

TRANSPORT LAYER

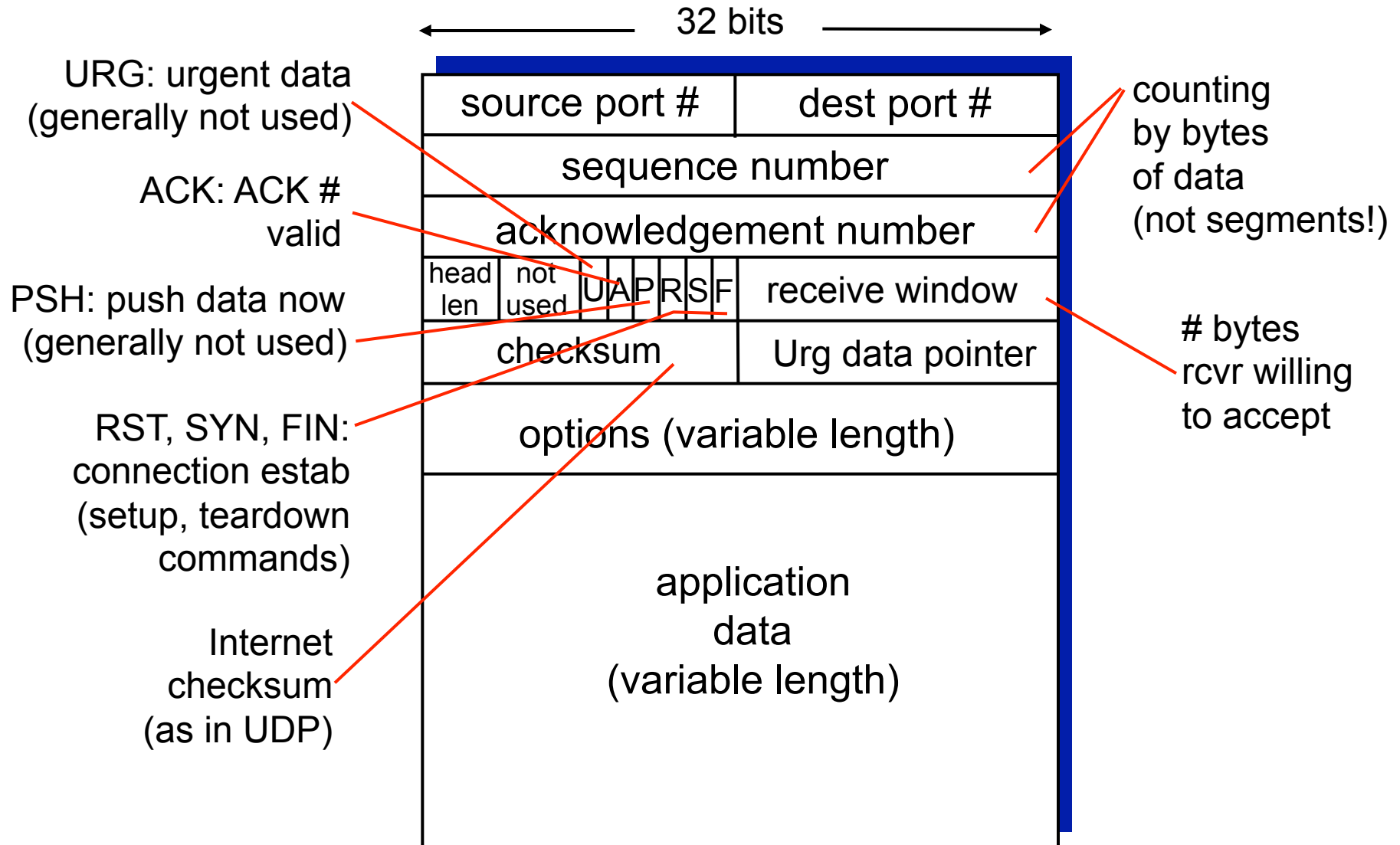
Choose your weapon

TCP: Overview

RFCs: 793, 1122, 1323, 2018, 2581

- **point-to-point:**
 - ▣ one sender, one receiver
- **reliable, in-order *byte stream*:**
 - ▣ no “message boundaries”
- **pipelined:**
 - ▣ TCP congestion and flow control set window size
- **full duplex data:**
 - ▣ bi-directional data flow in same connection
 - ▣ MSS: maximum segment size
- **connection-oriented:**
 - ▣ handshaking (exchange of control msgs) initializes sender, receiver state before data exchange
- **flow controlled:**
 - ▣ sender will not overwhelm receiver

TCP segment structure



TCP seq. numbers, ACKs

sequence numbers:

- byte stream “number” of first byte in segment’s data

acknowledgements:

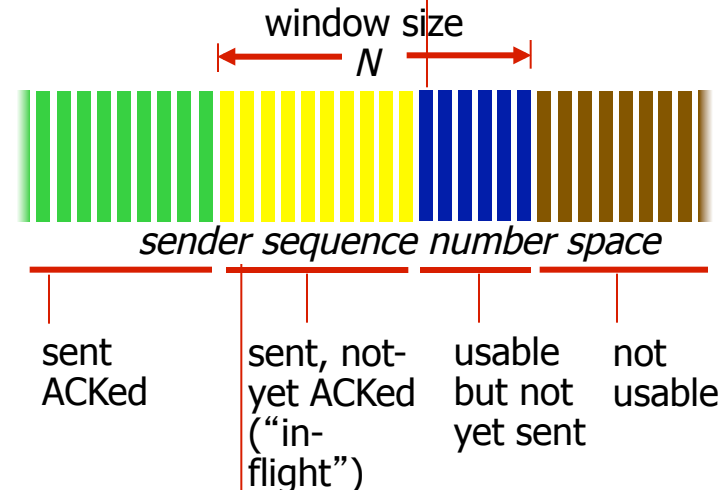
- seq # of next byte expected from other side
- cumulative ACK

