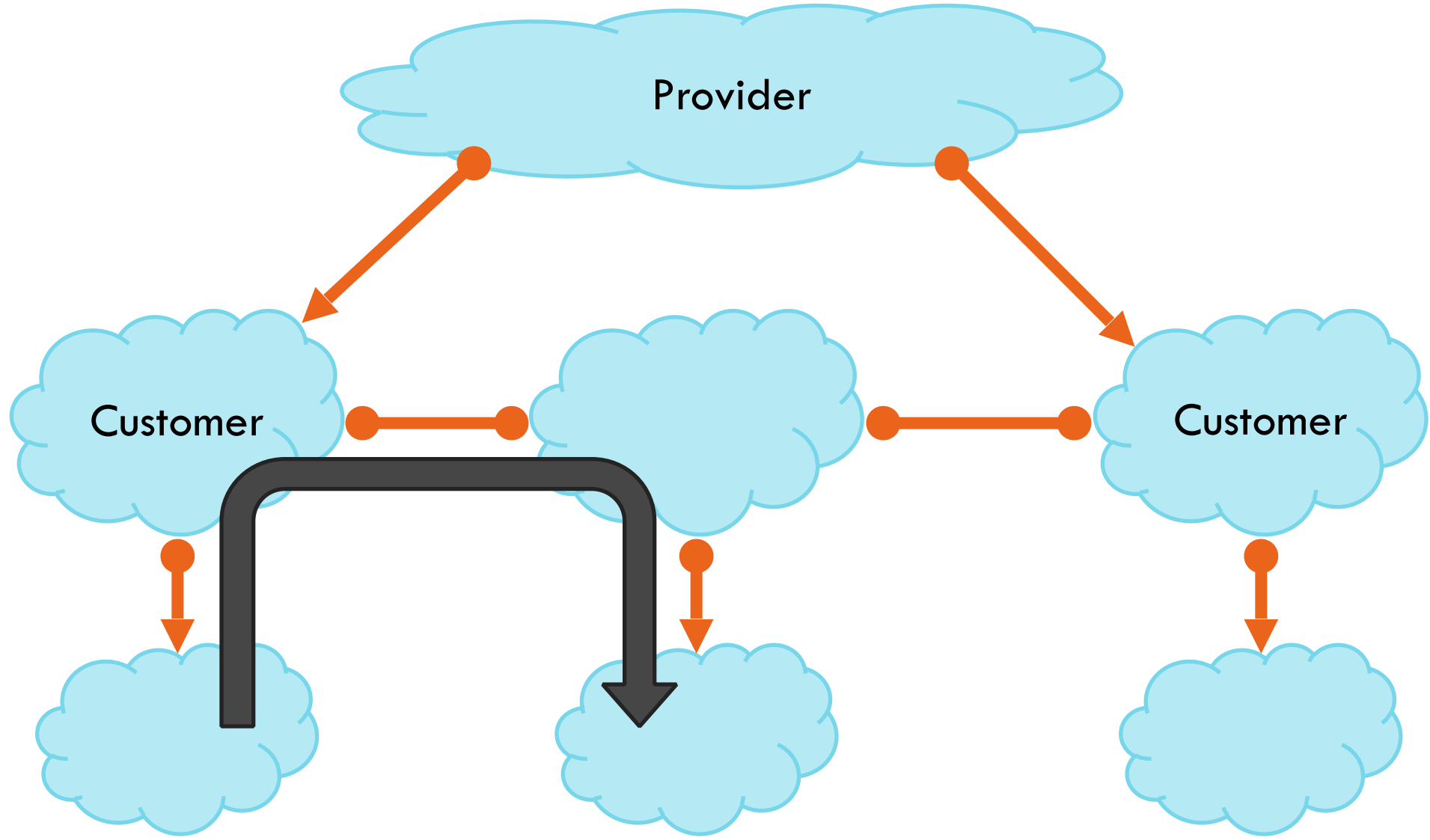
8

Peer 2 has no incentive to route 1 → 3



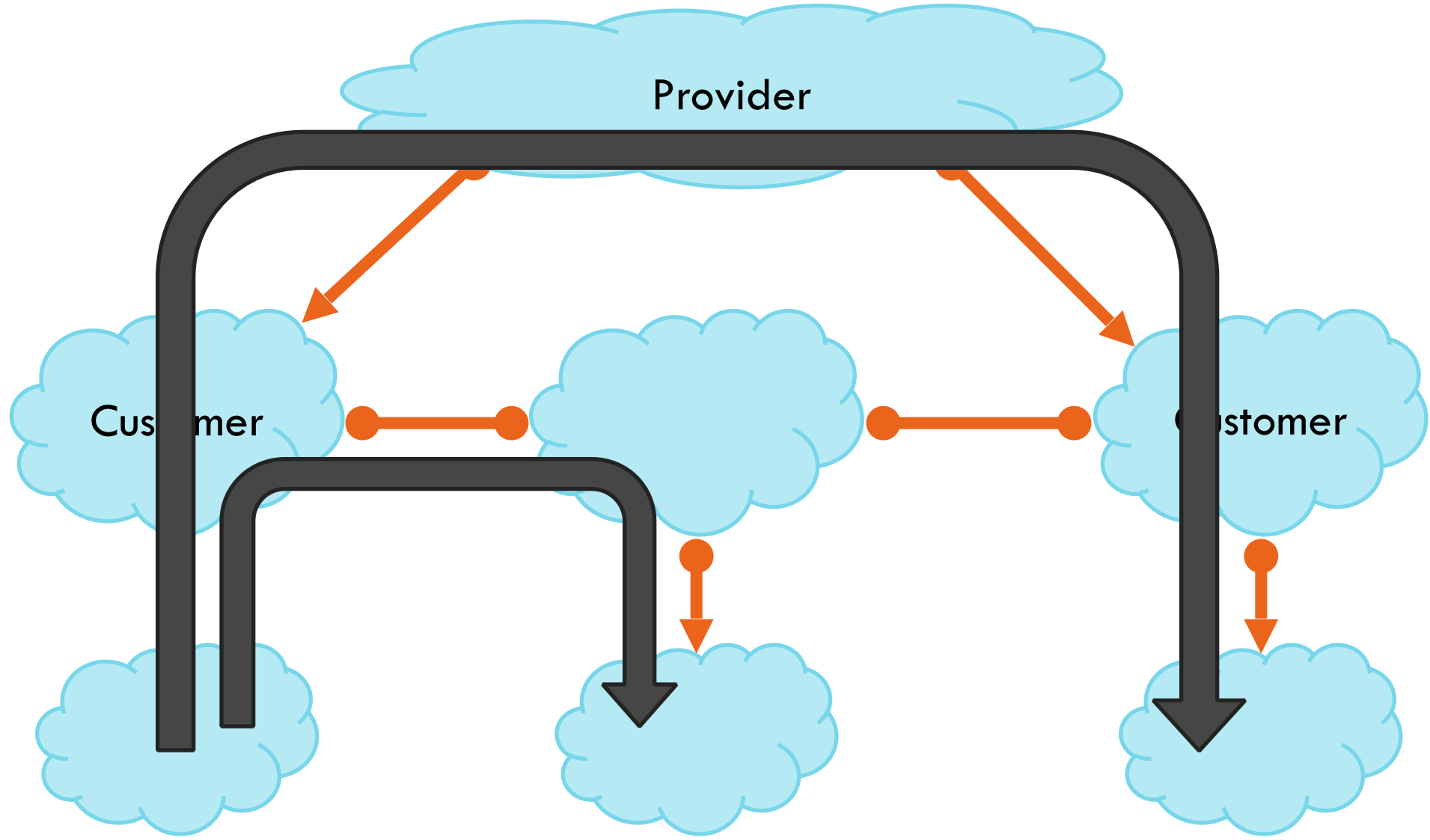
BGP Relationships

8



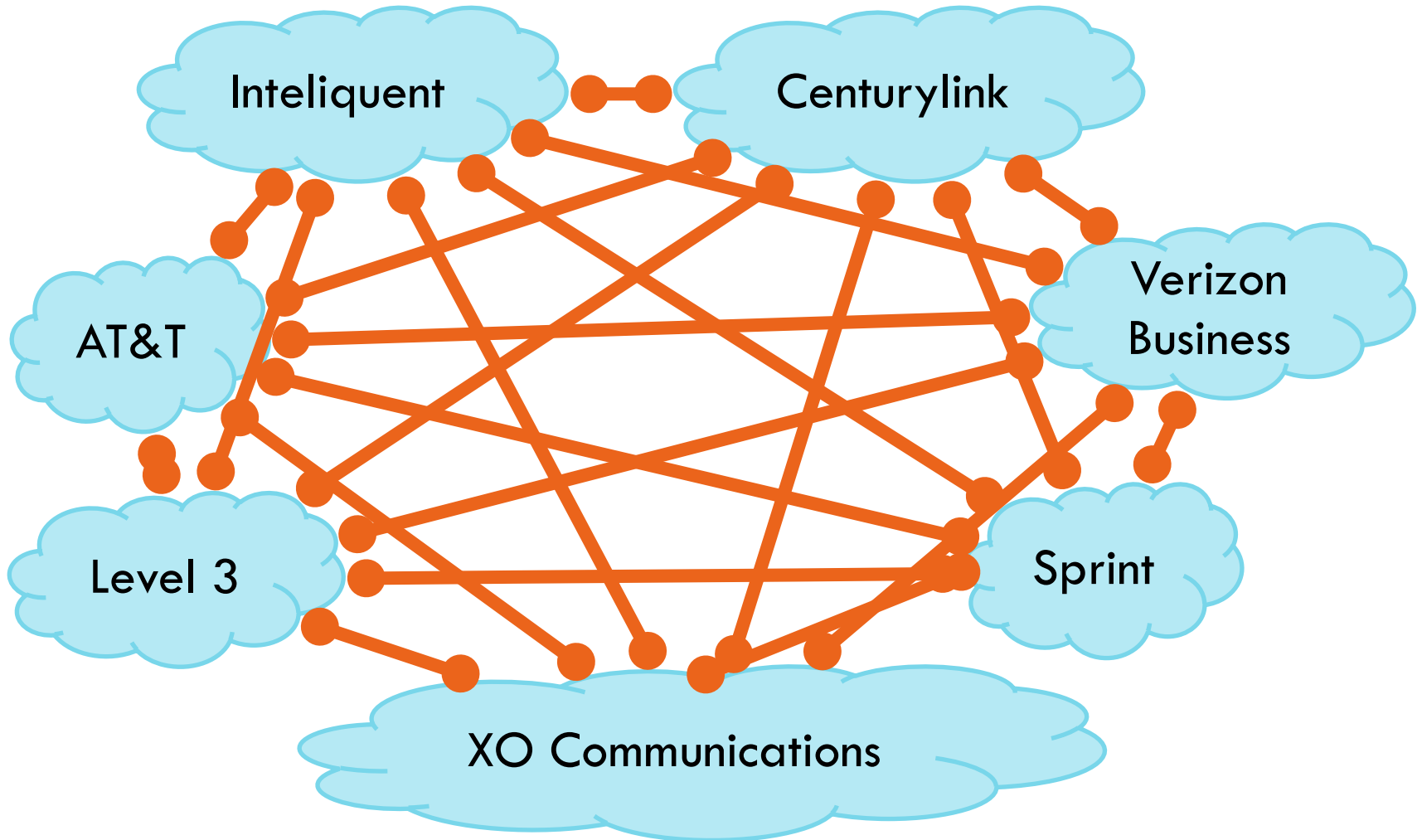
BGP Relationships

8



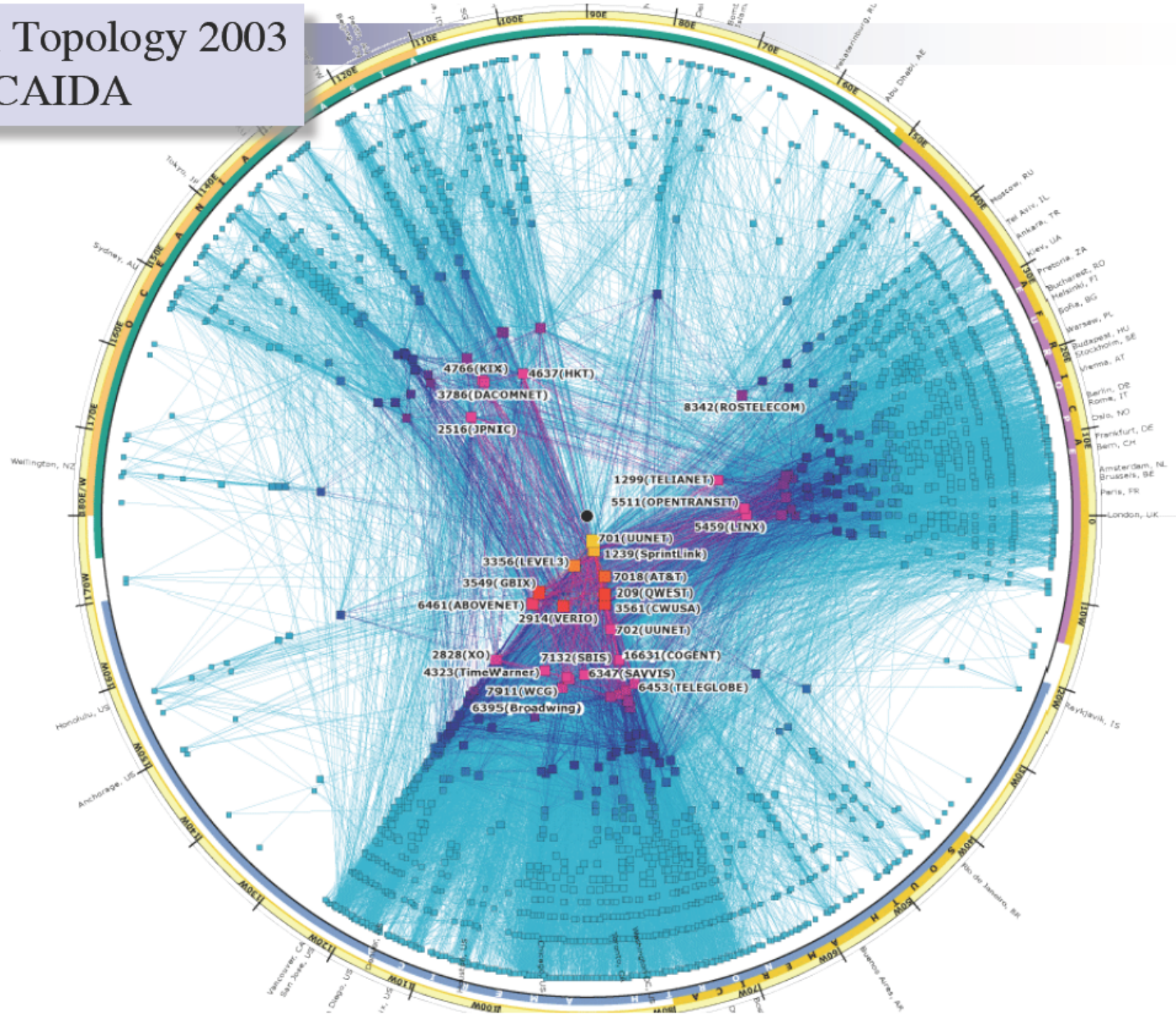
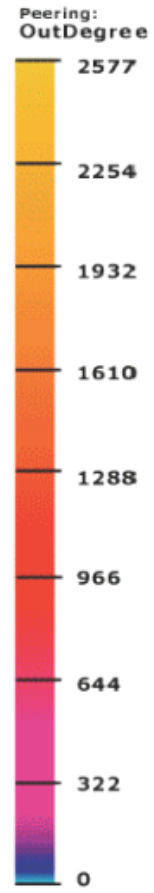
Tier-1 ISP Peering

9



AS-level Topology 2003

Source: CAIDA



Peering Wars

11

Peer

- Reduce upstream costs
- Improve end-to-end performance
- May be the only way to connect to parts of the Internet

Don't Peer

- You would rather have customers
- Peers are often competitors
- Peering agreements require periodic renegotiation

Peering Wars

11

Peer

- Reduce upstream costs
- Improve end-to-end performance
- May be the only way to connect to parts of the Internet

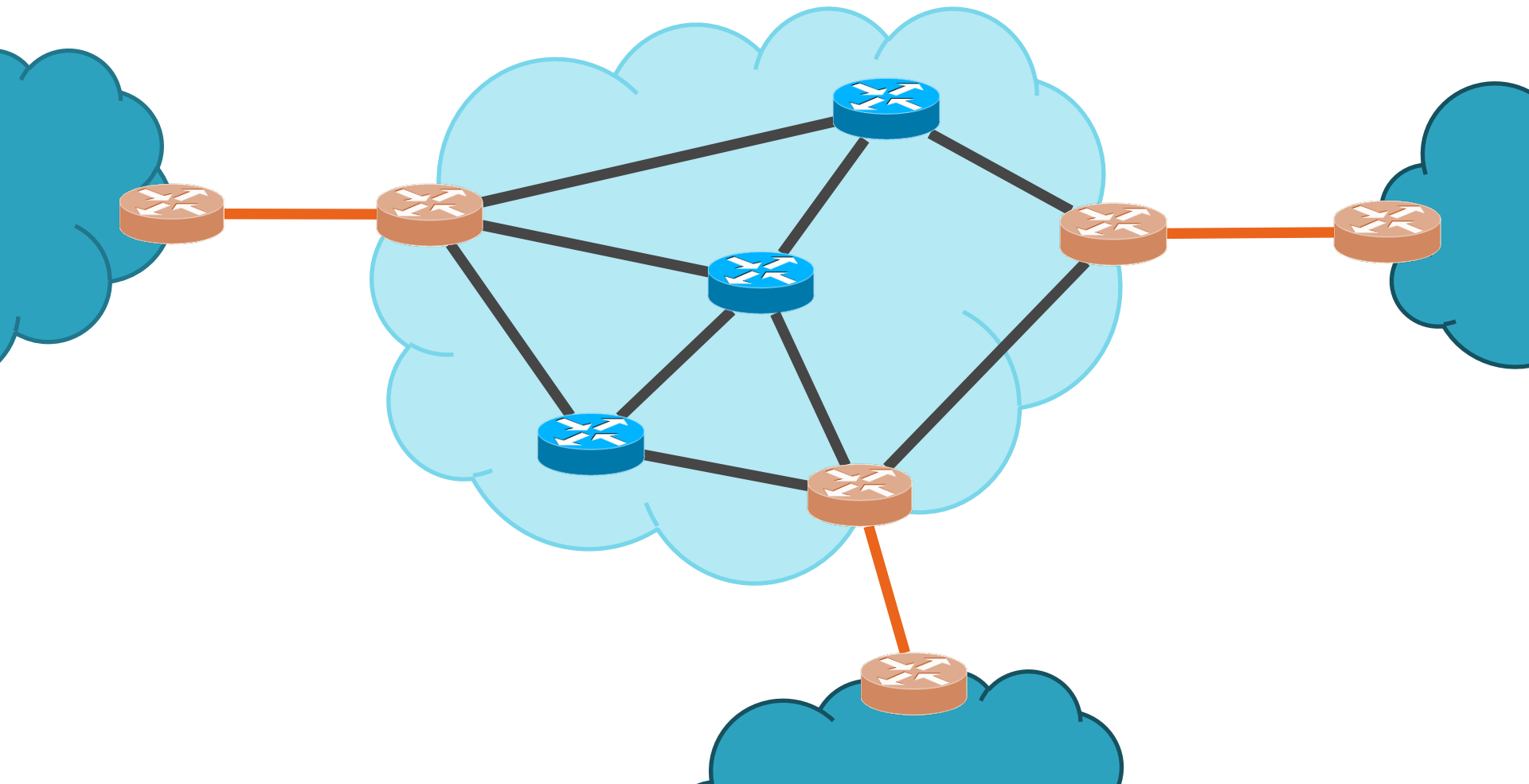
Don't Peer

- You would rather have customers
- Peers are often competitors
- Peering agreements require periodic renegotiation

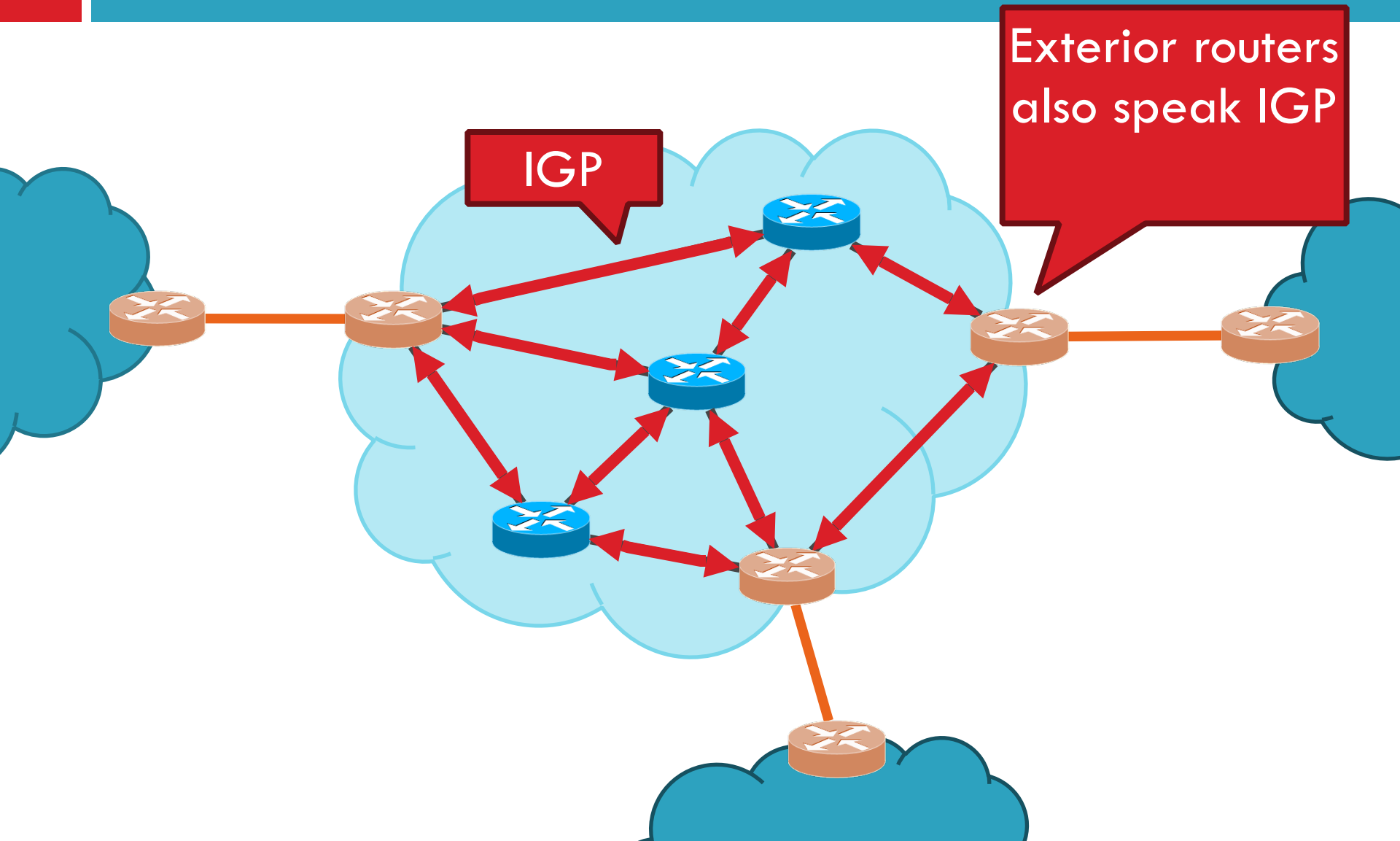
Peering struggles in the ISP world are extremely contentious, agreements are usually confidential

Two Types of BGP Neighbors

12

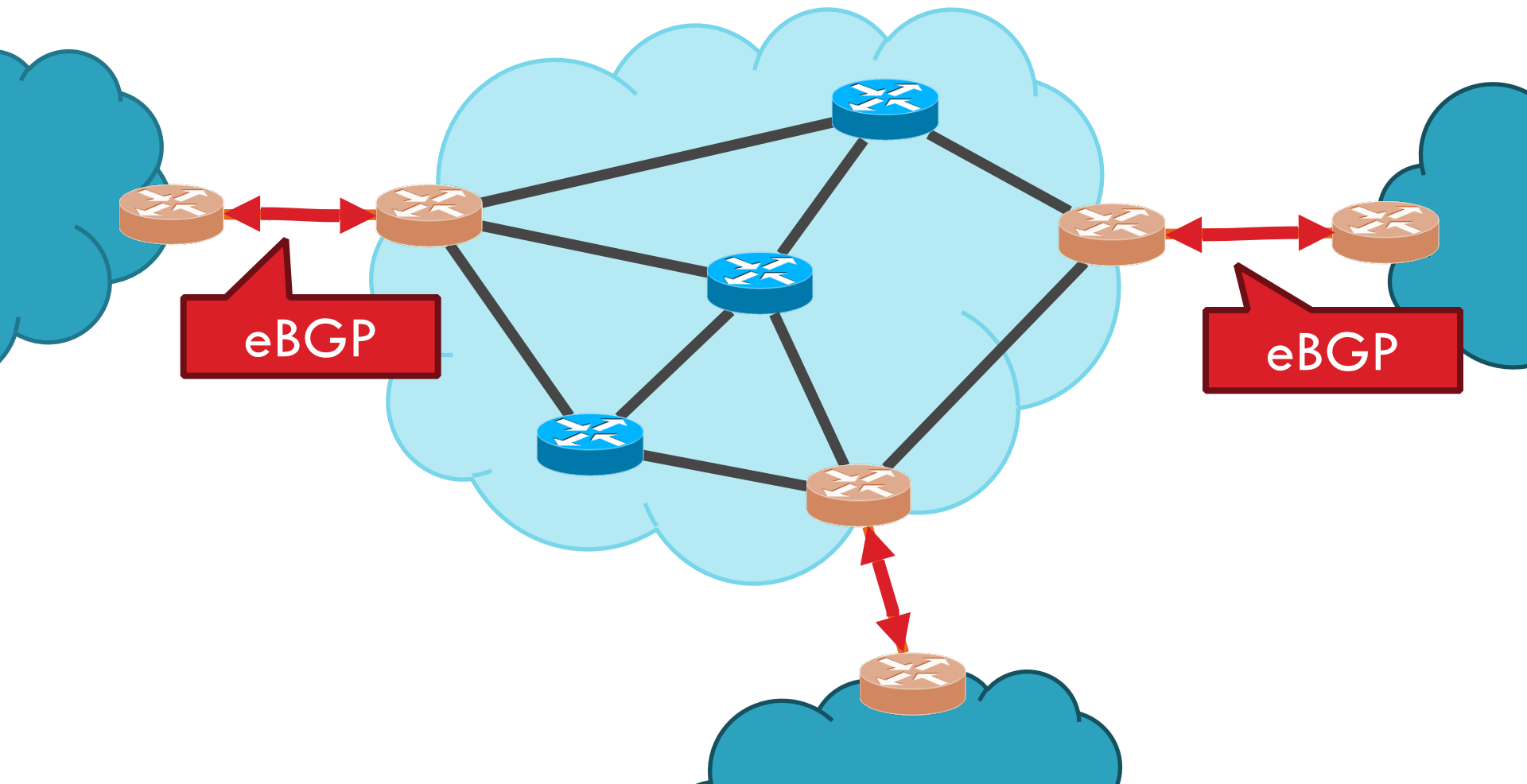


Two Types of BGP Neighbors

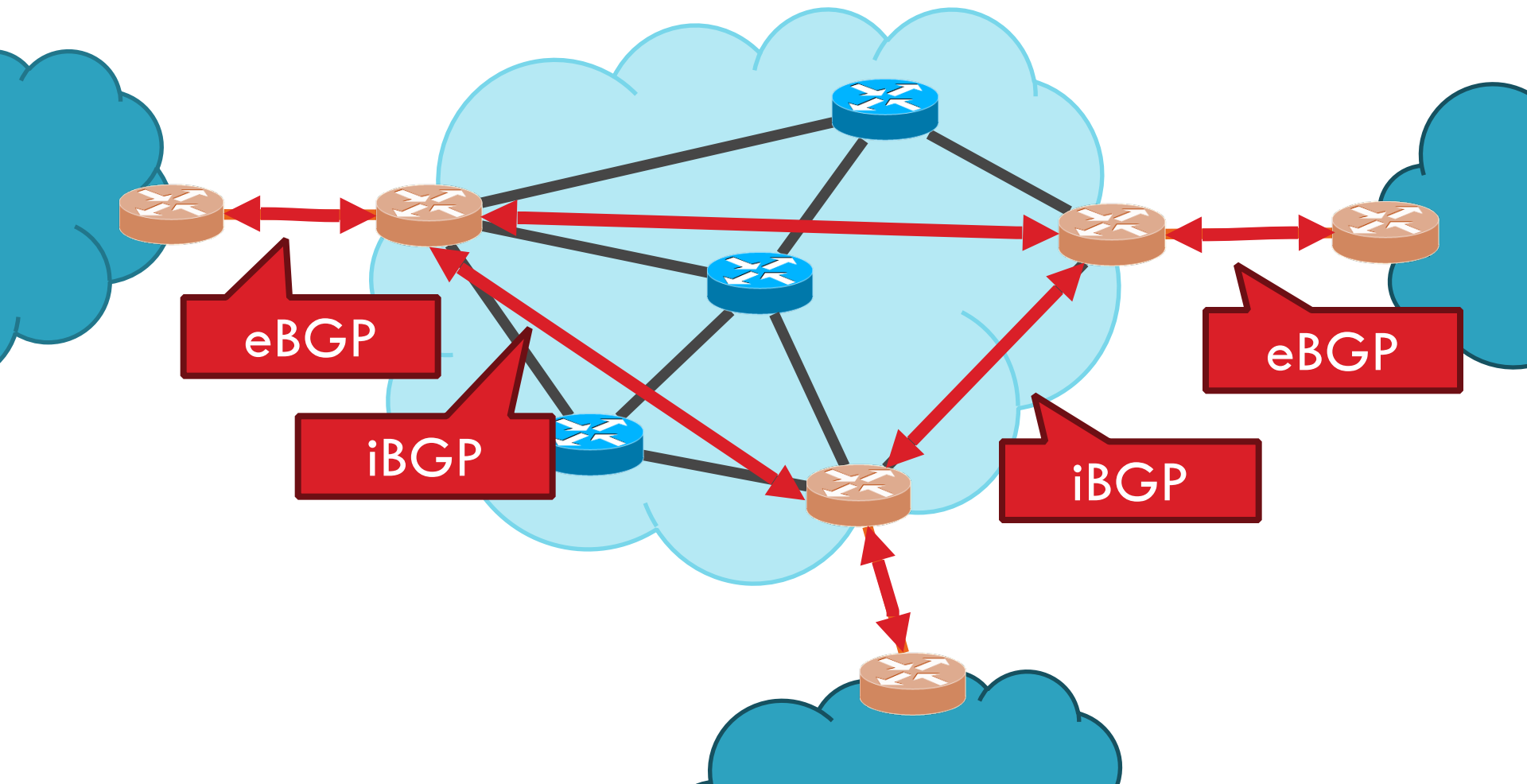


Two Types of BGP Neighbors

12

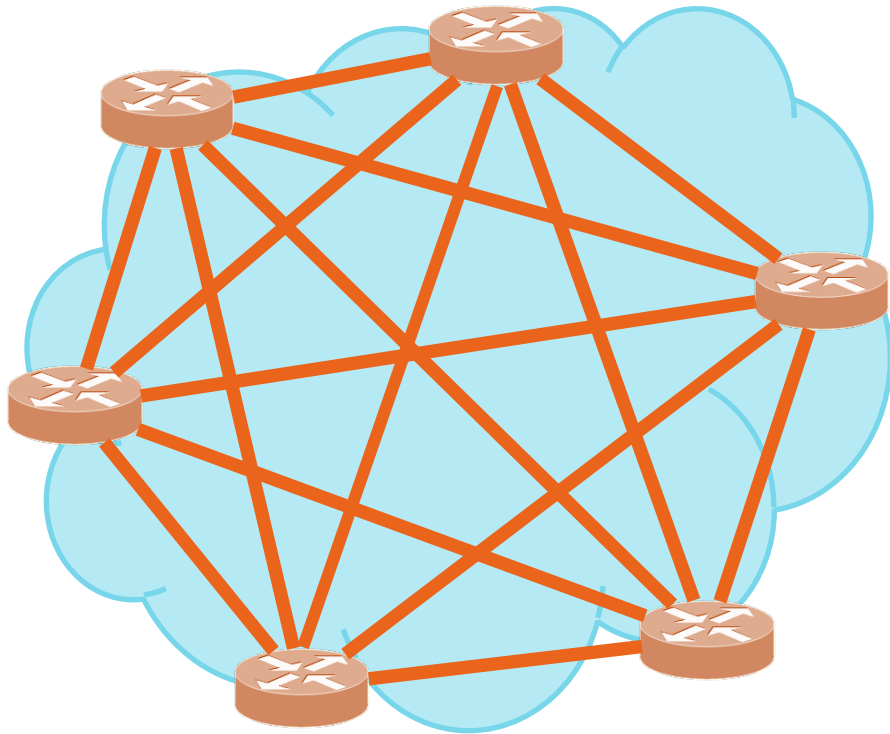


Two Types of BGP Neighbors



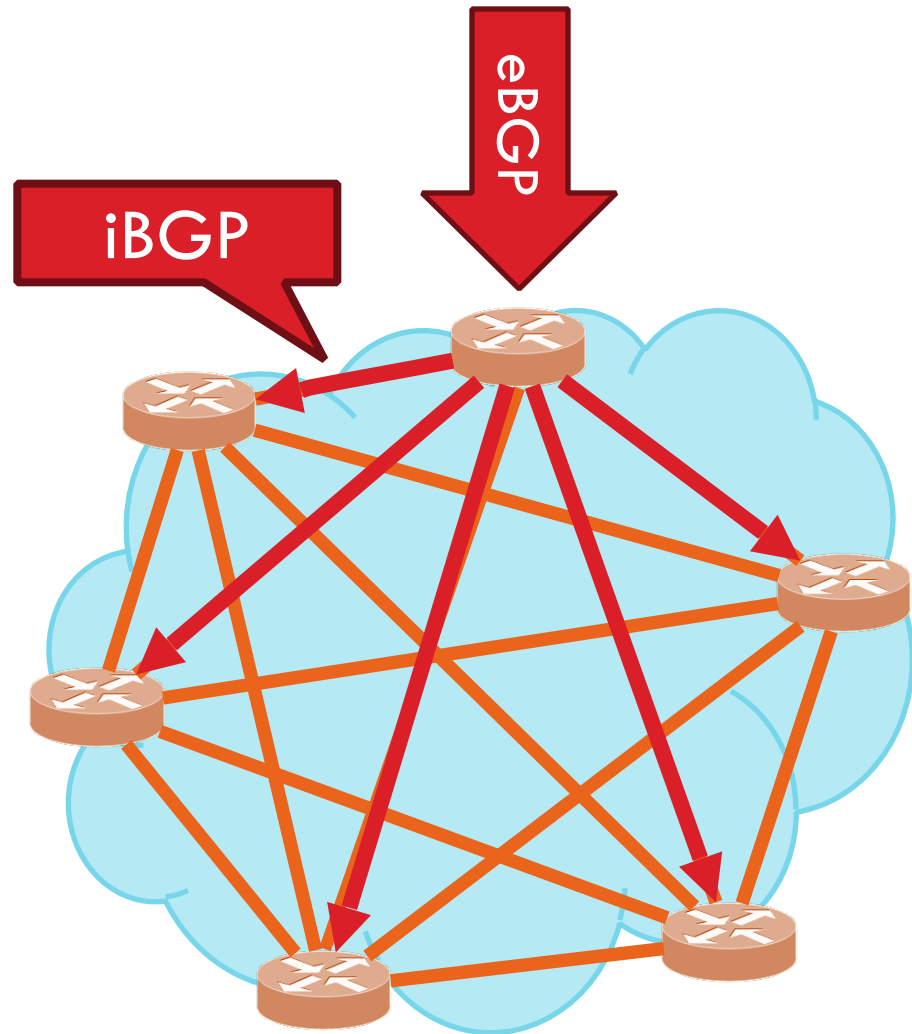
Full iBGP Meshes

13



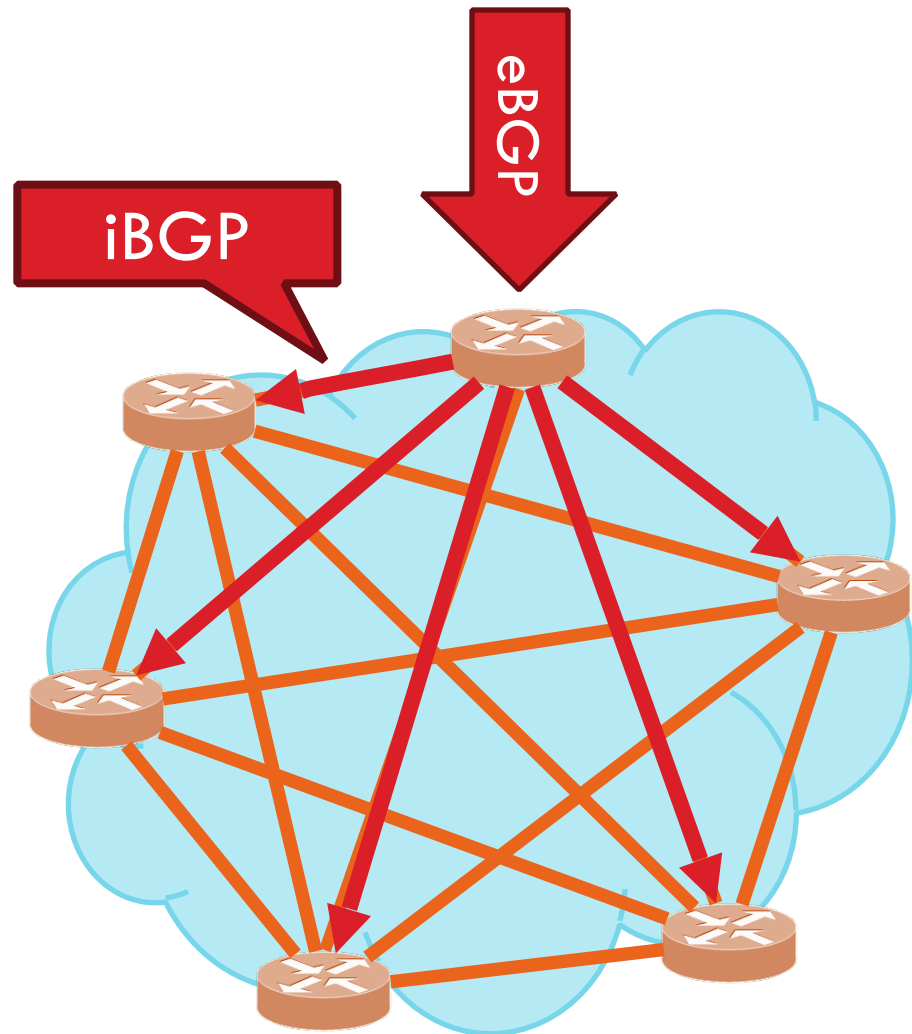
Full iBGP Meshes

13



Full iBGP Meshes

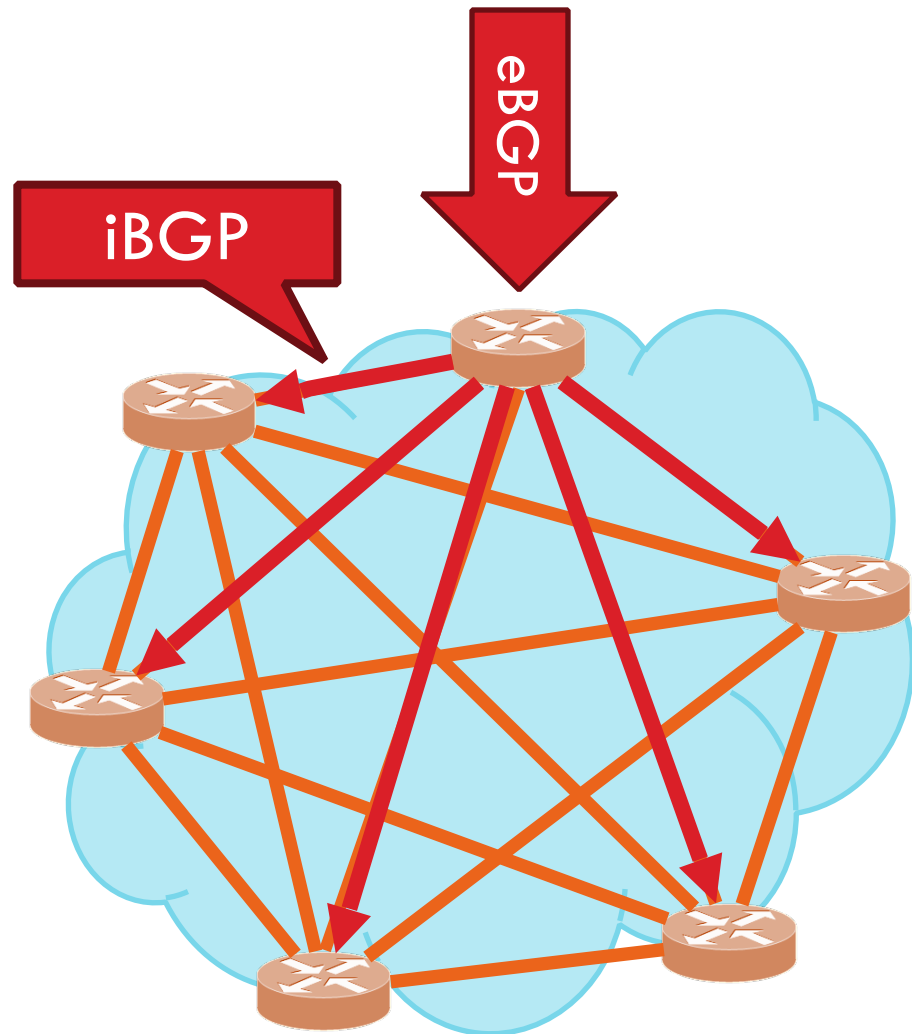
13



- Question: why do we need iBGP?
 - ▣ OSPF does not include BGP policy info
 - ▣ Prevents routing loops within the AS

Full iBGP Meshes

13

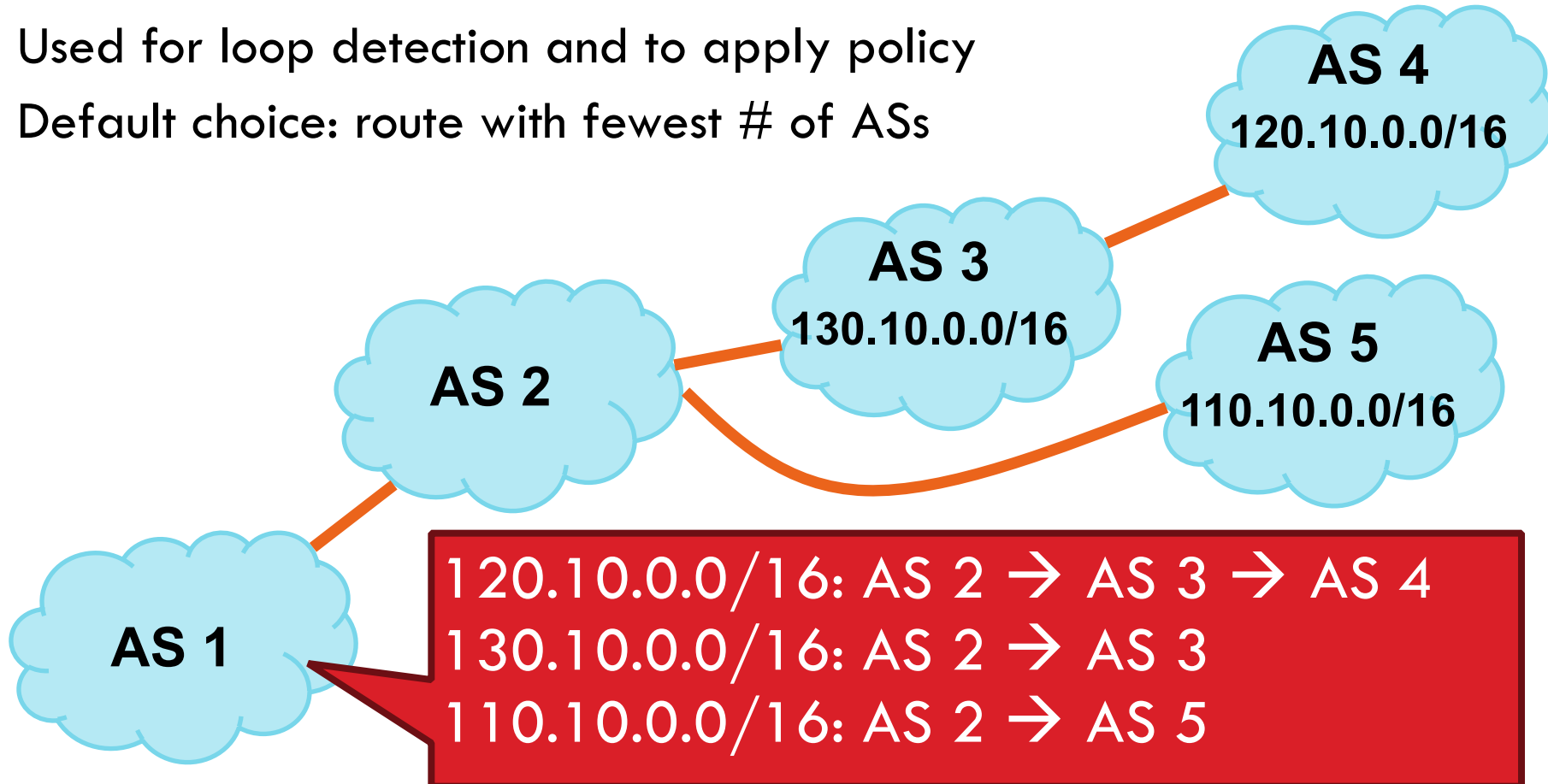


- Question: why do we need iBGP?
 - ▣ OSPF does not include BGP policy info
 - ▣ Prevents routing loops within the AS
- iBGP updates do not trigger announcements

Path Vector Protocol

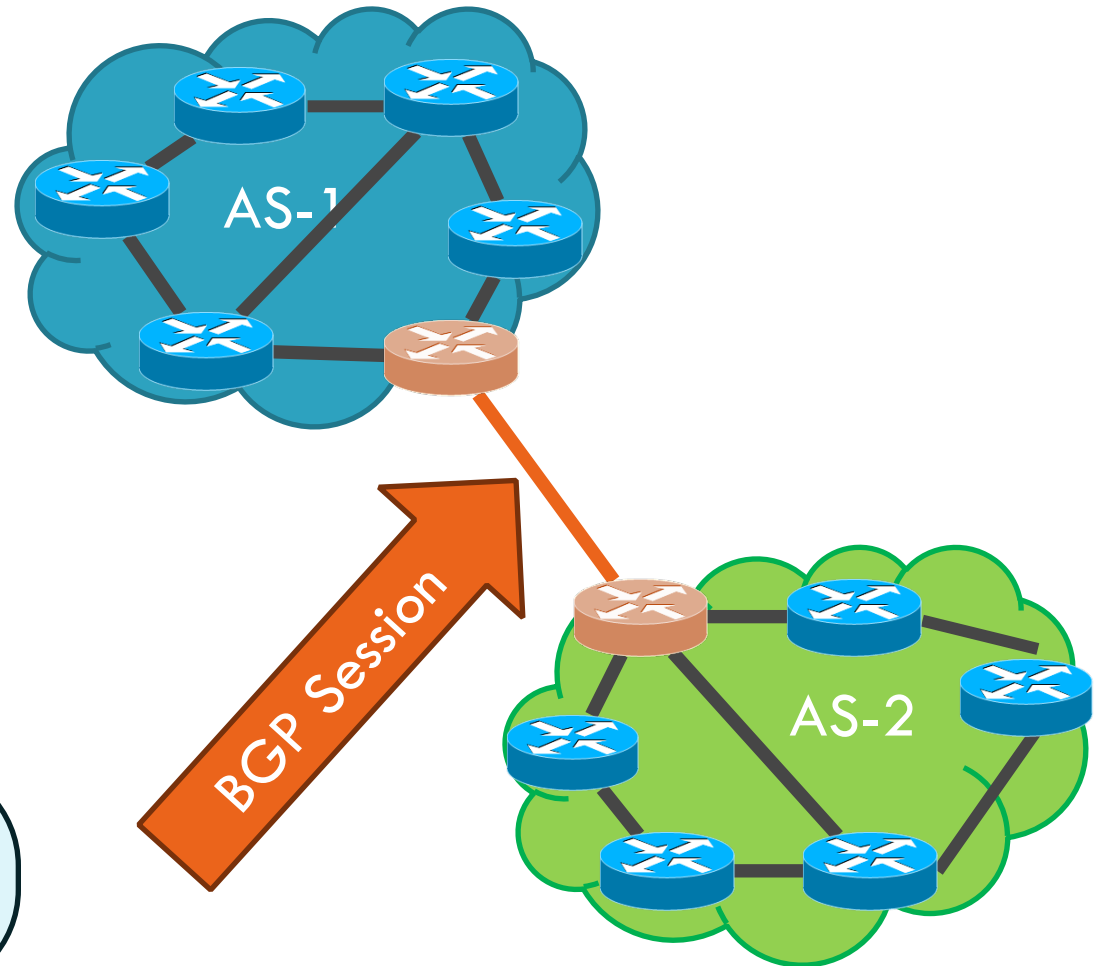
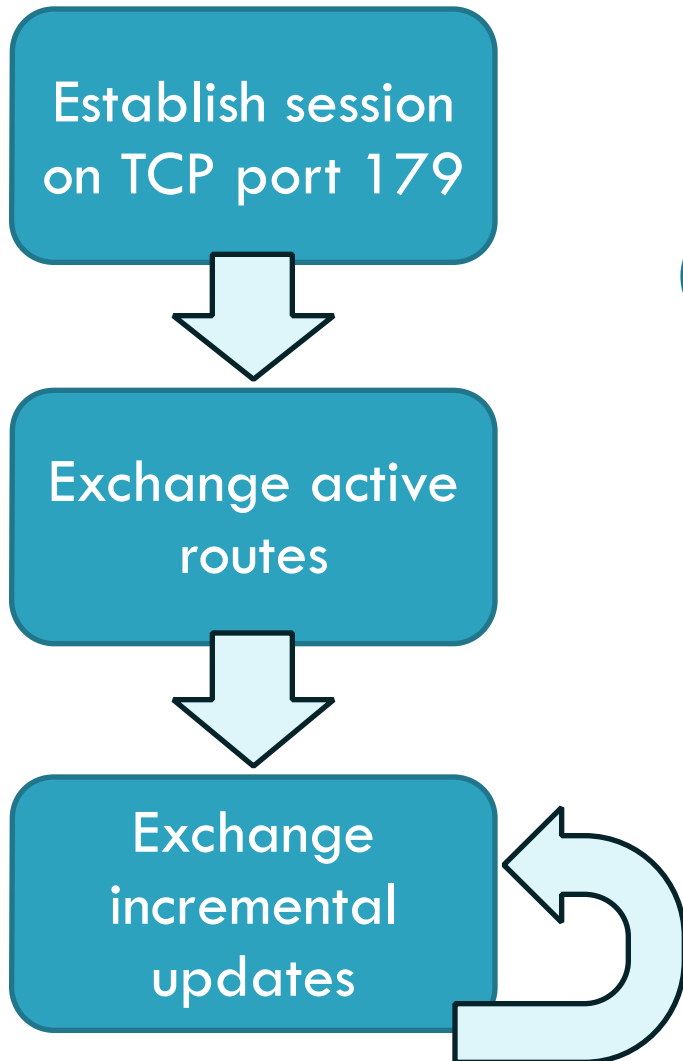
14

- AS-path: sequence of ASs a route traverses
 - ▣ Like distance vector, plus additional information
- Used for loop detection and to apply policy
- Default choice: route with fewest # of ASs



BGP Operations (Simplified)

15



Four Types of BGP Messages

16

- **Open:** Establish a peering session.
- **Keep Alive:** Handshake at regular intervals.
- **Notification:** Shuts down a peering session.
- **Update:** Announce new routes or withdraw previously announced routes.

