

# CSCI 466: Networks

TCP Flow Control, Timeout, Congestion Control

Reese Pearsall  
Fall 2024

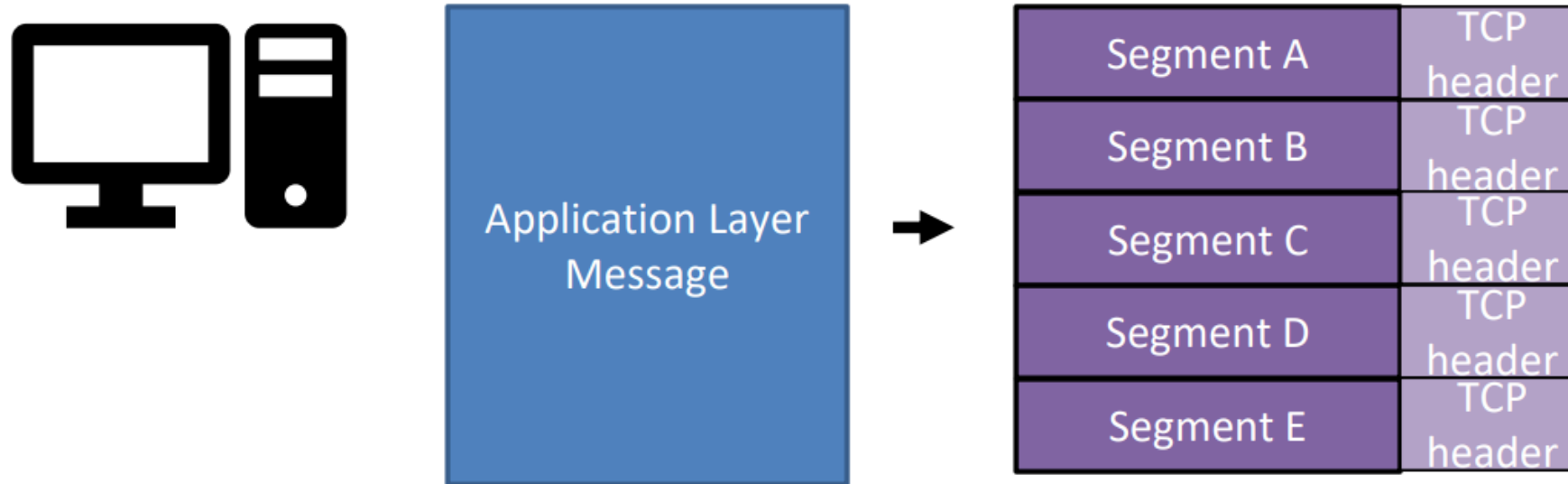
# Announcements

Wireshark Lab due on Wednesday

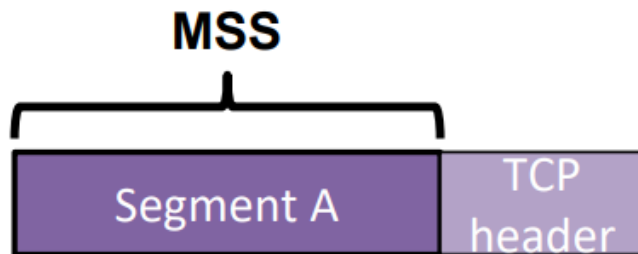
Have a good weekend



# TCP Flow Control

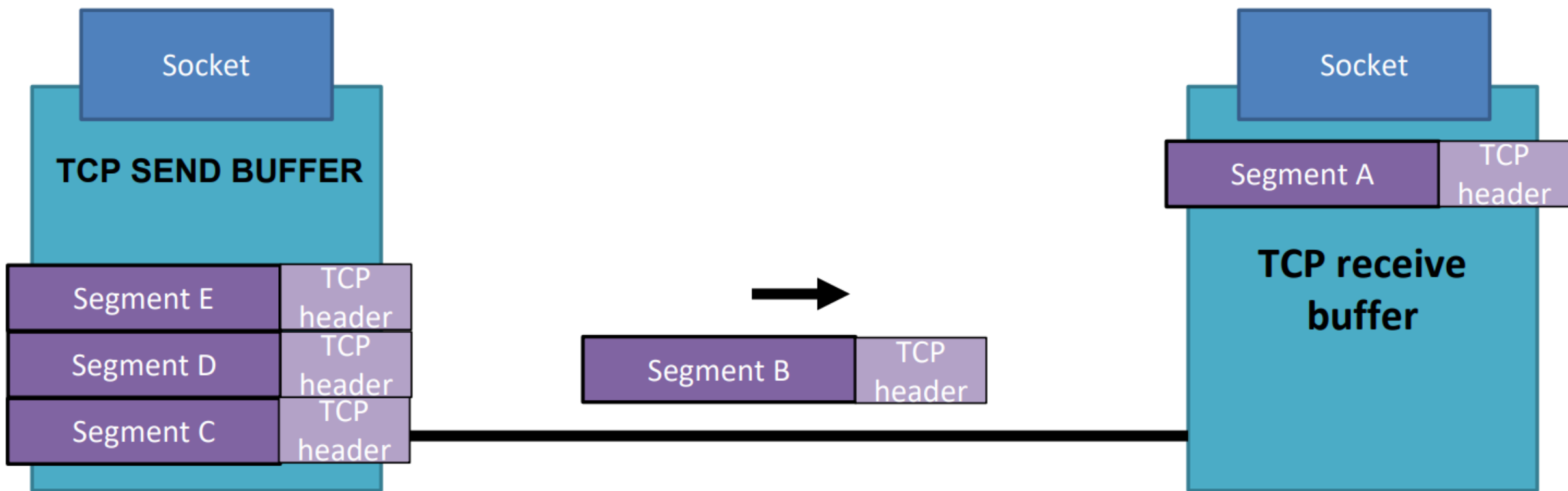


Application layer messages are split into smaller chunks called **segments**



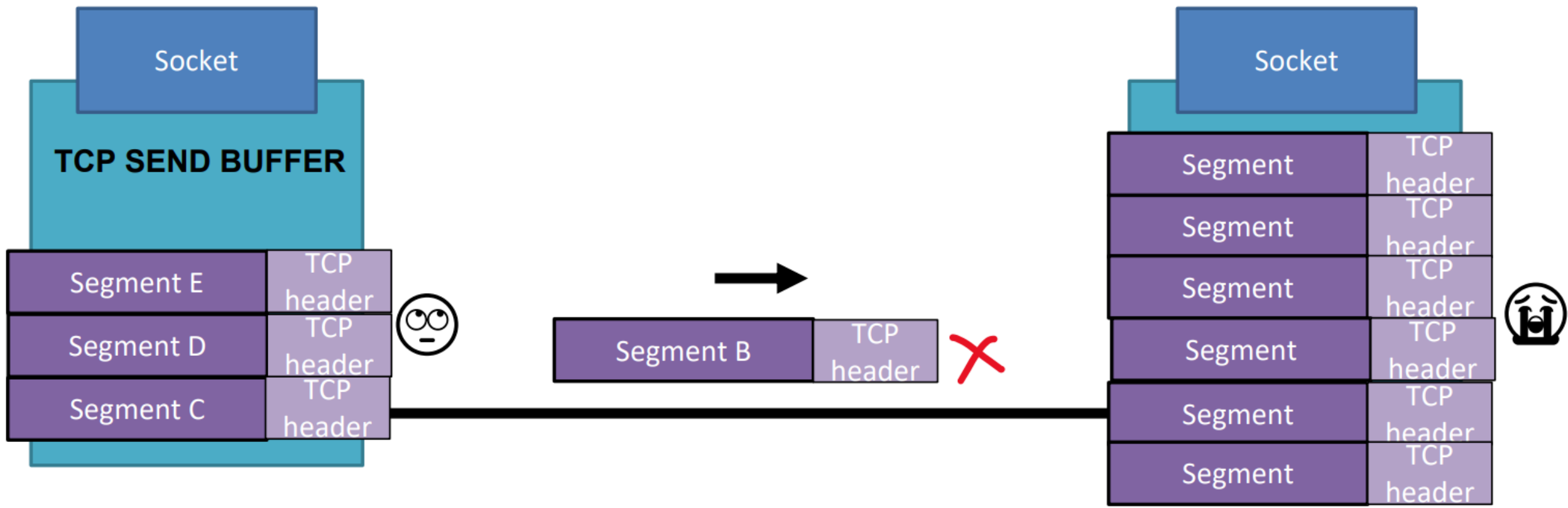
The size of these segments is determined by the **maximum segment size (MSS)**

# TCP Flow Control



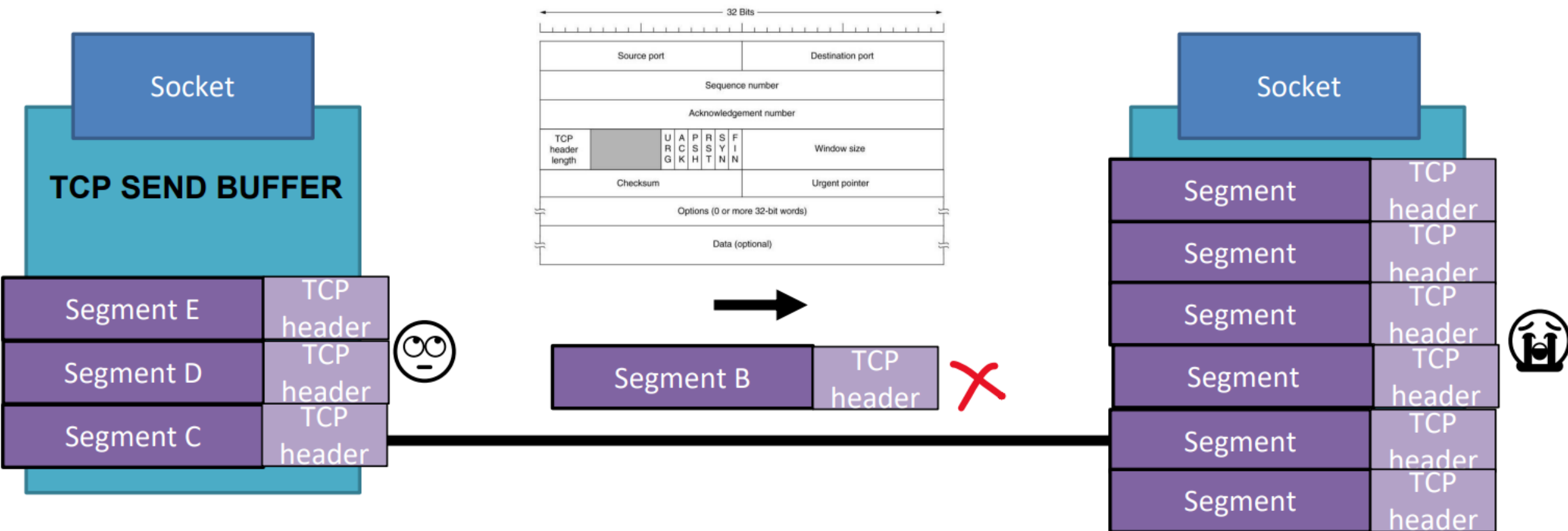
Applications read streams of data from a **TCP buffer**

