

Transport Control Protocol (TCP)

Richard T. B. Ma

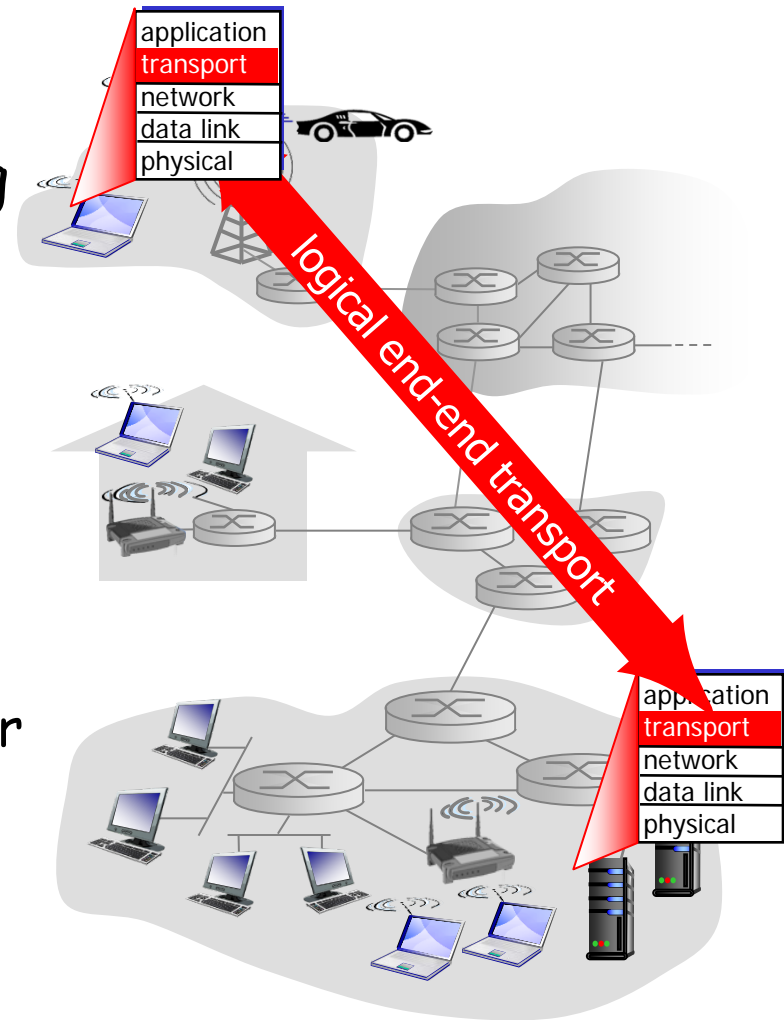
School of Computing

National University of Singapore

CS 3103: Compute Networks and Protocols

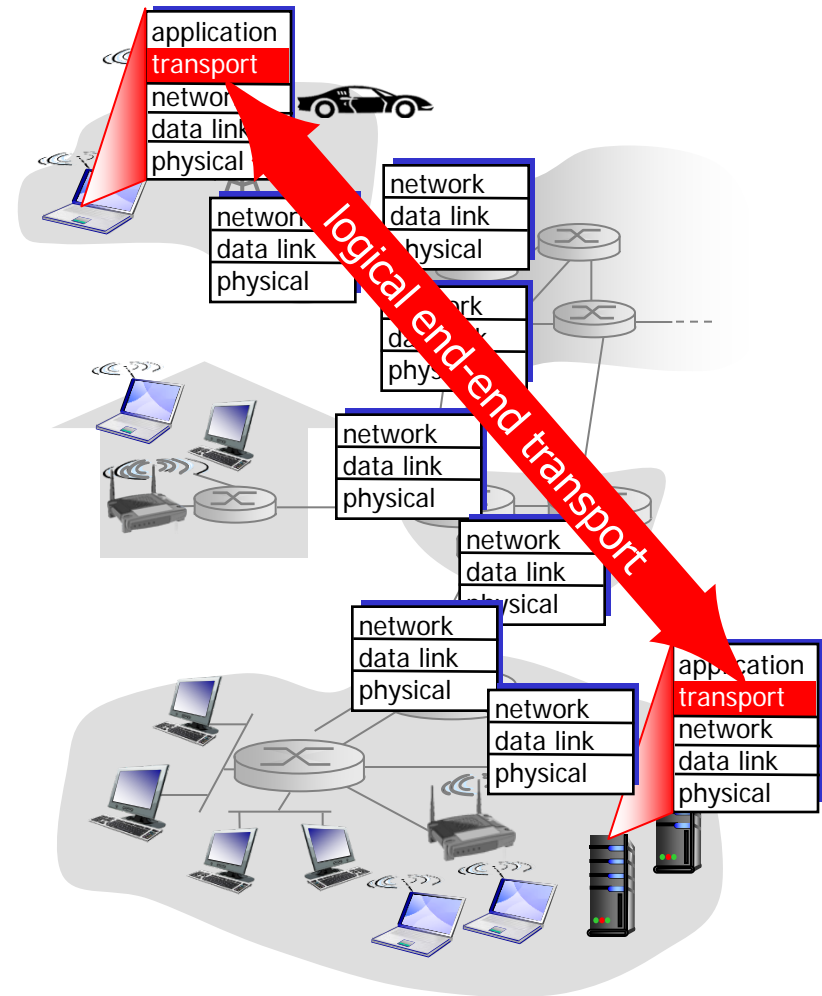
Transport services and protocols

- ❖ provide *logical communication* between app processes running on different hosts
- ❖ protocols run in end systems
 - send side: breaks app messages into *segments*, passes to network layer
 - rcv side: reassembles segments into messages, passes to app layer
- ❖ more than one transport protocol available to apps
 - Internet: TCP and UDP



Internet transport-layer protocols

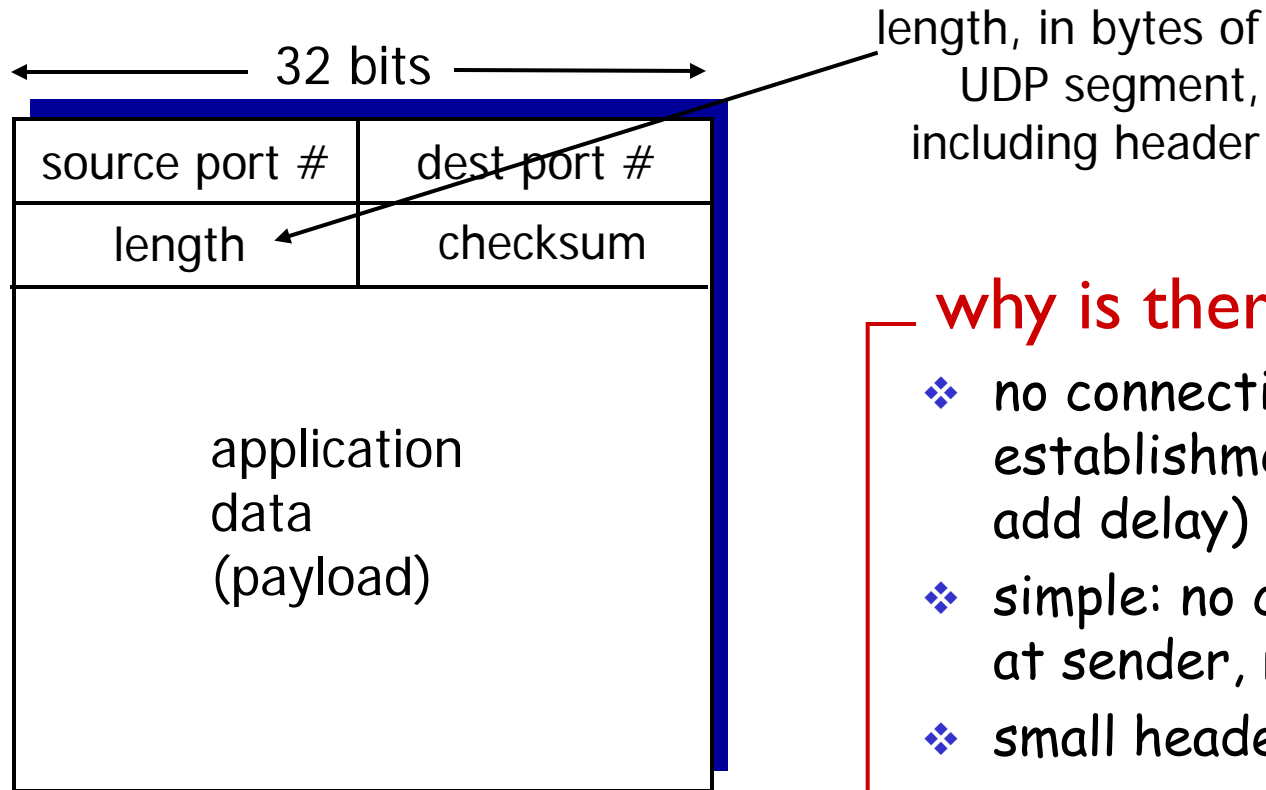
- ❑ reliable, in-order delivery (TCP)
 - ❖ congestion control
 - ❖ flow control
 - ❖ connection setup
- ❑ unreliable, unordered delivery: UDP
 - ❖ no-frills extension of "best-effort" IP
- ❑ services not available:
 - ❖ delay guarantees
 - ❖ bandwidth guarantees