Four Types of BGP Messages

16

- **Open**: Establish a peering session.
- **Keep Alive**: Handshake at regular intervals.
- **Notification**: Shuts down a peering session.
- **Update**: Announce new routes or withdraw previously announced routes.

announcement = IP prefix + attributes values

BGP Attributes

17

- Attributes used to select “best” path
 - ▣ LocalPref
 - Local preference policy to choose most preferred route
 - Overrides default fewest AS behavior

BGP Attributes

17

- Attributes used to select “best” path
 - LocalPref
 - Local preference policy to choose most preferred route
 - Overrides default fewest AS behavior
 - Multi-exit Discriminator (MED)
 - Specifies path for external traffic destined for an internal network
 - Chooses peering point for your network

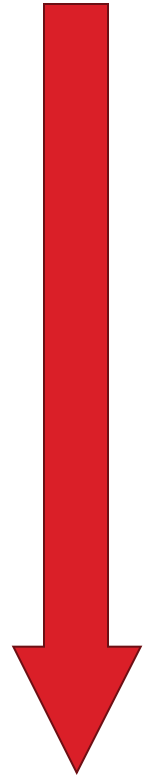
BGP Attributes

17

- Attributes used to select “best” path
 - LocalPref
 - Local preference policy to choose most preferred route
 - Overrides default fewest AS behavior
 - Multi-exit Discriminator (MED)
 - Specifies path for external traffic destined for an internal network
 - Chooses peering point for your network
 - Import Rules
 - What route advertisements do I accept?
 - Export Rules
 - Which routes do I forward to whom?