# TCP Flow Control



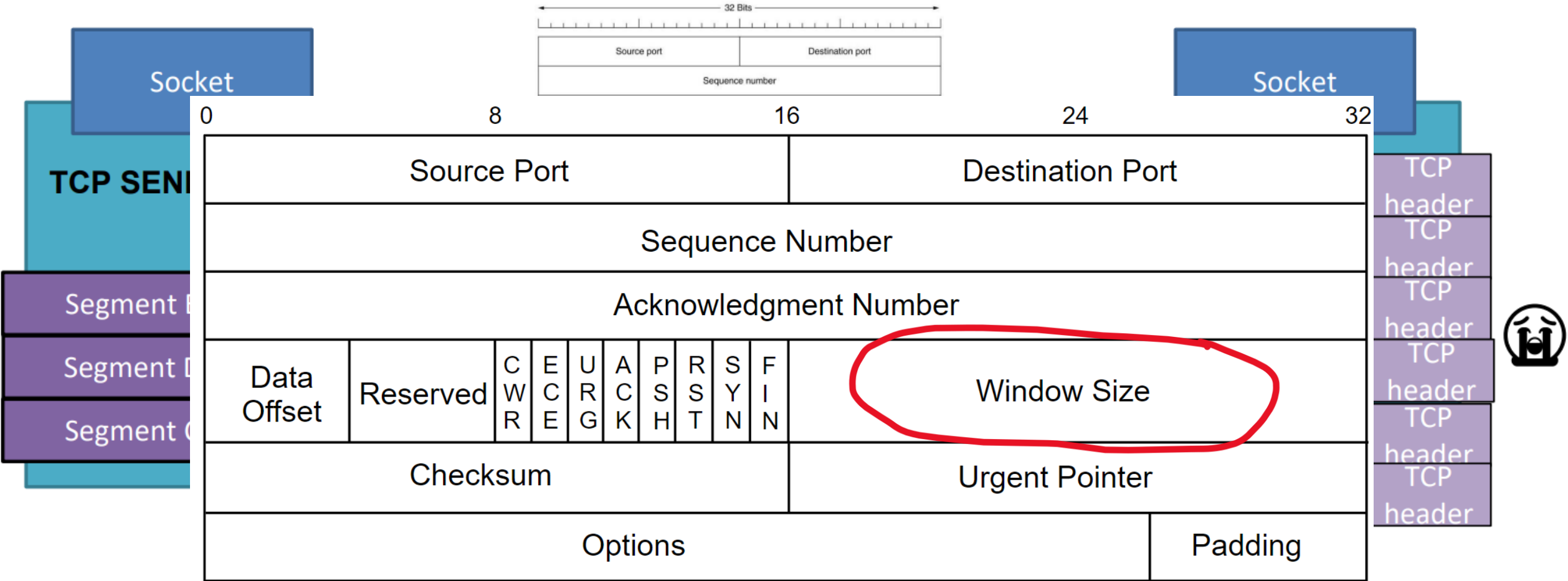
How could we prevent something like this from happening?

# TCP Flow Control



**We could send back to the sender how much available space we have in our buffer!**

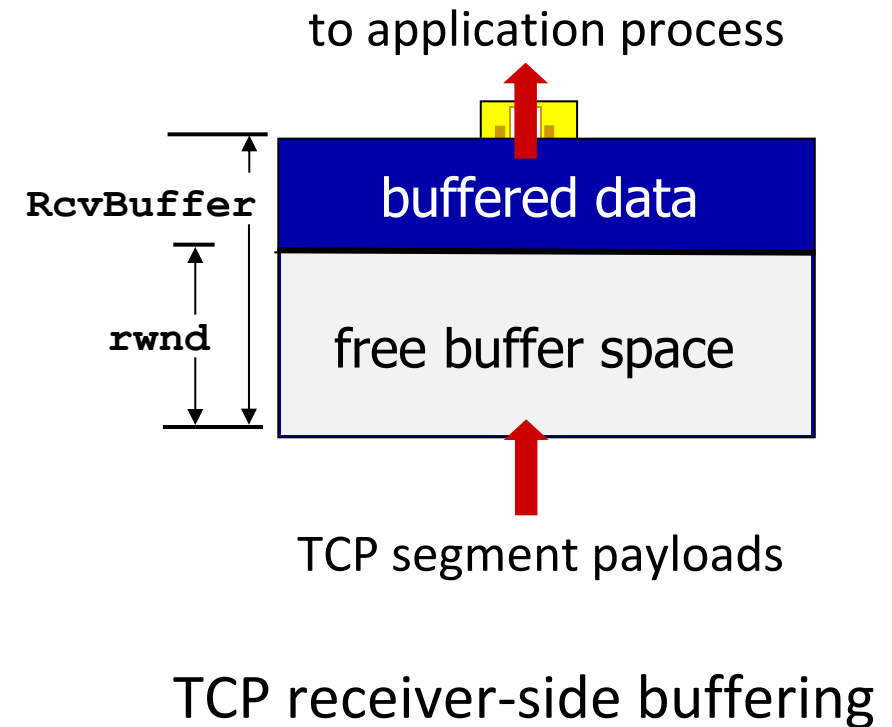
# TCP Flow Control



**sender how much available  
space we have in our buffer!**

# TCP Flow Control

- TCP receiver “advertises” free buffer space in **rwnd** field in TCP header
  - **RcvBuffer** size set via socket options (typical default is 4096 bytes)
  - many operating systems auto-adjust **RcvBuffer**
- sender limits amount of unACKed (“in-flight”) data to received **rwnd**
- guarantees receive buffer will not overflow





# TCP Flow Control

Wireshark ?

# TCP Flow Control

[https://media.pearsoncmg.com/aw/ecs\\_kurose\\_compnetwork\\_7/cw/content/interactiveanimations/flow-control/index.html](https://media.pearsoncmg.com/aw/ecs_kurose_compnetwork_7/cw/content/interactiveanimations/flow-control/index.html)

# TCP Timer

What is a good way to determine when to timeout? (aka the length of timer)

1. Too short: premature timeout, unnecessary retransmissions
2. Too long: slow reaction to segment loss

The TCP timeout value should around the same time it take to ....

# TCP Timer

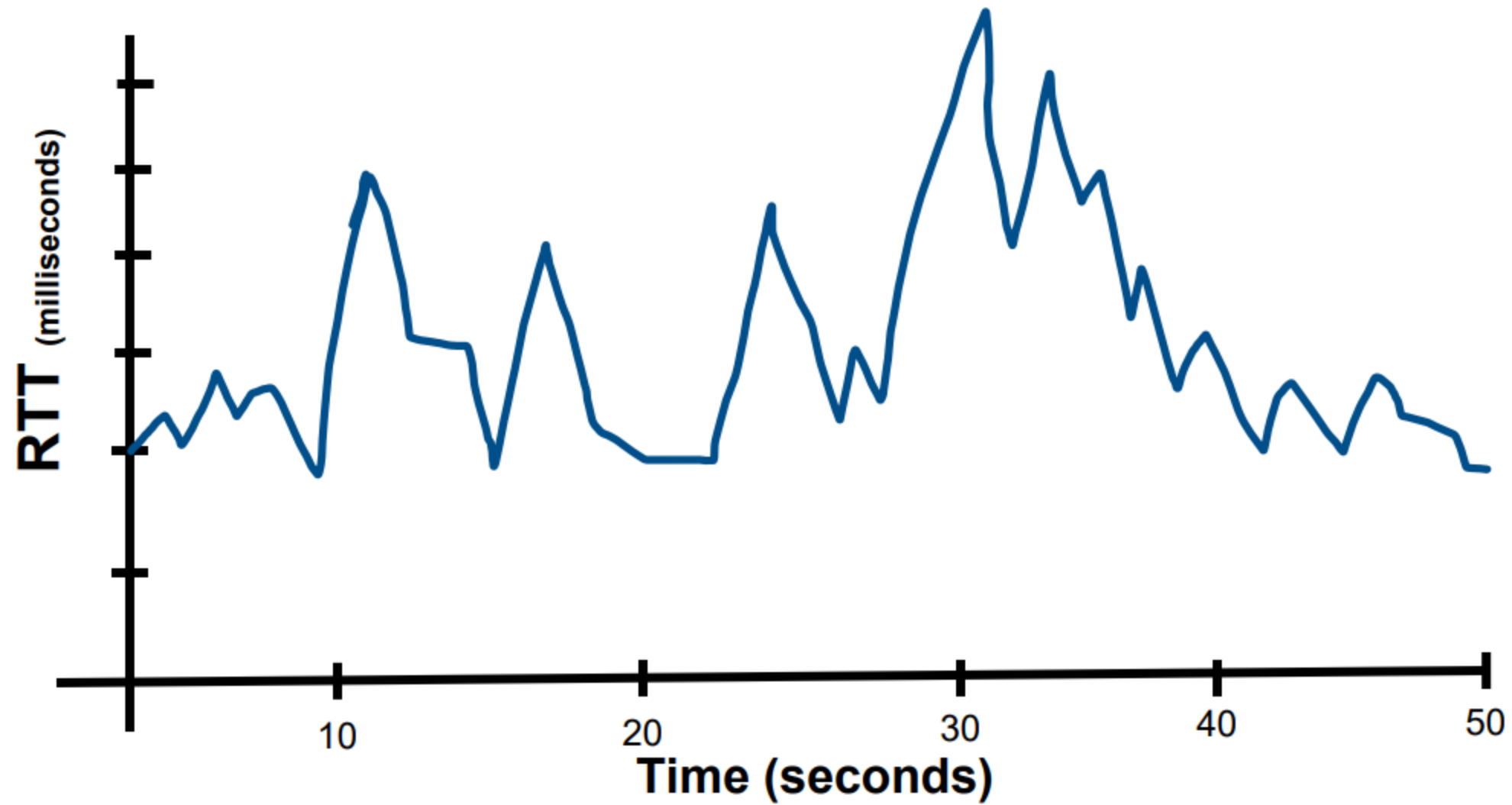
What is a good way to determine when to timeout? (aka the length of timer)

