

# CS 3700

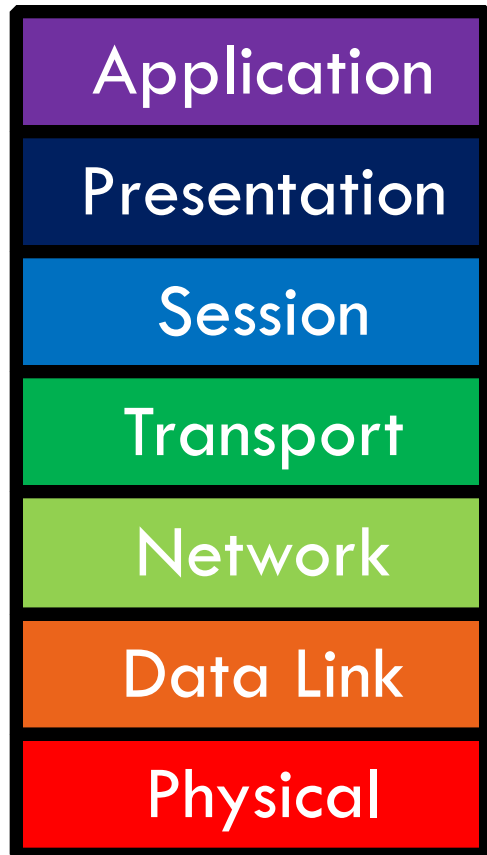
## Networks and Distributed Systems

### Lecture 9: UDP/TCP

Revised 2/9/2014

# Transport Layer

2



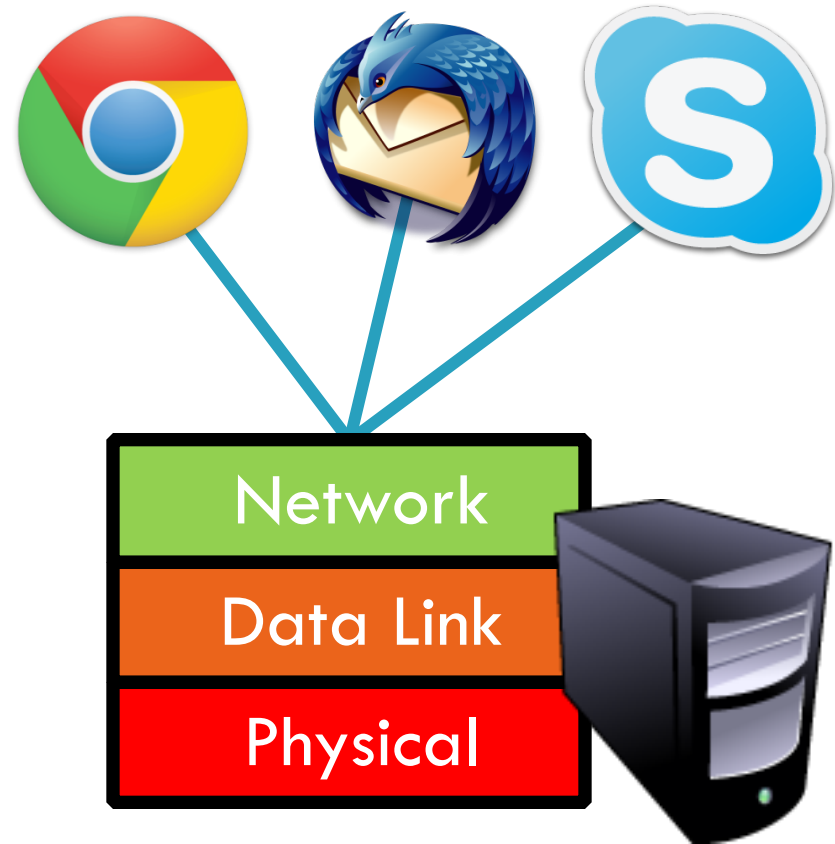
- Function:
  - ▣ Demultiplexing of data streams
- Optional functions:
  - ▣ Creating long lived connections
  - ▣ Reliable, in-order packet delivery
  - ▣ Error detection
  - ▣ Flow and congestion control
- Key challenges:
  - ▣ Detecting and responding to congestion
  - ▣ Balancing fairness against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

- ❑ UDP
- ❑ TCP

# The Case for Multiplexing

4

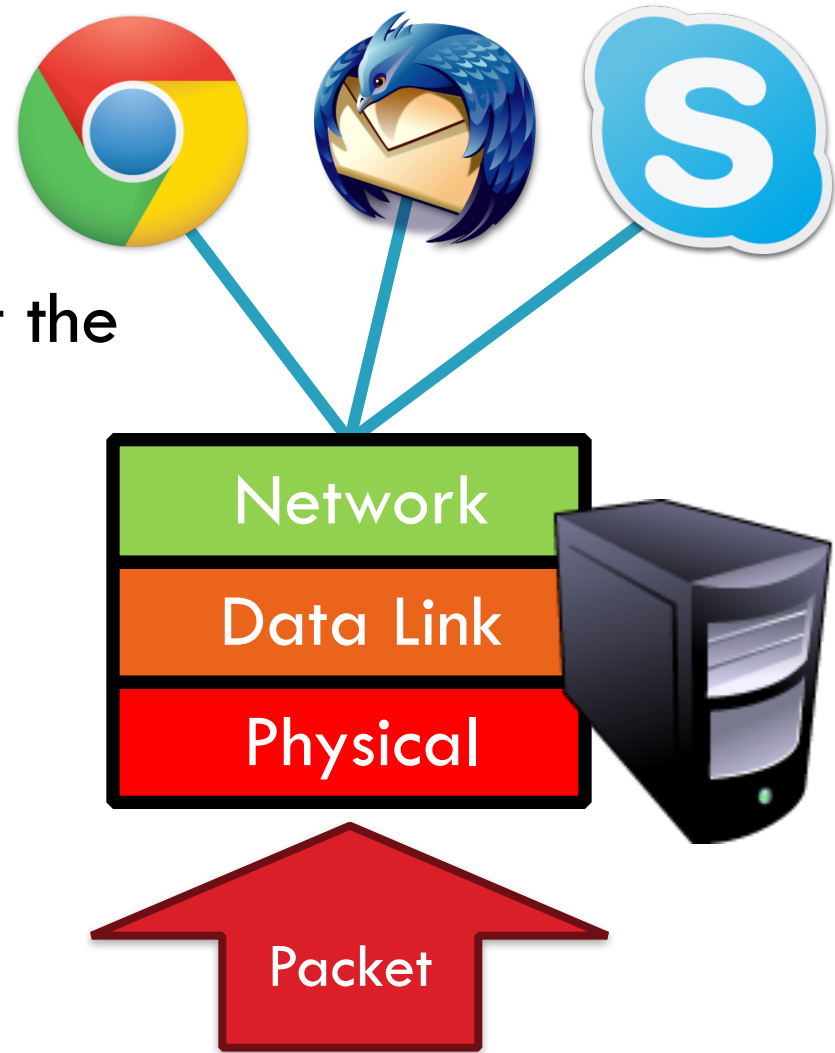
- Datagram network
  - No circuits
  - No connections



# The Case for Multiplexing

4

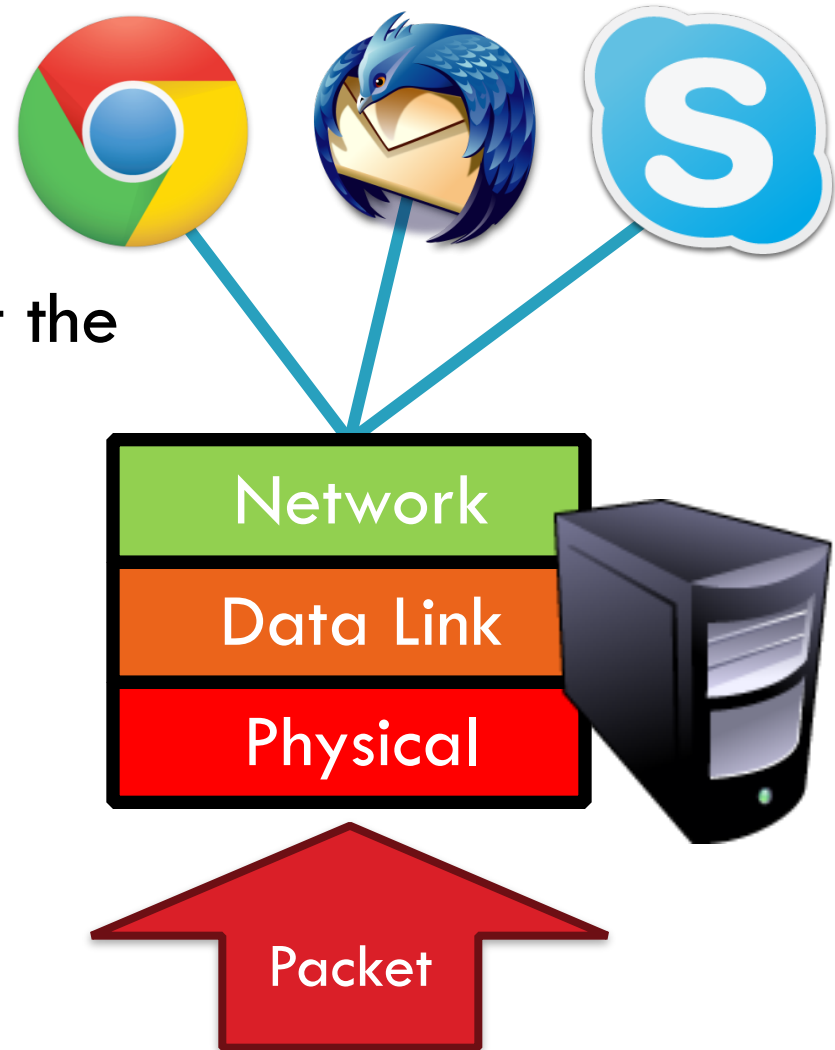
- Datagram network
  - ▣ No circuits
  - ▣ No connections
- Clients run many applications at the same time
  - ▣ Who to deliver packets to?



# The Case for Multiplexing

4

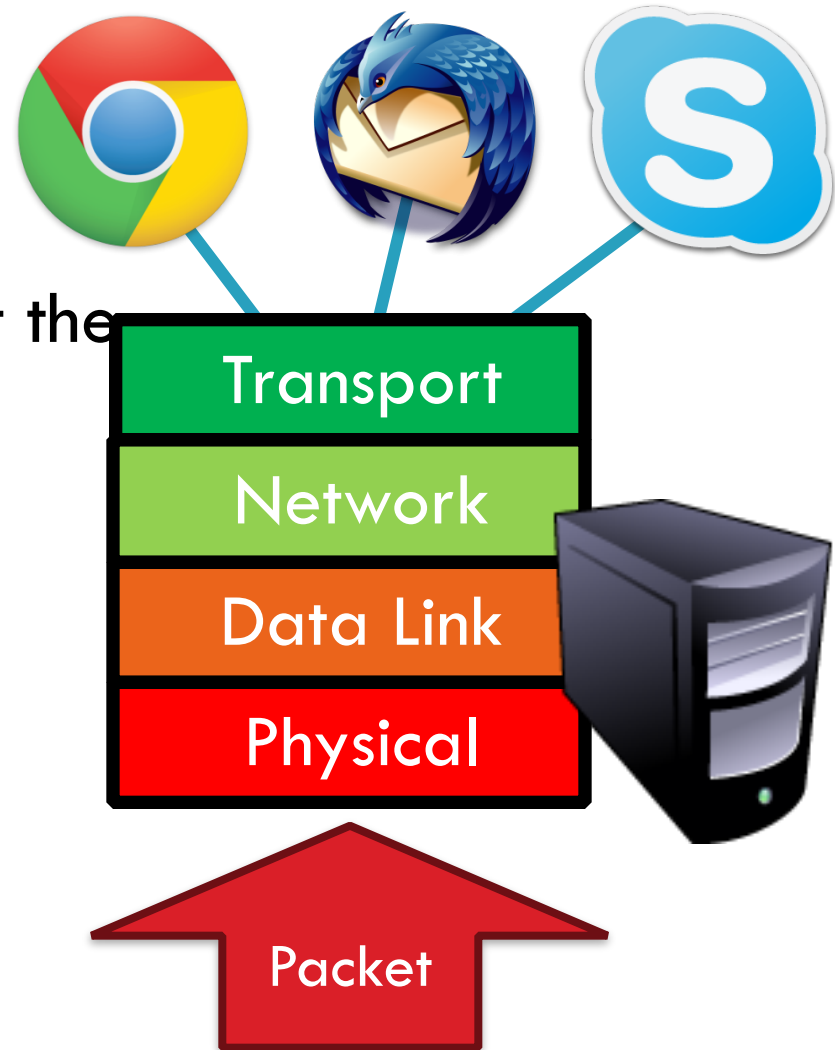
- Datagram network
  - ▣ No circuits
  - ▣ No connections
- Clients run many applications at the same time
  - ▣ Who to deliver packets to?
- IP header “protocol” field
  - ▣ 8 bits = 256 concurrent streams



# The Case for Multiplexing

4

- 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

5

Application

Host 1



Host 2



Host 3





# Demultiplexing Traffic

5

Host 1

Application



P1

P2

Transport



P3

P4

P5

Unique port for each application

Host 3

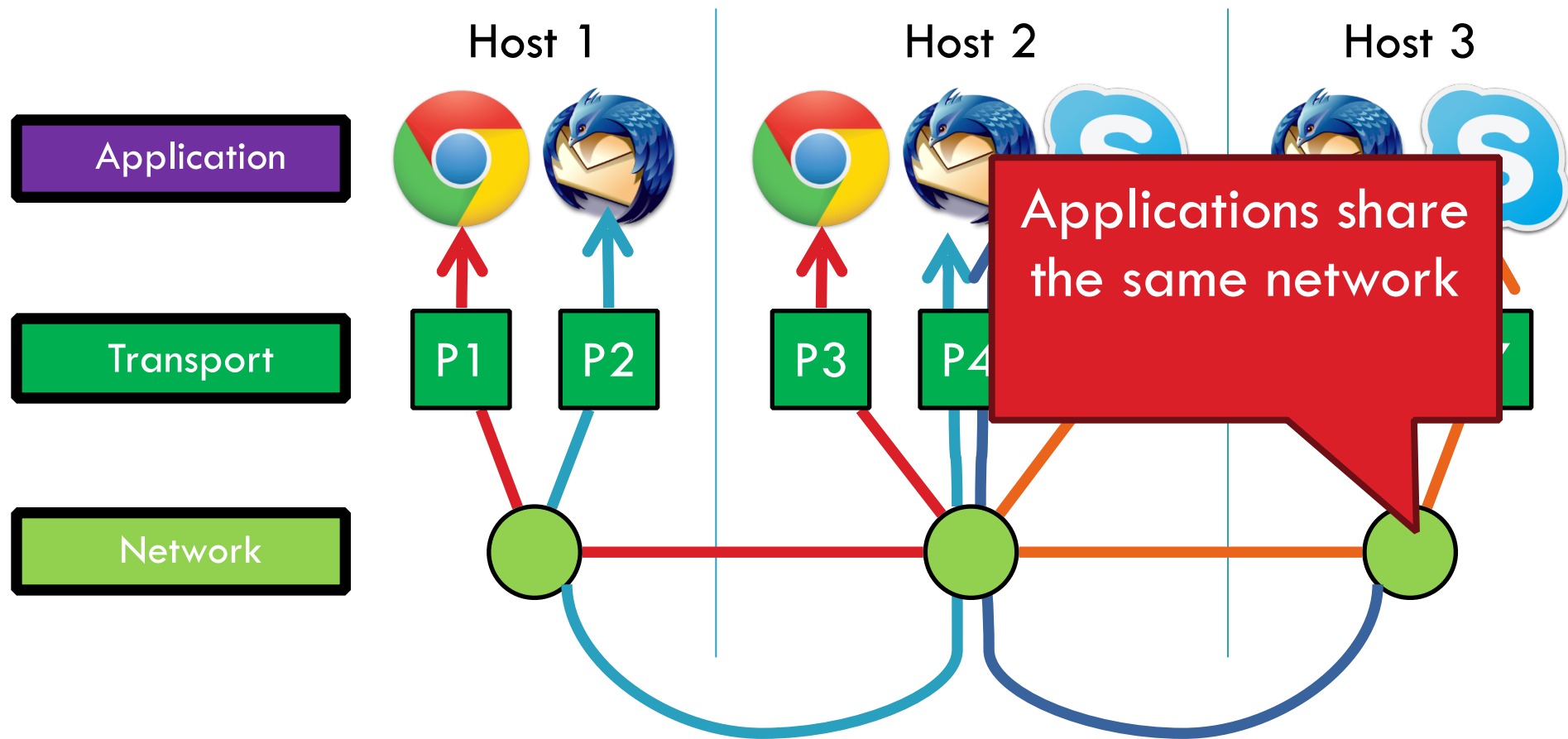


P6

P7

# Demultiplexing Traffic

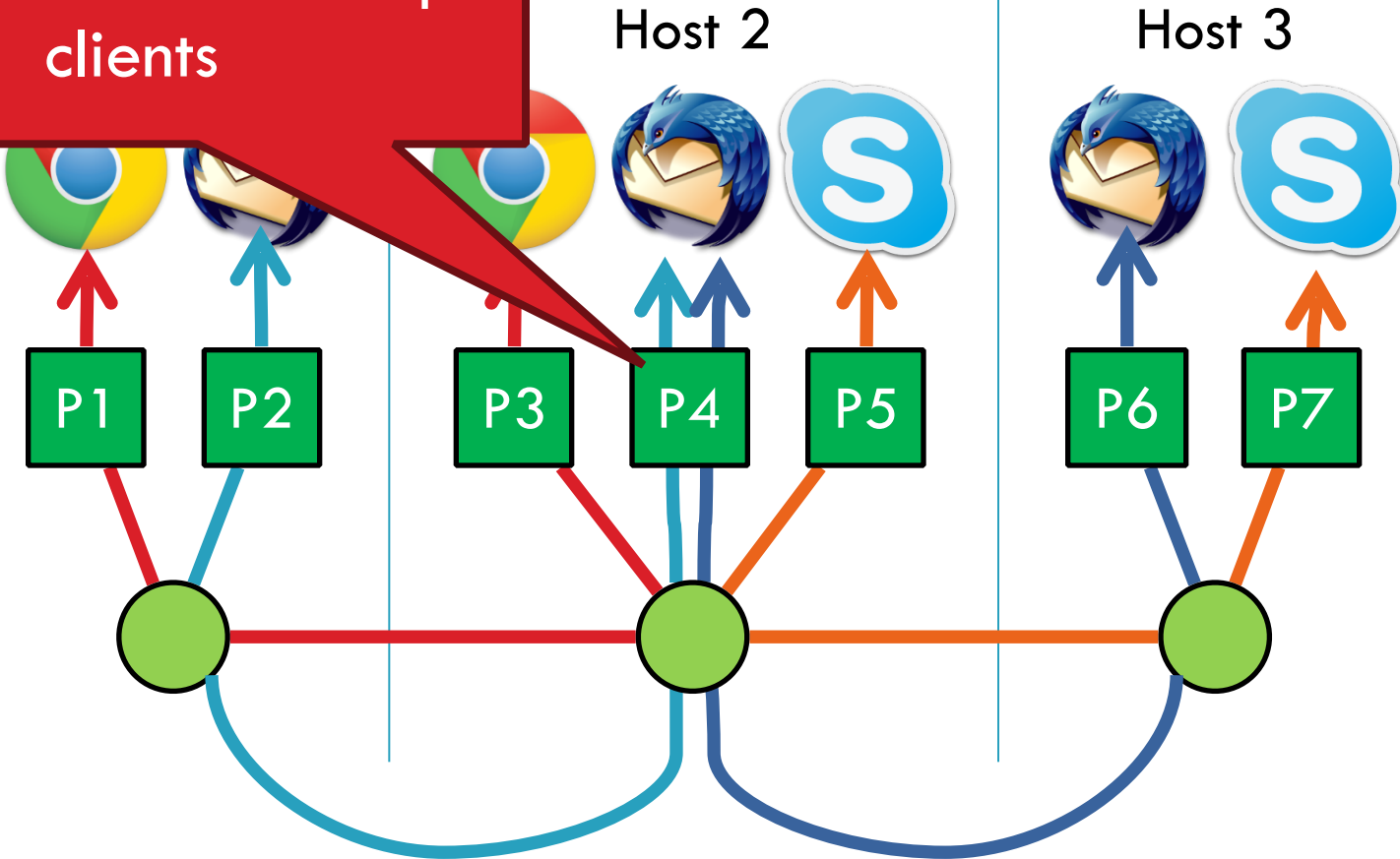
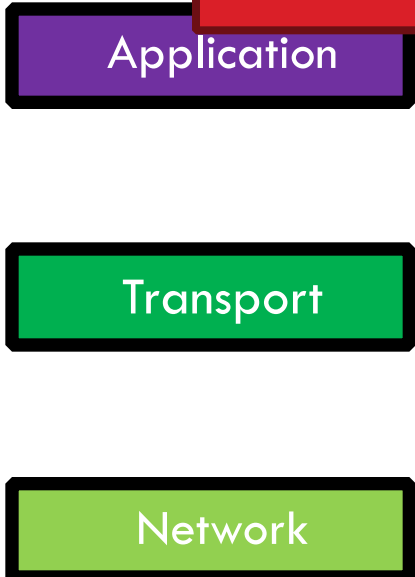
5