Q: how receiver handles out-of-order segments

- A:** TCP spec doesn’t say, - up to implementor

outgoing segment from sender

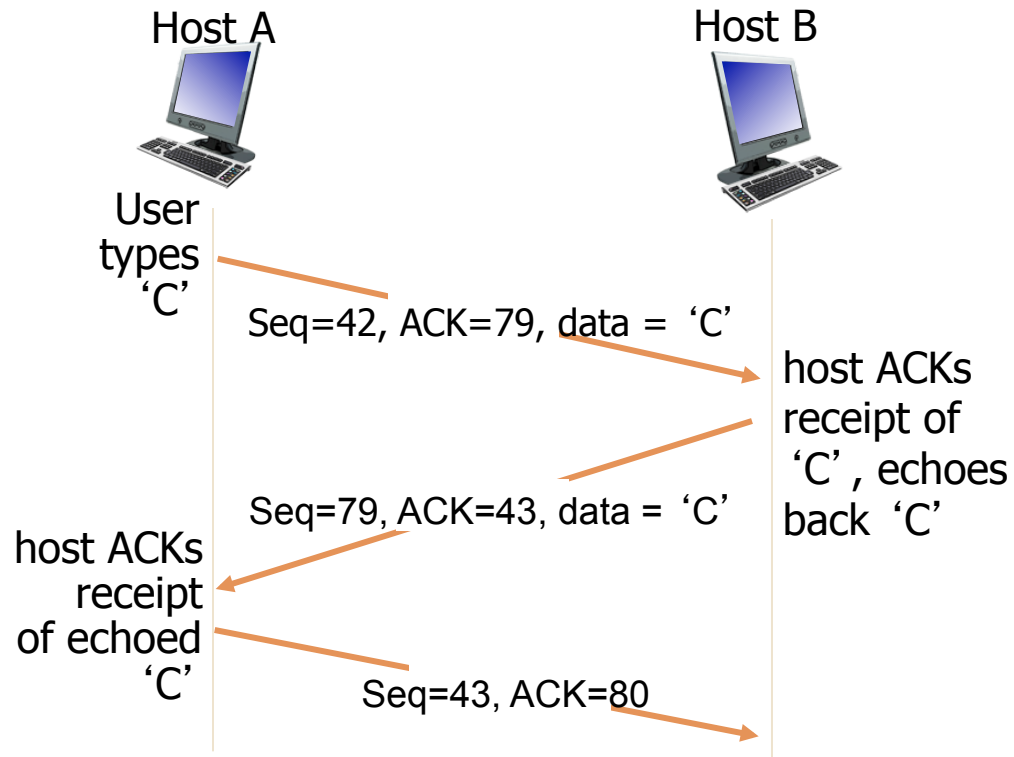
source port #	dest port #
sequence number	
acknowledgement number	
	rwnd
checksum	urg pointer



incoming segment to sender

source port #	dest port #
sequence number	
acknowledgement number	
	A
checksum	urg pointer

TCP seq. numbers, ACKs



simple telnet scenario

TCP round trip time, timeout

Q: how to set TCP timeout value?

- longer than RTT
 - ▣ but RTT varies
- *too short*: premature timeout, unnecessary retransmissions
- *too long*: slow reaction to segment loss

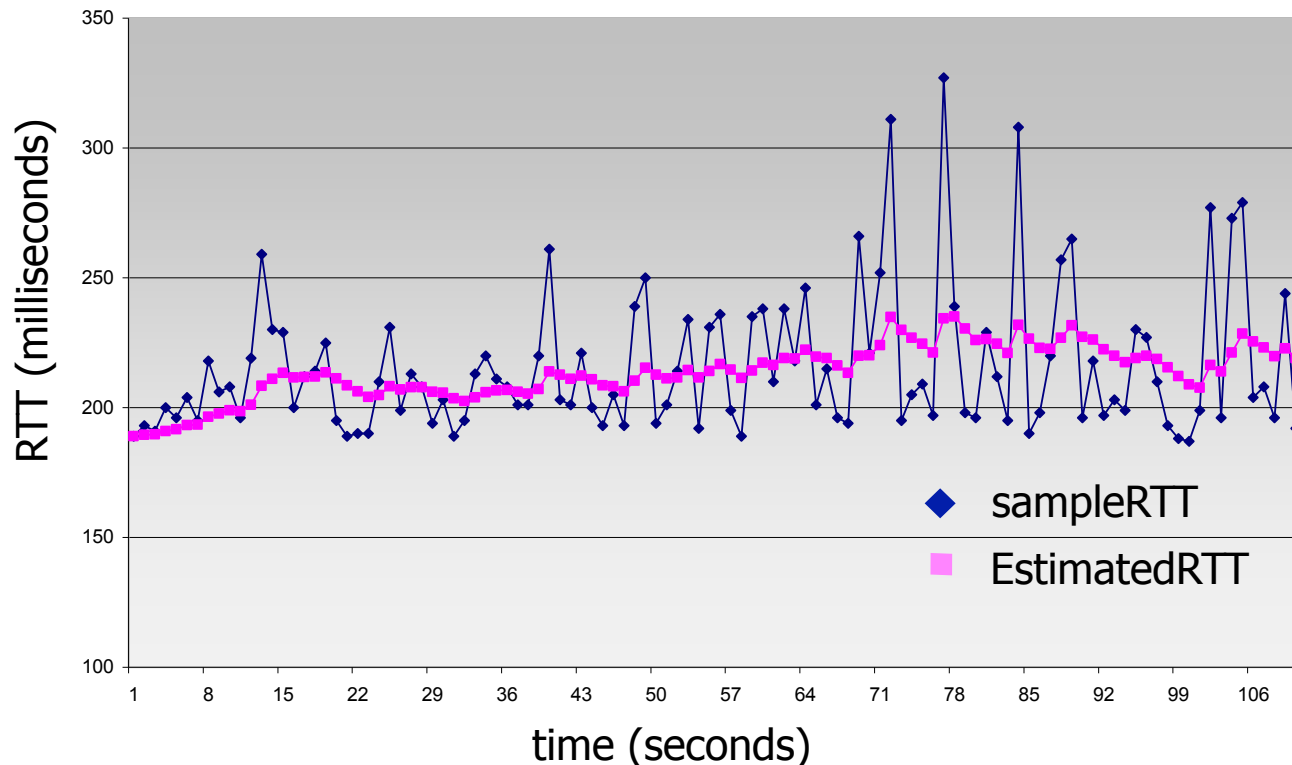
Q: how to estimate RTT?