1. Too short: premature timeout, unnecessary retransmissions
2. Too long: slow reaction to segment loss

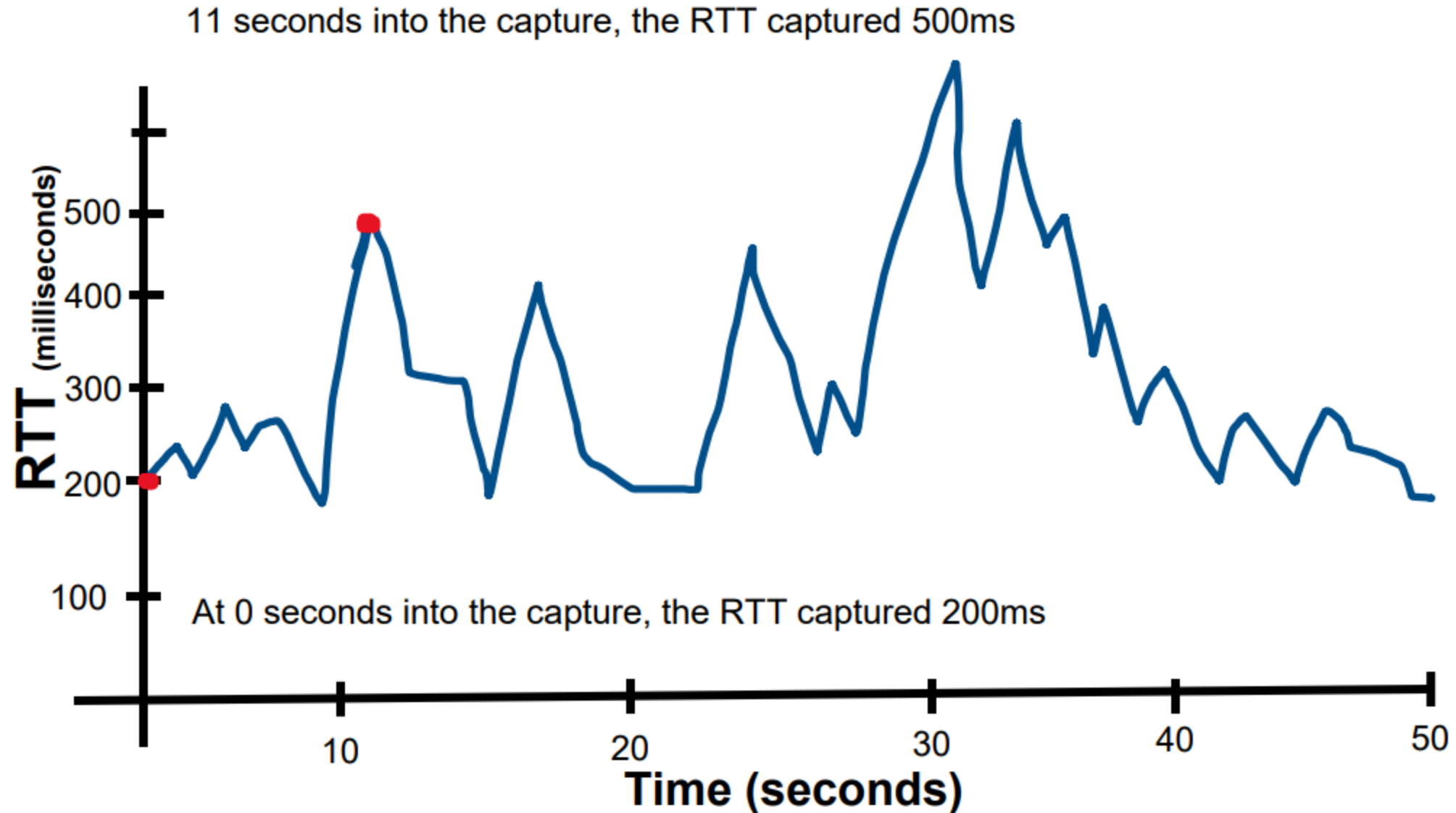
The TCP timeout value should around the same time it take to receive an acknowledgement on a sent packet (on average)

Let's consider setting it to be a dynamic value!