# Demultiplexing Traffic

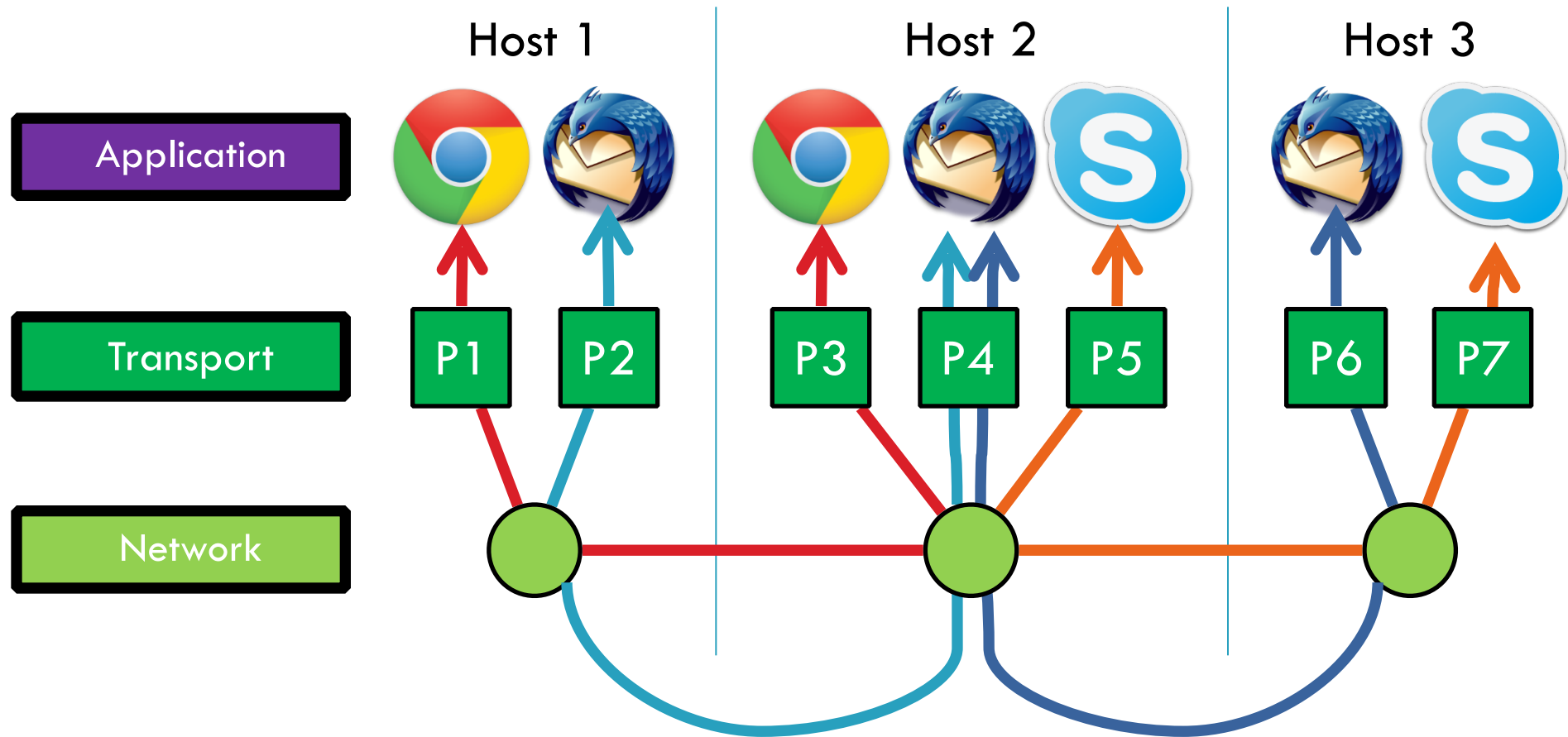
5

Server applications communicate with multiple clients



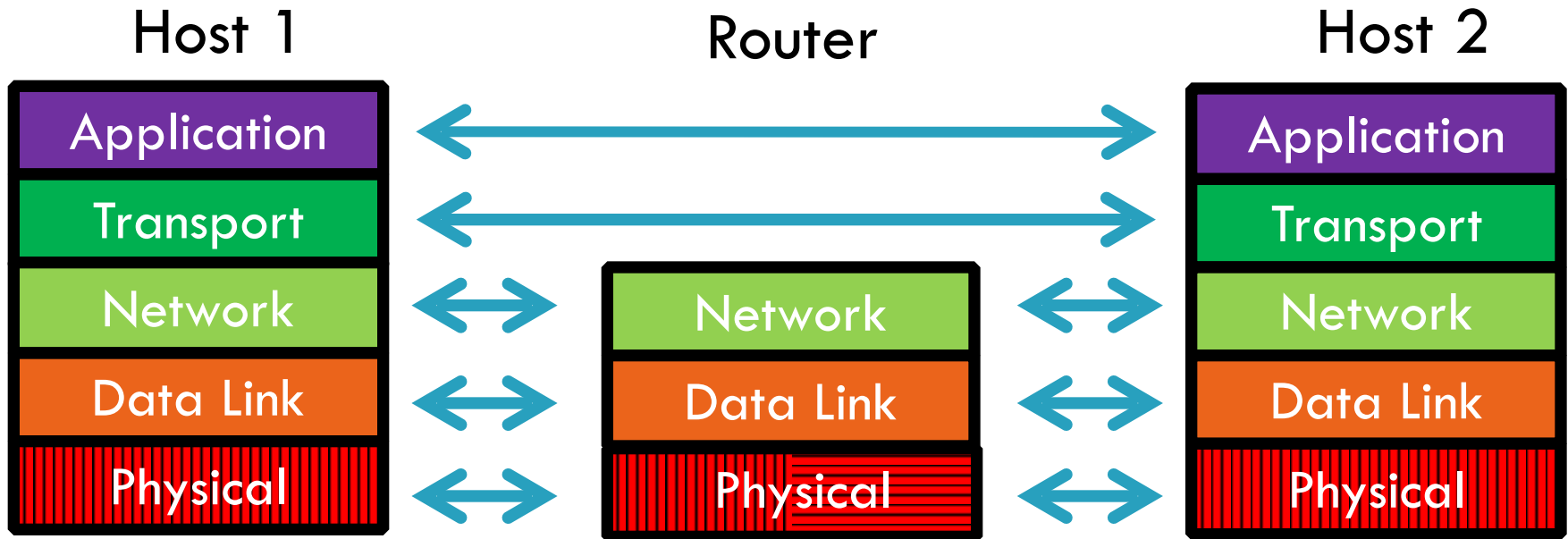
# Demultiplexing Traffic

5



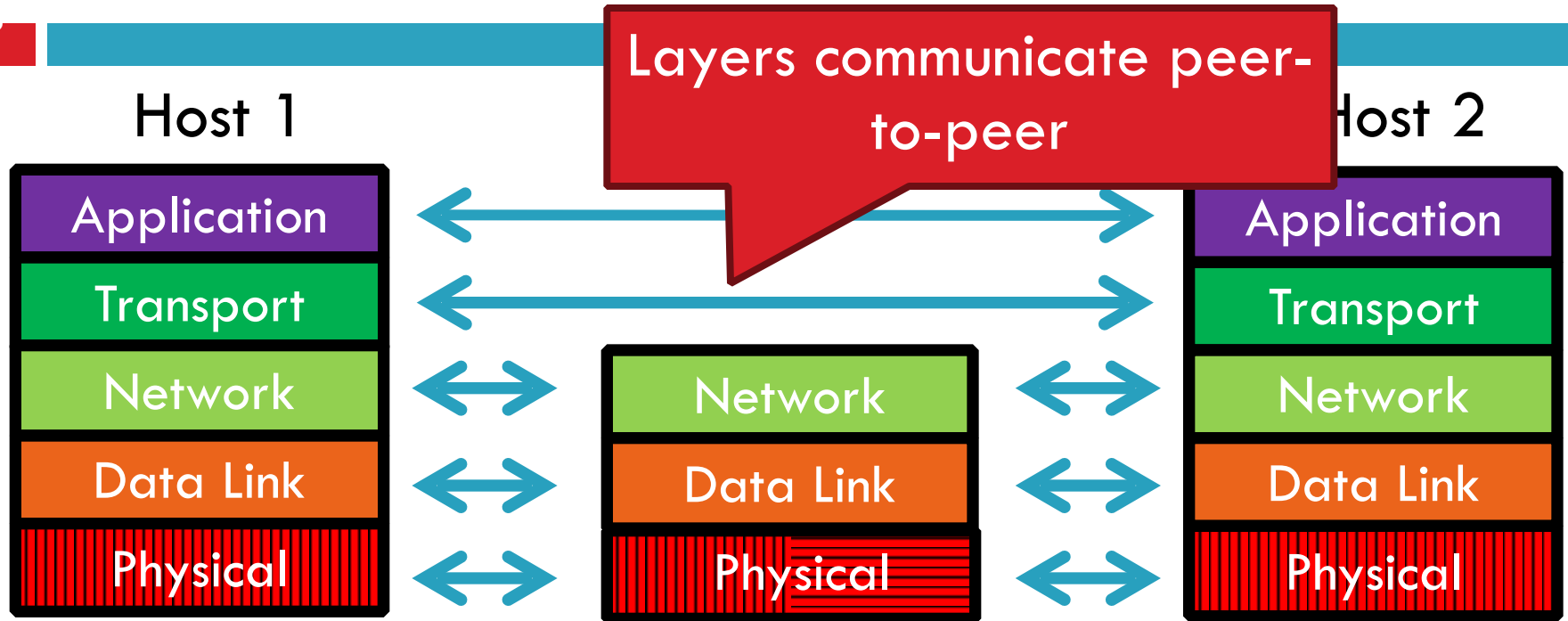
Endpoints identified by  $\langle src\_ip, src\_port, dest\_ip, dest\_port \rangle$

# Layering, Revisited



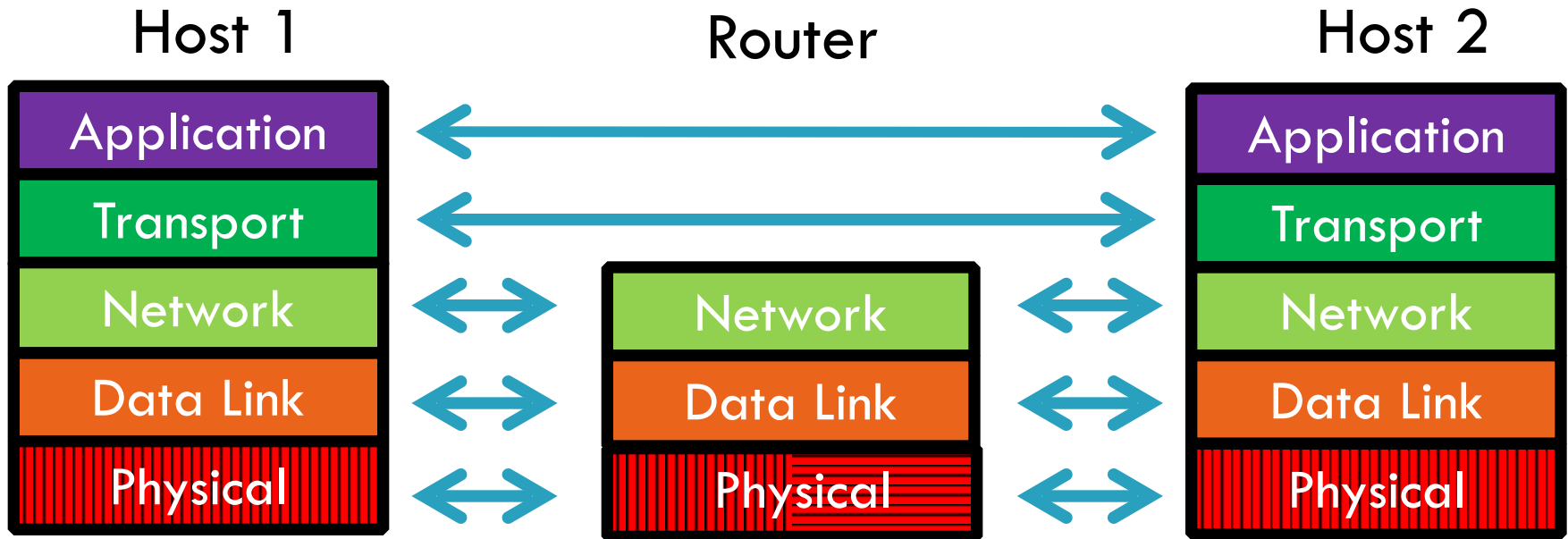
- Lowest level end-to-end protocol (in theory)
  - ▣ Transport header only read by source and destination
  - ▣ Routers view transport header as payload

# Layering, Revisited



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

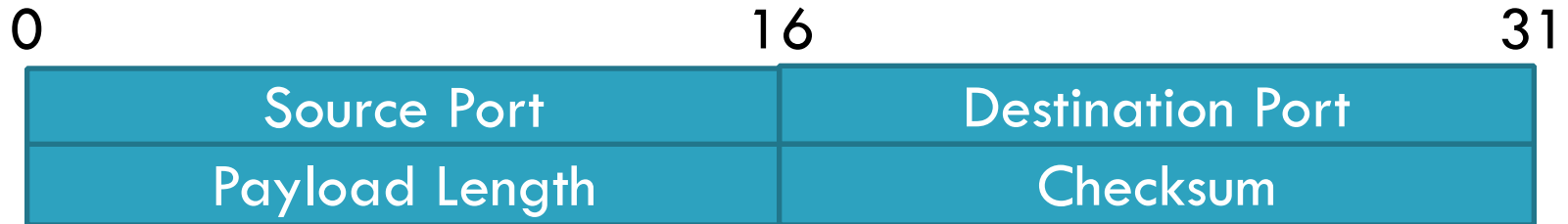
# Layering, Revisited



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

# User Datagram Protocol (UDP)

7



- 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

8

- Invented after TCP
  - Why?

# Uses for UDP

8

- Invented after TCP
  - Why?
- Not all applications can tolerate TCP

# Uses for UDP

8

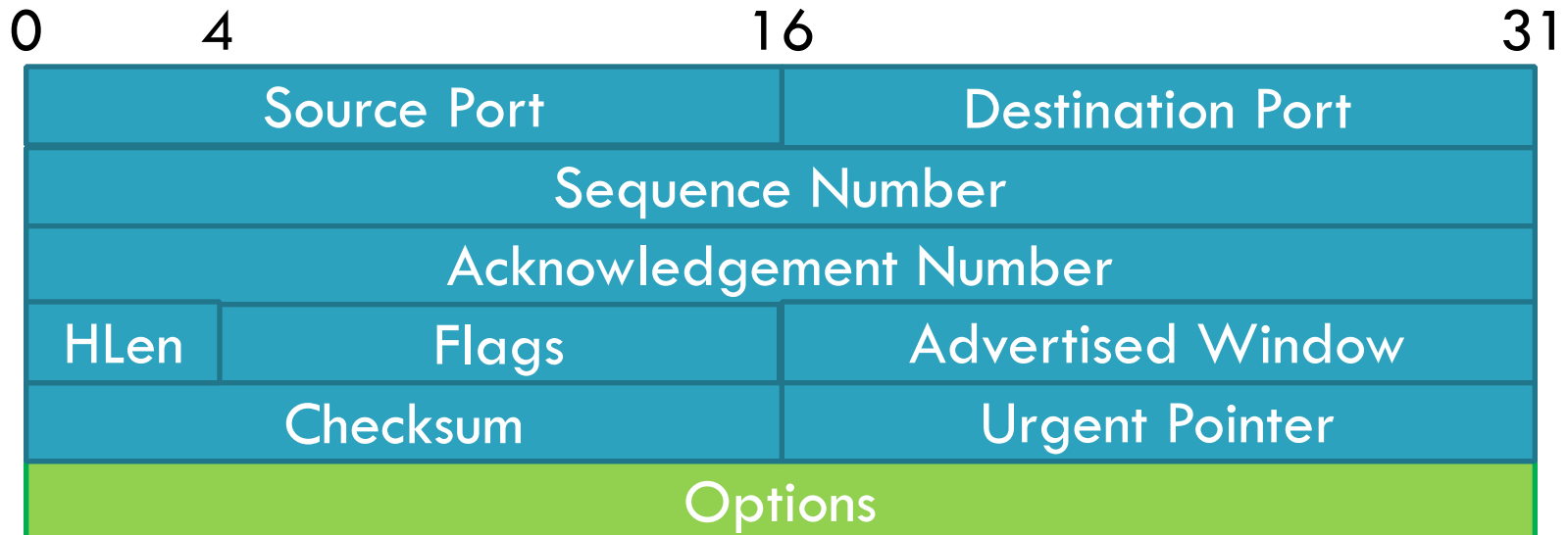
- Invented after TCP
  - Why?
- Not all applications can tolerate TCP
- Custom protocols can be built on top of UDP
  - Reliability? Strict ordering?
  - Flow control? Congestion control?
- Examples
  - RTMP, real-time media streaming (e.g. voice, video)
  - Facebook datacenter protocol

- ❑ UDP
- ❑ TCP

# Transmission Control Protocol

10

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

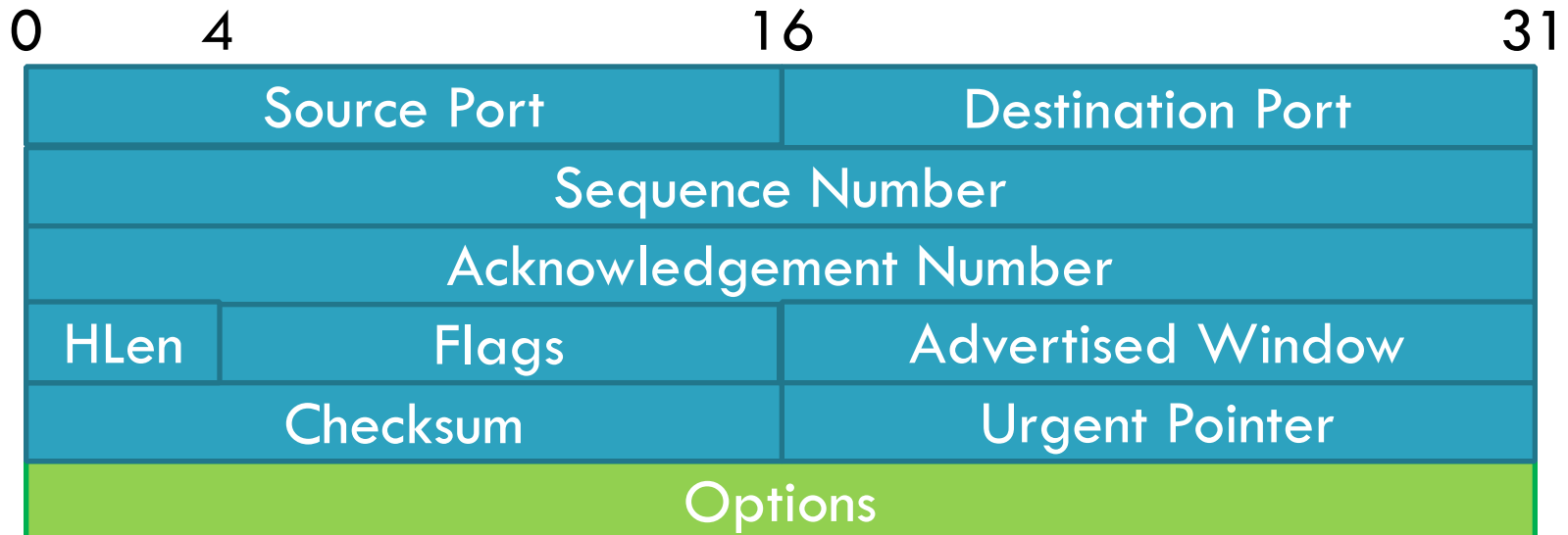


# Transmission Control Protocol

10

- 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

11

- Why do we need connection setup?

