

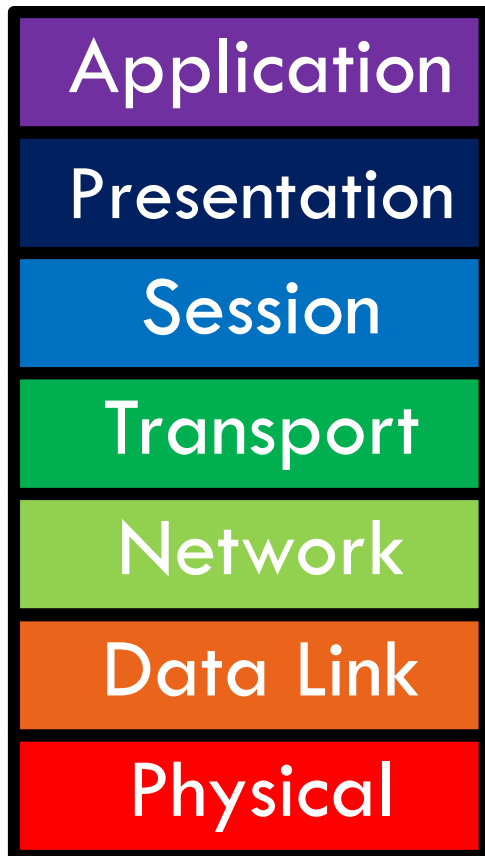
Computer Networks

Lecture 10: Network Layer – Part II

Based on slides from D. Choffnes Northeastern U. and P. Gill from StonyBrook University
Revised Autumn 2015 by S. Laki

Network Layer

2



□ Function:

- ▣ Route packets end-to-end on a network, through multiple hops

□ Key challenge:

- ▣ How to represent addresses
- ▣ How to route packets
 - Scalability
 - Convergence

Intra-domain Routing Protocols

3

- Distance vector
 - ▣ Routing Information Protocol (RIP), based on Bellman-Ford
 - ▣ Routers periodically exchange reachability information with neighbors
- Link state
 - ▣ Open Shortest Path First (OSPF), based on Dijkstra
 - ▣ Each network periodically **floods** immediate reachability information to all other routers
 - ▣ Per router local computation to determine full routes

- ❑ Distance Vector Routing
 - ❑ RIP
- ❑ Link State Routing
 - ❑ OSPF
 - ❑ IS-IS

Distance Vector Routing

5

- What is a distance vector?
 - ▣ Current best known cost to reach a destination
- Idea: exchange vectors among neighbors to learn about lowest cost paths

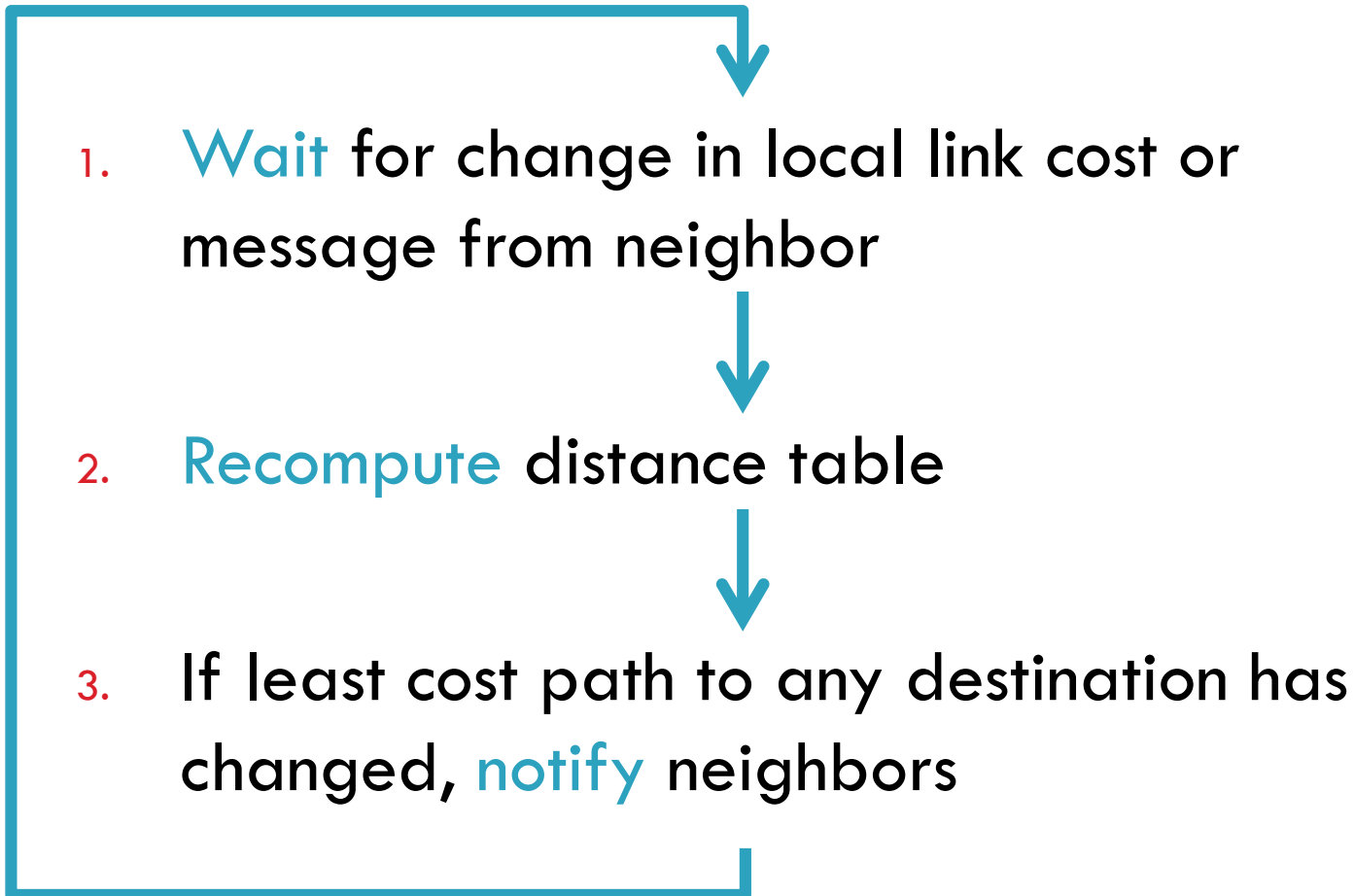
DV Table
at Node C

Destination	Cost
A	7
B	1
D	2
E	5
F	1

- No entry for C
 - Initially, only has info for immediate neighbors
 - ▣ Other destinations cost = ∞
 - Eventually, vector is filled
- Routing Information Protocol (RIP)

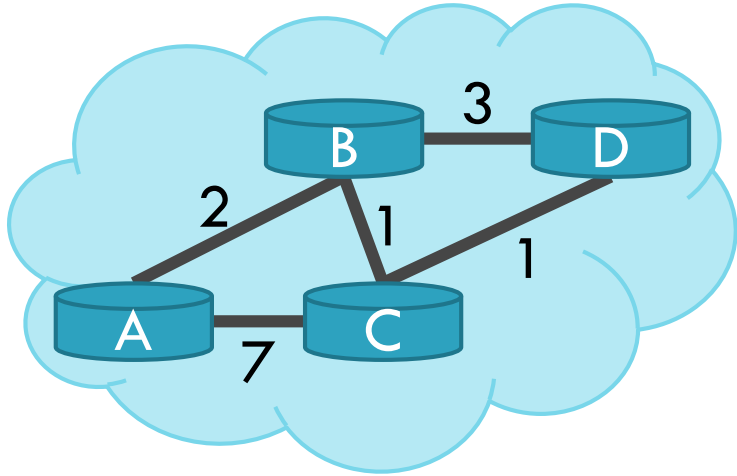
Distance Vector Routing Algorithm

6

- 
- ```
graph TD; A[] --> B[1. Wait for change in local link cost or message from neighbor]; B --> C[2. Recompute distance table]; C --> D[3. If least cost path to any destination has changed, notify neighbors]; D --> A;
```
1. **Wait** for change in local link cost or message from neighbor
  2. **Recompute** distance table
  3. If least cost path to any destination has changed, **notify** neighbors

# Distance Vector Initialization

7



Node A

| Dest. | Cost     | Next |
|-------|----------|------|
| B     | 2        | B    |
| C     | 7        | C    |
| D     | $\infty$ |      |

Node B

| Dest. | Cost | Next |
|-------|------|------|
| A     | 2    | A    |
| C     | 1    | C    |
| D     | 3    | D    |

1. Initialization:
2. for all neighbors  $V$  do
3. if  $V$  adjacent to  $A$
4.  $D(A, V) = c(A, V)$ ;
5. else
6.  $D(A, V) = \infty$ ;
- ...

Node C

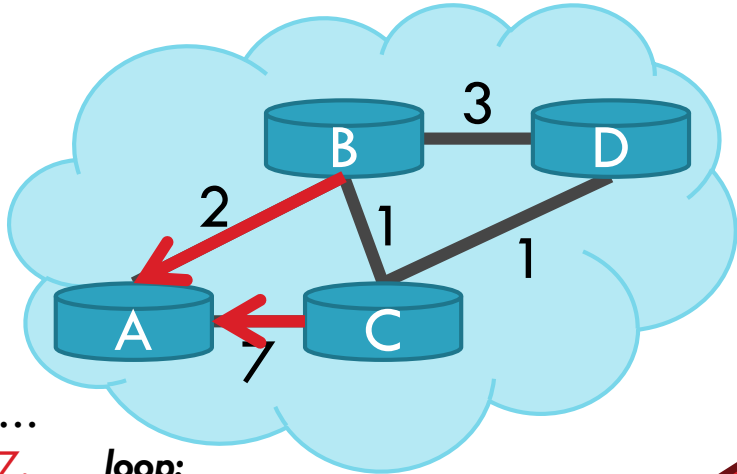
| Dest. | Cost | Next |
|-------|------|------|
| A     | 7    | A    |
| B     | 1    | B    |
| D     | 1    | D    |

Node D

| Dest. | Cost     | Next |
|-------|----------|------|
| A     | $\infty$ |      |
| B     | 3        | B    |
| C     | 1        | C    |

# Distance Vector: 1<sup>st</sup> Iteration

8



Node A

| Dest. | Cost | Next |
|-------|------|------|
| B     | 2    | B    |
| C     | 3    | B    |
| D     | 5    | B    |

Node B

| Dest. | Cost | Next |
|-------|------|------|
| A     | 2    | A    |
| C     | 1    | C    |
| D     | 2    | C    |



...  
7. loop:  
...  
12. else if (update D(V, Y) received)  
13. for all destinations Y  
14. if (destination Y is not A)  
15.  $D(A, Y) = \min(D(A, Y), D(A, B) + D(B, Y))$   
16. else  
17.  $D(A, Y) = \min(D(A, Y), D(A, C) + D(C, Y))$   
18. if (there is a new min. for dest. Y)  
19. send D(A, Y) to all neighbors  
20. forever

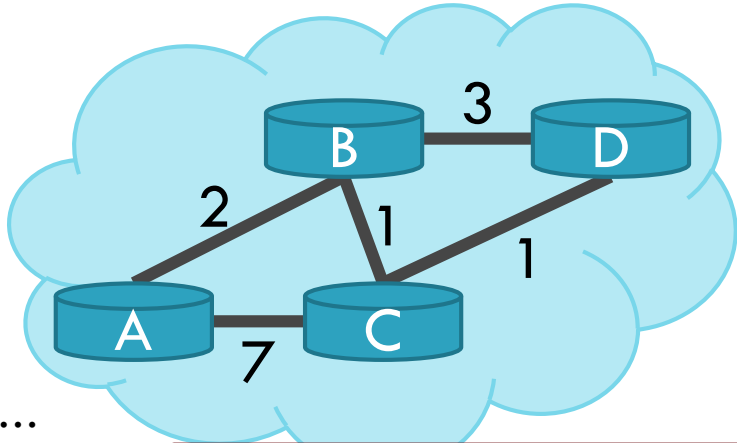
$D(A, C) = \min(D(A, C), D(A, B) + D(B, C))$   
 $D(A, D) = \min(D(A, D), D(A, B) + D(B, D))$   
 $= \min(8, 3 + 3) = 5$

| Dest. | Cost | Next |
|-------|------|------|
| B     | 1    | B    |
| D     | 1    | D    |
| B     | 3    | B    |
| C     | 1    | C    |



# Distance Vector: End of 3<sup>rd</sup> Iteration

9



Node A

| Dest. | Cost | Next |
|-------|------|------|
| B     | 2    | B    |
| C     | 3    | B    |
| D     | 4    | B    |

Node B

| Dest. | Cost | Next |
|-------|------|------|
| A     | 2    | A    |
| C     | 1    | C    |
| D     | 2    | C    |

- Nothing changes, algorithm terminates
- Until something changes...

| Dest. | Cost | Next |
|-------|------|------|
| A     | 3    | B    |
| B     | 1    | B    |
| D     | 1    | D    |

| Dest. | Cost | Next |
|-------|------|------|
| A     | 4    | C    |
| B     | 2    | C    |
| C     | 1    | C    |

```

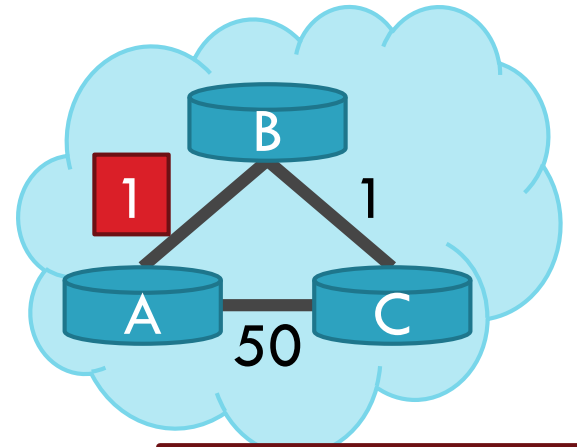
...
7. loop
...
12. else
13. for
14.
15.
16. else
17. D(A, Y) =
 min(D(A, Y),
 D(A, V) + D(V, Y));
18. if (there is a new min. for dest. Y)
19. send D(A, Y) to all neighbors
20. forever