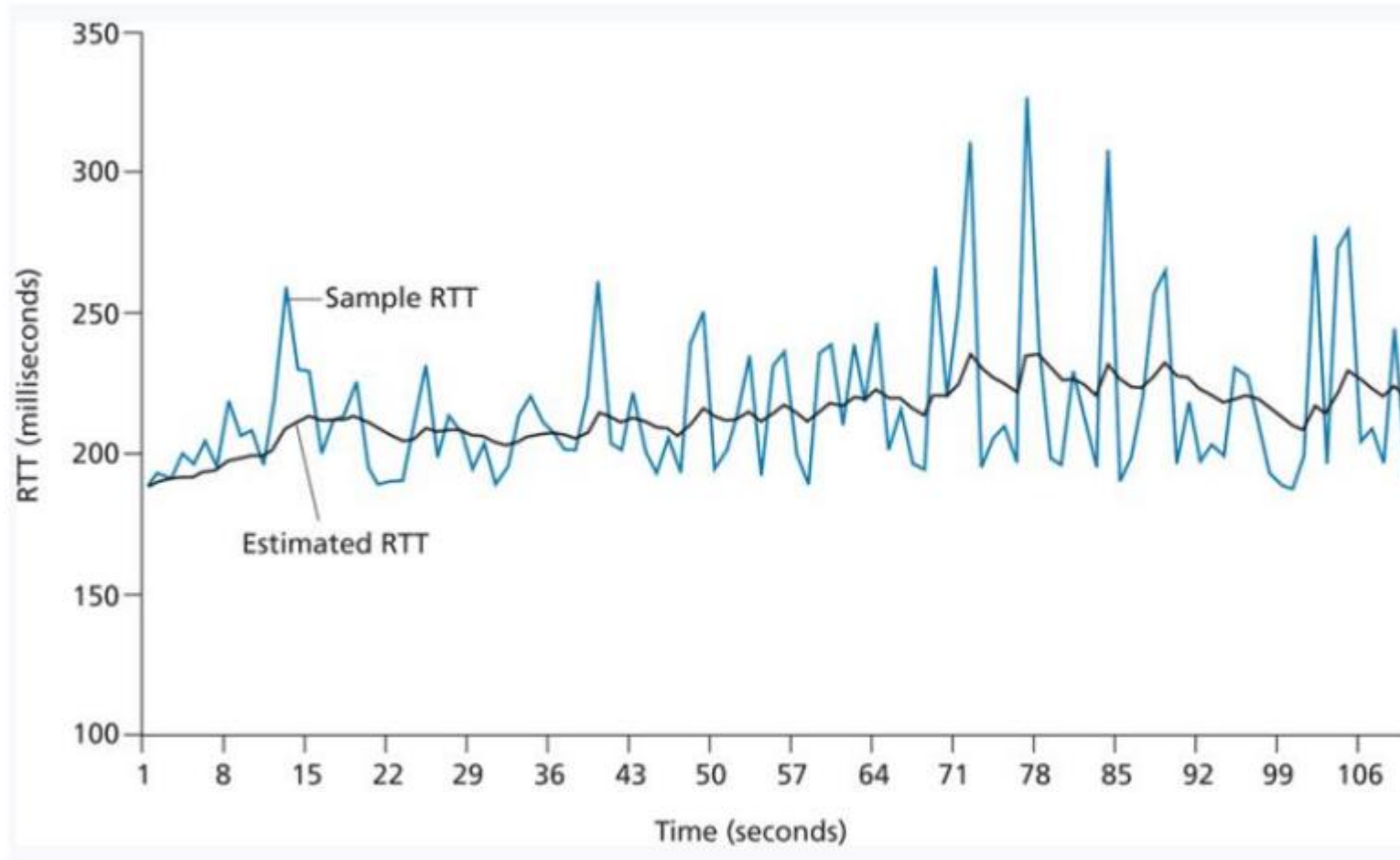
# TCP Timeout



# TCP Timeout



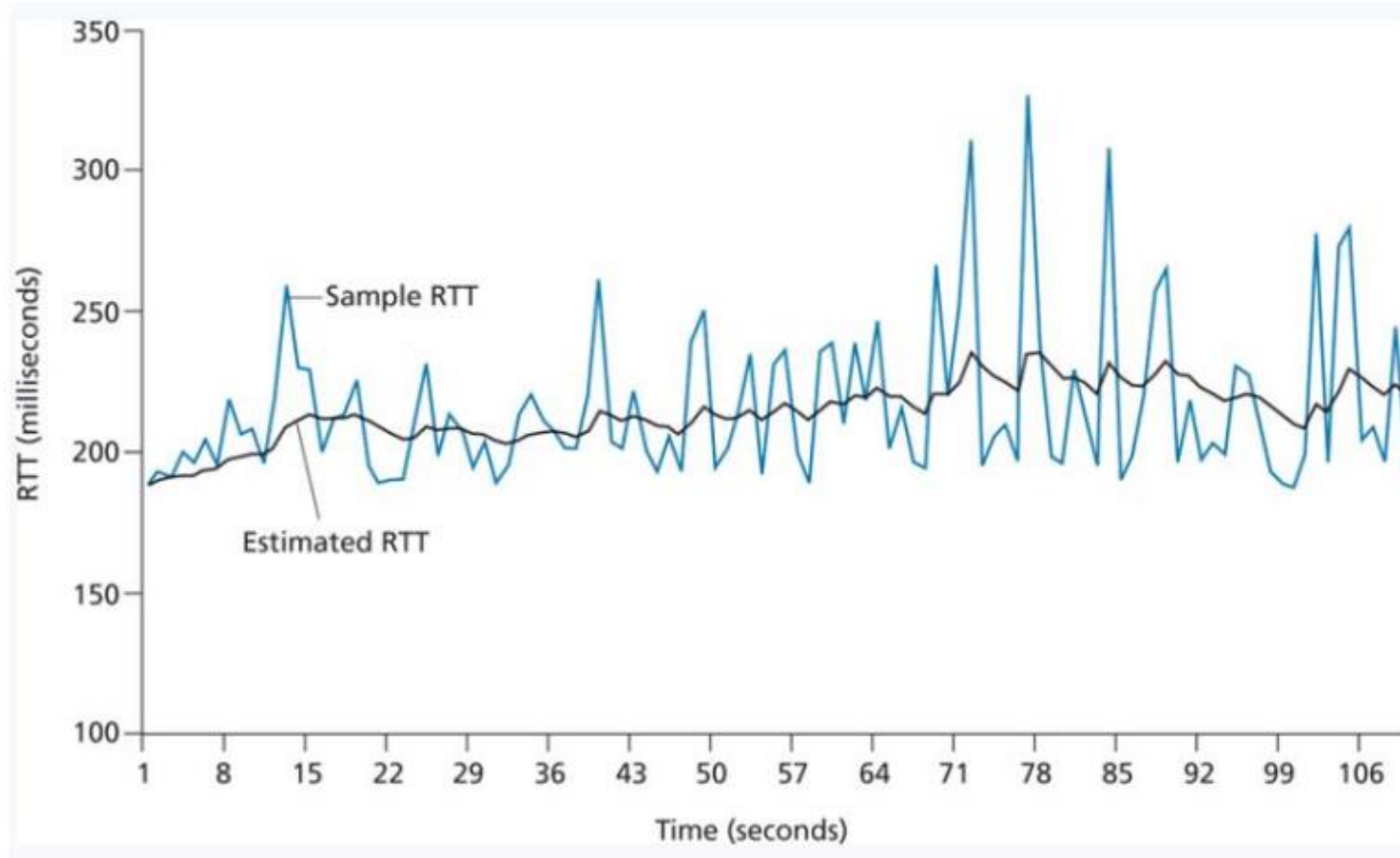
# TCP Timeout



$$\text{EstimatedRTT} = (1 - \alpha) \cdot \text{EstimatedRTT} + \alpha \cdot \text{SampleRTT}$$

$$\alpha = 0.125$$

# TCP Timeout



In addition, we also want  
some kind of safety margin

$$\text{DevRTT} = (1-\beta) * \text{DevRTT} + \beta * |\text{SampleRTT} - \text{EstimatedRTT}|$$

(typically,  $\beta = 0.25$ )

TimeoutInterval =

$$\text{EstimatedRTT} + 4 * \text{DevRTT}$$

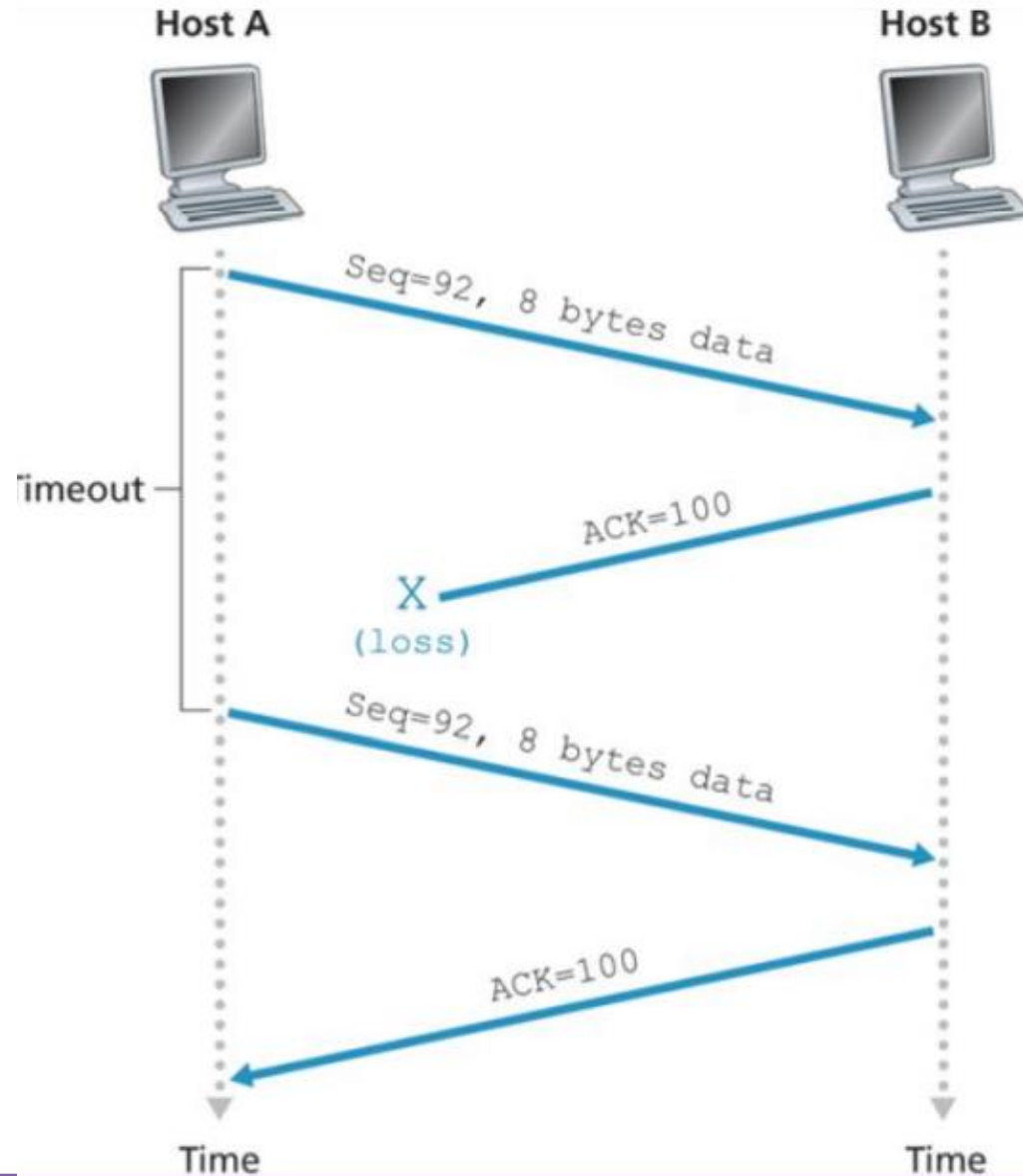
(safety margin)

$$\text{EstimatedRTT} = (1 - \alpha) \cdot \text{EstimatedRTT} + \alpha \cdot \text{SampleRTT}$$