# Connection Setup

11

- 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?



# Connection Setup

11

- 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

12

**Client**



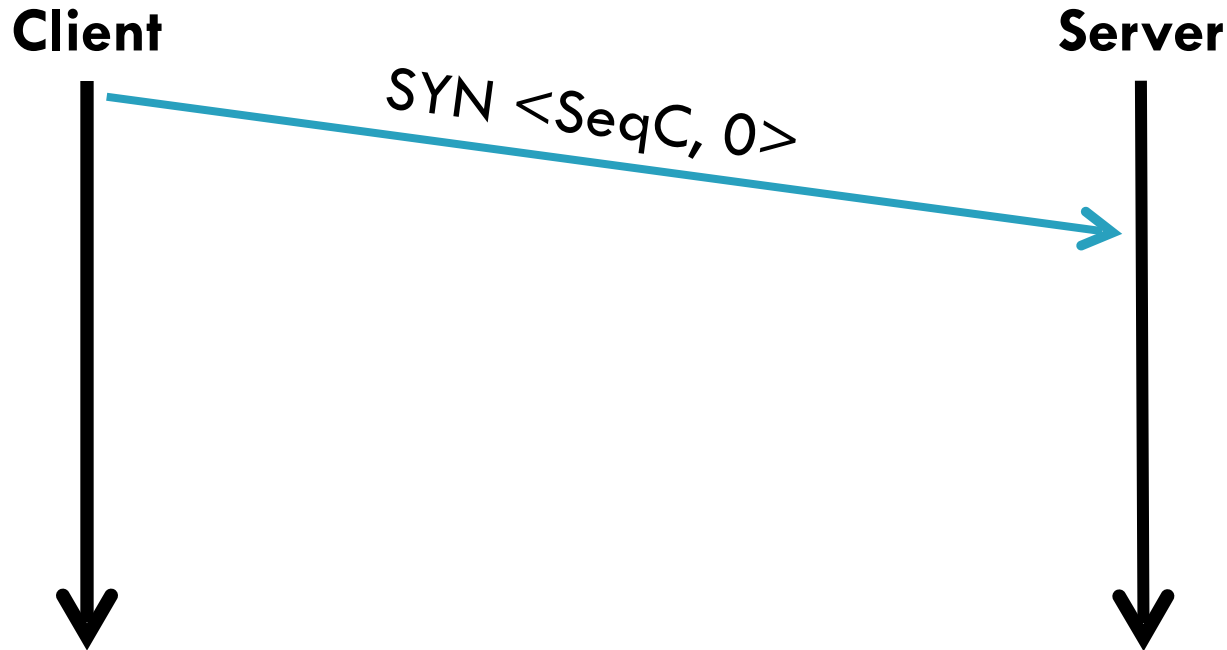
**Server**



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

# Three Way Handshake

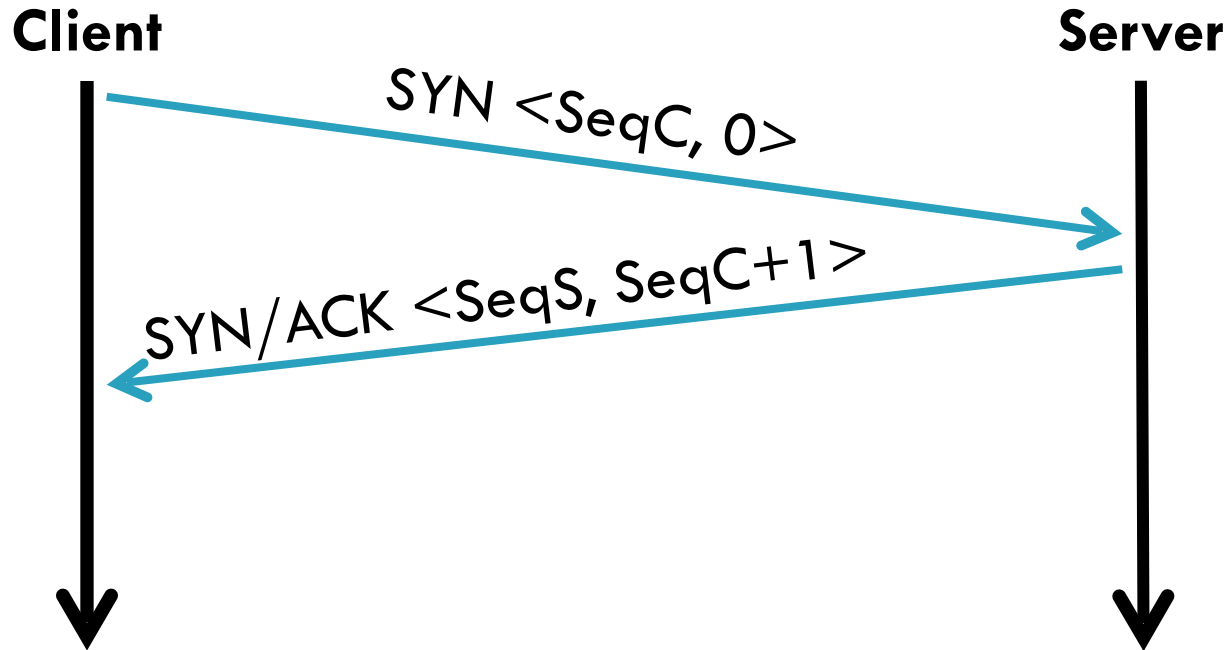
12



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

# Three Way Handshake

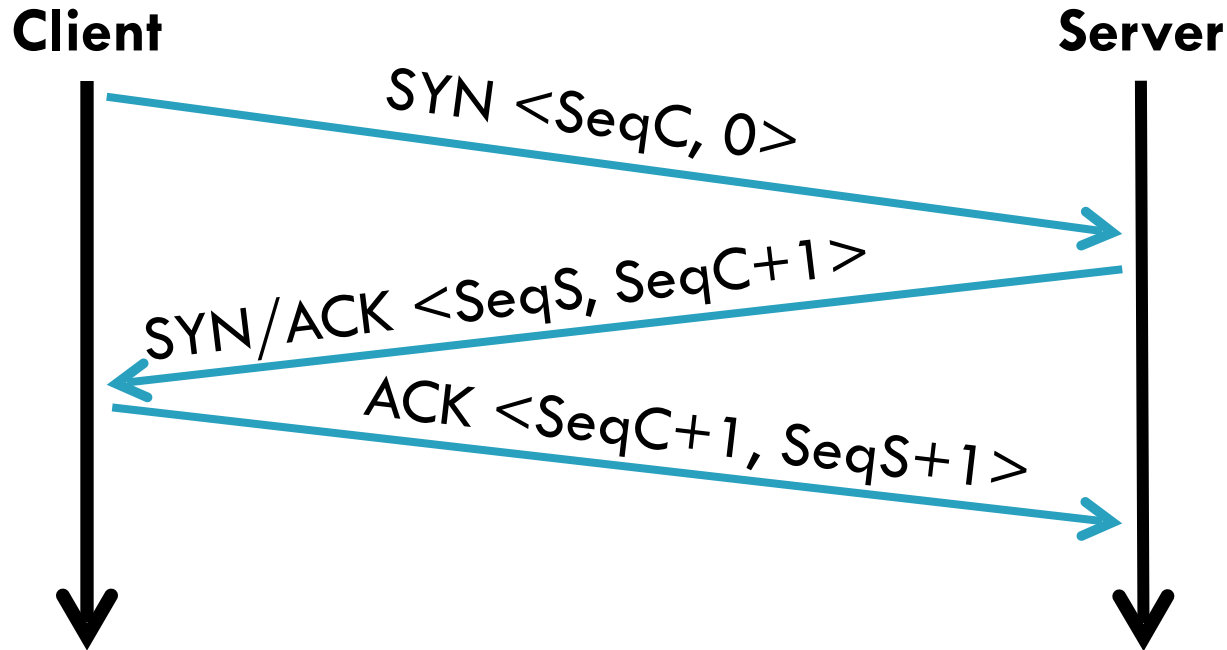
12



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

# Three Way Handshake

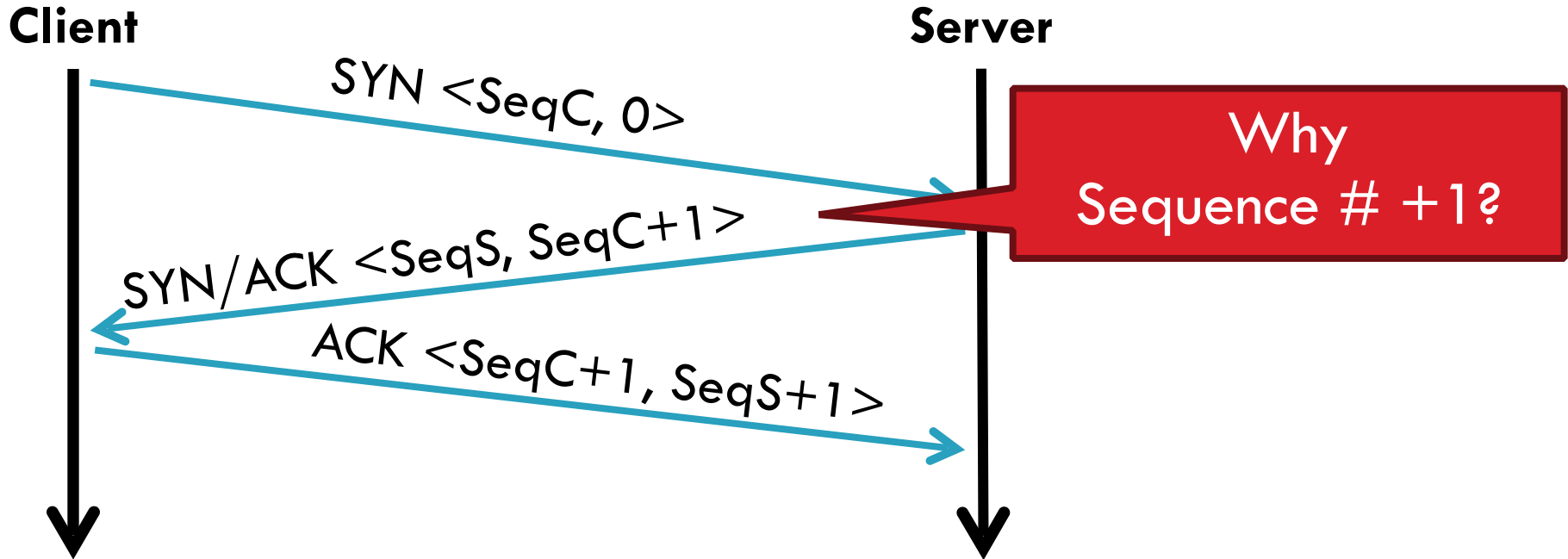
12



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

# Three Way Handshake

12



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

# Connection Setup Issues

13

- Connection confusion
  - How to disambiguate connections from the same host?
  - Random sequence numbers

# Connection Setup Issues

13

- Connection confusion
  - How to disambiguate connections from the same host?
  - Random sequence numbers
- Source spoofing
  - Kevin Mitnick
  - Need good random number generators!



# Connection Setup Issues

13

- Connection confusion
  - ▣ How to disambiguate connections from the same host?
  - ▣ Random sequence numbers
- Source spoofing
  - ▣ Kevin Mitnick
  - ▣ Need good random number generators!
- Connection state management
  - ▣ Each SYN allocates state on the server
  - ▣ SYN flood = denial of service attack
  - ▣ Solution: SYN cookies