Route Selection Summary

18



Route Selection Summary

18



Highest Local Preference

Enforce relationships

Route Selection Summary

18



Highest Local Preference

Enforce relationships

Shortest AS Path

Lowest MED

Lowest IGP Cost to BGP Egress

Traffic engineering

Route Selection Summary

18



Highest Local Preference

Enforce relationships

Shortest AS Path

Lowest MED

Lowest IGP Cost to BGP Egress

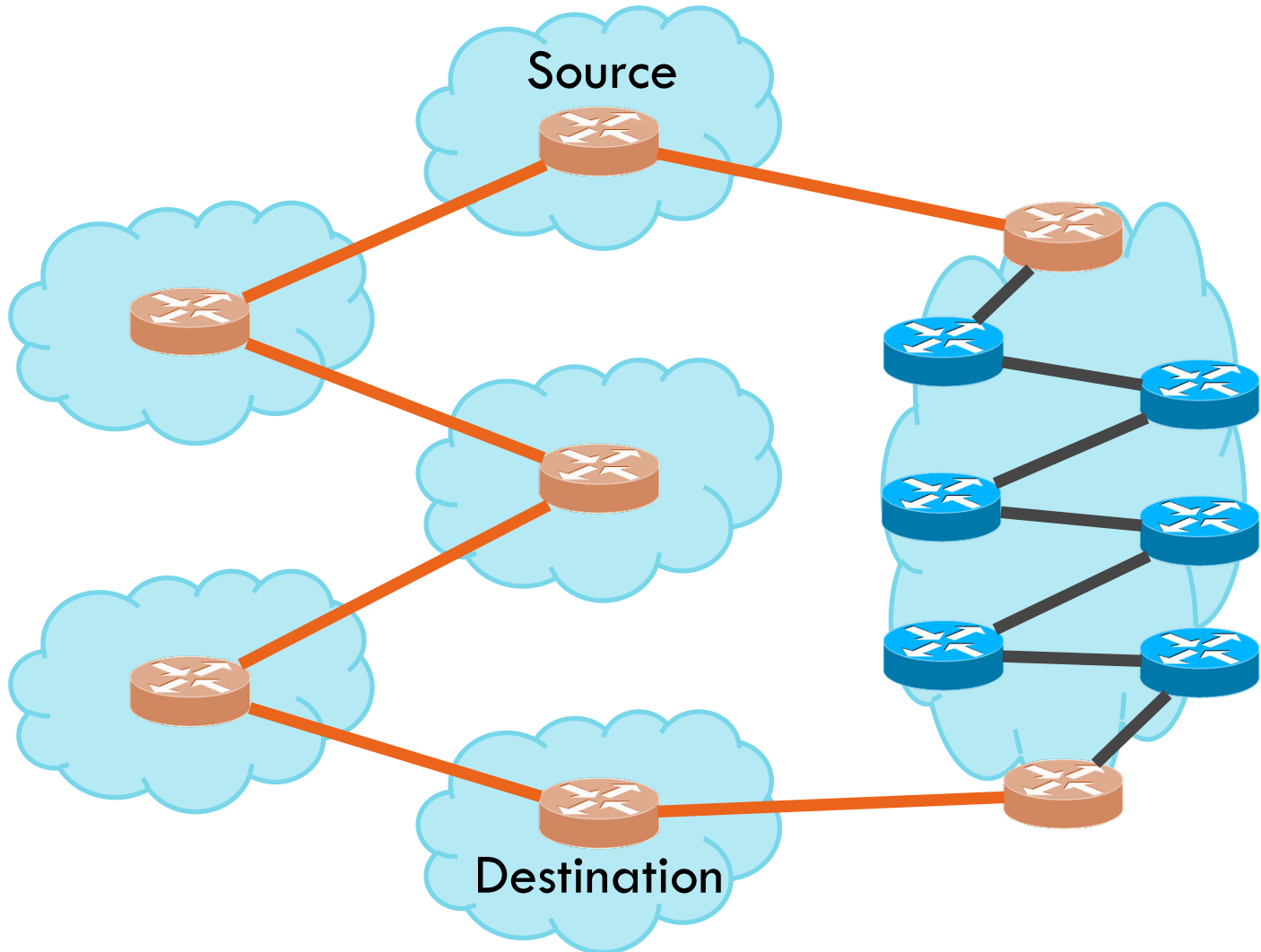
Traffic engineering

Lowest Router ID

**When all else fails,
break ties**

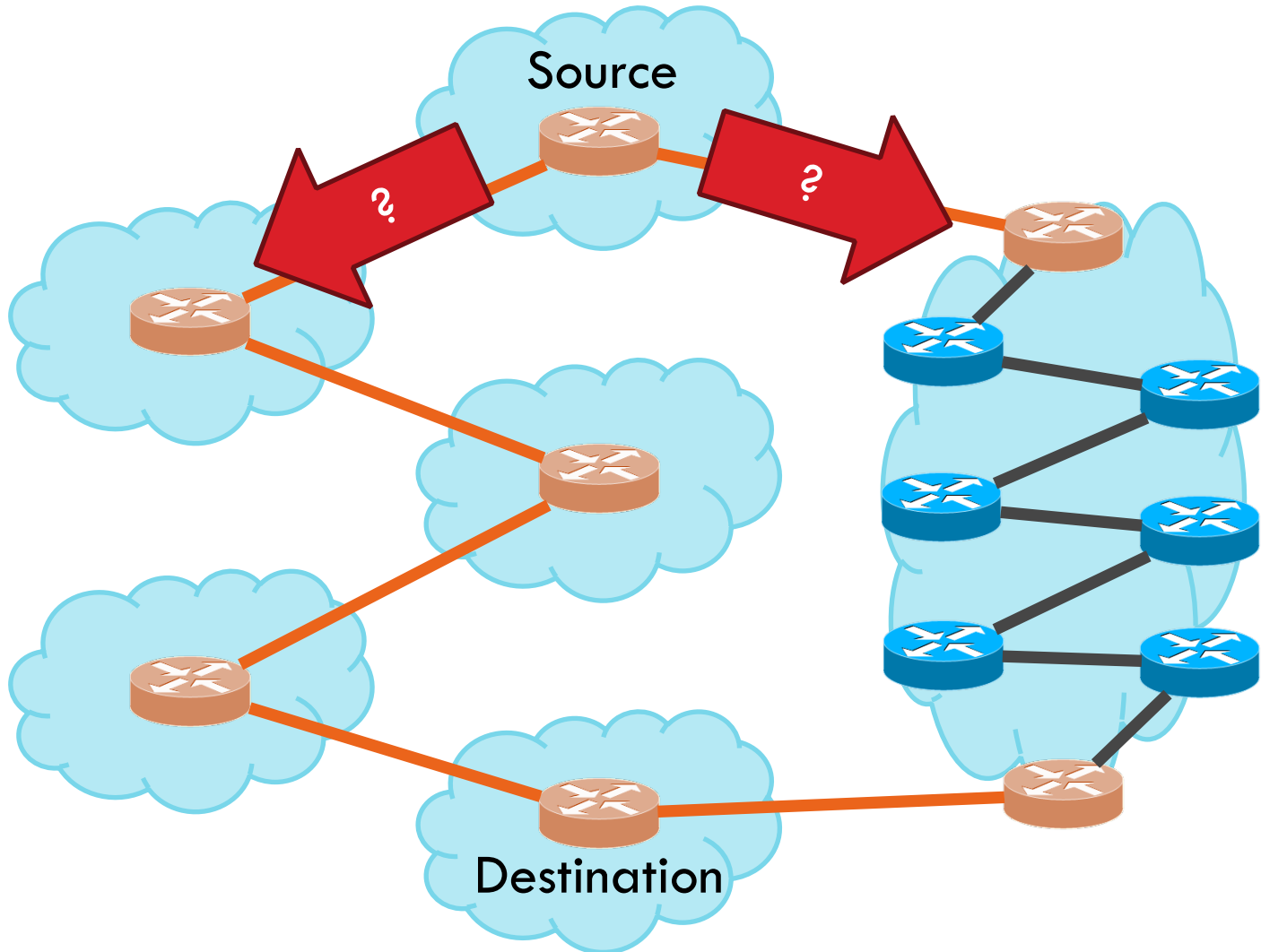
Shortest AS Path \neq Shortest Path

19



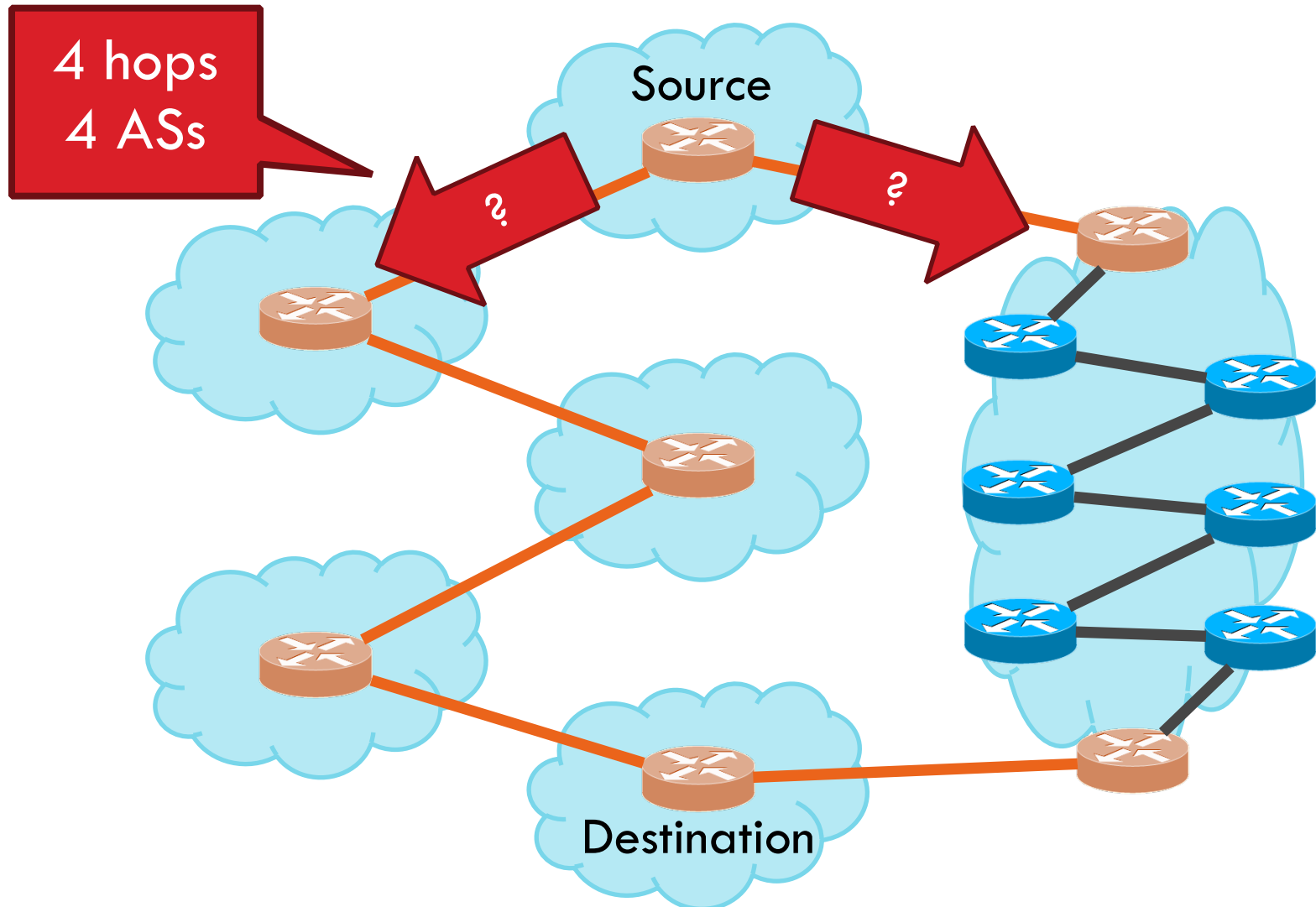
Shortest AS Path \neq Shortest Path

19



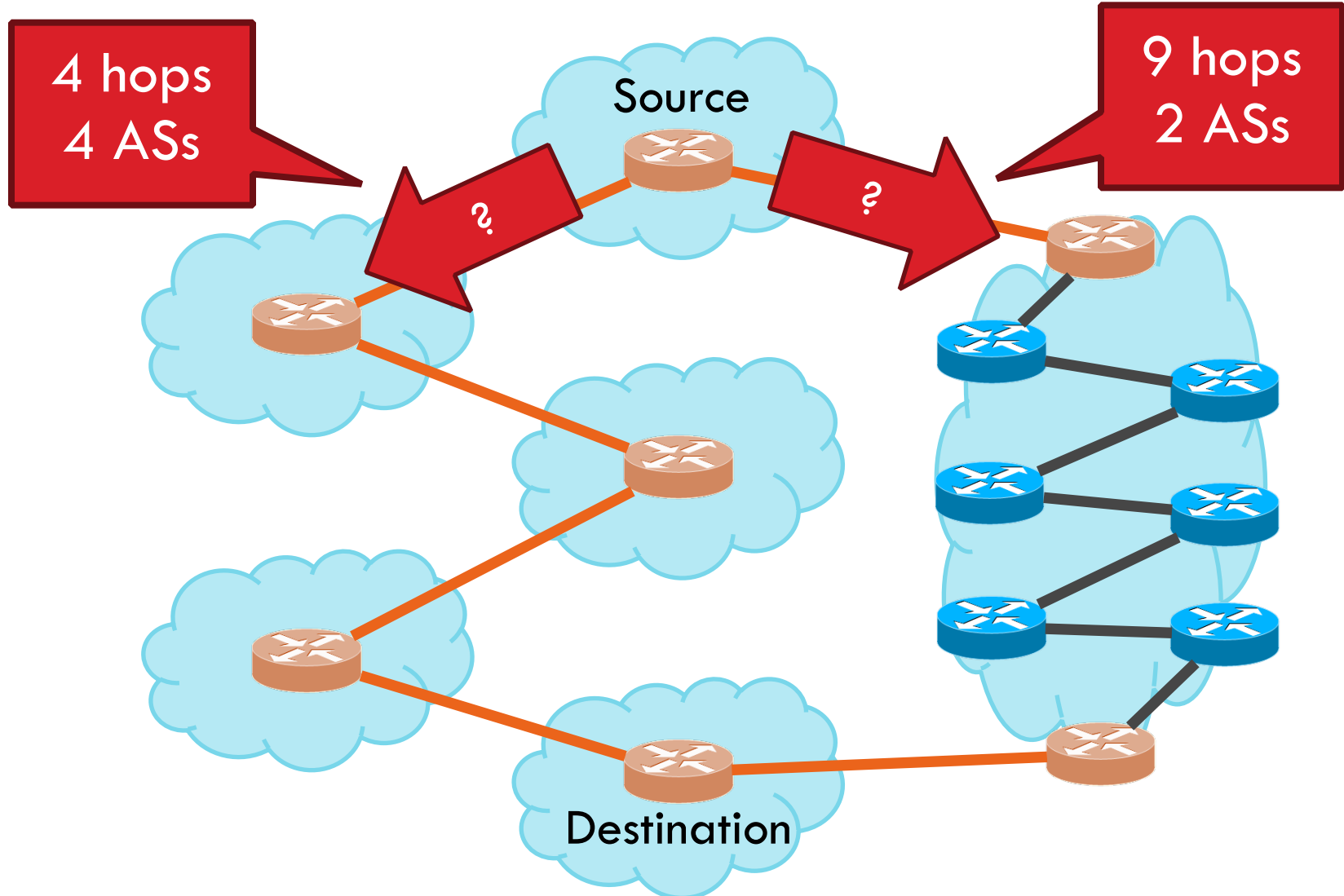
Shortest AS Path \neq Shortest Path

19



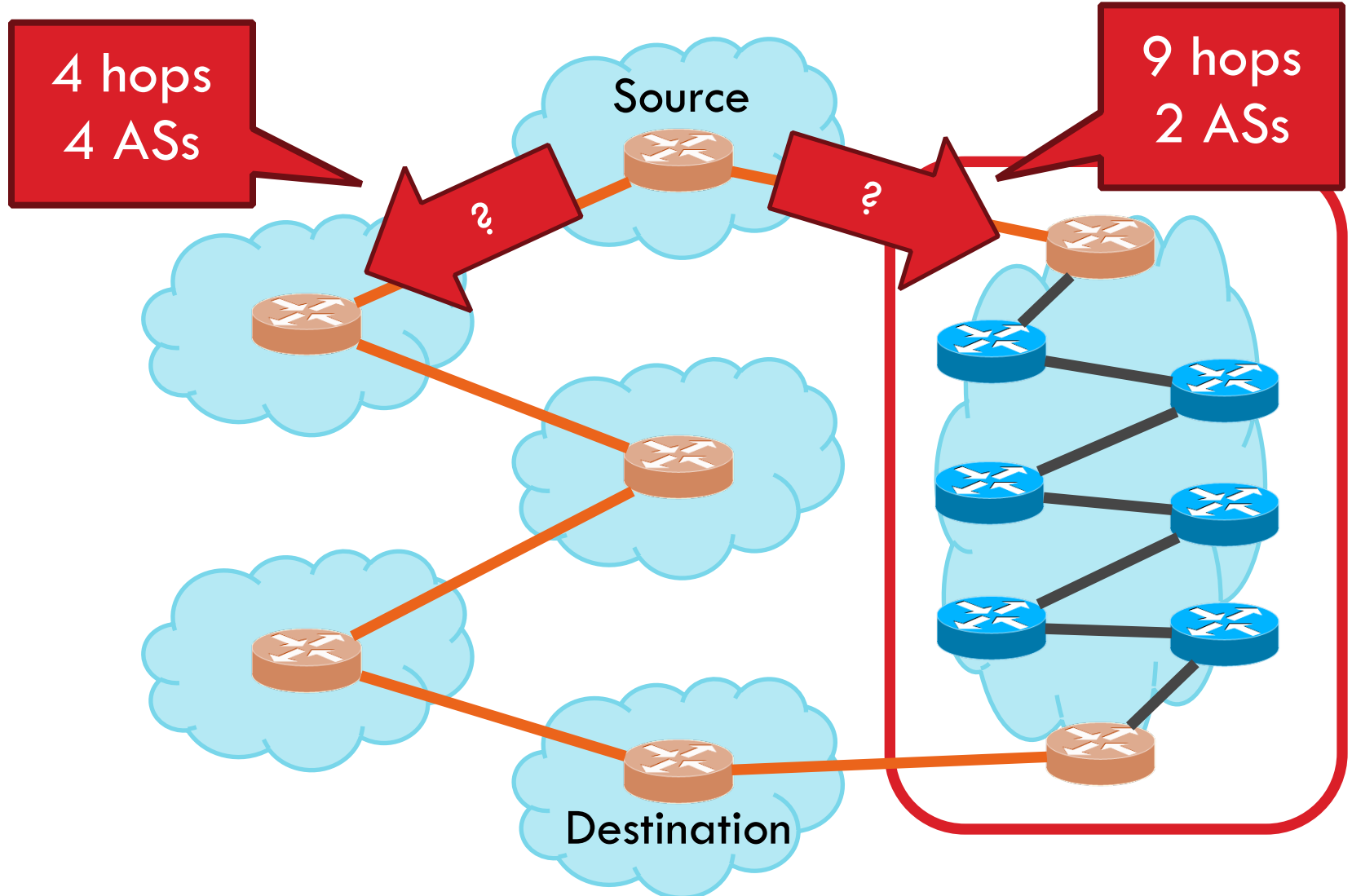
Shortest AS Path \neq Shortest Path

19



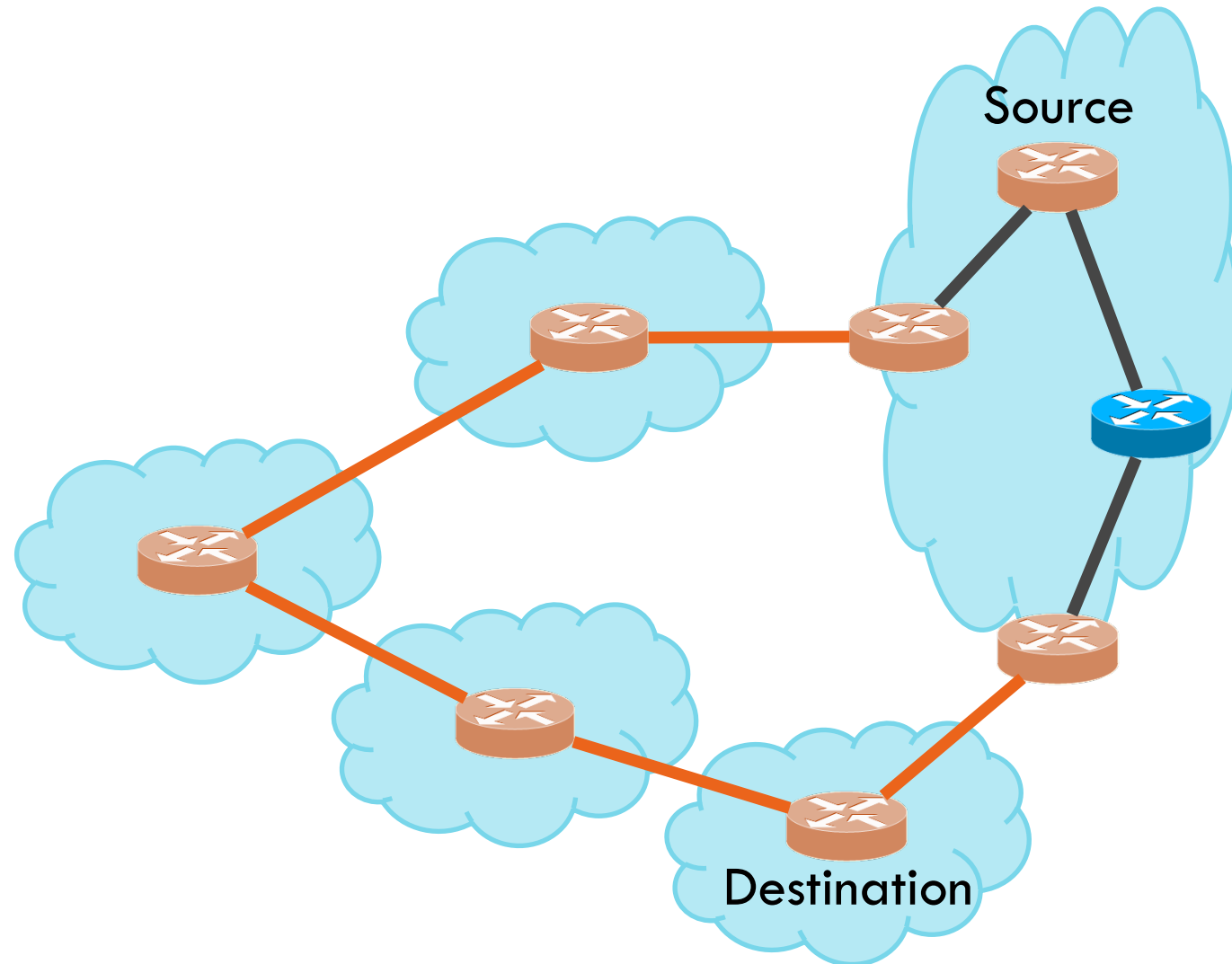
Shortest AS Path \neq Shortest Path

19



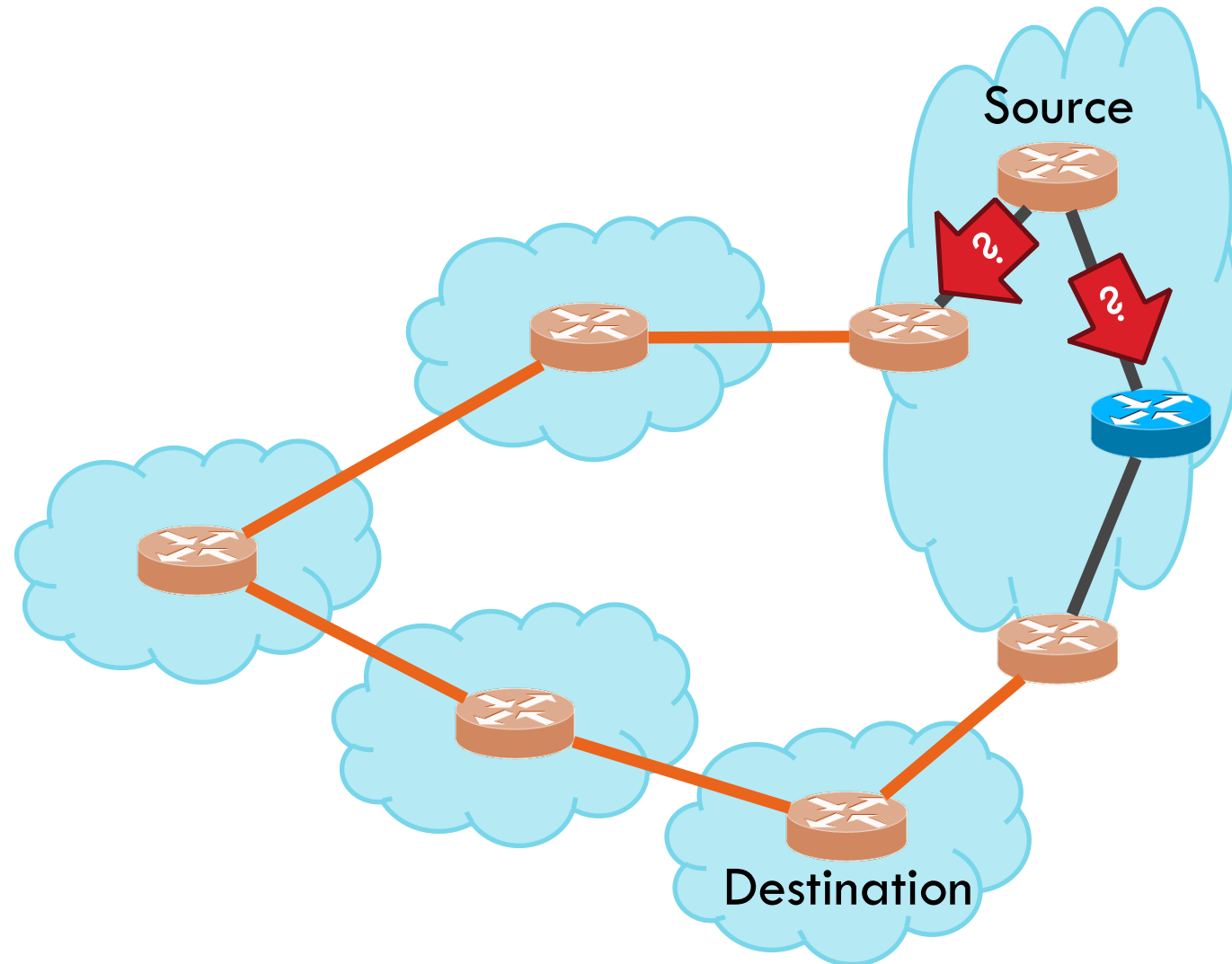
Hot Potato Routing

20



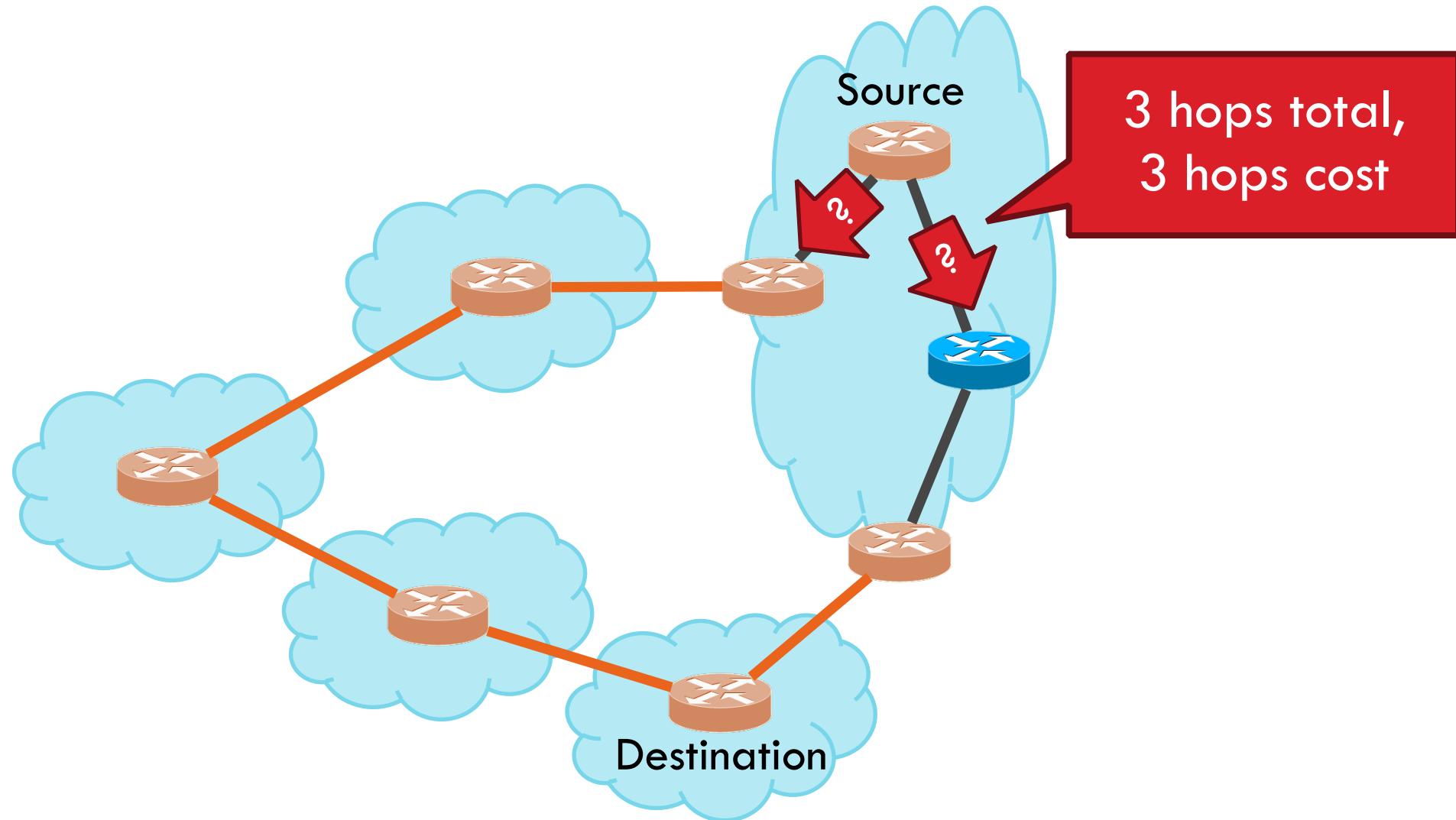
Hot Potato Routing

20



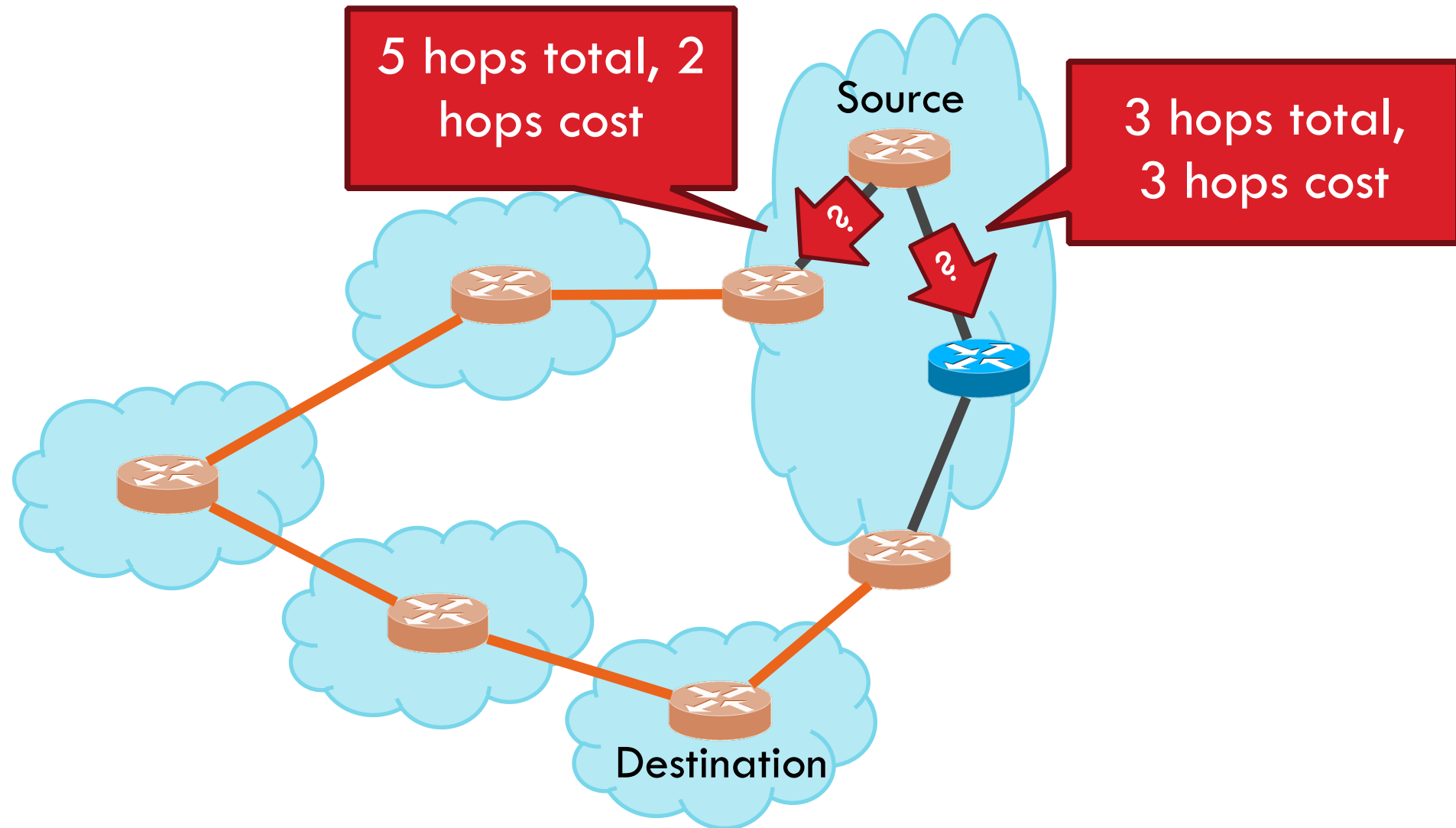
Hot Potato Routing

20



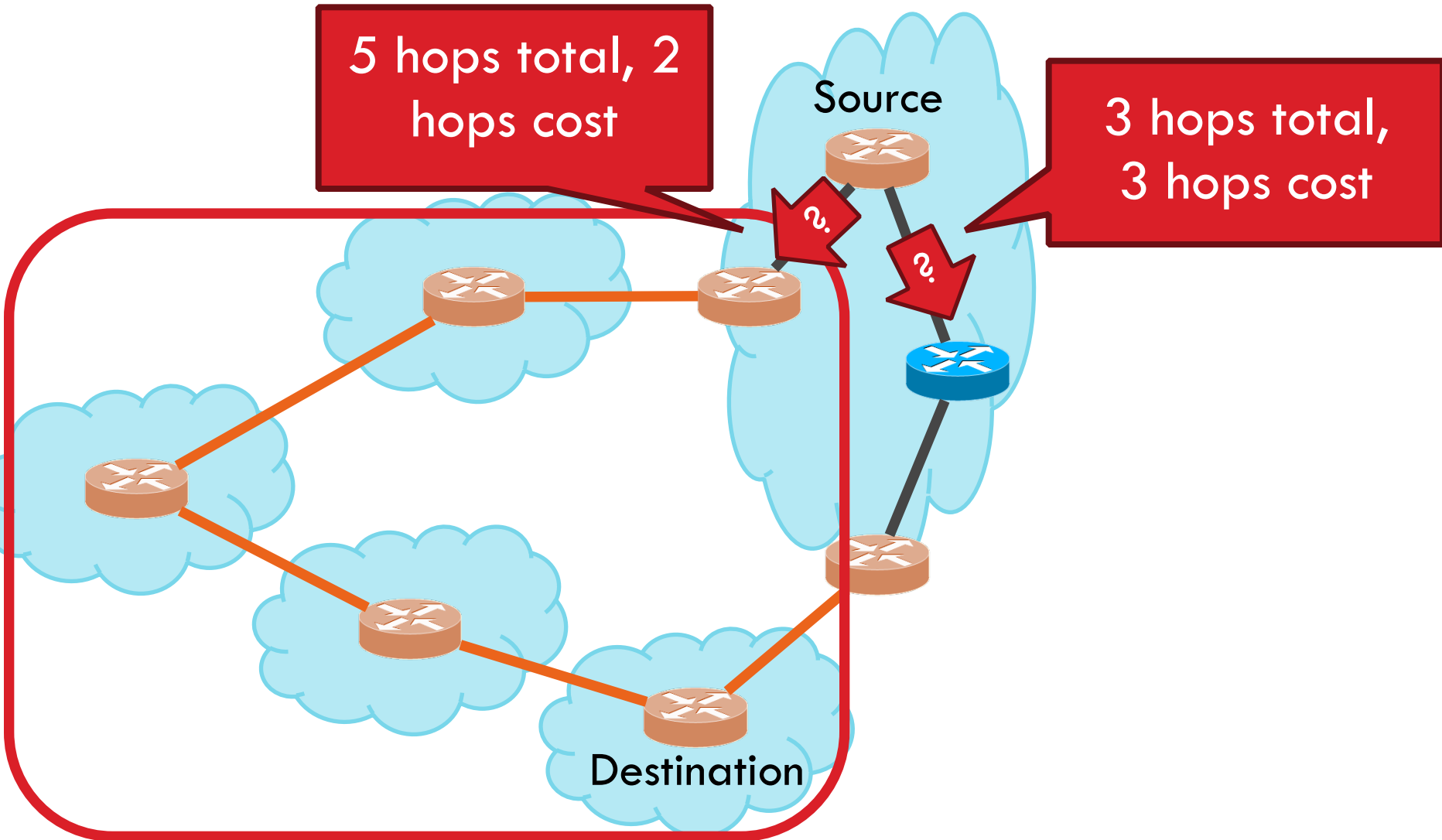
Hot Potato Routing

20



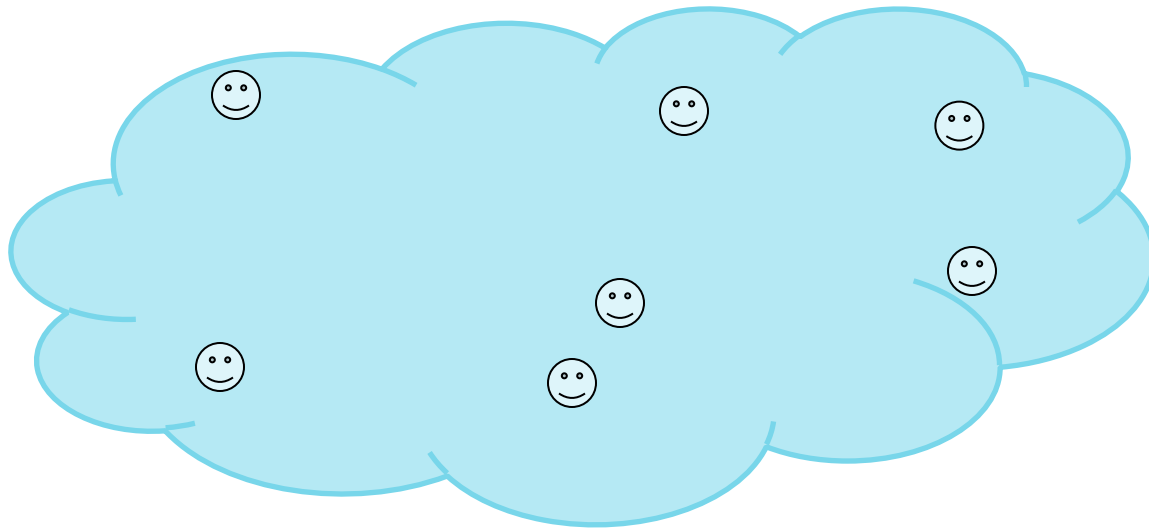
Hot Potato Routing

20



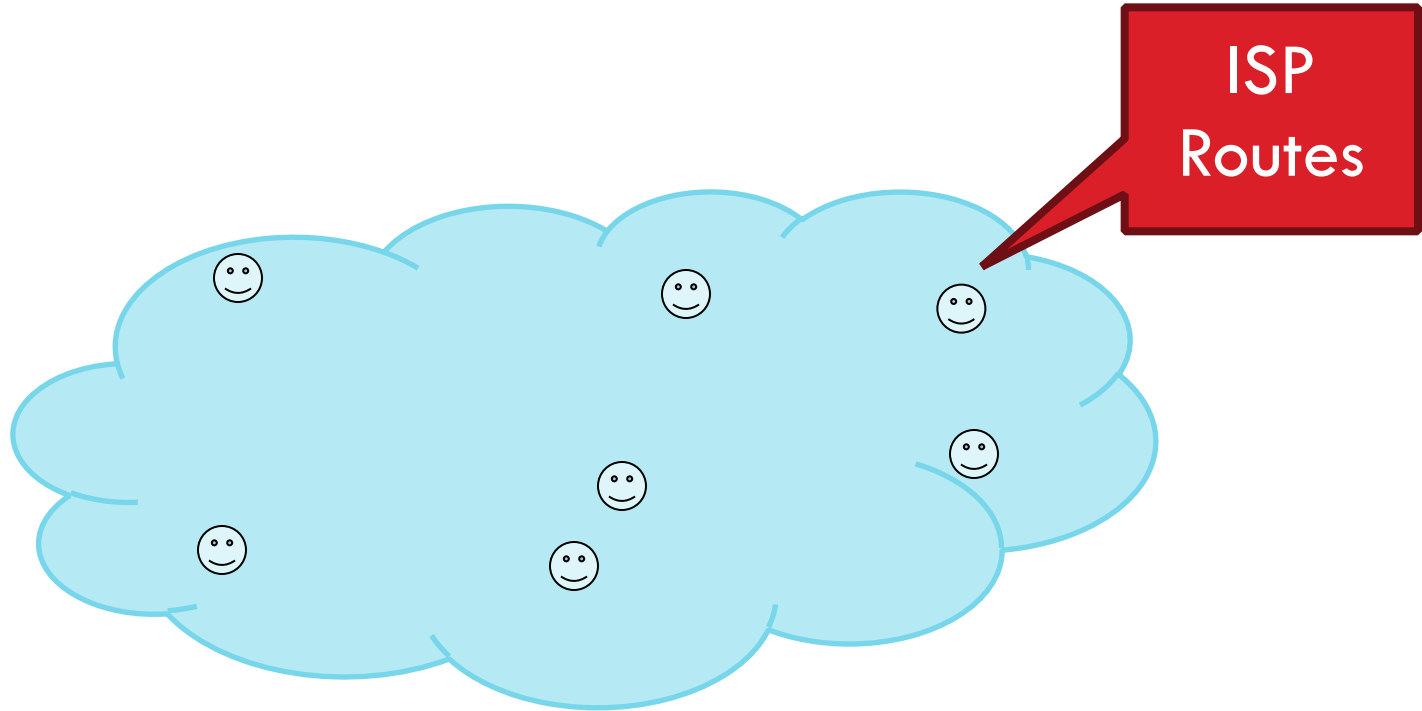
Importing Routes

21



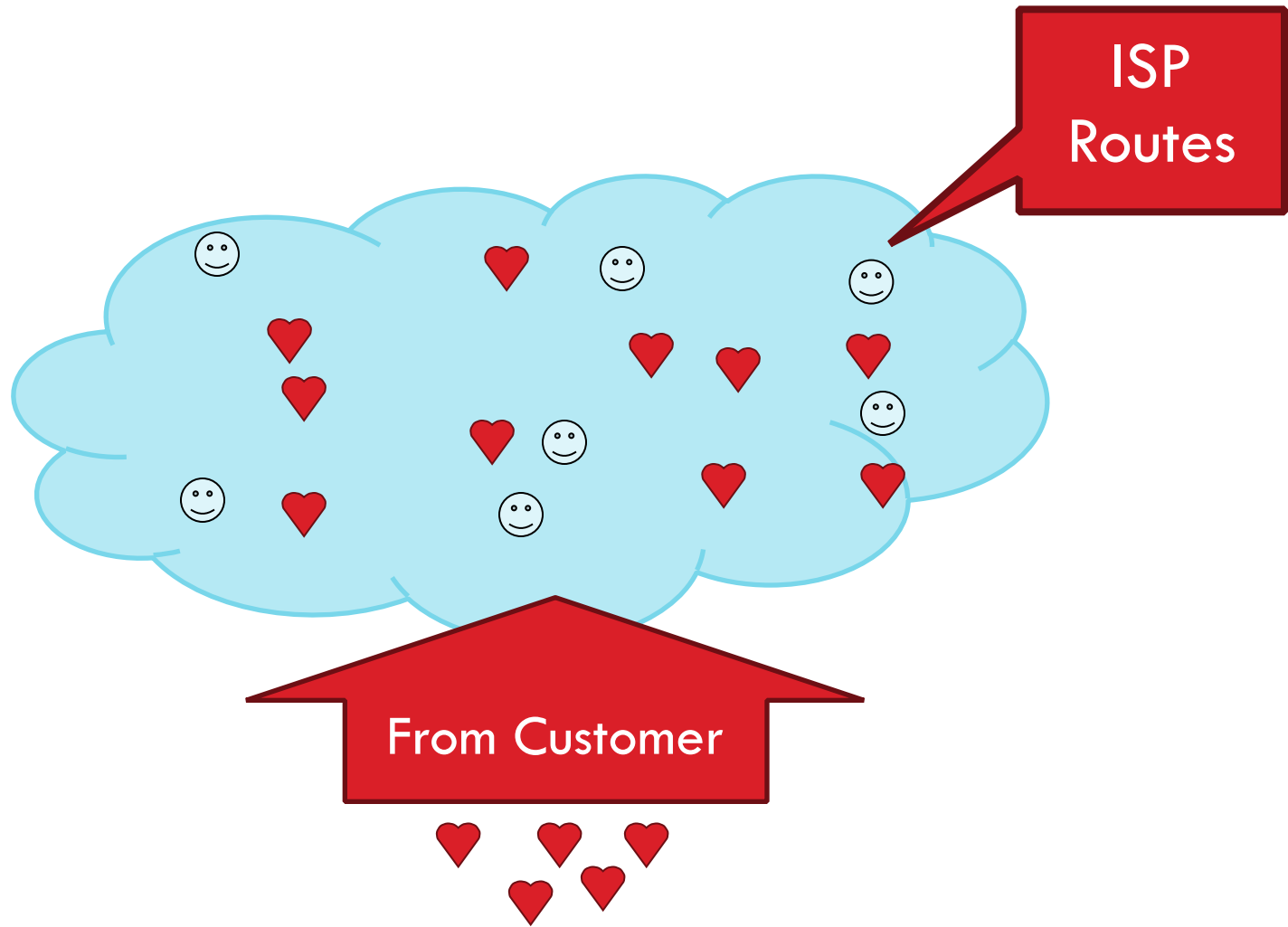
Importing Routes

21



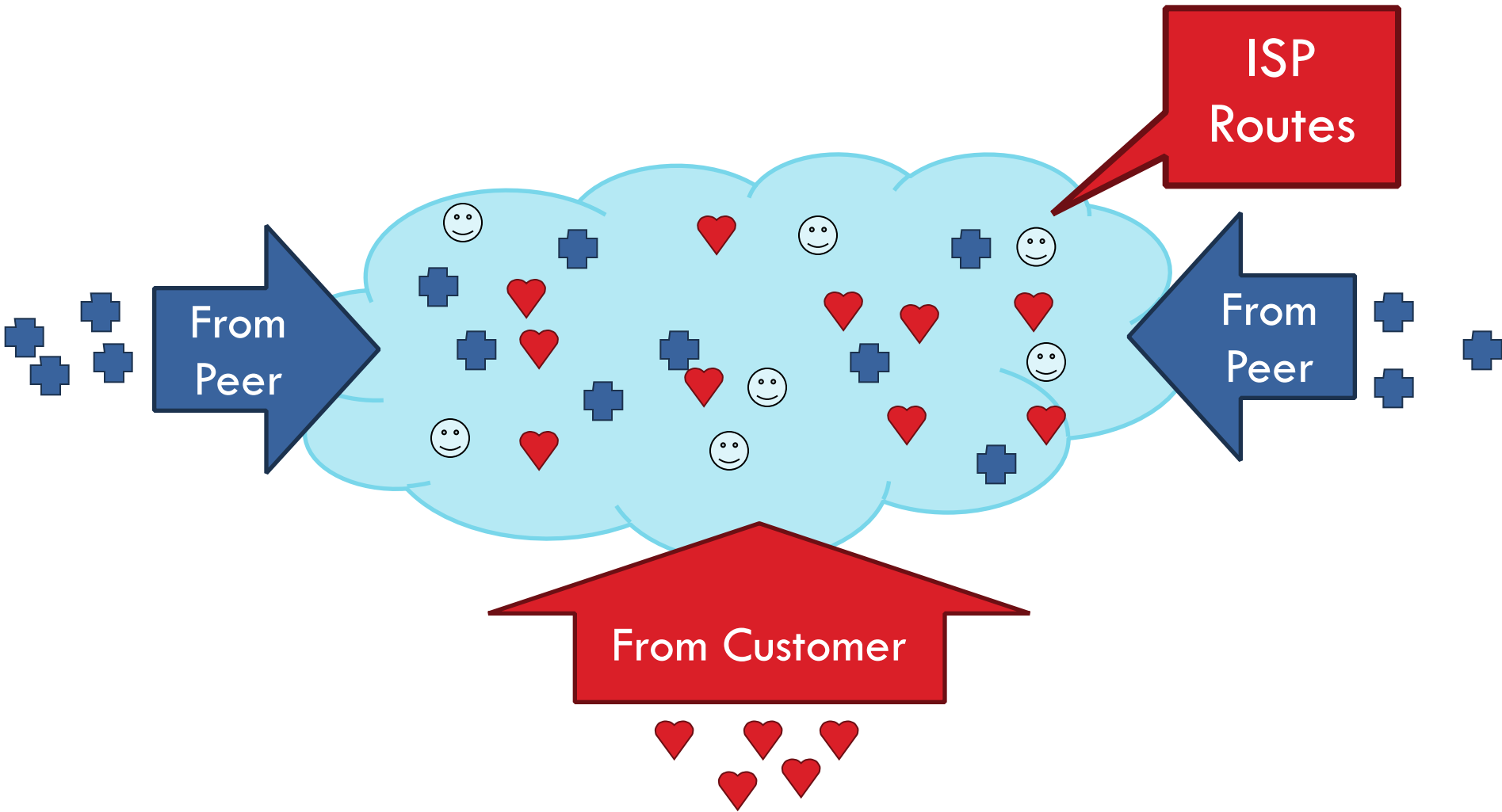
Importing Routes

21



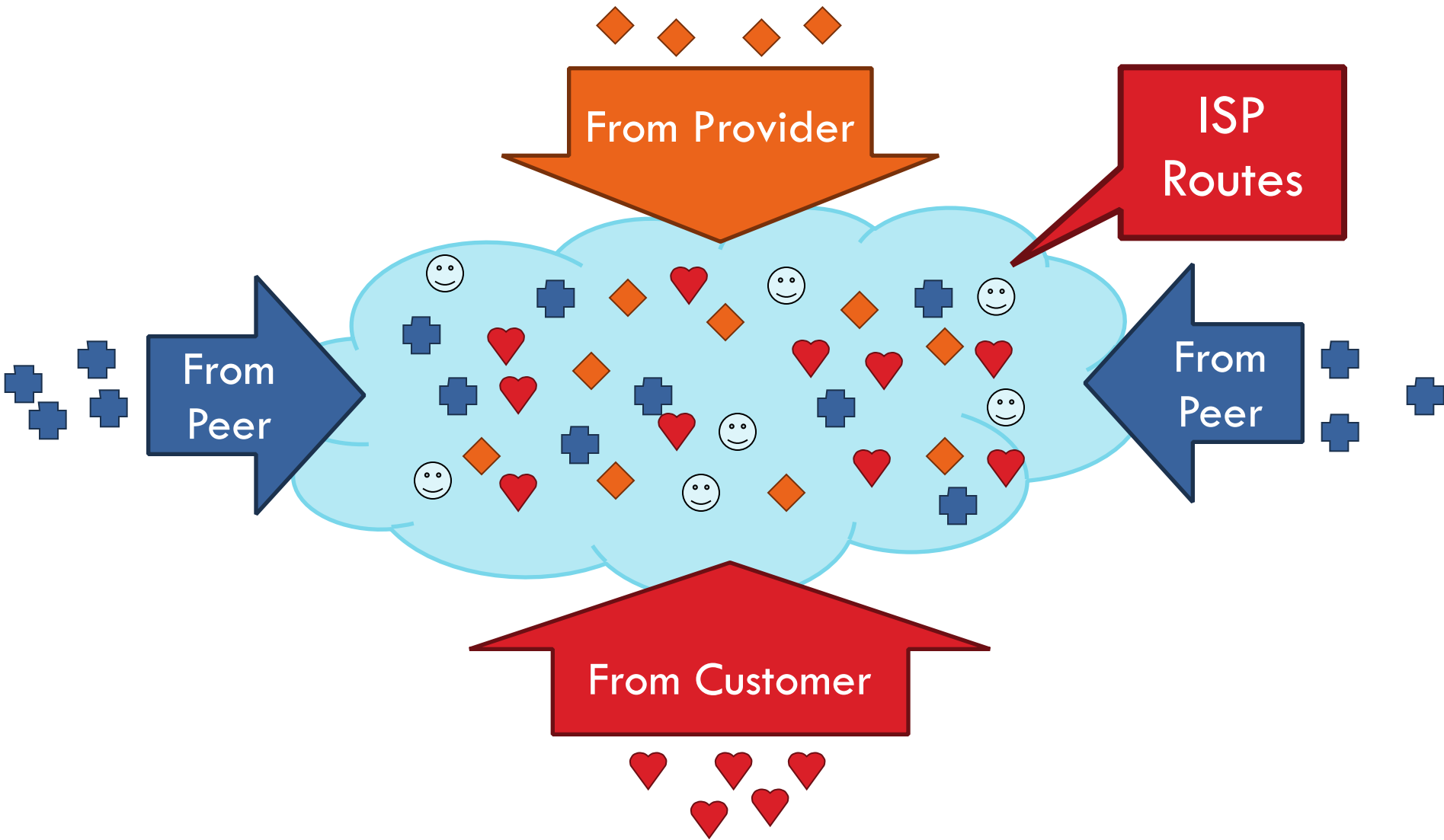
Importing Routes

21



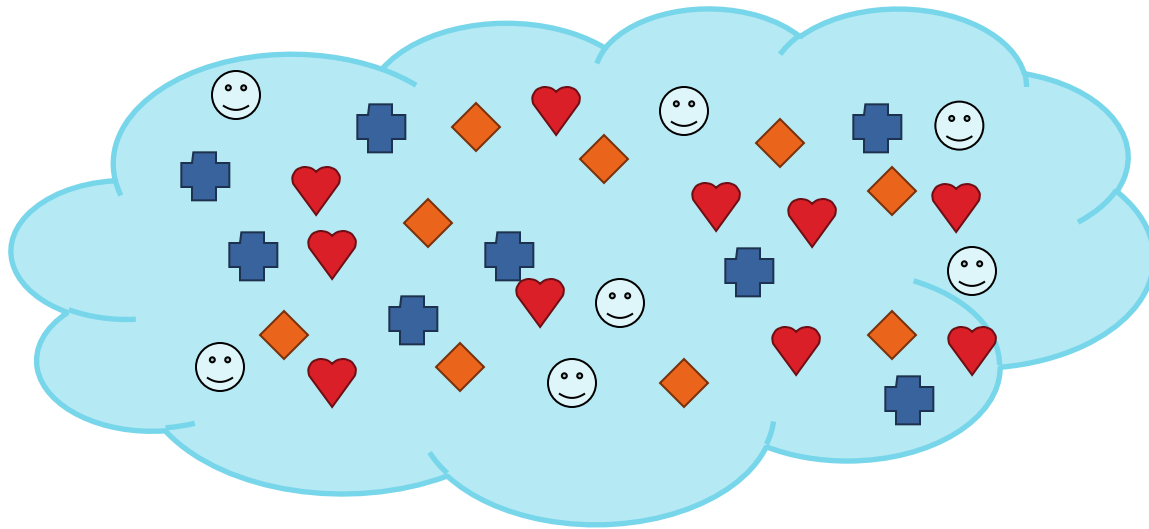
Importing Routes

21



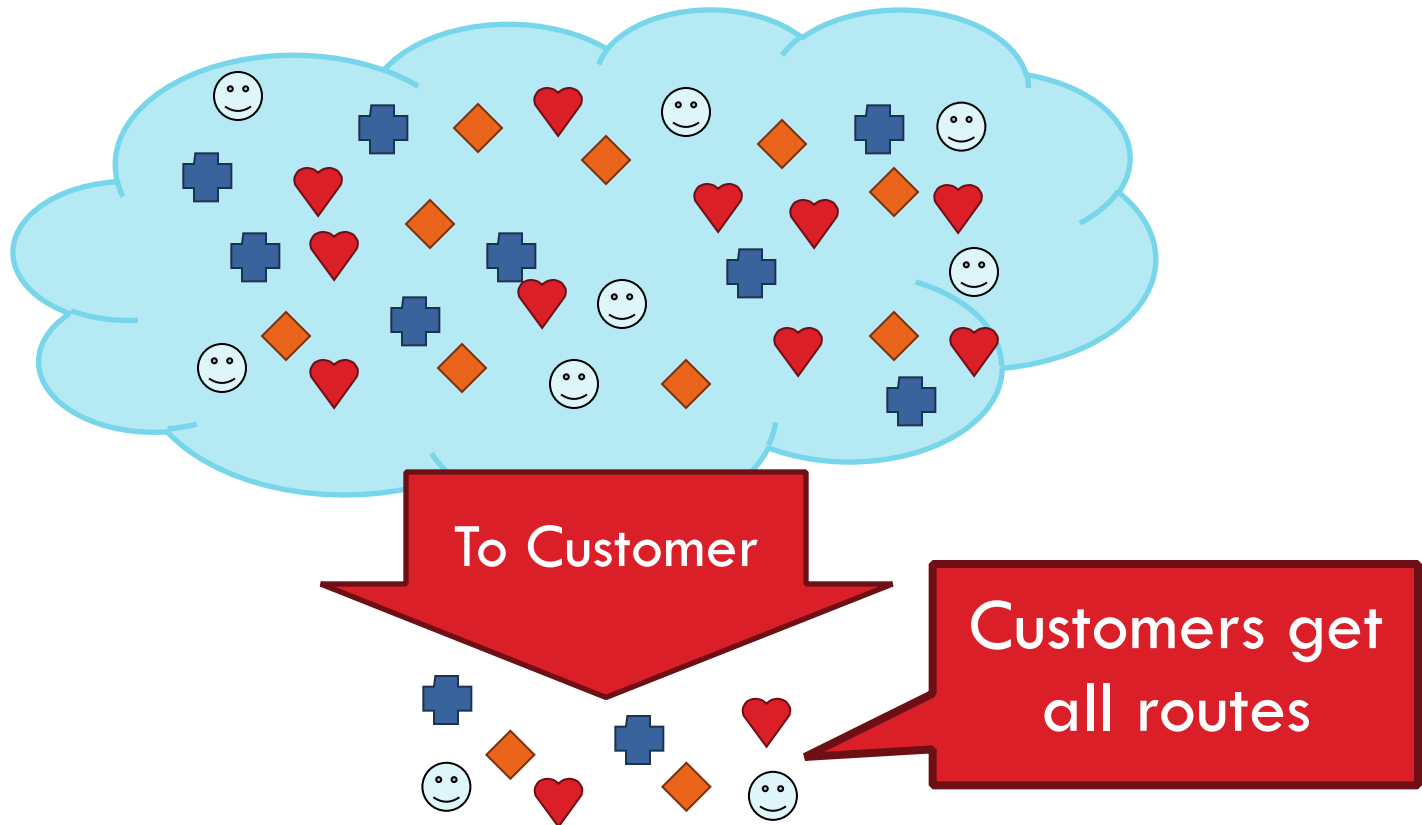
Exporting Routes

22



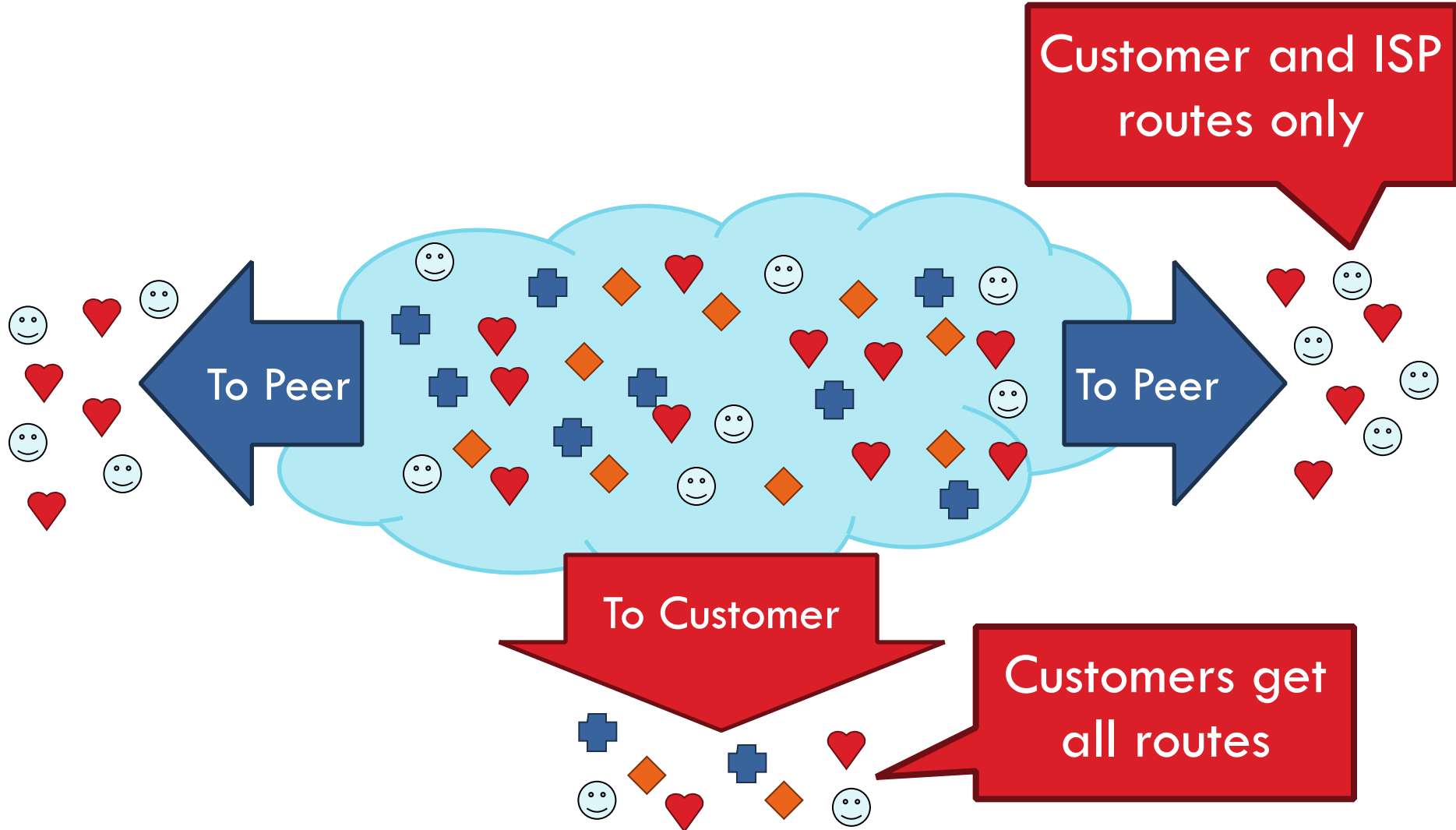
Exporting Routes

22



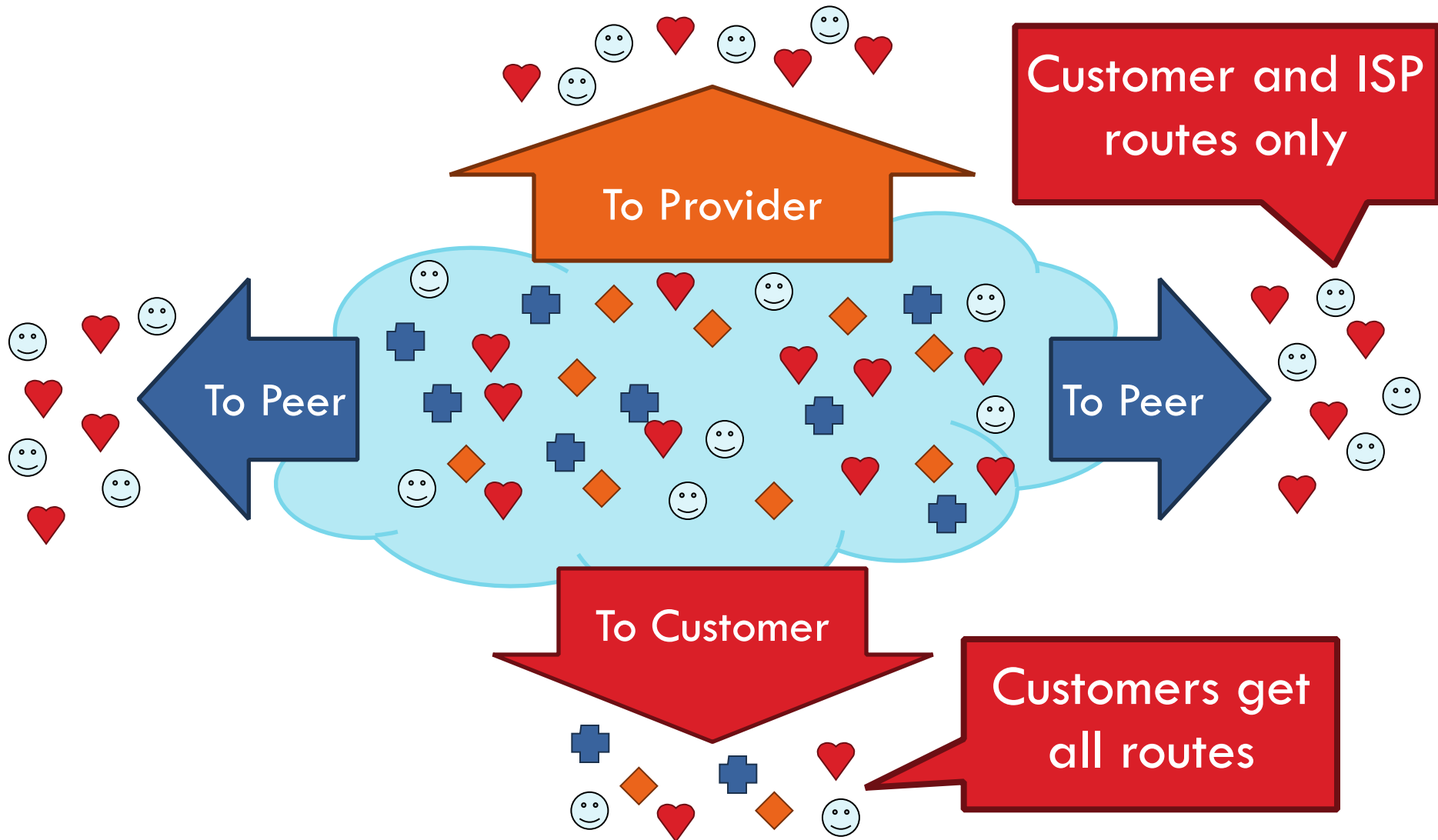
Exporting Routes

22



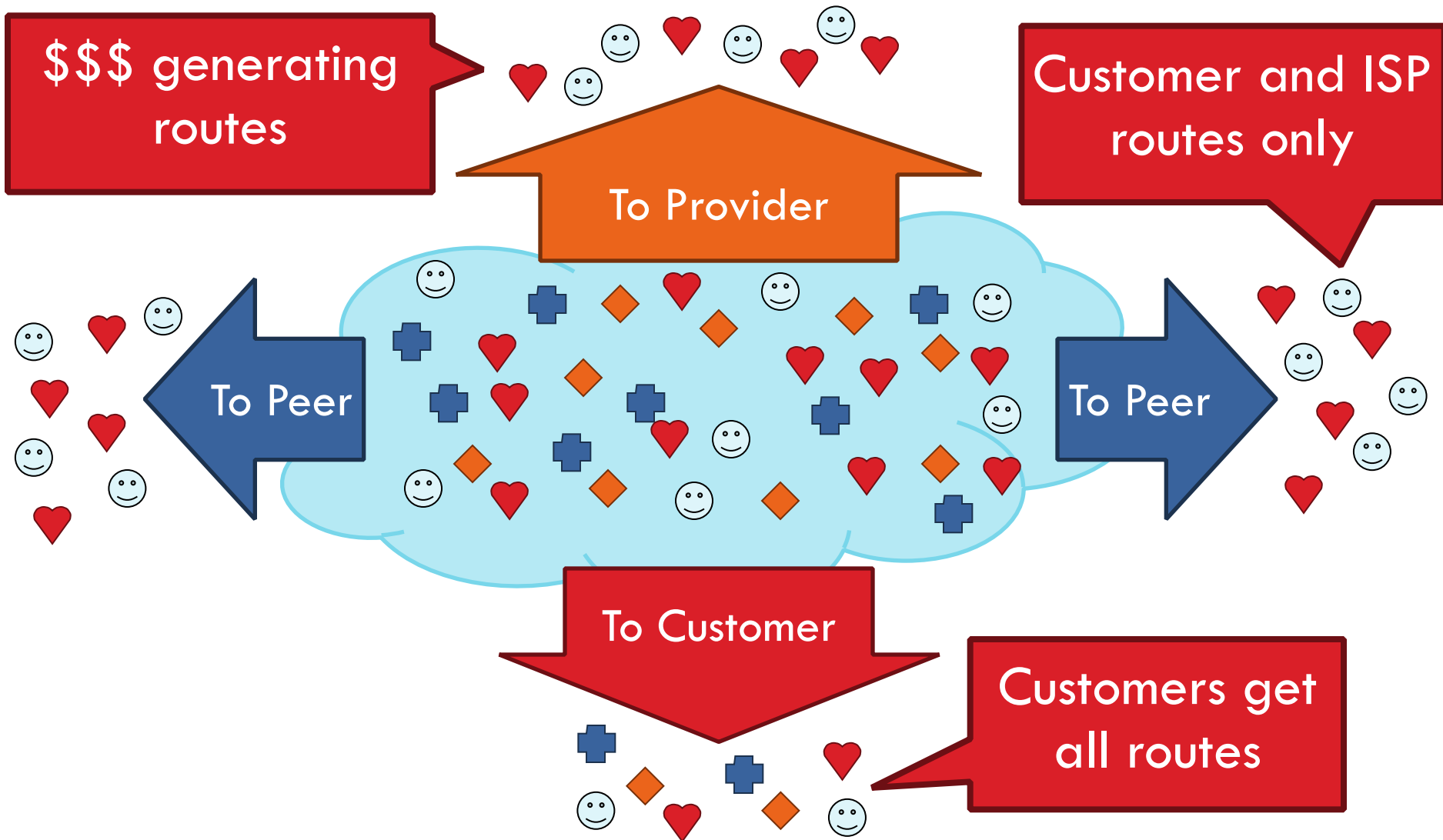
Exporting Routes

22



Exporting Routes

22



Modeling BGP

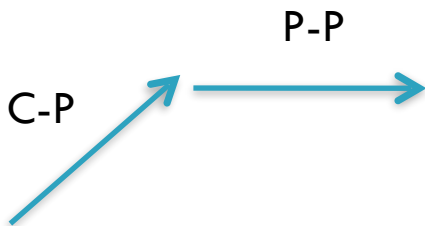
23

- AS relationships
 - Customer/provider
 - Peer
 - Sibling, IXP
- Gao-Rexford model
 - AS prefers to use customer path, then peer, then provider
 - Follow the money!
 - Valley-free routing
 - Hierarchical view of routing (incorrect but frequently used)

Modeling BGP

23

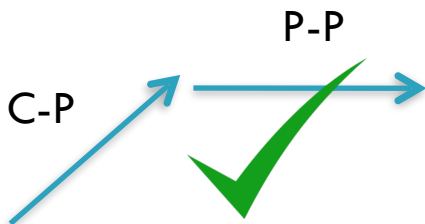
- AS relationships
 - Customer/provider
 - Peer
 - Sibling, IXP
- Gao-Rexford model
 - AS prefers to use customer path, then peer, then provider
 - Follow the money!
 - Valley-free routing
 - Hierarchical view of routing (incorrect but frequently used)



Modeling BGP

23

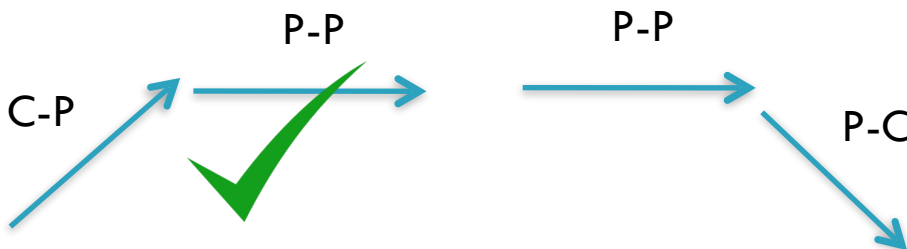
- AS relationships
 - Customer/provider
 - Peer
 - Sibling, IXP
- Gao-Rexford model
 - AS prefers to use customer path, then peer, then provider
 - Follow the money!
 - Valley-free routing
 - Hierarchical view of routing (incorrect but frequently used)



Modeling BGP

23

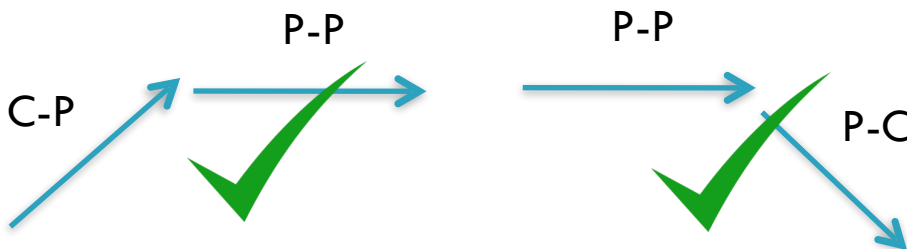
- AS relationships
 - Customer/provider
 - Peer
 - Sibling, IXP
- Gao-Rexford model
 - AS prefers to use customer path, then peer, then provider
 - Follow the money!
 - Valley-free routing
 - Hierarchical view of routing (incorrect but frequently used)



Modeling BGP

23

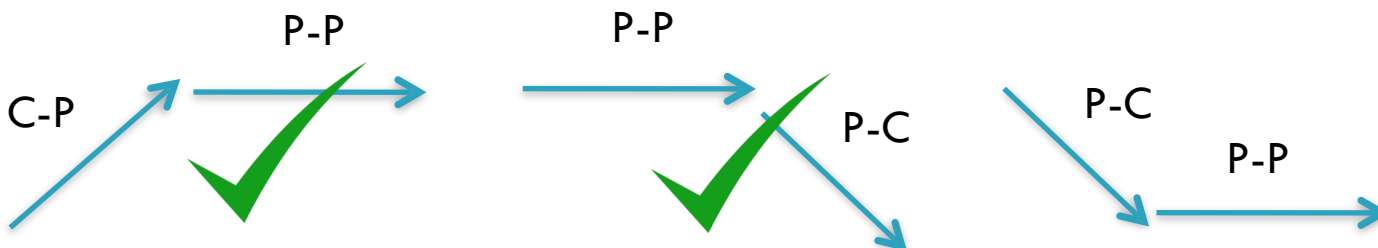
- AS relationships
 - Customer/provider
 - Peer
 - Sibling, IXP
- Gao-Rexford model
 - AS prefers to use customer path, then peer, then provider
 - Follow the money!
 - Valley-free routing
 - Hierarchical view of routing (incorrect but frequently used)



Modeling BGP

23

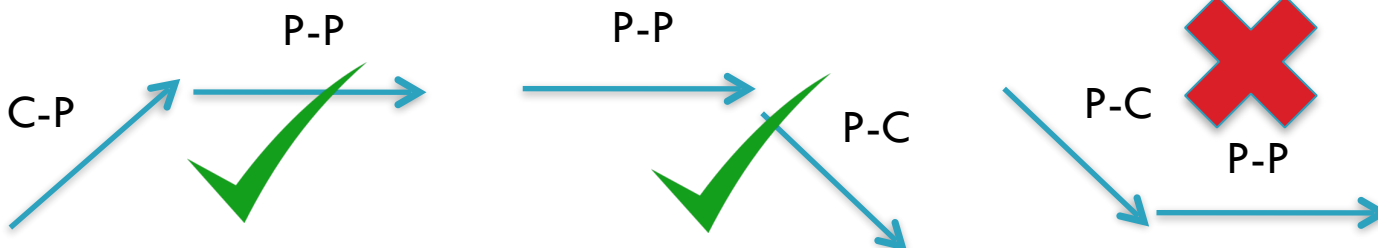
- AS relationships
 - Customer/provider
 - Peer
 - Sibling, IXP
- Gao-Rexford model
 - AS prefers to use customer path, then peer, then provider
 - Follow the money!
 - Valley-free routing
 - Hierarchical view of routing (incorrect but frequently used)



Modeling BGP

23

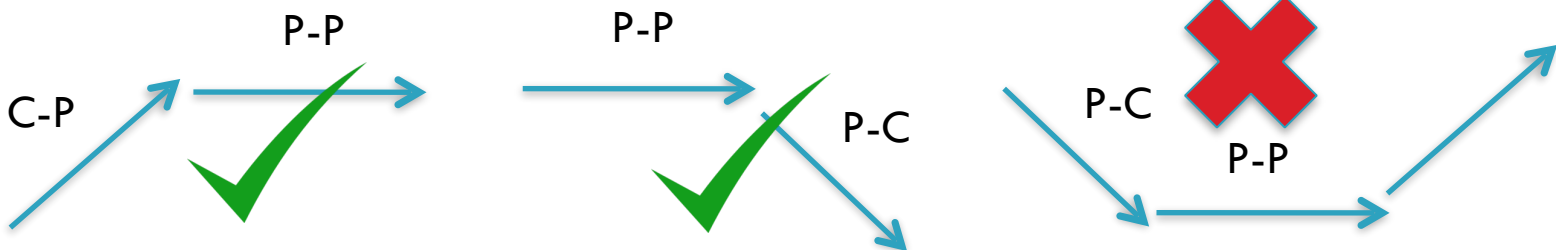
- AS relationships
 - Customer/provider
 - Peer
 - Sibling, IXP
- Gao-Rexford model
 - AS prefers to use customer path, then peer, then provider
 - Follow the money!
 - Valley-free routing
 - Hierarchical view of routing (incorrect but frequently used)



Modeling BGP

23

- AS relationships
 - Customer/provider
 - Peer
 - Sibling, IXP
- Gao-Rexford model
 - AS prefers to use customer path, then peer, then provider
 - Follow the money!
 - Valley-free routing
 - Hierarchical view of routing (incorrect but frequently used)



AS Relationships: It's Complicated

24

- GR Model is strictly hierarchical
 - ▣ Each AS pair has exactly one relationship
 - ▣ Each relationship is the same for all prefixes

AS Relationships: It's Complicated

24

- GR Model is strictly hierarchical
 - ▣ Each AS pair has exactly one relationship
 - ▣ Each relationship is the same for all prefixes
- In practice it's much more complicated
 - ▣ Rise of widespread peering
 - ▣ Regional, per-prefix peerings
 - ▣ Tier-1's being shoved out by "hypergiants"
 - ▣ IXPs dominating traffic volume

AS Relationships: It's Complicated

24

- GR Model is strictly hierarchical
 - ▣ Each AS pair has exactly one relationship
 - ▣ Each relationship is the same for all prefixes
- In practice it's much more complicated
 - ▣ Rise of widespread peering
 - ▣ Regional, per-prefix peerings
 - ▣ Tier-1's being shoved out by "hypergiants"
 - ▣ IXPs dominating traffic volume
- Modeling is very hard, very prone to error
 - ▣ Huge potential impact for understanding Internet behavior

Other BGP Attributes

25

- AS_SET
 - ▣ Instead of a single AS appearing at a slot, it's a set of Ases
 - ▣ Why?

Other BGP Attributes

25

- AS_SET
 - ▣ Instead of a single AS appearing at a slot, it's a set of Ases
 - ▣ Why?
- Communities
 - ▣ Arbitrary number that is used by neighbors for routing decisions
 - Export this route only in Europe
 - Do not export to your peers
 - ▣ Usually stripped after first interdomain hop
 - ▣ Why?

Other BGP Attributes

25

- AS_SET
 - ▣ Instead of a single AS appearing at a slot, it's a set of Ases
 - ▣ Why?
- Communities
 - ▣ Arbitrary number that is used by neighbors for routing decisions
 - Export this route only in Europe
 - Do not export to your peers
 - ▣ Usually stripped after first interdomain hop
 - ▣ Why?
- Prepending
 - ▣ Lengthening the route by adding multiple instances of ASN
 - ▣ Why?

- ❑ BGP Basics
- ❑ Stable Paths Problem
- ❑ BGP in the Real World

What Problem is BGP Solving?

27

Underlying Problem	Distributed Solution
<i>Shortest Paths</i>	<i>RIP, OSPF, IS-IS, etc.</i>
<i>???</i>	<i>BGP</i>

What Problem is BGP Solving?