UDP: User Datagram Protocol [RFC 768]

- ❑ “bare bones” Internet transport protocol
- ❑ “best effort” service, UDP segments may be:
 - ❖ lost
 - ❖ delivered out-of-order
- ❑ *connectionless*:
 - ❖ no handshaking between UDP sender, receiver
 - ❖ each UDP segment handled independently of others
- ❖ UDP use:
 - streaming multimedia apps (loss tolerant, rate sensitive)
 - DNS
- ❖ reliable transfer over UDP:
 - add reliability at application layer
 - application-specific error recovery!

UDP: segment header

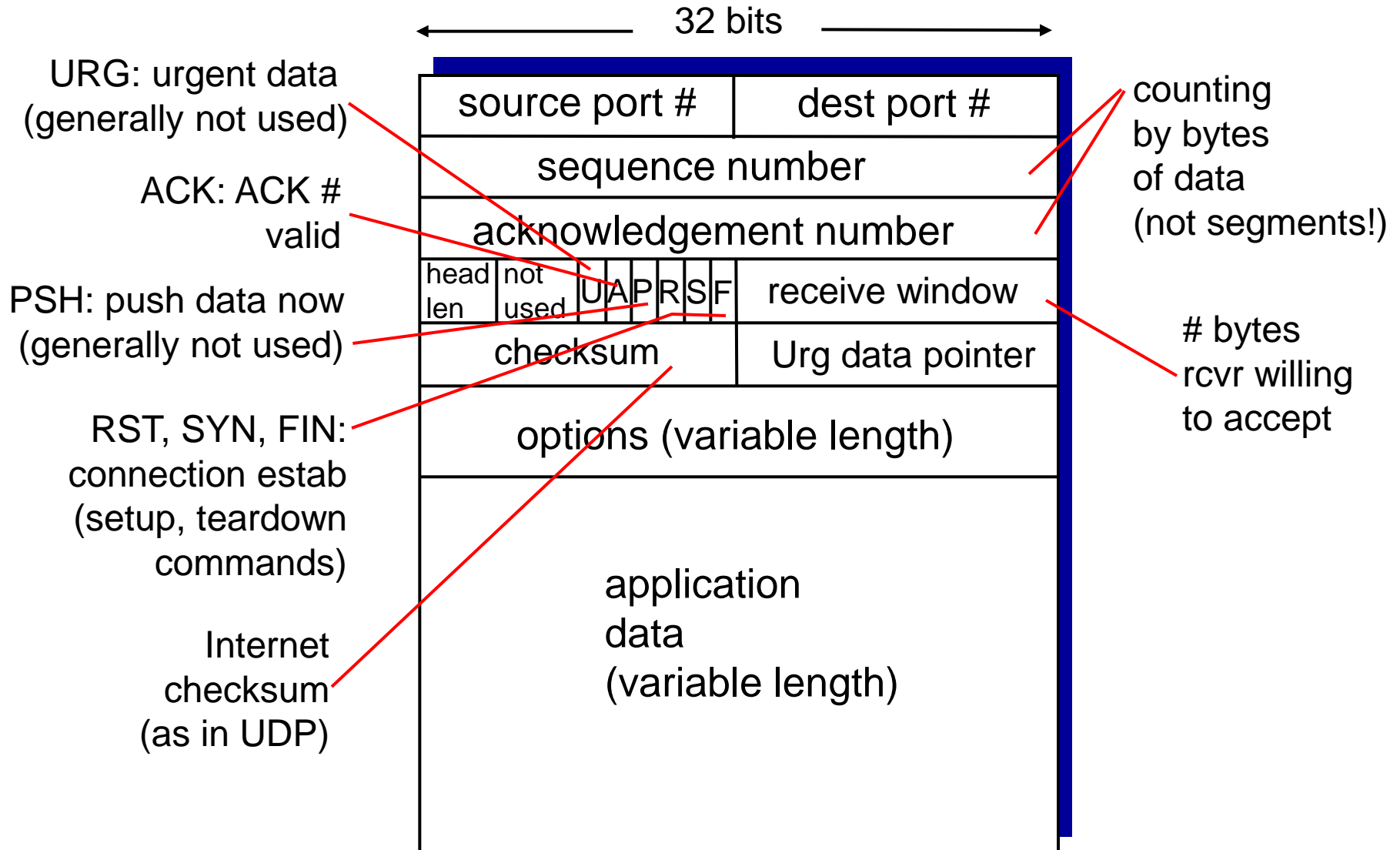


UDP segment format

— why is there a UDP? —

- ❖ no connection establishment (which can add delay)
- ❖ simple: no connection state at sender, receiver
- ❖ small header size
- ❖ no congestion control: UDP can blast away as fast as desired

TCP segment structure



TCP seq. numbers, ACKs

sequence numbers:

- ❖ byte stream “number” of the first byte in segment's data

acknowledgements:

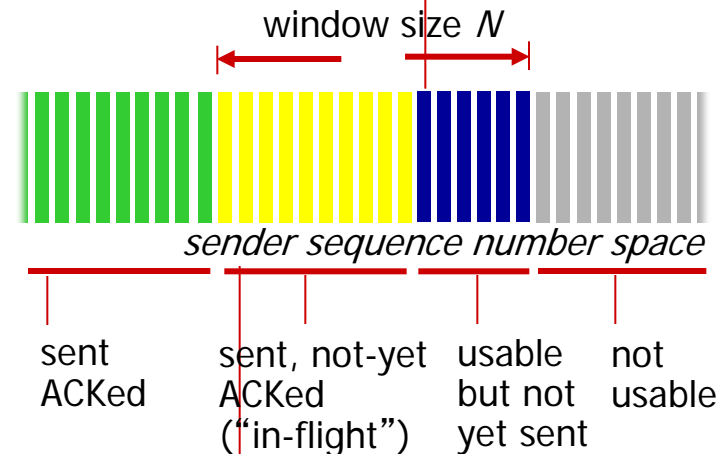
- ❖ seq # of next byte expected from other side
- ❖ cumulative ACK

Q: how receiver handles out-of-order segments

- ❖ A: TCP spec doesn't say, - up to implementor

outgoing segment from sender

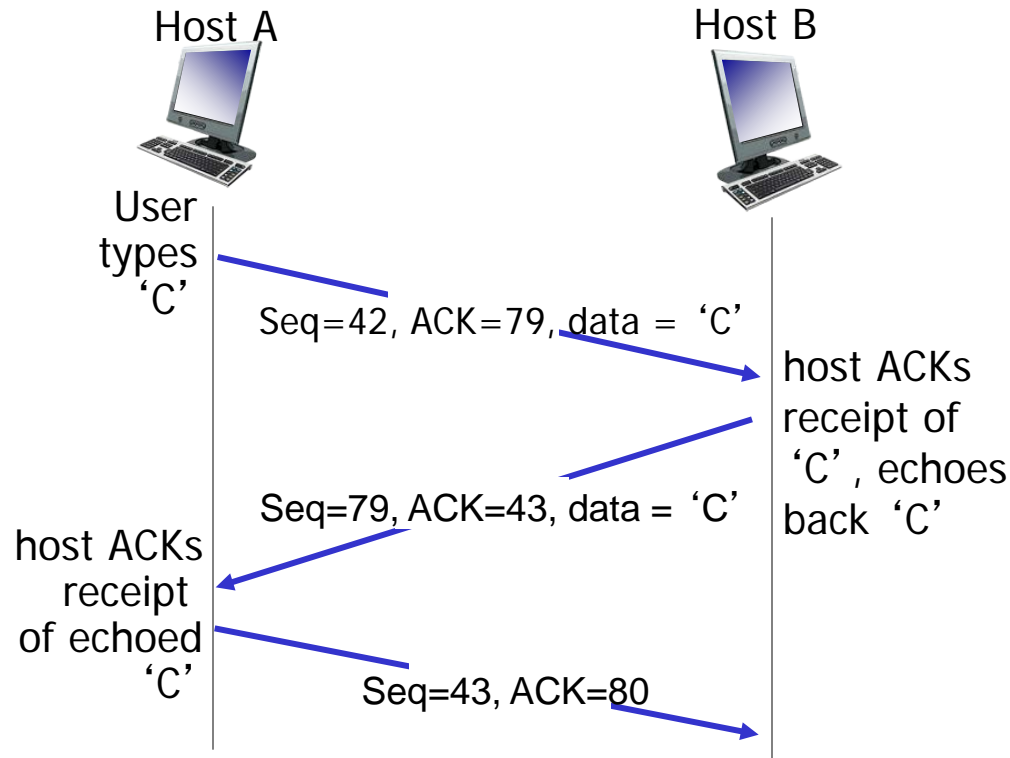
source port #	dest port #
sequence number	
acknowledgement number	
	rwnd
checksum	urg pointer