# Connection Tear Down

14

- Either side can initiate tear down

**Client**



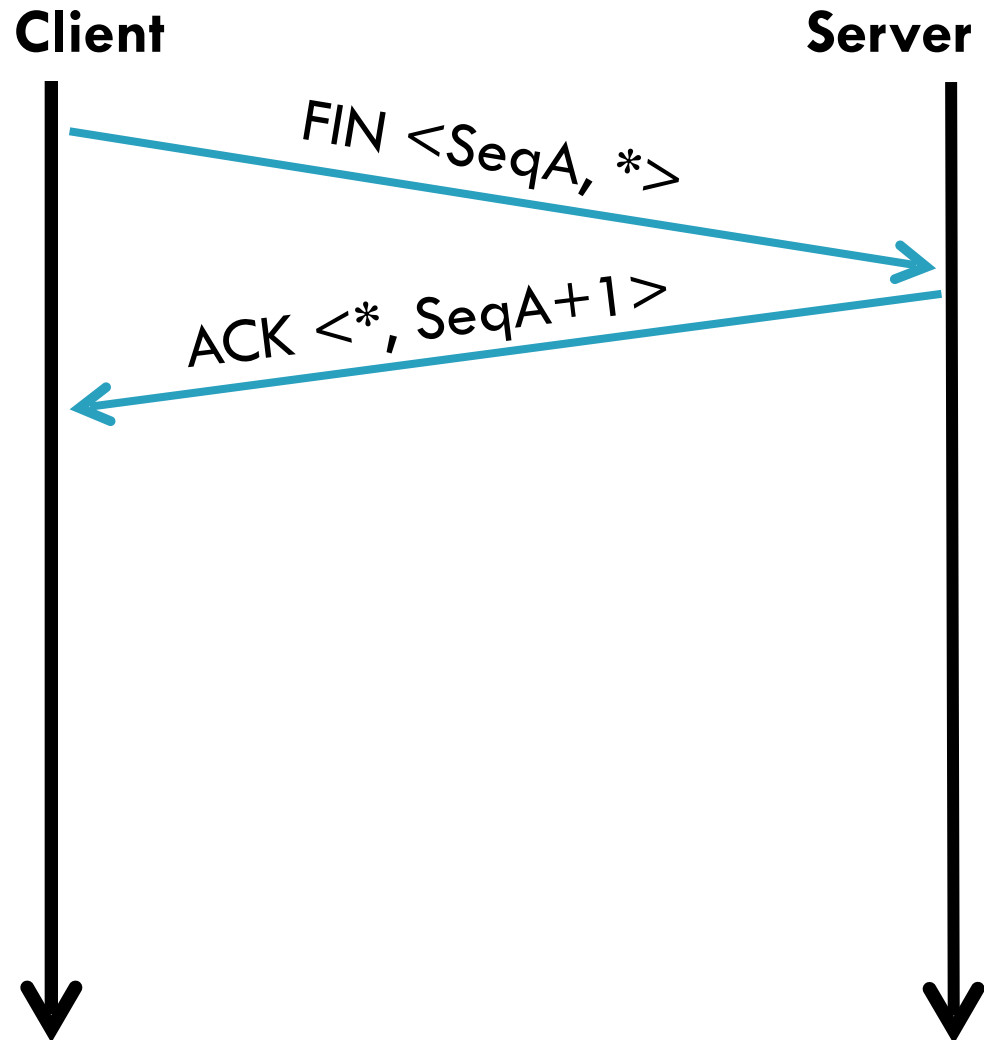
**Server**



# Connection Tear Down

14

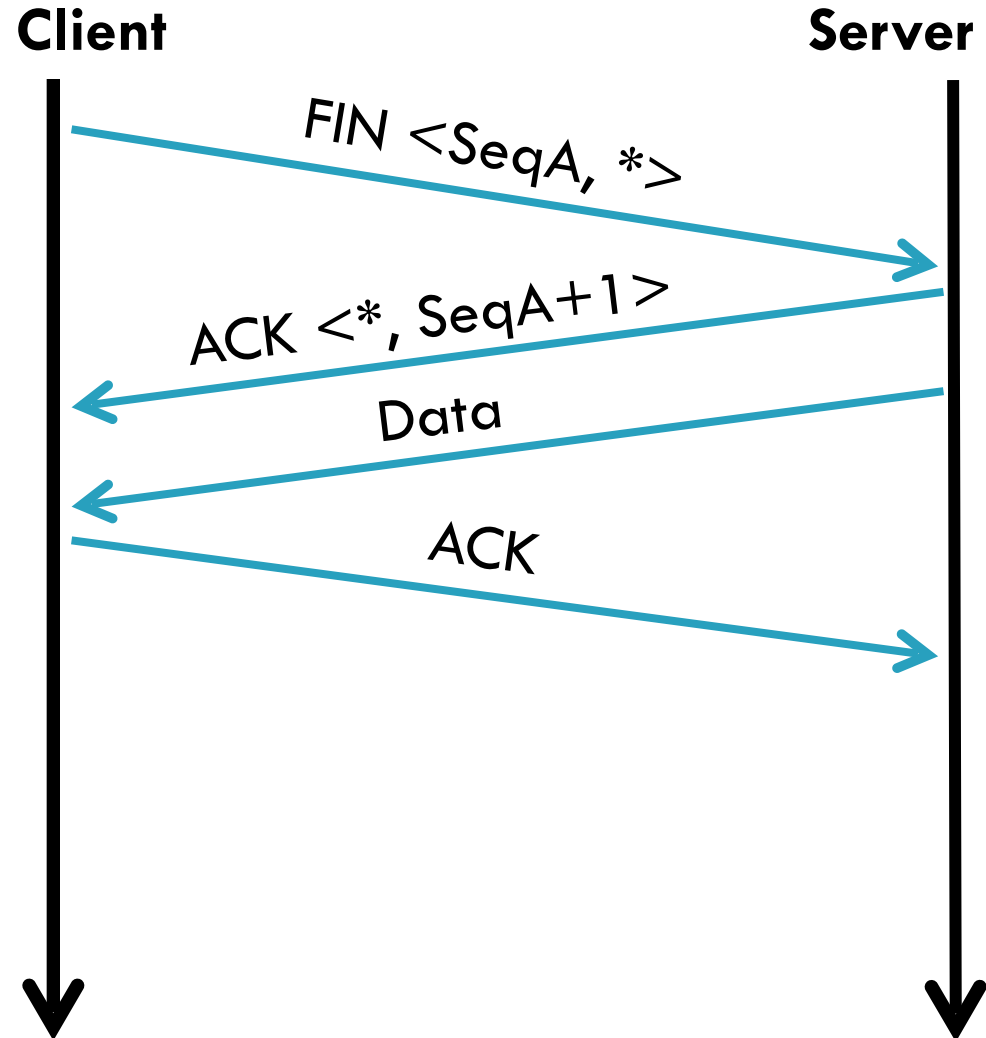
- Either side can initiate tear down



# Connection Tear Down

14

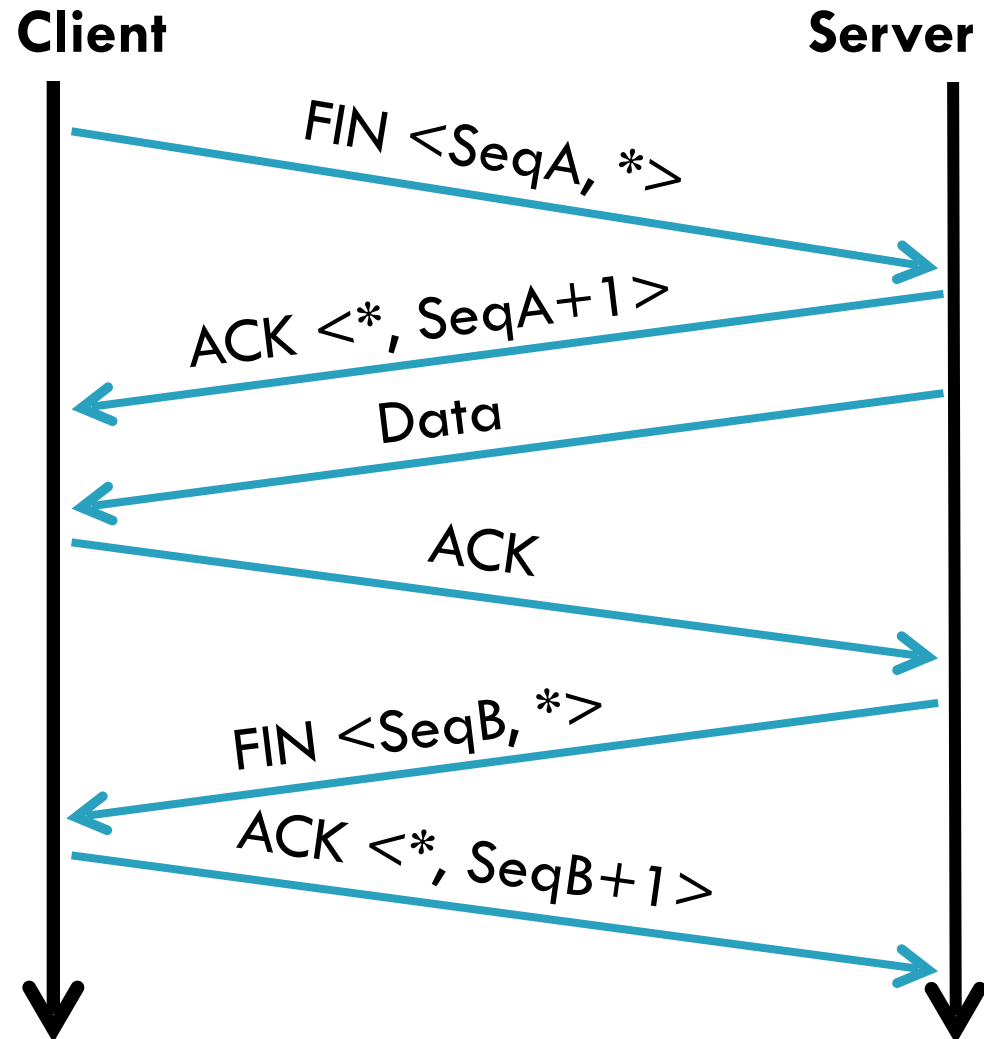
- Either side can initiate tear down
- Other side may continue sending data
  - ▣ Half open connection
  - ▣ *shutdown()*



# Connection Tear Down

14

- Either side can initiate tear down
- Other side may continue sending data
  - ▣ Half open connection
  - ▣ *shutdown()*
- Acknowledge the last FIN
  - ▣ Sequence number + 1



# Sequence Number Space

15

- TCP uses a byte stream abstraction
  - Each byte in each stream is numbered
  - 32-bit value, wraps around
  - Initial, random values selected during setup

# Sequence Number Space

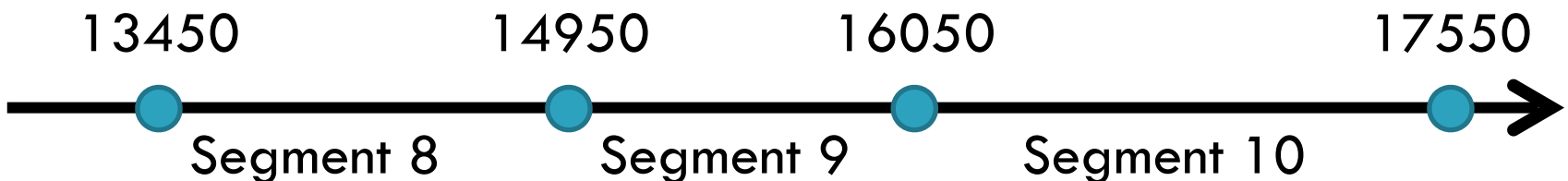
15

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

# Sequence Number Space

15

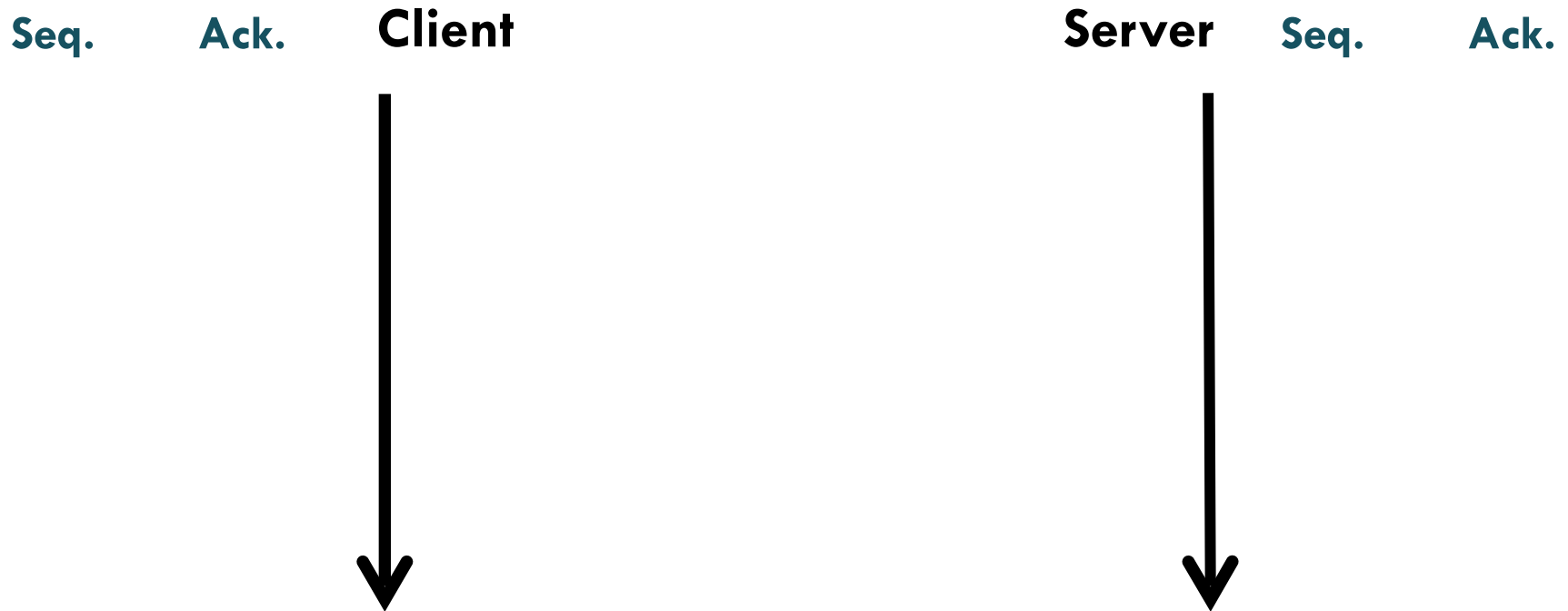
- TCP uses a byte stream abstraction
  - ▣ Each byte in each stream is numbered
  - ▣ 32-bit value, wraps around
  - ▣ Initial, random values selected during setup
- Byte stream broken down into segments (packets)
  - ▣ Size limited by the Maximum Segment Size (MSS)
  - ▣ Set to limit fragmentation
-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





# Bidirectional Communication

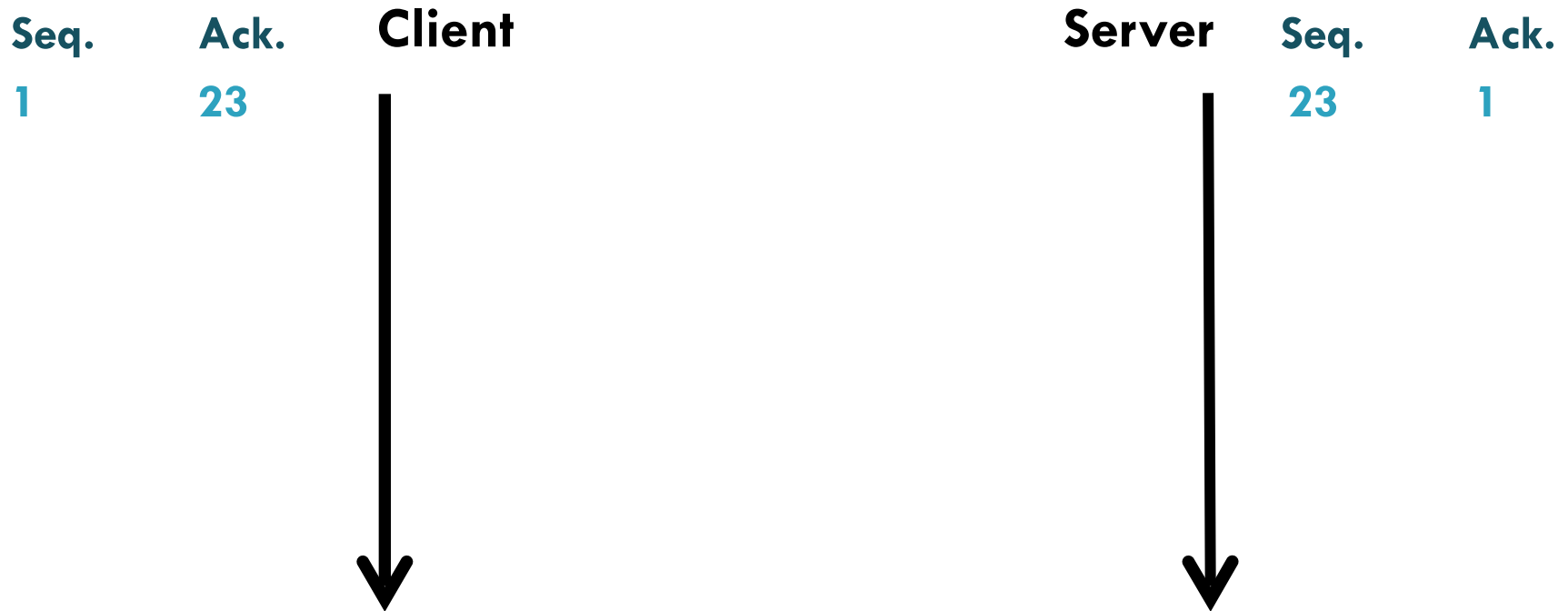
16



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

# Bidirectional Communication

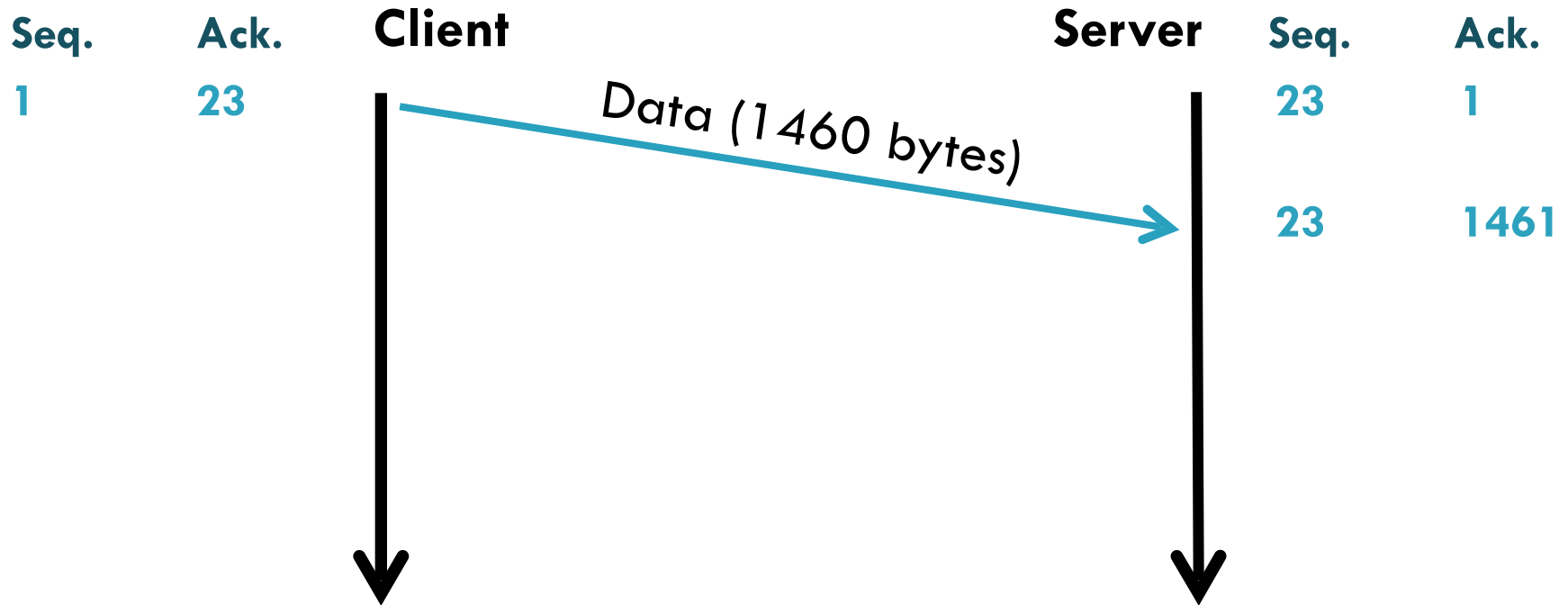
16



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

# Bidirectional Communication

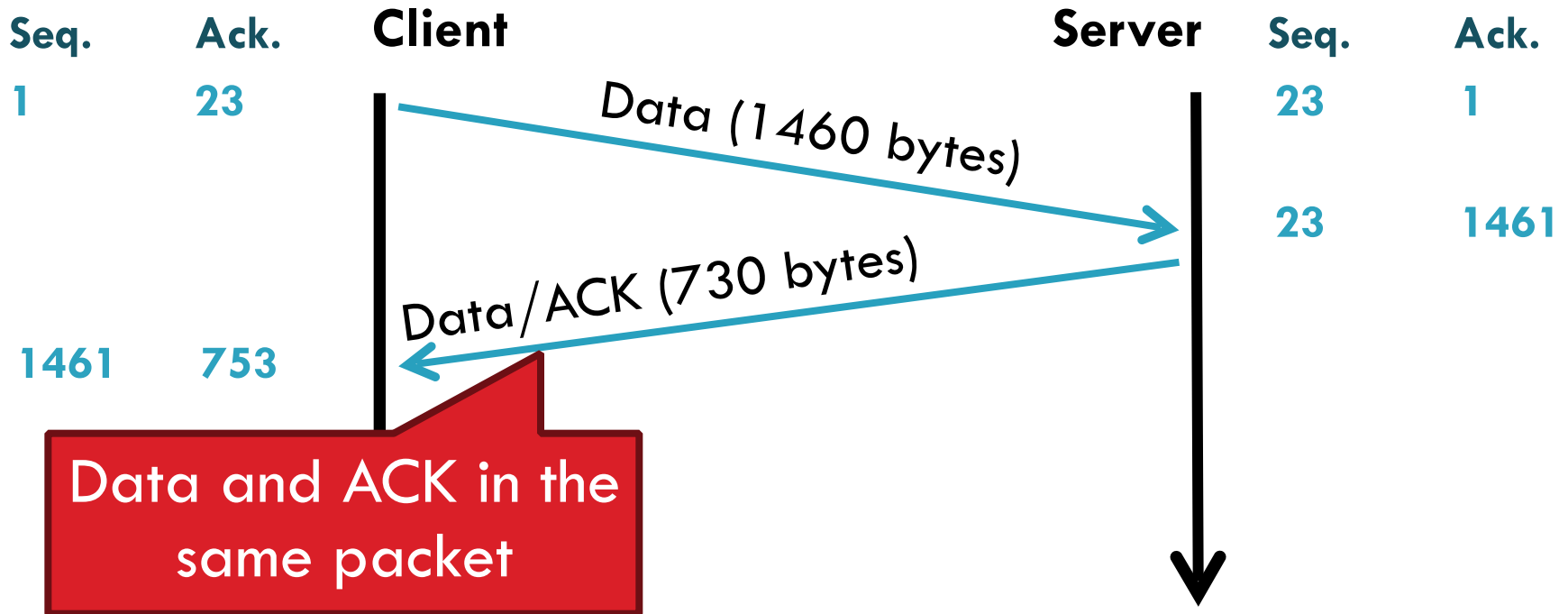
16



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

# Bidirectional Communication

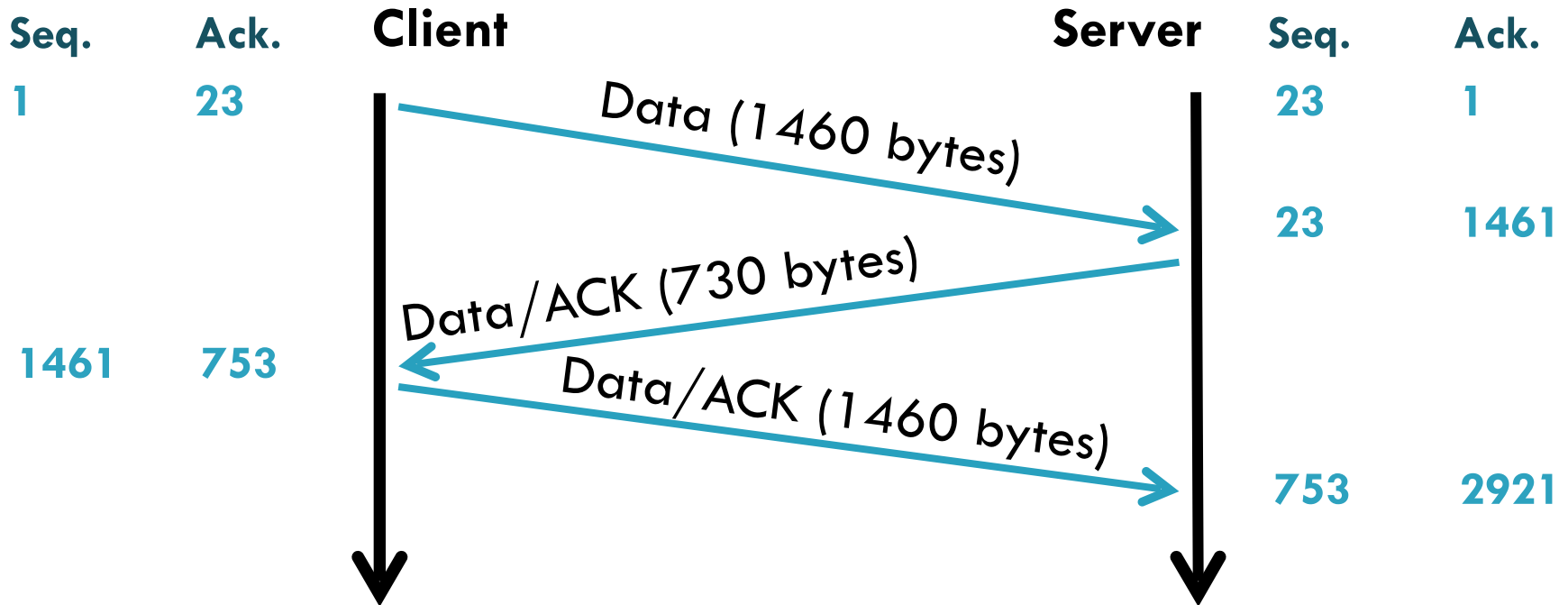
16



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

# Bidirectional Communication

16



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

# Flow Control

17

- Problem: how many packets should a sender transmit?
  - Too many packets may overwhelm the receiver
  - Size of the receivers buffers may change over time