27

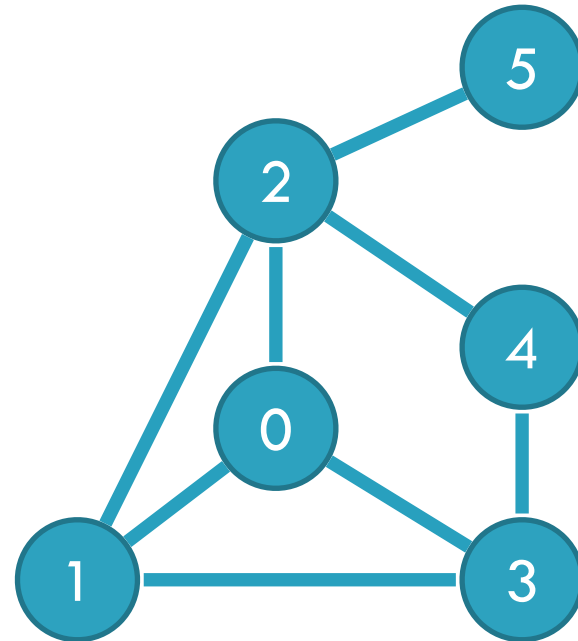
Underlying Problem	Distributed Solution
<i>Shortest Paths</i>	<i>RIP, OSPF, IS-IS, etc.</i>
<i>???</i>	<i>BGP</i>

- Knowing ??? can:
 - Aid in the analysis of BGP policy
 - Aid in the design of BGP extensions
 - Help explain BGP routing anomalies
 - Give us a deeper understanding of the protocol

The Stable Paths Problem

28

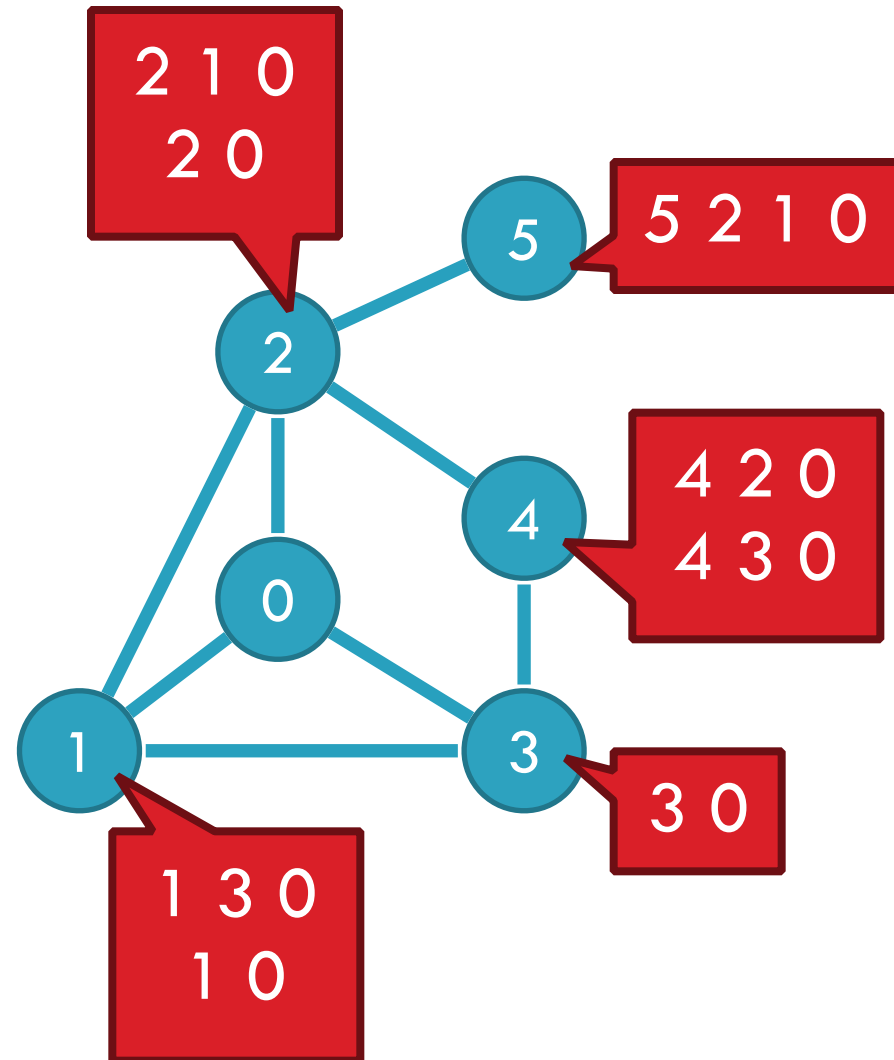
- An instance of the SPP:
 - Graph of nodes and edges
 - Node 0, called the origin



The Stable Paths Problem

28

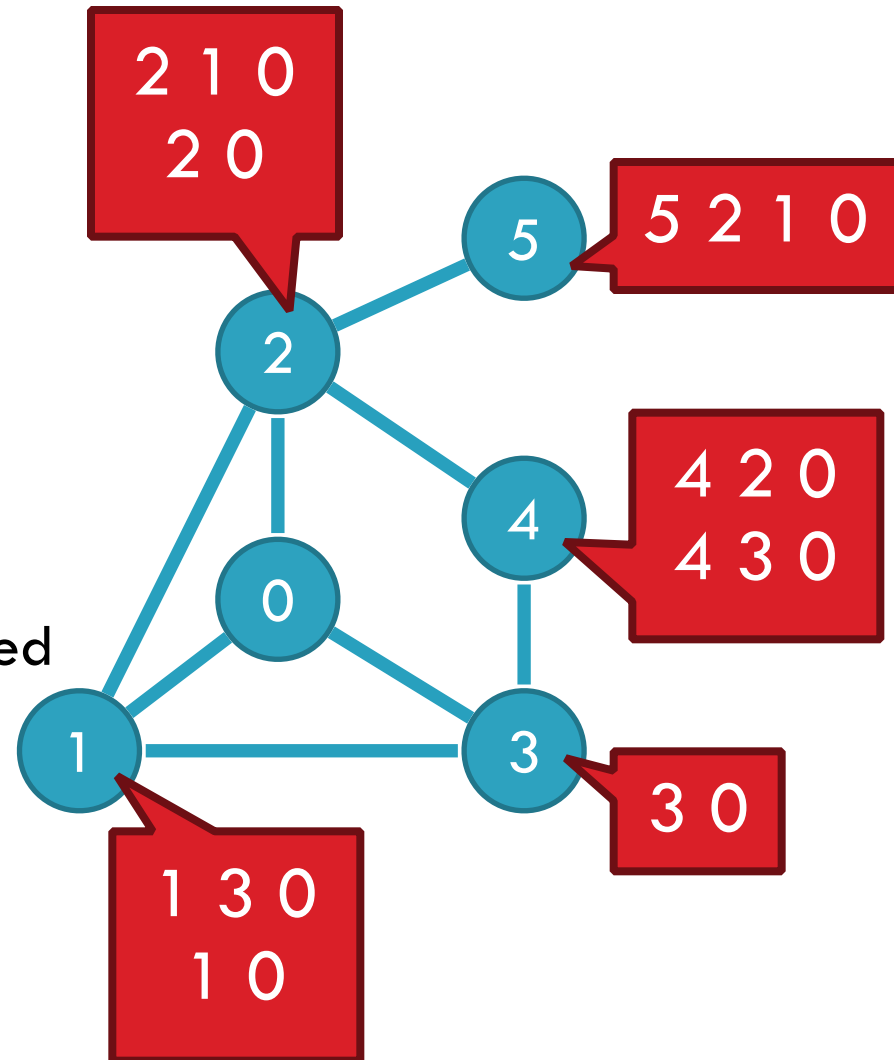
- An instance of the SPP:
 - Graph of nodes and edges
 - Node 0, called the origin
 - A set of permitted paths from each node to the origin
 - Each set contains the null path



The Stable Paths Problem

28

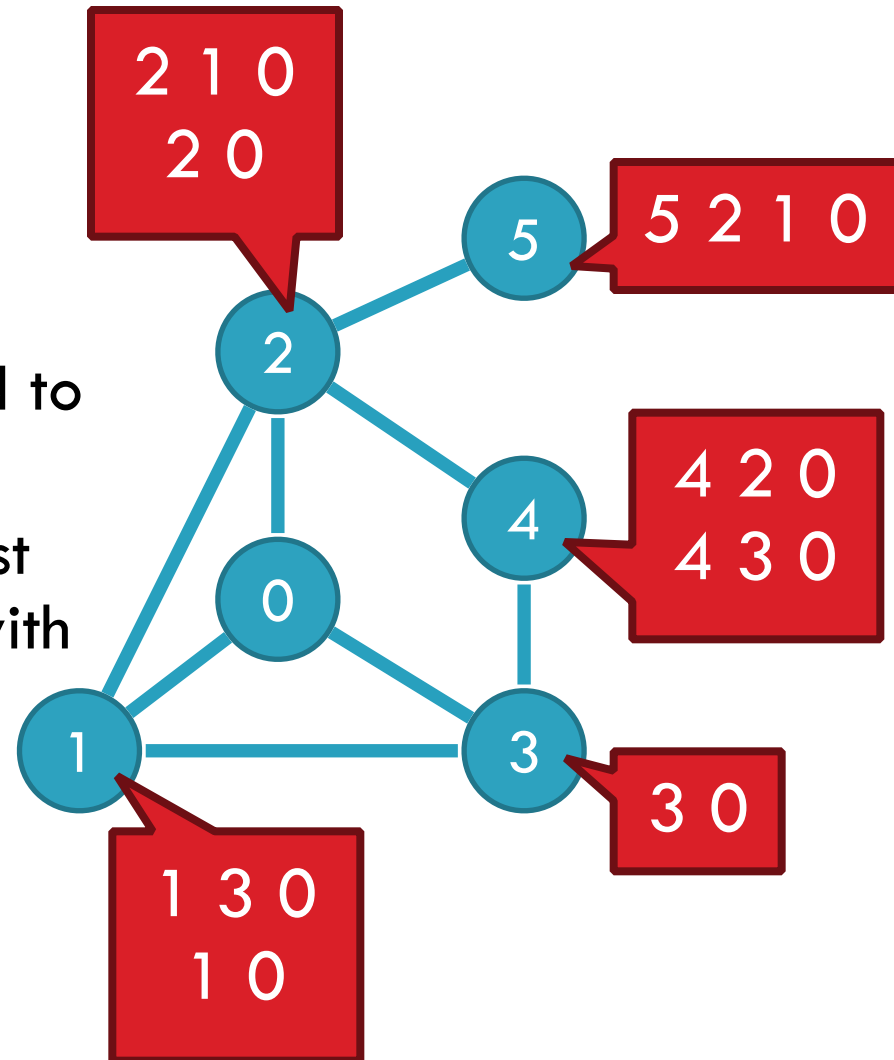
- An instance of the SPP:
 - Graph of nodes and edges
 - Node 0, called the origin
 - A set of permitted paths from each node to the origin
 - Each set contains the null path
 - Each set of paths is ranked
 - Null path is always least preferred



A Solution to the SPP

29

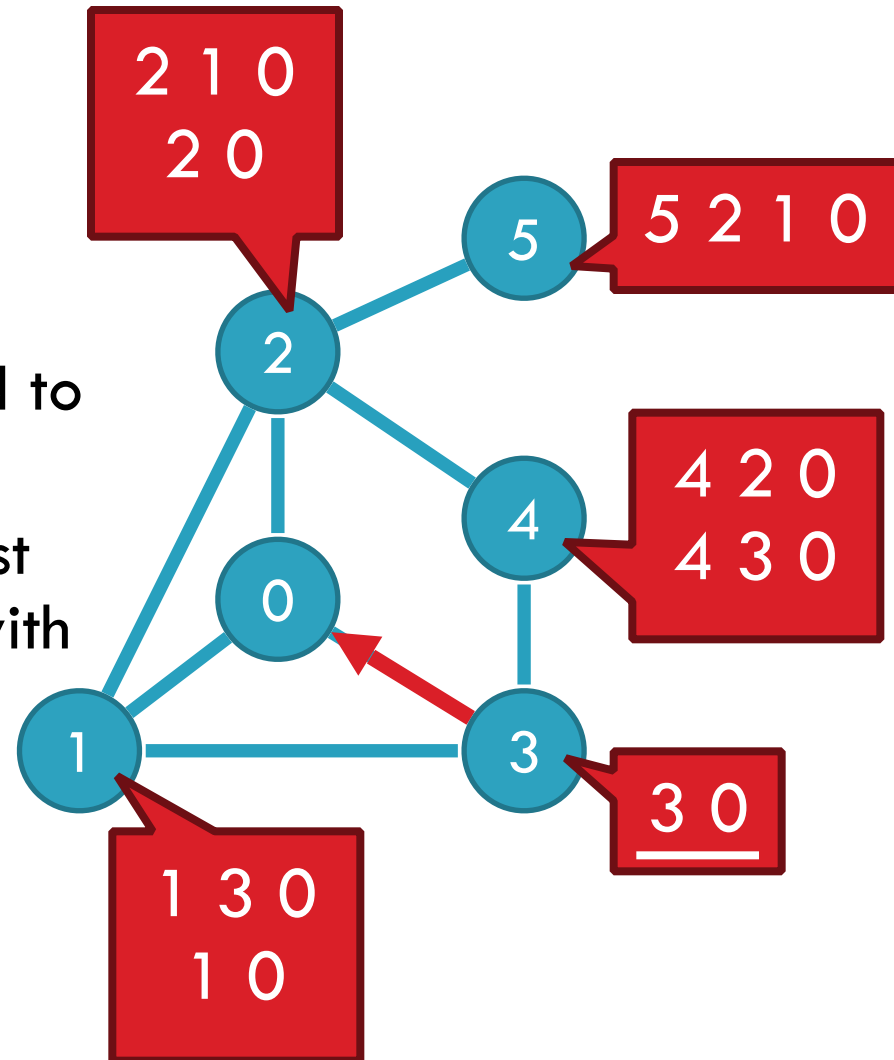
- A solution is an assignment of permitted paths to each node such that:
 - ▣ Node u 's path is either null or uwP , where path w is assigned to node w and edge $u \rightarrow w$ exists
 - ▣ Each node is assigned the highest ranked path that is consistent with their neighbors



A Solution to the SPP

29

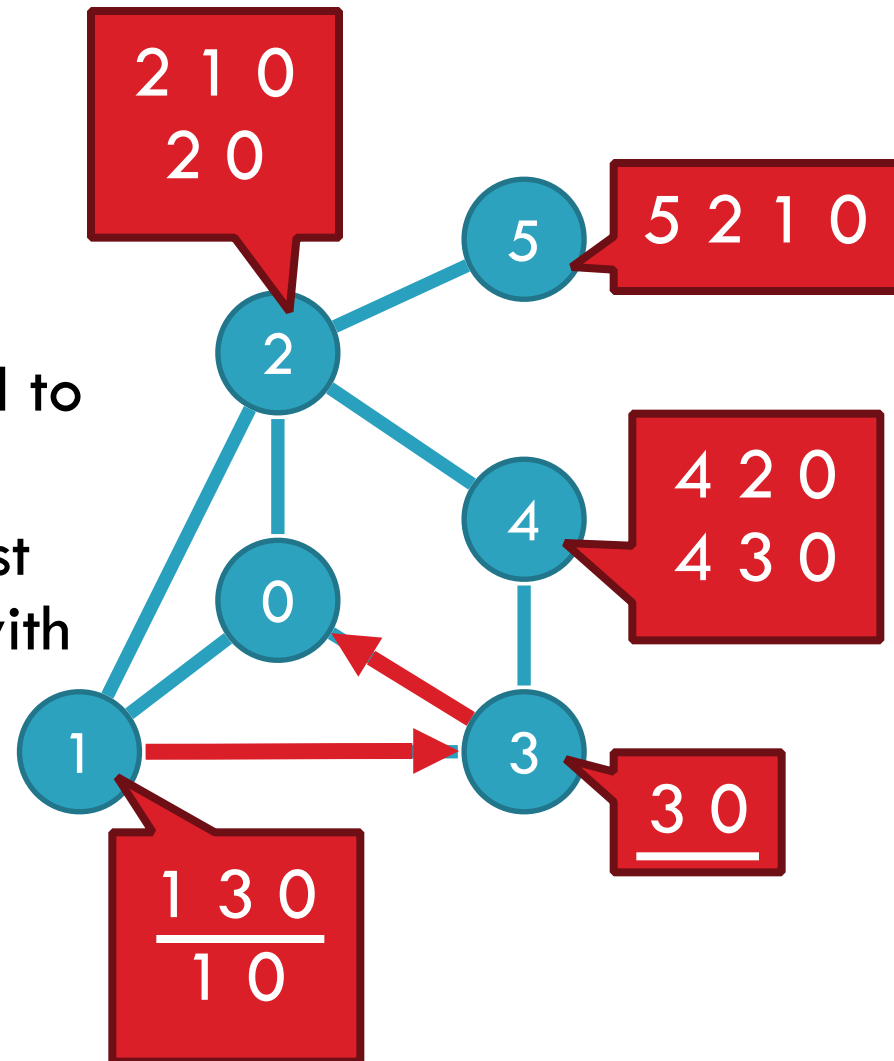
- A solution is an assignment of permitted paths to each node such that:
 - ▣ Node u 's path is either null or uwP , where path w is assigned to node w and edge $u \rightarrow w$ exists
 - ▣ Each node is assigned the highest ranked path that is consistent with their neighbors



A Solution to the SPP

29

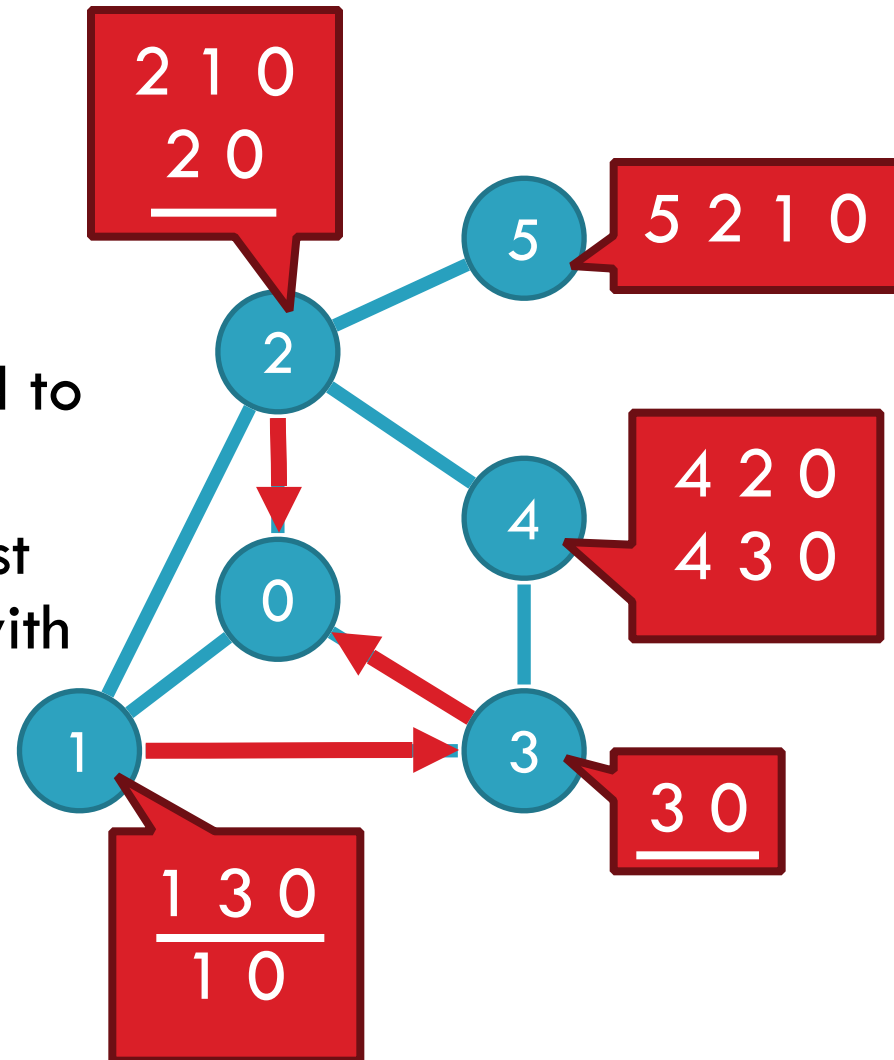
- A solution is an assignment of permitted paths to each node such that:
 - ▣ Node u 's path is either null or uwP , where path w is assigned to node w and edge $u \rightarrow w$ exists
 - ▣ Each node is assigned the highest ranked path that is consistent with their neighbors



A Solution to the SPP

29

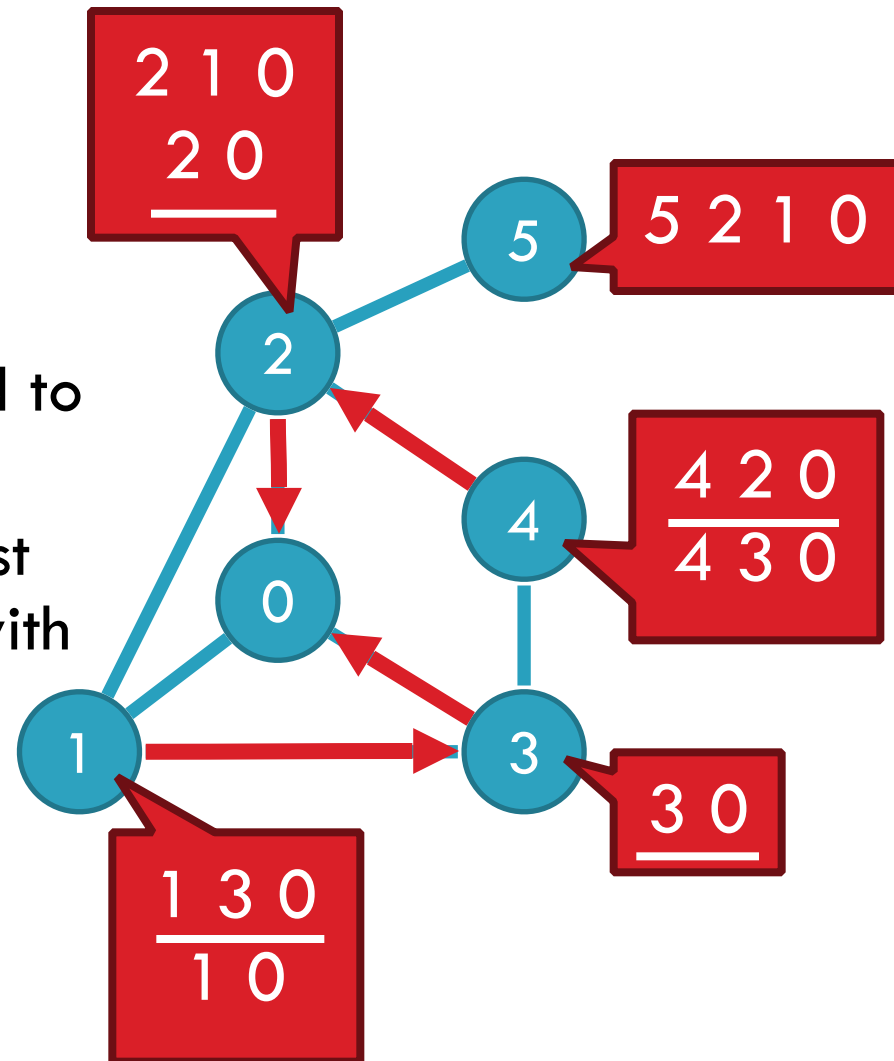
- A solution is an assignment of permitted paths to each node such that:
 - ▣ Node u 's path is either null or uwP , where path w is assigned to node w and edge $u \rightarrow w$ exists
 - ▣ Each node is assigned the highest ranked path that is consistent with their neighbors



A Solution to the SPP

29

- A solution is an assignment of permitted paths to each node such that:
 - ▣ Node u 's path is either null or uwP , where path w is assigned to node w and edge $u \rightarrow w$ exists
 - ▣ Each node is assigned the highest ranked path that is consistent with their neighbors

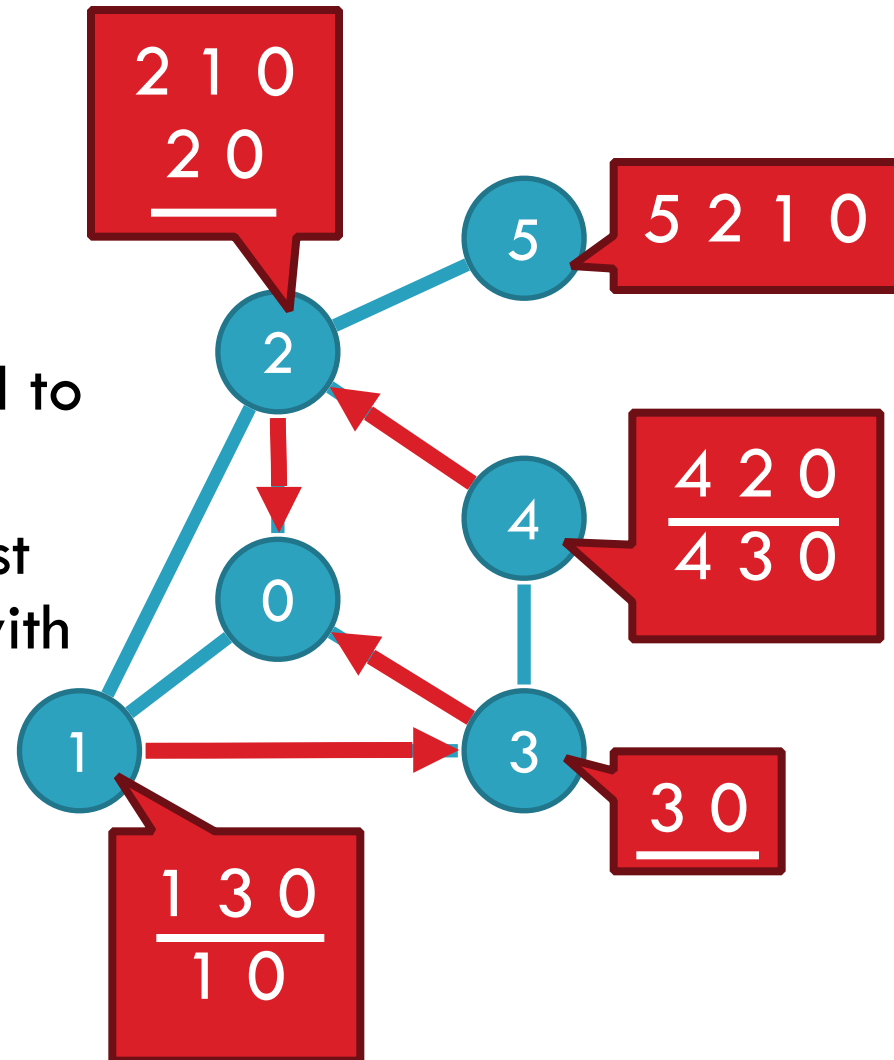


A Solution to the SPP

29

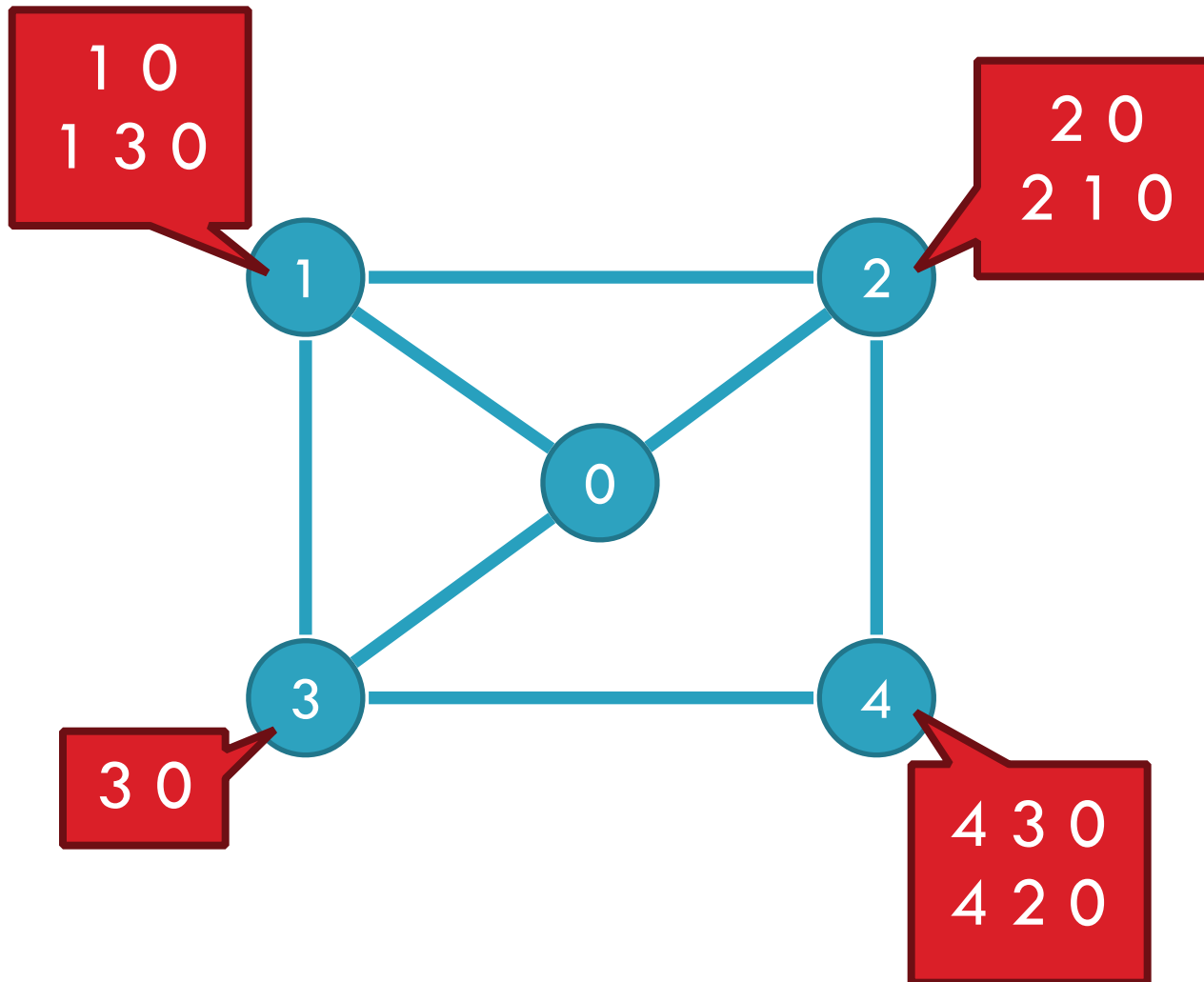
- A solution is an assignment of permitted paths to each node such that