```

```

7. loop:
8. wait (link cost update or update message)
9. if (c(A,V) changes by d)
10. for all destinations Y through V do
11. D(A,Y) = D(A,Y) + d
12. else if (update D(V, Y) received from V)
13. for all destinations Y do
14. if (destination Y through V)
15. D(A,Y) = D(A,V) + D(V, Y);
16. else
17. D(A, Y) = min(D(A, Y), D(A, V) + D(V, Y));

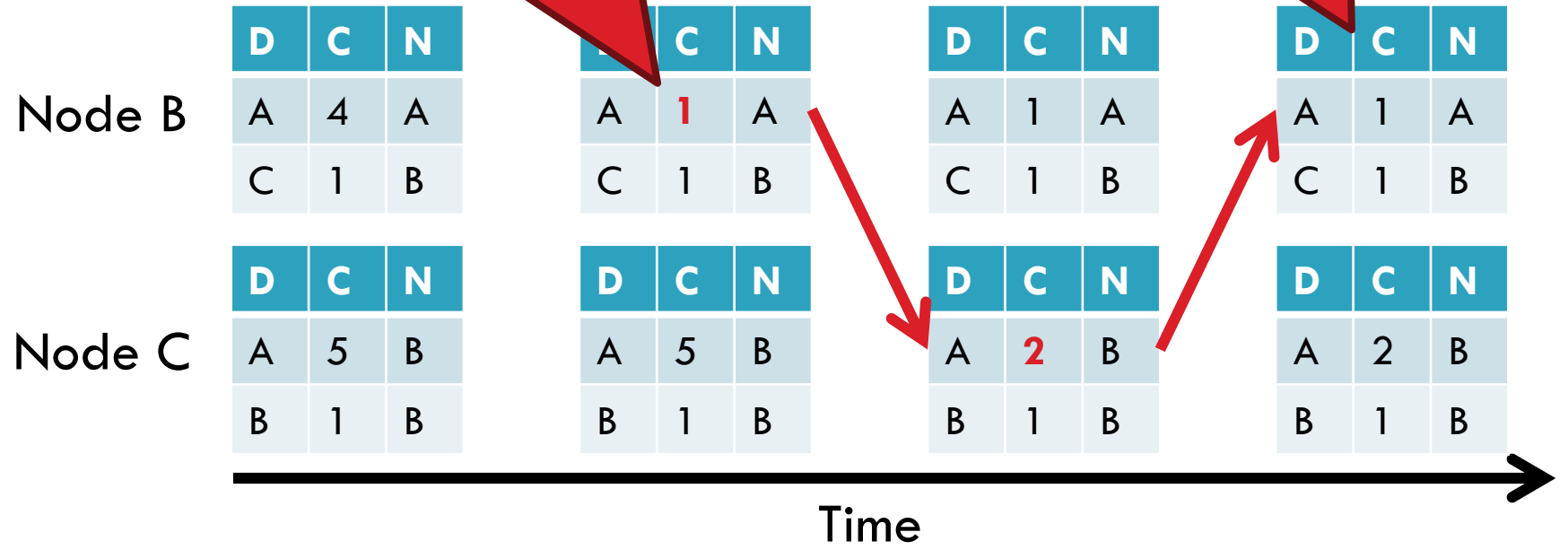
```



Link Cost Algorithm

Good news travels fast

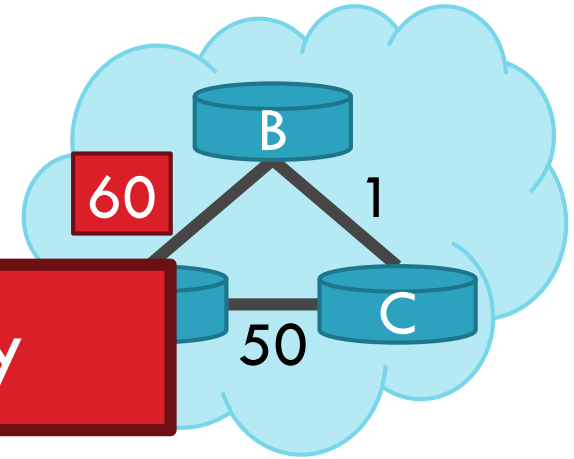
Algorithm terminates



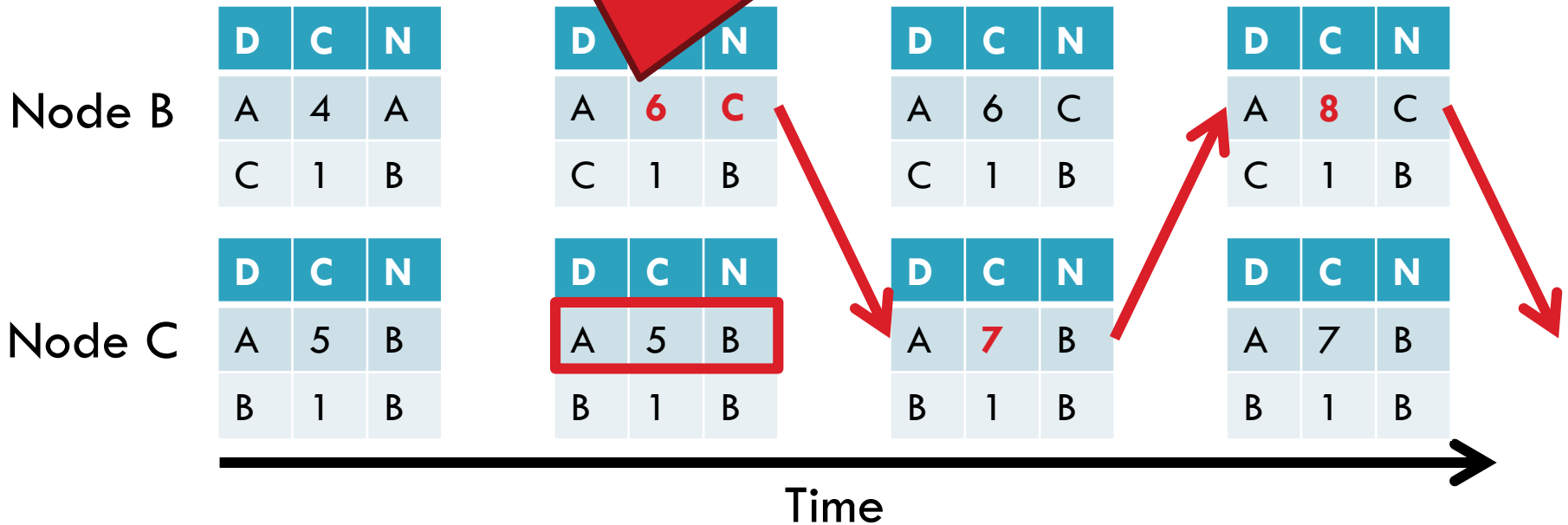
# Count to Infinity Problem

11

- Node B knows  $D(C, A) = 5$
- However, B does not know the path is  $C \rightarrow B \rightarrow A$
- Thus,  $D(B, A) = 60$



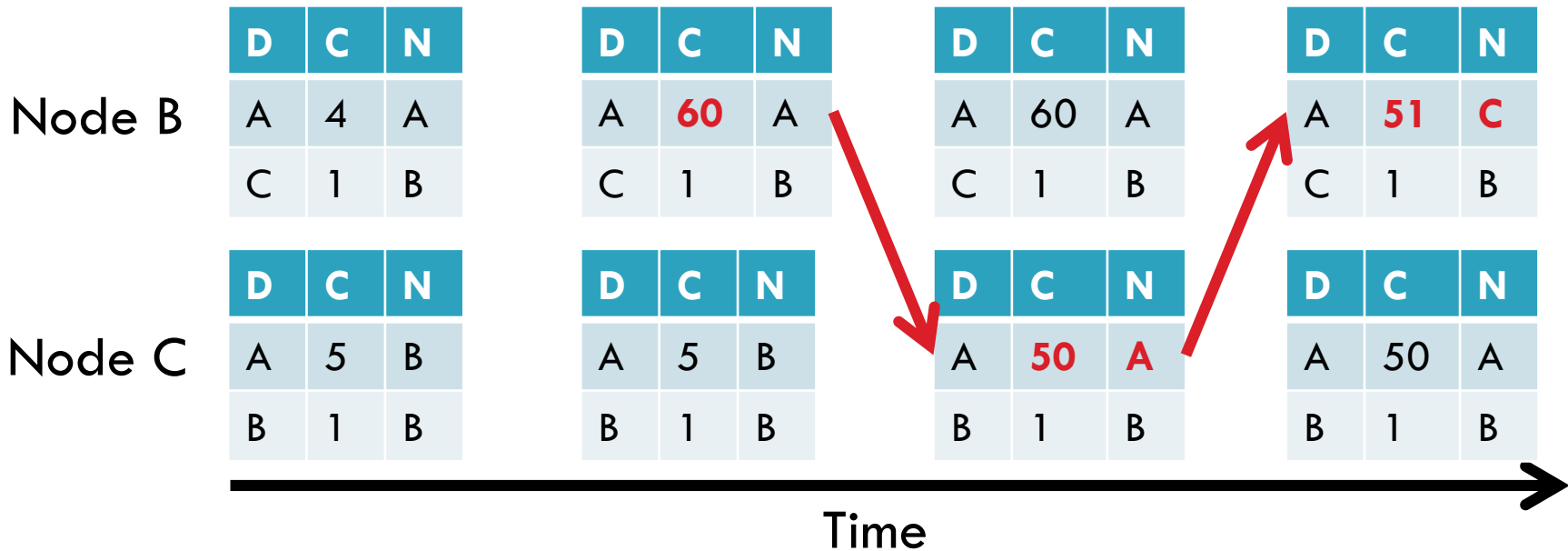
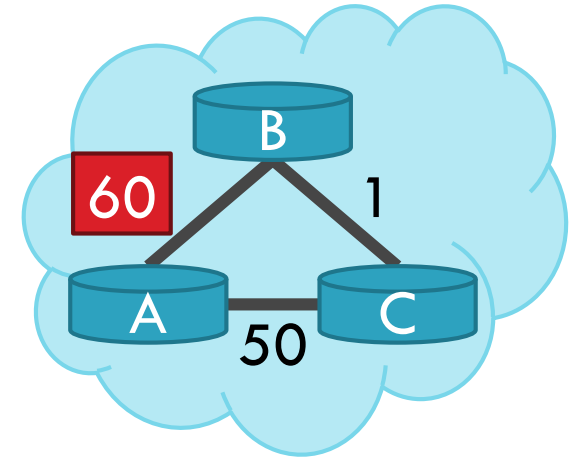
Bad news travels slowly



# Poisoned Reverse

12

- If C routes through B to get to A
  - ▣ C tells B that  $D(C, A) = \infty$
  - ▣ Thus, B won't route to A via C

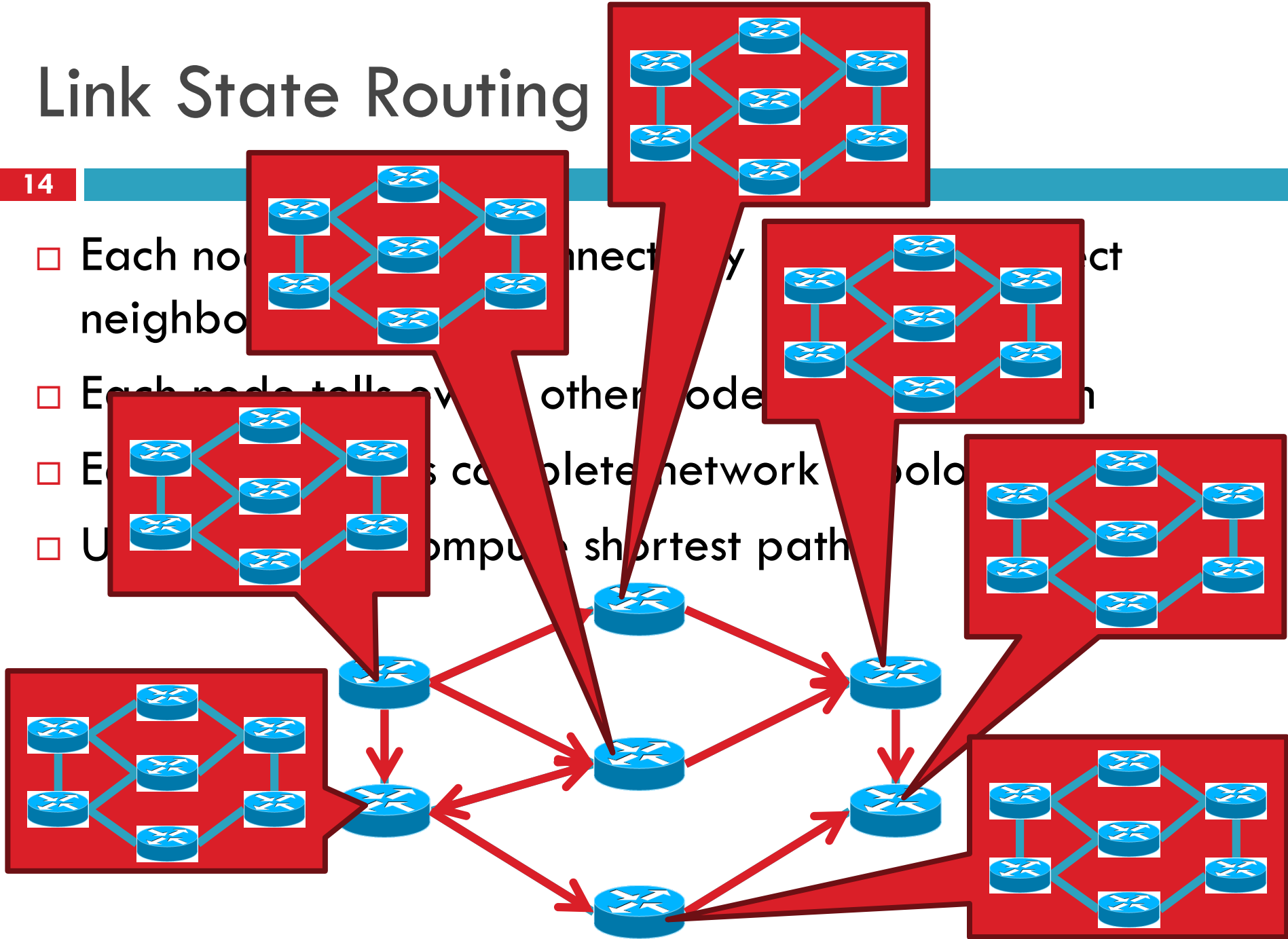


- ❑ Distance Vector Routing
  - ❑ RIP
- ❑ Link State Routing
  - ❑ OSPF
  - ❑ IS-IS

# Link State Routing

14

- Each node connects to its neighbors
- Each node tells every other node about its neighbors
- Each node has complete network topology
- Use Dijkstra's algorithm to compute shortest path



# Flooding Details

15

- Each node periodically generates Link State Packet
  - ▣ ID of node generating the LSP
  - ▣ List of direct neighbors and costs
  - ▣ Sequence number (64-bit, assumed to never wrap)
  - ▣ Time to live
- Flood is reliable (ack + retransmission)
- Sequence number “versions” each LSP
- Receivers flood LSPs to their own neighbors
  - ▣ Except whoever originated the LSP
- LSPs also generated when link states change

# OSPF vs. IS-IS

16

- Two different implementations of link-state routing

## OSPF

- Favored by companies, datacenters
- More optional features
- Built on top of IPv4
  - ▣ LSAs are sent via IPv4
  - ▣ OSPFv3 needed for IPv6

## IS-IS

- Favored by ISPs
- Less “chatty”
  - ▣ Less network overhead
  - ▣ Supports more devices
- Not tied to IP
  - ▣ Works with IPv4 or IPv6