# Flow Control

17

- 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

# Flow Control

17

- 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

18

Packet Sent

Src. Port		Dest. Port	
Sequence Number			
Acknowledgement Number			
HL	Flags	Window	
Checksum		Urgent Pointer	

Packet Received

Src. Port		Dest. Port	
Sequence Number			
Acknowledgement Number			
HL	Flags	Window	
Checksum		Urgent Pointer	



# Flow Control: Sender Side

18

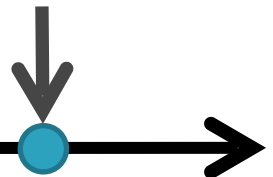
## Packet Sent

Src. Port		Dest. Port	
Sequence Number			
Acknowledgement Number			
HL	Flags	Window	
Checksum		Urgent Pointer	

## Packet Received

Src. Port		Dest. Port	
Sequence Number			
Acknowledgement Number			
HL	Flags	Window	
Checksum		Urgent Pointer	

**App Write**



# Flow Control: Sender Side

18

Packet Sent

Src. Port		Dest. Port	
Sequence Number			
Acknowledgement Number			
HL	Flags	Window	
Checksum		Urgent Pointer	

Packet Received

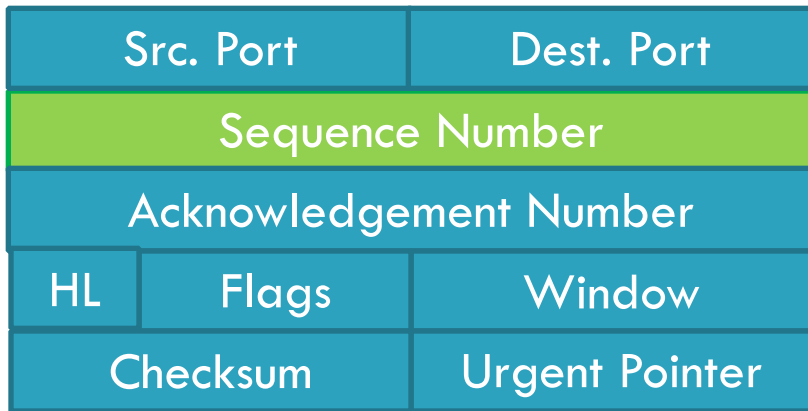
Src. Port		Dest. Port	
Sequence Number			
Acknowledgement Number			
HL	Flags	Window	
Checksum		Urgent Pointer	



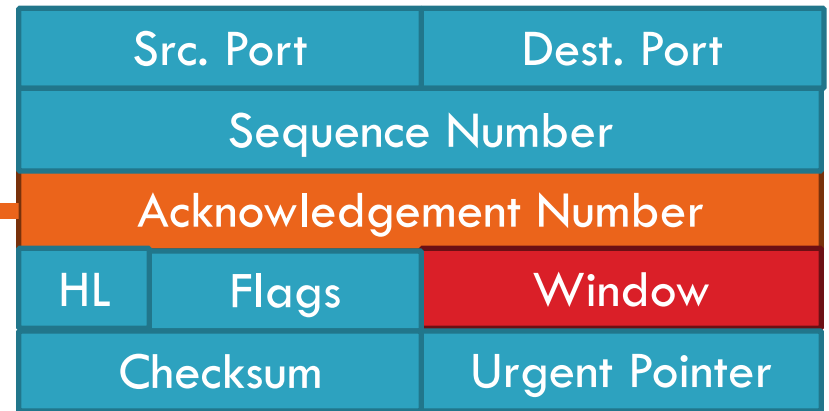
# Flow Control: Sender Side

18

Packet Sent



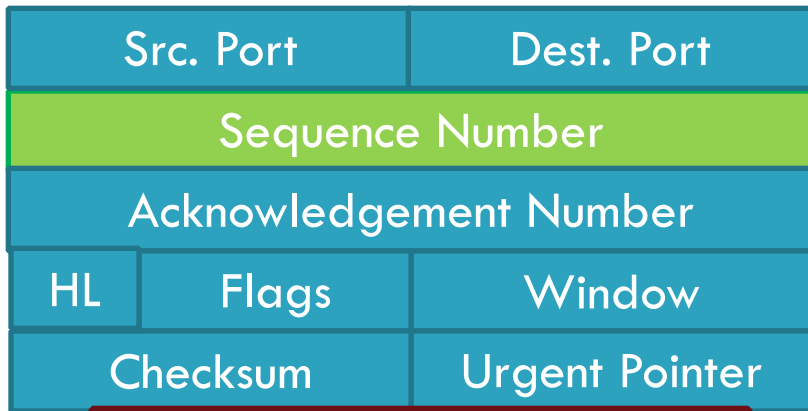
Packet Received



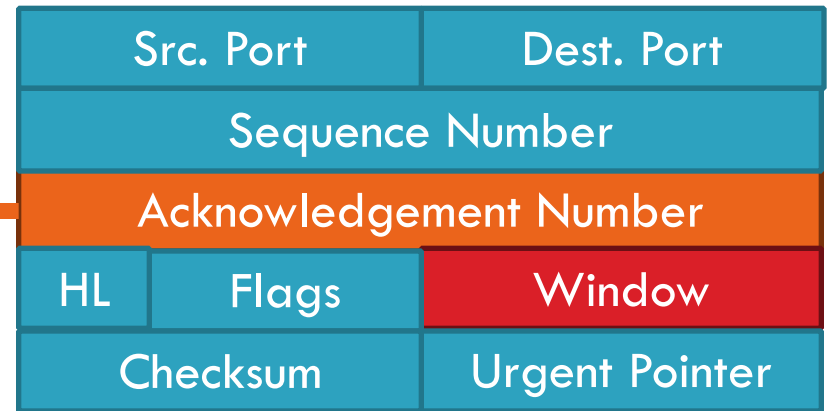
# Flow Control: Sender Side

18

## Packet Sent



## Packet Received



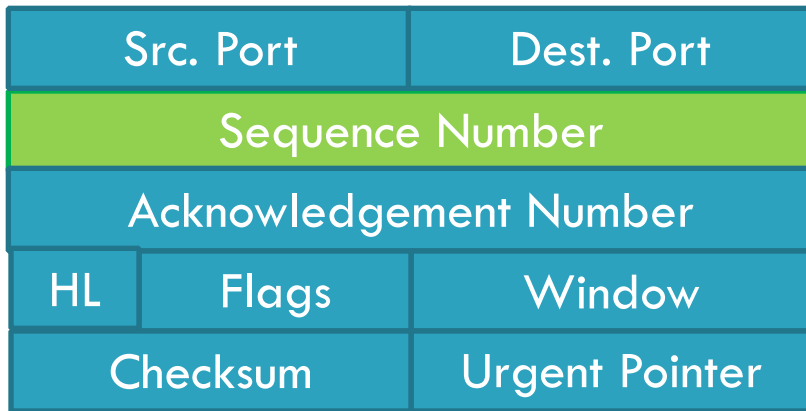
**Must be buffered until ACKed**



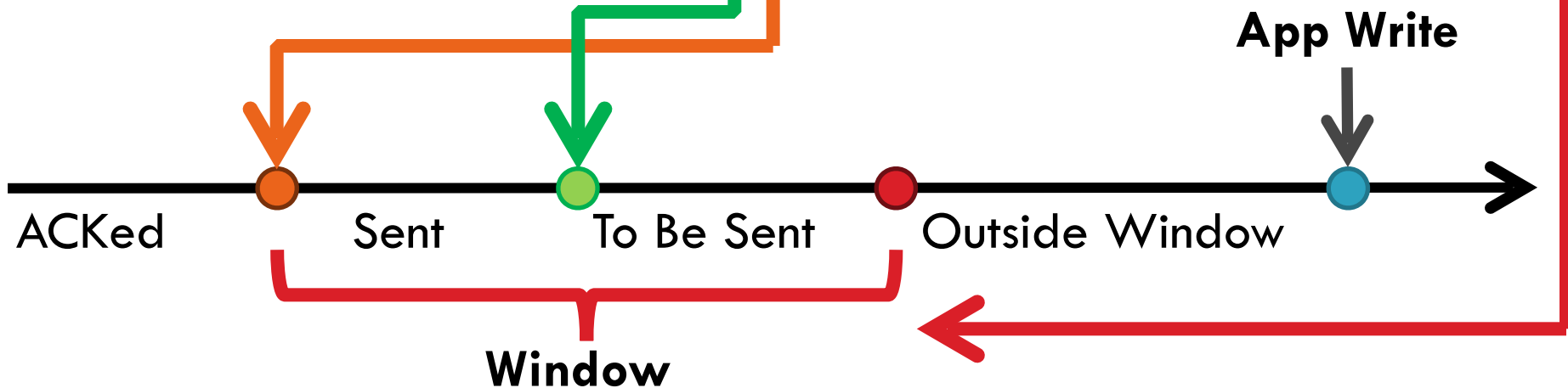
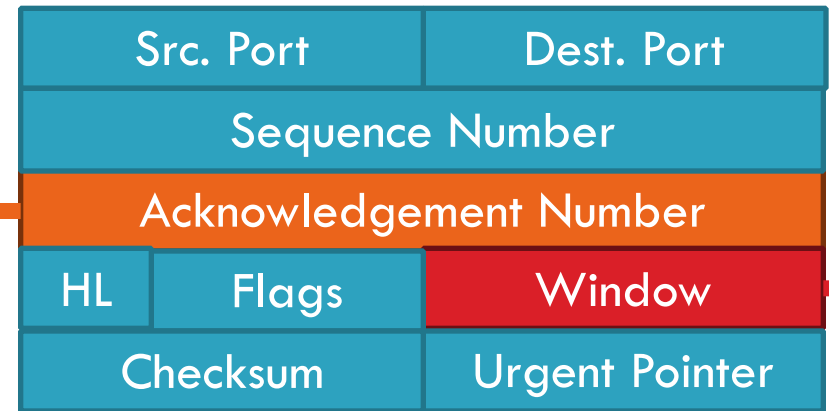
# Flow Control: Sender Side