$$\alpha = 0.125$$

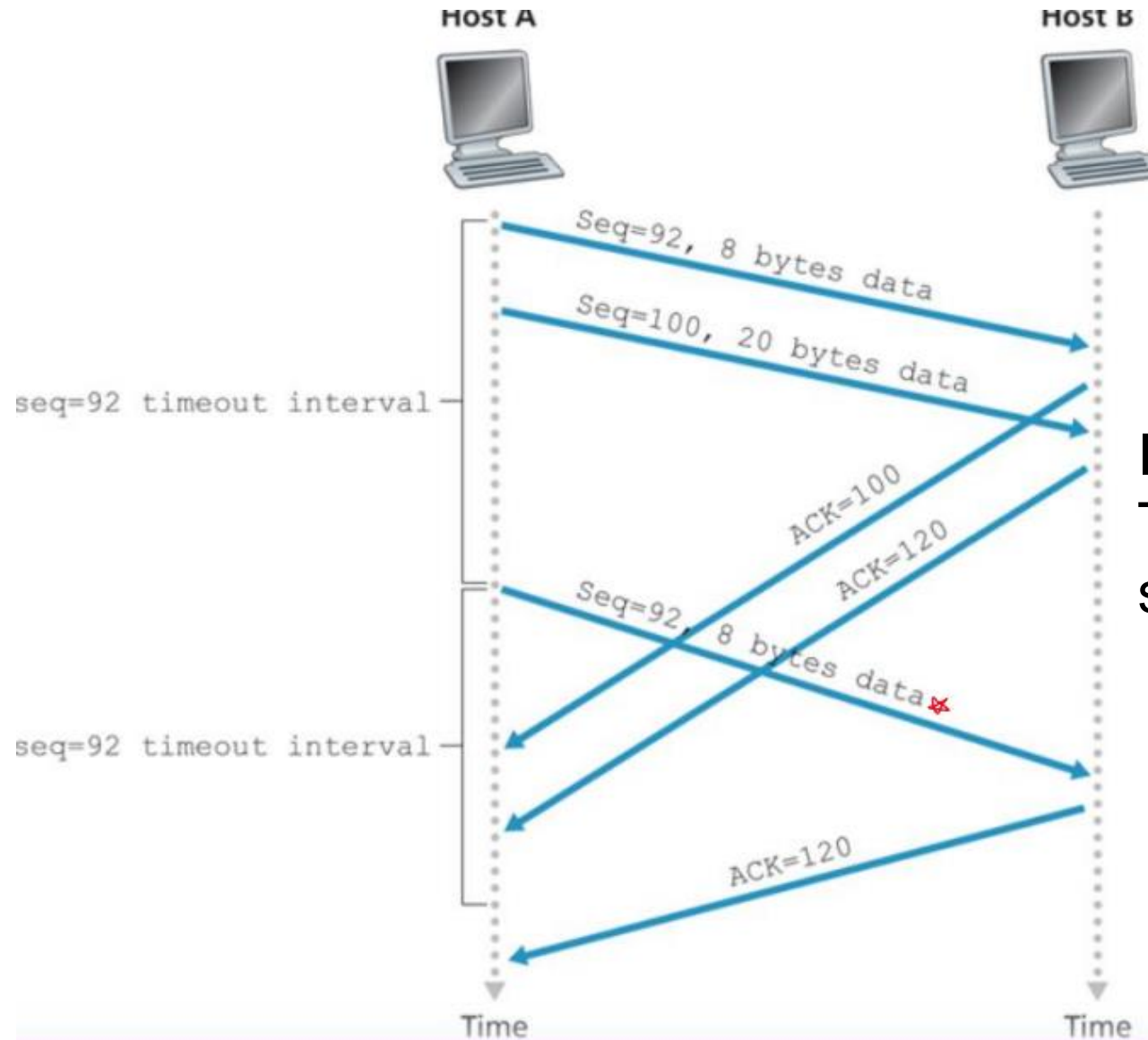


# TCP Timeout



TCP retransmits on ACK loss

# TCP Timeout



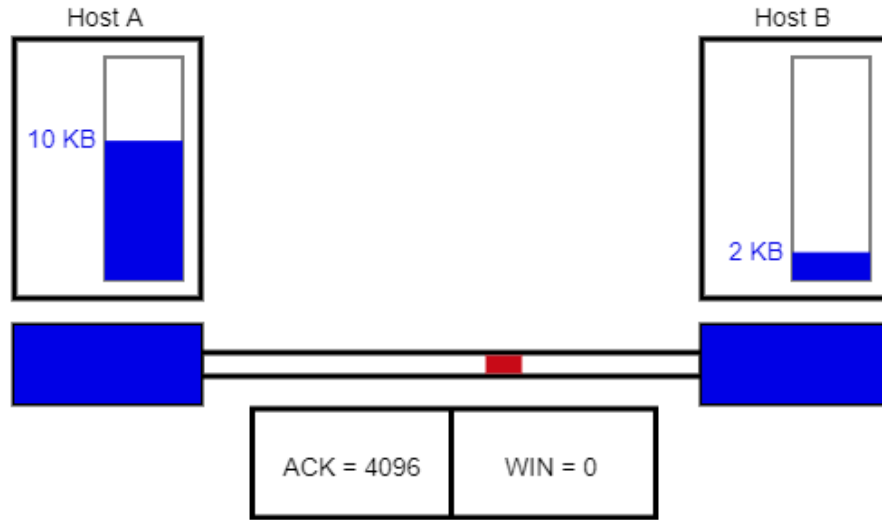
If multiple ACKS are lost/late, TCP only resends the first segment in the sequence

# TCP Timeout

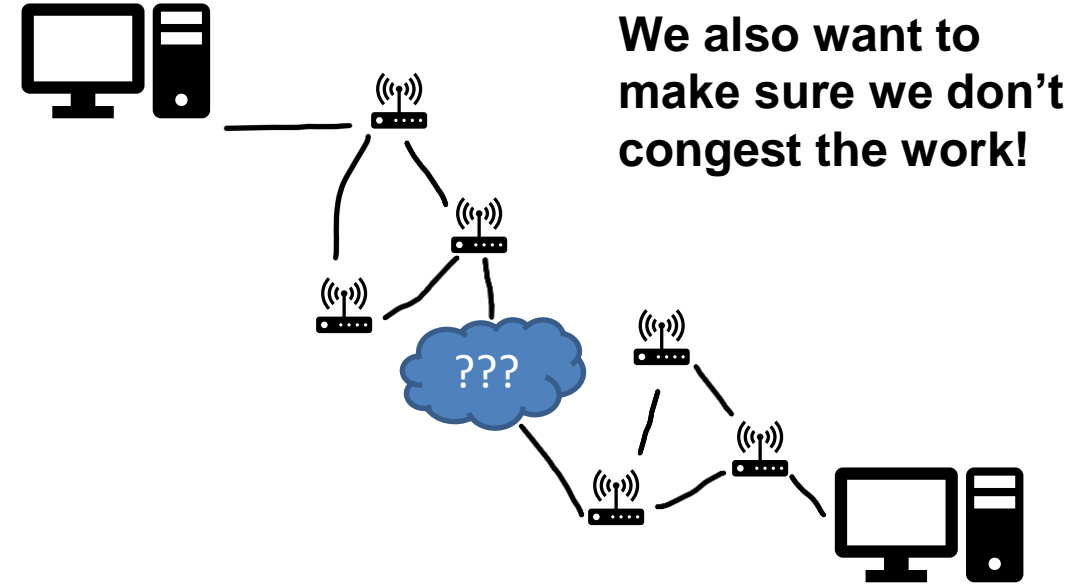
Event	TCP Receiver Action
Arrival of in-order segment with expected sequence number. All data up to expected sequence number already acknowledged.	Delayed ACK. Wait up to 500 msec for arrival of another in-order segment. If next in-order segment does not arrive in this interval, send an ACK.
Arrival of in-order segment with expected sequence number. One other in-order segment waiting for ACK transmission.	Immediately send single cumulative ACK, ACKing both in-order segments.
Arrival of out-of-order segment with higher-than-expected sequence number. Gap detected.	Immediately send duplicate ACK, indicating sequence number of next expected byte (which is the lower end of the gap).
Arrival of segment that partially or completely fills in gap in received data.	Immediately send ACK, provided that segment starts at the lower end of gap.

Actions for specific TCP events are described well in TCP documentation

# TCP Congestion Control



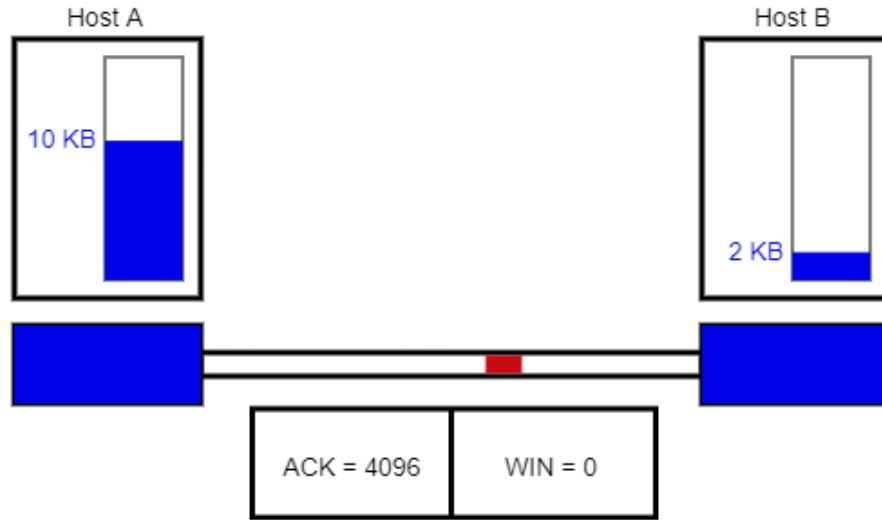
TCP sends back amount of available buffer space in the receiver  
This helps make sure we don't overwhelm the receiver



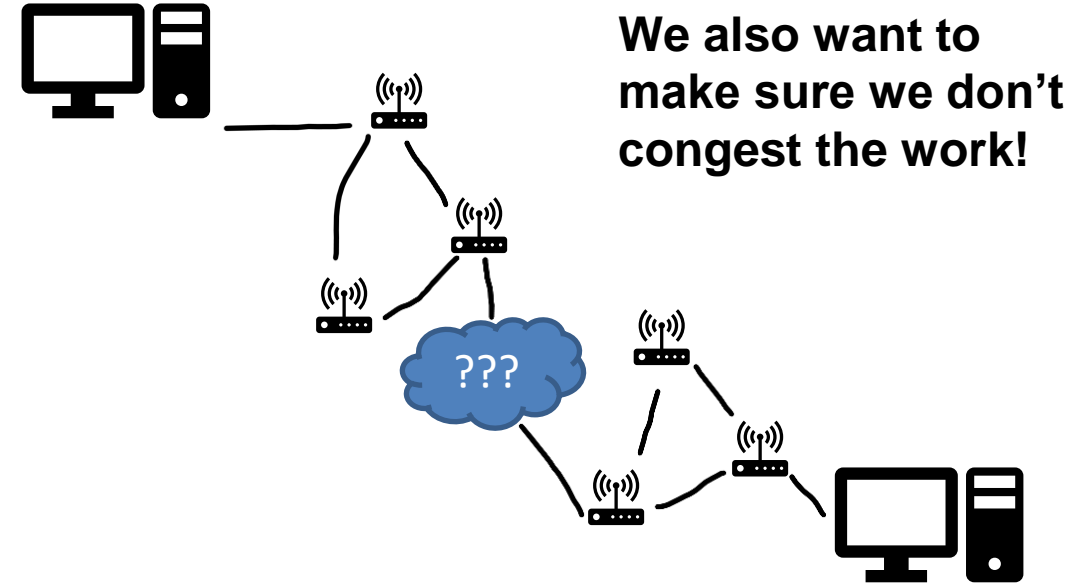
Issues:

- If the network is congested, we want to slow down our sending rate
- If the network is not congested, we should try to send more stuff

# TCP Congestion Control



TCP sends back amount of available buffer space in the receiver  
This helps make sure we don't overwhelm the receiver

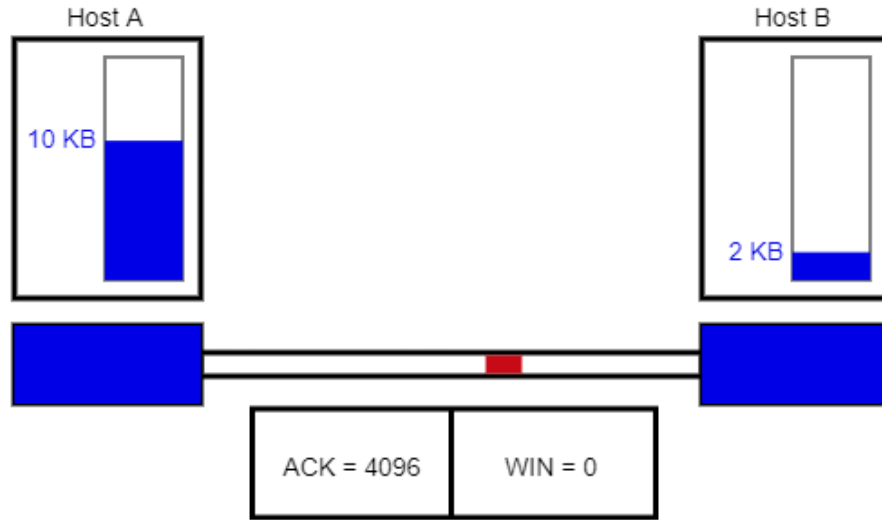


Issues:

- If the network is congested, we want to slow down our sending rate
- If the network is not congested, we should try to send more stuff

From the sender perspective, how could we measure how congested the network is?