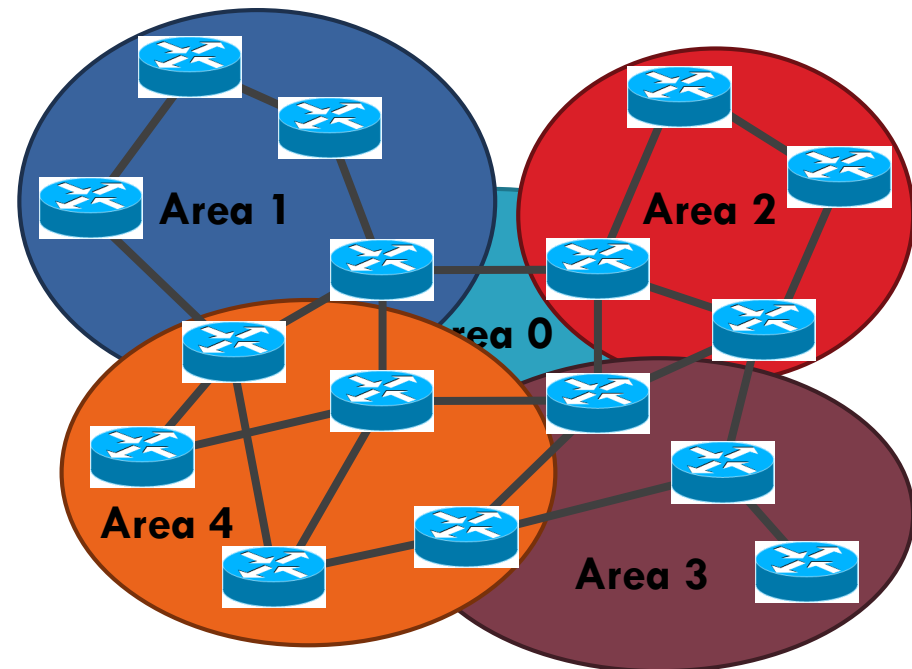


# Different Organizational Structure

17

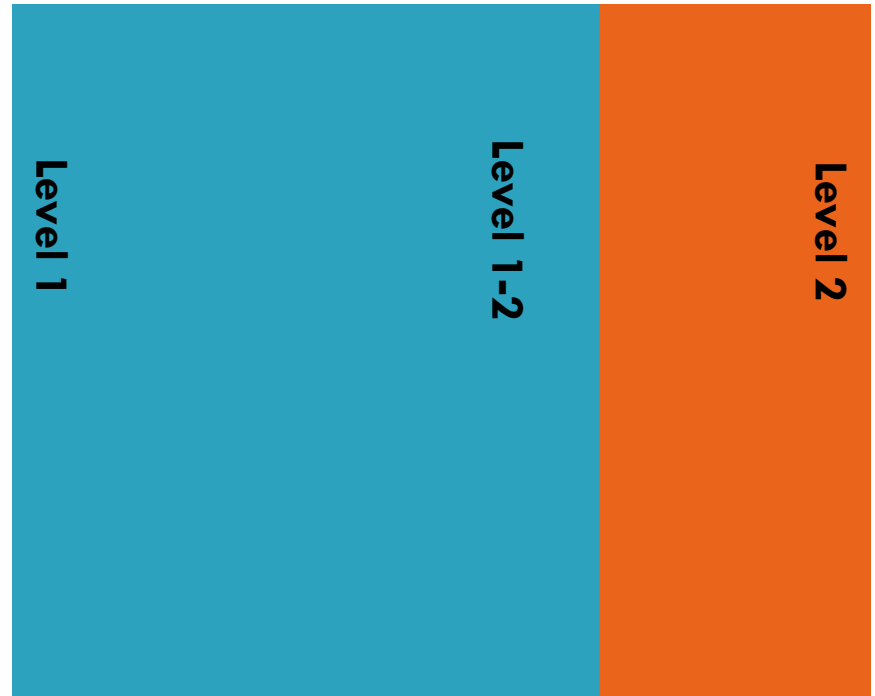
## OSPF

- ❑ Organized around overlapping areas
- ❑ Area 0 is the core network



## IS-IS

- ❑ Organized as a 2-level hierarchy
- ❑ Level 2 is the backbone



# Network Layer, Control Plane

18

Data Plane

Application

Presentation

Session

Transport

Network

Data Link

Physical

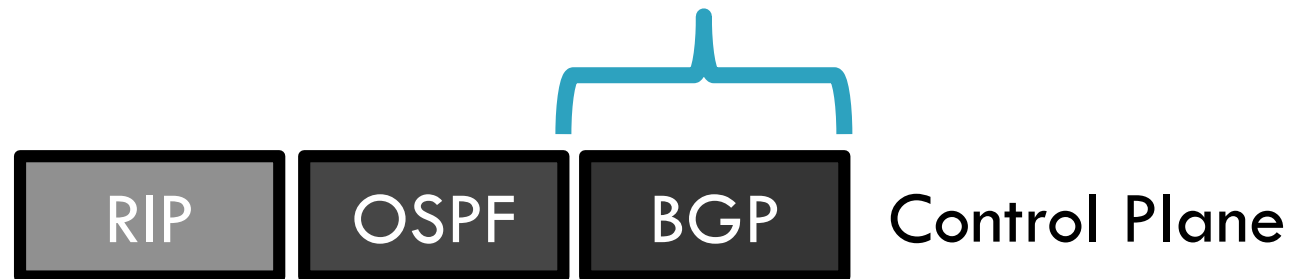
□ Function:

▣ Set up routes between networks

□ Key challenges:

▣ Implementing provider policies

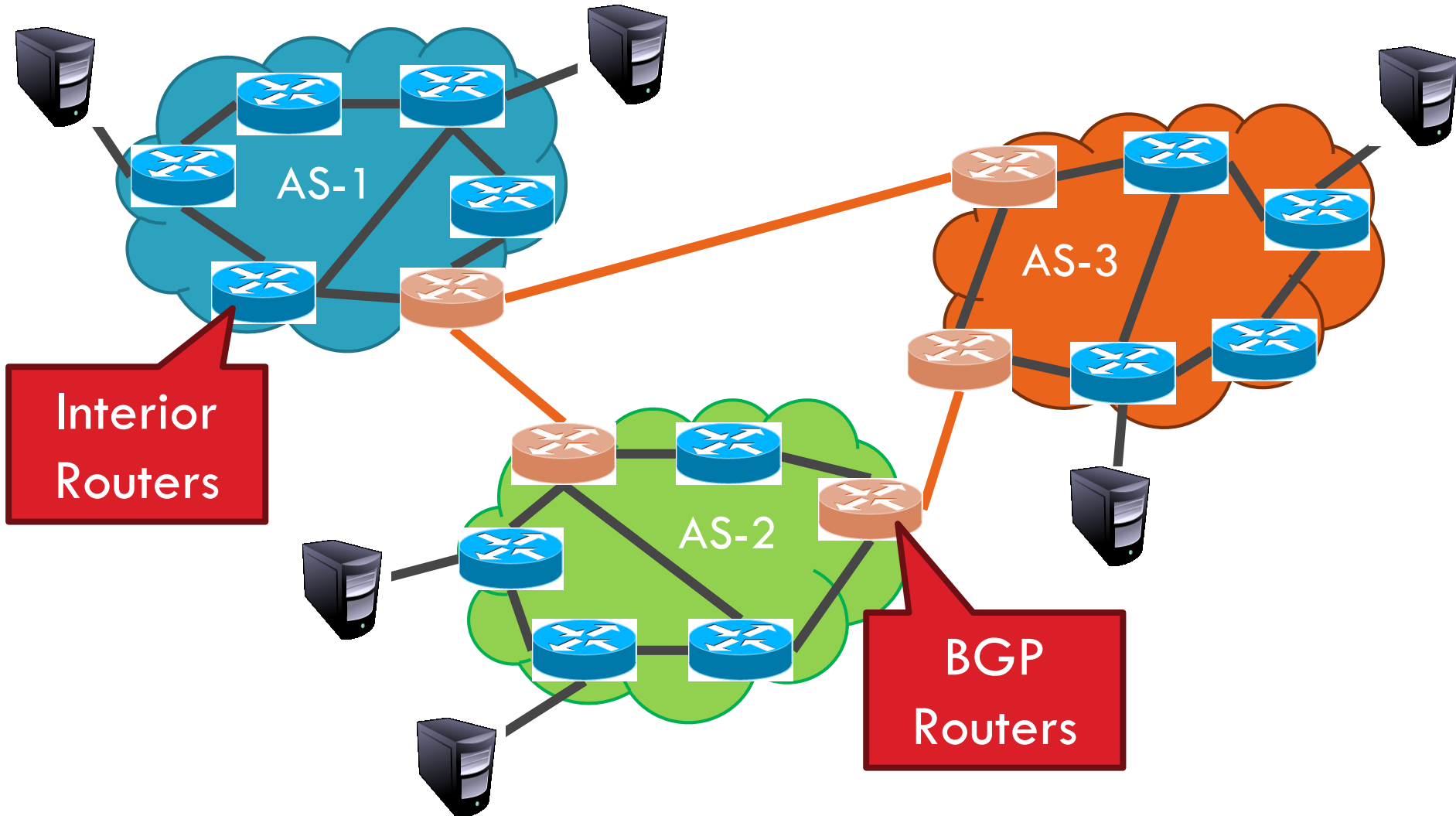
▣ Creating stable paths



- ❑ BGP Basics
- ❑ Stable Paths Problem
- ❑ BGP in the Real World
- ❑ Debugging BGP Path Problems

# ASs, Revisited

20



# AS Numbers

21

- ❑ Each AS identified by an ASN number
  - ❑ 16-bit values (latest protocol supports 32-bit ones)
  - ❑ 64512 – 65535 are reserved
- ❑ Currently, there are ~ 40000 ASNs
  - ❑ AT&T: 5074, 6341, 7018, ...
  - ❑ Sprint: 1239, 1240, 6211, 6242, ...
  - ❑ ELTE: 2012
  - ❑ Google 15169, 36561 (formerly YT), + others
  - ❑ Facebook 32934
  - ❑ North America ASs → <ftp://ftp.arin.net/info/asn.txt>

# Inter-Domain Routing

22

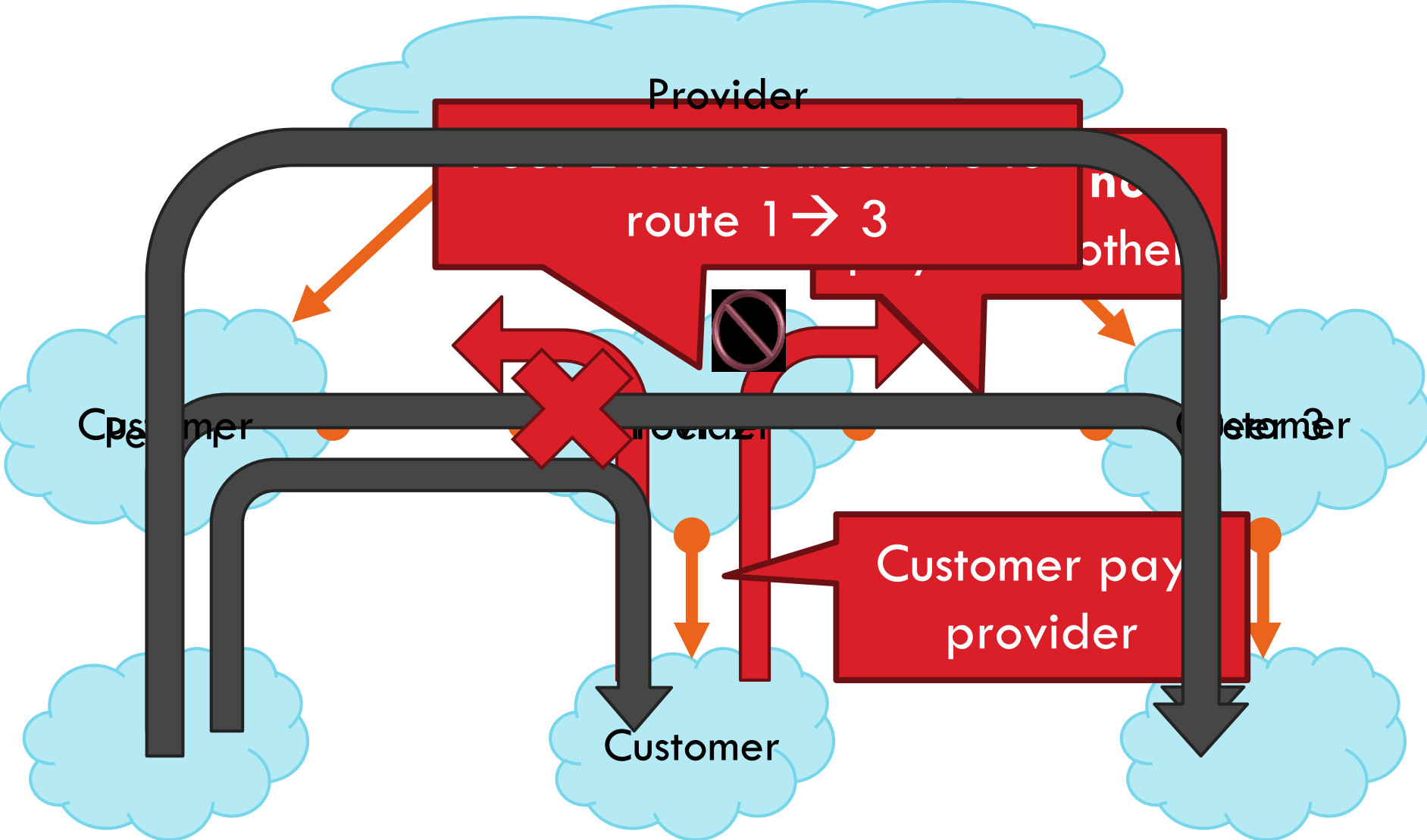
- Global connectivity is at stake!
  - ▣ Thus, all ASs must use the same protocol
  - ▣ Contrast with intra-domain routing
- What are the requirements?
  - ▣ Scalability
  - ▣ Flexibility in choosing routes
    - Cost
    - Routing around failures
- Question: link state or distance vector?
  - ▣ Trick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

# BGP

23

- 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)