18

Packet Sent

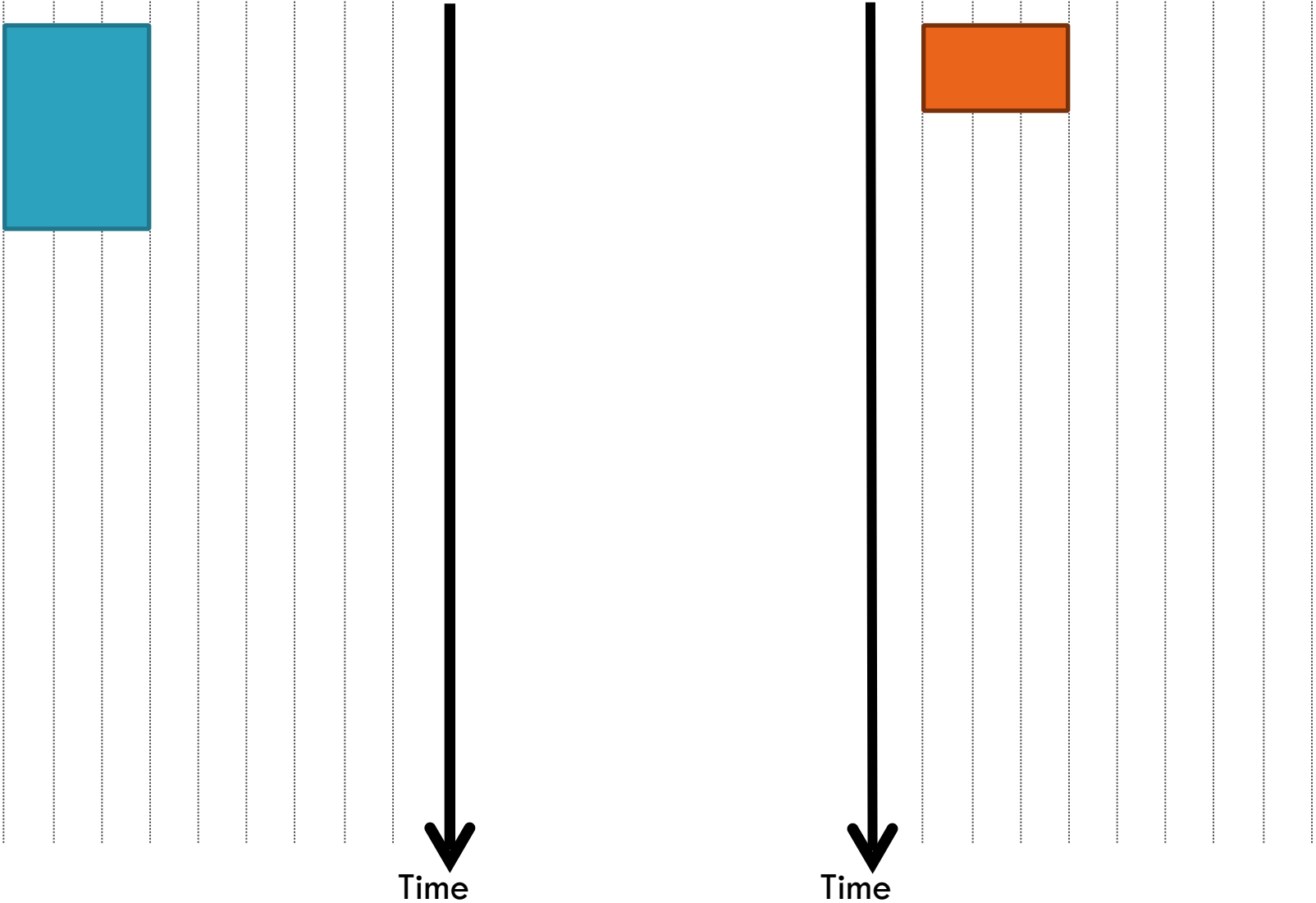


Packet Received



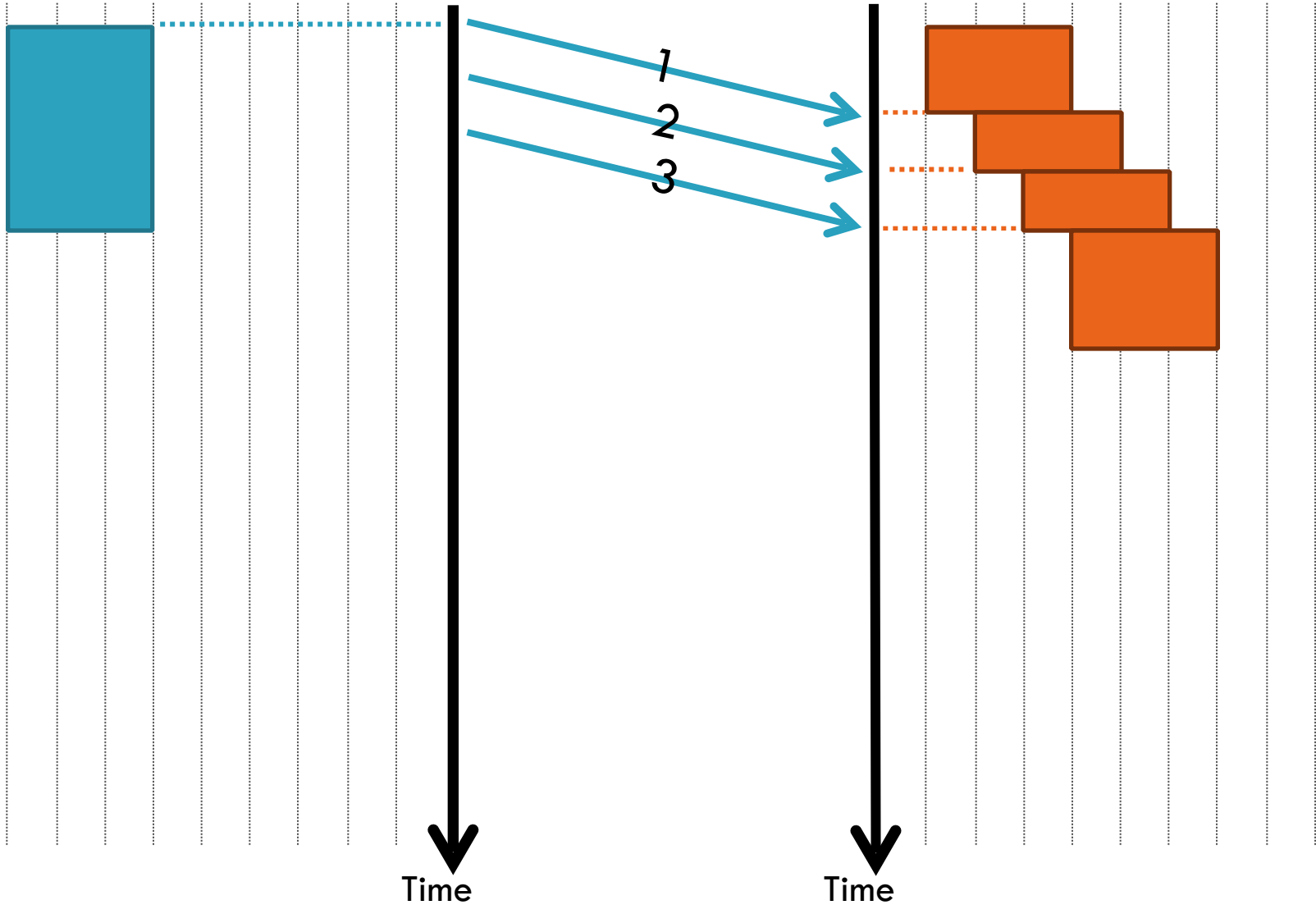
# Sliding Window Example

19



# Sliding Window Example

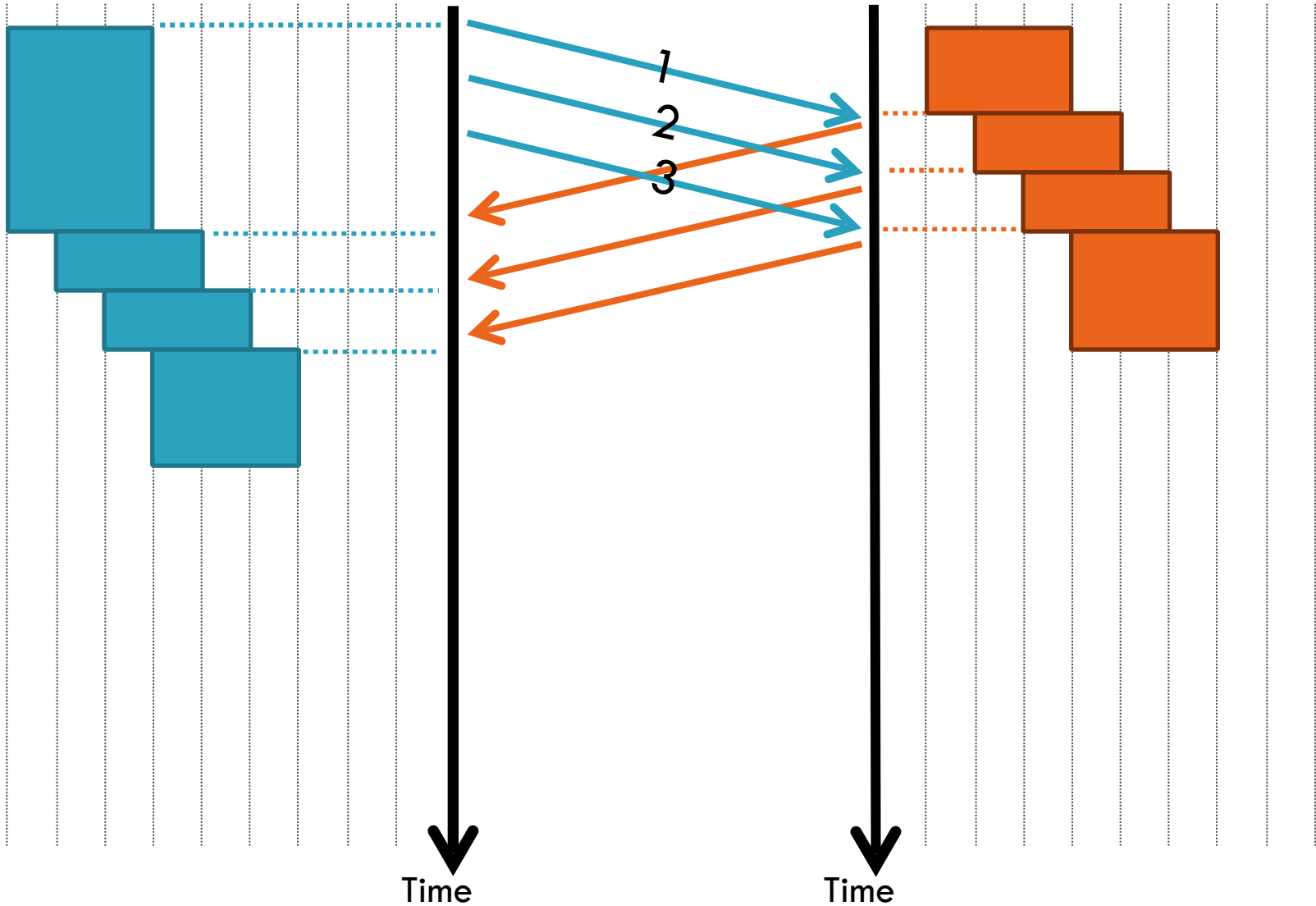
19





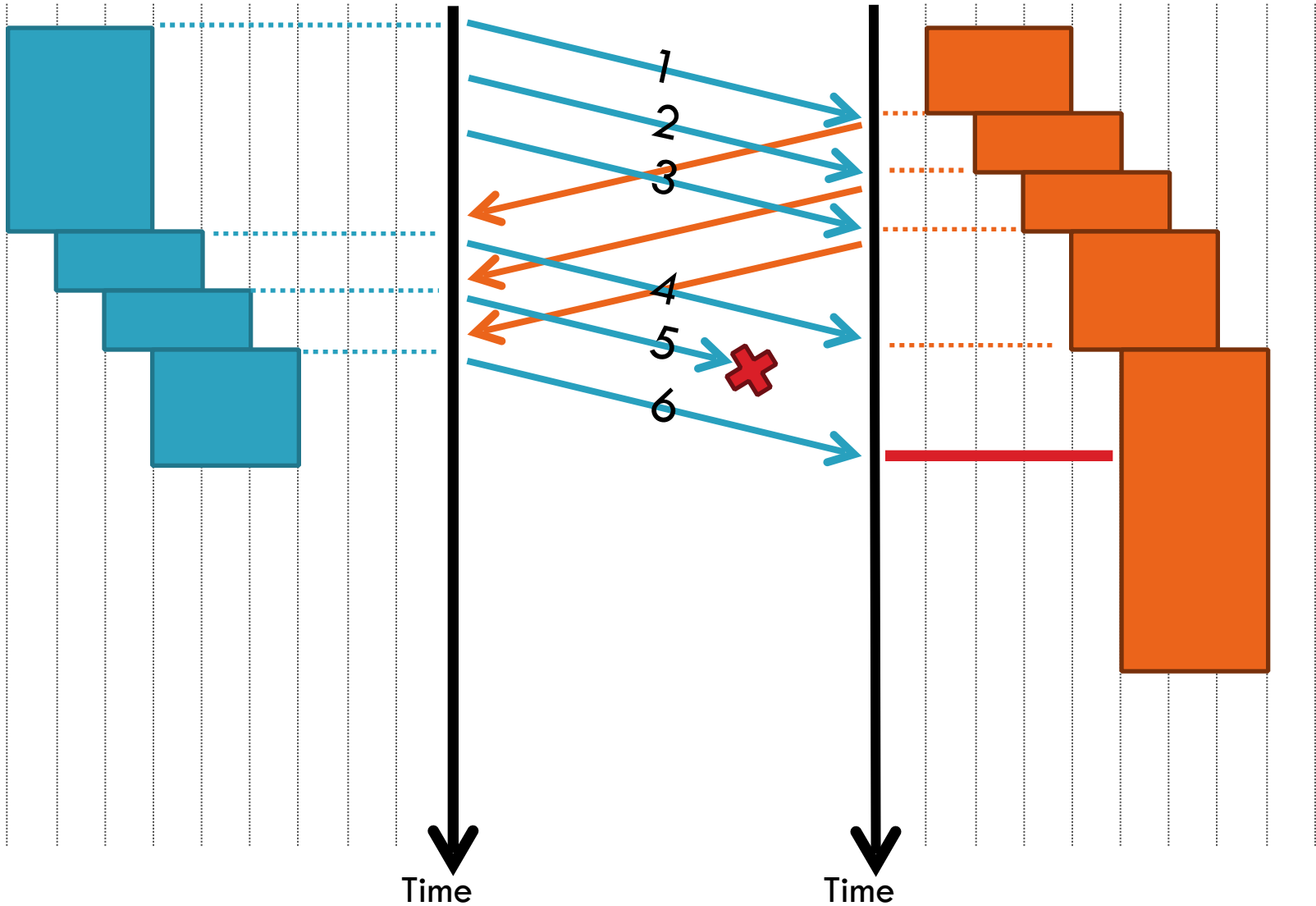
# Sliding Window Example

19



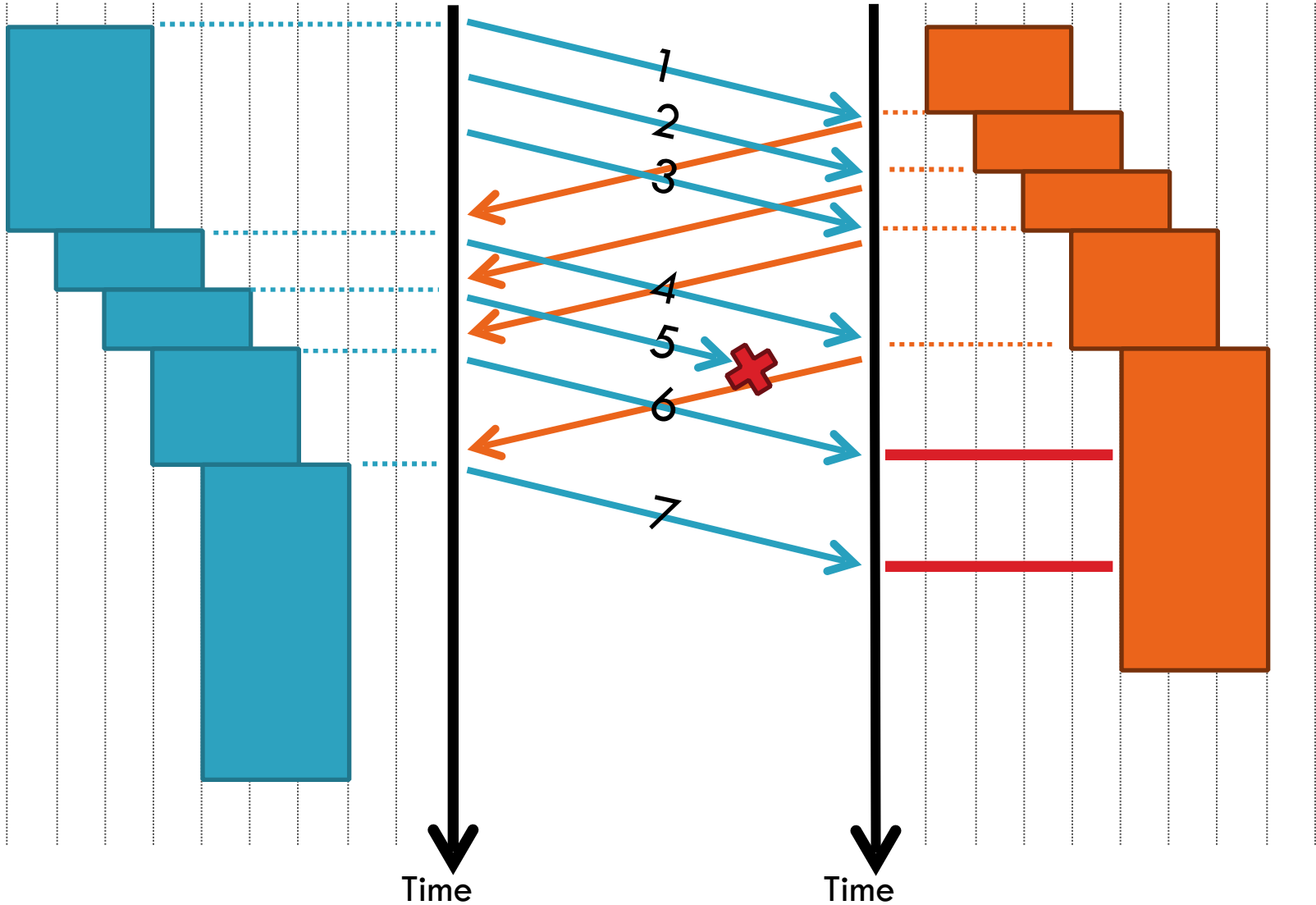
# Sliding Window Example

19



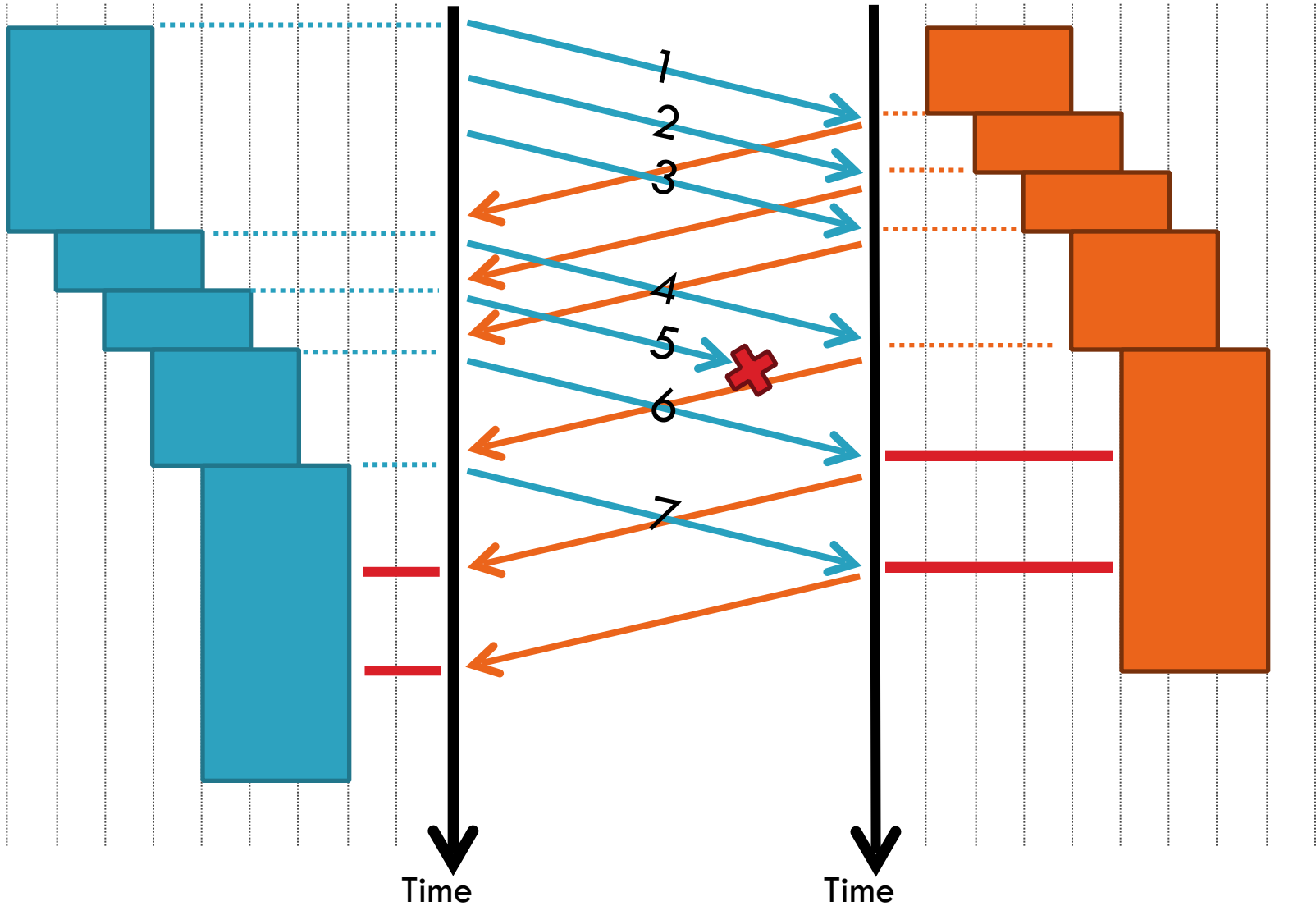
# Sliding Window Example

19



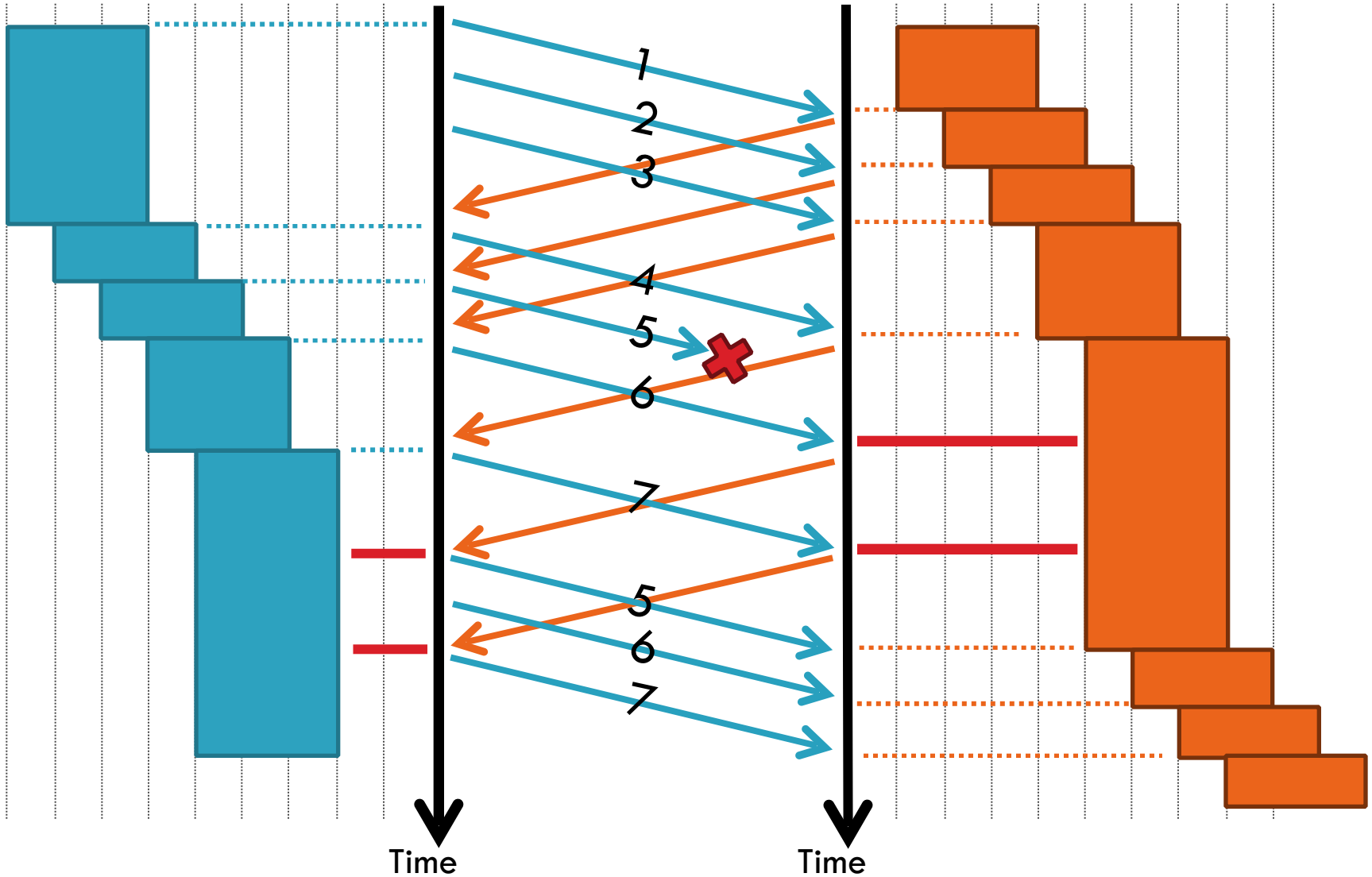
# Sliding Window Example

19



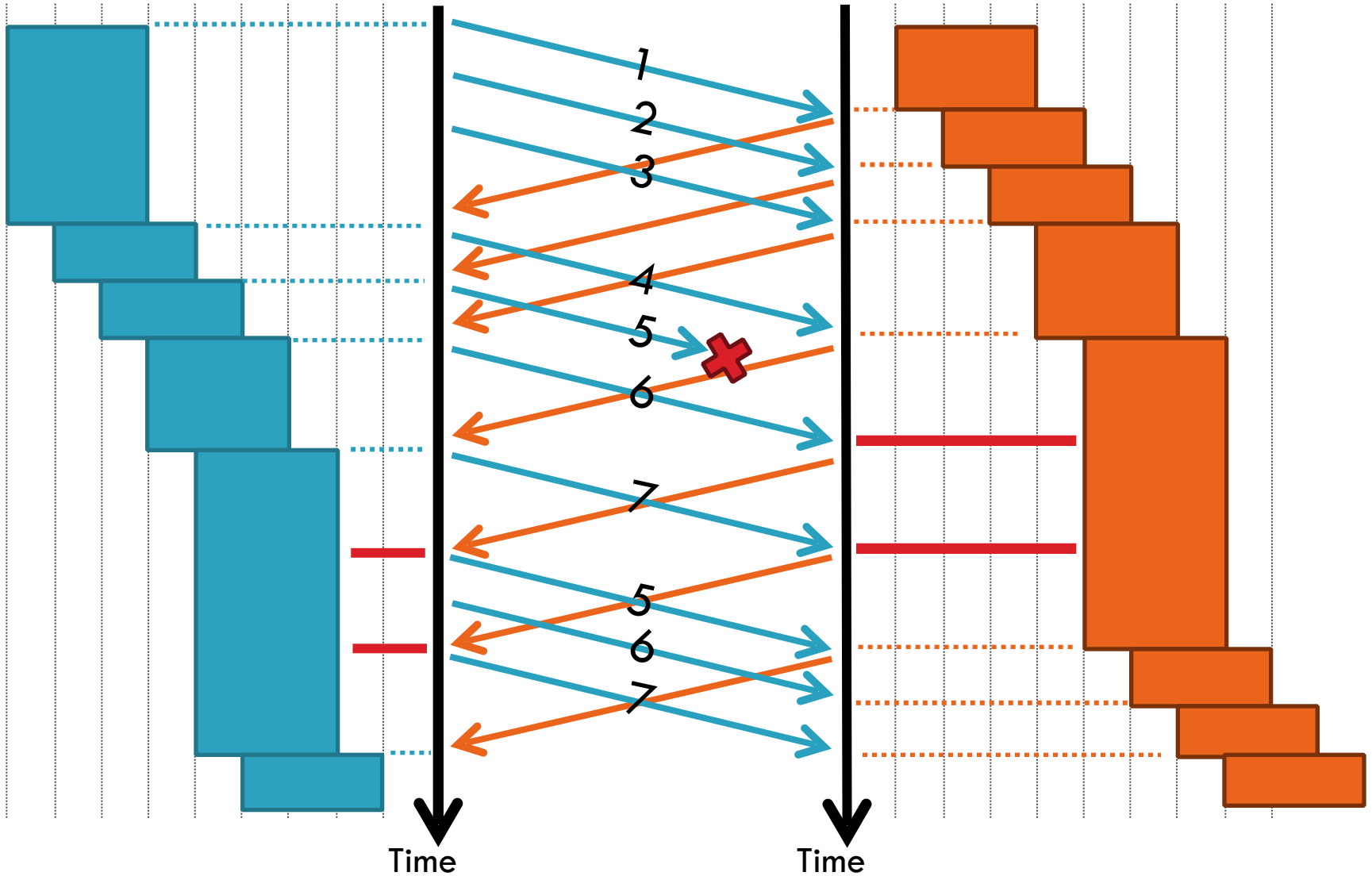
# Sliding Window Example

19



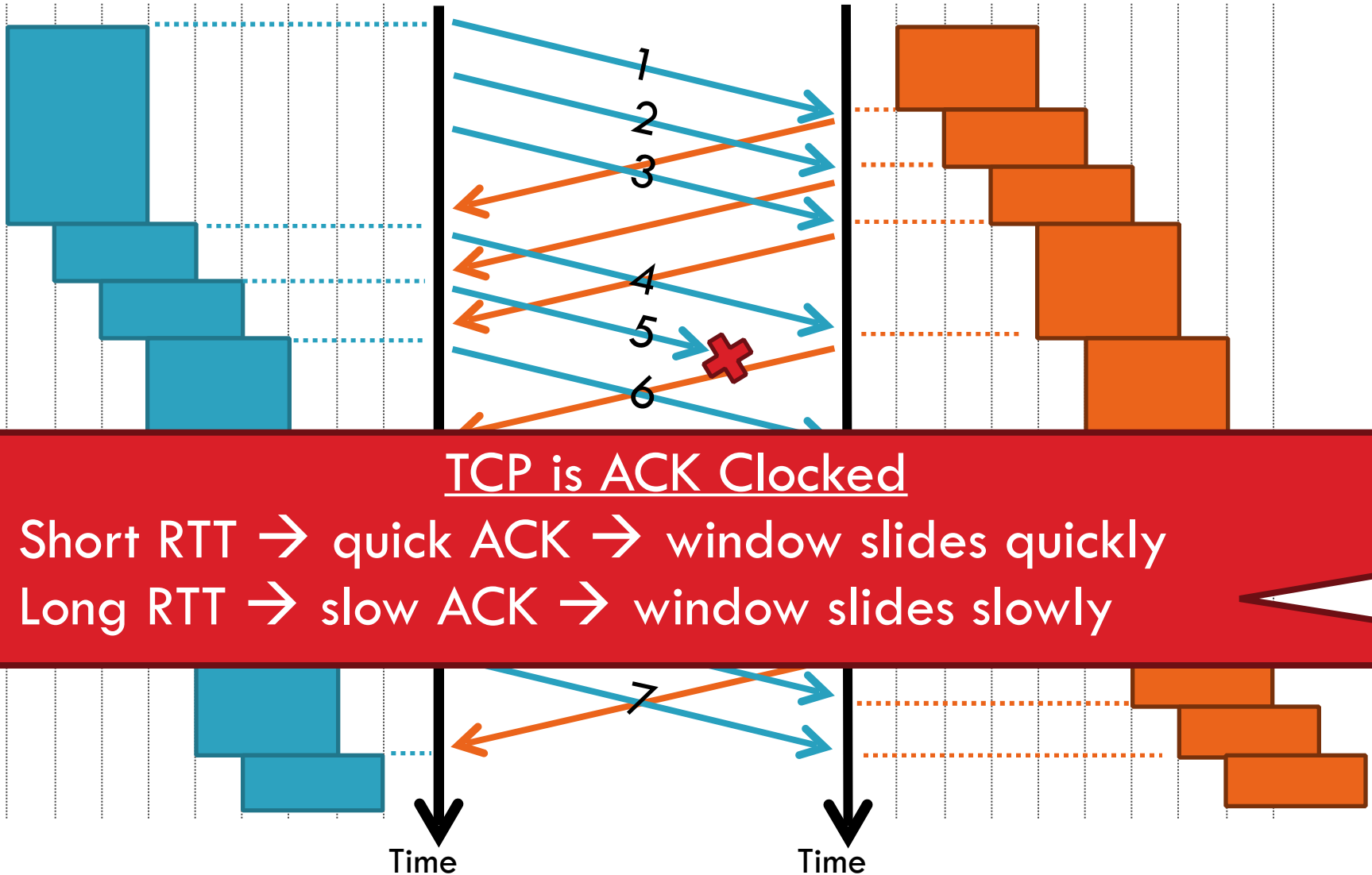
# Sliding Window Example

19



# Sliding Window Example

19



# What Should the Receiver ACK?

20

1. ACK every packet



# What Should the Receiver ACK?

20

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

# What Should the Receiver ACK?

20

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

# What Should the Receiver ACK?

20

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

# What Should the Receiver ACK?

20

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

21

- 32 bits, unsigned
  - ▣ Why so big?

# Sequence Numbers, Revisited

21

- 32 bits, unsigned
  - Why so big?
- For the sliding window you need...
  - $|\text{Sequence \# Space}| > 2 * |\text{Sending Window Size}|$
  - $2^{32} > 2 * 2^{16}$

# Sequence Numbers, Revisited

21

- 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 3 minutes
  - Sequence number would wrap around at 286Mbps

# Sequence Numbers, Revisited

21

- 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 3 minutes
  - Sequence number would wrap around at 286Mbps
    - What about GigE? PAWS algorithm + TCP options



# Silly Window Syndrome

22

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

# Silly Window Syndrome

22

- Problem: what if the window size is very small?
  - Multiple, small packets, headers dominate data



# Silly Window Syndrome

22

- 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

23

1. If the window  $\geq$  MSS and available data  $\geq$  MSS:  
Send the data
2. Elif there is unACKed data:  
Enqueue data in a buffer (send after a timeout)
3. Else: send the data