incoming segment to sender

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

TCP seq. numbers, ACKs



simple telnet scenario

TCP 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 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

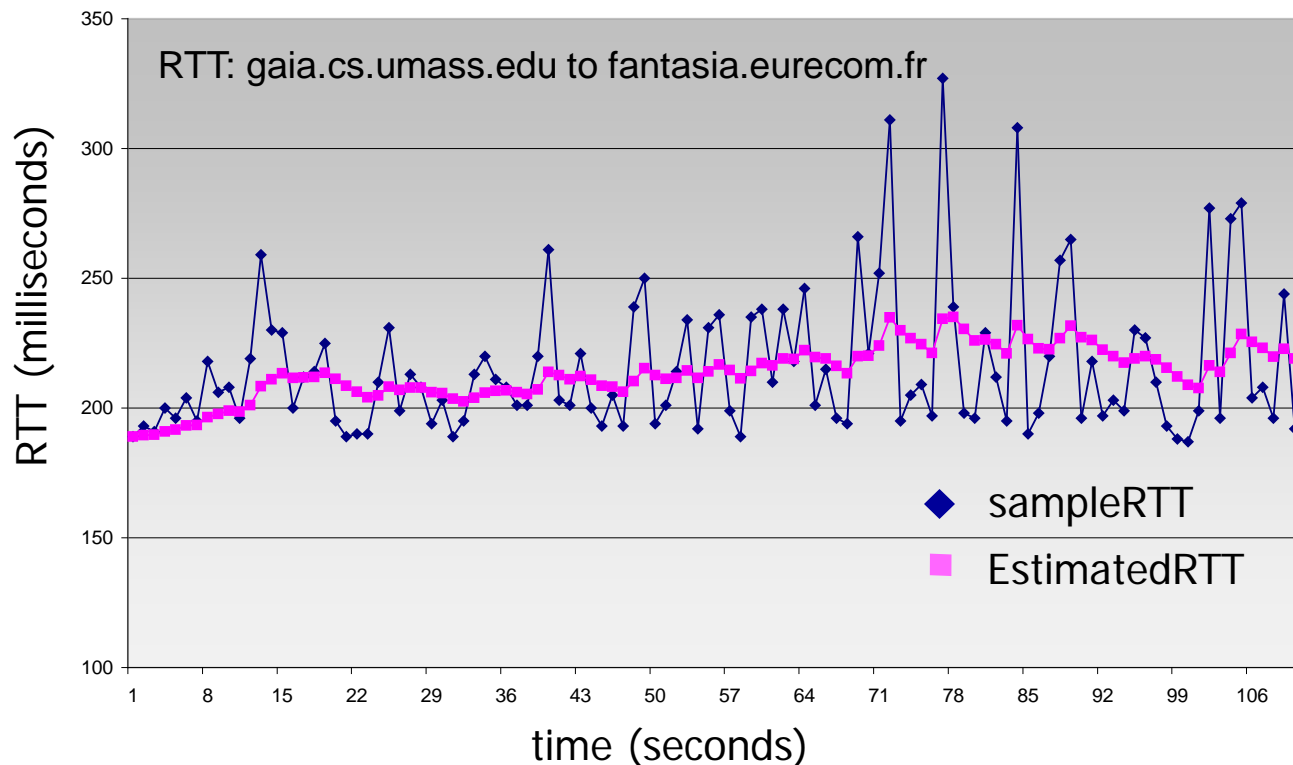
Q: how to estimate RTT?

- **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 round trip time, timeout

- **timeout interval:** EstimatedRTT plus “safety margin”
 - ❖ large variation in EstimatedRTT -> larger safety margin
- **estimate SampleRTT deviation from EstimatedRTT:**

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

(typically, $\beta = 0.25$)

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



↑
estimated RTT

↑
“safety margin”

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`

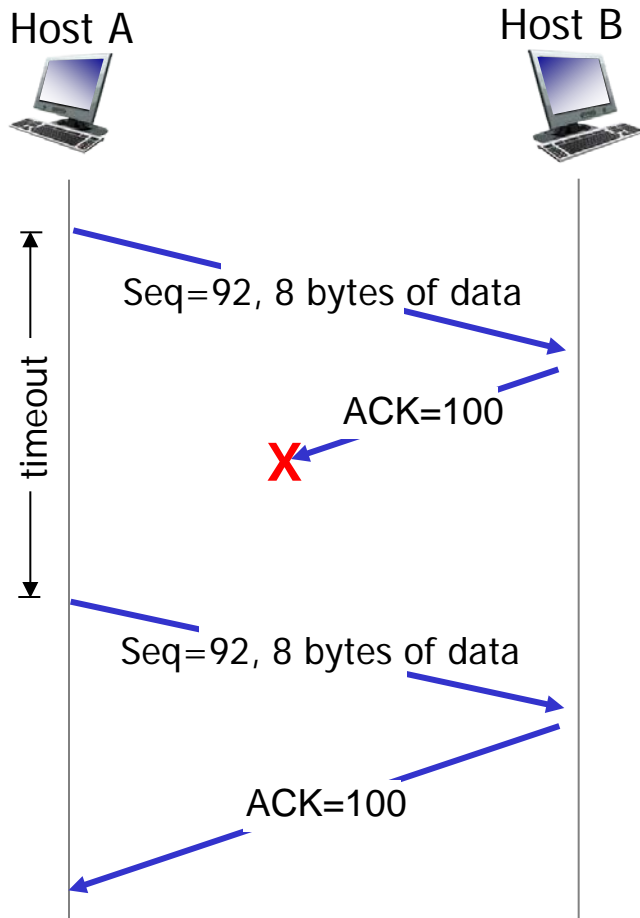
timeout:

- ❑ retransmit segment that caused timeout
- ❑ restart timer

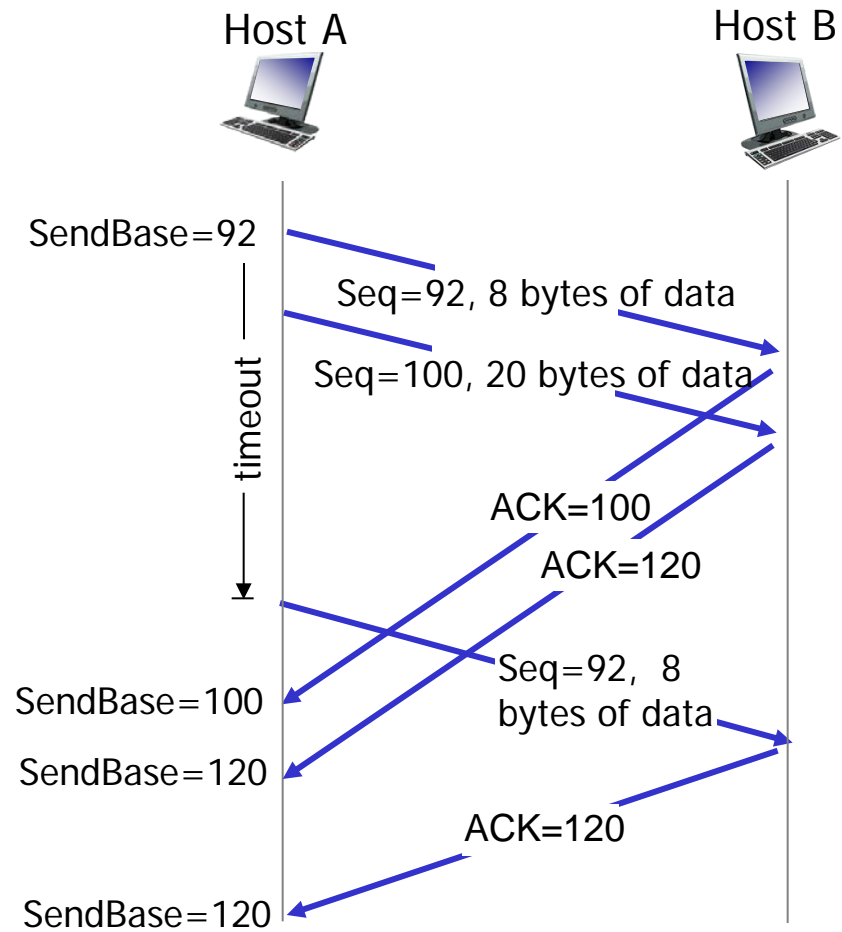
ack rcvd:

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

TCP: retransmission scenarios

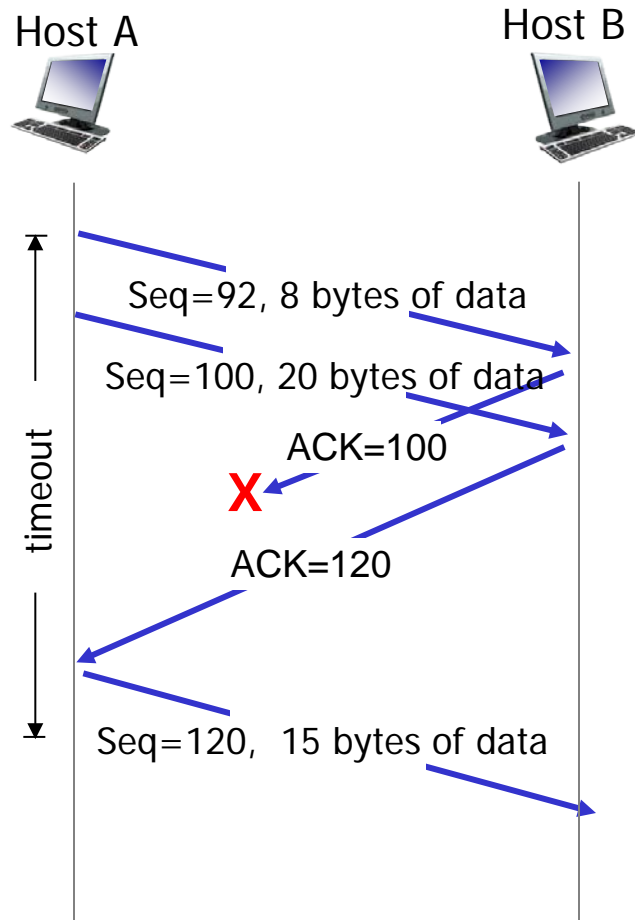


lost ACK scenario



premature timeout

TCP: retransmission scenarios



cumulative ACK

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

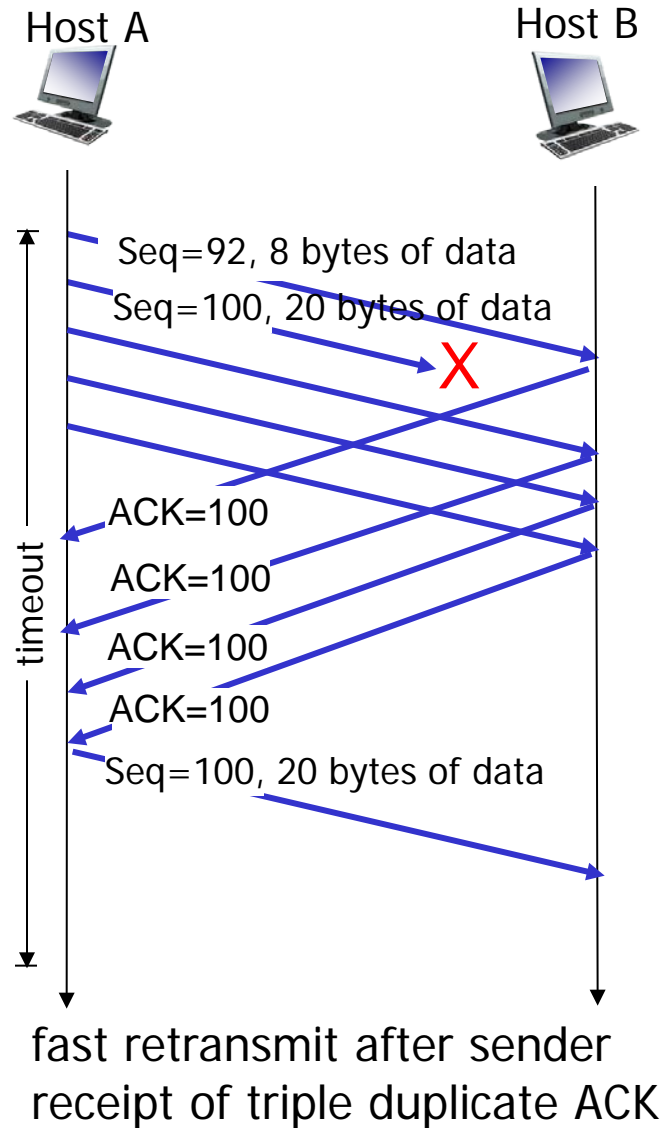
- ❑ time-out period often relatively long:
 - ❖ long delay before resending lost packet
- ❑ detect lost segments via duplicate ACKs.
 - ❖ sender often sends many segments back-to-back
 - ❖ if segment is lost, there will likely be many duplicate ACKs.