- **SampleRTT**: measured time from segment transmission until ACK receipt
 - ▣ ignore retransmissions
- **SampleRTT** will vary, want estimated RTT “smoother”
 - ▣ average several *recent* measurements, not just current **SampleRTT**

TCP round trip time, timeout

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

- ❖ exponential weighted moving average
- ❖ influence of past sample decreases exponentially fast
- ❖ typical value: $\alpha = 0.125$



TCP reliable data transfer

- TCP creates rdt service on top of IP's unreliable service
 - ▣ pipelined segments
 - ▣ cumulative acks
 - ▣ single retransmission timer
- retransmissions triggered by:
 - ▣ timeout events
 - ▣ duplicate acks

let's initially consider simplified TCP sender:

- ▣ ignore duplicate acks
- ▣ ignore flow control, congestion control

TCP sender events:

data rcvd from app:

- create segment with seq #
- seq # is byte-stream number of first data byte in segment
- start timer if not already running
 - ▣ think of timer as for oldest unacked segment
 - ▣ expiration interval: `TimeoutInterval`

timeout:

- retransmit segment that caused timeout
- restart timer

ack rcvd:

- if ack acknowledges previously unacked segments
 - ▣ update what is known to be ACKed
 - ▣ start timer if there are still unacked segments

TCP ACK generation [RFC 1122, RFC 2581]

<i>event at receiver</i>	<i>TCP receiver action</i>
arrival of in-order segment with expected seq #. All data up to expected seq # already ACKed	delayed ACK. Wait up to 500ms for next segment. If no next segment, send ACK
arrival of in-order segment with expected seq #. One other segment has ACK pending	immediately send single cumulative ACK, ACKing both in-order segments
arrival of out-of-order segment higher-than-expect seq. # . Gap detected	immediately send <i>duplicate ACK</i> , indicating seq. # of next expected byte
arrival of segment that partially or completely fills gap	immediate send ACK, provided that segment starts at lower end of gap

TCP fast retransmit

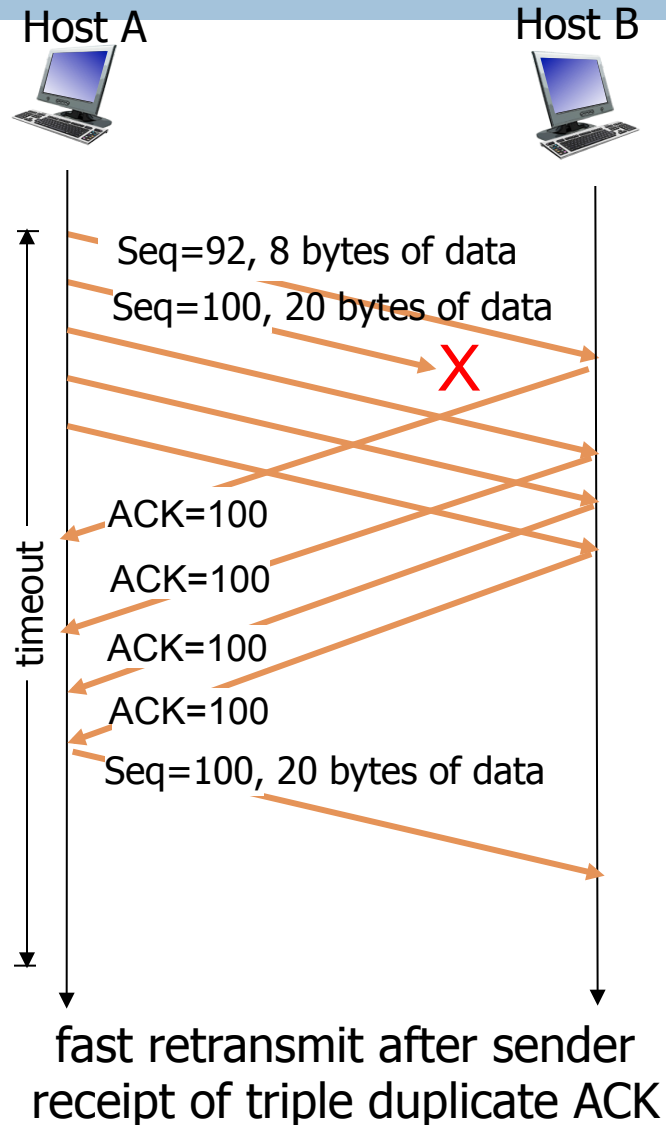
- time-out period often relatively long:
 - ▣ long delay before resending lost packet
- detect lost segments via duplicate ACKs.
 - ▣ sender often sends many segments back-to-back
 - ▣ if segment is lost, there will likely be many duplicate ACKs.

TCP fast retransmit

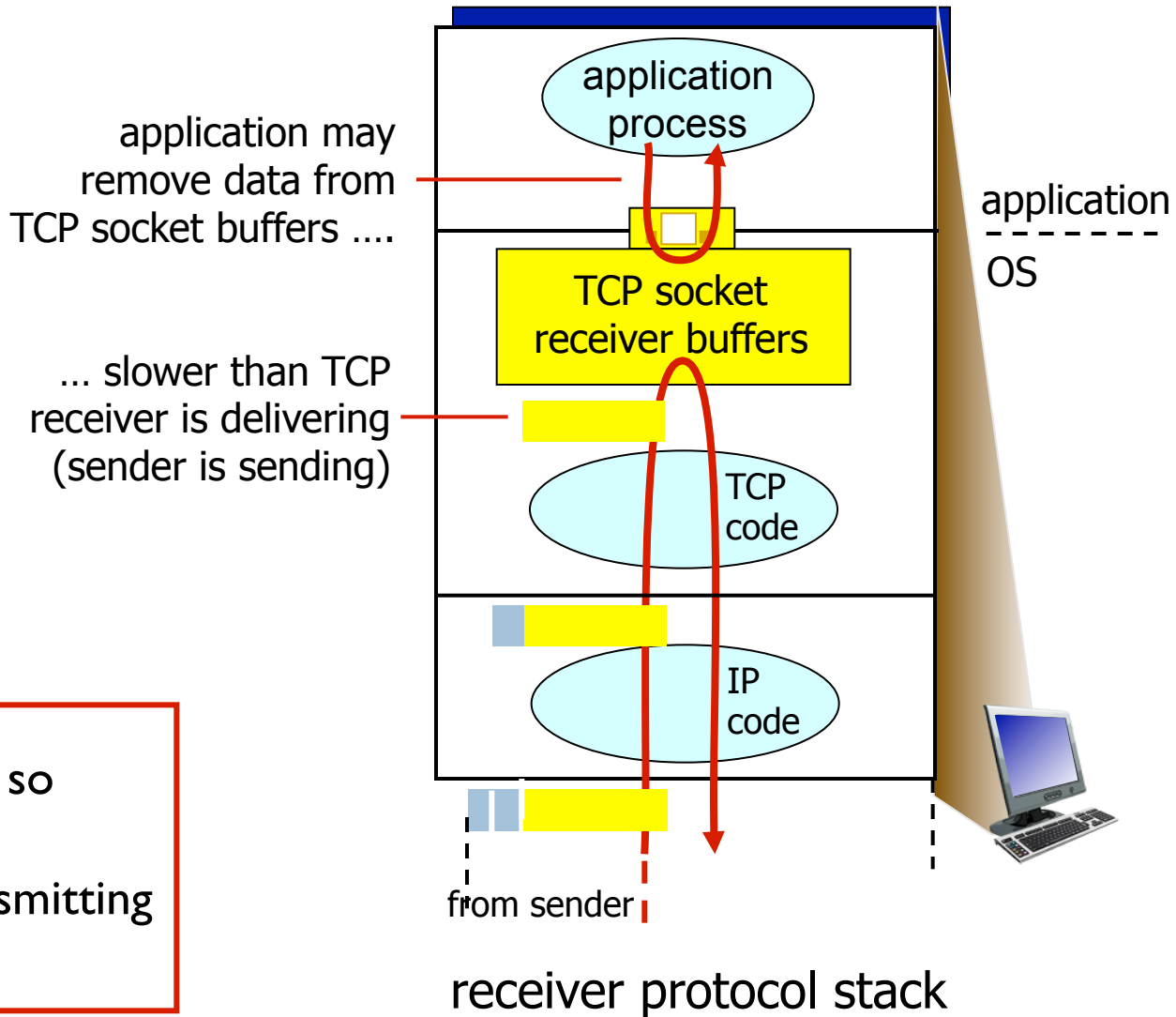
if sender receives 3 ACKs for same data (“triple duplicate ACKs”), resend unacked segment with smallest seq #