# Nagle's Algorithm

23



1. If the window  $\geq$  MSS and available data  $\geq$  MSS:  
Send the data
2. Elif there is unACKed data:  
Enqueue data in a buffer (send after a timeout)
3. Else: send the data

A red callout box with a black border and a white shadow, pointing to the third step of the algorithm. It contains the text: "Send a non-full packet if nothing else is happening".

Send a non-full packet if nothing else is happening



# Nagle's Algorithm

23

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 (send after a timeout)
3. Else: send the data  Send a non-full packet if nothing else is happening

# Nagle's Algorithm

23

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 (send after a timeout)
  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));`

# Error Detection

24

- Checksum detects (some) packet corruption
  - ▣ Computed over IP header, TCP header, and data



# Error Detection

24

- Checksum detects (some) packet corruption
  - ▣ Computed over IP header, TCP header, and data
- Sequence numbers catch sequence problems
  - ▣ Duplicates are ignored
  - ▣ Out-of-order packets are reordered or dropped
  - ▣ Missing sequence numbers indicate lost packets

# Error Detection

24

- Checksum detects (some) packet corruption
  - ▣ Computed over IP header, TCP header, and data
- Sequence numbers catch sequence problems
  - ▣ Duplicates are ignored
  - ▣ Out-of-order packets are reordered or dropped
  - ▣ Missing sequence numbers indicate lost packets
- Lost segments detected by sender
  - ▣ Use **timeout** to detect missing ACKs
  - ▣ Need to estimate RTT to calibrate the timeout
  - ▣ Sender must keep copies of all data until ACK

# Retransmission Time Outs (RTO)

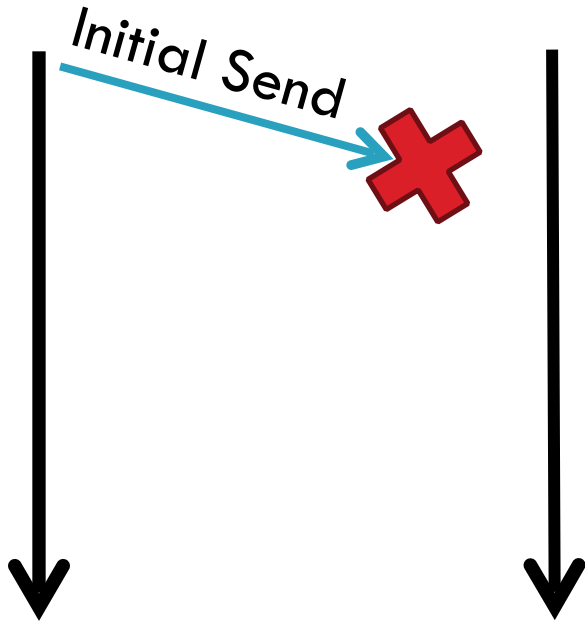
25

- Problem: time-out is linked to round trip time

# Retransmission Time Outs (RTO)

25

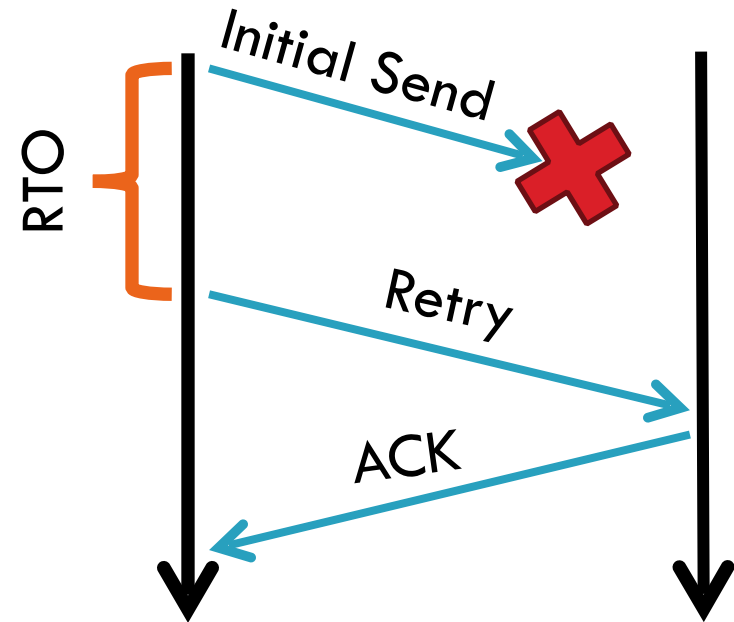
- Problem: time-out is linked to round trip time



# Retransmission Time Outs (RTO)

25

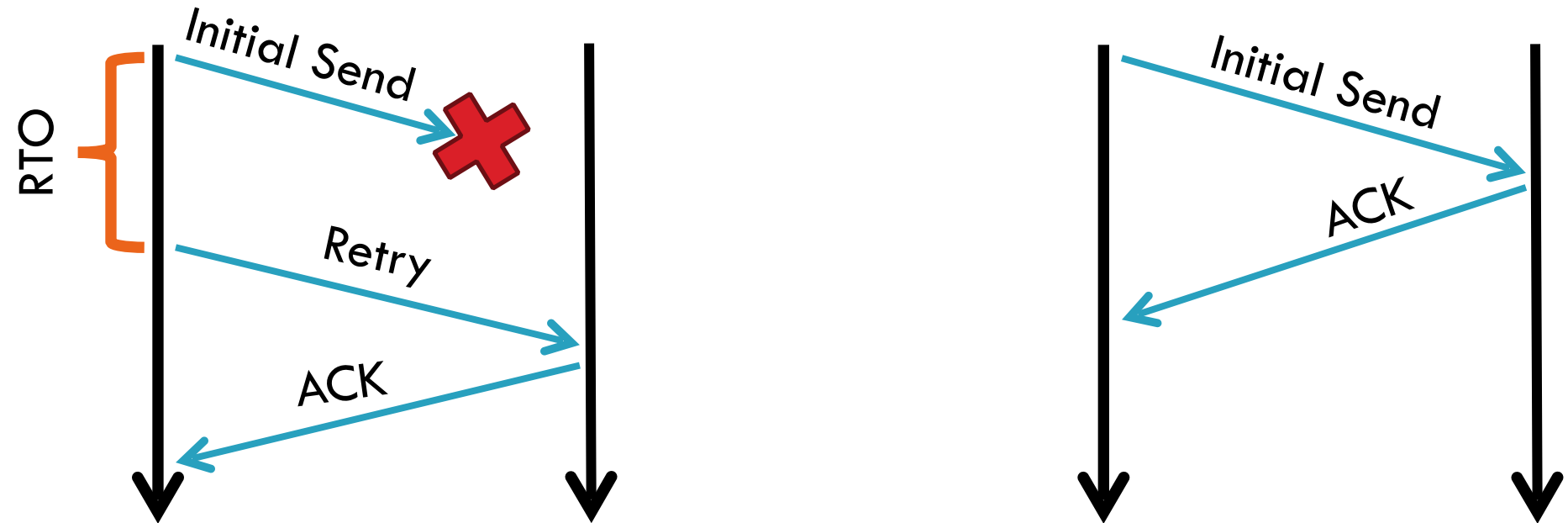
- Problem: time-out is linked to round trip time



# Retransmission Time Outs (RTO)

25

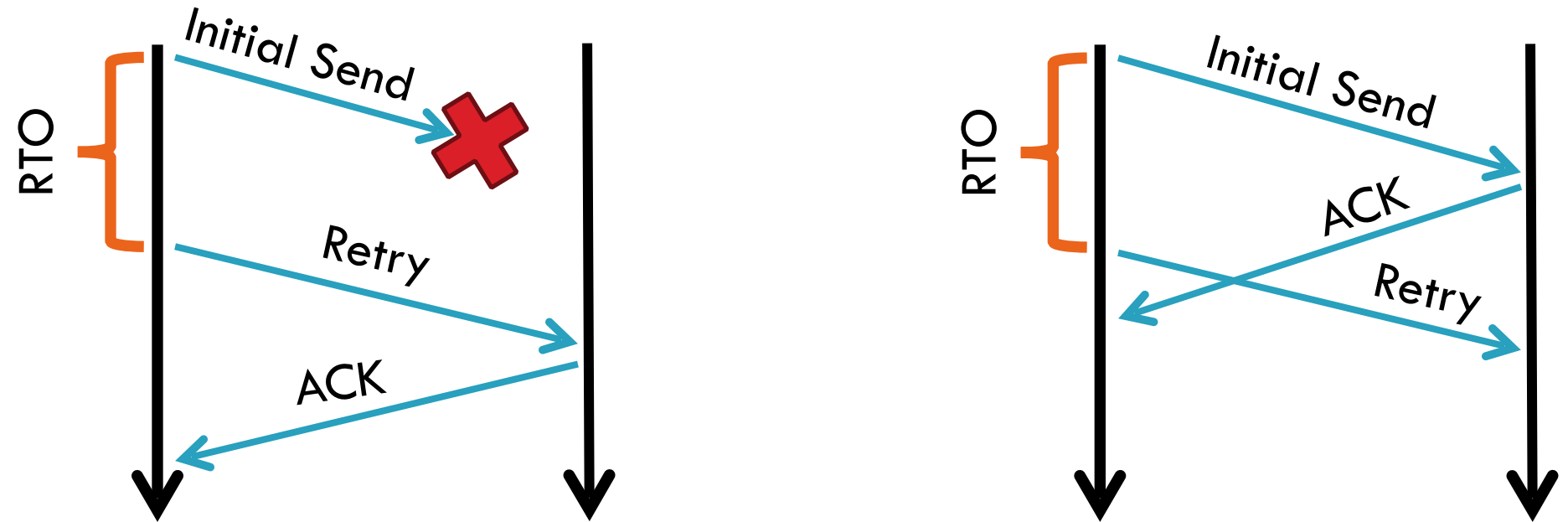
- Problem: time-out is linked to round trip time