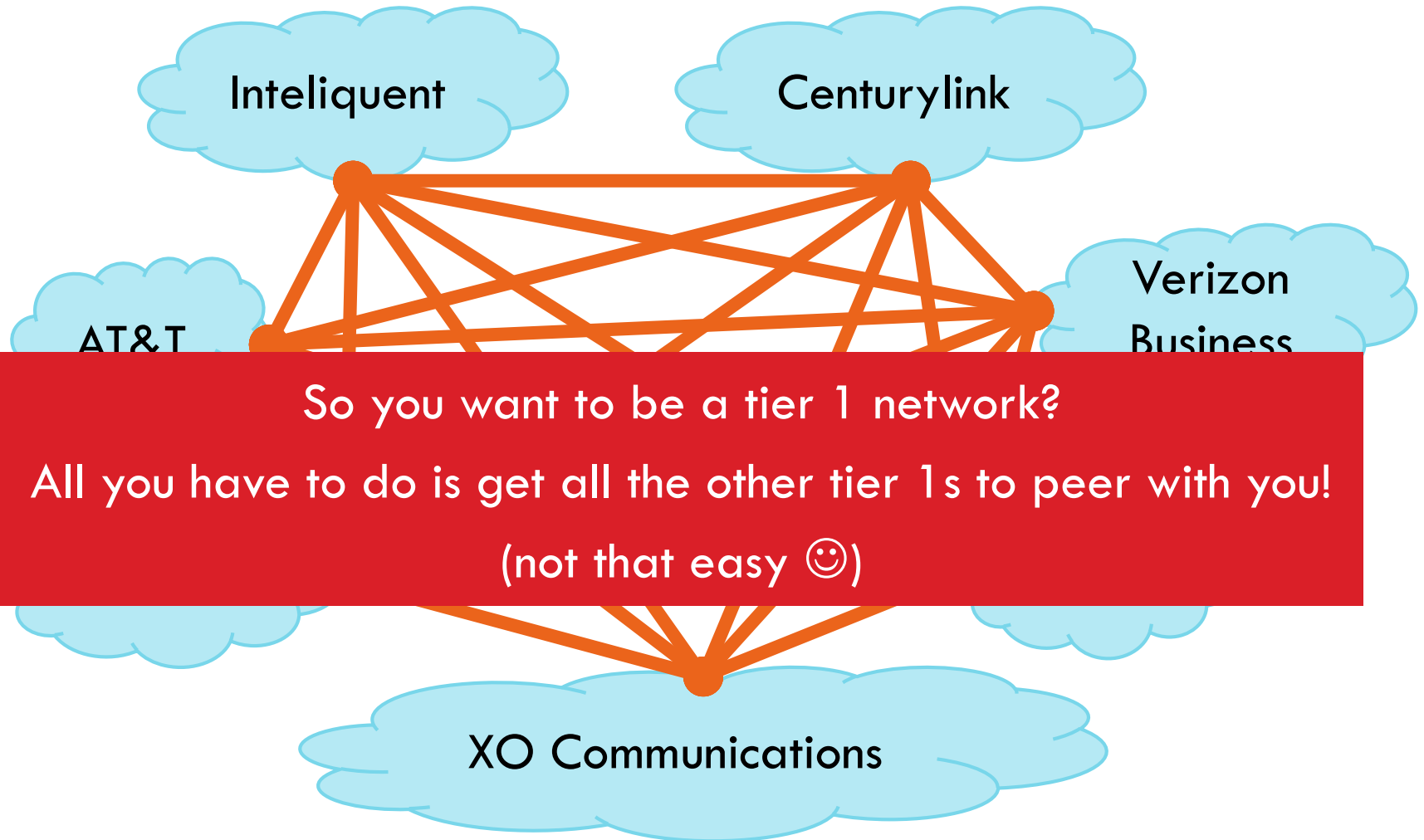
# BGP Relationships





# Tier-1 ISP Peering

25





# Peering Wars

27

Peer

Don't Peer

- Reduce upstream costs
- You would rather have

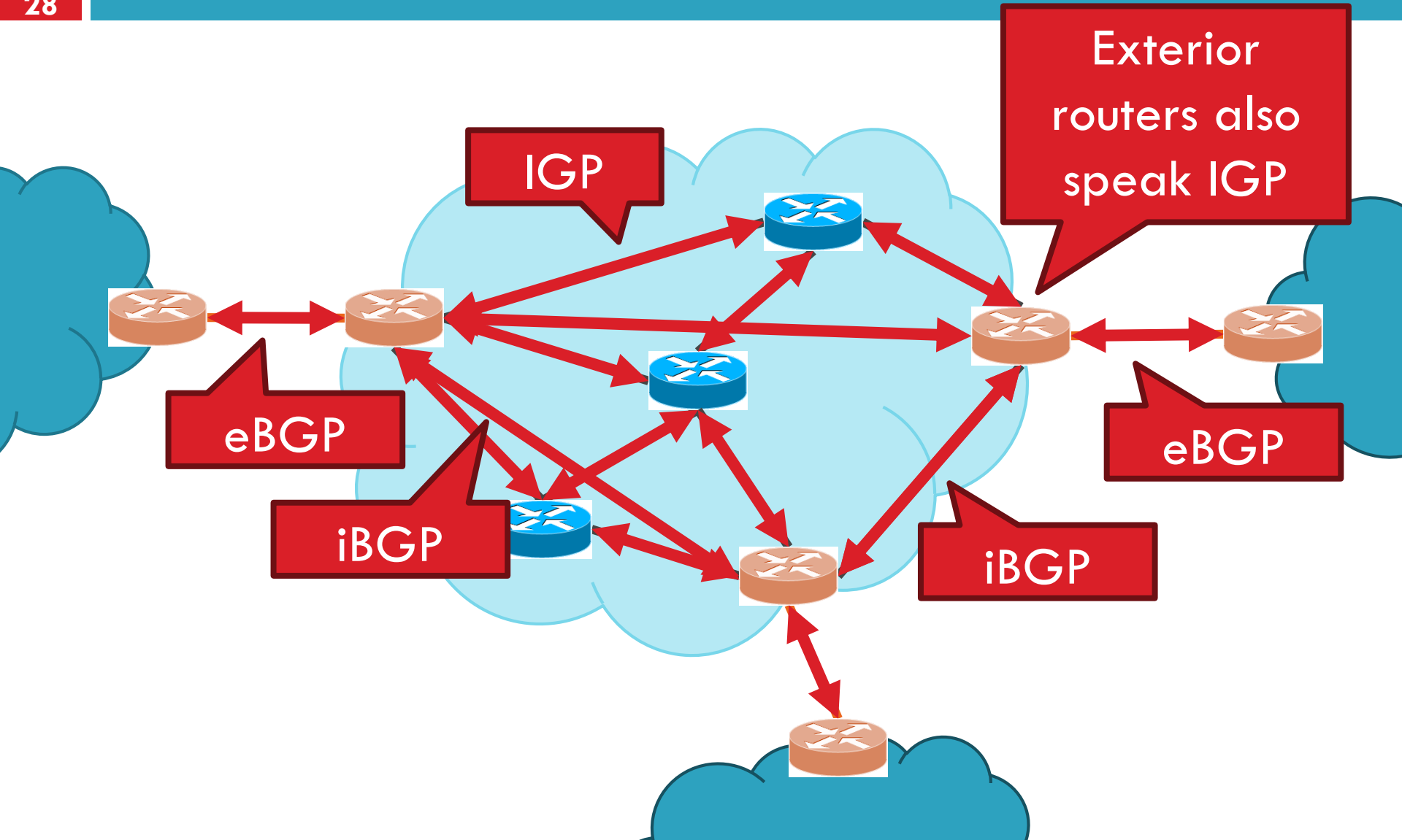
- Peering struggles in the ISP world are extremely contentious
- agreements are usually confidential

- Example: If you are a customer of my peer why should I peer with you? You should pay me too!

Incentive to keep relationships private!

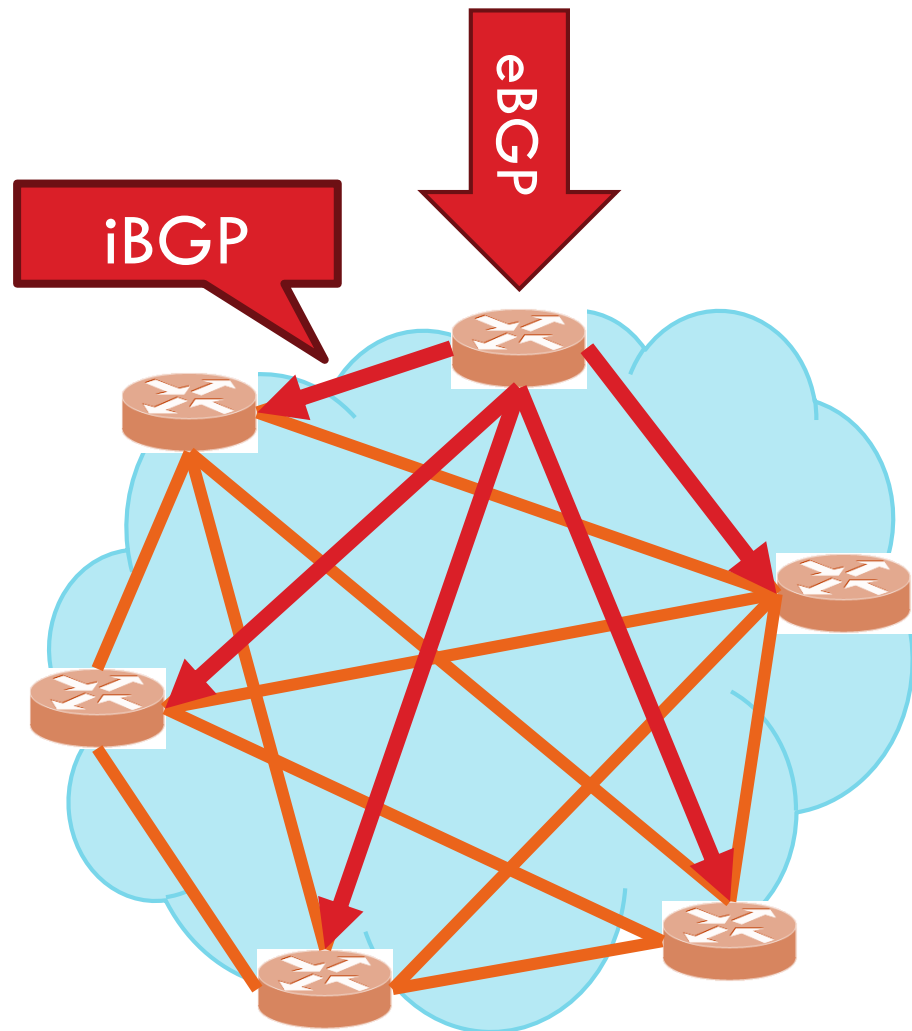
# Two Types of BGP Neighbors

28



# Full iBGP Meshes

29

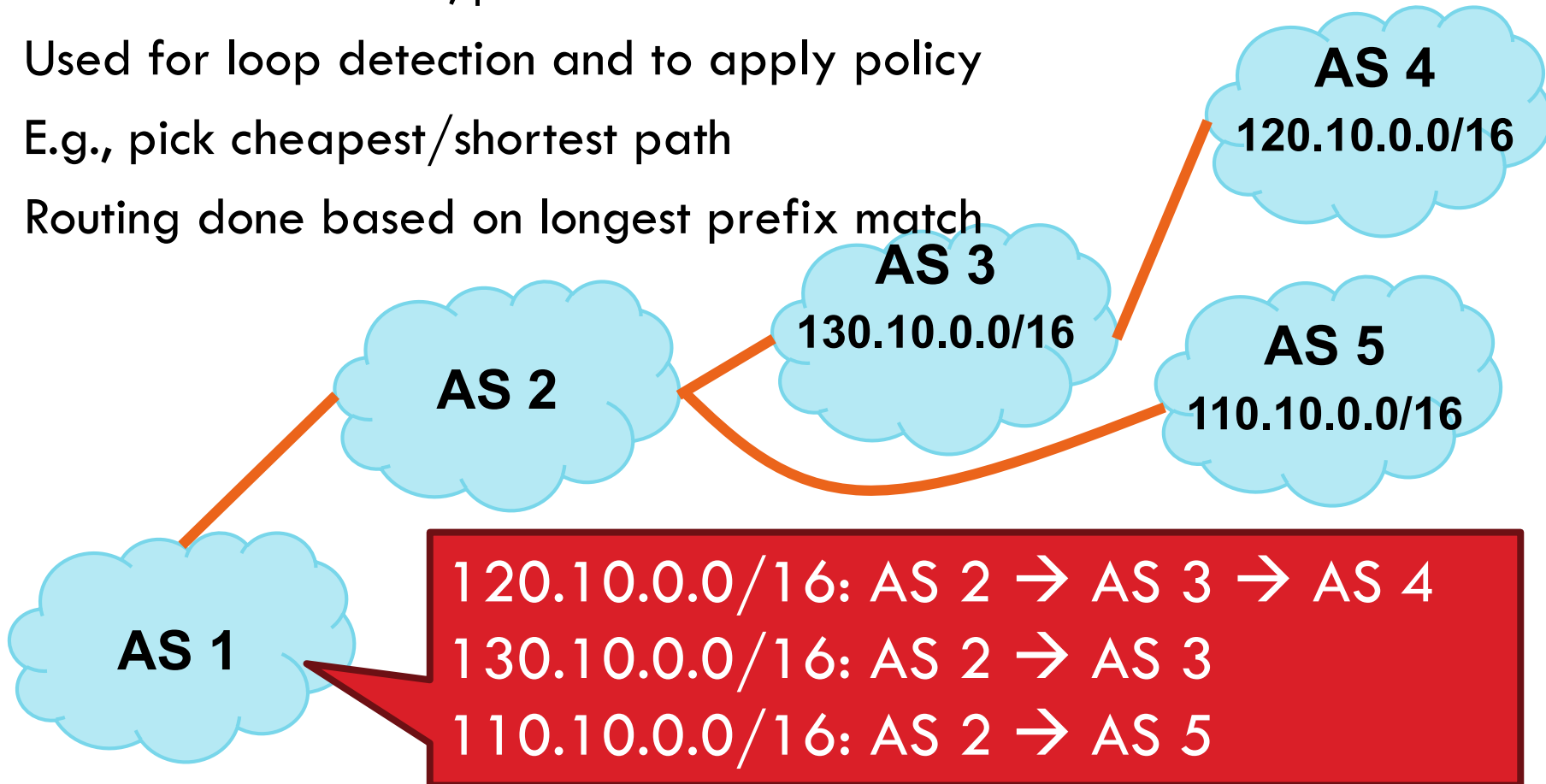


- Question: why do we need iBGP?
  - ▣ OSPF does not include BGP policy info
  - ▣ Prevents routing loops within the AS
- iBGP updates do not trigger announcements

# Path Vector Protocol

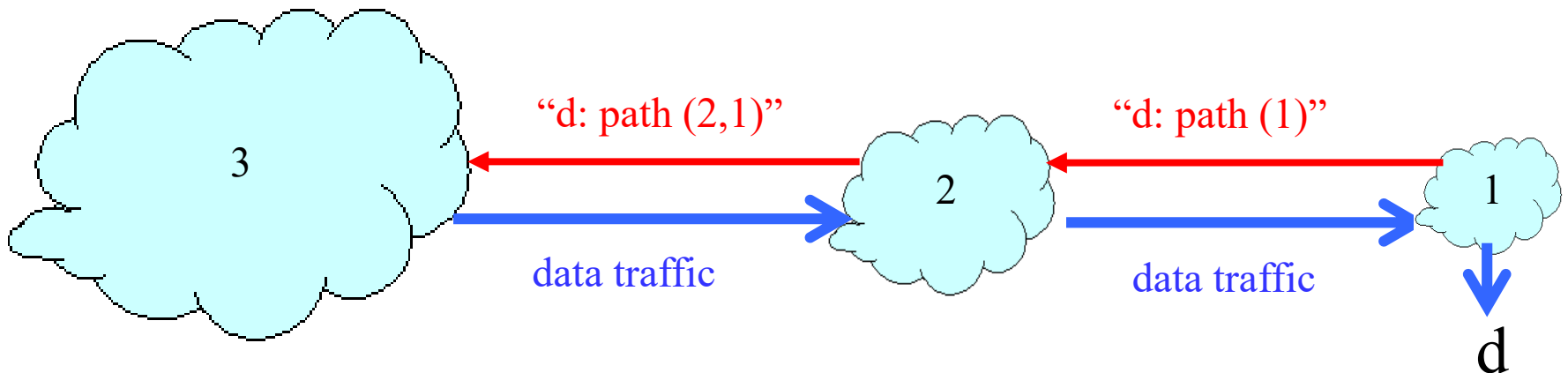
30

- AS-path: sequence of ASs a route traverses
  - ▣ Like distance vector, plus additional information
- Used for loop detection and to apply policy
- E.g., pick cheapest/shortest path
- Routing done based on longest prefix match