Solutions need not use the shortest paths, or form a spanning tree



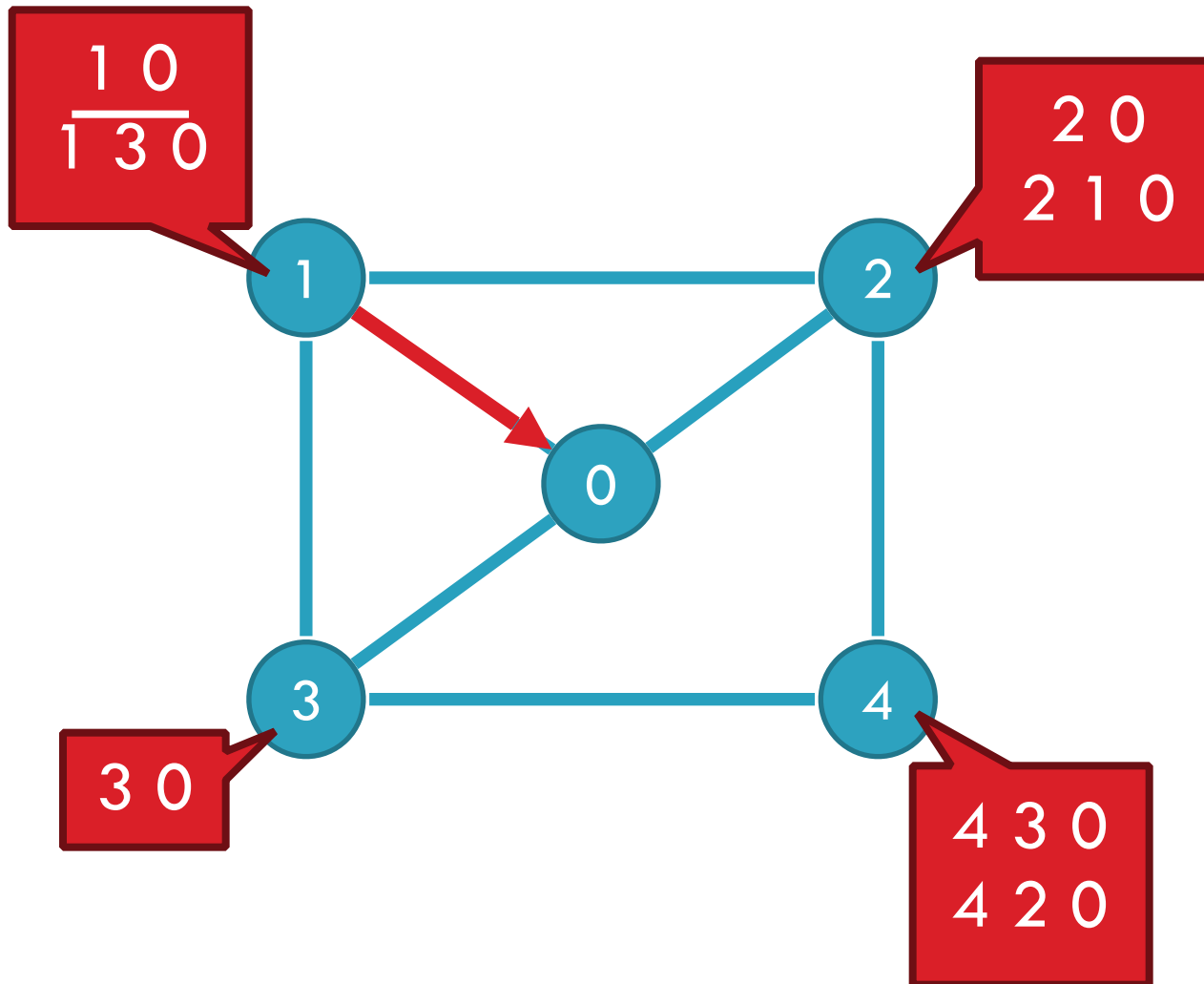
Simple SPP Example

30



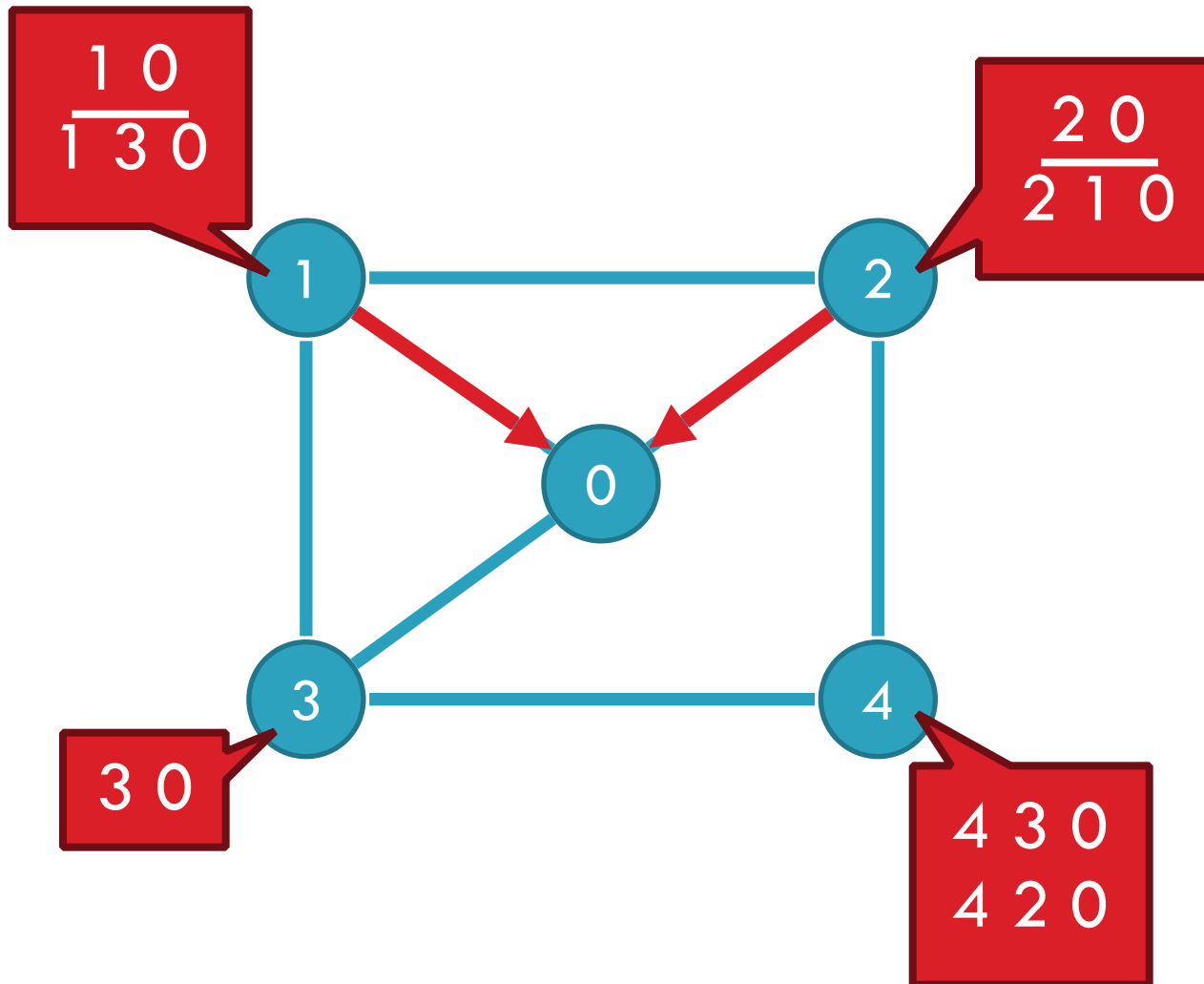
Simple SPP Example

30



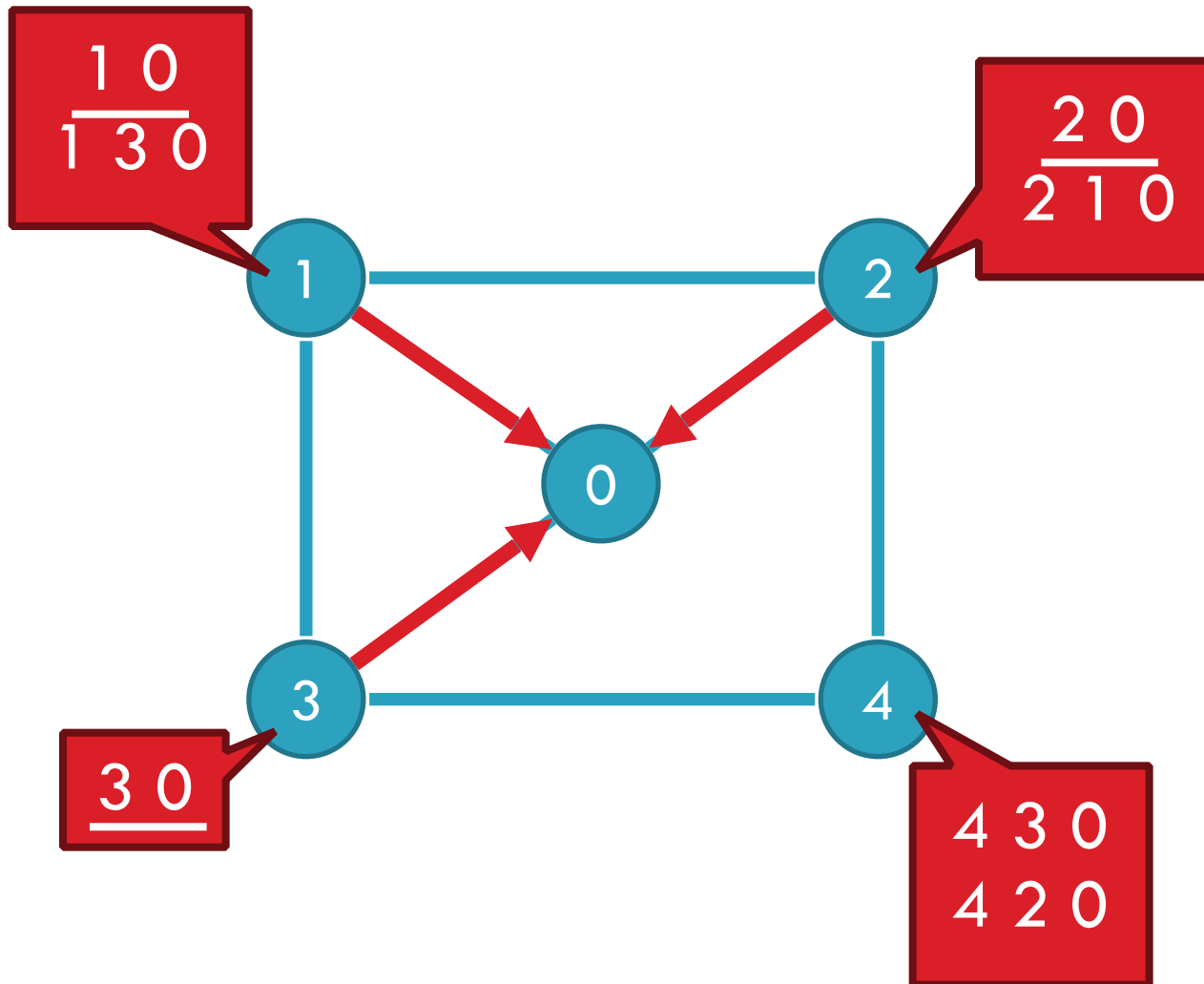
Simple SPP Example

30



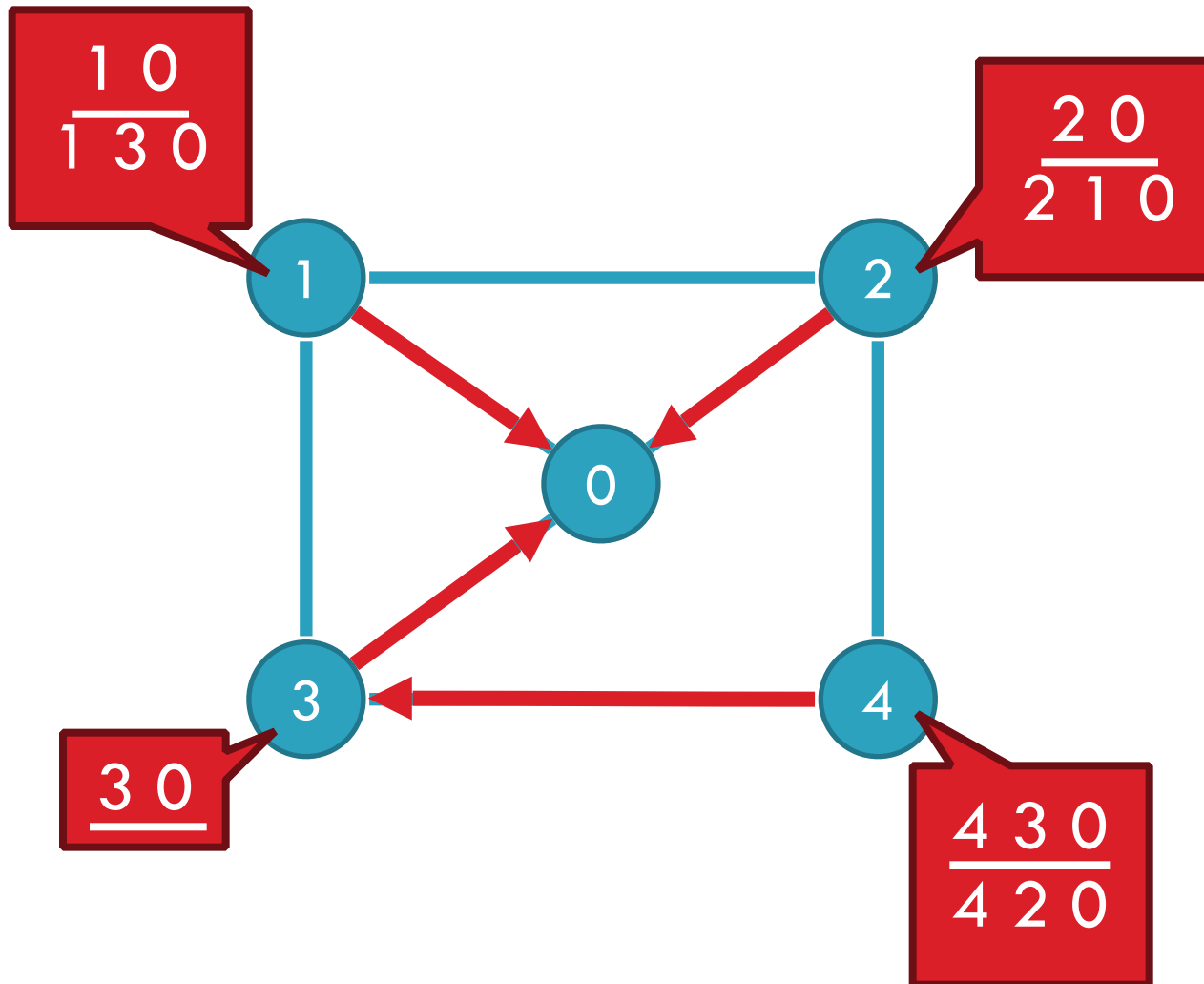
Simple SPP Example

30



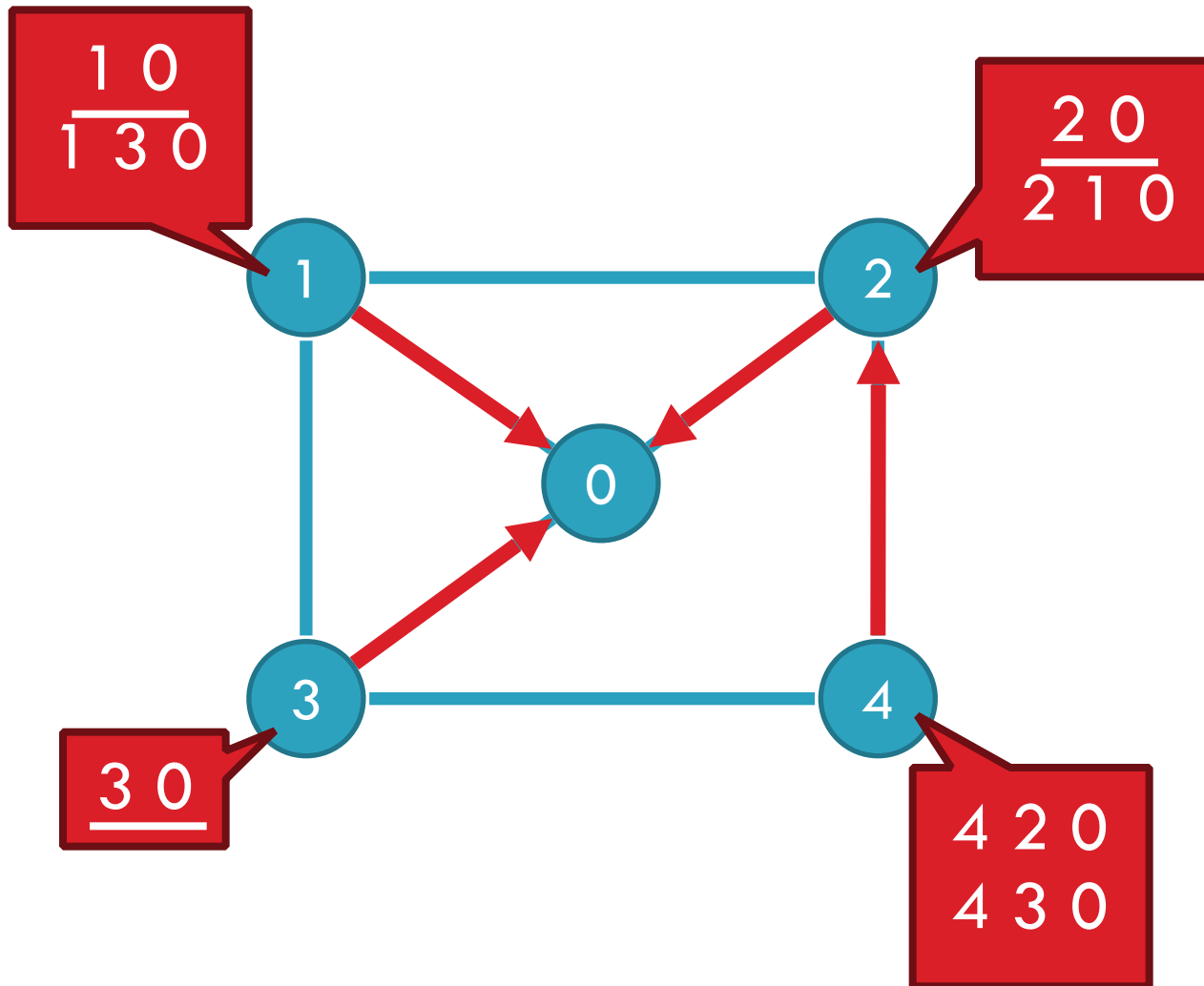
Simple SPP Example

30



Simple SPP Example

30



Simple SPP Example

30

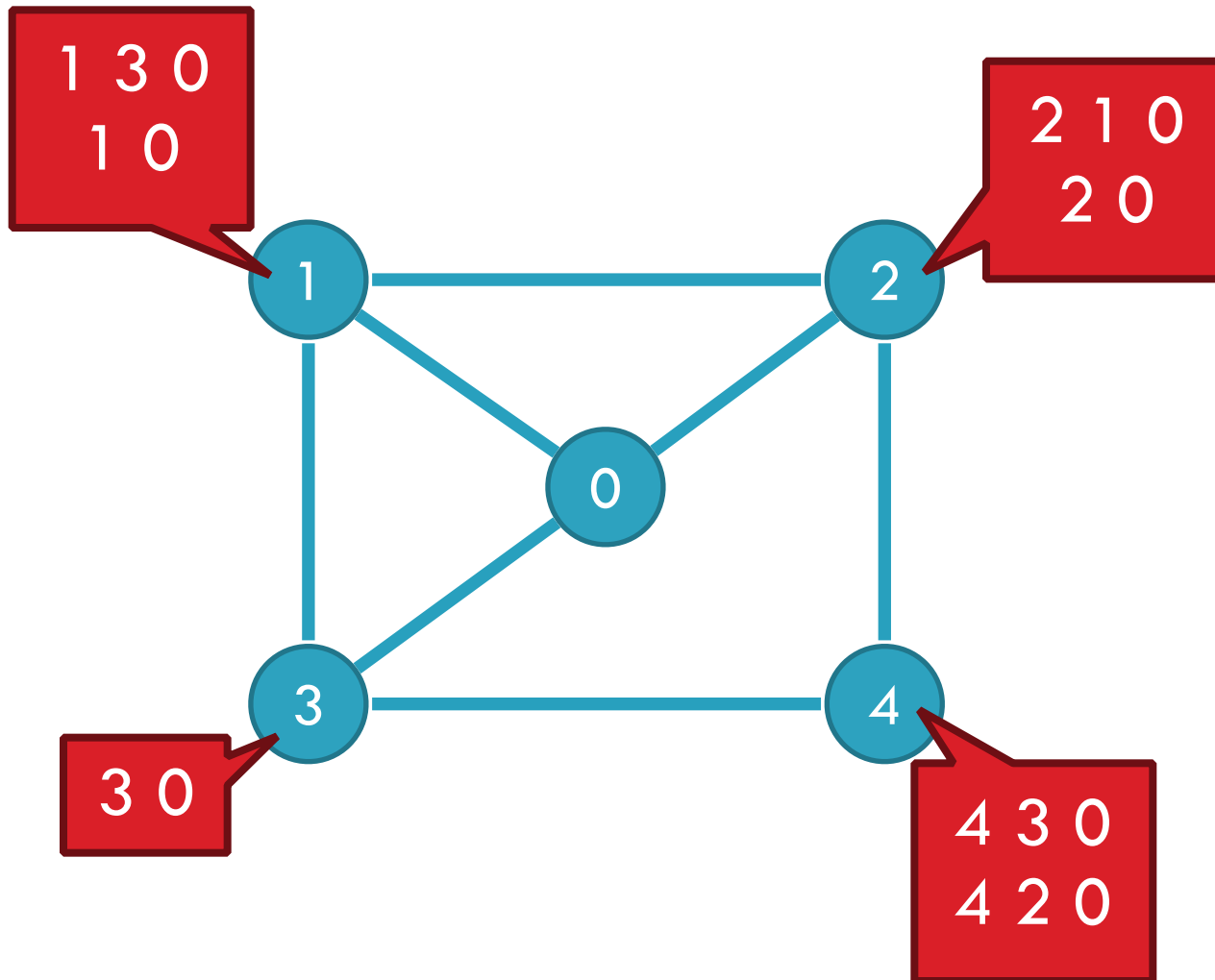


- Each node gets its preferred route
- Totally stable topology



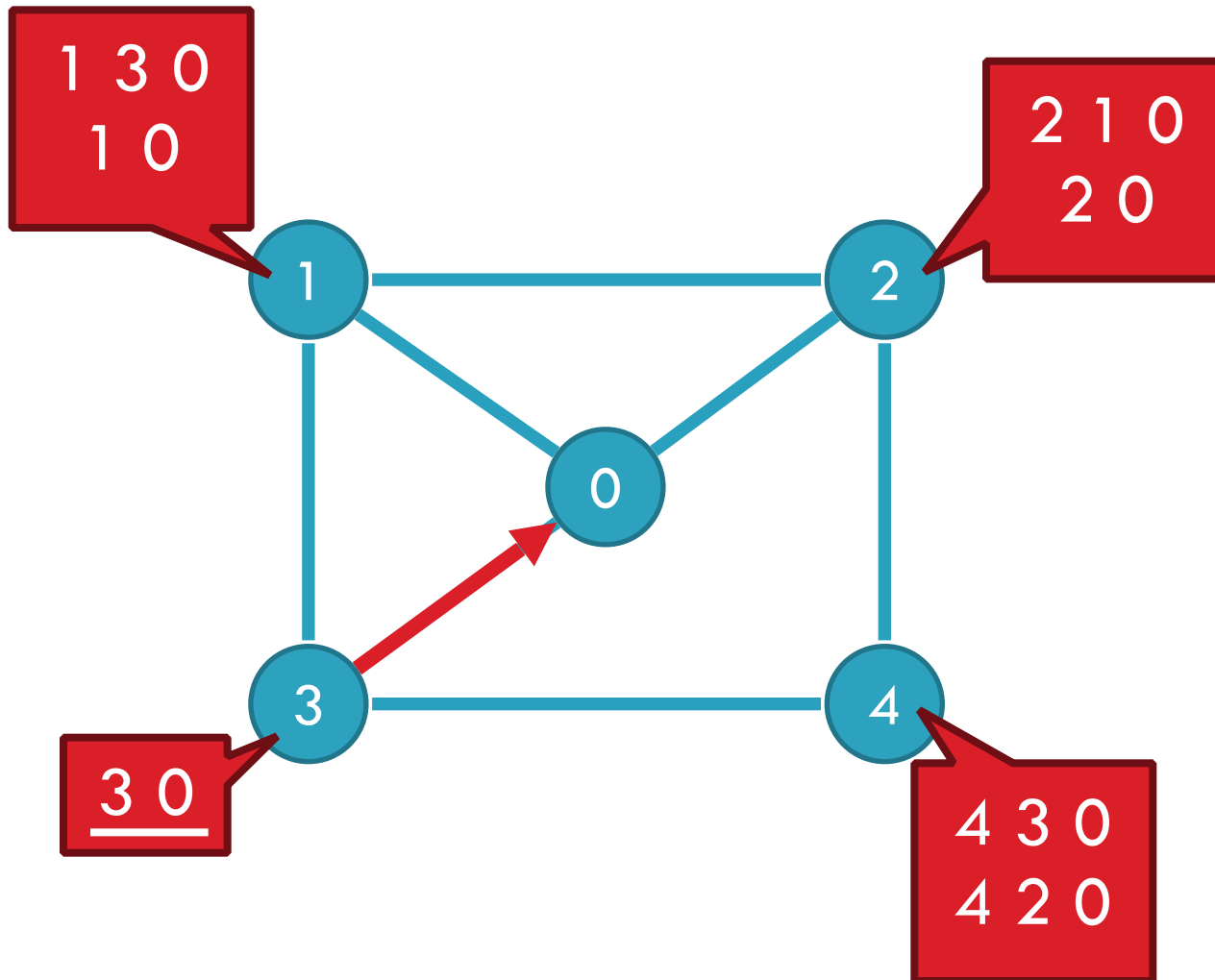
Good Gadget

31



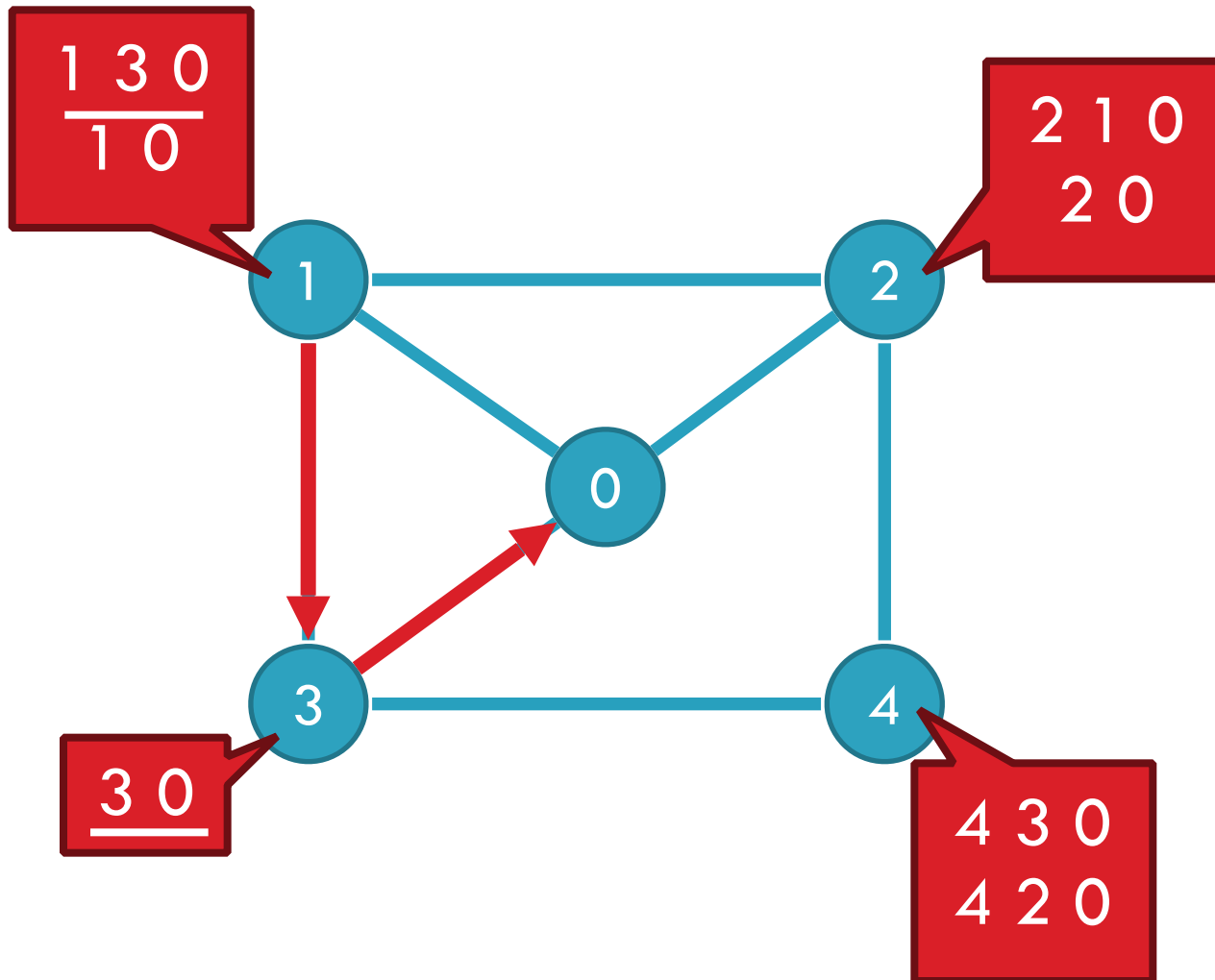
Good Gadget

31



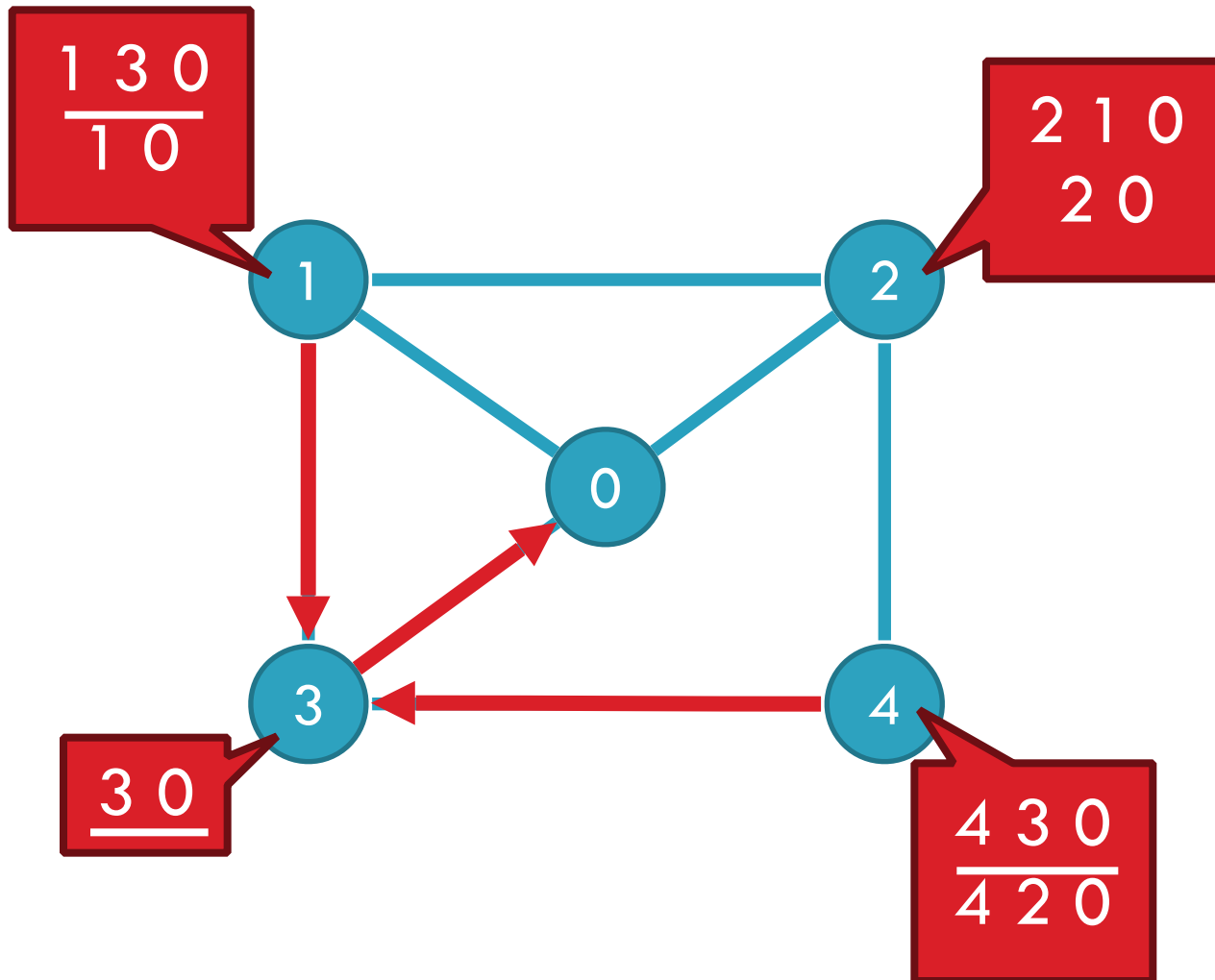
Good Gadget

31



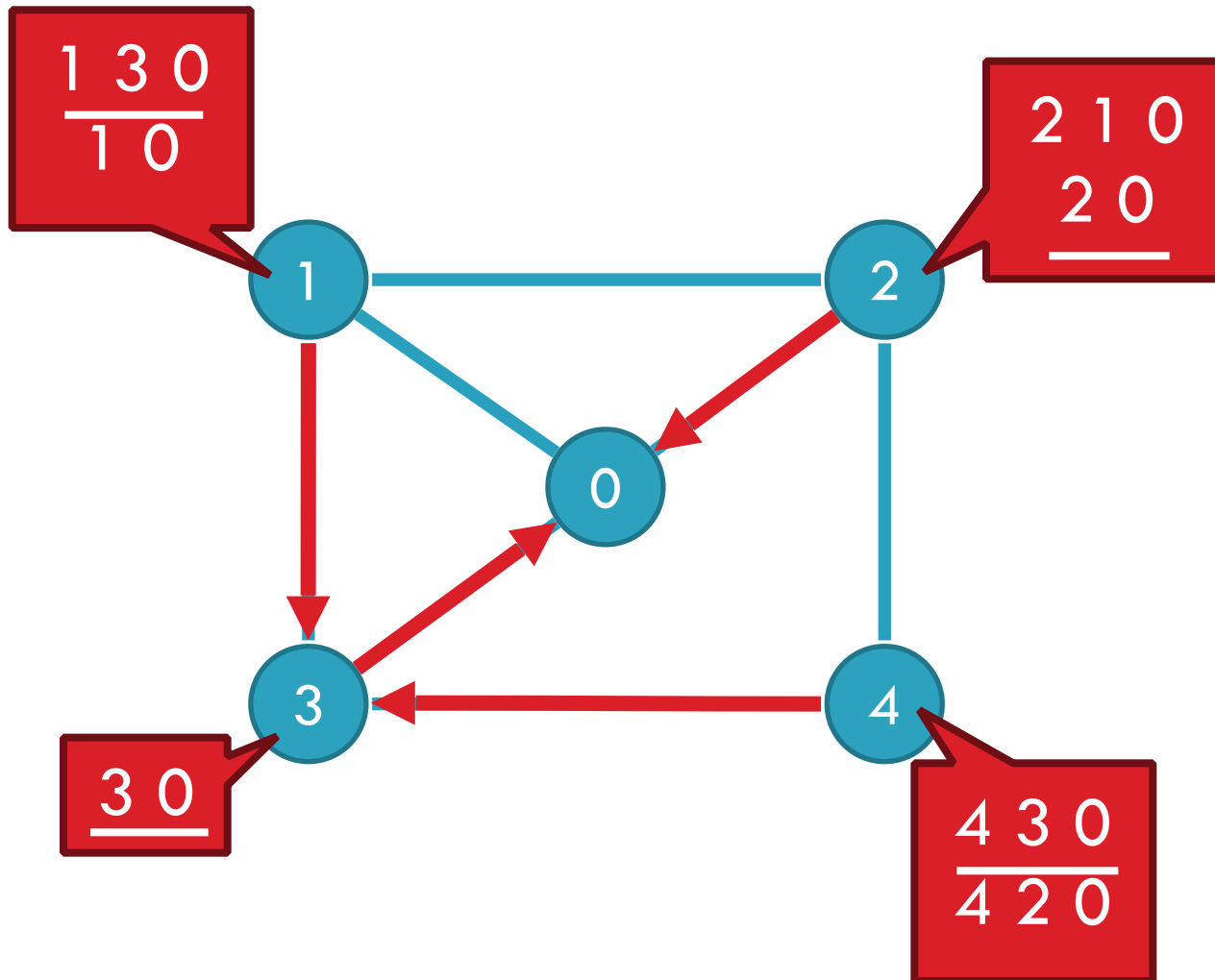
Good Gadget

31



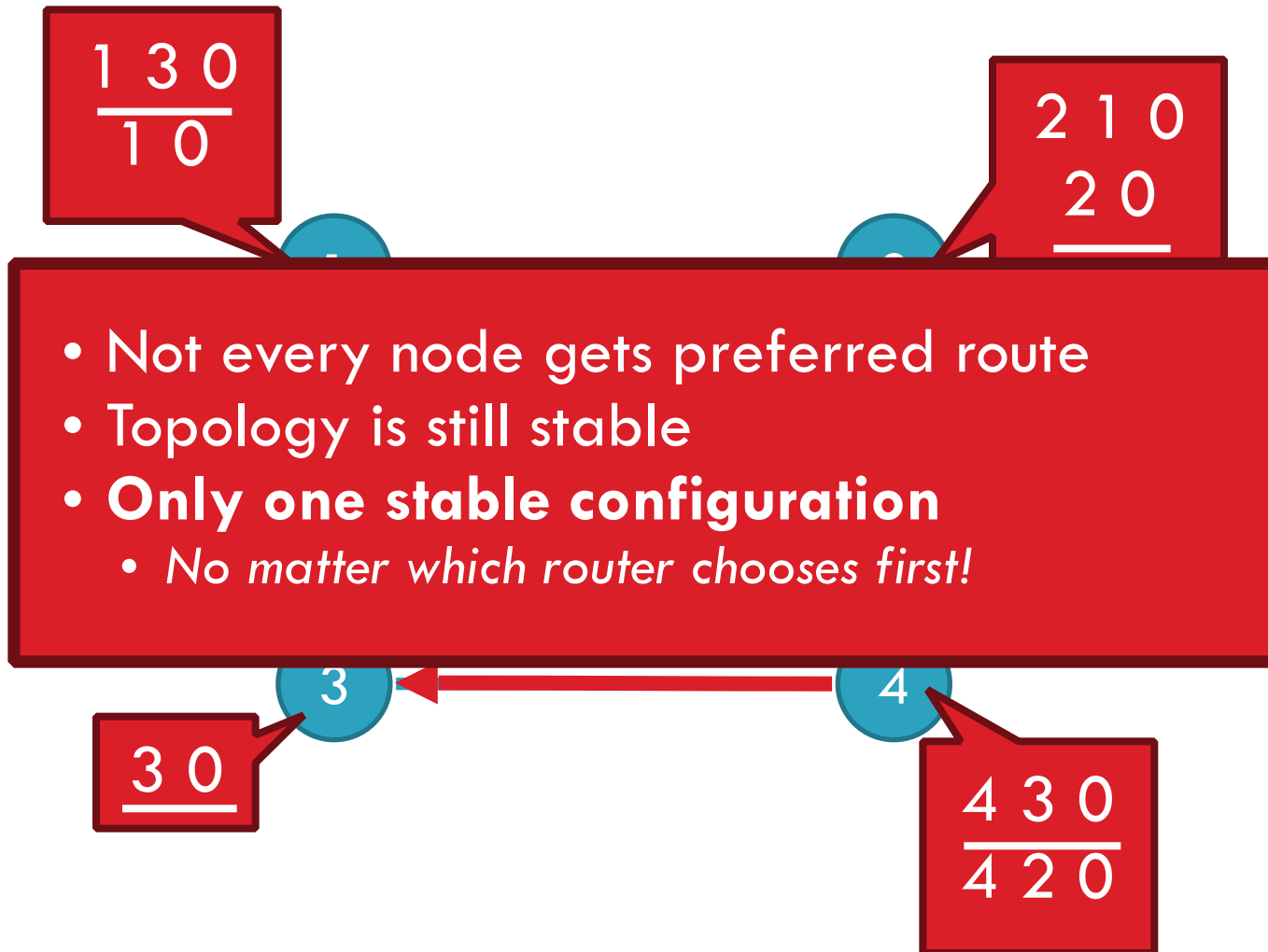
Good Gadget

31



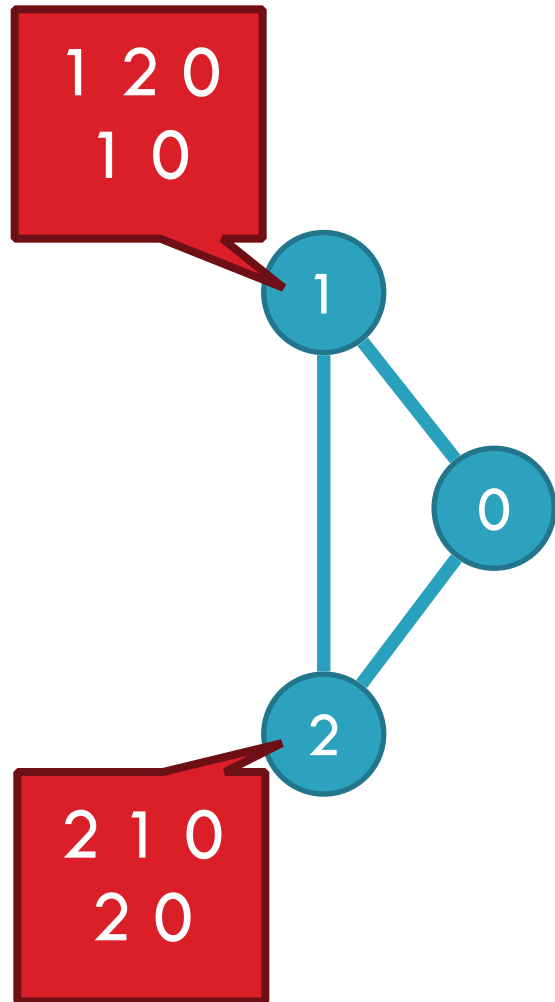
Good Gadget

31



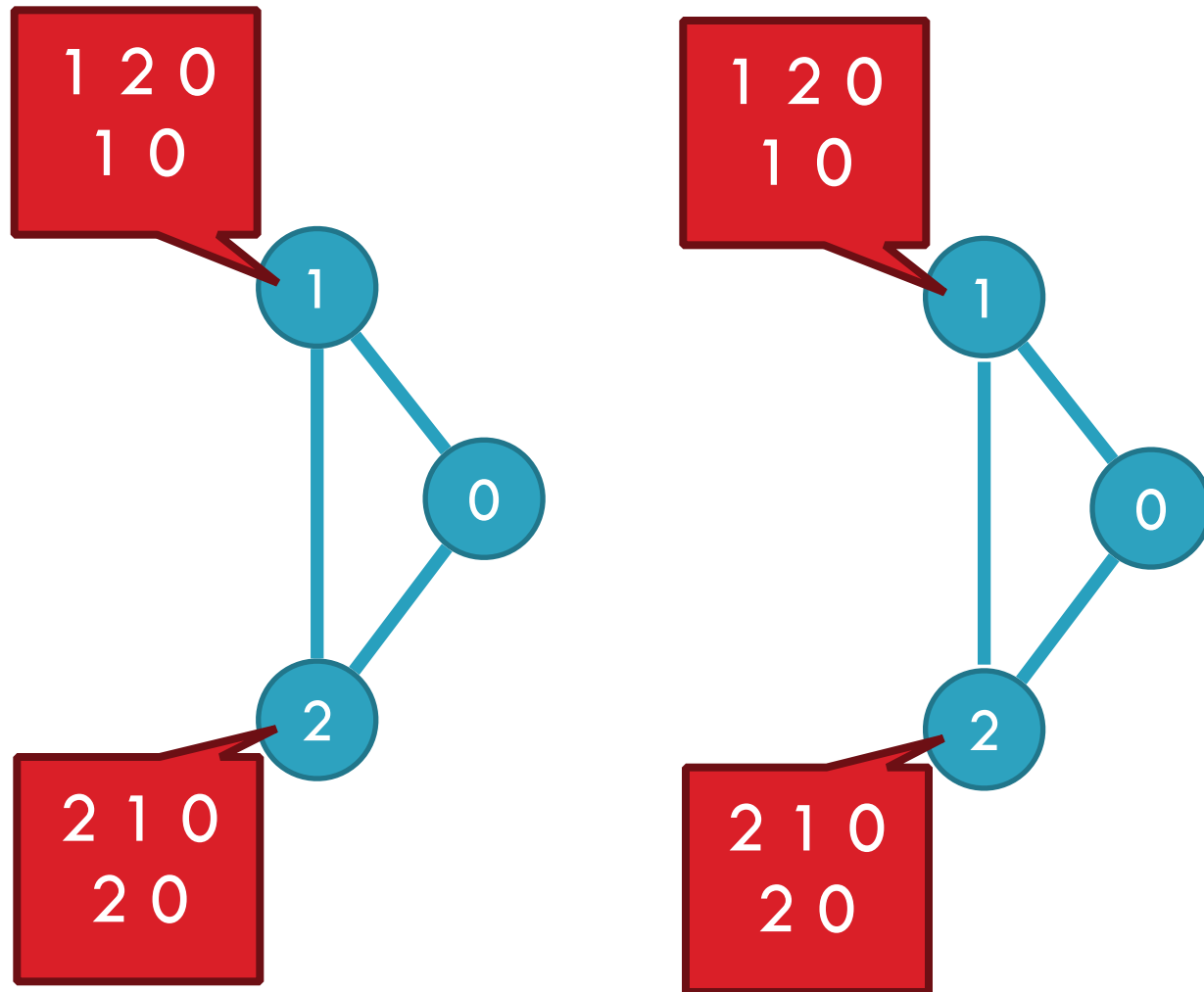
SPP May Have Multiple Solutions

32



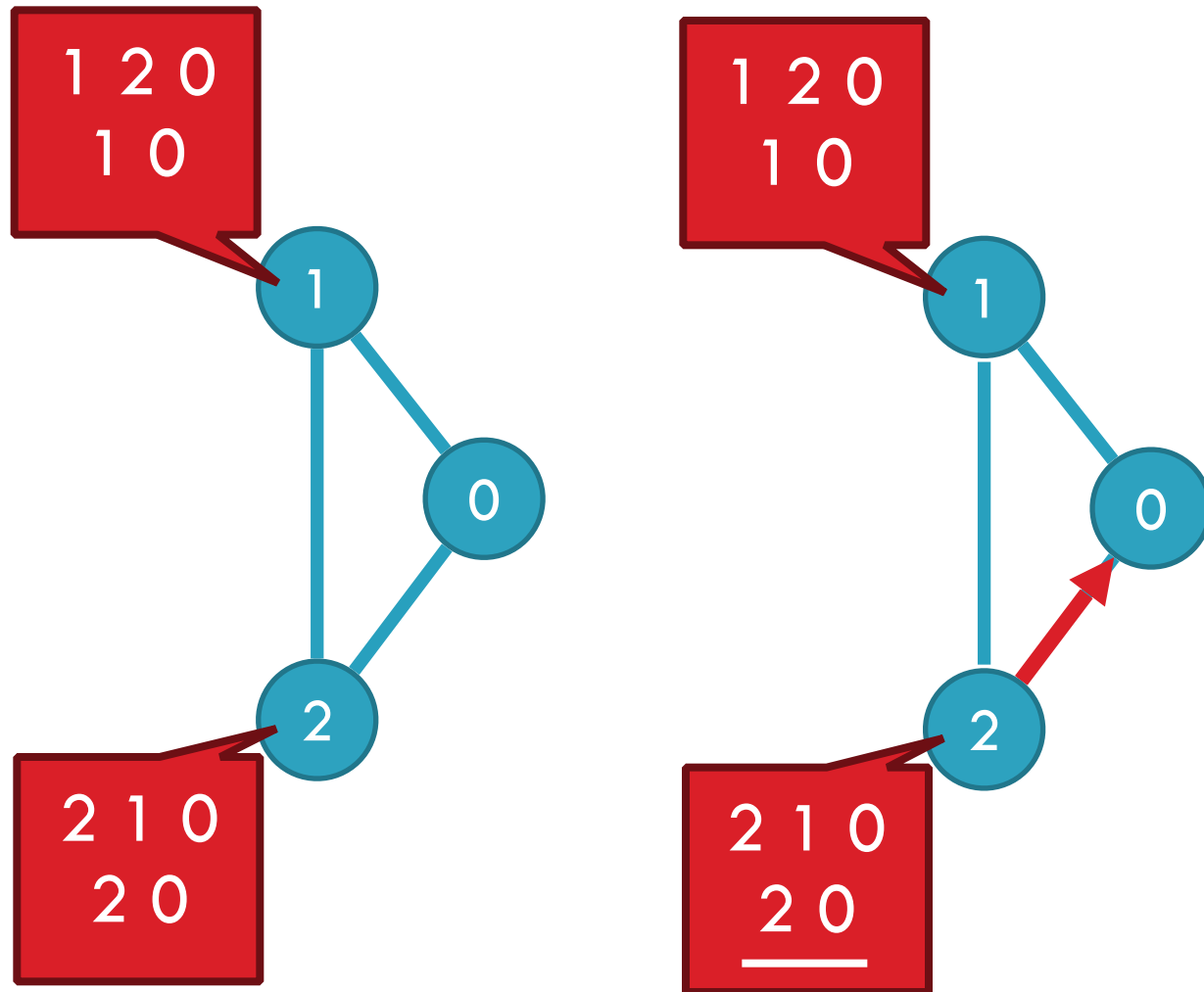
SPP May Have Multiple Solutions

32



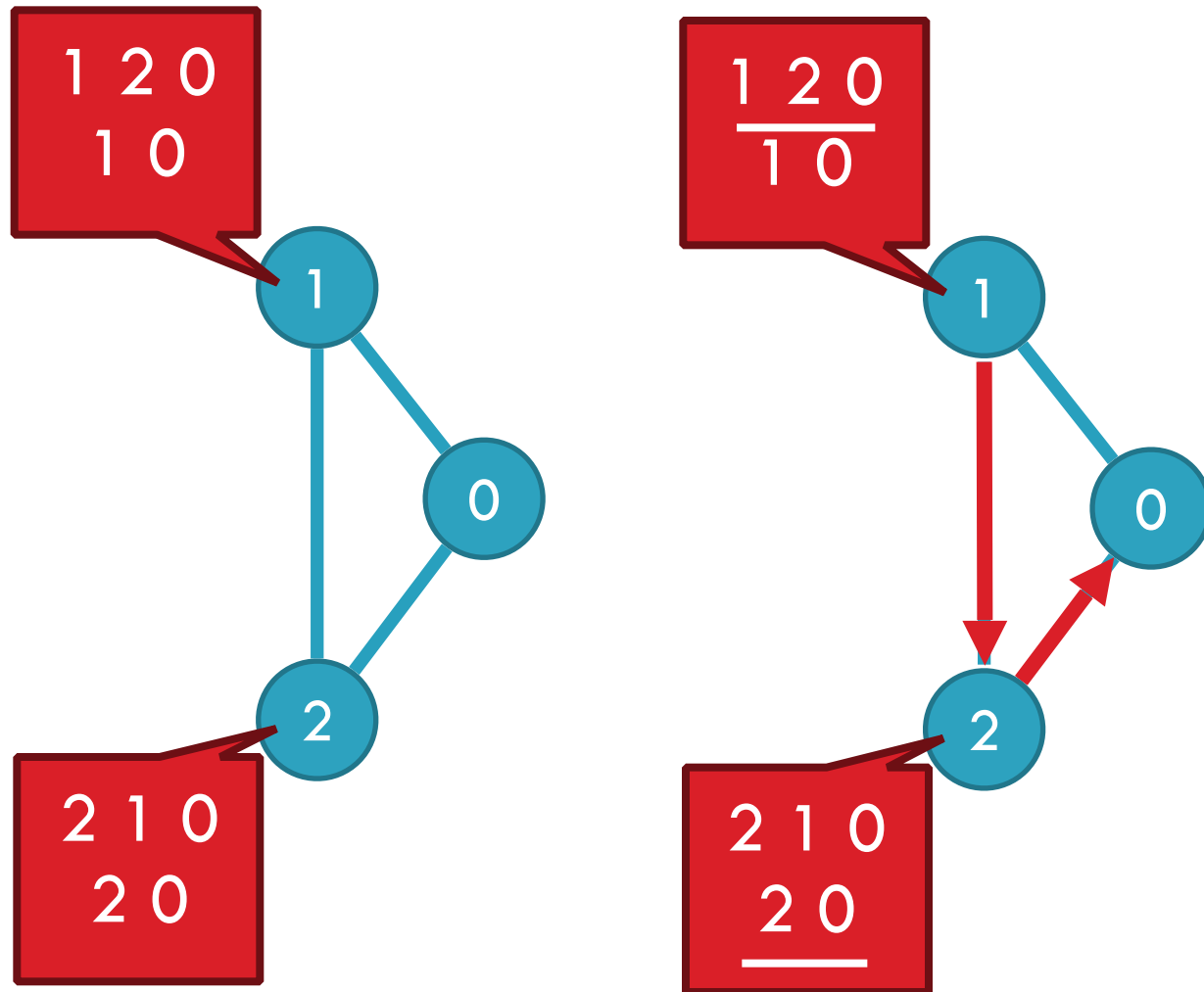
SPP May Have Multiple Solutions

32



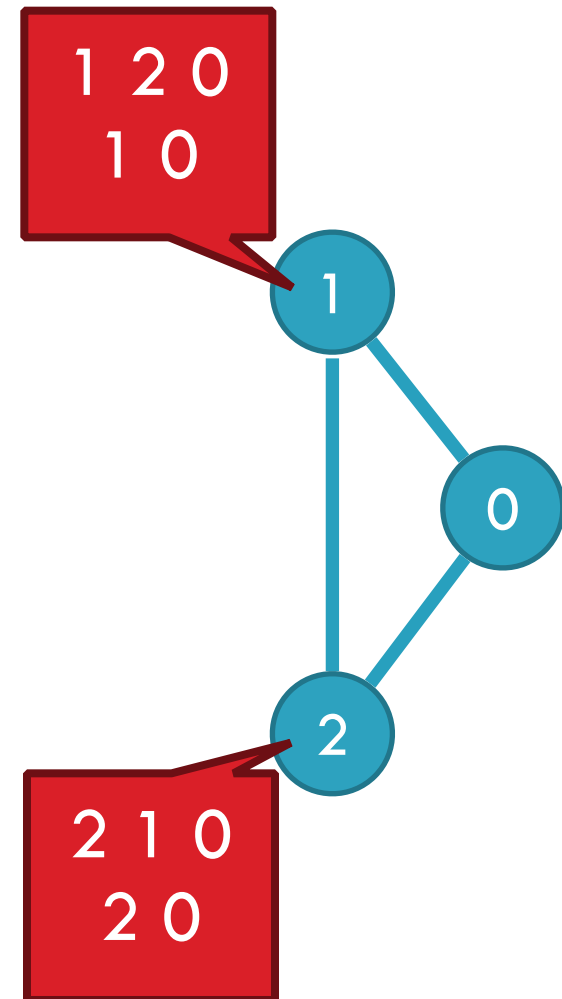
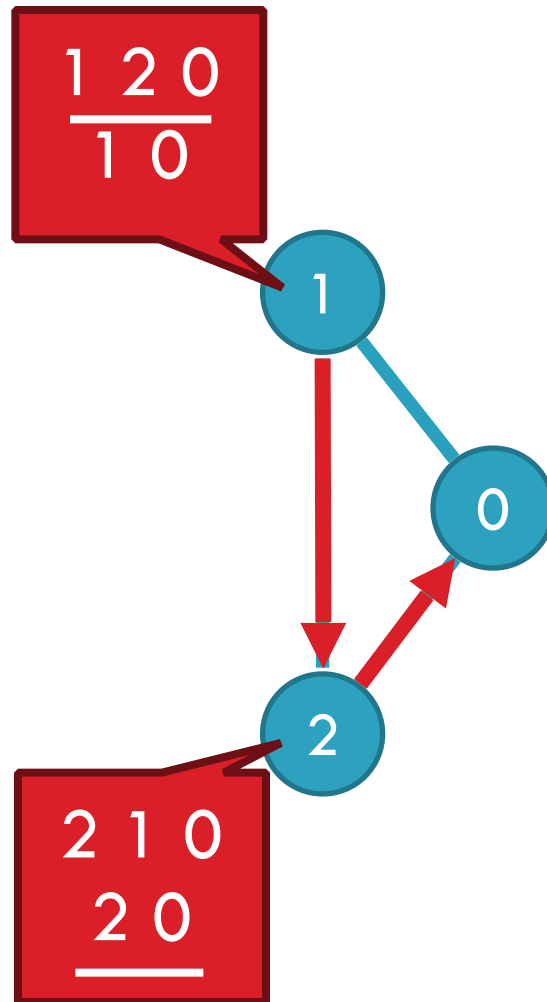
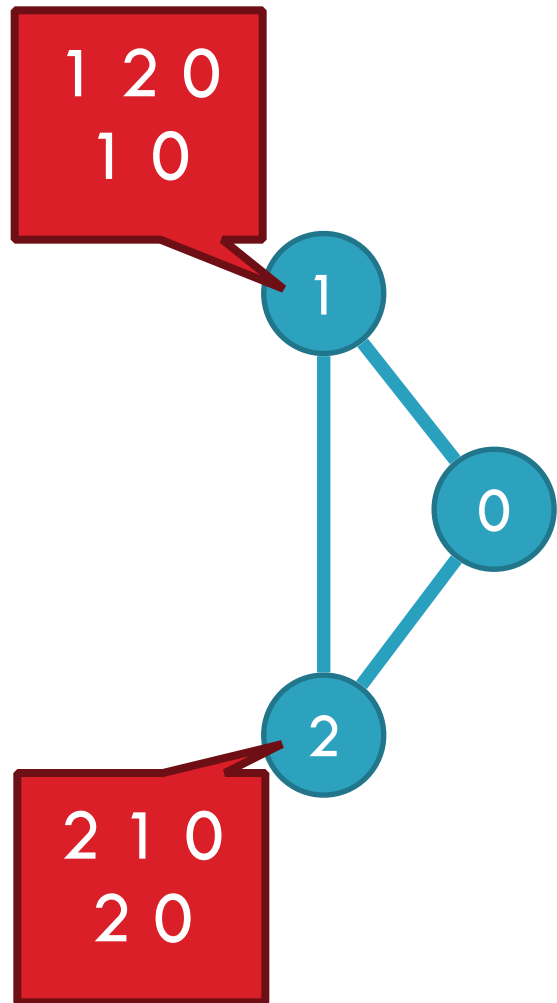
SPP May Have Multiple Solutions