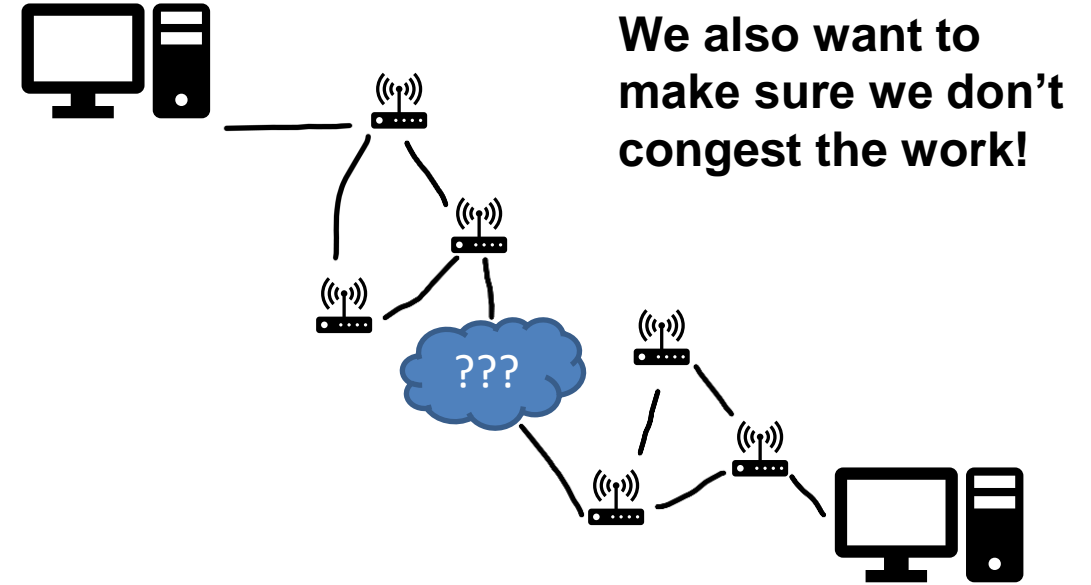
# TCP Congestion Control



TCP sends back amount of available buffer space in the receiver  
This helps make sure we don't overwhelm the receiver

Some ways we could measure how congested the network is

- See how many dropped packets we are getting
- Amount of duplicate ACKs received
- Amount of UnAcked packets



Issues:

- If the network is congested, we want to slow down our sending rate
- If the network is not congested, we should try to send more stuff

# TCP Congestion Control

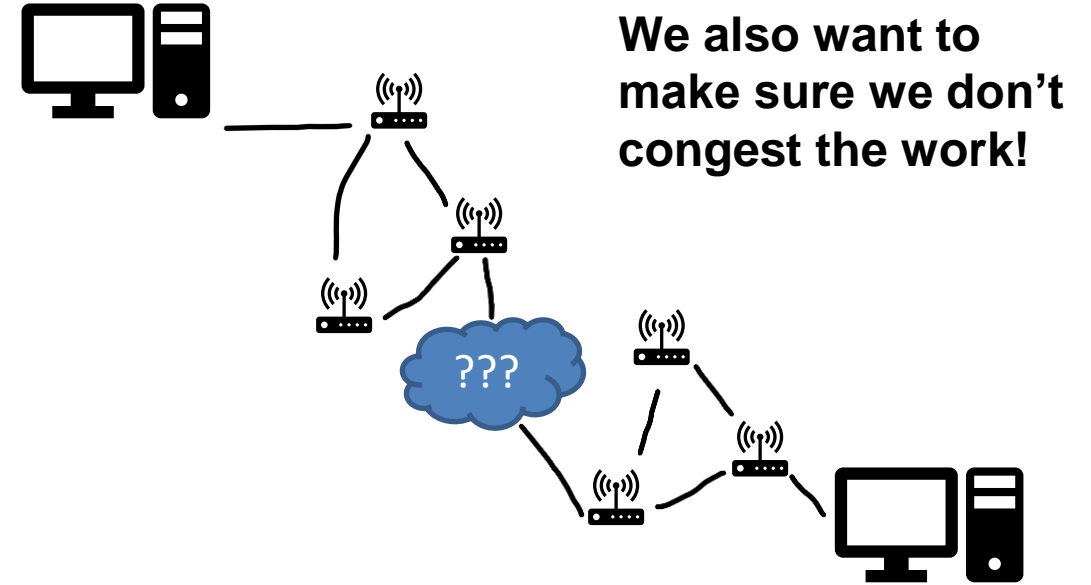
TCP sender also has a **congestion window (cwnd)**, which controls the amount of unAck'd that can be sent out

TCP is **self-clocking**

(It uses acknowledgements to trigger, or clock, its increase in congestion window size)

The amount of unacknowledged data at a sender may not exceed the *minimum* of the congestion window and receiving window

$$\text{LastByteSent} - \text{LastByteAcked} \leq \min\{\text{cwnd}, \text{rwnd}\}$$