- likely that unacked segment lost, so don't wait for timeout

TCP fast retransmit



TCP flow control

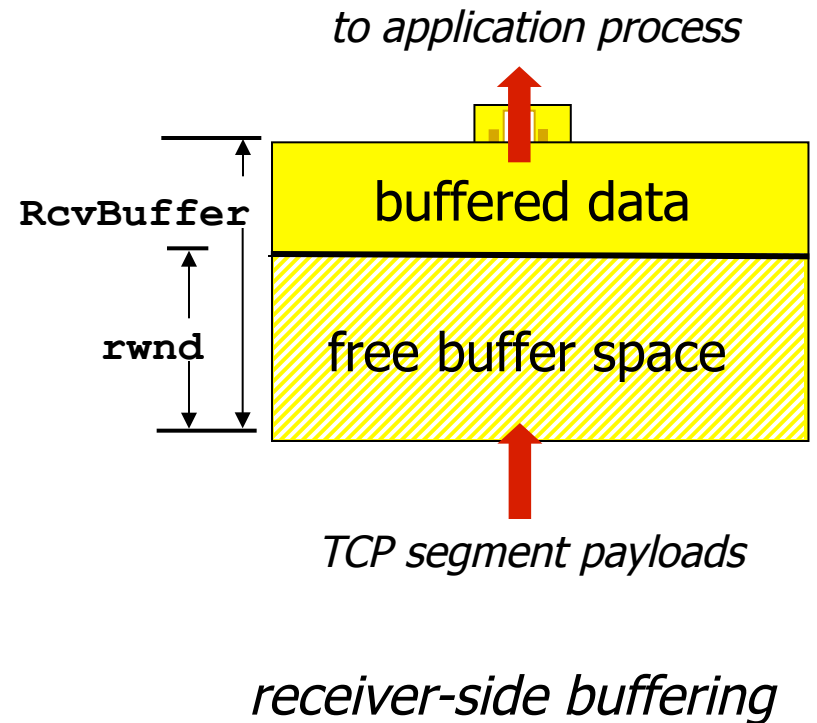


flow control

receiver controls sender, so
sender won't overflow
receiver's buffer by transmitting
too much, too fast

TCP flow control

- receiver “advertises” free buffer space by including **rwnd** value in TCP header of receiver-to-sender segments
 - ▣ **RcvBuffer** size set via socket options (typical default is 4096 bytes)
 - ▣ many operating systems autoadjust **RcvBuffer**
- sender limits amount of unacked (“in-flight”) data to receiver’s **rwnd** value
- guarantees receive buffer will not overflow

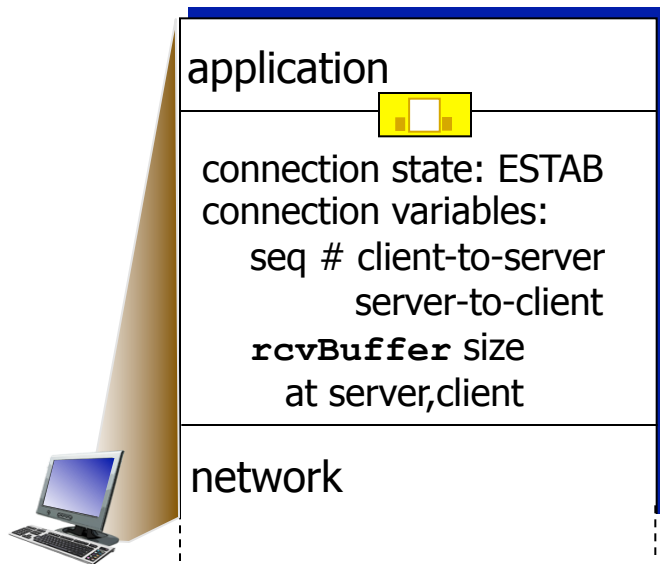


rwnd = Receiver Window

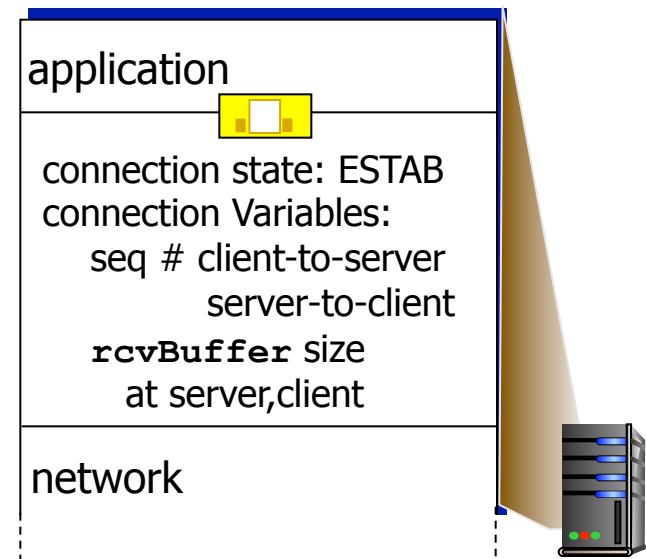
Connection Management

before exchanging data, sender/receiver “handshake”:

- agree to establish connection (each knowing the other willing to establish connection)
- agree on connection parameters



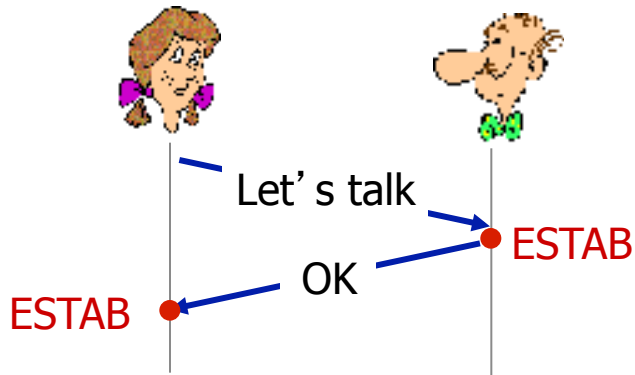
```
Socket clientSocket =  
    newSocket("hostname", "port  
    number");
```



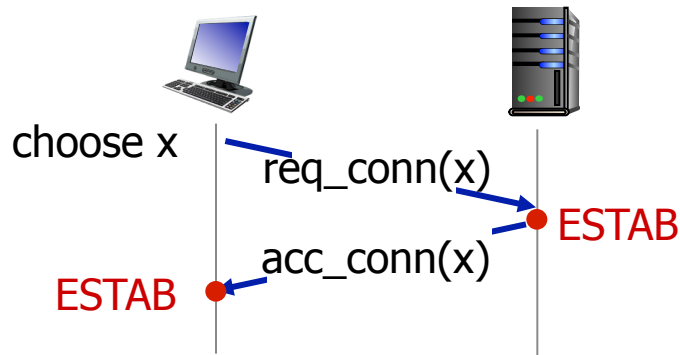
```
Socket connectionSocket =  
    welcomeSocket.accept();
```

Establish a verbal connection

2-way handshake:



- Q: will 2-way handshake always work in network?
- variable delays
 - retransmitted messages (e.g. `req_conn(x)`) due to message loss
 - message reordering
 - can't "see" other side



TCP 3-way handshake

client state

LISTEN

SYNSENT

ESTAB

choose init seq num, x
send TCP SYN msg



SYNbit=1, Seq=x

SYNbit=1, Seq=y
ACKbit=1; ACKnum=x+1

received SYNACK(x)
indicates server is live;
send ACK for SYNACK;
this segment may contain
client-to-server data

ACKbit=1, ACKnum=y+1



choose init seq num, y
send TCP SYNACK
msg, acking SYN

received ACK(y)
indicates client is live

server state

LISTEN

SYN RCVD

ESTAB

TCP: closing a connection

- client, server each close their side of connection
 - ▣ send TCP segment with FIN bit = 1
- respond to received FIN with ACK
 - ▣ on receiving FIN, ACK can be combined with own FIN
- simultaneous FIN exchanges can be handled

TCP: closing a connection

client state

ESTAB

`clientSocket.close()`

FIN_WAIT_1

can no longer
send but can
receive data

FIN_WAIT_2

wait for server
close

TIMED_WAIT

timed wait
for $2 * \text{max}$
segment lifetime

CLOSED



FINbit=1, seq=x

ACKbit=1; ACKnum=x+1

FINbit=1, seq=y

ACKbit=1; ACKnum=y+1

can still
send data

can no longer
send data

server state

ESTAB

CLOSE_WAIT

LAST_ACK

CLOSED

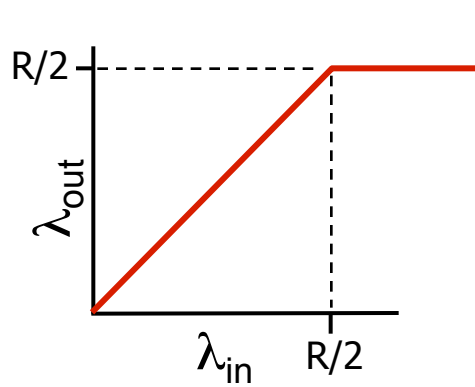
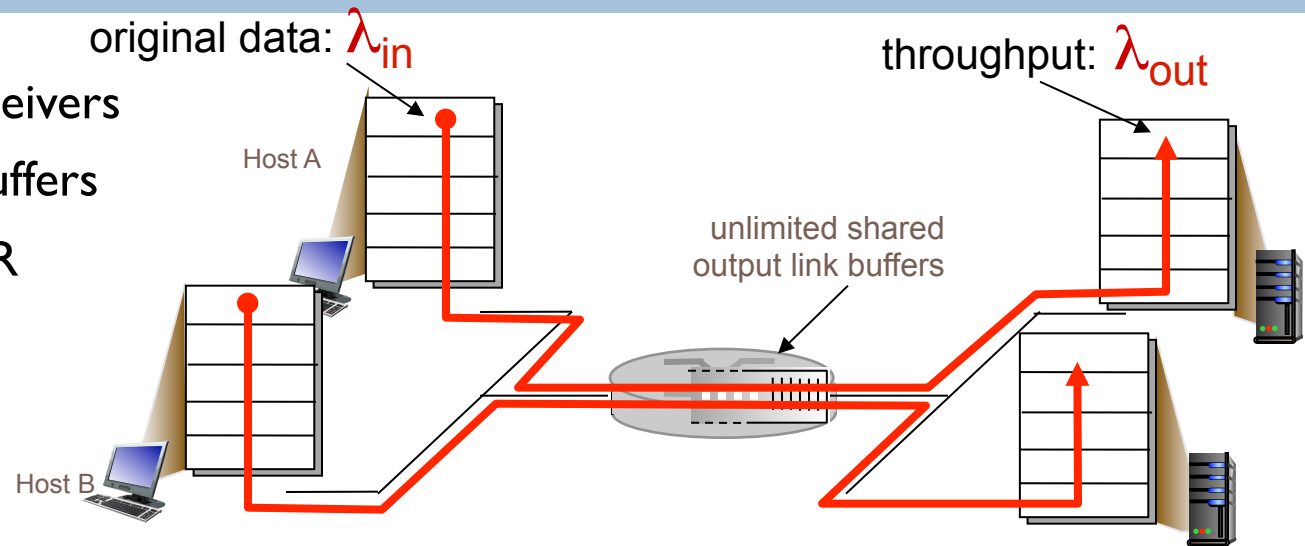
Principles of congestion control

congestion:

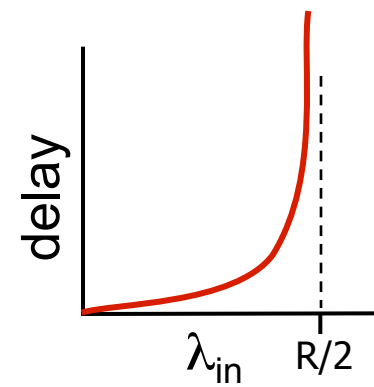
- informally: “too many sources sending too much data too fast for *network* to handle”
- different from flow control!
- manifestations:
 - ▣ lost packets (buffer overflow at routers)
 - ▣ long delays (queueing in router buffers)
- a top-10 problem!

Causes/costs of congestion: Infinite Buffer

- two senders, two receivers
- one router, infinite buffers
- output link capacity: R
- no retransmission



- maximum per-connection throughput: $R/2$



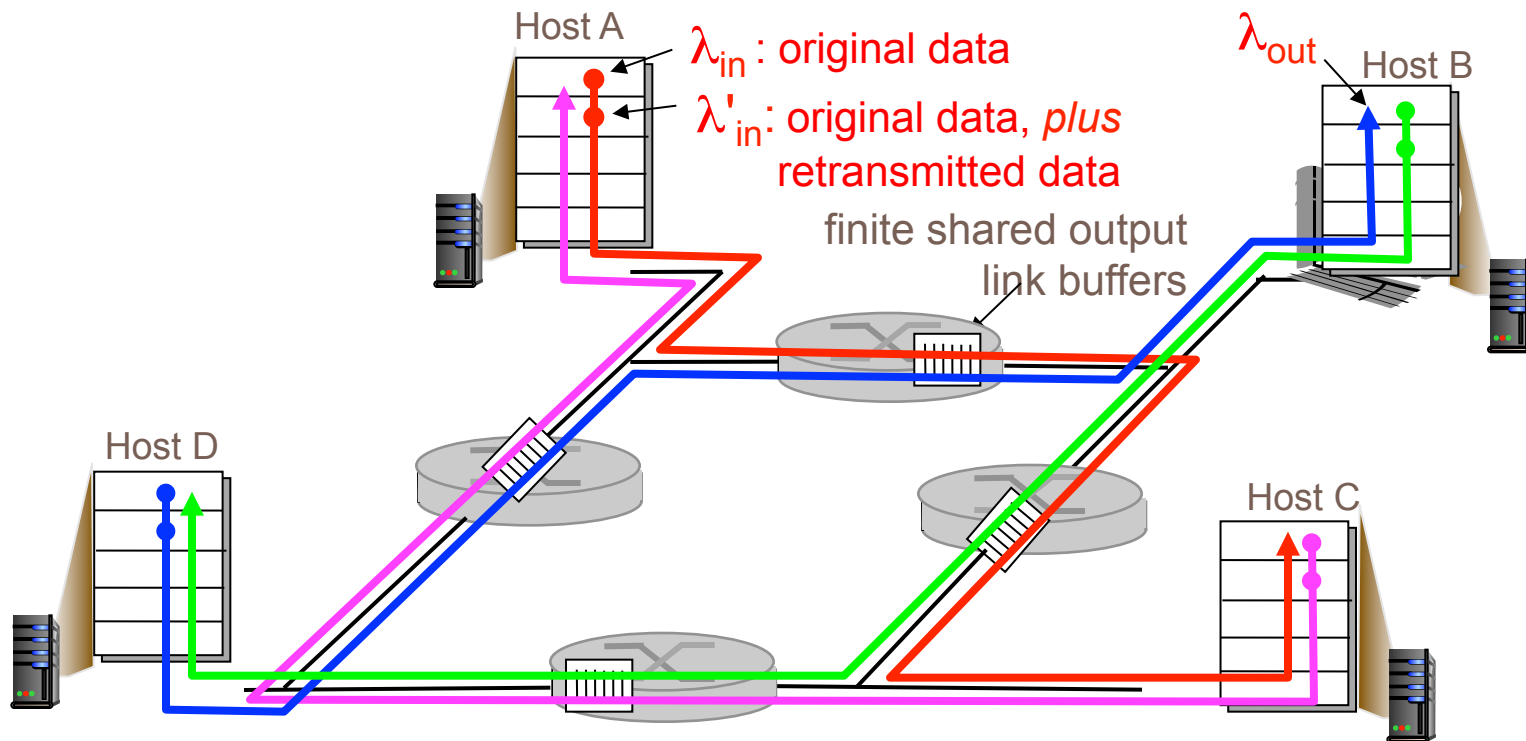
- ❖ large delays as arrival rate, λ_{in} , approaches capacity

Causes/costs of congestion: Finite Buffer

- four senders
- multihop paths
- timeout/retransmit

Q: what happens as λ_{in} and λ'_{in} increase ?

A: as red λ'_{in} increases, all arriving blue pkts at upper queue are dropped, blue throughput $\rightarrow 0$



Approaches towards congestion control

two broad approaches towards congestion control:

end-end congestion control:

- no explicit feedback from network
- congestion inferred from end-system observed loss, delay
- approach taken by TCP