32



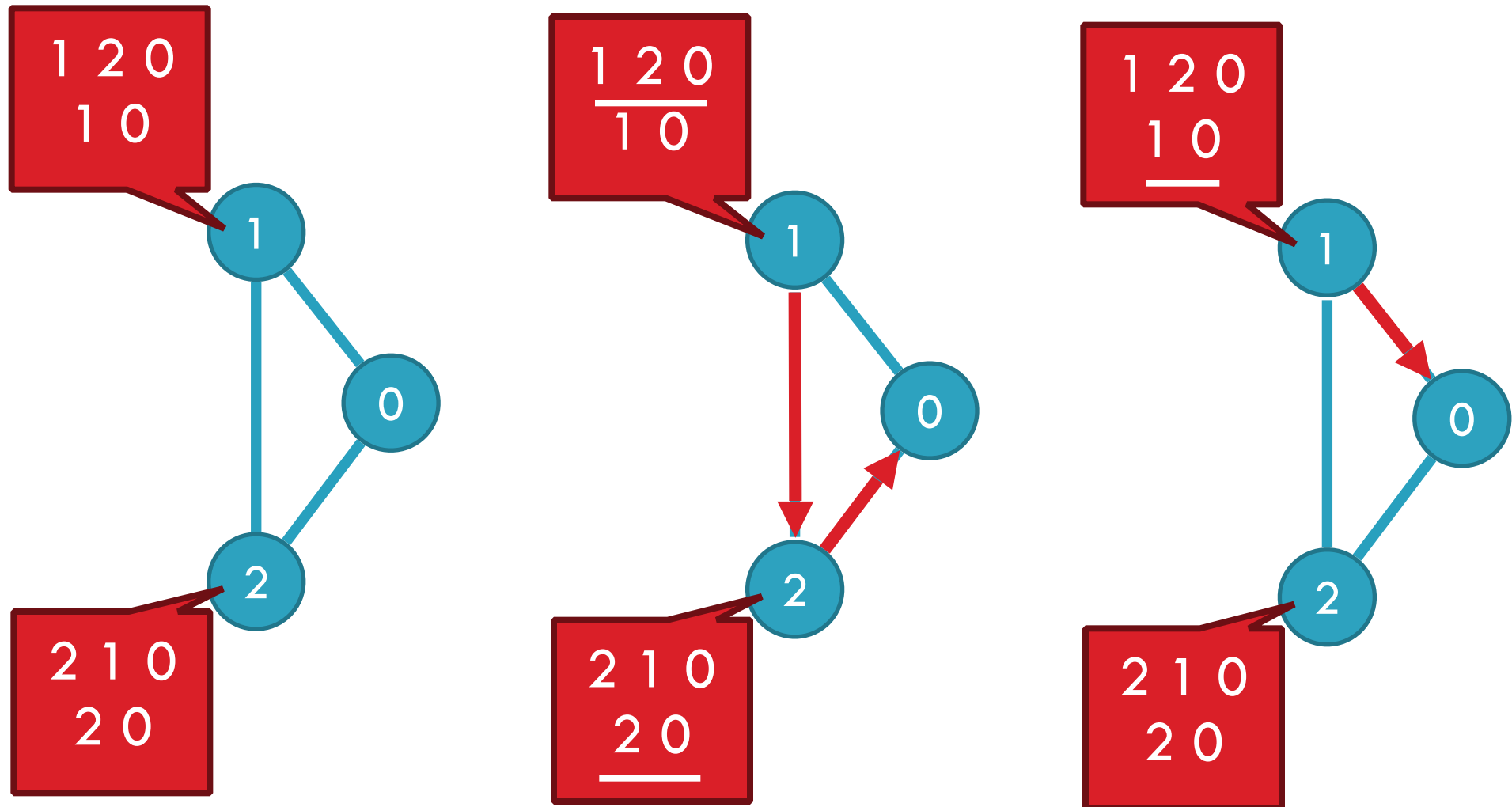
SPP May Have Multiple Solutions

32



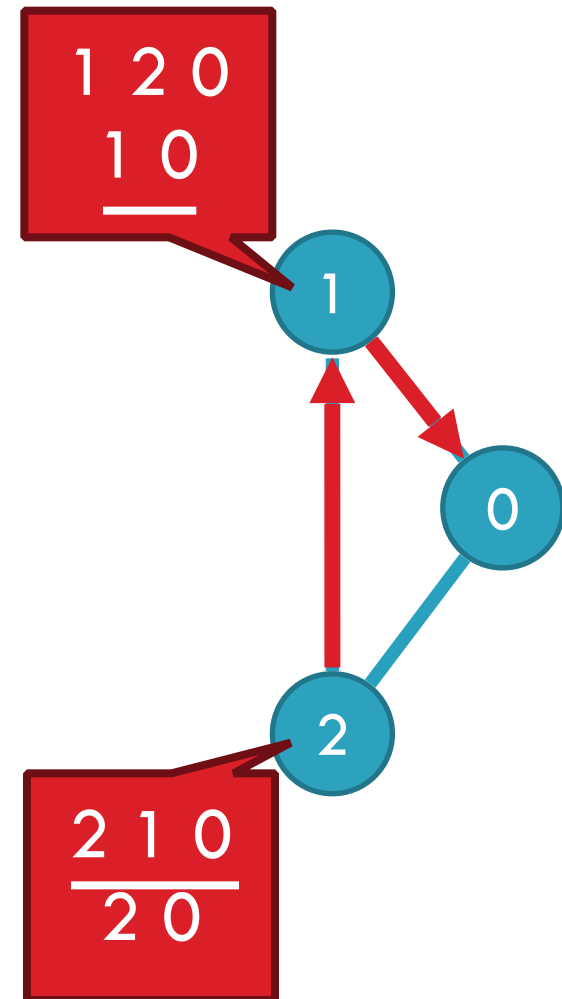
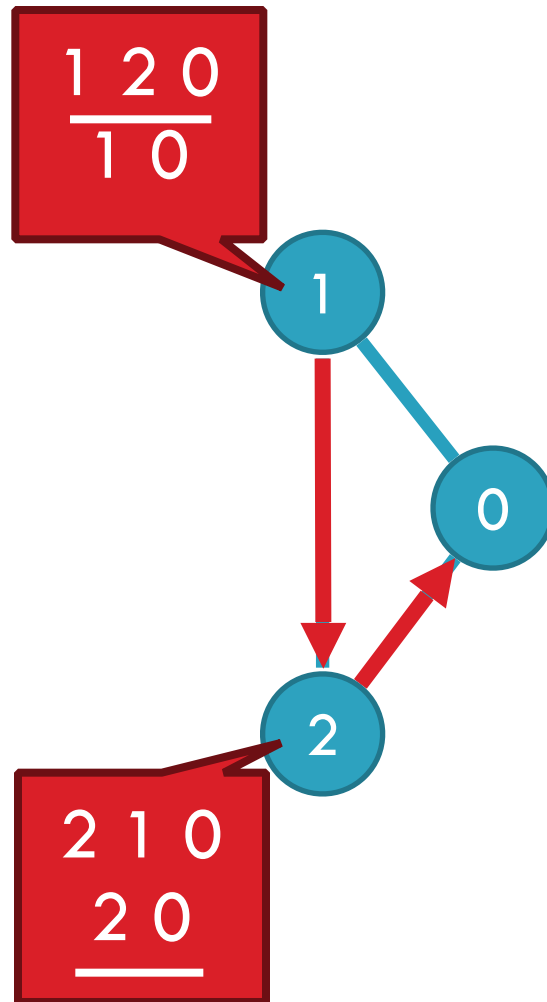
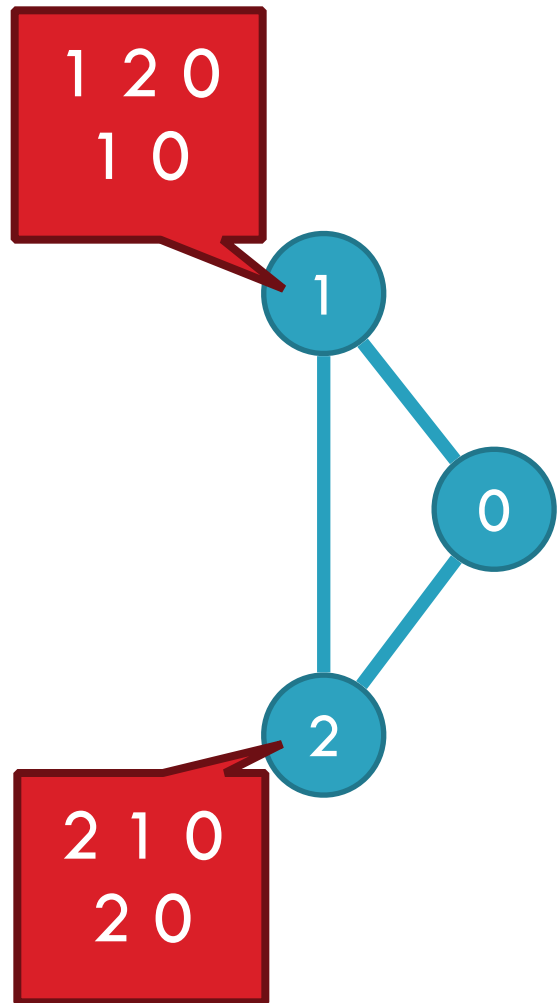
SPP May Have Multiple Solutions