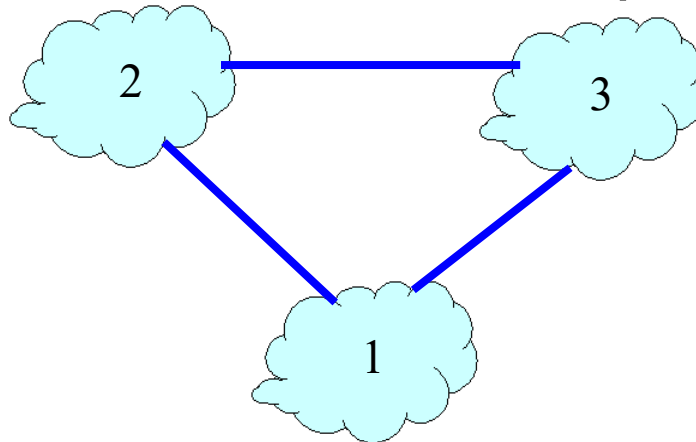
# Path-Vector Routing

- 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*



# Flexible Policies

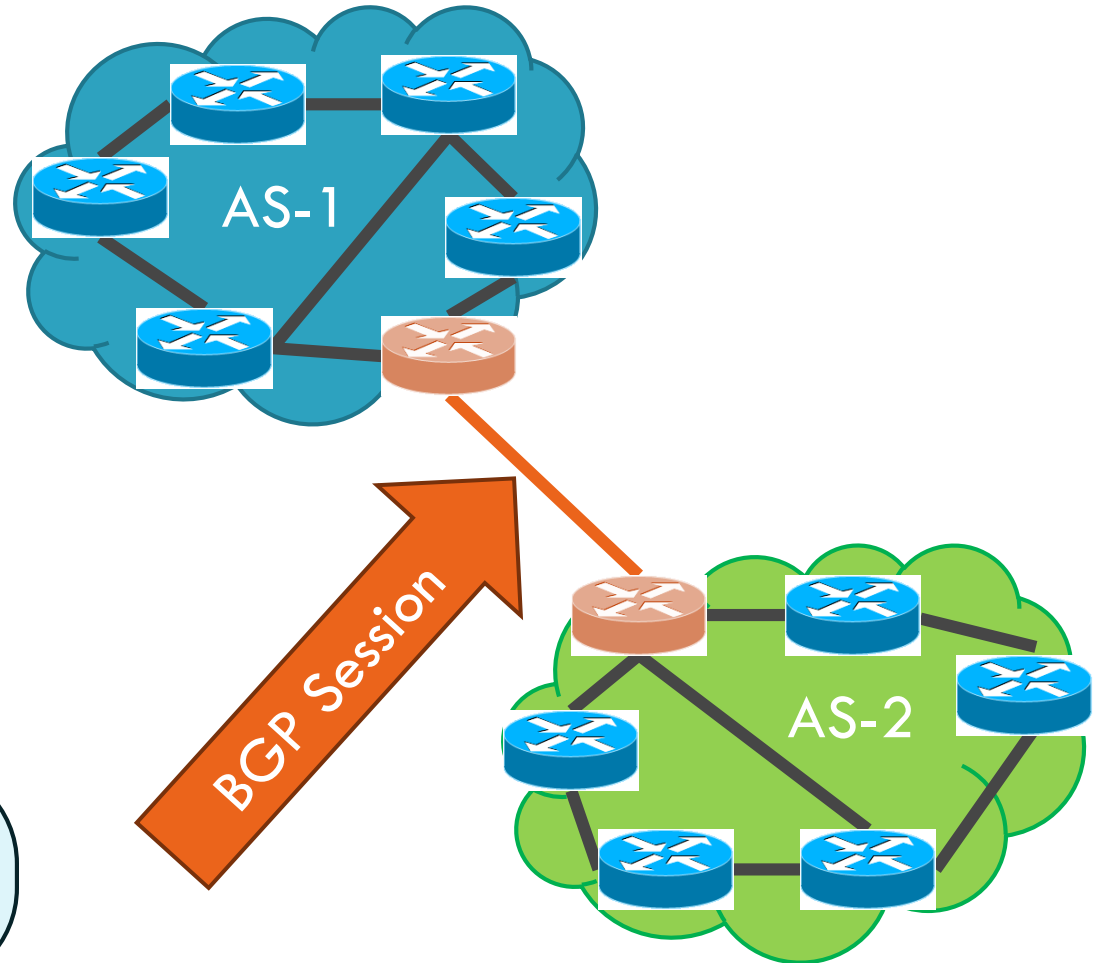
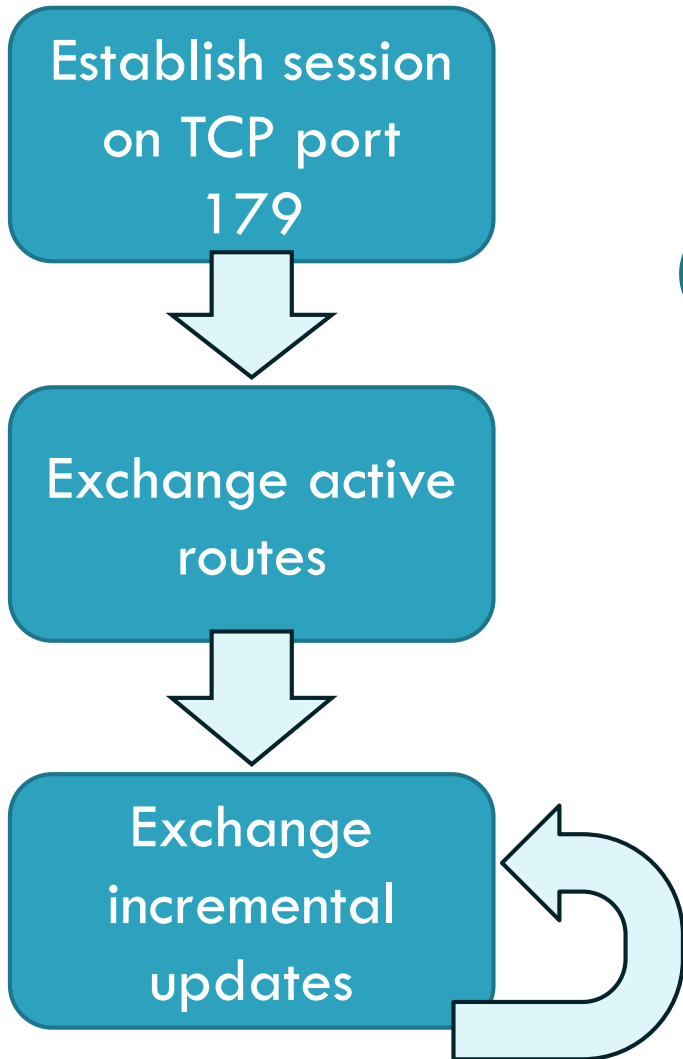
- Each node can apply local policies
  - ▣ Path selection: Which path to use?
  - ▣ Path export: Which paths to advertise?
- Examples
  - ▣ Node 2 may prefer the path “2, 3, 1” over “2, 1”
  - ▣ Node 1 may not let node 3 hear the path “1, 2”





# BGP Operations (Simplified)

33



# Four Types of BGP Messages

34

- ❑ **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

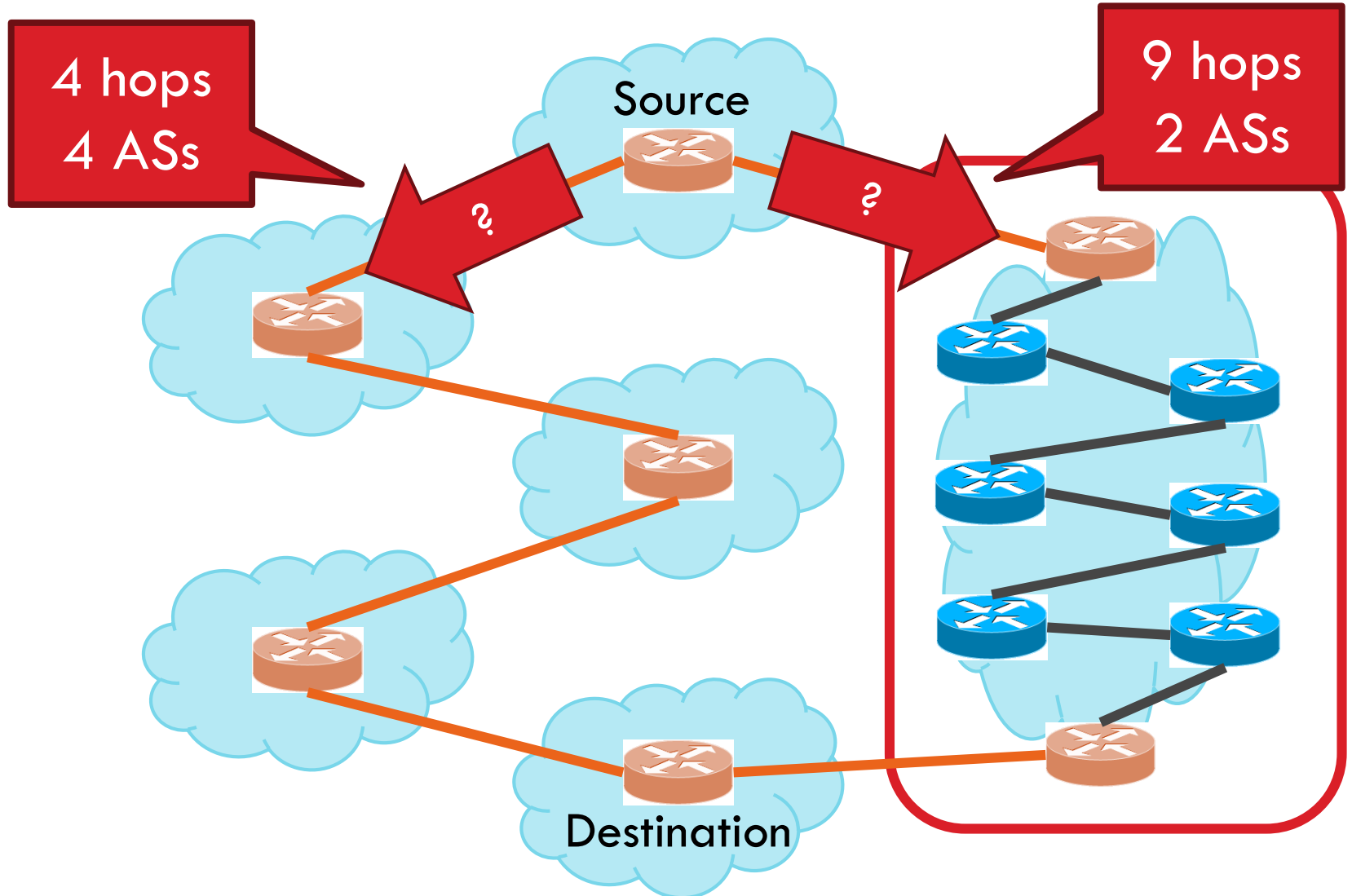
# BGP Attributes

35

- 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?