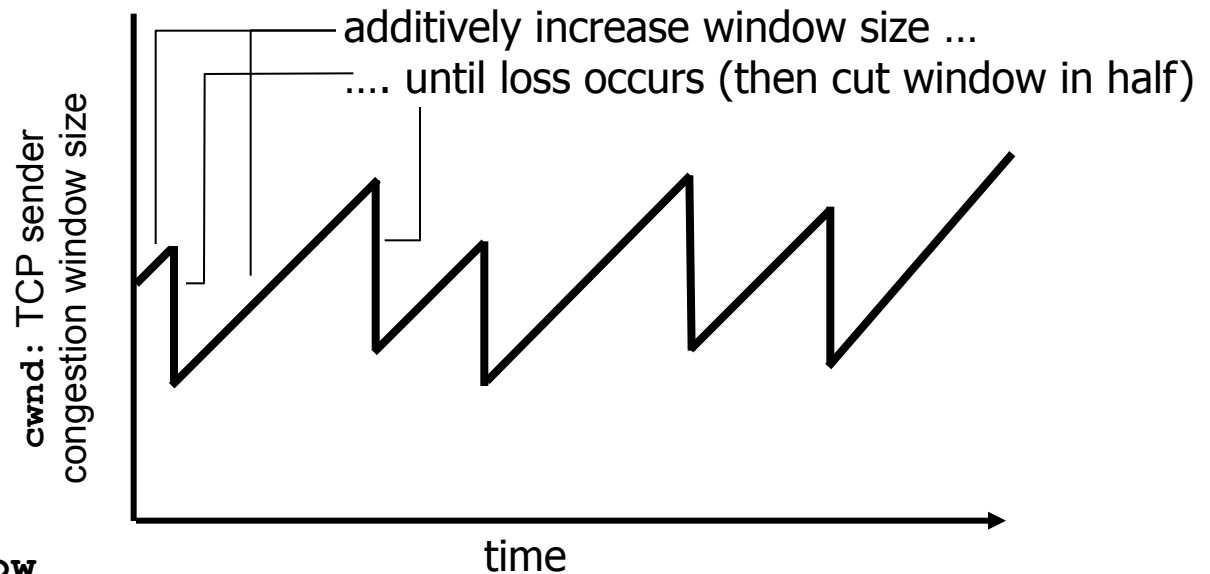
network-assisted congestion control:

- routers provide feedback to end systems
 - ▣ single bit indicating congestion (SNA, DECbit, TCP/IP ECN, ATM)
 - ▣ explicit rate for sender to send at

TCP congestion control: additive increase multiplicative decrease

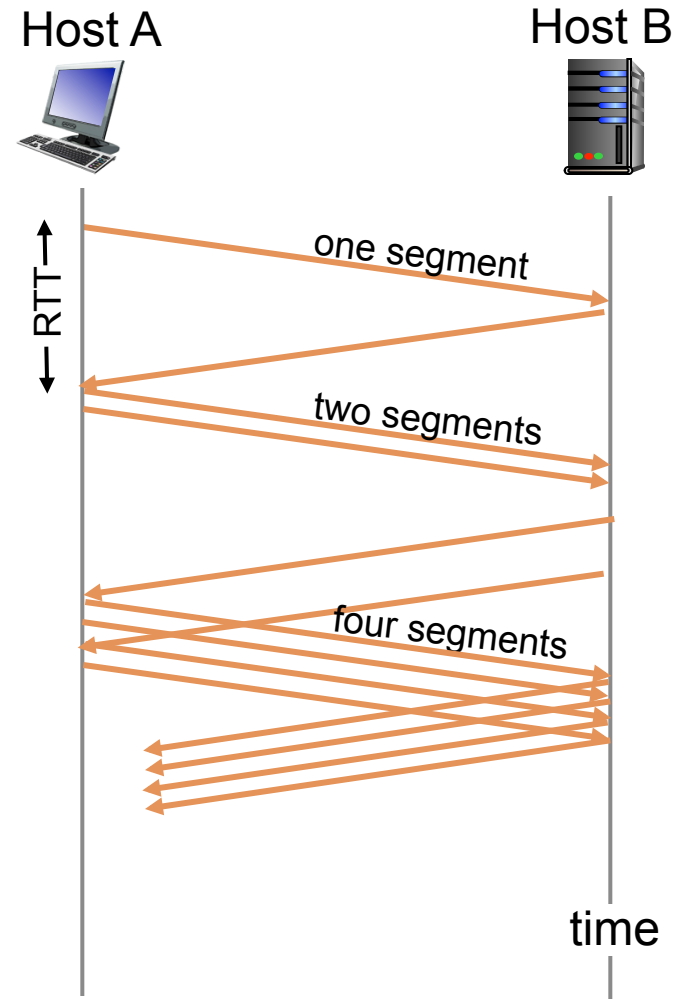
- ❖ *approach*: sender increases transmission rate (window size), probing for usable bandwidth, until loss occurs
 - *additive increase*: increase **cwnd** by 1 MSS every RTT until loss detected
 - *multiplicative decrease*: cut **cwnd** in half after loss