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

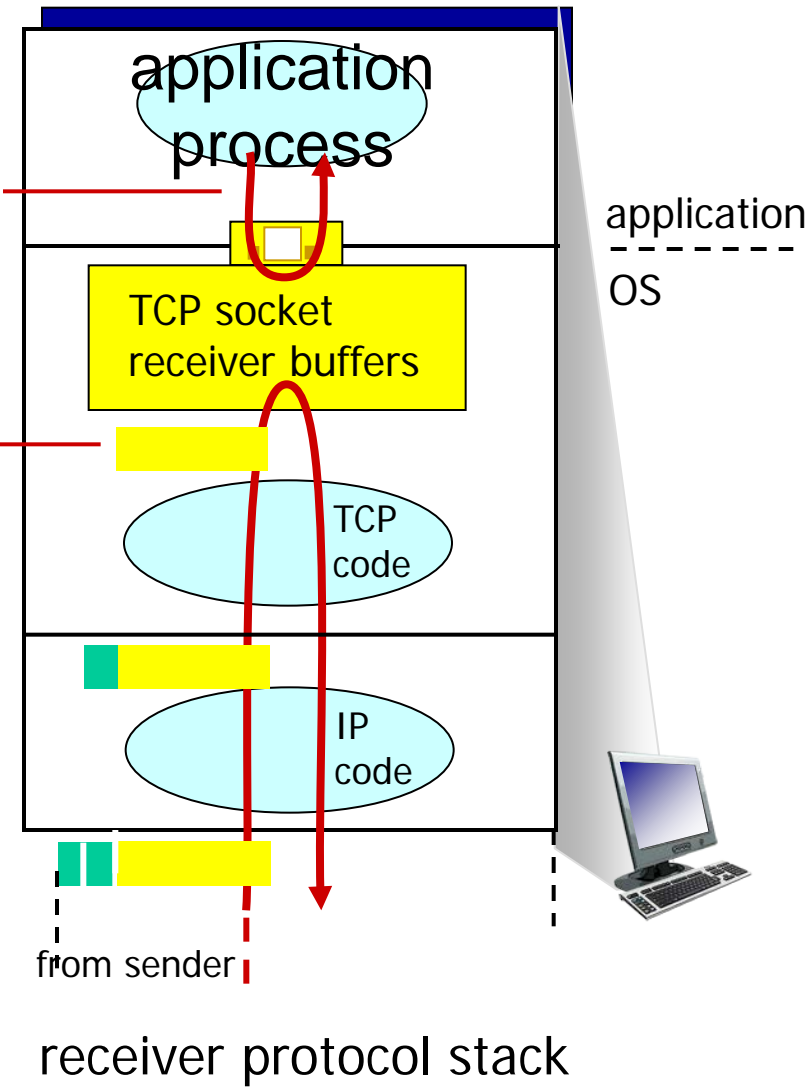


TCP flow control

application may
remove data from
TCP socket buffers

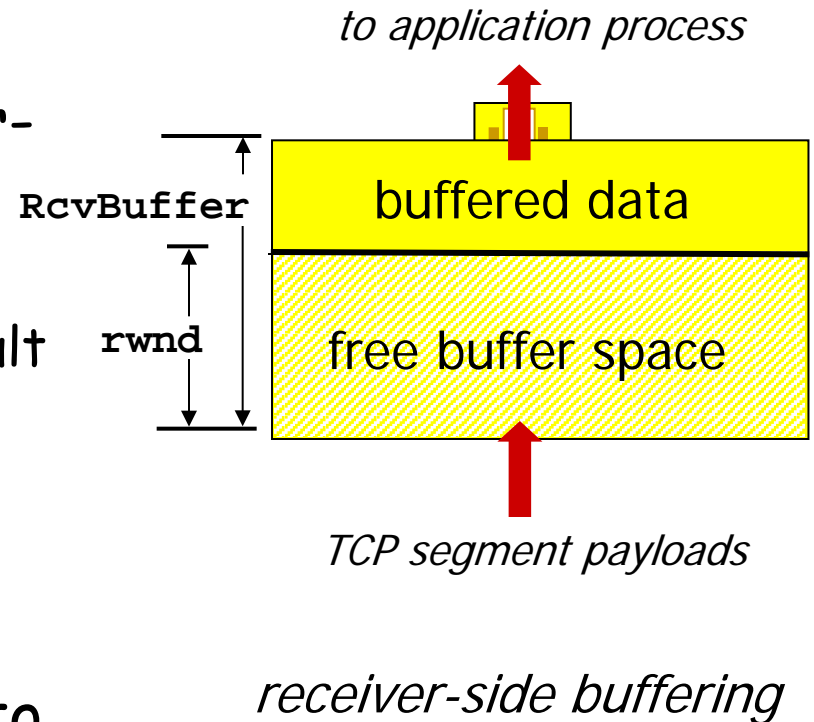
... slower than TCP
receiver is delivering
(sender is sending)

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



TCP flow control

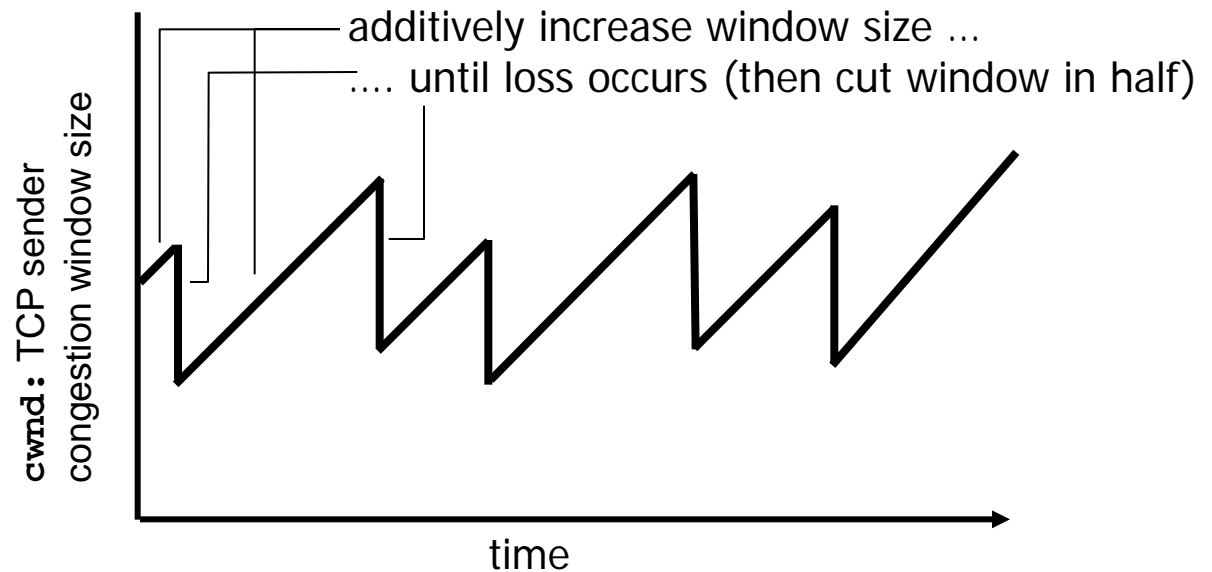
- ❑ receiver “advertises” free buffer space by including `rwnd` in header of receiver-to-sender segments
 - ❖ `RcvBuffer` size set via socket options (typical default is 4096 bytes)
 - ❖ many operating systems autoadjust `RcvBuffer`
- ❑ sender limits amount of unacked (“in-flight”) data to receiver’s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



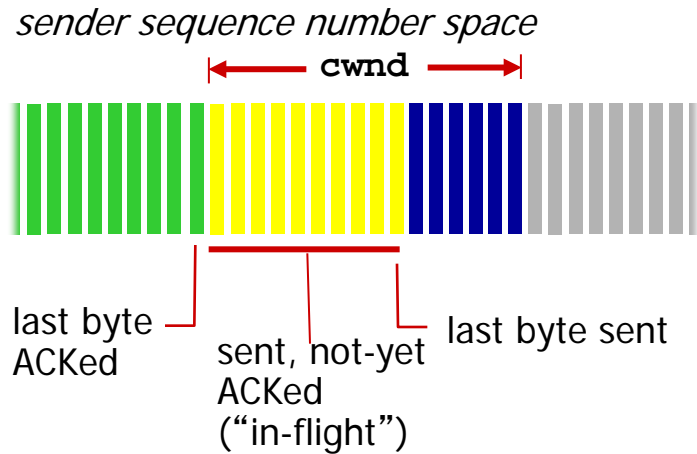
TCP congestion control: additive increase multiplicative decrease

- ❖ *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 Congestion Control: details



TCP sending rate:

❖ *roughly: send cwnd bytes, wait RTT for ACKS, then send more bytes*

□ sender limits transmission:

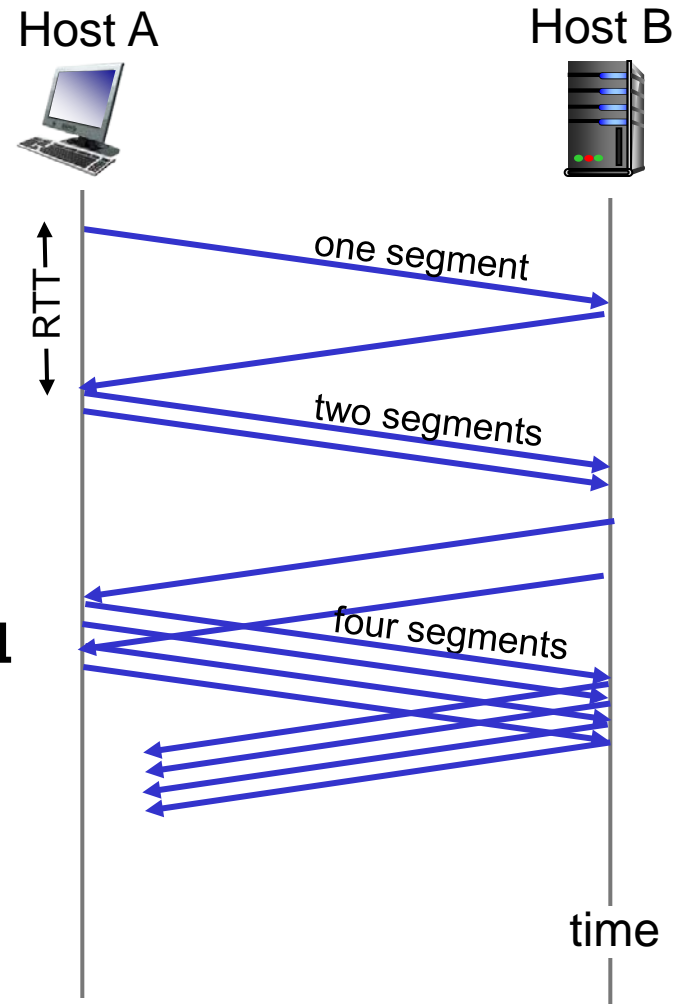
$$\text{LastByteSent} - \text{LastByteAcked} \leq \text{cwnd}$$