# Shortest AS Path $\neq$ Shortest Path

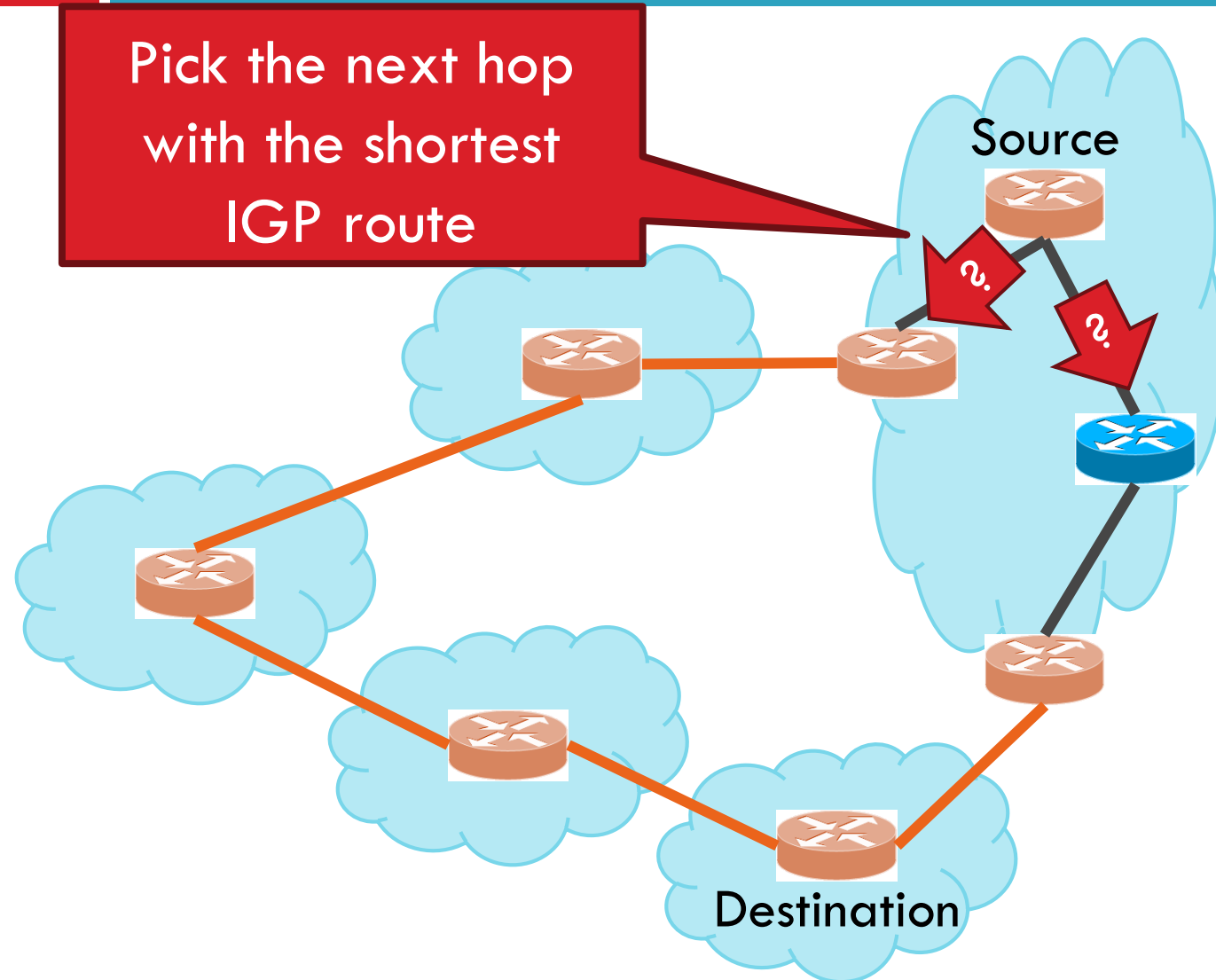
36



# Hot Potato Routing

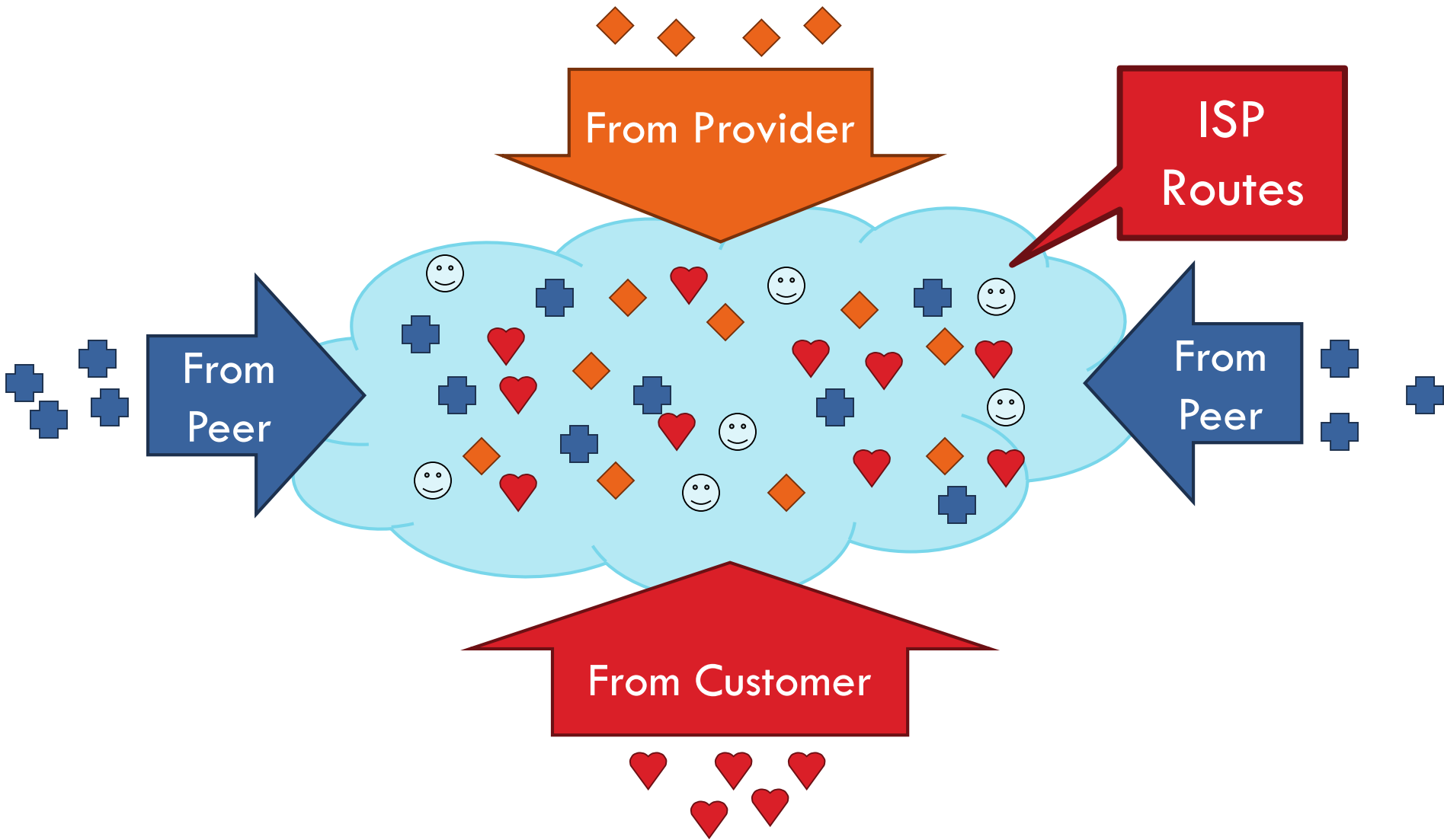
37

Pick the next hop with the shortest IGP route



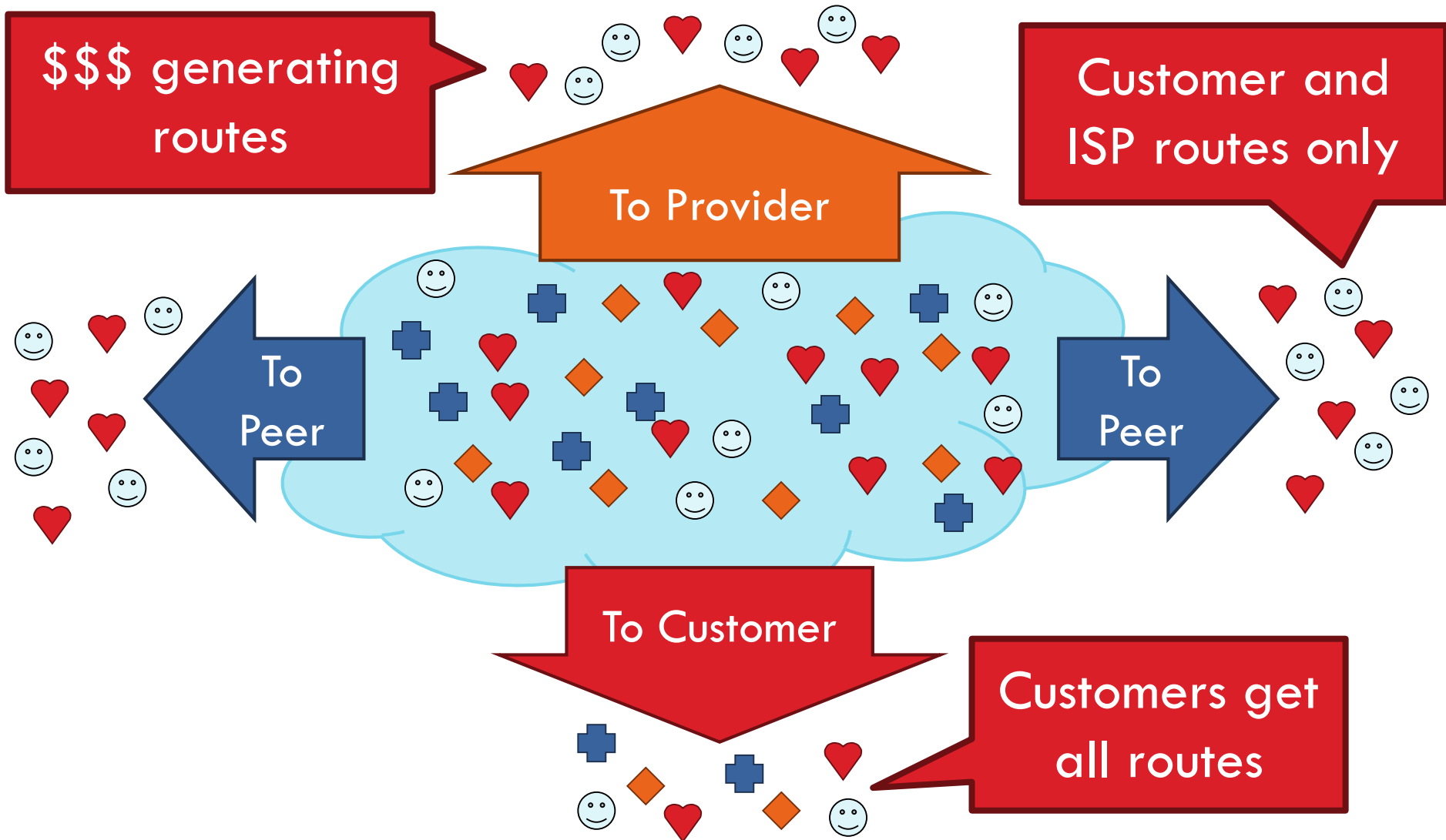
# Importing Routes

38



# Exporting Routes

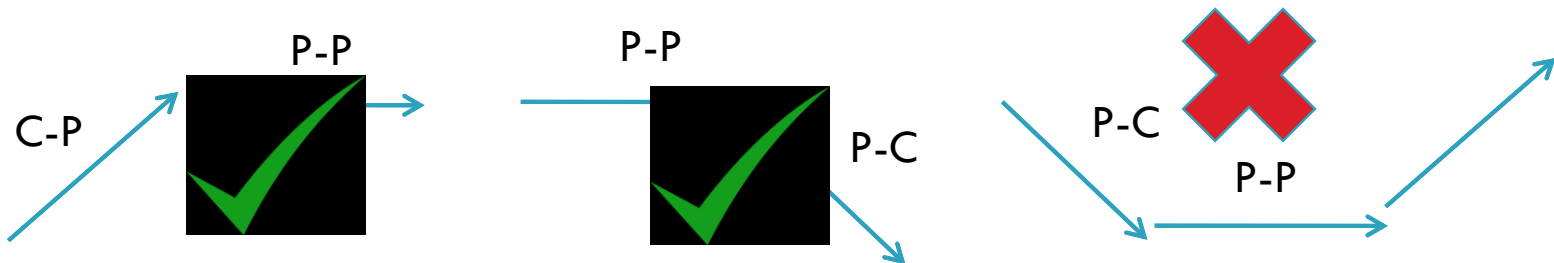
39



# Modeling BGP

40

- 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

41

- 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

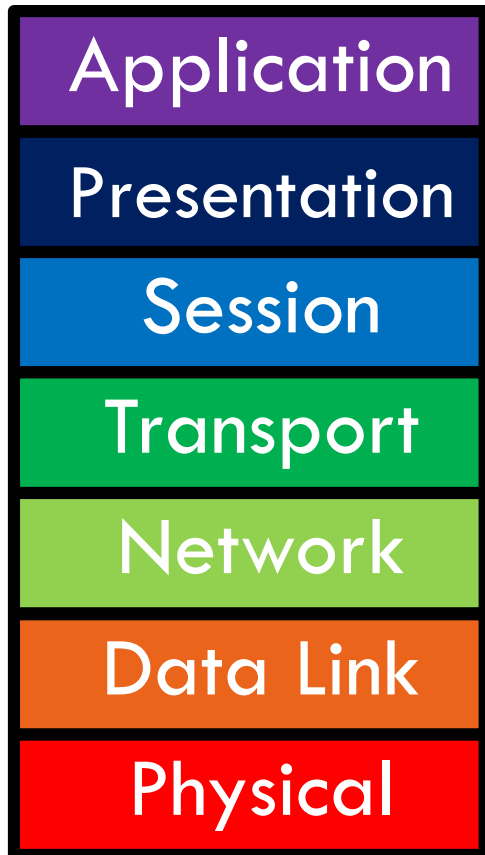
# Other BGP Attributes

42

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

# Transport Layer

43



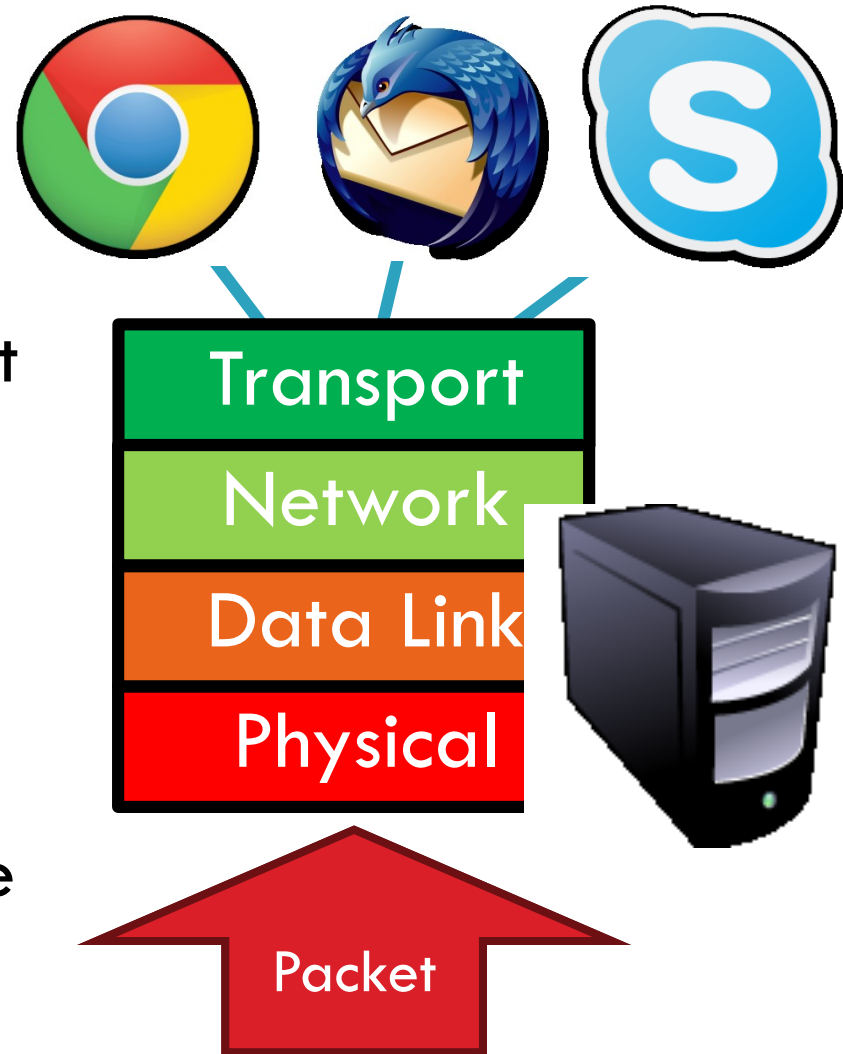
- Function:
  - ▣ Demultiplexing of data streams
- Optional functions:
  - ▣ Creating long lived connections
  - ▣ Reliable, in-order packet delivery
  - ▣ Error detection
  - ▣ Flow and congestion control
- Key challenges:
  - ▣ Detecting and responding to congestion
  - ▣ Balancing fairness against high utilization

- UDP
- TCP
- Congestion Control
- Evolution of TCP
- Problems with TCP

# The Case for Multiplexing

45

- ❑ Datagram network
  - ❑ No circuits
  - ❑ No connections
- ❑ Clients run many applications at the same time
  - ❑ Who to deliver packets to?
- ❑ IP header “protocol” field
  - ❑ 8 bits = 256 concurrent streams
- ❑ Insert Transport Layer to handle demultiplexing