32



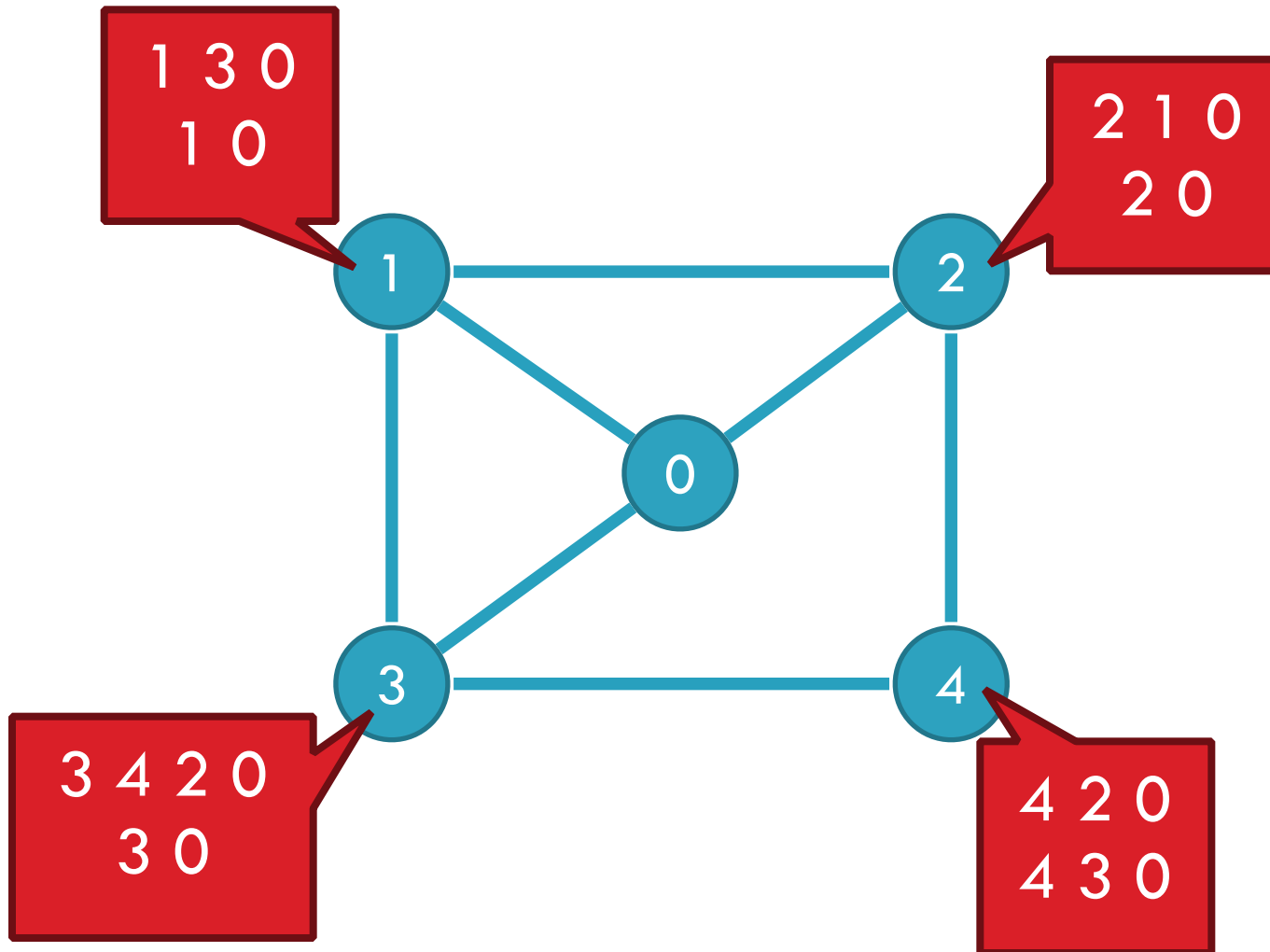
SPP May Have Multiple Solutions

32



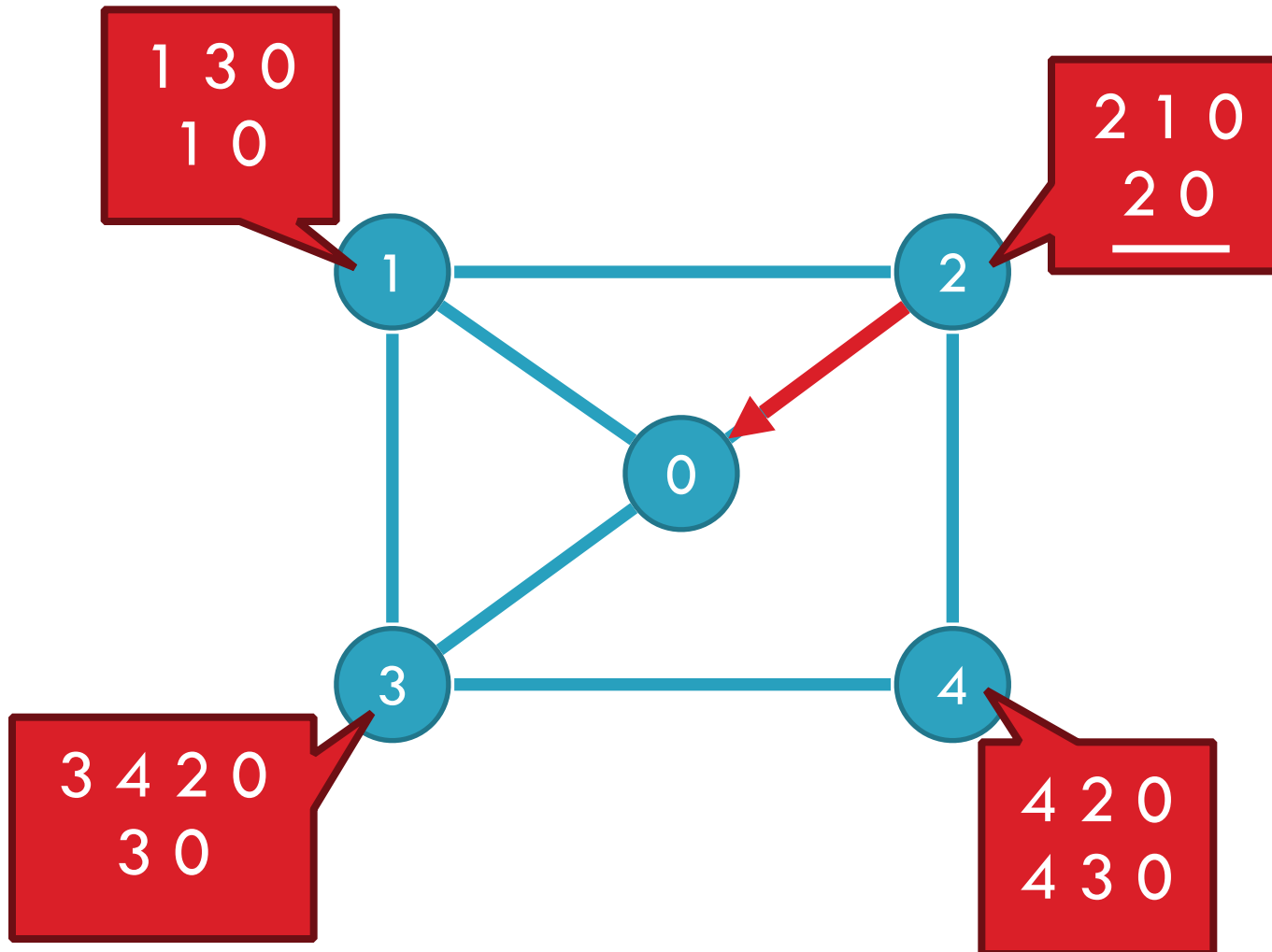
Bad Gadget

33



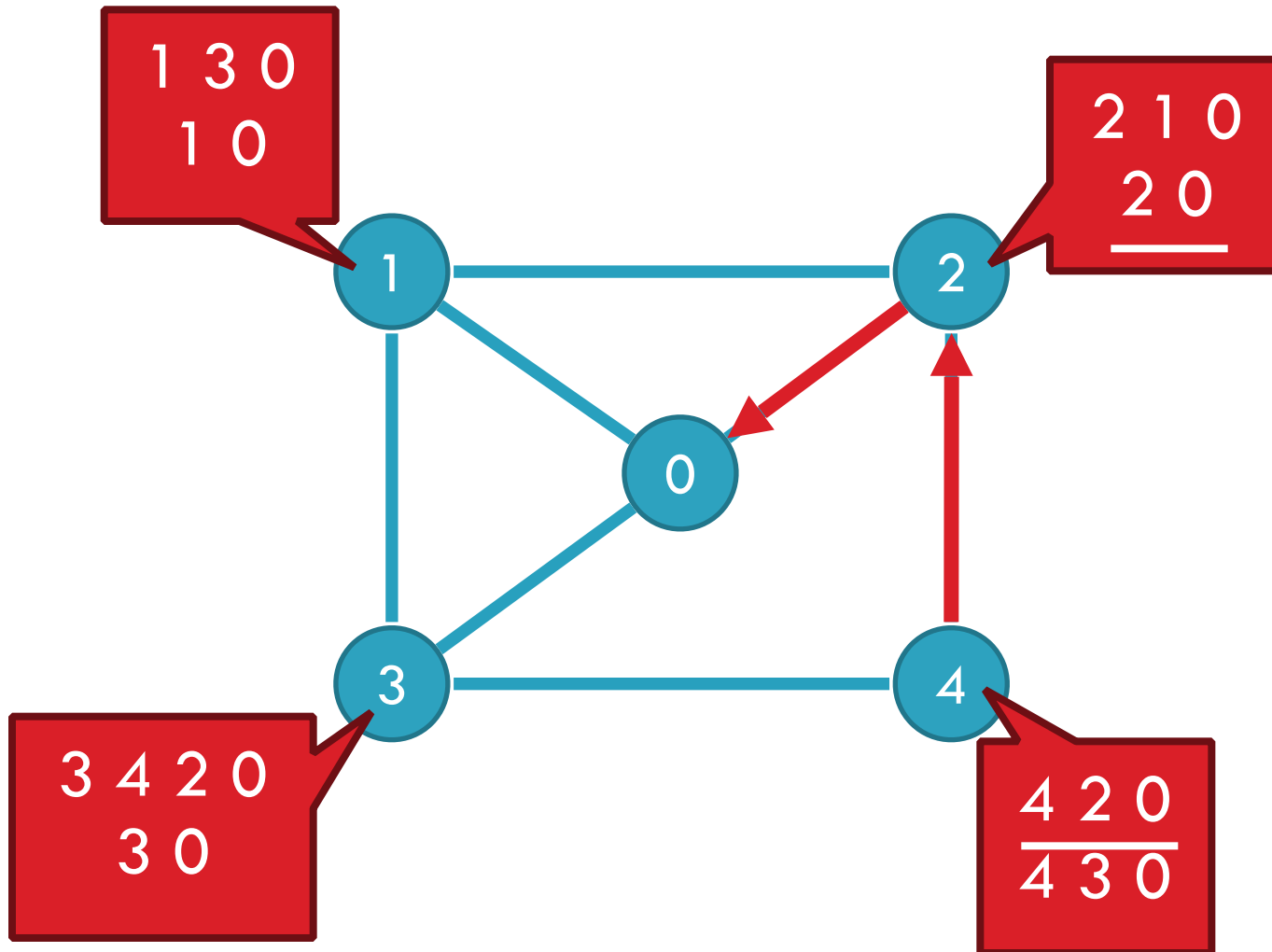
Bad Gadget

33



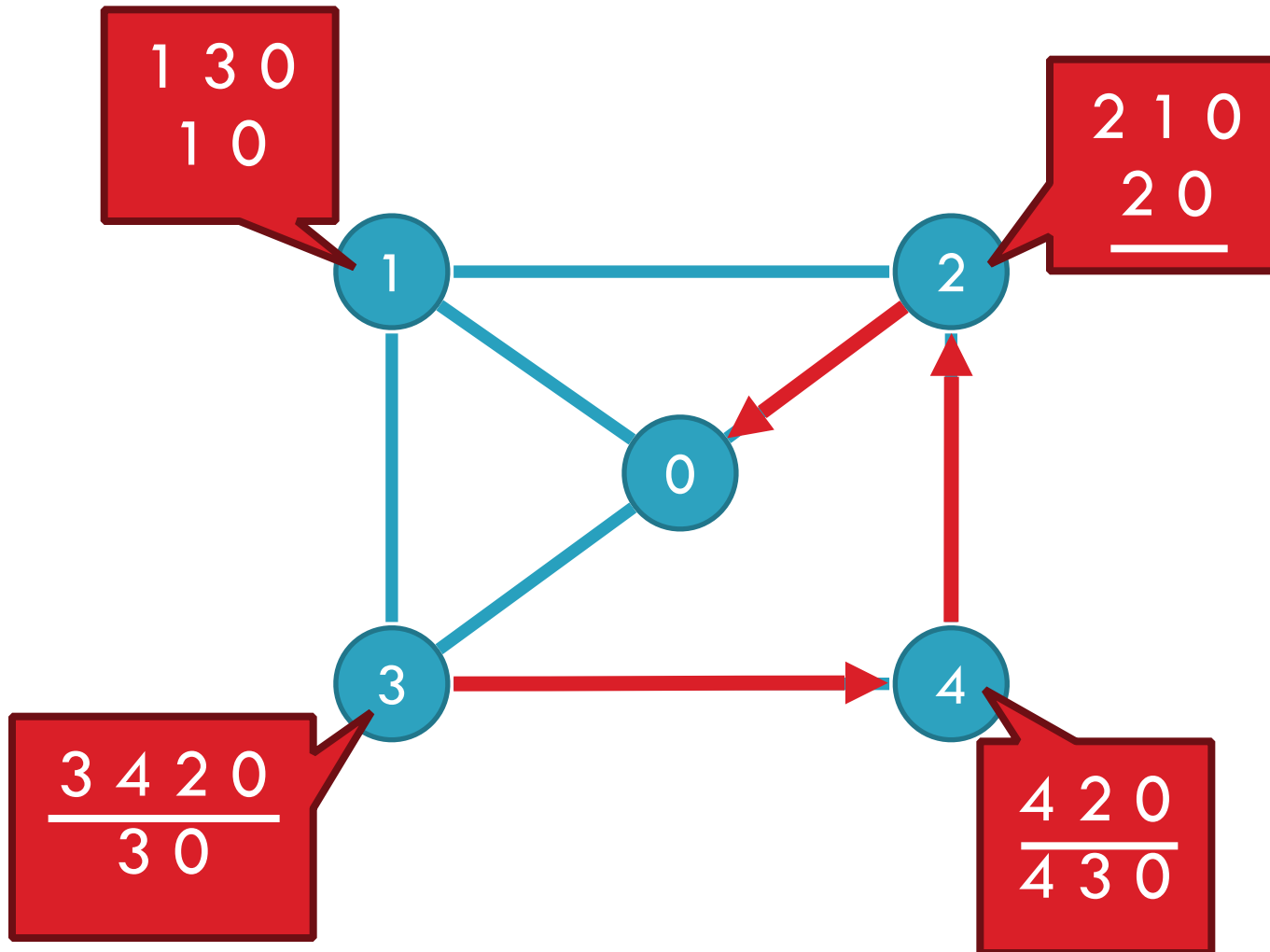
Bad Gadget

33



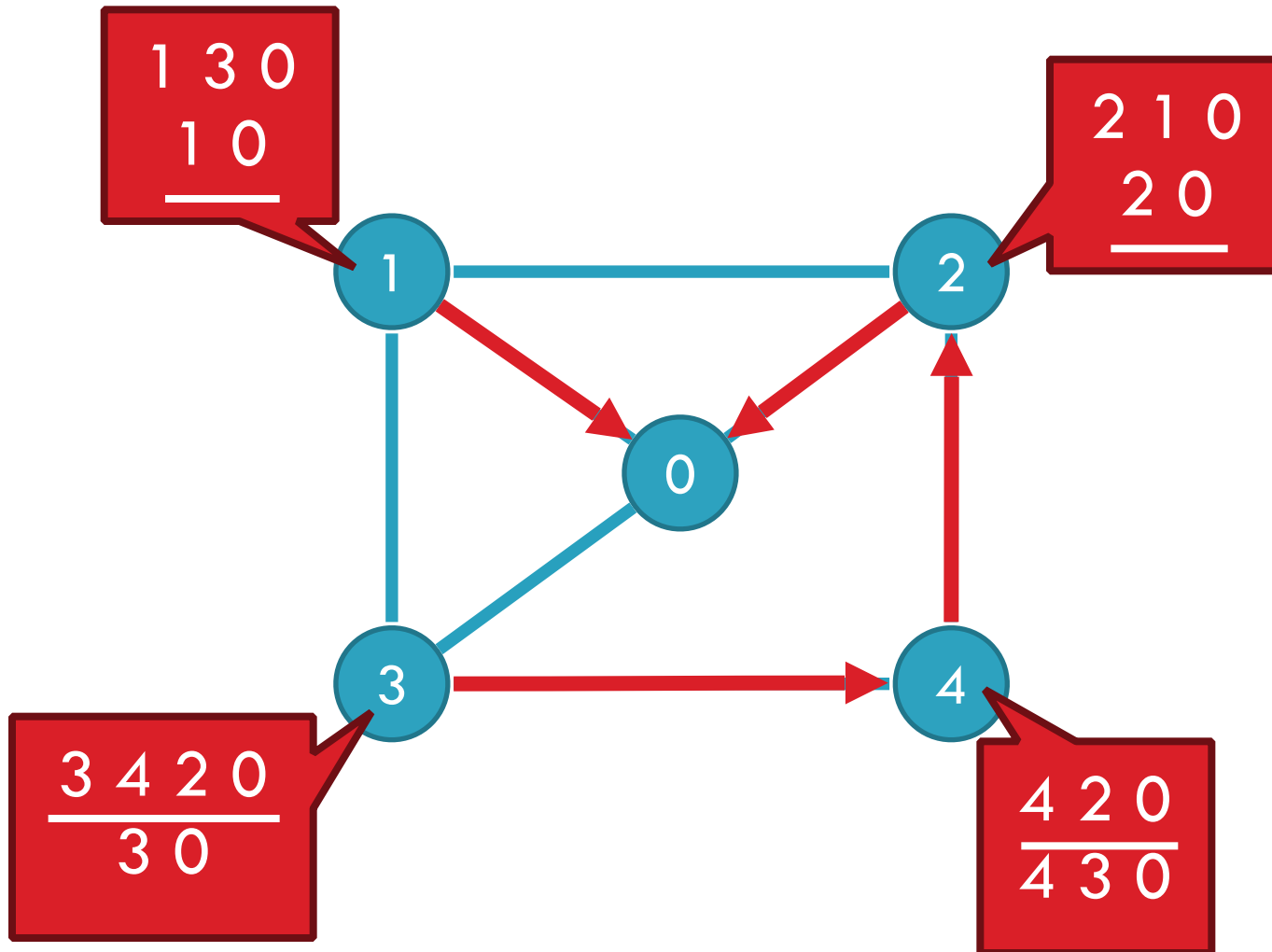
Bad Gadget

33



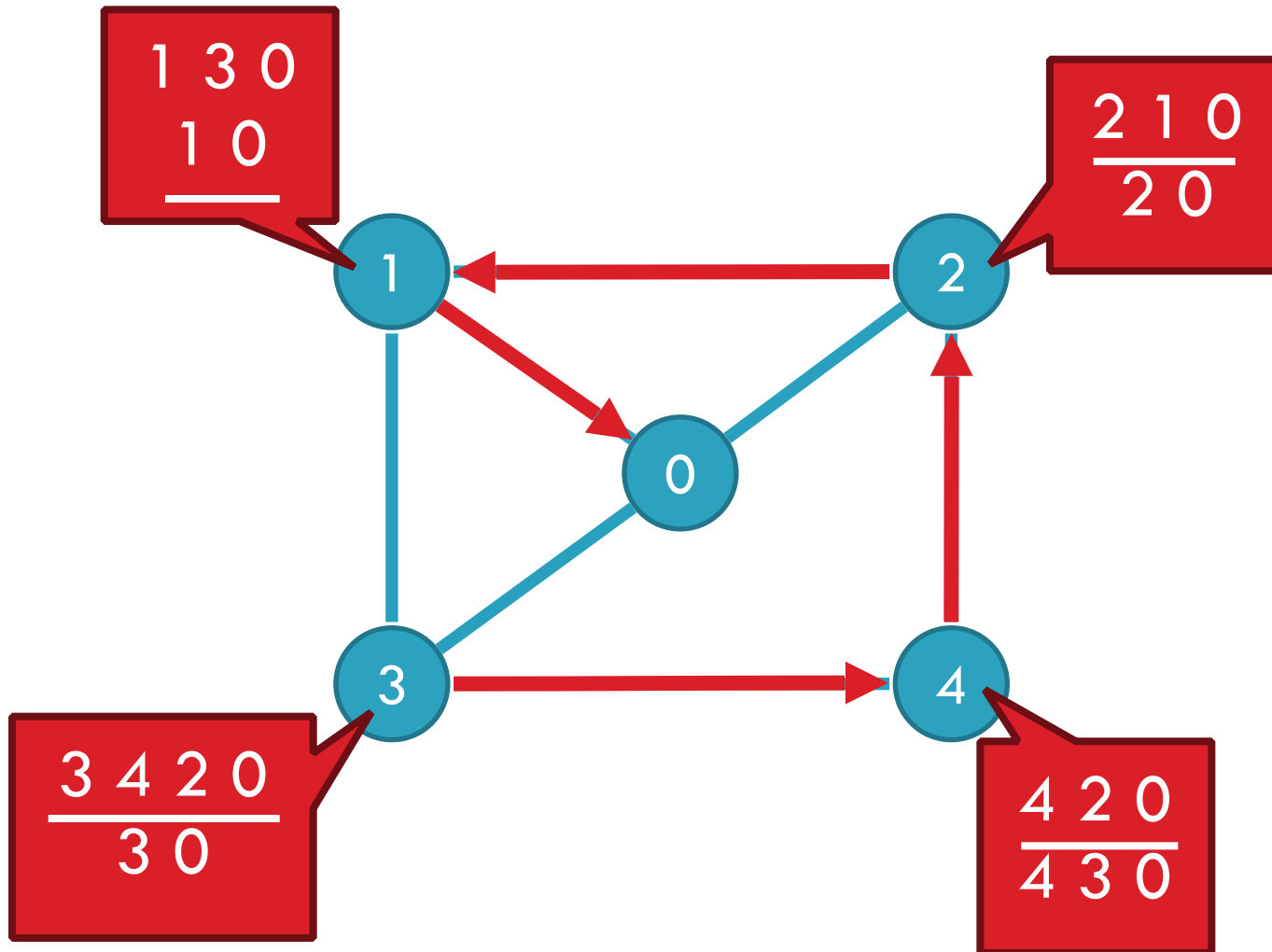
Bad Gadget

33



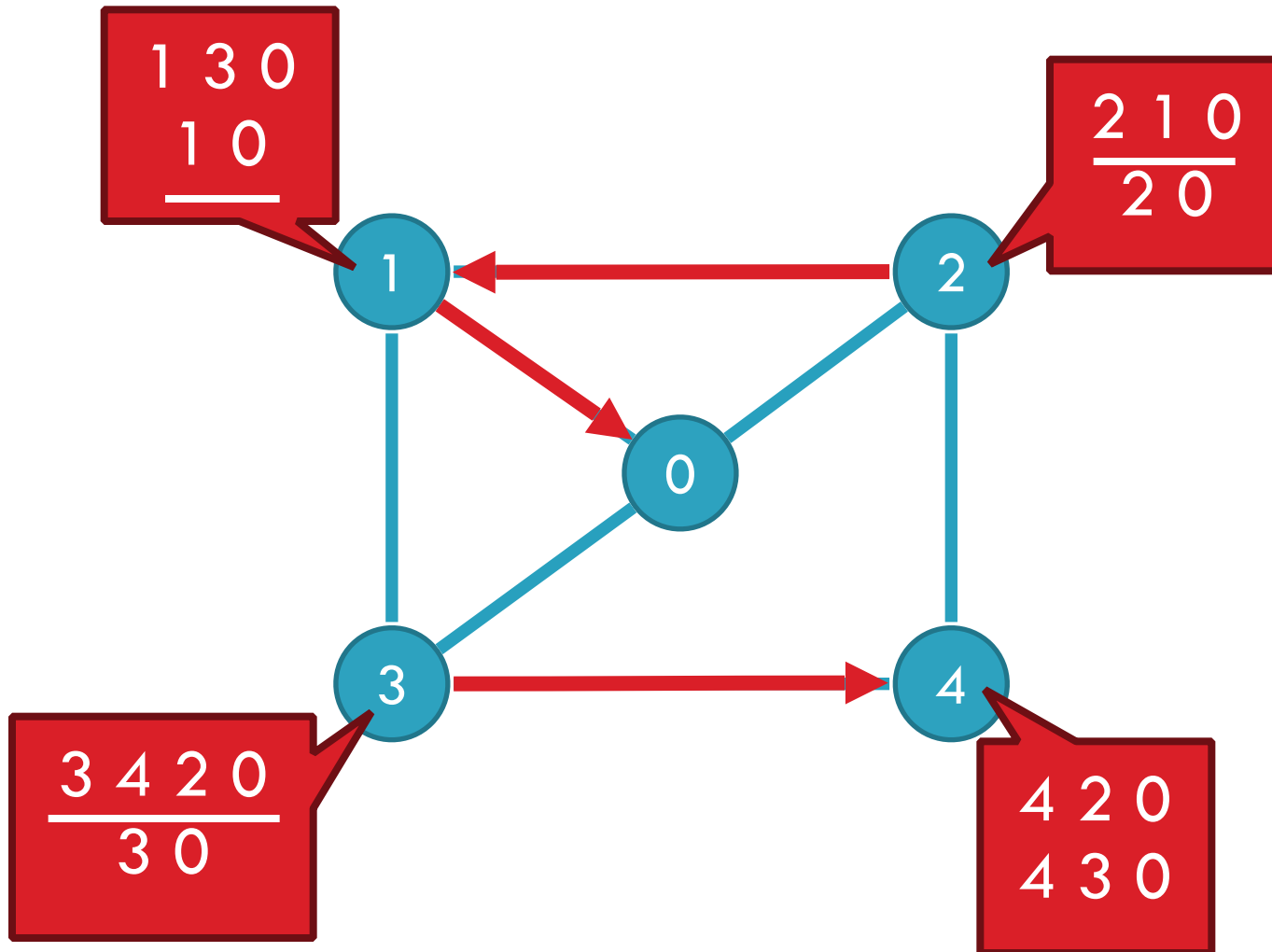
Bad Gadget

33



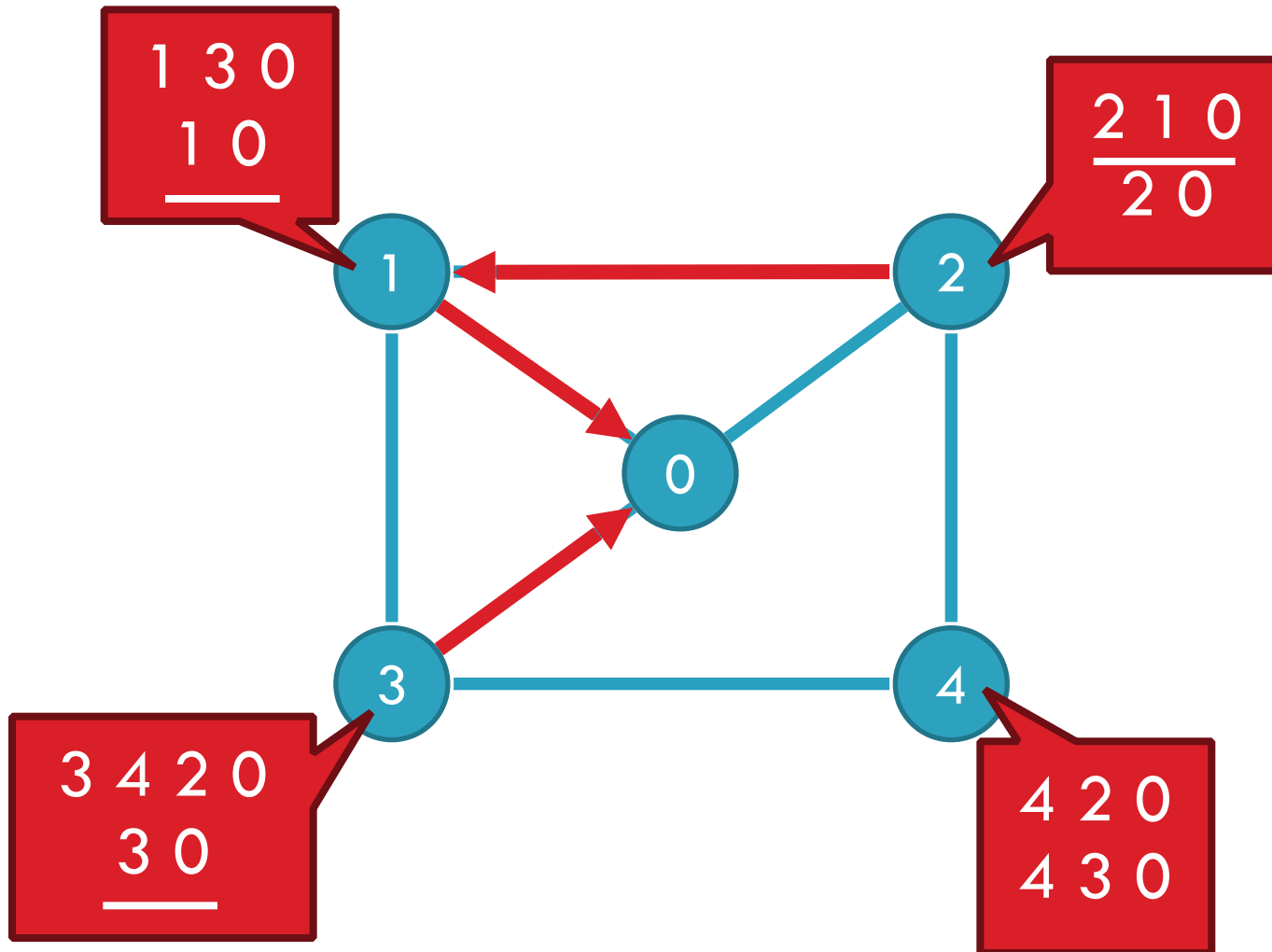
Bad Gadget

33



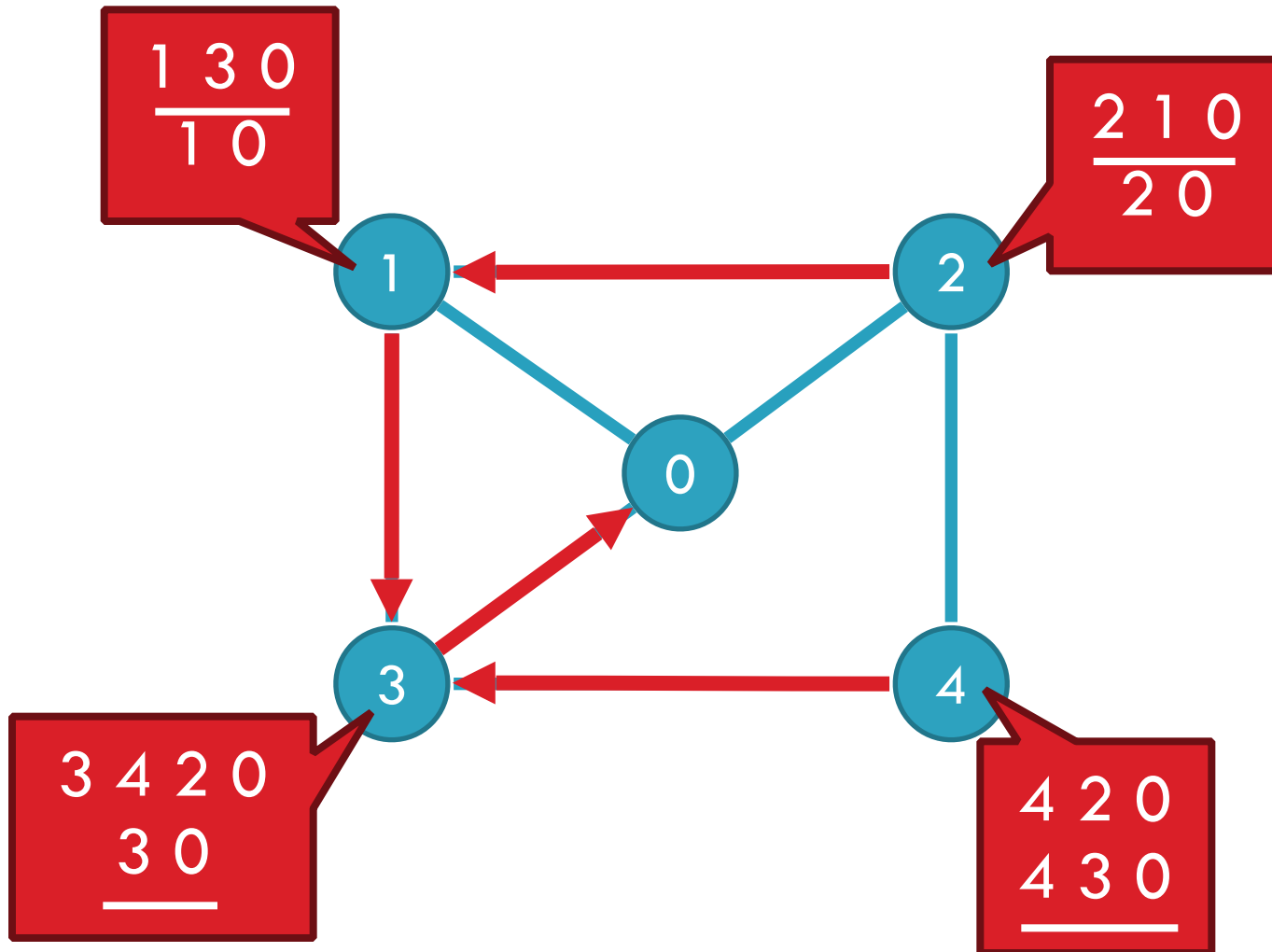
Bad Gadget

33



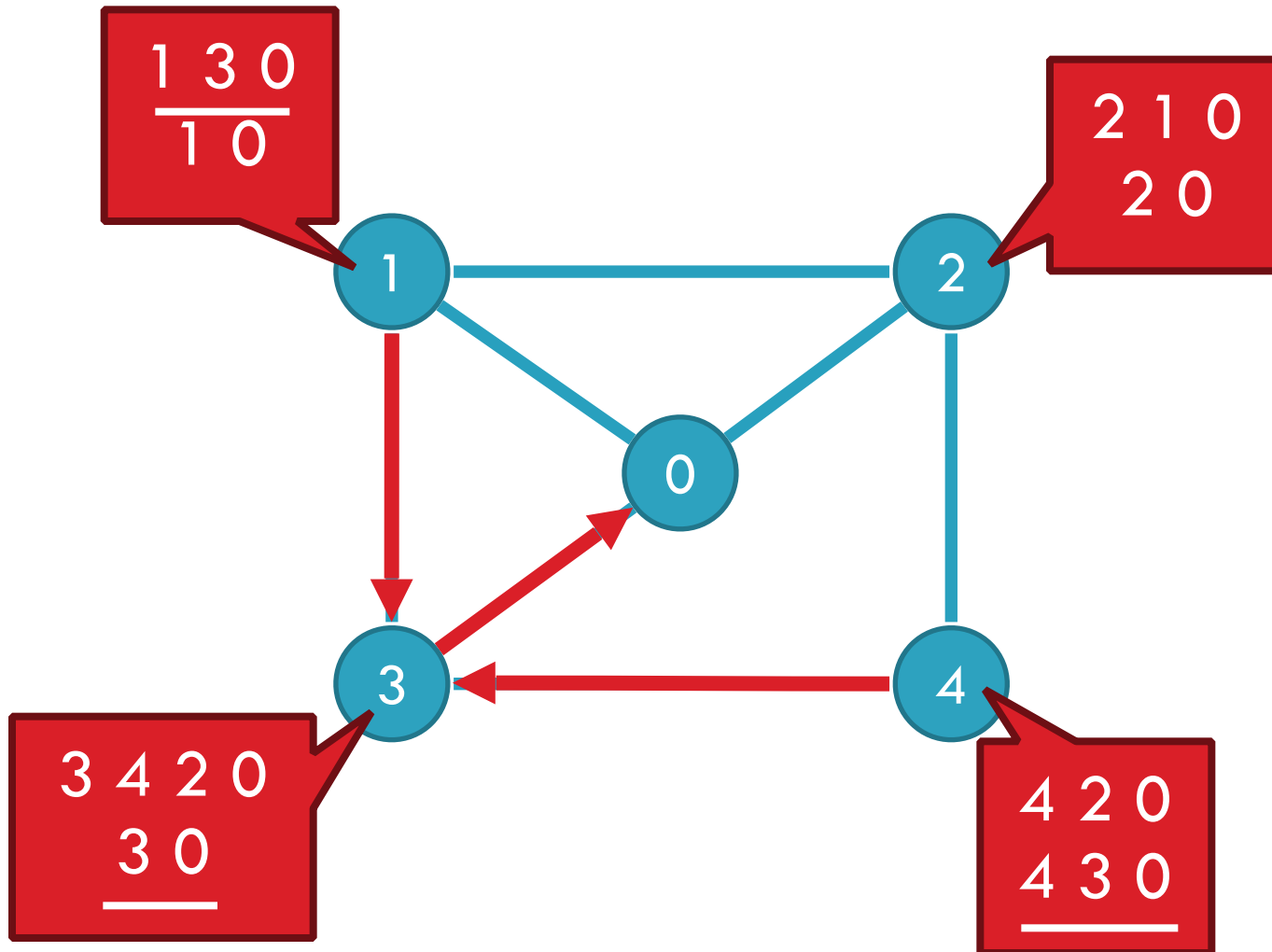
Bad Gadget

33



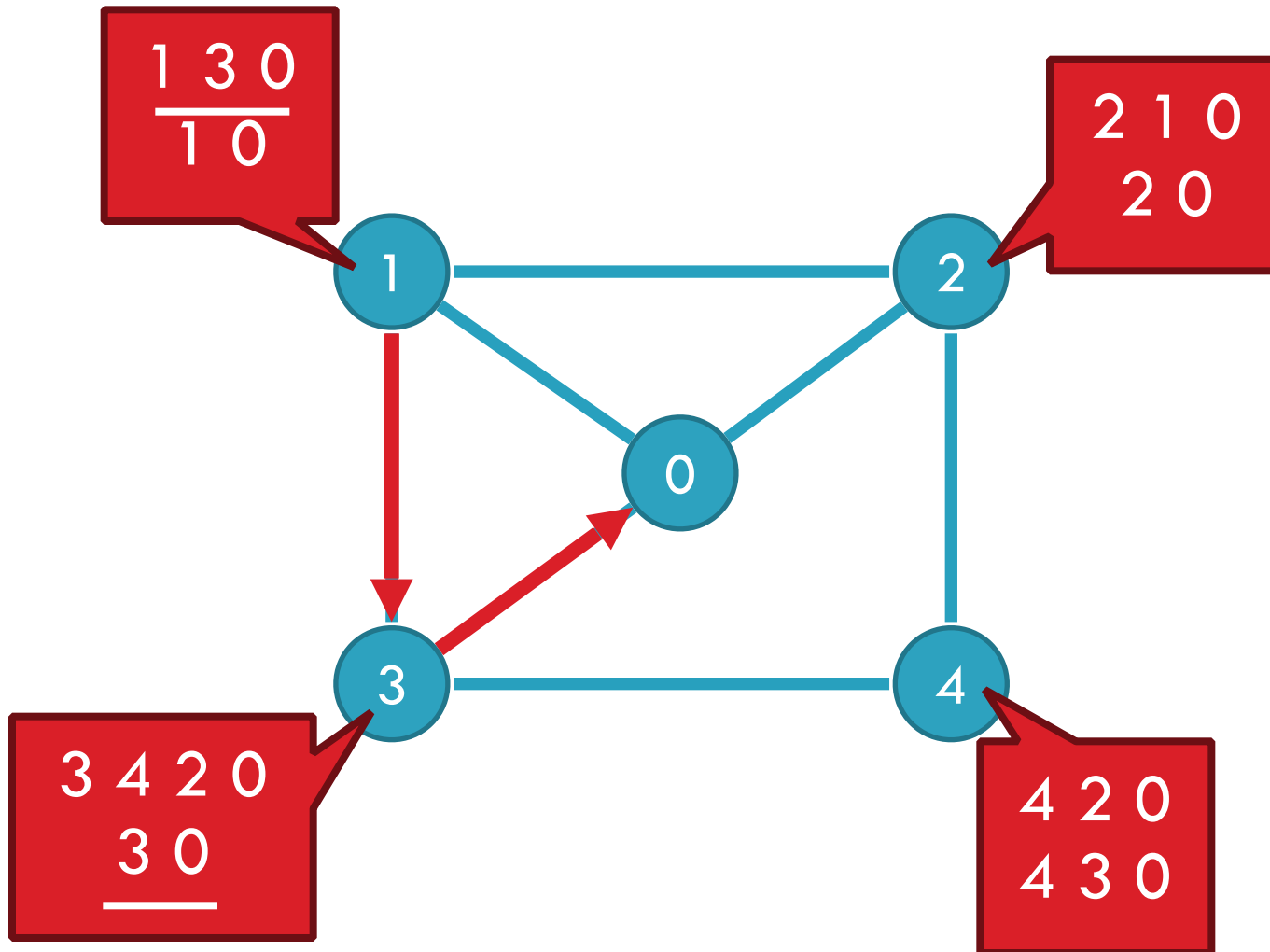
Bad Gadget

33



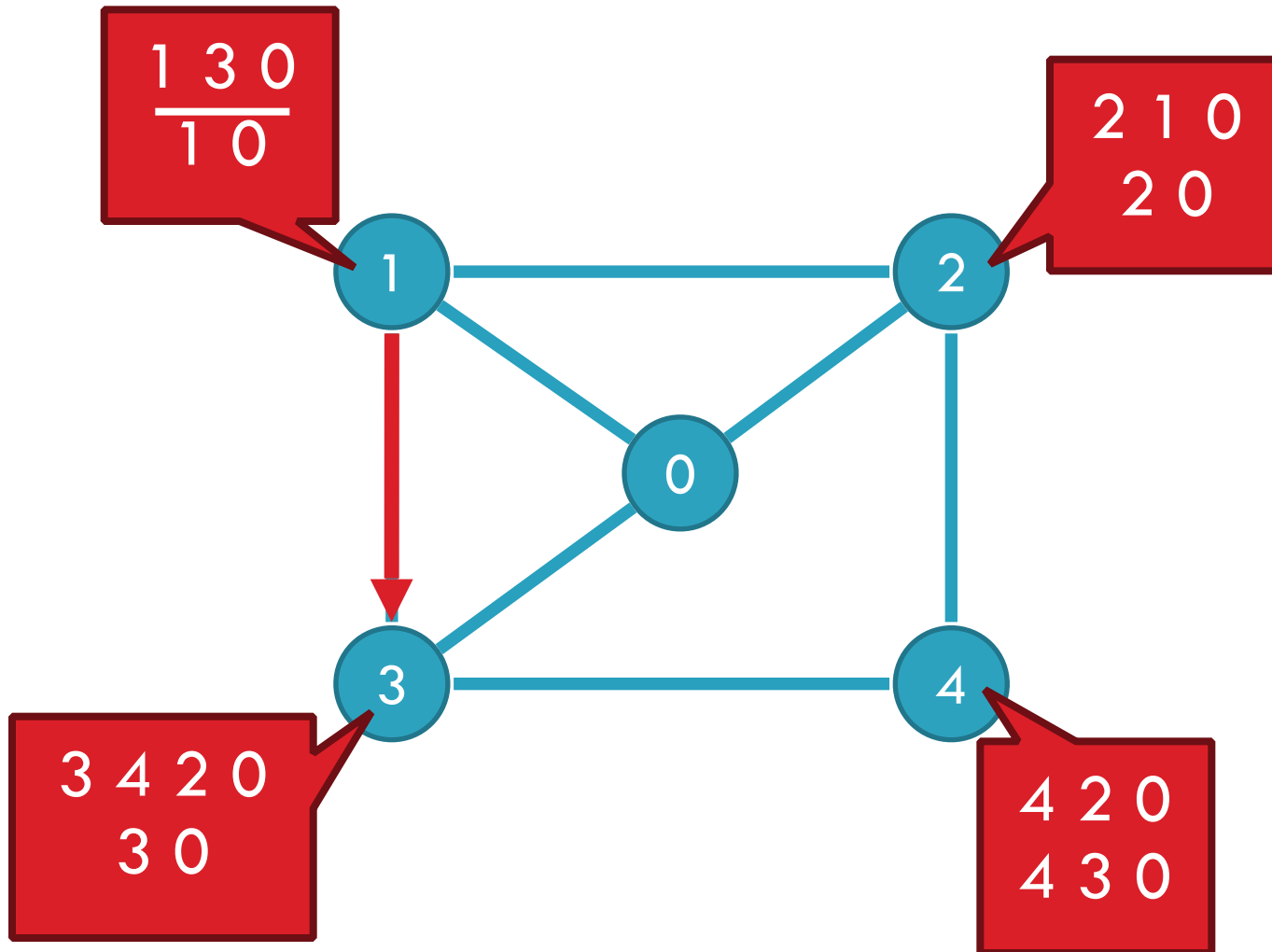
Bad Gadget

33



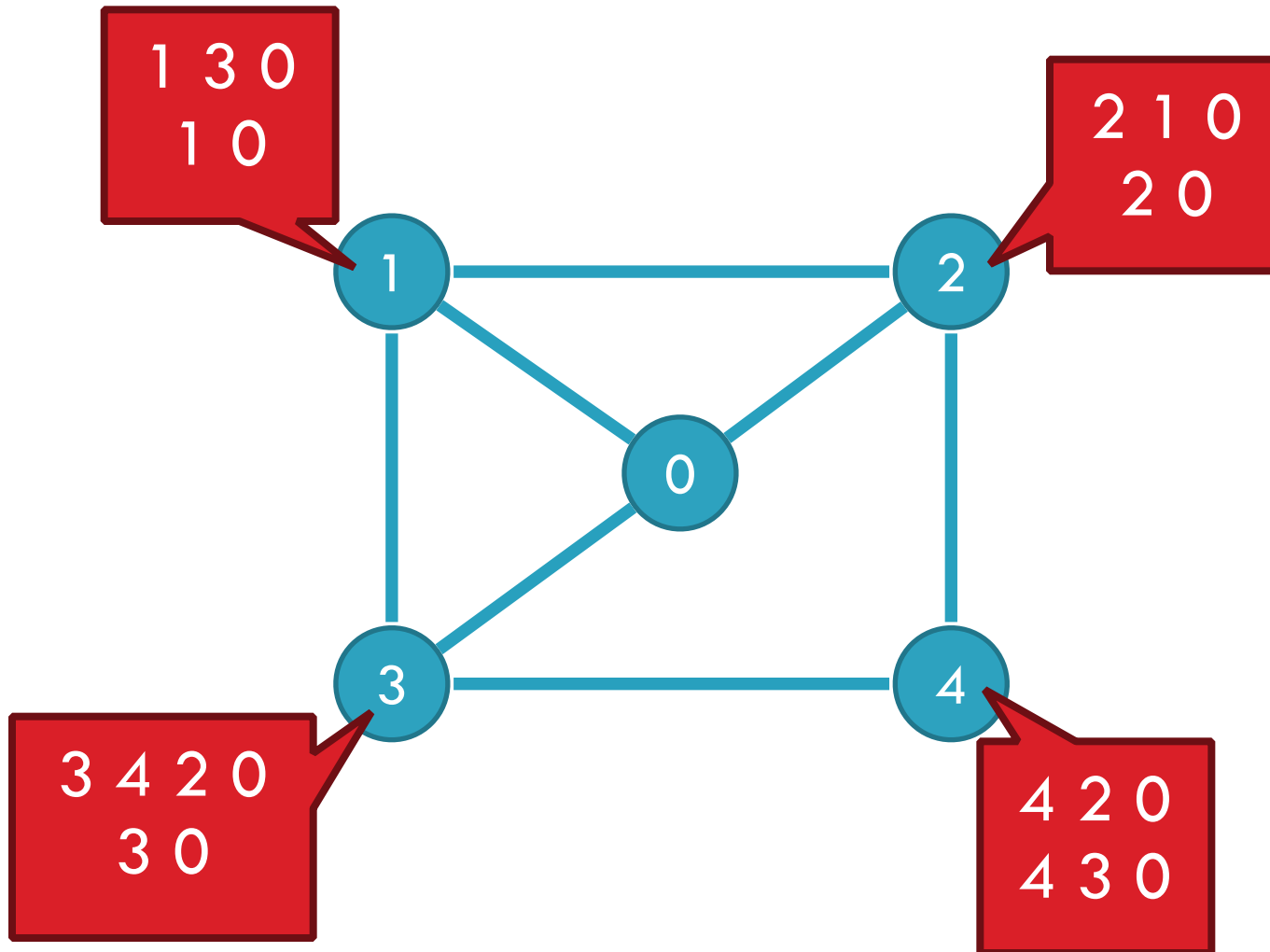
Bad Gadget

33



Bad Gadget

33



Bad Gadget

33

1 3 0

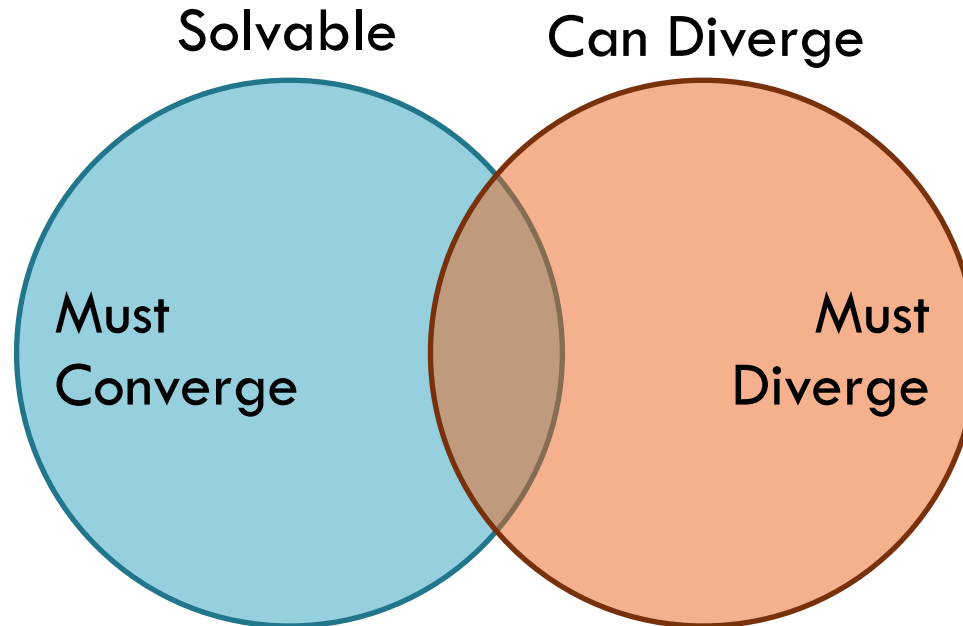
- That was only one round of oscillation!
- This keeps going, infinitely
- Problem stems from:
 - Local (not global) decisions
 - Ability of one node to improve its path selection

4 3 0

SPP Explains BGP Divergence

34

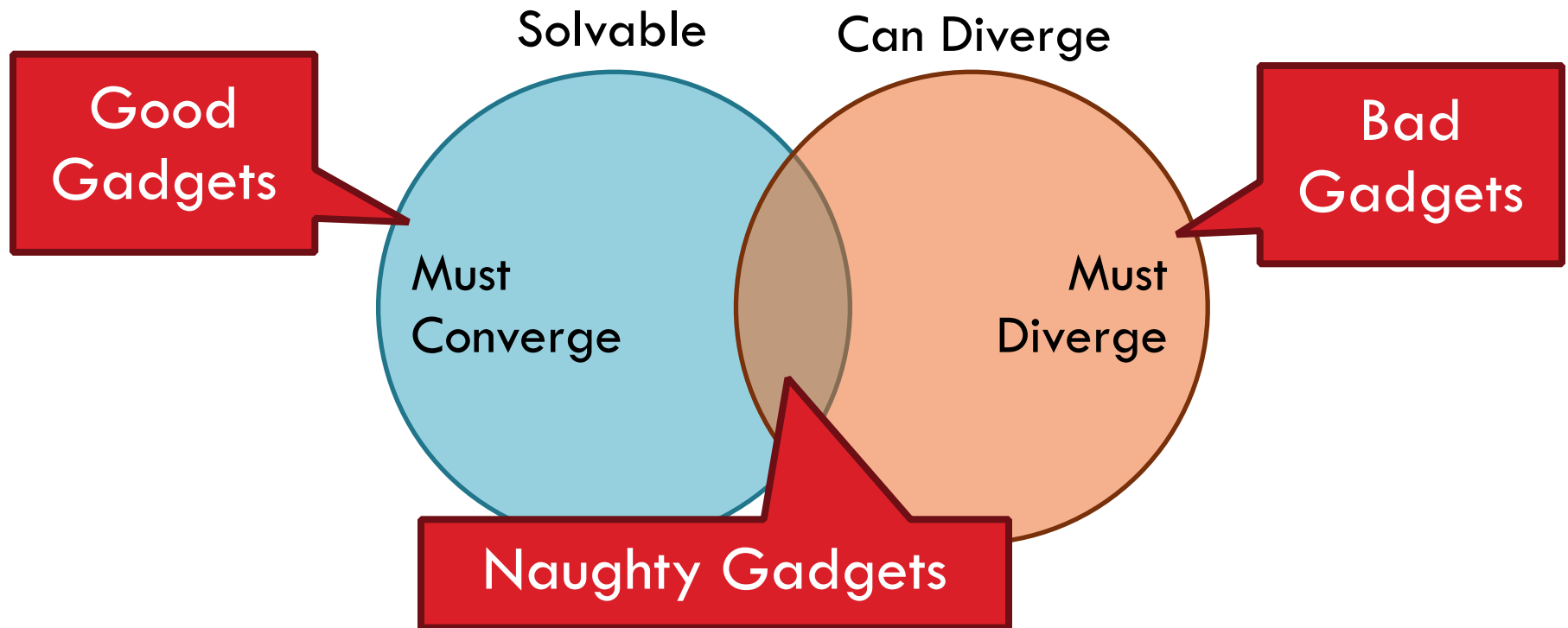
- BGP is **not** guaranteed to converge to stable routing
 - ▣ Policy inconsistencies may lead to “livelock”
 - ▣ Protocol oscillation



SPP Explains BGP Divergence

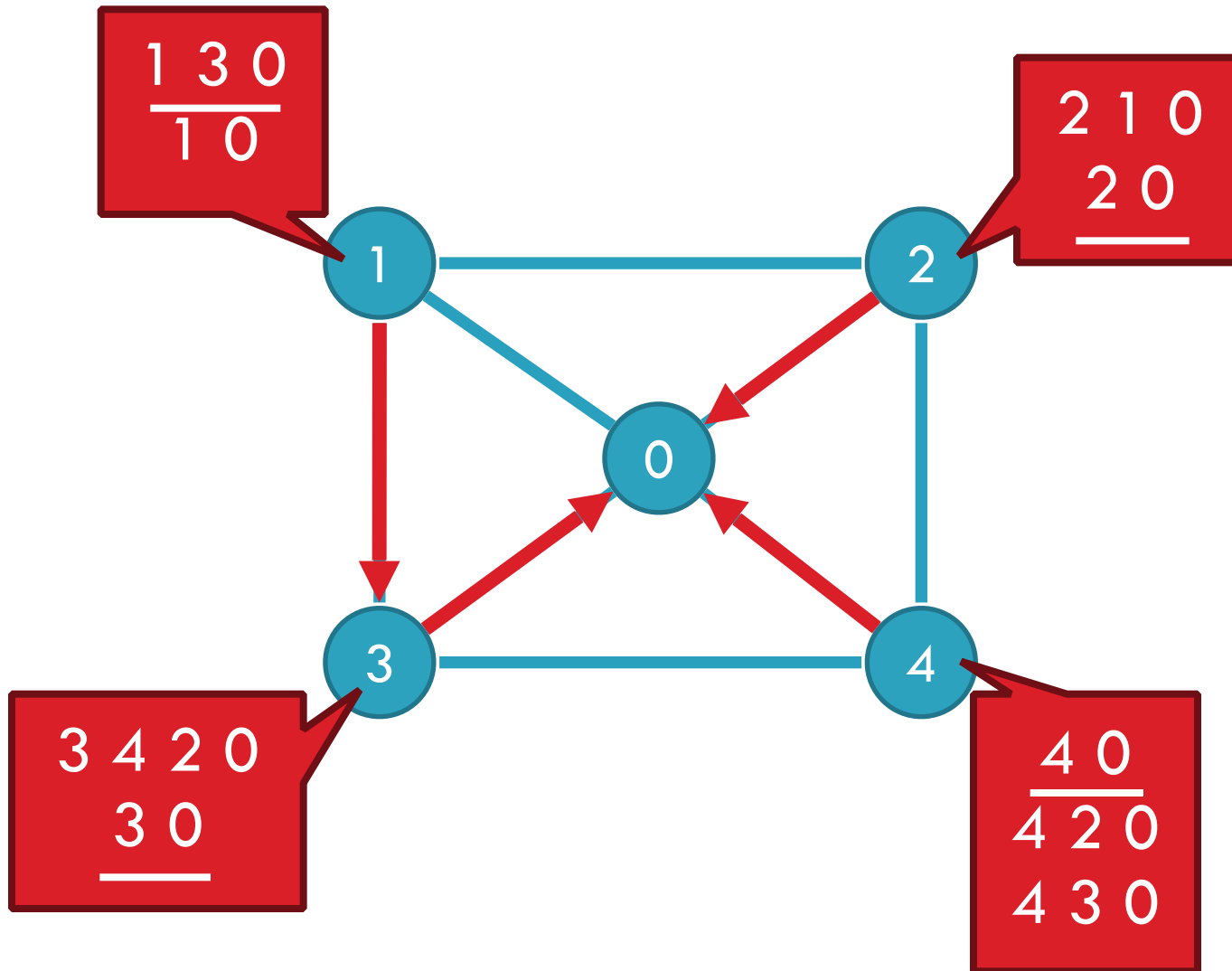
34

- BGP is **not** guaranteed to converge to stable routing
 - ▣ Policy inconsistencies may lead to “livelock”
 - ▣ Protocol oscillation



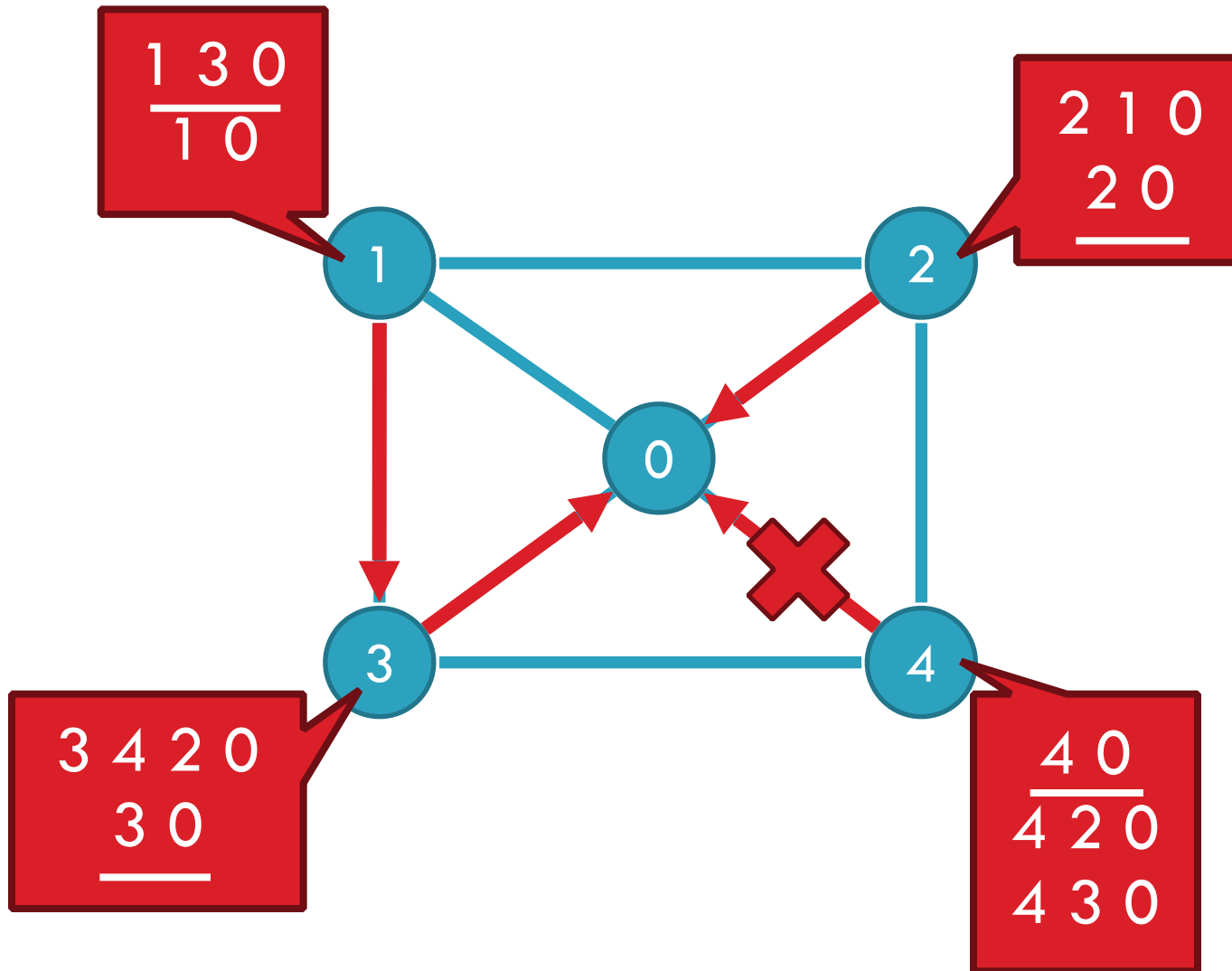
Beware of Backup Policies

35



Beware of Backup Policies

35

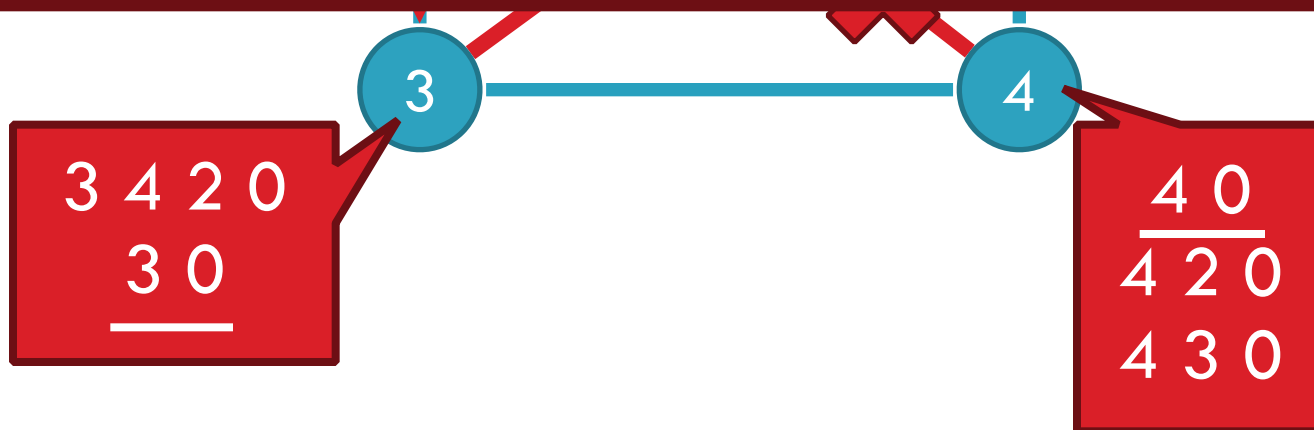


Beware of Backup Policies

35

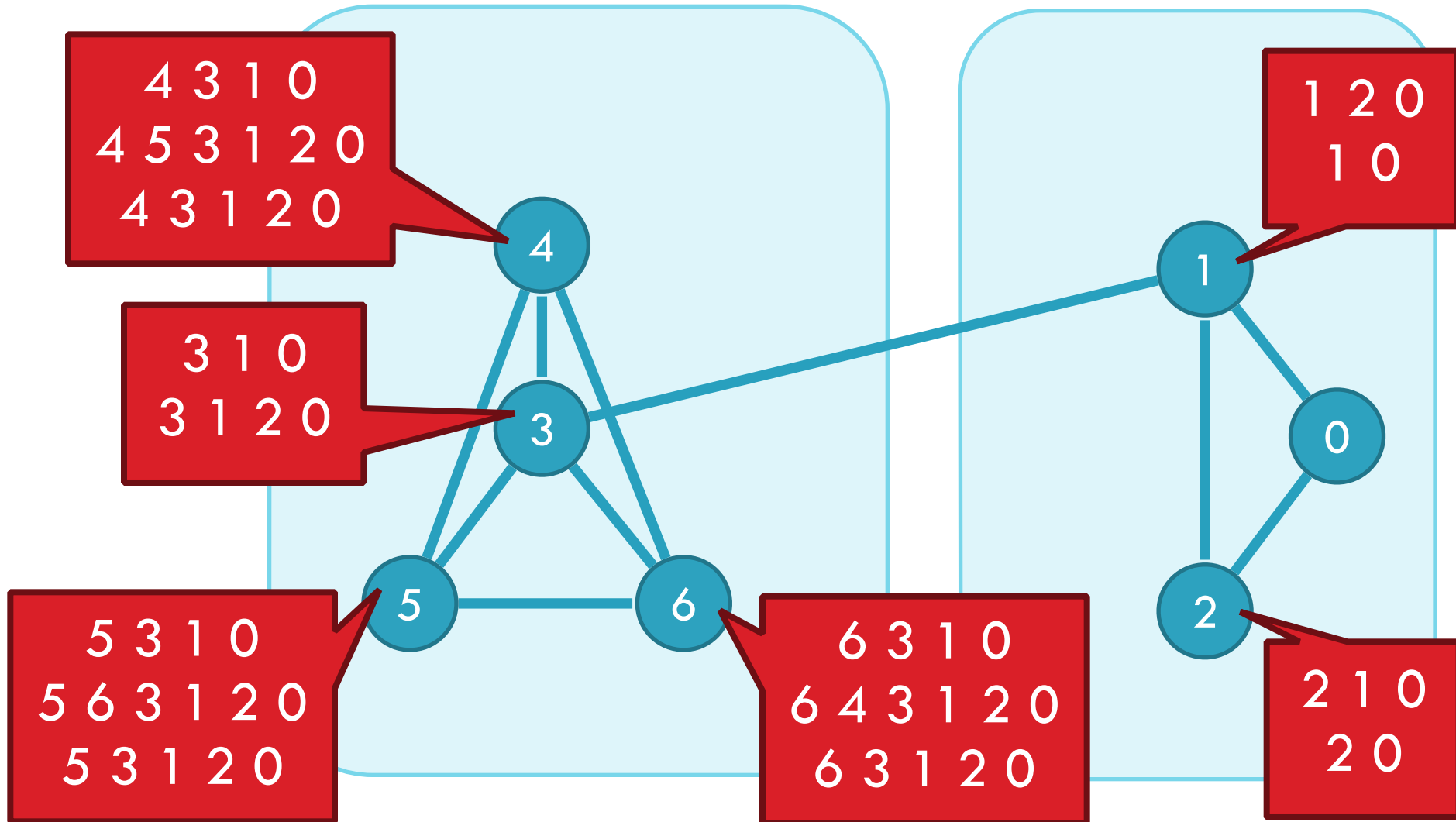


- BGP is not robust
- It may not recover from link failure



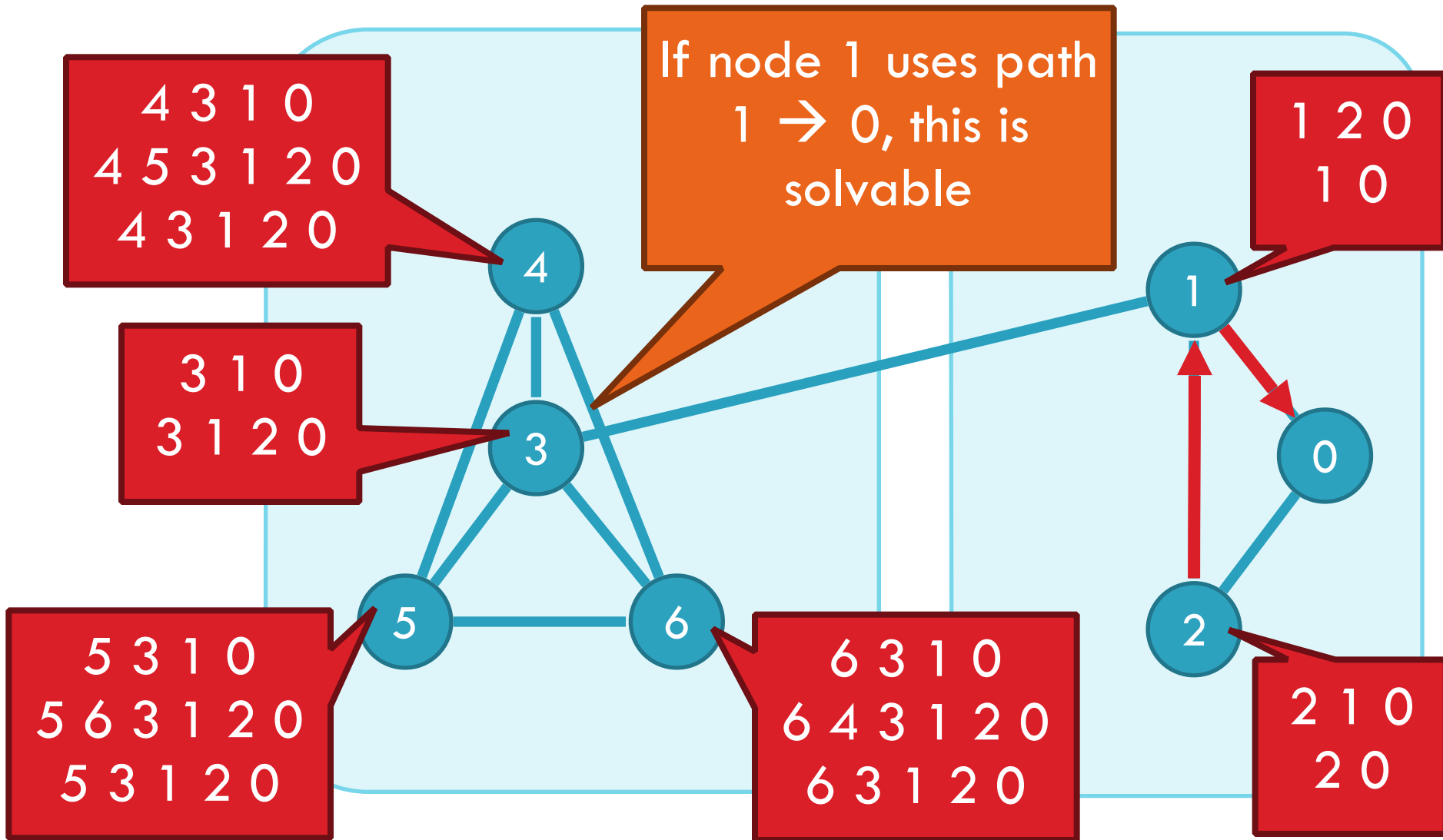
BGP is Precarious

36

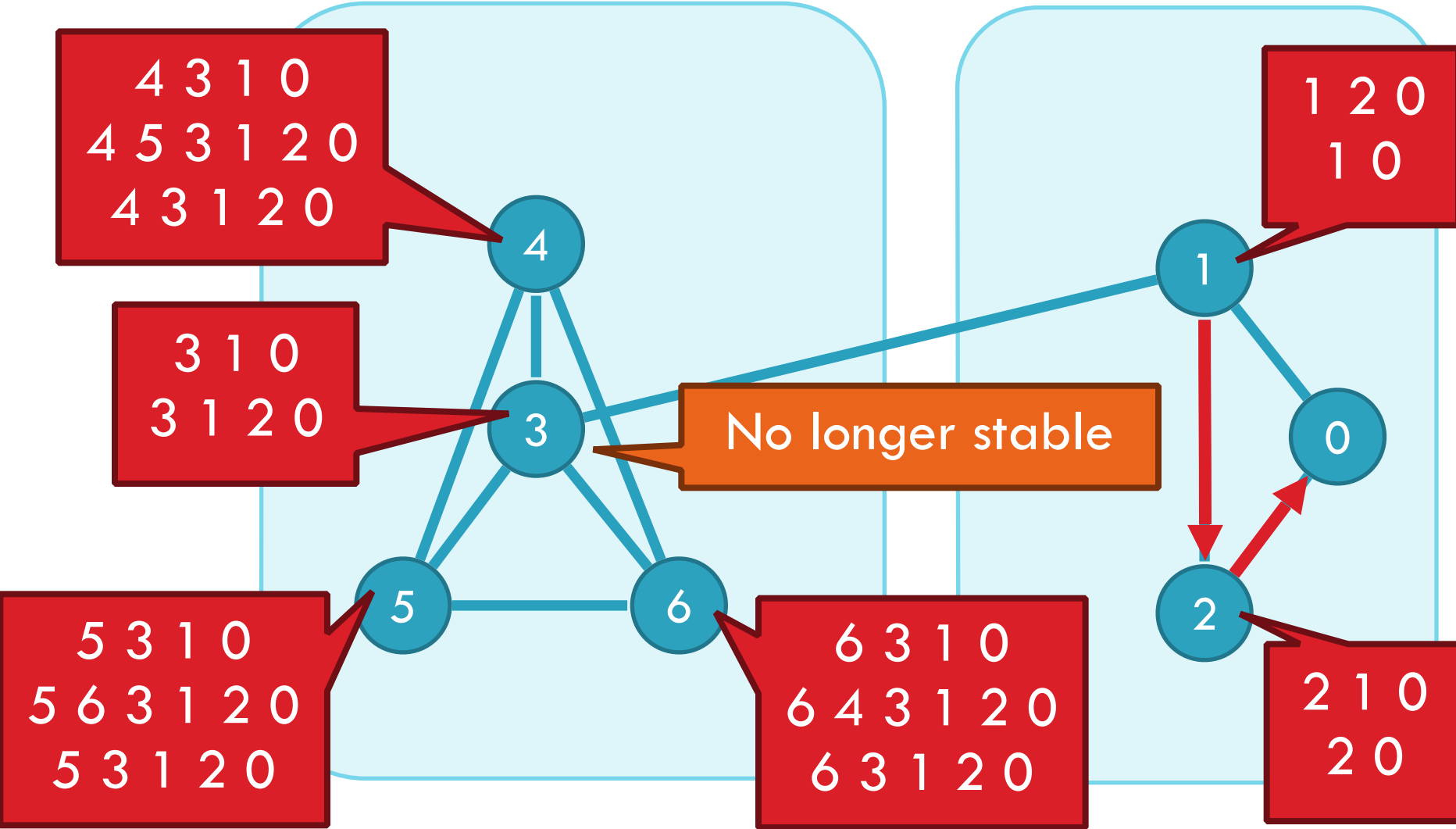


BGP is Precarious

36



BGP is Precarious



Can BGP Be Fixed?

- Unfortunately, SPP is NP-complete

Can BGP Be Fixed?

- Unfortunately, SPP is NP-complete

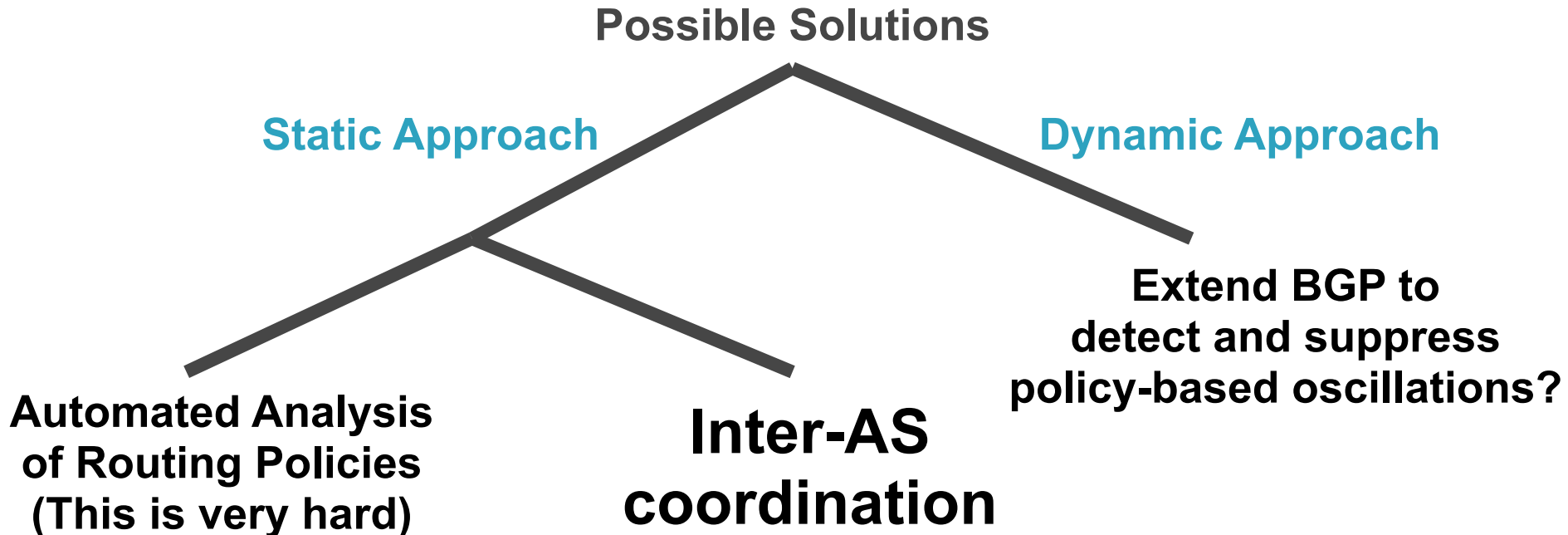
Possible Solutions

Dynamic Approach

Extend BGP to
detect and suppress
policy-based oscillations?

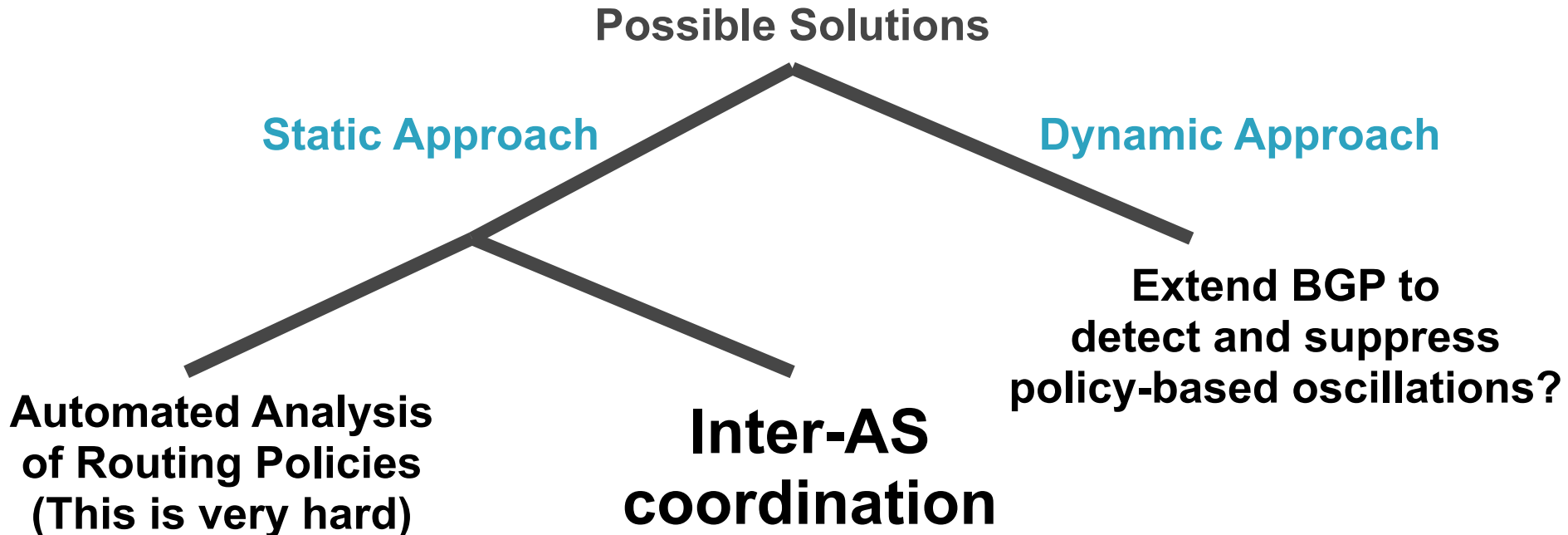
Can BGP Be Fixed?

- Unfortunately, SPP is NP-complete



Can BGP Be Fixed?

- Unfortunately, SPP is NP-complete



These approaches are complementary

- ❑ BGP Basics
- ❑ Stable Paths Problem
- ❑ BGP in the Real World

Motivation

39

- Routing reliability/fault-tolerance on small time scales (minutes) not previously a priority
- Transaction oriented and interactive applications (e.g. Internet Telephony) will require higher levels of end-to-end network reliability
- How well does the Internet routing infrastructure tolerate faults?