□ cwnd is dynamic, a function of perceived network congestion

$$\text{rate} \approx \frac{\text{cwnd}}{\text{RTT}} \text{ bytes/sec}$$

TCP Slow Start

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



TCP: detecting, reacting to loss

- ❑ loss indicated by timeout:
 - ❖ `cwnd` set to 1 MSS;
 - ❖ window then grows exponentially (as in slow start) to threshold, then grows linearly
- ❑ loss indicated by 3 duplicate ACKs: TCP RENO
 - ❖ dup ACKs indicate network capable of delivering some segments
 - ❖ `cwnd` is cut in half window then grows linearly
- ❑ TCP Tahoe always sets `cwnd` to 1 (timeout or 3 duplicate acks)

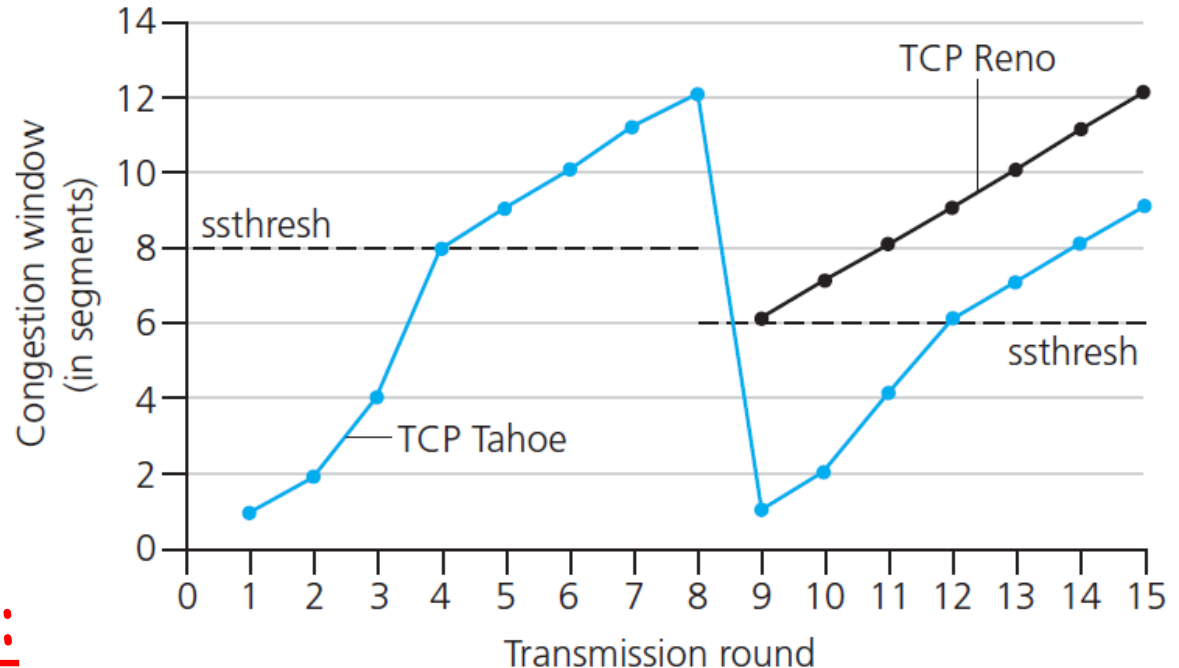
TCP: from slow start to CA

Q: when should the exponential increase switch to linear?

A: when `cwnd` gets to 1/2 of its value before timeout.

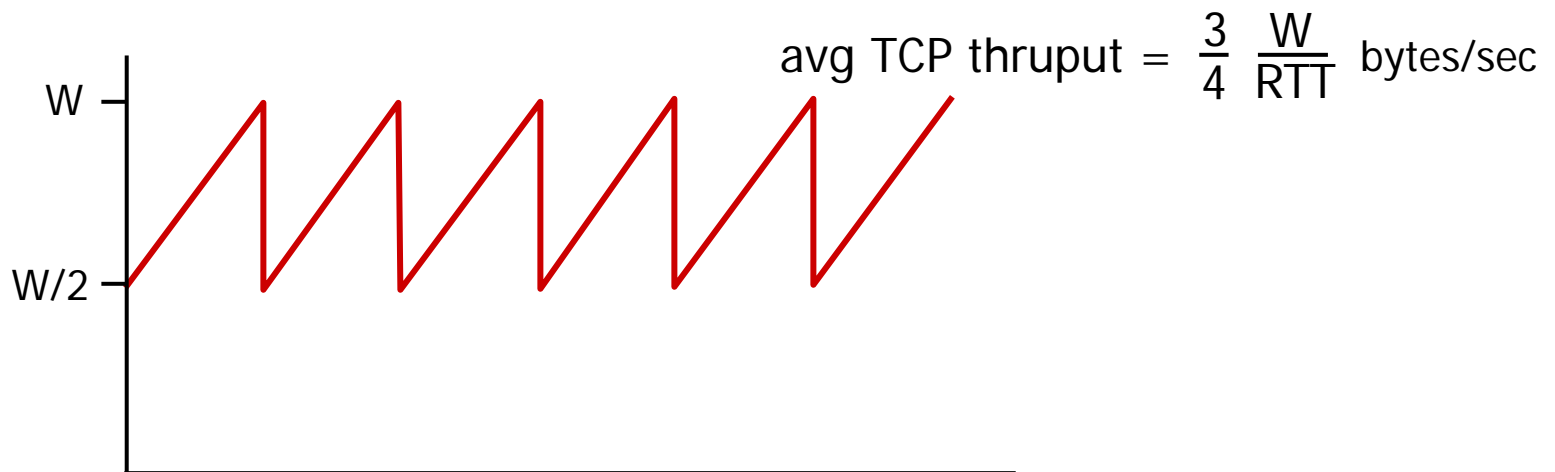
Implementation:

- ❑ variable `ssthresh`
- ❑ on loss event, `ssthresh` is set to 1/2 of `cwnd` just before loss event



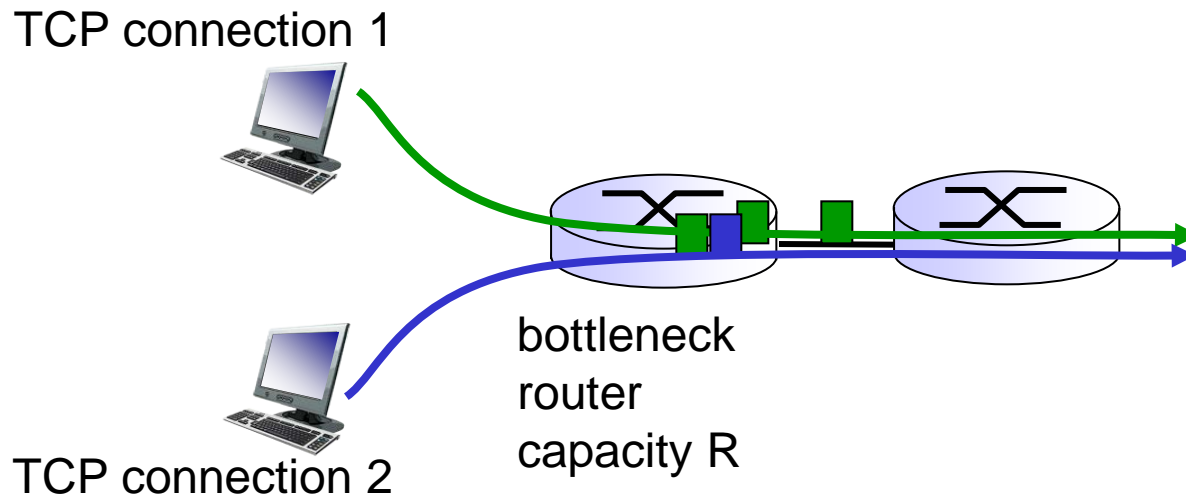
TCP throughput

- avg. TCP throughput 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 $3/4W$ per RTT