# Demultiplexing Traffic

46

Server applications communicate with multiple clients

Unique port for each application

Applications share the same network

Application

Transport

Network

P1

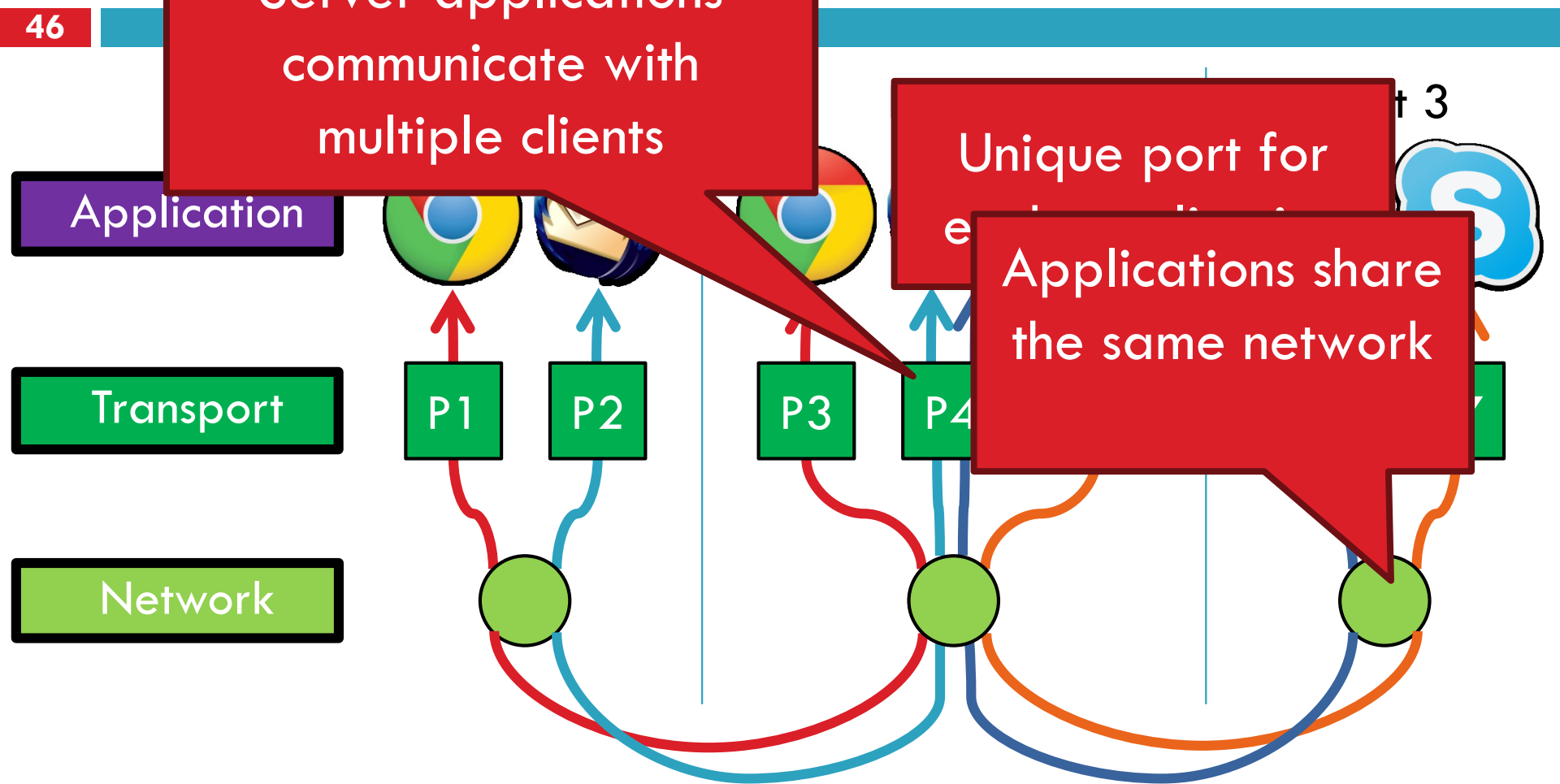
P2

P3

P4

Port 3

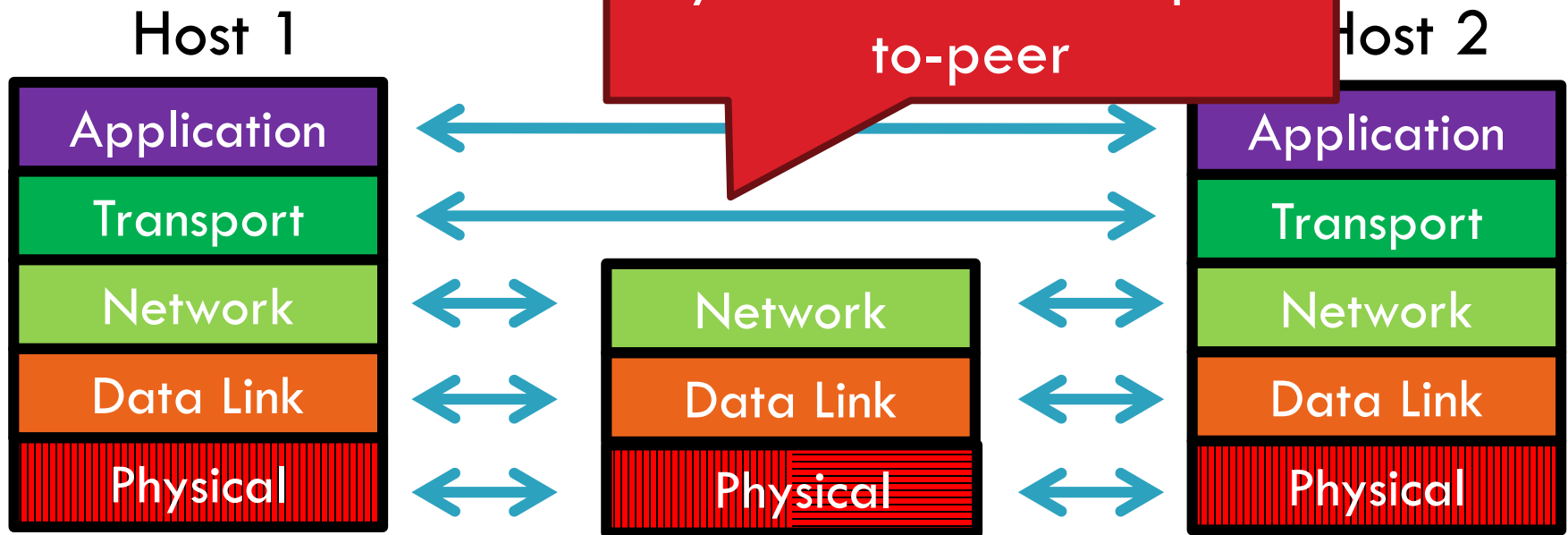
Endpoints identified by  $\langle src\_ip, src\_port, dest\_ip, dest\_port \rangle$



# Layering, Revisited

47

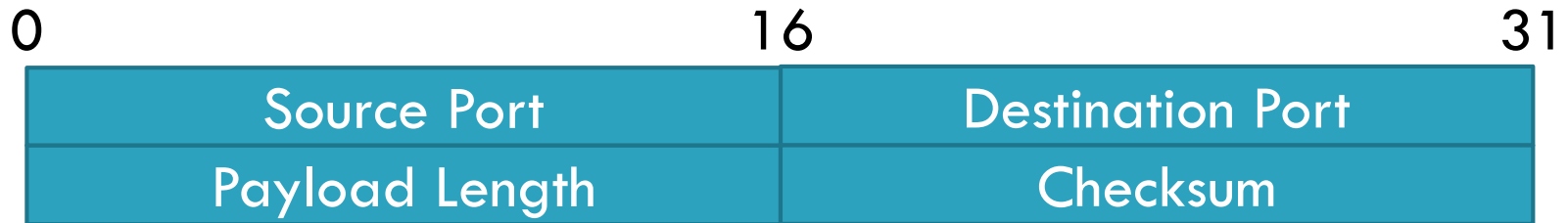
Layers communicate peer-to-peer



- Lowest level end-to-end protocol
  - ▣ Transport header only read by source and destination
  - ▣ Routers view transport header as payload

# User Datagram Protocol (UDP)

48



- ❑ Simple, connectionless datagram
  - ▣ C sockets: `SOCK_DGRAM`
- ❑ Port numbers enable demultiplexing
  - ▣ 16 bits = 65535 possible ports
  - ▣ Port 0 is invalid
- ❑ Checksum for error detection
  - ▣ Detects (some) corrupt packets
  - ▣ Does not detect dropped, duplicated, or reordered packets



# Uses for UDP

49

- ❑ Invented after TCP
  - ❑ Why?
- ❑ Not all applications can tolerate TCP
- ❑ Custom protocols can be built on top of UDP
  - ❑ Reliability? Strict ordering?
  - ❑ Flow control? Congestion control?
- ❑ Examples
  - ❑ RTMP, real-time media streaming (e.g. voice, video)
  - ❑ Facebook datacenter protocol

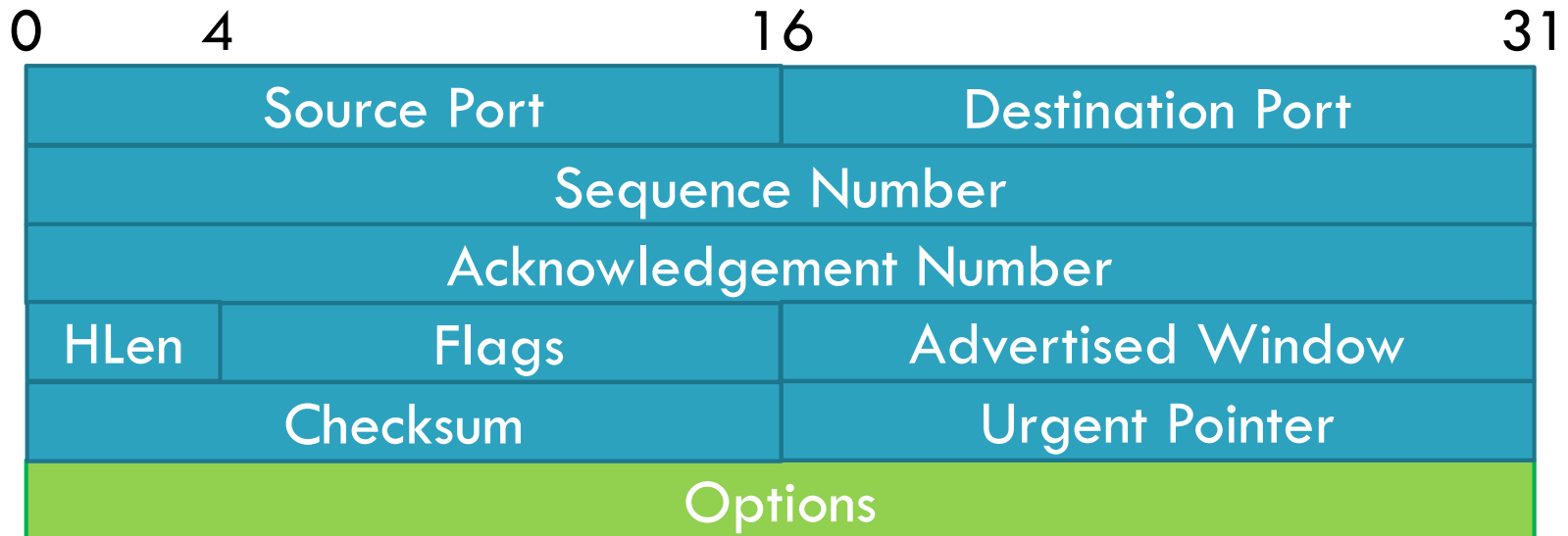
- ❑ UDP – already discussed
- ❑ TCP
- ❑ Congestion Control
- ❑ Evolution of TCP
- ❑ Problems with TCP

# Transmission Control Protocol

51

- Reliable, in-order, bi-directional byte streams
  - ▣ Port numbers for demultiplexing
  - ▣ Virtual circuits (connections)
  - ▣ Flow control
  - ▣ Congestion control, approximate fairness

Why these features?



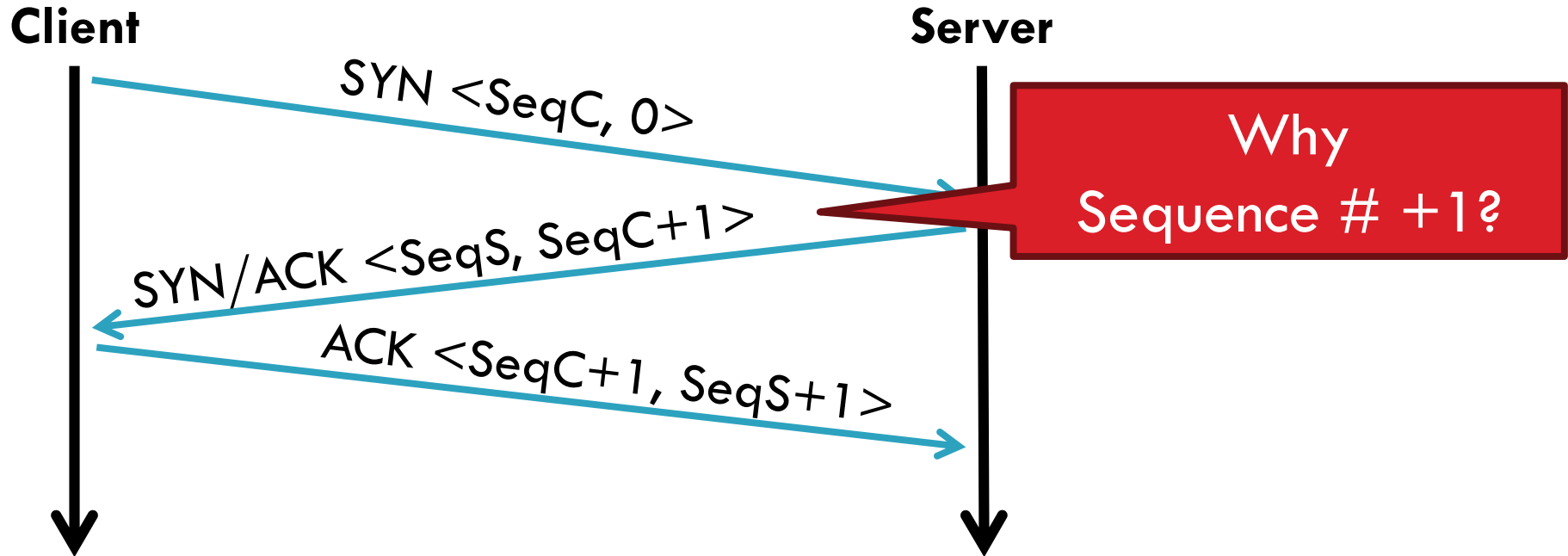
# Connection Setup

52

- Why do we need connection setup?
  - ▣ To establish state on both hosts
  - ▣ Most important state: sequence numbers
    - Count the number of bytes that have been sent
    - Initial value chosen at random
    - Why?
- Important TCP flags (1 bit each)
  - ▣ SYN – synchronization, used for connection setup
  - ▣ ACK – acknowledge received data
  - ▣ FIN – finish, used to tear down connection

# Three Way Handshake

53

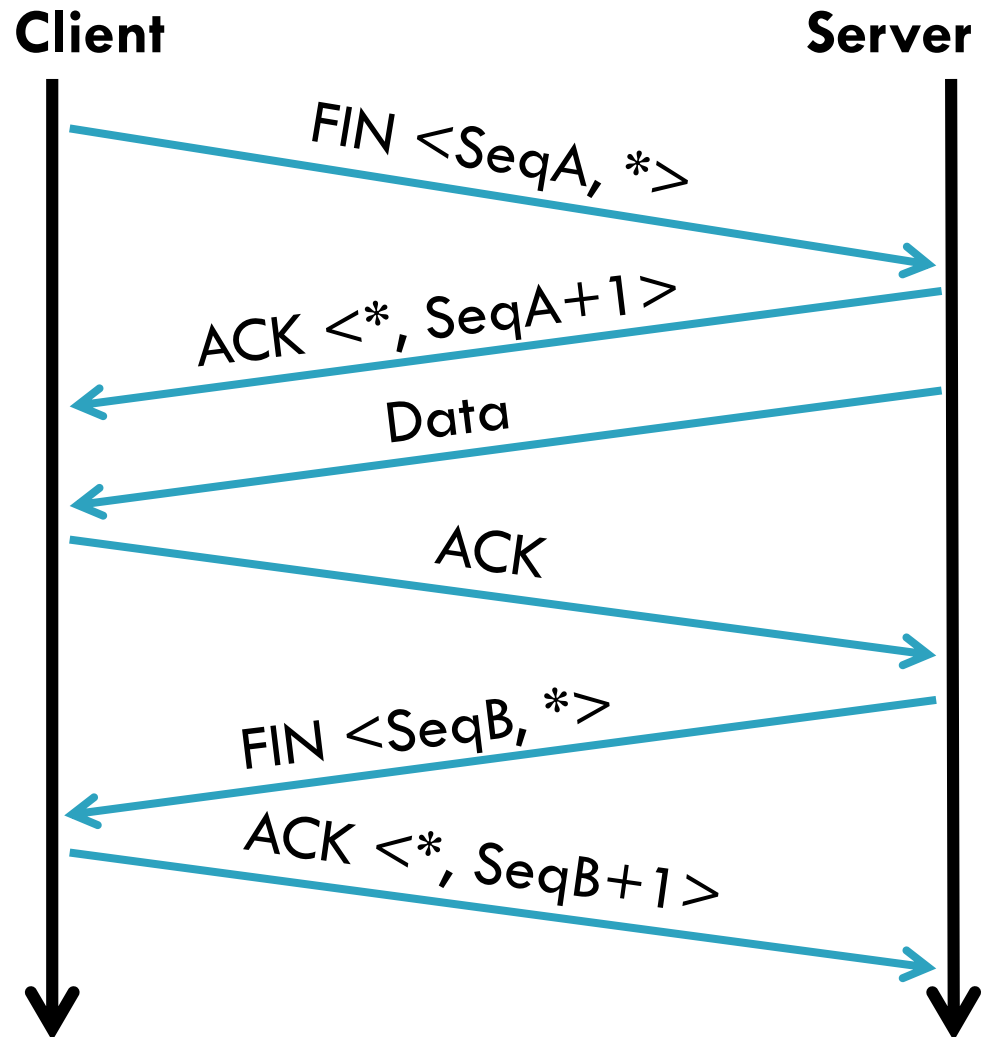


- Each side:
  - ▣ Notifies the other of starting sequence number
  - ▣ ACKs the other side's starting sequence number

# Connection Tear Down

55

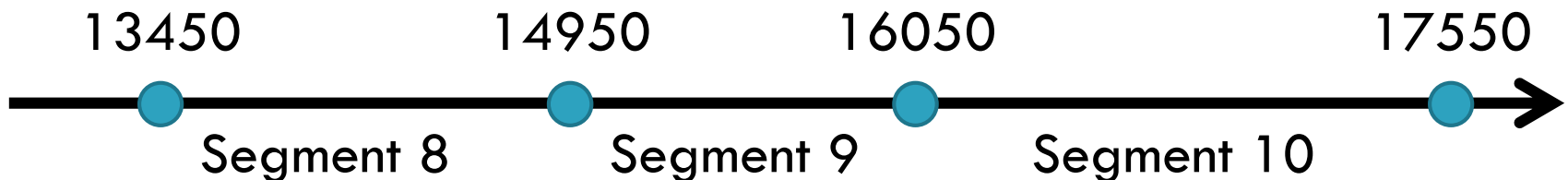
- Either side can initiate tear down
- Other side may continue sending data
  - ▣ Half open connection
  - ▣ `shutdown()`
- Acknowledge the last FIN
  - ▣ Sequence number + 1
- What happens if 2<sup>nd</sup> FIN is lost?