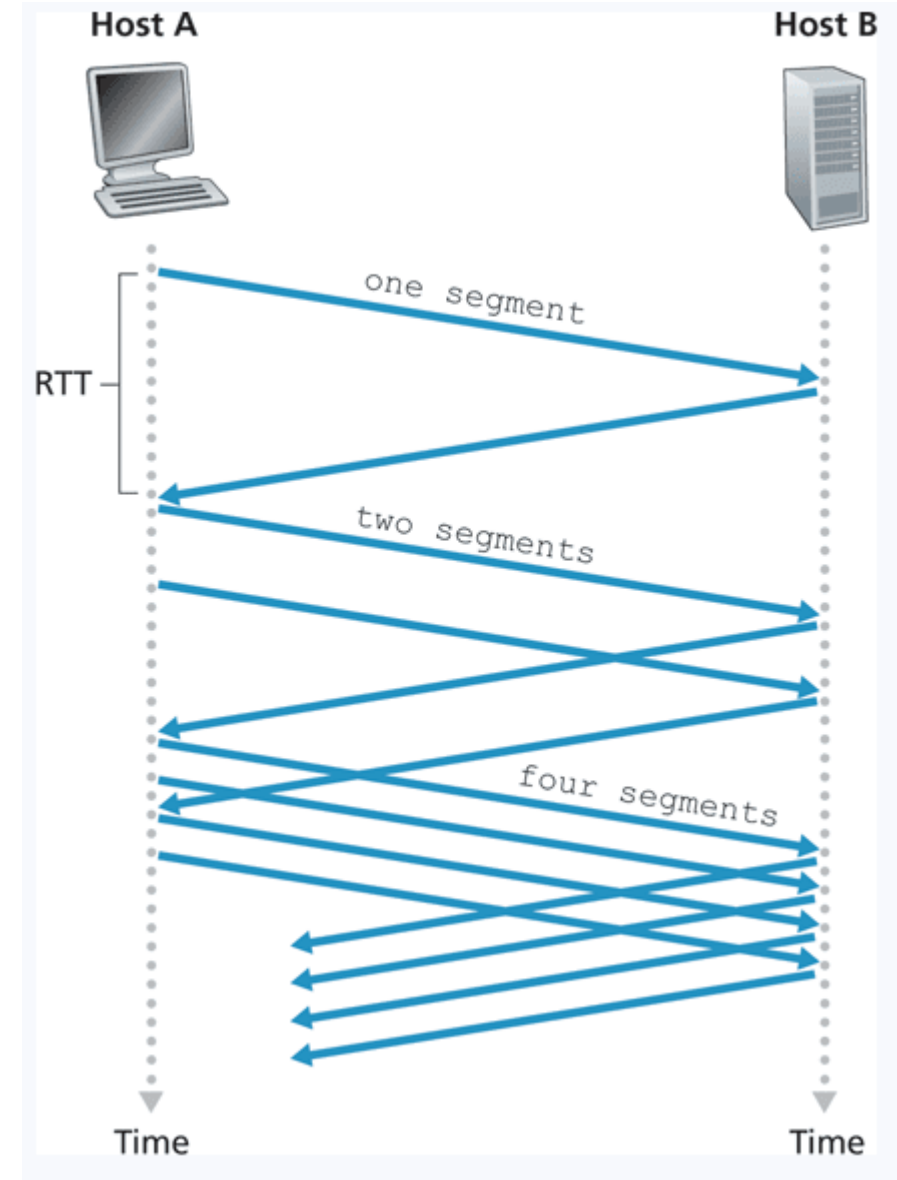


# TCP Congestion Control Algorithm

TCP Algorithm to prevent network congestion

- **Slow Start**
- Congestion Avoidance
- Fast recovery

Start sending slow, but exponentially grows up to a *threshold*





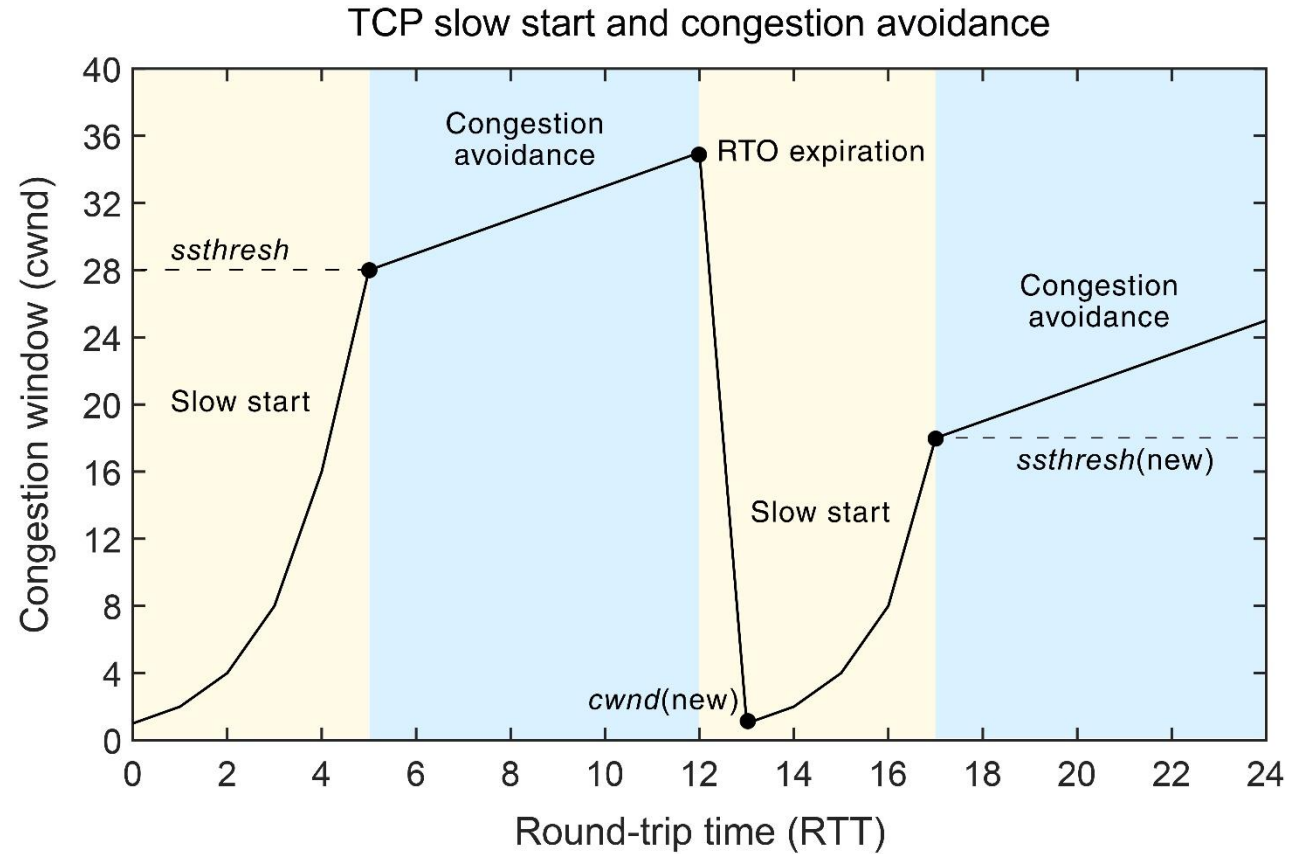
# TCP Congestion Control Algorithm

## TCP Algorithm to prevent network congestion

- Slow Start
- **Congestion Avoidance**
- Fast recovery

Linearly increase congestion window for each ACK received

When a loss event occurs, significantly decrease congestion window and slow down transmission rate, and enter **fast recovery**

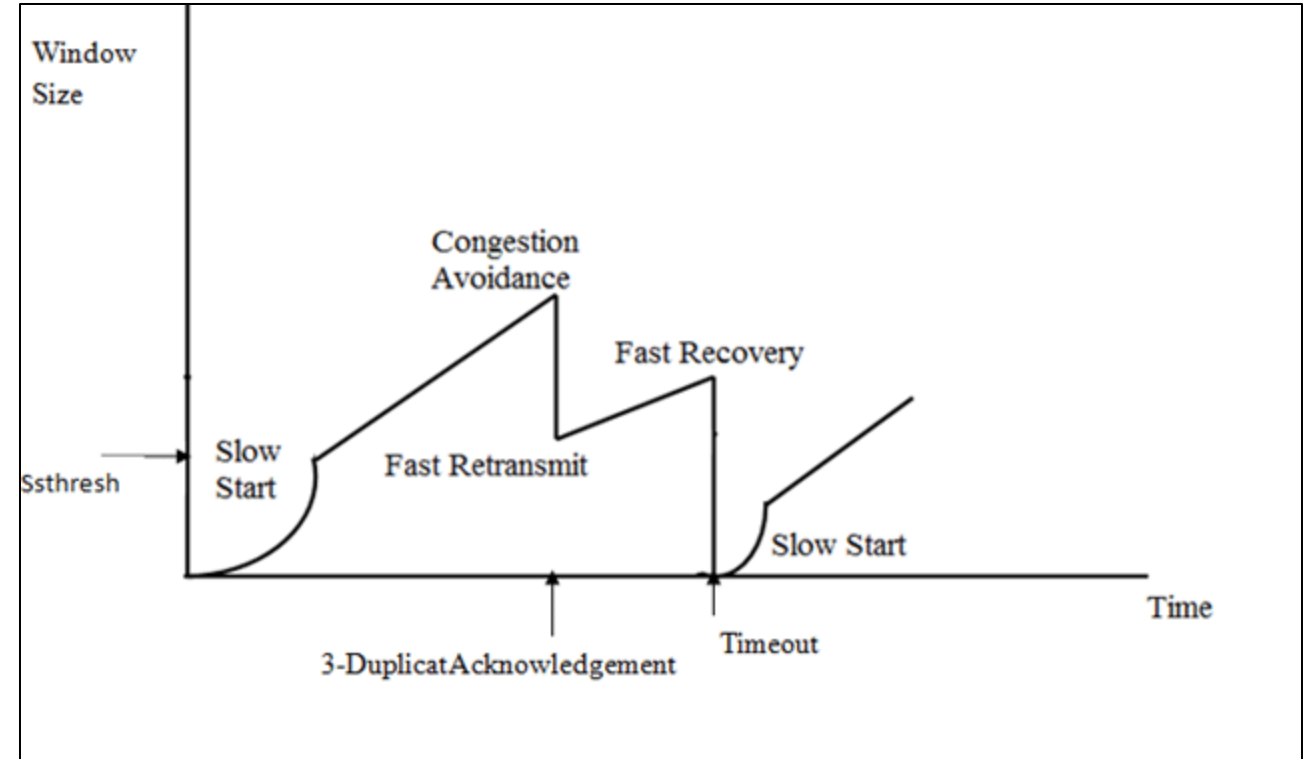


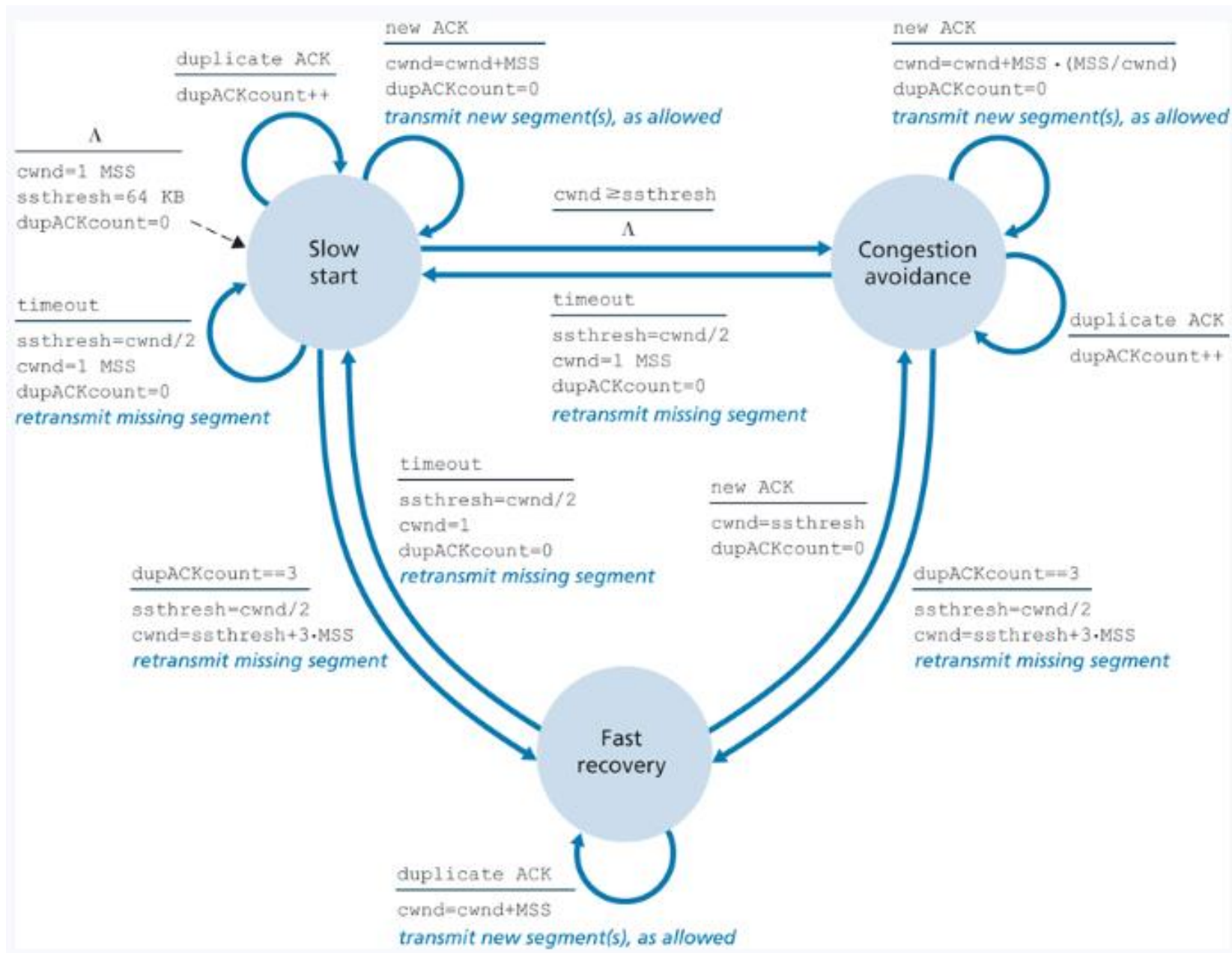
# TCP Congestion Control Algorithm

TCP Algorithm to prevent network congestion

- Slow Start
- Congestion Avoidance
- **Fast recovery**

Upon knowledge of packet loss, throttle the TCP connection and start off slow again