Conventional Wisdom

40

- Internet routing is robust under faults
 - ▣ Supports path re-routing
 - ▣ Path restoration on the order of seconds
- BGP has good convergence properties
 - ▣ Does not exhibit looping/bouncing problems of RIP
- Internet fail-over will improve with faster routers and faster links
- More redundant connections (multi-homing) will always improve fault-tolerance

Delayed Routing Convergence

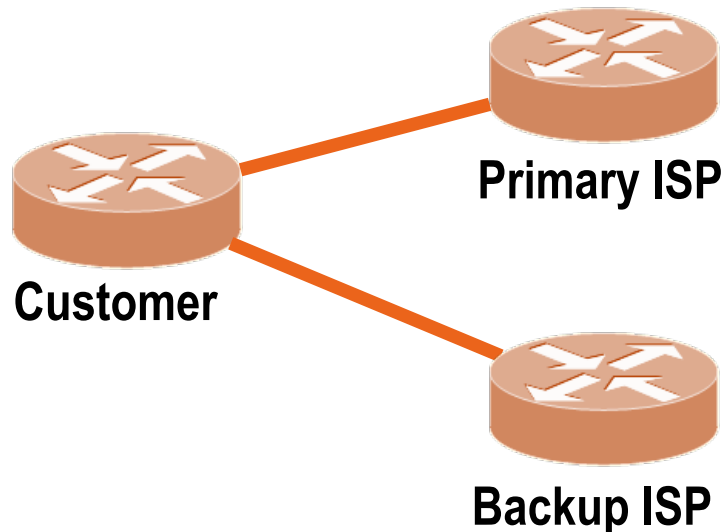
41

- Conventional wisdom about routing convergence is not accurate
 - Measurement of BGP convergence in the Internet
 - Analysis/intuition behind delayed BGP routing convergence
 - Modifications to BGP implementations which would improve convergence times

Open Question

42

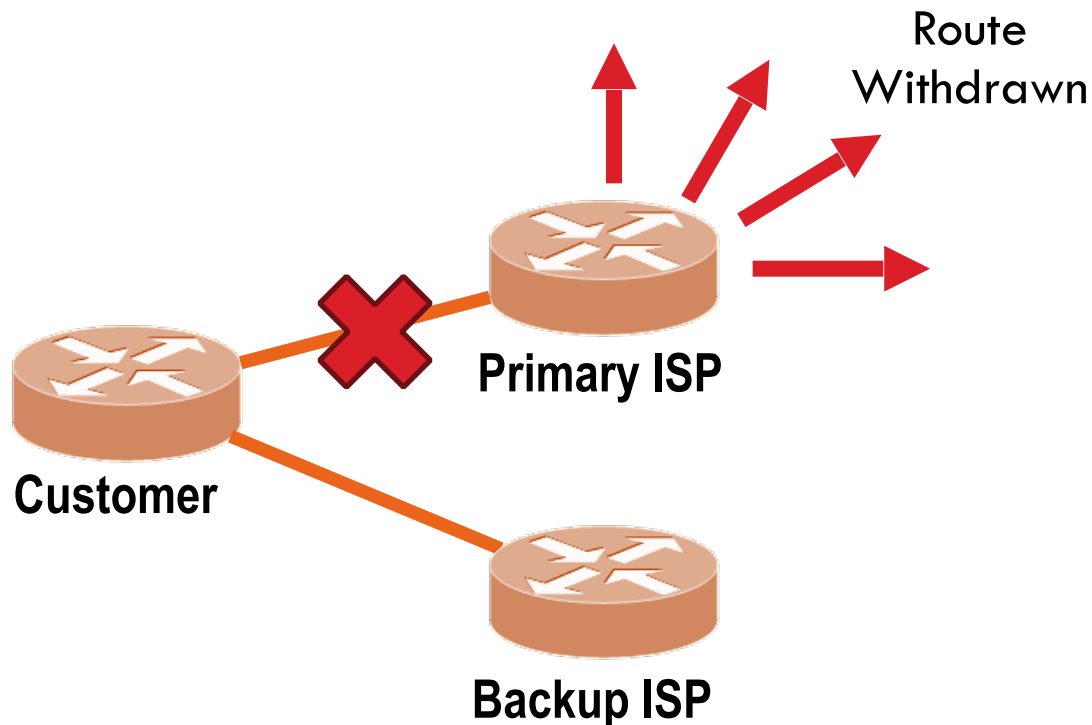
- After a fault in a path to multi-homed site, how long does it take for majority of Internet routers to fail-over to secondary path?



Open Question

42

- After a fault in a path to multi-homed site, how long does it take for majority of Internet routers to fail-over to secondary path?

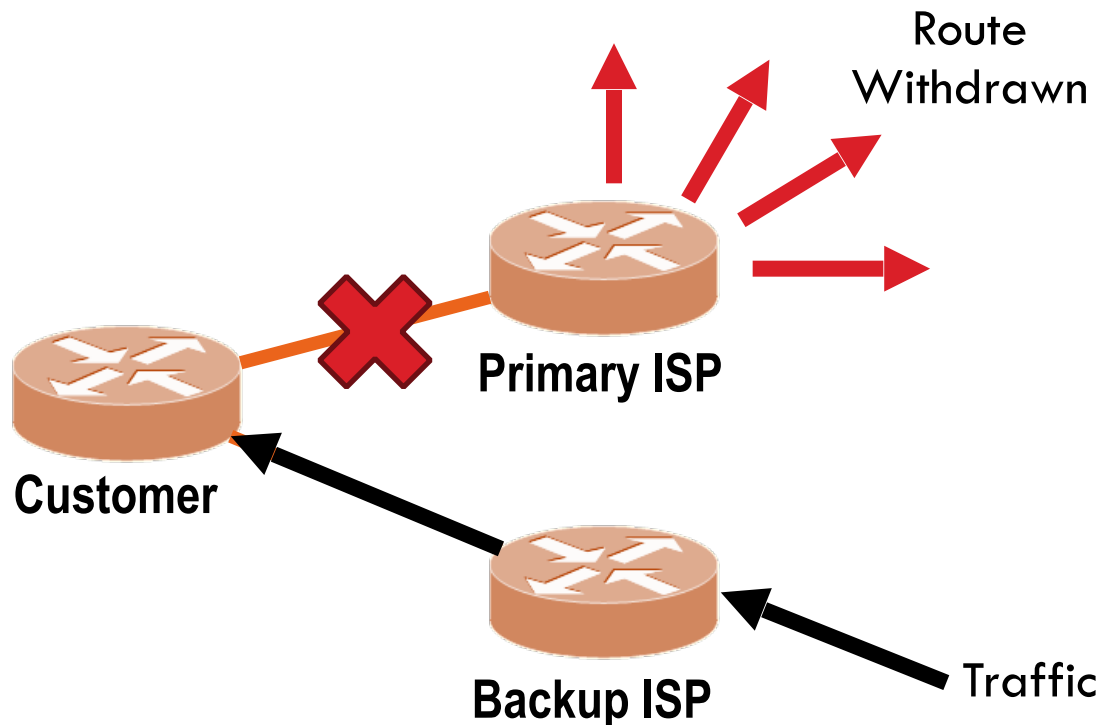


- Routing table convergence

Open Question

42

- After a fault in a path to multi-homed site, how long does it take for majority of Internet routers to fail-over to secondary path?



- Routing table convergence
- Stable end-to-end paths

Bad News

43

- With unconstrained policies:
 - Divergence
 - Possible create unsatisfiable policies
 - NP-complete to identify these policies
 - Happening today?

Bad News

43

- With unconstrained policies:
 - Divergence
 - Possible create unsatisfiable policies
 - NP-complete to identify these policies
 - Happening today?
- With constrained policies (e.g. shortest path first)
 - Transient oscillations
 - BGP usually converges
 - It may take a very long time...

Bad News

43

- With unconstrained policies:
 - Divergence
 - Possible create unsatisfiable policies
 - NP-complete to identify these policies
 - Happening today?
- With constrained policies (e.g. shortest path first)
 - Transient oscillations
 - BGP usually converges
 - It may take a very long time...
- BGP Beacons: focuses on constrained policies

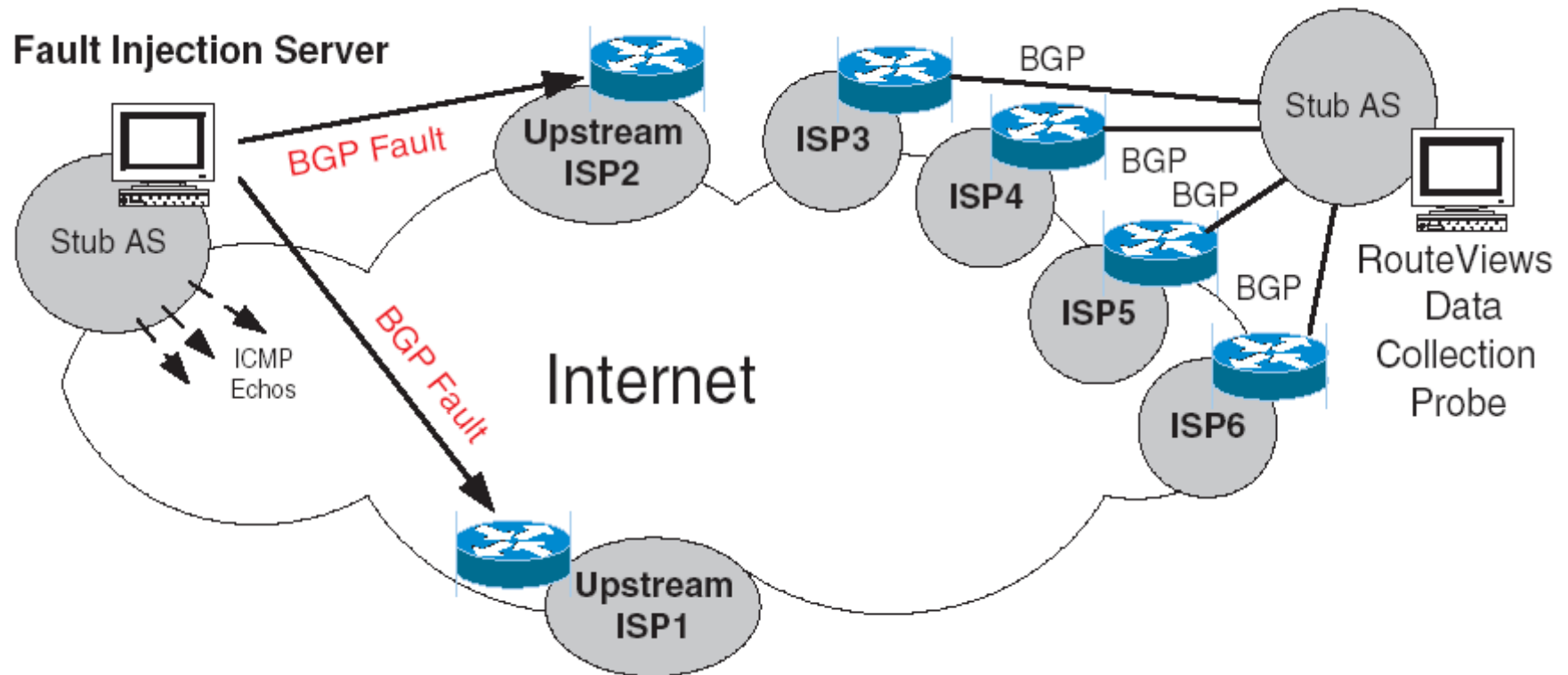
16 Month Study of Convergence

44

- Instrument the Internet
 - Inject BGP faults (announcements/withdrawals) of varied prefix and AS path length into topologically and geographically diverse ISP peering sessions
 - Monitor impact faults through
 - Recording BGP peering sessions with 20 tier1/tier2 ISPs
 - Active ICMP measurements (512 byte/second to 100 random web sites)
 - Wait two years (and 250,000 faults)

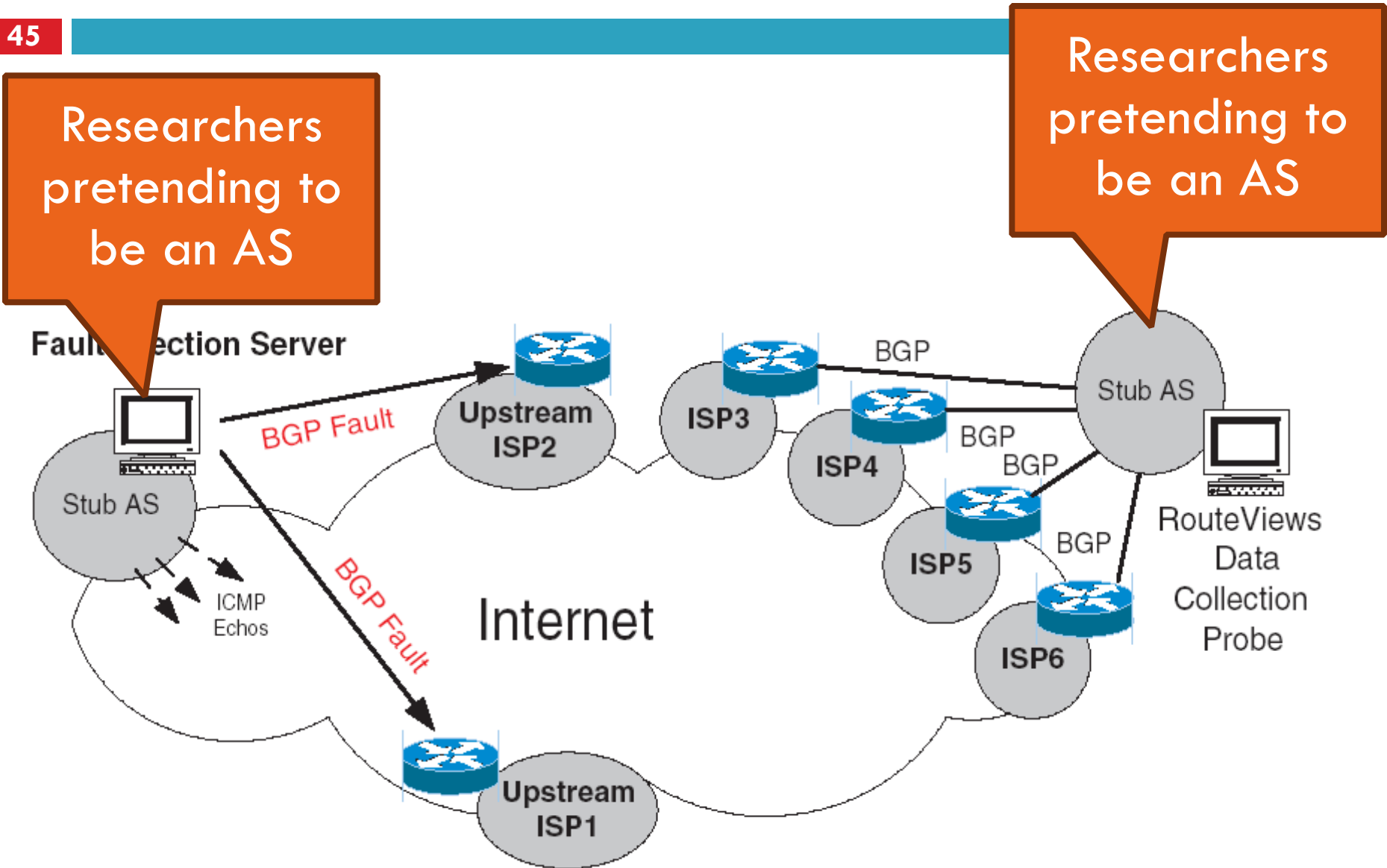
Measurement Architecture

45



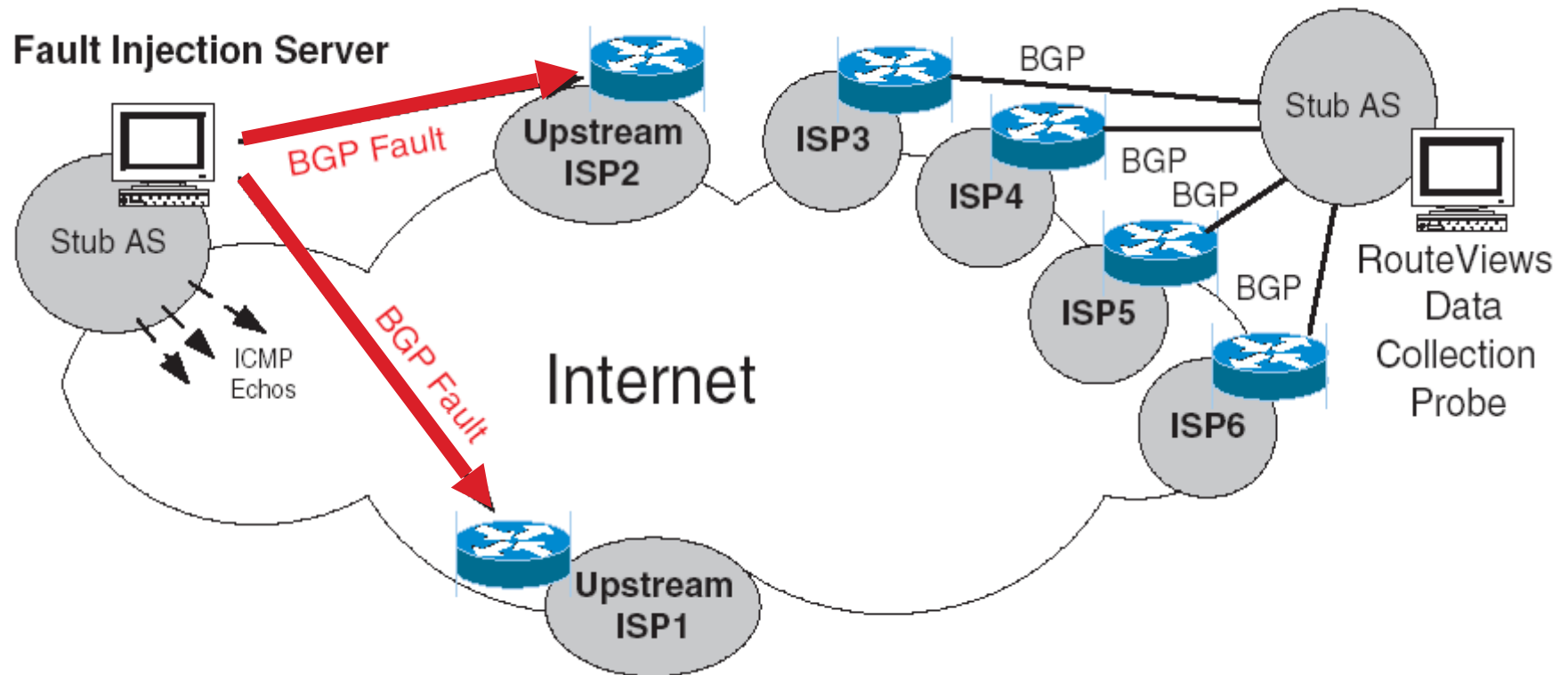
Measurement Architecture

45



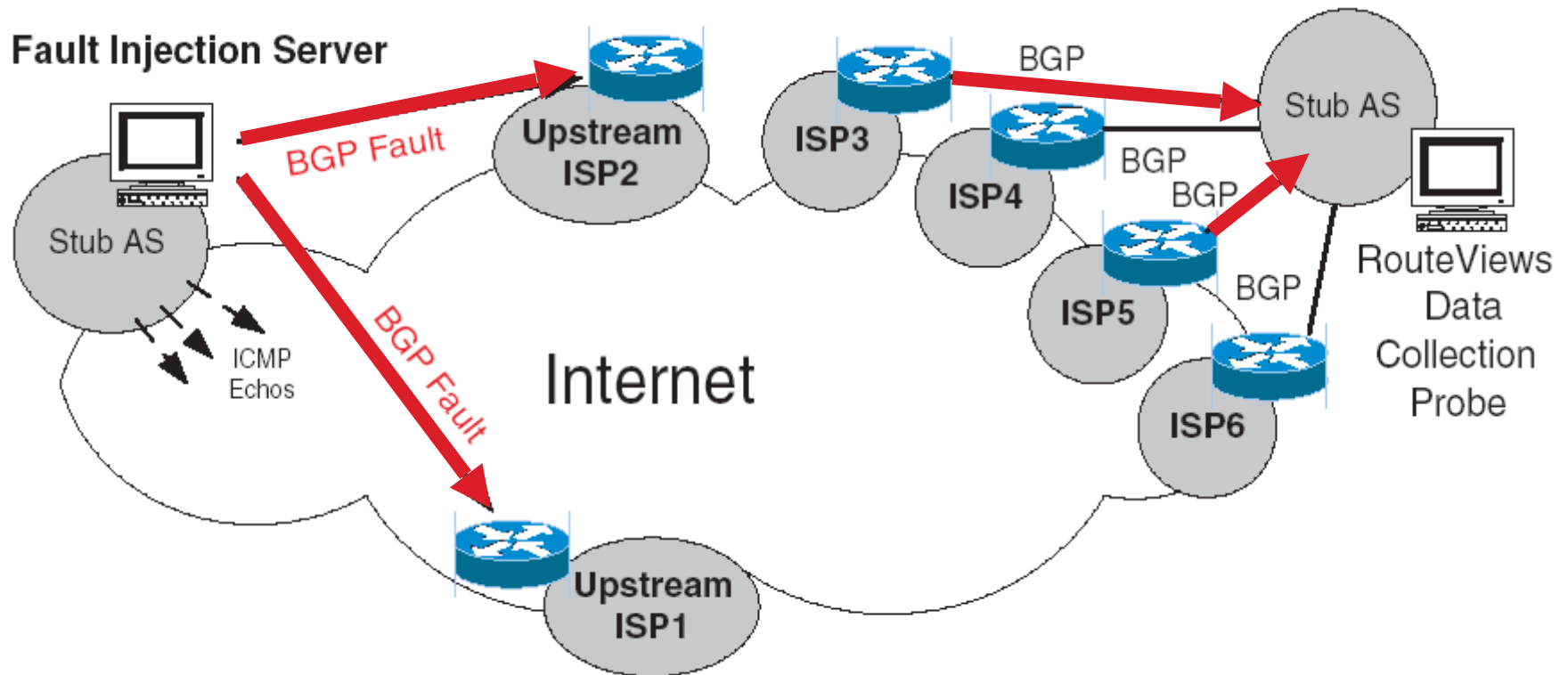
Measurement Architecture

45



Measurement Architecture

45



Announcement Scenarios

46

- T_{up} – a new route is advertised
- T_{down} – A route is withdrawn
 - i.e. single-homed failure
- T_{short} – Advertise a shorter/better AS path
 - i.e. primary path repaired
- T_{long} – Advertise a longer/worse AS path
 - i.e. primary path fails

Major Convergence Results

47

- Routing convergence requires an order of magnitude longer than expected
 - ▣ 10s of minutes
- Routes converge more quickly following T_{up}/Repair than T_{down}/Failure events
 - ▣ Bad news travels more slowly
- Withdrawals (T_{down}) generate several more announcements than new routes (T_{up})

Example

48

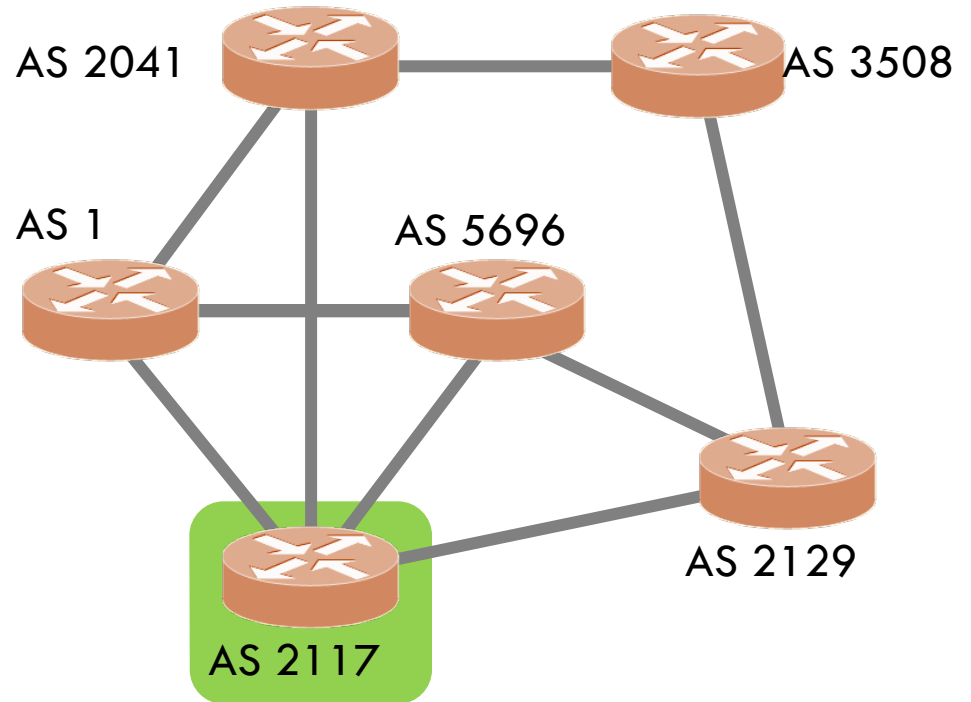
<u>TIME</u>	<u>BGP Message/Event</u>
10:40:30	<i>Route Fails/Withdrawn by AS2129</i>
10:41:08	<i>2117 announce 5696 2129</i>
10:41:32	<i>2117 announce 1 5696 2129</i>
10:41:50	<i>2117 announce 2041 3508 3508 4540 7037 1239 5696 2129</i>
10:42:17	<i>2117 announce 1 2041 3508 3508 4540 7037 1239 5696 2129</i>
10:43:05	<i>2117announce 2041 3508 3508 4540 7037 1239 6113 5696 2129</i>
10:43:35	<i>2117 announce 1 2041 3508 3508 4540 7037 1239 6113 5696 2129</i>
10:43:59	<i>2117 sends withdraw</i>

- ❑ BGP log of updates from AS2117 for route via AS2129
- ❑ One withdrawal triggers 6 announcements and one withdrawal from 2117
- ❑ Increasing AS path length until final withdrawal

Why So Many Announcements?

49

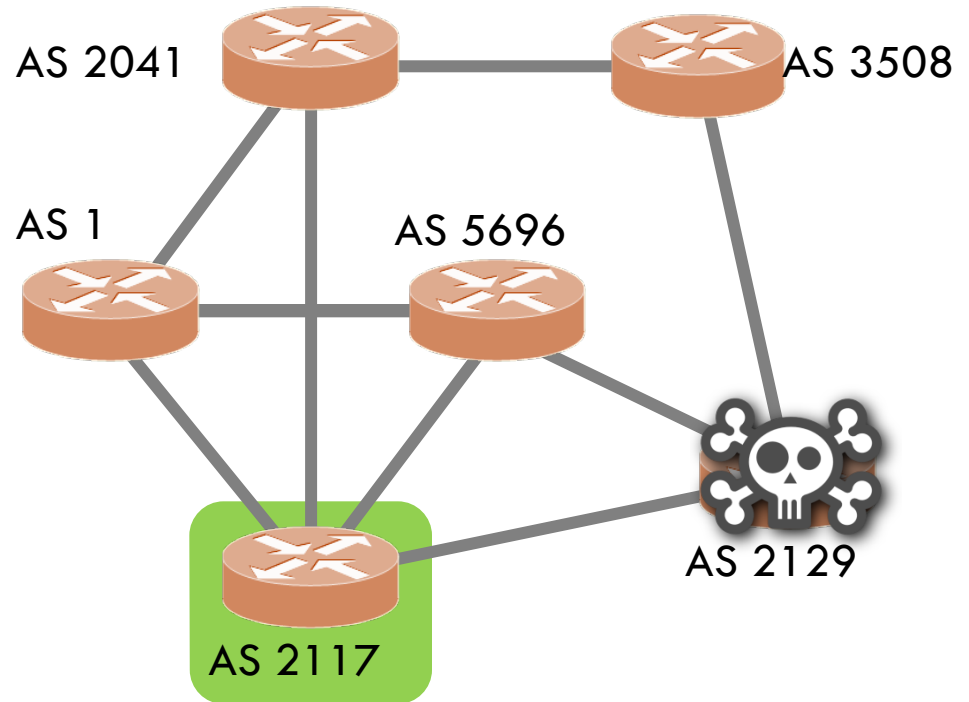
Events from AS 2177



Why So Many Announcements?

49

Events from AS 2177

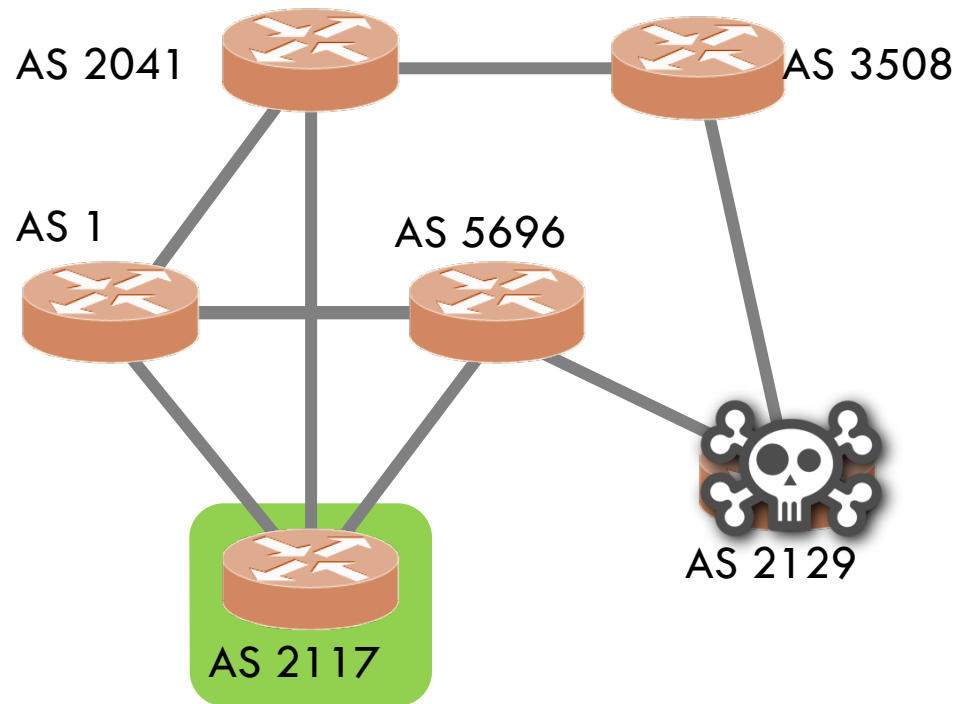


Why So Many Announcements?

49

Events from AS 2177

1. **Route Fails: AS 2129**

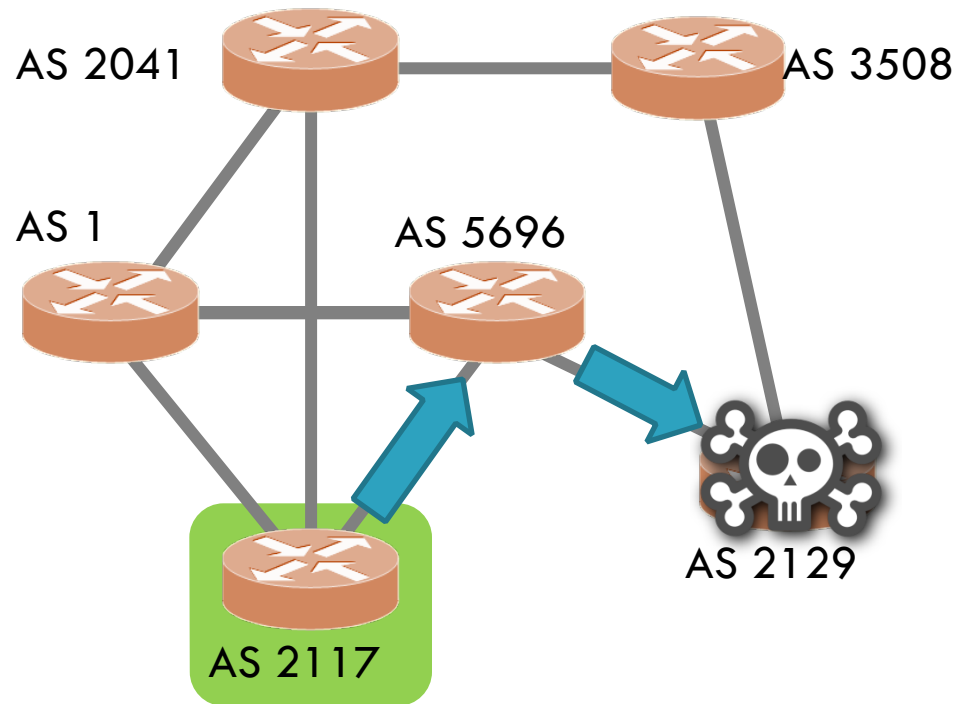


Why So Many Announcements?

49

Events from AS 2177

1. **Route Fails: AS 2129**
2. **Announce: 5696 2129**

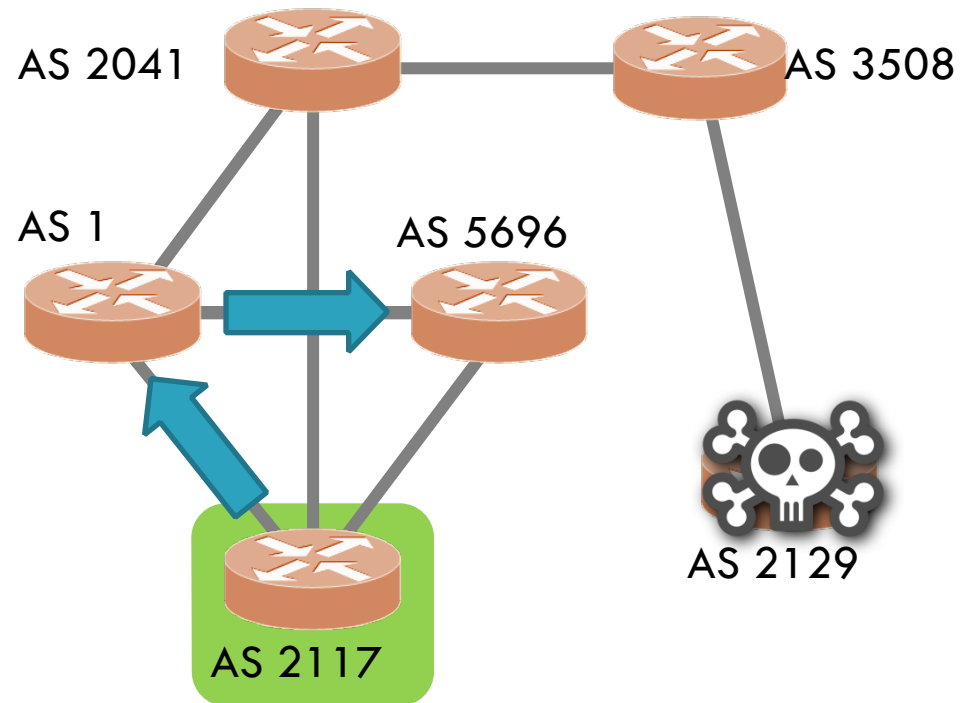


Why So Many Announcements?

49

Events from AS 2177

1. **Route Fails: AS 2129**
2. **Announce: 5696 2129**
3. **Announce: 1 5696 2129**

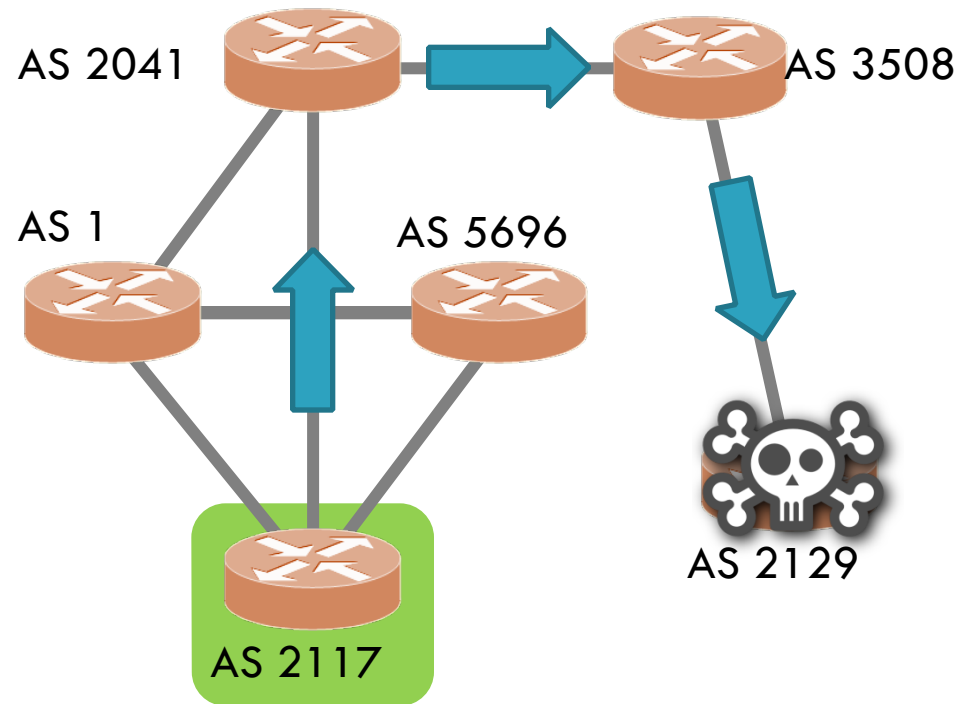


Why So Many Announcements?

49

Events from AS 2177

1. **Route Fails: AS 2129**
2. **Announce: 5696 2129**
3. **Announce: 1 5696 2129**
4. **Announce: 2041 3508 2129**

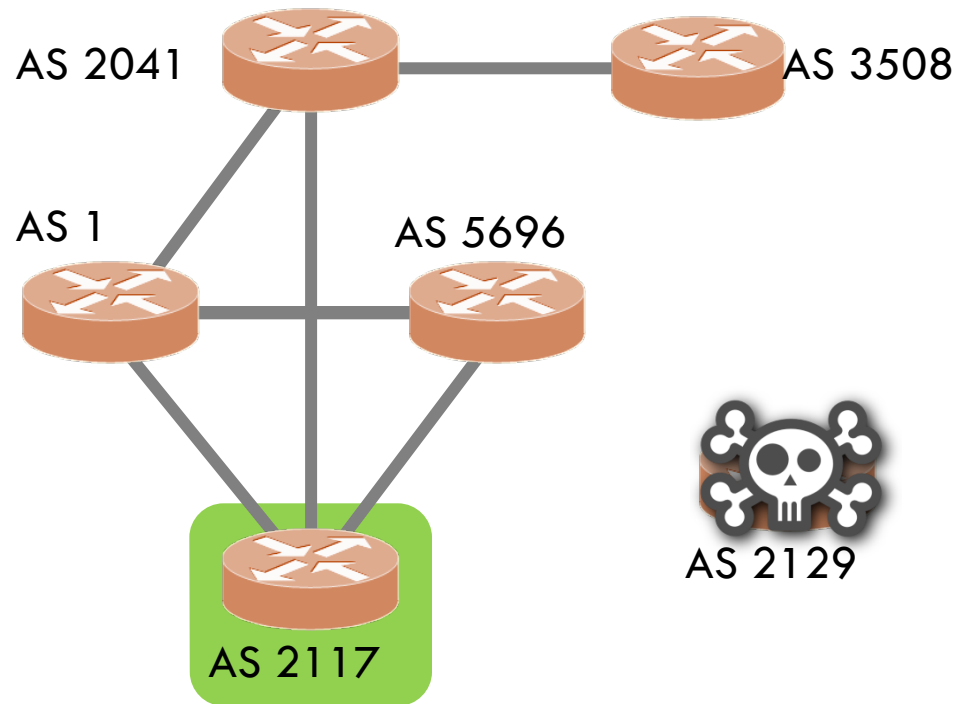


Why So Many Announcements?

49

Events from AS 2177

1. **Route Fails: AS 2129**
2. **Announce: 5696 2129**
3. **Announce: 1 5696 2129**
4. **Announce: 2041 3508 2129**
5. **Announce: 1 2041 3508 2129**
6. **Route Withdrawn: 2129**



How Many Announcements Does it Take For an AS to Withdraw a Route?

50

```
7/5 19:33:25      Route R is withdrawn
7/5 19:34:15      AS6543 anno unce R 6543 66665 8918 1 5696 999
7/5 19:35:00      AS6543 anno unce R 6543 66665 8918 67455 6461 5696 999
7/5 19:35:37      AS6543 anno unce R 6543 66665 4332 6461 5696 999
7/5 19:35:39      AS6543 anno unce R 6543 66665 5378 6660 67455 6461 5696 999
7/5 19:35:39      AS6543 anno unce R 6543 66665 65 6461 5696 999
7/5 19:35:52      AS6543 anno unce R 6543 66665 6461 5696 999
7/5 19:36:00      AS6543 anno unce R 6543 66665 5378 6765 6660 67455 6461 5696 999
...
7/5 19:38:22      AS6543 withdraw R
```

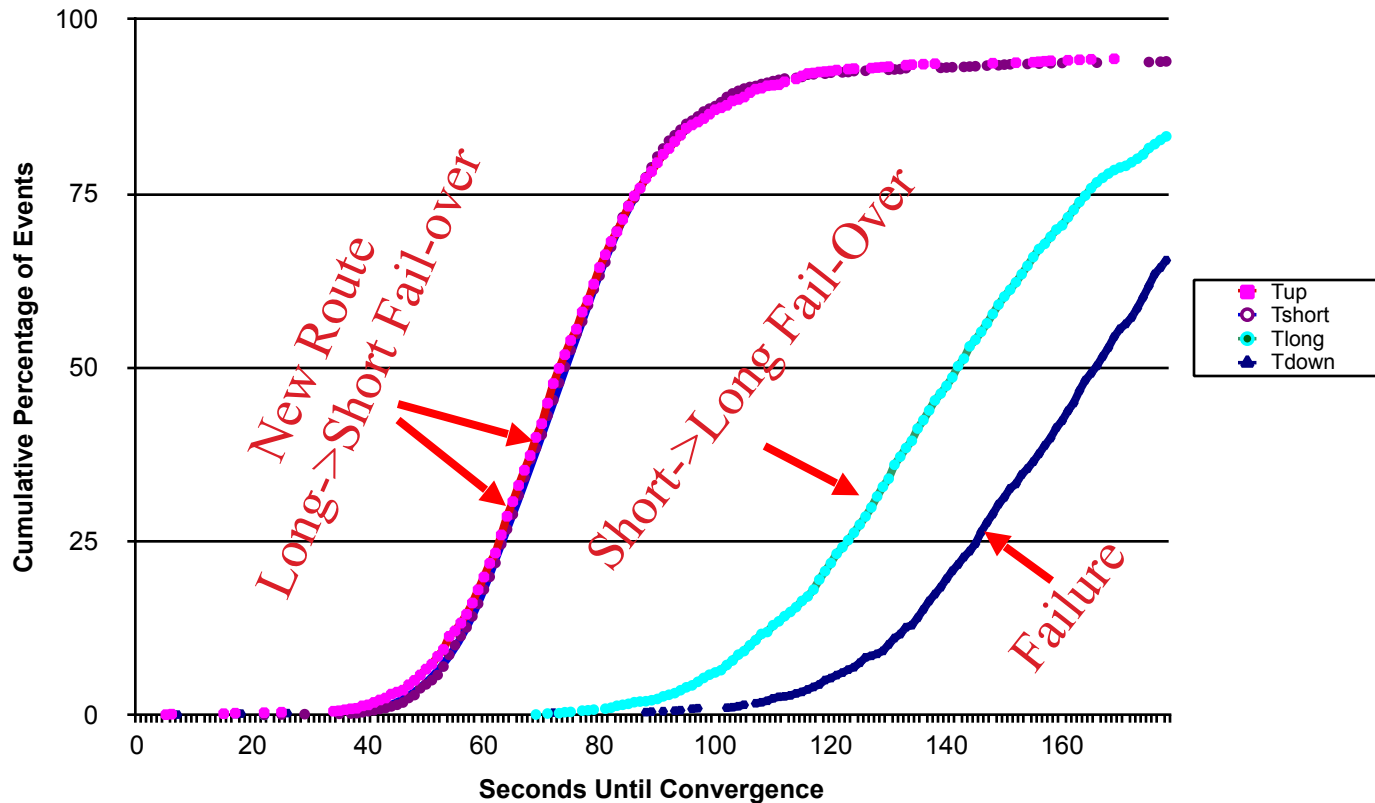
How Many Announcements Does it Take For an AS to Withdraw a Route?

50

```
7/5 19:33:25      Route R is withdrawn
7/5 19:34:15      AS6543 announce R 6543 66665 8918 1 5696 999
7/5 19:35:00      AS6543 announce R 6543 66665 8918 67455 6461 5696 999
7/5 19:35:37      AS6543 announce R 6543 66665 4332 6461 5696 999
7/5 19:35:39      AS6543 announce R 6543 66665 5378 6660 67455 6461 5696 999
7/5 19:35:39      AS6543 announce R 6543 66665 65 6461 5696 999
7/5 19:35:52      AS6543 announce R 6543 66665 6461 5696 999
7/5 19:36:00      AS6543 announce R 6543 66665 5378 6765 6660 67455 6461 5696 999
...
7/5 19:38:22      AS6543 withdraw R
```

Answer: *up to 19*

BGP Routing Table Convergence Times



- ❑ Less than half of Tdown events converge within two minutes
- ❑ T-up/T-short and T-down/T-long form equivalence classes
- ❑ Long tailed distribution (up to 15 minutes)

Failures, Fail-overs and Repairs

52

- Bad news does not travel fast...
- Repairs (Tup) exhibit similar convergence as long-short AS path fail-over
- Failures (Tdown) and short-long fail-overs (e.g. primary to secondary path) also similar
 - ▣ Slower than Tup (e.g. a repair)
 - ▣ 80% take longer than two minutes
 - ▣ Fail-over times degrade the greater the degree of multi-homing

Intuition for Delayed Convergence

- There exists possible ordering of messages such that BGP will explore ALL possible AS paths of ALL possible lengths
- BGP is $O(N!)$, where N number of default-free BGP routers in a complete graph with default policy

Impact of Delayed Convergence

- Why do we care about routing table convergence?
 - ▣ It impacts end-to-end connectivity for Internet paths
- ICMP experiment results
 - ▣ Loss of connectivity, packet loss, latency, and packet re-ordering for an average of 3-5 minutes after a fault
- Why?
 - ▣ Routers drop packets when next hop is unknown
 - ▣ Path switching spikes latency/delay
 - ▣ Multi-pathing causes reordering

In real life ...

- Discussed worst case BGP behavior
- In practice, BGP policy prevents worst case from happening
- BGP timers also provide synchronization and limits possible orderings of messages

Inter-Domain Routing Summary

- BGP4 is the only inter-domain routing protocol currently in use world-wide
- Issues?
 - Lack of security
 - Ease of misconfiguration
 - Poorly understood interaction between local policies
 - Poor convergence
 - Lack of appropriate information hiding
 - Non-determinism
 - Poor overload behavior

Lots of research into how to fix this

57

- Security
 - ▣ BGPSEC, RPKI
- Misconfigurations, inflexible policy
 - ▣ SDN
- Policy Interactions
 - ▣ PoiRoot (root cause analysis)
- Convergence
 - ▣ Consensus Routing
- Inconsistent behavior
 - ▣ LIFEGUARD, among others

Why are these still issues?

58

- Backward compatibility
- Buy-in / incentives for operators
- Stubbornness

Why are these still issues?

58

- Backward compatibility
- Buy-in / incentives for operators
- Stubbornness

Very similar issues to IPv6 deployment