# Sequence Number Space

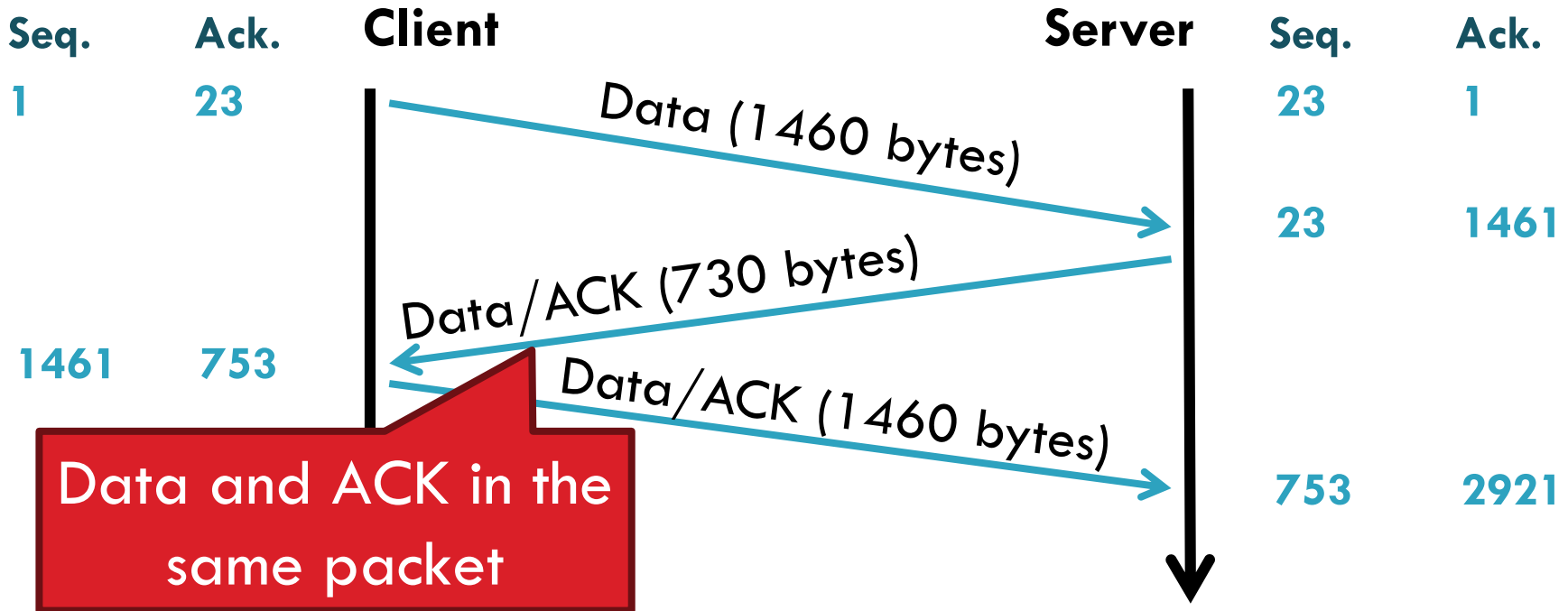
56

- TCP uses a byte stream abstraction
  - ▣ Each byte in each stream is numbered
  - ▣ 32-bit value, wraps around
  - ▣ Initial, random values selected during setup. Why?
- Byte stream broken down into segments (packets)
  - ▣ Size limited by the Maximum Segment Size (MSS)
  - ▣ Set to limit fragmentation
- Each segment has a sequence number



# Bidirectional Communication

57



- Each side of the connection can send and receive
  - ▣ Different sequence numbers for each direction



# Flow Control

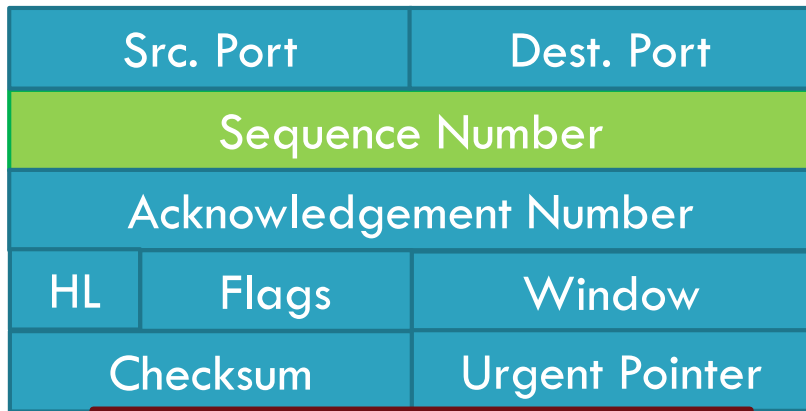
58

- Problem: how many packets should a sender transmit?
  - ▣ Too many packets may overwhelm the receiver
  - ▣ Size of the receivers buffers may change over time
- Solution: sliding window
  - ▣ Receiver tells the sender how big their buffer is
  - ▣ Called the **advertised window**
  - ▣ For window size  $n$ , sender may transmit  $n$  bytes without receiving an ACK
  - ▣ After each ACK, the window slides forward
- Window may go to zero!

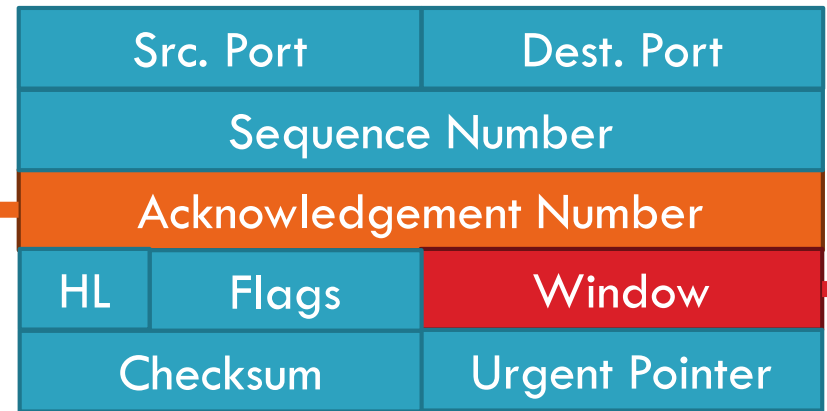
# Flow Control: Sender Side

59

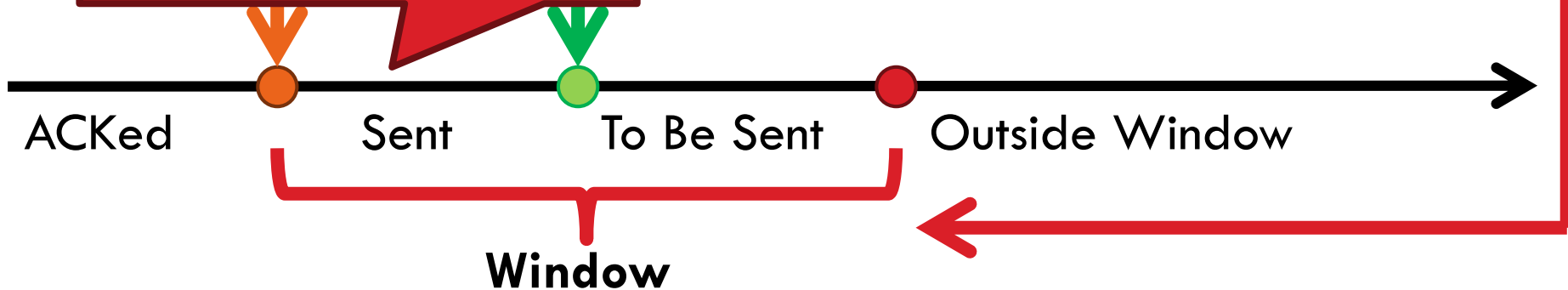
Packet Sent



Packet Received

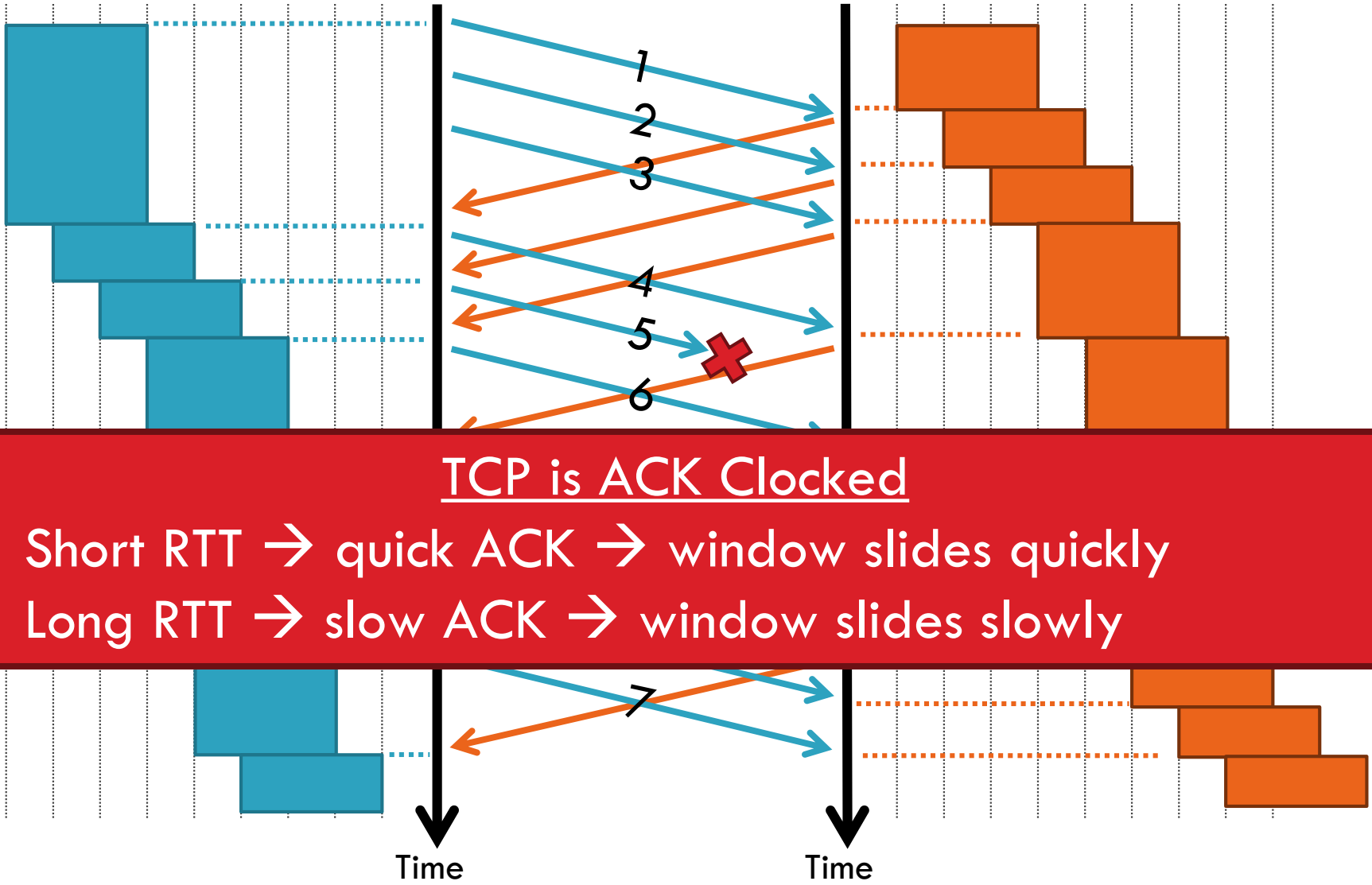


Must be buffered until ACKed



# Sliding Window Example

60



# Observations

61

- Throughput is  $\sim w/\text{RTT}$
- Sender has to buffer all unacknowledged packets, because they may require retransmission
- Receiver may be able to accept out-of-order packets, but only up to buffer limits

# What Should the Receiver ACK?

62

1. ACK every packet
2. Use *cumulative ACK*, where an ACK for sequence  $n$  implies ACKS for all  $k < n$
3. Use *negative ACKs* (NACKs), indicating which packet did not arrive
4. Use *selective ACKs* (SACKs), indicating those that did arrive, even if not in order
  - ▣ SACK is an actual TCP extension

# Sequence Numbers, Revisited

63

- 32 bits, unsigned
  - ▣ Why so big?
- For the sliding window you need...
  - ▣  $|\text{Sequence \# Space}| > 2 * |\text{Sending Window Size}|$
  - ▣  $2^{32} > 2 * 2^{16}$
- Guard against stray packets
  - ▣ IP packets have a maximum segment lifetime (MSL) of 120 seconds
    - i.e. a packet can linger in the network for 2 minutes

# Silly Window Syndrome

64

□ Problem: what if the window size is very small?

□ Multiple, small packets, headers dominate data





□ Equivalent problem: sender transmits packets one byte at a time

1. `for (int x = 0; x < strlen(data); ++x)`
2. `write(socket, data + x, 1);`

# Nagle's Algorithm

65

1. If the window  $\geq$  MSS and available data  $\geq$  MSS:  
Send the data  Send a full packet
  2. Elif there is unACKed data:  
Enqueue data in a buffer until an ACK is received
  3. Else: send the data  Send a non-full packet if nothing else is happening
- Problem: Nagle's Algorithm delays transmissions
- ▣ What if you need to send a packet immediately?
    1. `int flag = 1;`
    2. `setsockopt(sock, IPPROTO_TCP, TCP_NODELAY, (char *) &flag, sizeof(int));`