AIMD saw tooth
behavior: probing
for bandwidth

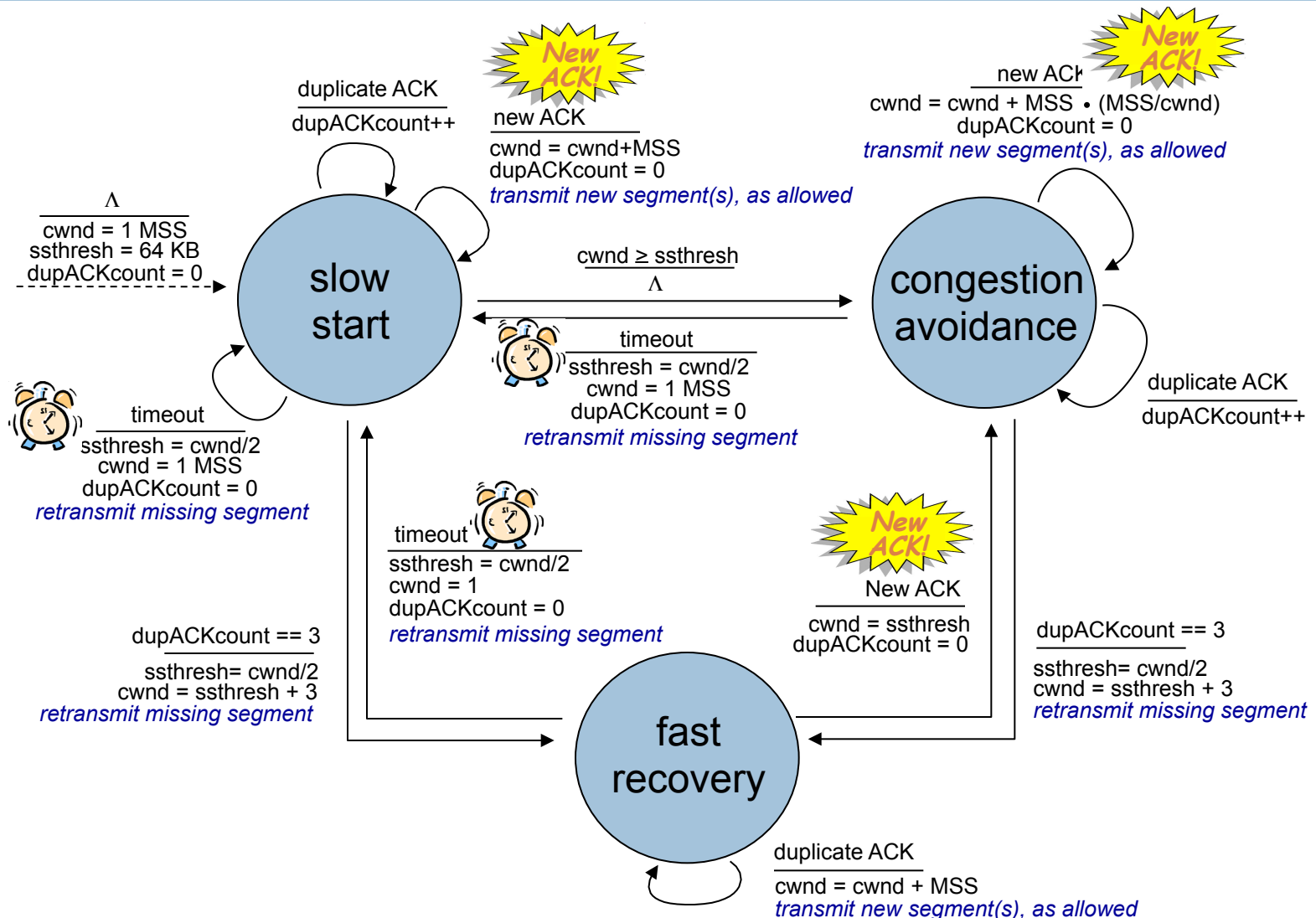


TCP Slow Start

- when connection begins, increase rate exponentially until first loss event:
 - ▣ initially **cwnd** = 1 MSS
 - ▣ double **cwnd** every RTT
 - ▣ done by incrementing **cwnd** for every ACK received
- summary: initial rate is slow but ramps up exponentially fast

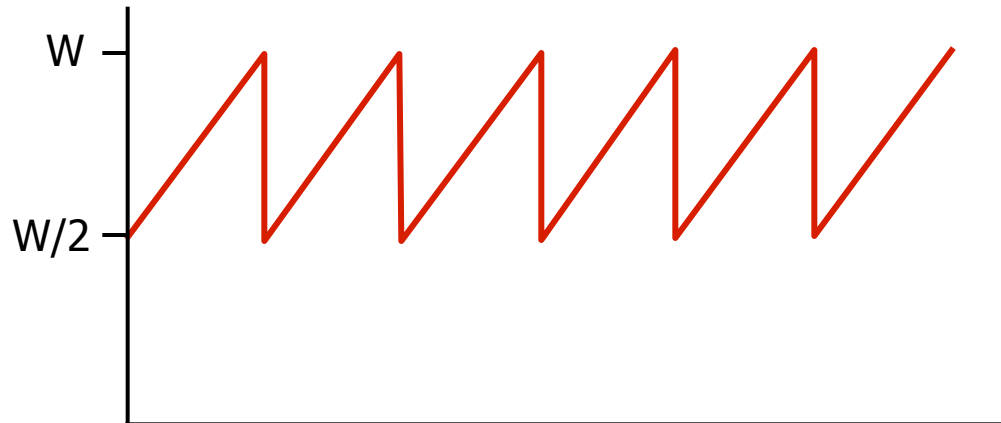


Summary: TCP Congestion Control



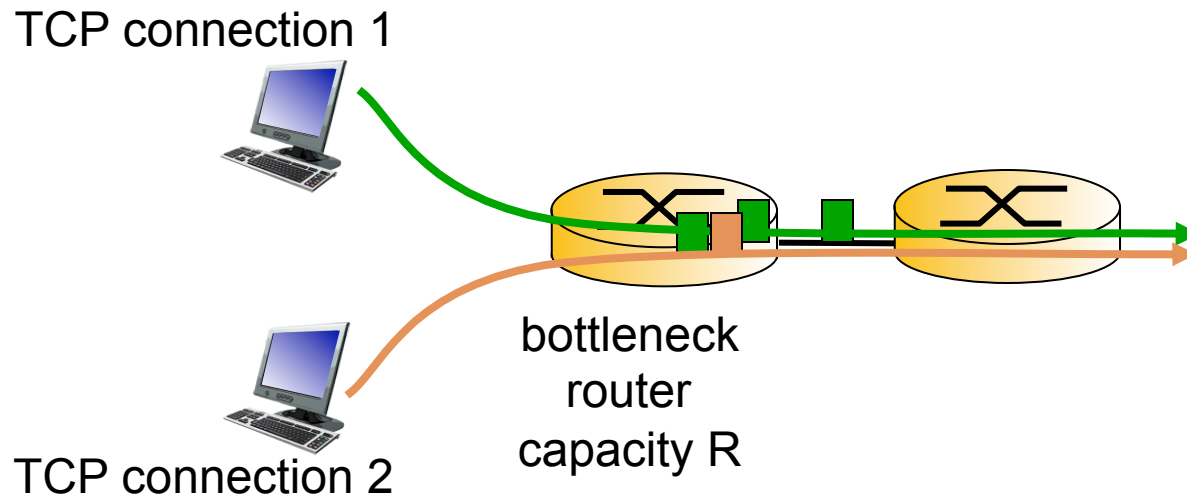
TCP throughput

- avg. TCP thruput as function of window size, RTT?
 - ignore slow start, assume always data to send
- **W: window size** (measured in bytes) **where loss occurs**
 - avg. window size (# in-flight bytes) is $\frac{3}{4} W$
 - avg. thruput is $\frac{3}{4} W$ per RTT
$$\text{avg TCP thruput} = \frac{3}{4} \frac{W}{\text{RTT}} \text{ bytes/sec}$$



TCP Fairness

fairness goal: if K TCP sessions share same bottleneck link of bandwidth R , each should have average rate of R/K



Fairness @ Transport Layer

Fairness and UDP

- multimedia apps often do not use TCP
 - ▣ do not want rate throttled by congestion control
- instead use UDP:
 - ▣ send audio/video at constant rate, tolerate packet loss

Fairness, parallel TCP connections

- application can open multiple parallel connections between two hosts
- web browsers do this
- e.g., link of rate R with 9 existing connections:
 - ▣ new app asks for 1 TCP, gets rate $R/10$
 - ▣ new app asks for 11 TCPs, gets $R/2$