# Retransmission Time Outs (RTO)

25

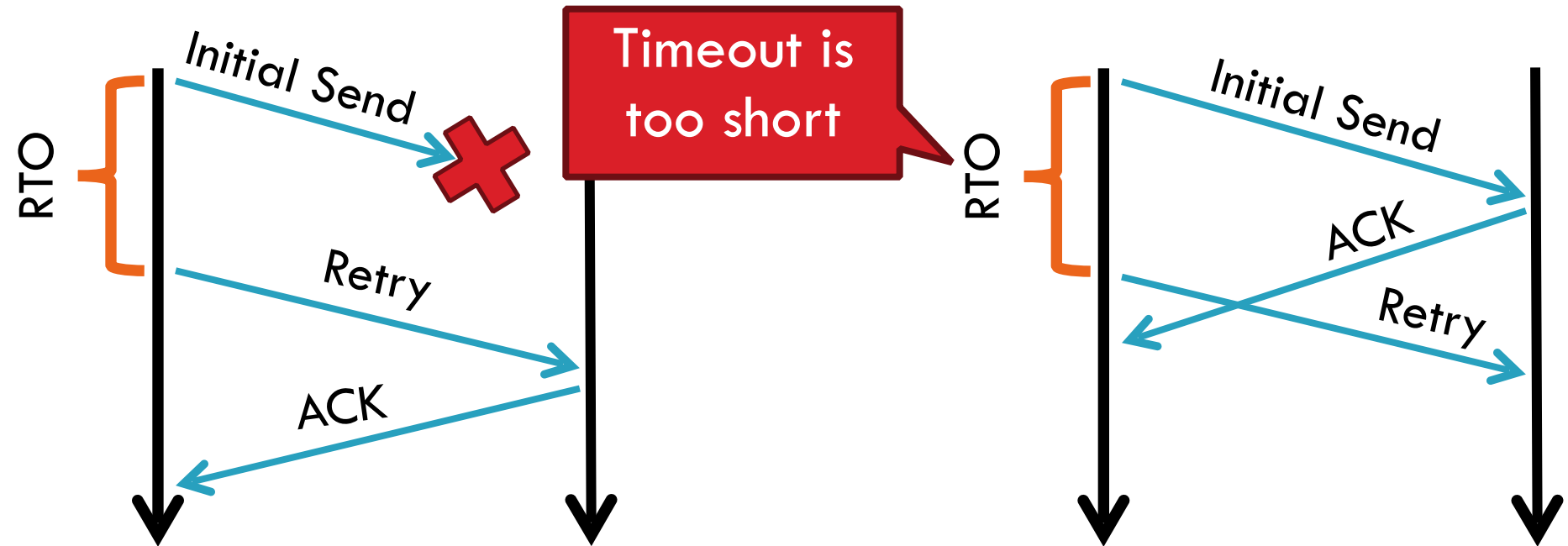
- Problem: time-out is linked to round trip time



# Retransmission Time Outs (RTO)

25

- Problem: time-out is linked to round trip time

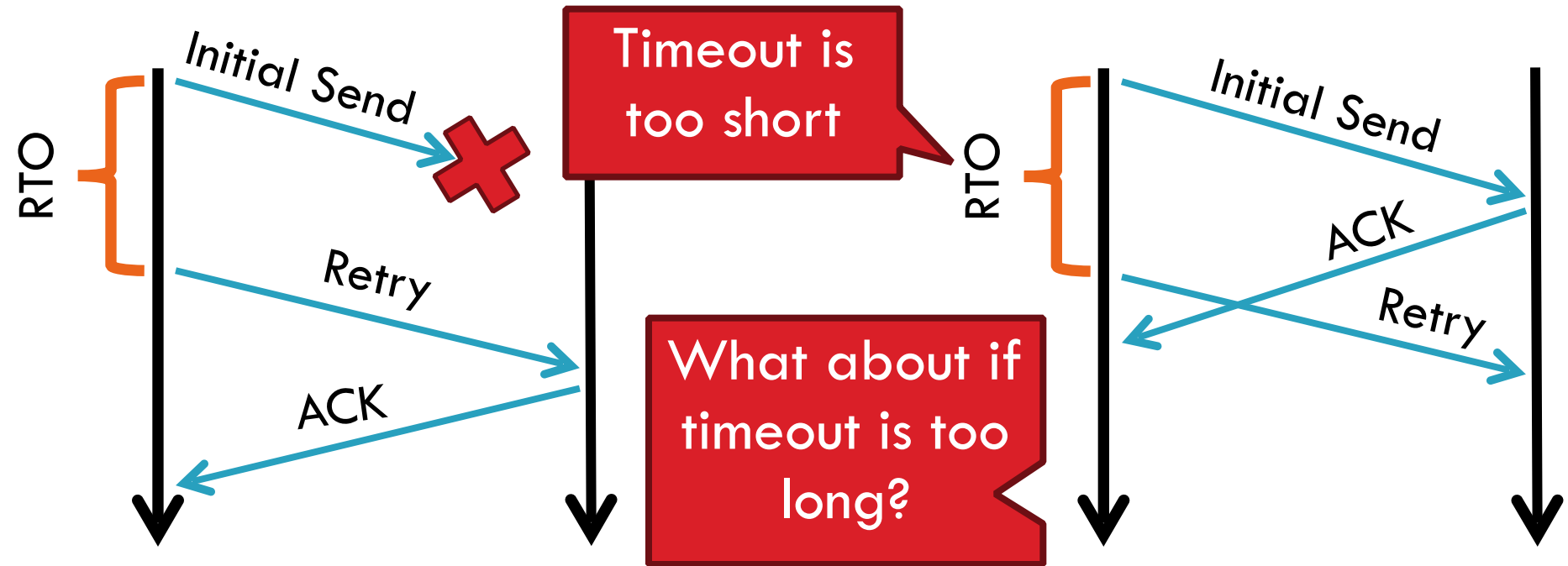




# Retransmission Time Outs (RTO)

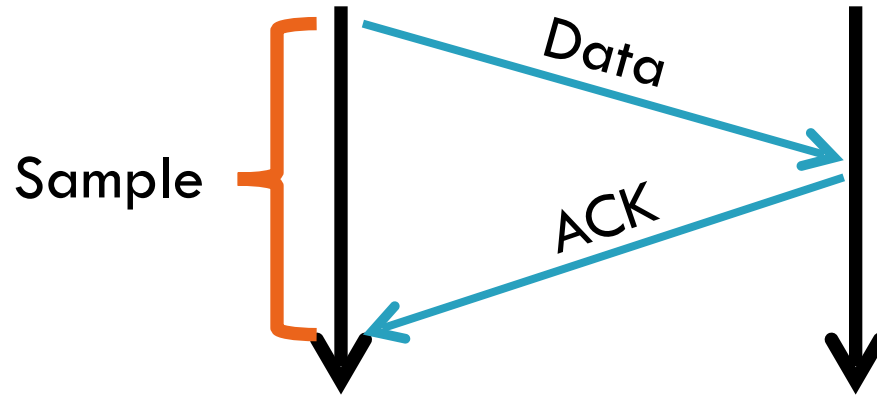
25

- Problem: time-out is linked to round trip time



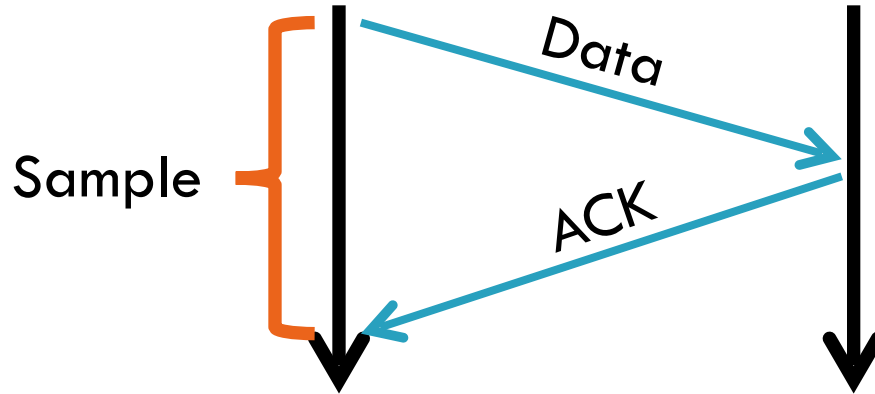
# Round Trip Time Estimation

26



# Round Trip Time Estimation

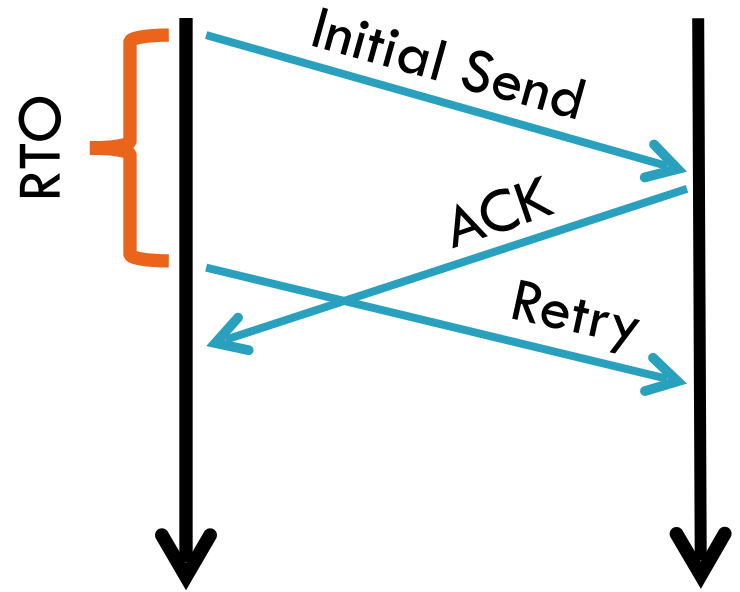
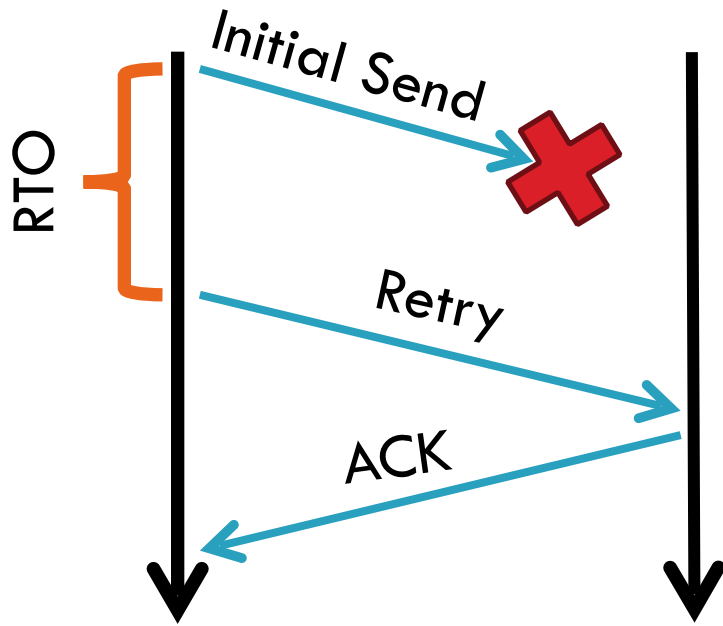
26



- Original TCP round-trip estimator
  - ▣ RTT estimated as a moving average
  - ▣  $\text{new\_rtt} = \alpha (\text{old\_rtt}) + (1 - \alpha)(\text{new\_sample})$
  - ▣ Recommended  $\alpha$ : 0.8-0.9 (0.875 for most TCPs)
- $\text{RTO} = 2 * \text{new\_rtt}$  (i.e. TCP is conservative)

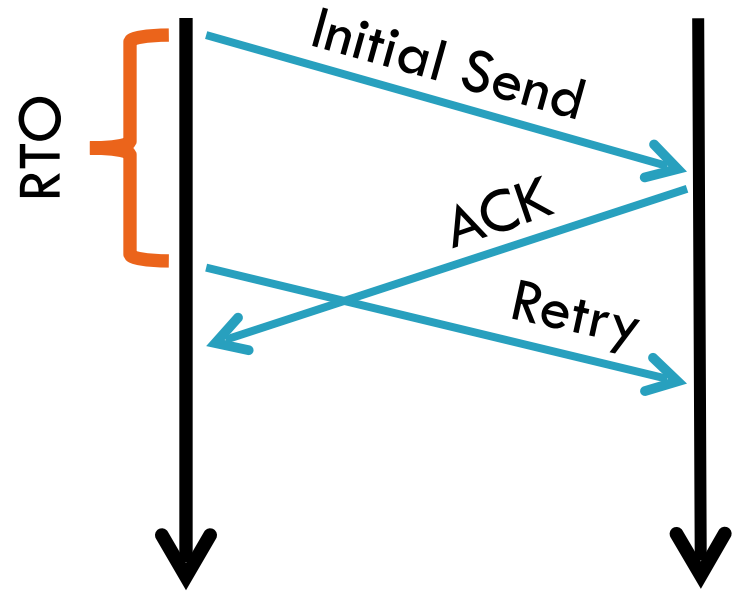
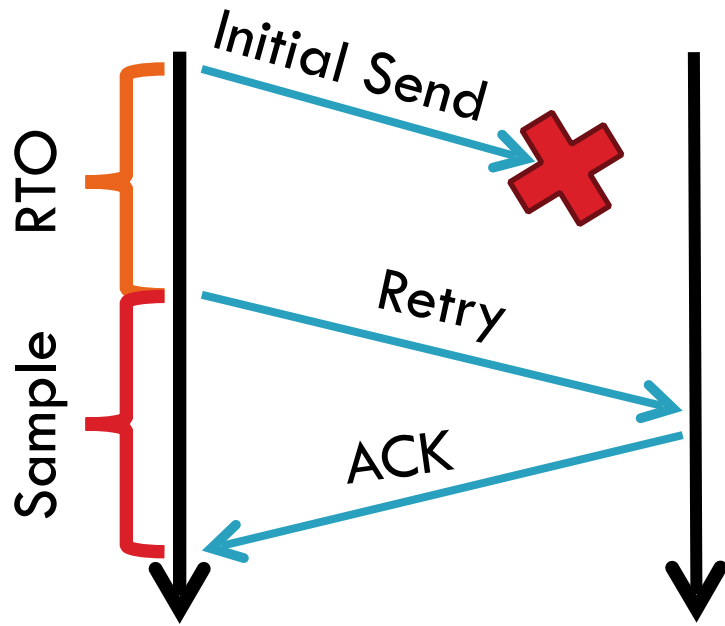
# RTT Sample Ambiguity

27



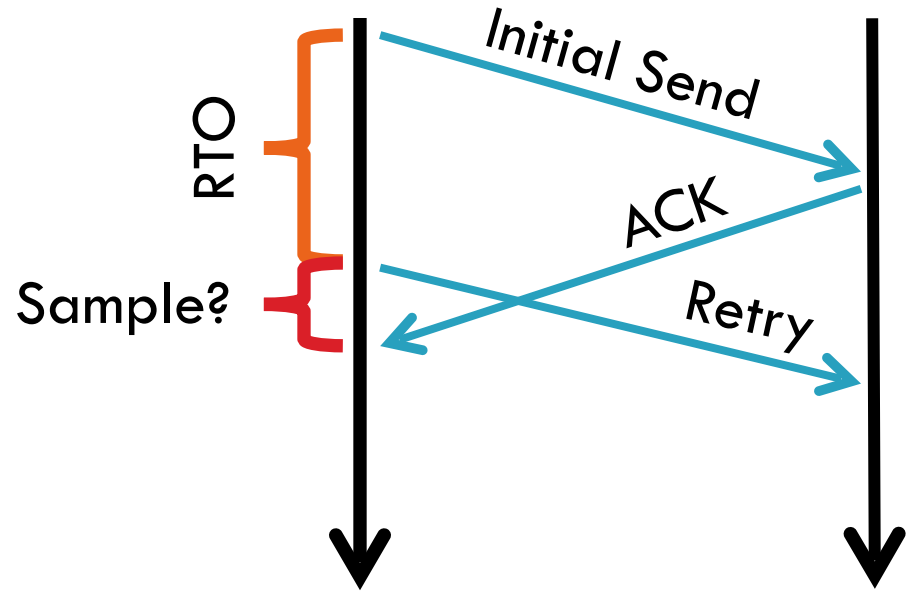
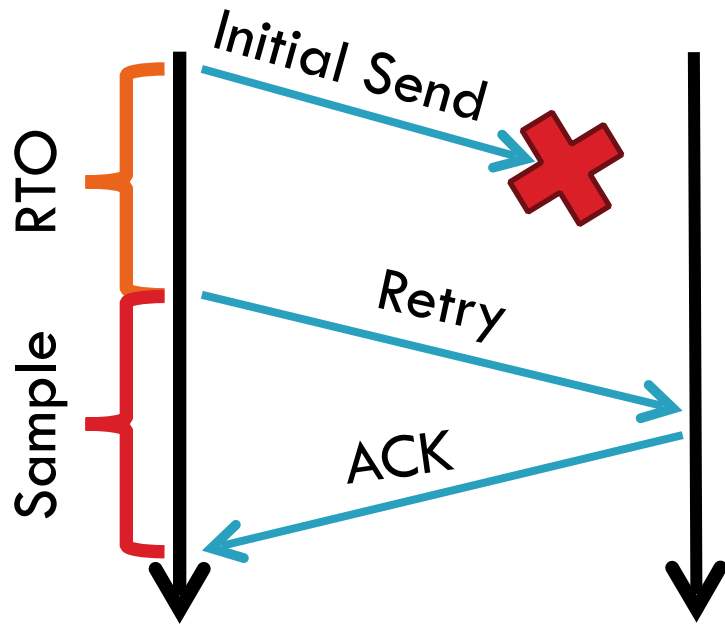
# RTT Sample Ambiguity

27



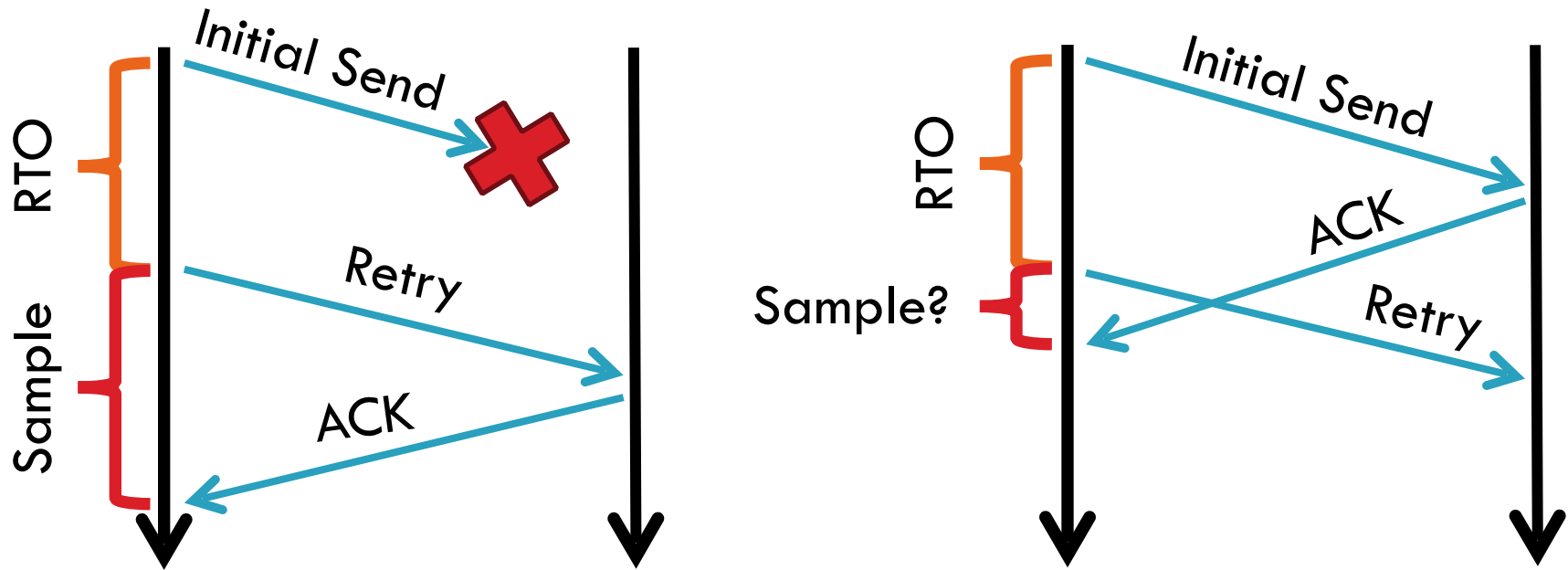
# RTT Sample Ambiguity

27



# RTT Sample Ambiguity

27



- Karn's algorithm: ignore samples for retransmitted segments