# TCP Congestion Control Algorithm

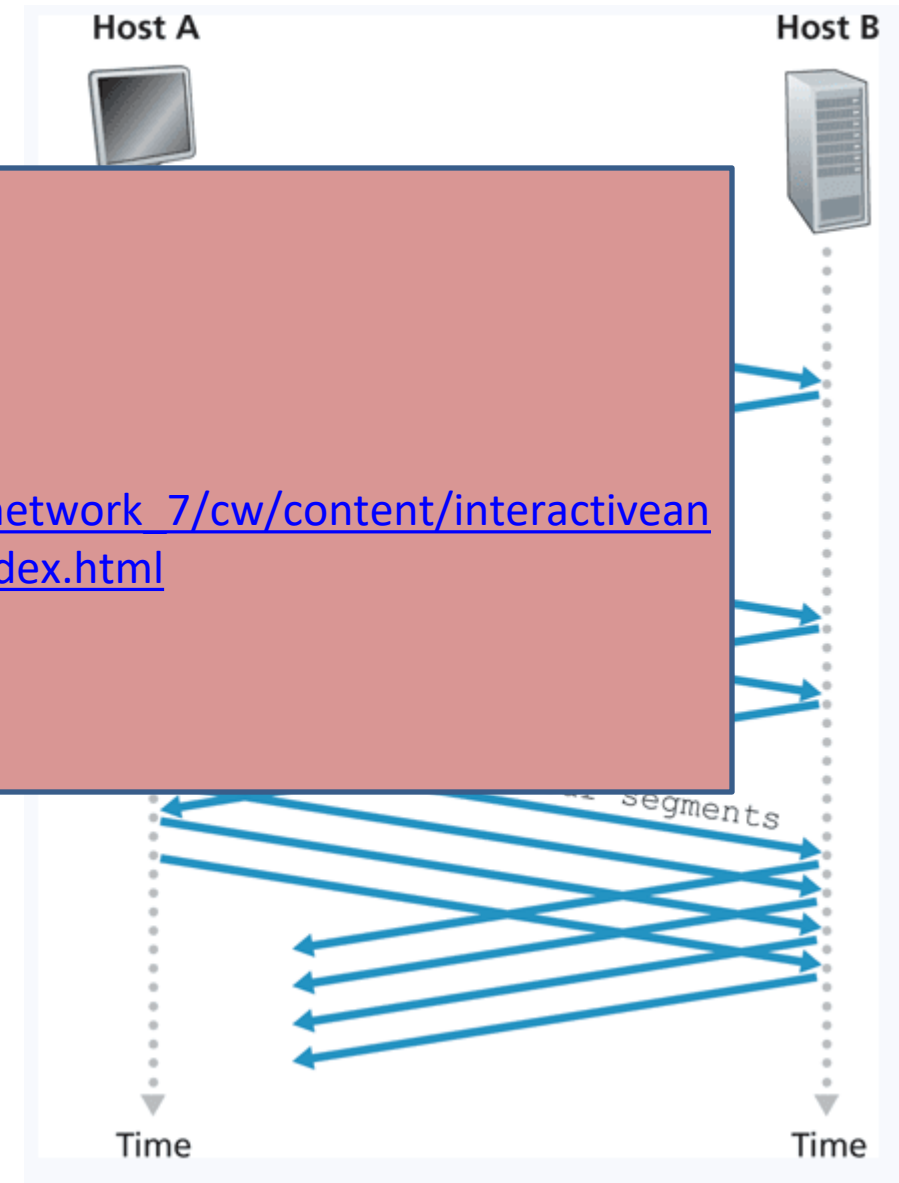
TCP Alg

- Slow
- Cong
- Fast

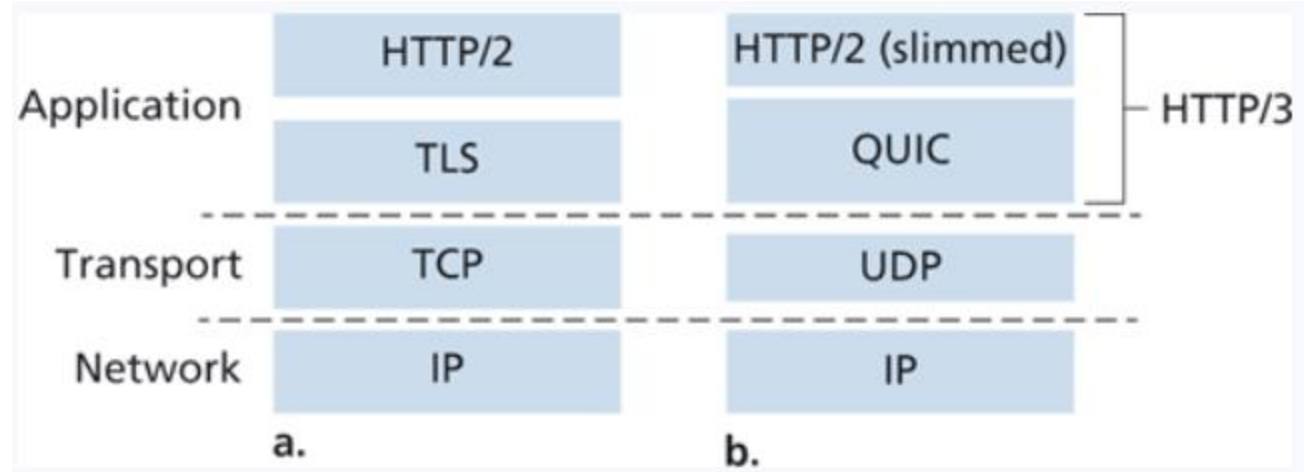
Upon knowledge of packet loss, throttle the TCP connection and start off slow again

Animation time!

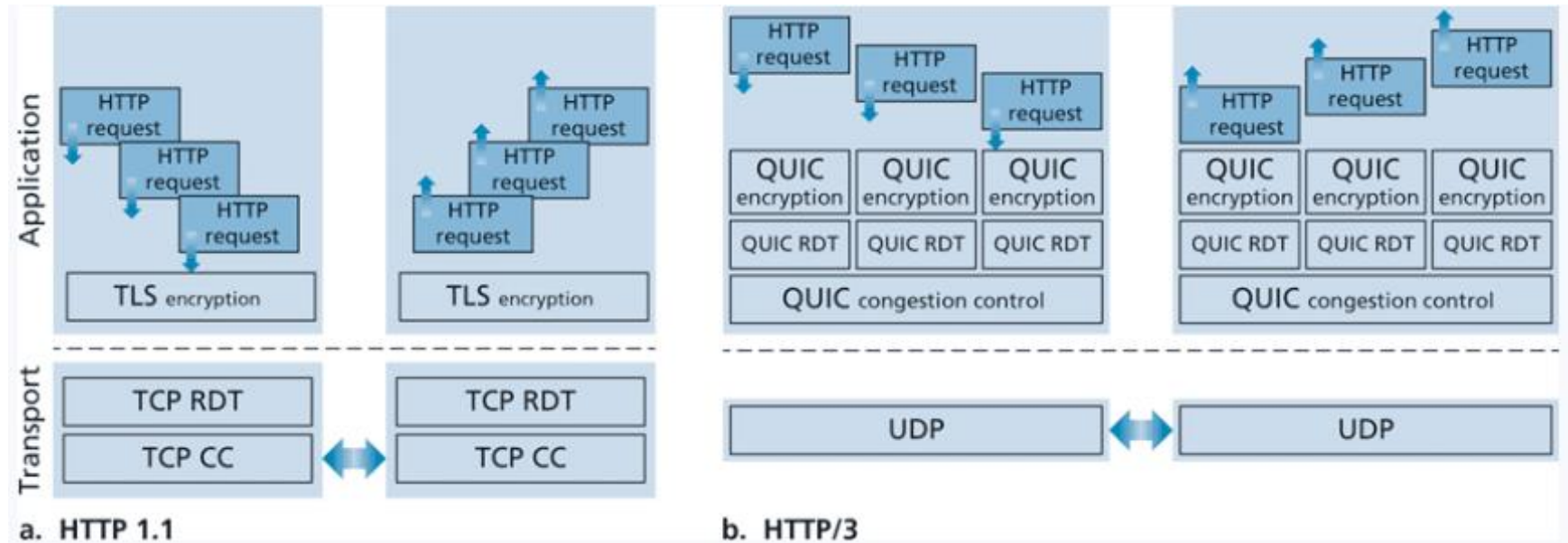
[https://media.pearsoncmg.com/aw/ecs\\_kurose\\_compnetwork\\_7/cw/content/interactiveanimations/tcp-congestion/index.html](https://media.pearsoncmg.com/aw/ecs_kurose_compnetwork_7/cw/content/interactiveanimations/tcp-congestion/index.html)



# Current transport layer implementation



Transport layer protocols and congestion control is still a heavily researched area!



RFCs (Request for Comments) documents and describes the details and standards of how internet protocols (such as HTTP, TCP, UDP) should work

## TCP- RFC 793

TRANSMISSION CONTROL PROTOCOL

DARPA INTERNET PROGRAM

PROTOCOL SPECIFICATION

September 1981

## UDP- RFC 768

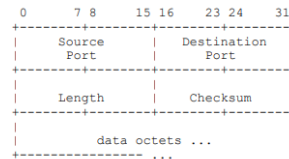
User Datagram Protocol

Introduction

This User Datagram Protocol (UDP) is defined to make available a datagram mode of packet-switched computer communication in the environment of an interconnected set of computer networks. This protocol assumes that the Internet Protocol (IP) [1] is used as the underlying protocol.

This protocol provides a procedure for application programs to send messages to other programs with a minimum of protocol mechanism. The protocol is transaction oriented, and delivery and duplicate protection are not guaranteed. Applications requiring ordered reliable delivery of streams of data should use the Transmission Control Protocol (TCP) [2].

Format



User Datagram Header Format

## DNS- RFC 1035

DOMAIN NAMES - IMPLEMENTATION AND SPECIFICATION

### 1. STATUS OF THIS MEMO

This RFC describes the details of the domain system and protocol, and assumes that the reader is familiar with the concepts discussed in a companion RFC, "Domain Names - Concepts and Facilities" [RFC-1034].

The domain system is a mixture of functions and data types which are an official protocol and functions and data types which are still experimental. Since the domain system is intentionally extensible, new data types and experimental behavior should always be expected in parts of the system beyond the official protocol. The official protocol parts include standard queries, responses and the Internet class RR data formats (e.g., host addresses). Since the previous RFC set, several definitions have changed, so some previous definitions are obsolete.



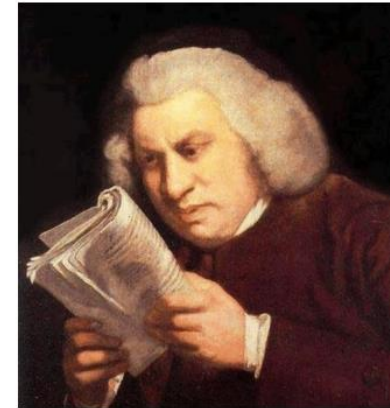
## A Standard for the Transmission of IP Datagrams on Avian Carriers

### Status of this Memo

This memo describes an experimental method for the encapsulation of IP datagrams in avian carriers. This specification is primarily useful in Metropolitan Area Networks. This is an experimental, not recommended standard. Distribution of this memo is unlimited.

### Overview and Rational

[Avian carriers can provide high delay, low throughput, and low altitude service.] The connection topology is limited to a single point-to-point path for each carrier, used with standard carriers, but many carriers can be used without significant interference with each other, outside of early spring. This is because of the 3D ether space available to the carriers, in contrast to the 1D ether used by IEEE802.3. The carriers have an intrinsic collision avoidance system, which increases availability. Unlike some network technologies, such as packet radio, communication is not limited to line-of-sight distance. Connection oriented service is available in some cities, usually based upon a central hub topology.



# FIN