TCP Fairness

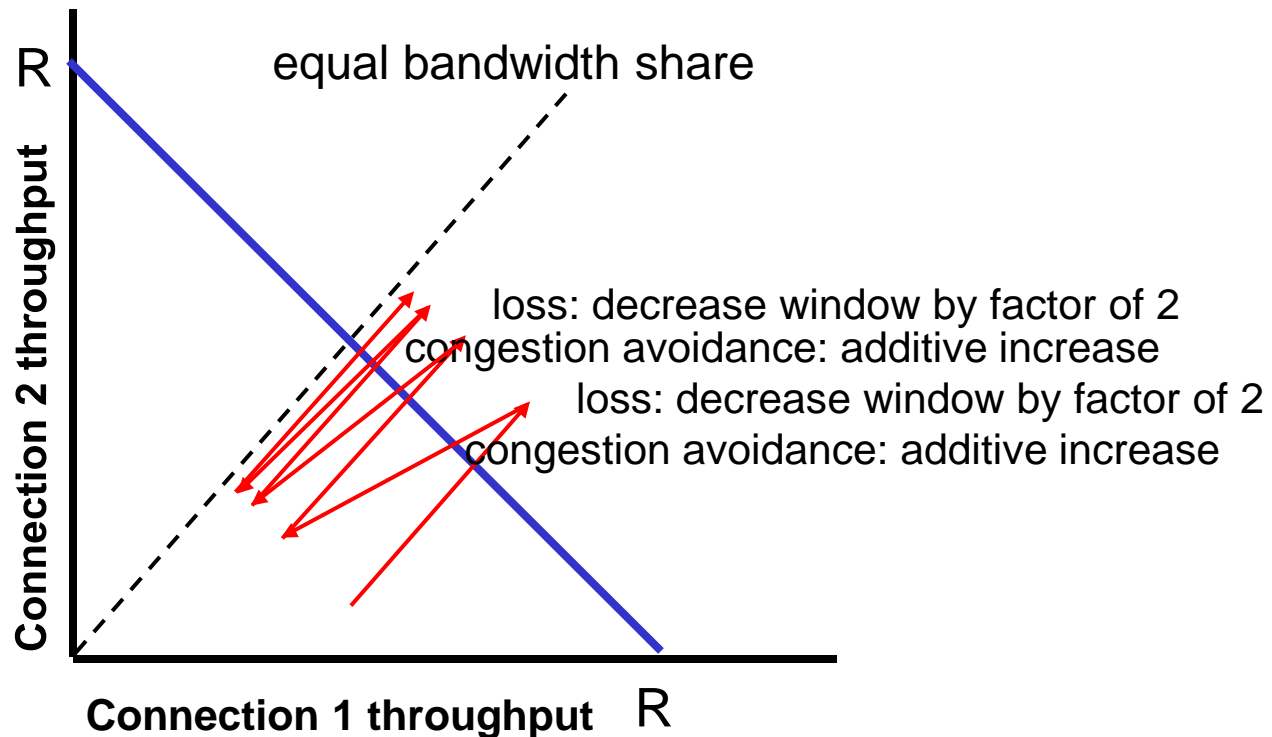
fairness goal: if K TCP sessions share same bottleneck link of bandwidth R , each should have average rate of R/K



Why is TCP fair?

two competing sessions:

- ❖ additive increase gives slope of 1, as throughput increases
- ❖ multiplicative decrease decreases throughput proportionally



Fairness (more)

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

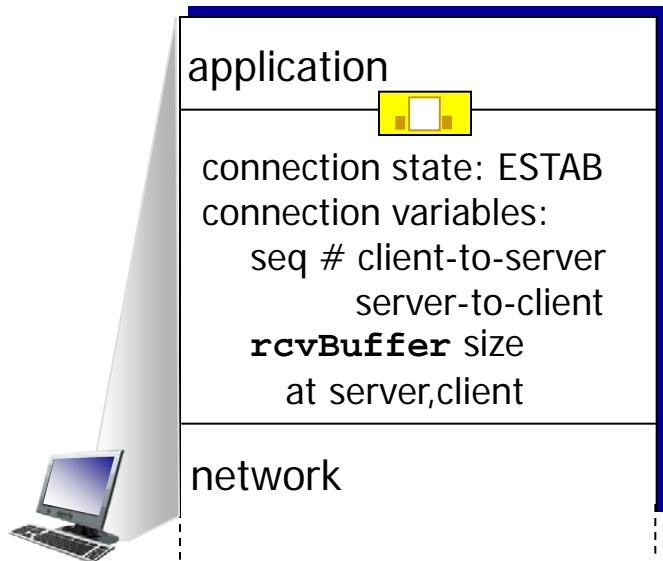
Fairness, parallel TCP connections

- ❑ 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 more than $R/2$

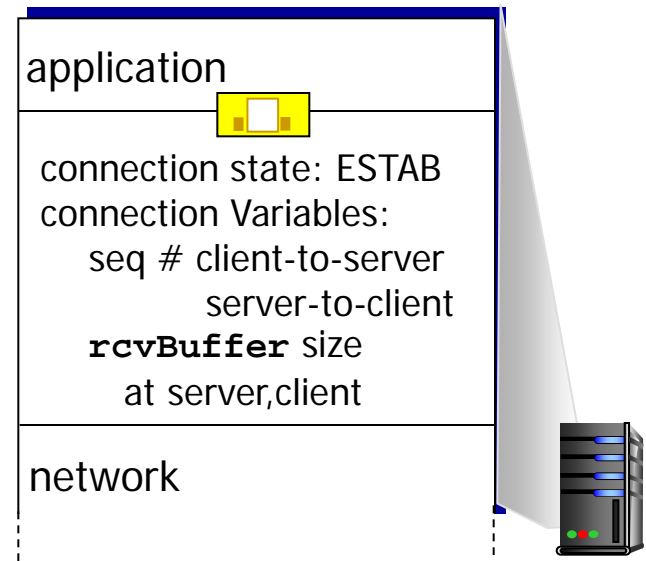
Connection Management

before exchanging data, sender/receiver "handshake":

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



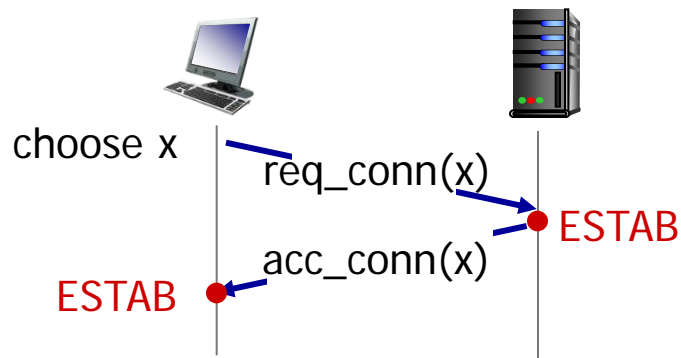
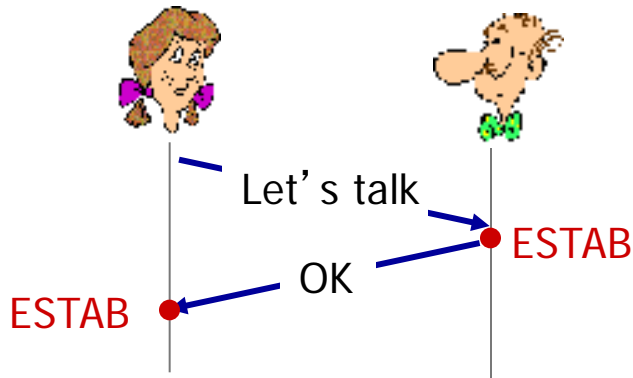
```
Socket clientSocket =  
    newSocket("hostname", "port  
    number");
```



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

Agreeing to establish a connection

2-way handshake:

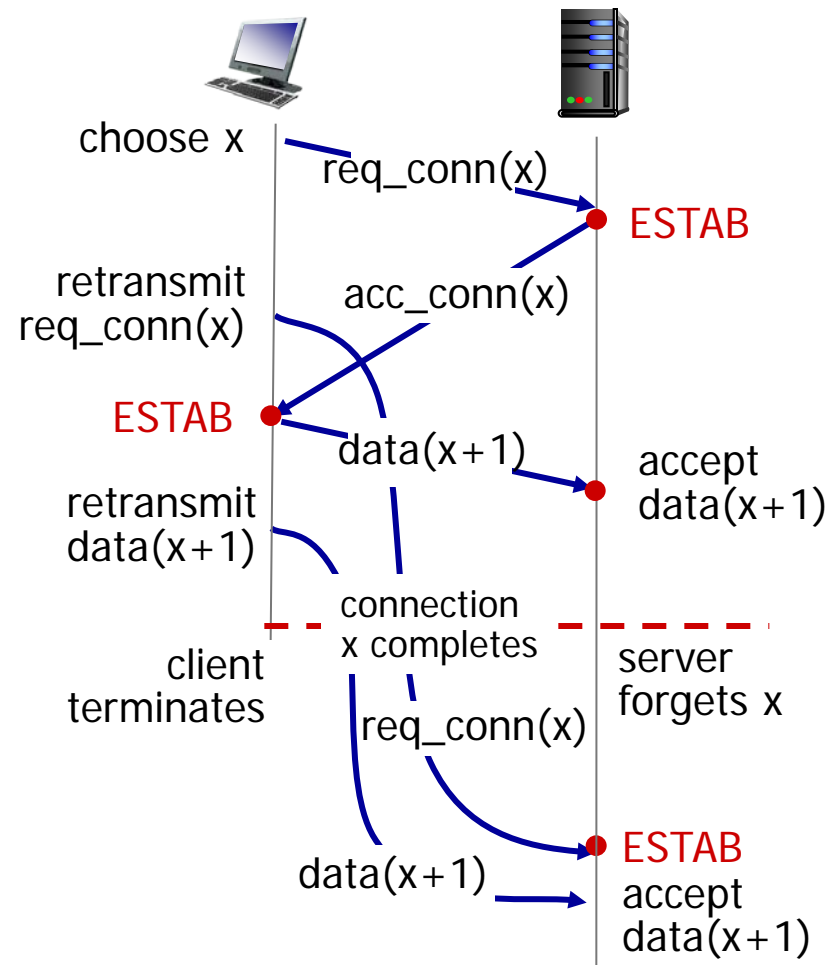
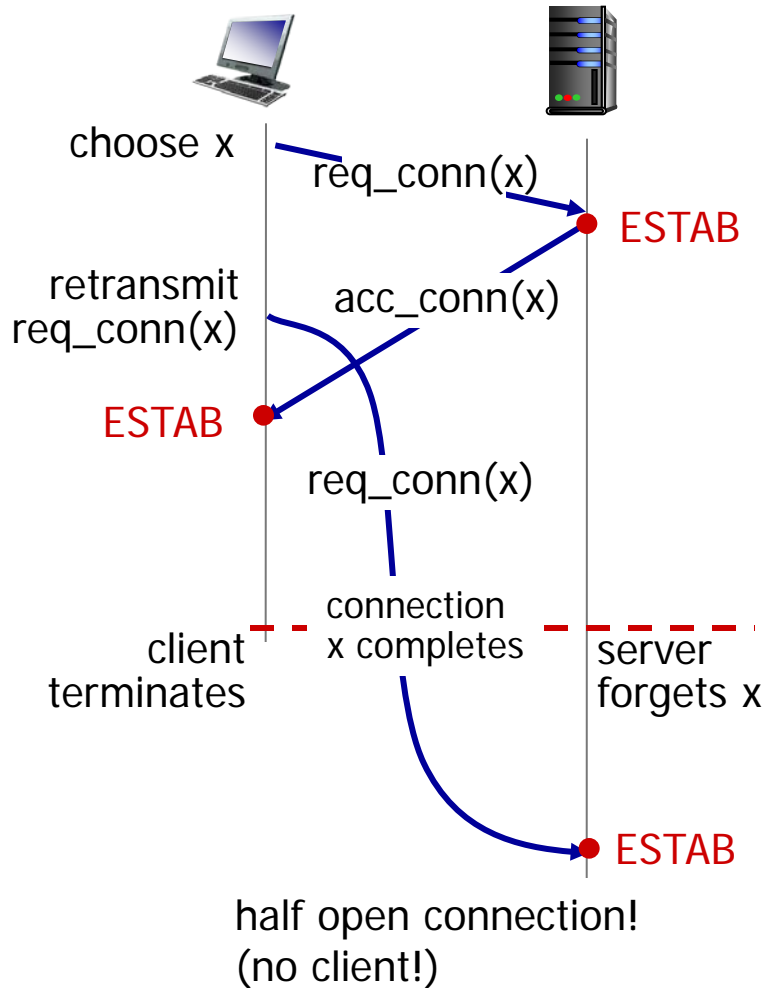


Q: will 2-way handshake always work in network?

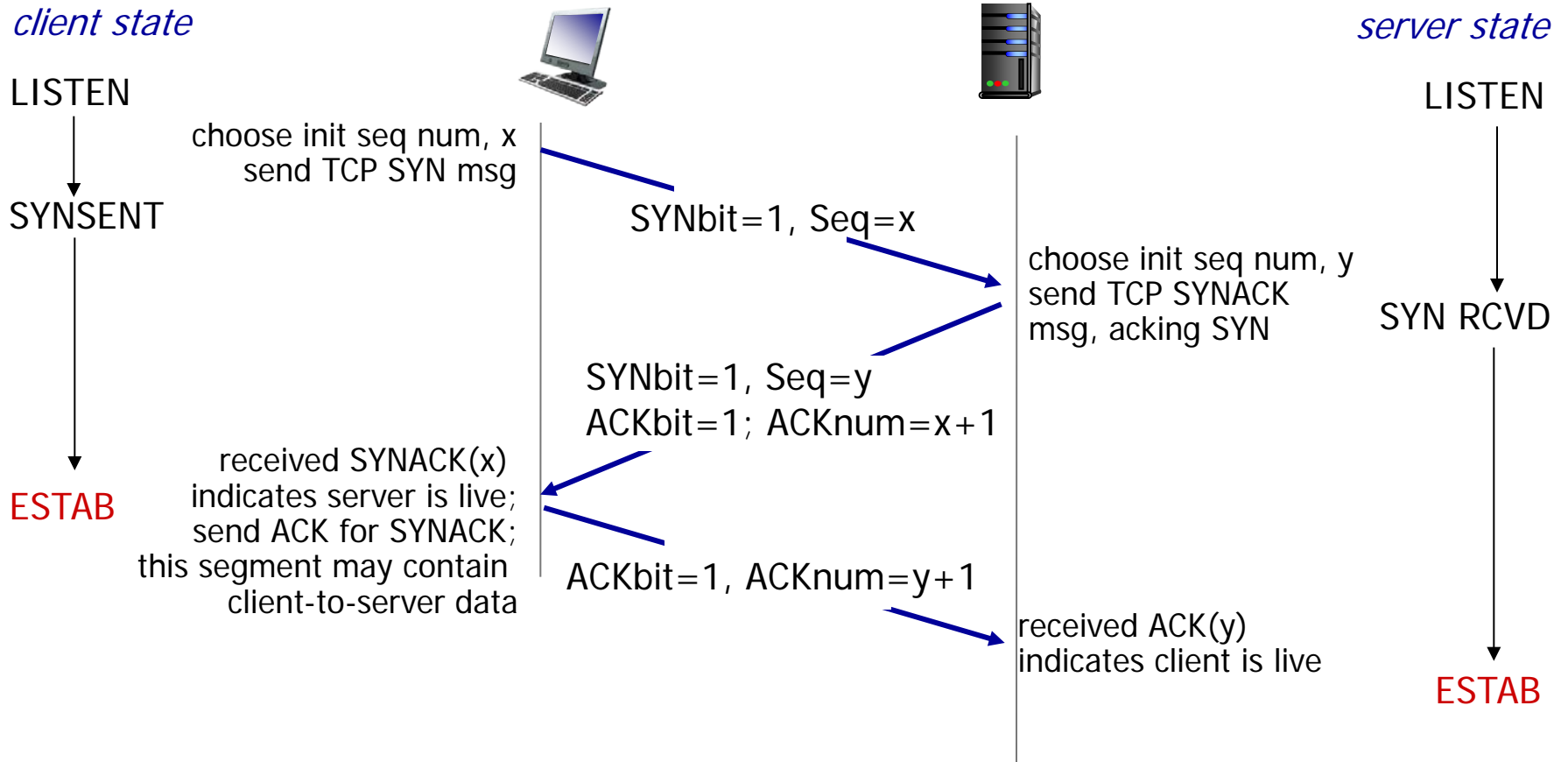
- ☐ variable delays
- ☐ retransmitted messages (e.g. `req_conn(x)`) due to message loss
- ☐ message reordering
- ☐ can't "see" other side

Agreeing to establish a connection

2-way handshake failure scenarios:



TCP 3-way handshake



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