# CS 456/656
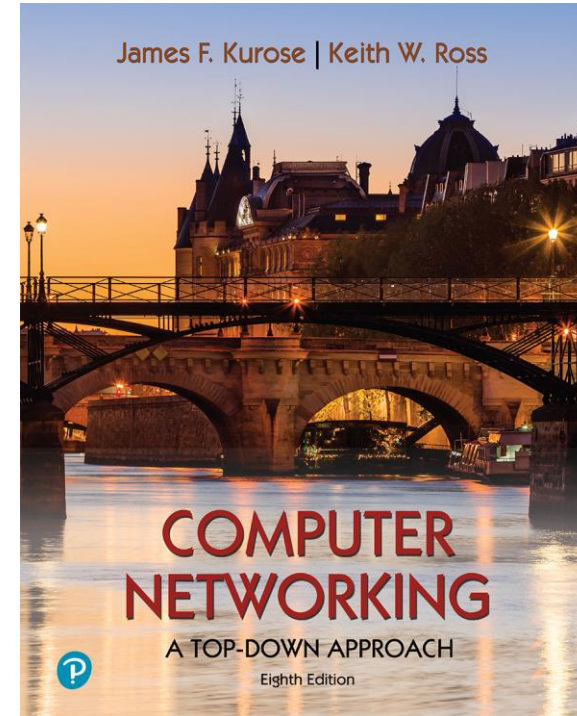# Computer Networks

## Lecture 7: Transport Layer – Part 3

Mina Tahmasbi Arashloo and Bo Sun

Fall 2024

# A note on the slides

Adapted from the slides that accompany this book.

*Computer Networking: A Top-Down Approach*
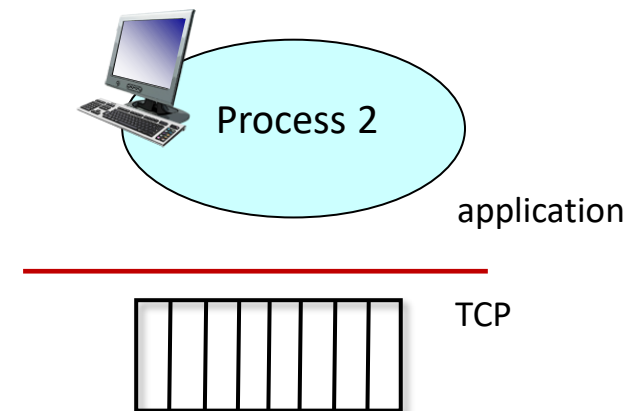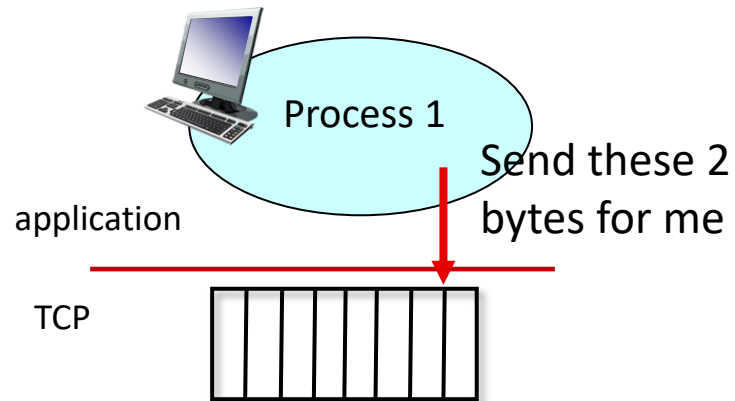8th edition
Jim Kurose, Keith Ross
Pearson, 2020

# Transport layer: roadmap

- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- **Connection-oriented transport: TCP**
    - segment structure
    - connection management
    - reliable data transfer
    - flow control
- Principles of congestion control
- TCP congestion control

# TCP: a widely-used reliable transport protocol

- **Guarantees reliable, in-order byte steam:**
  - no "message boundaries"

# TCP: a widely-used reliable transport protocol

- **Guarantees** reliable, in-order *byte steam:*
  - no "message boundaries"

# TCP: a widely-used reliable transport protocol

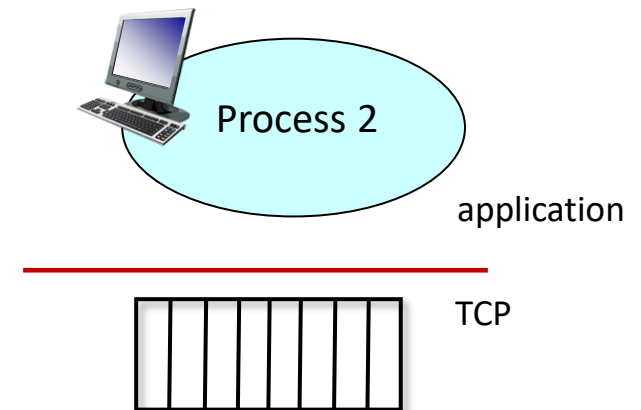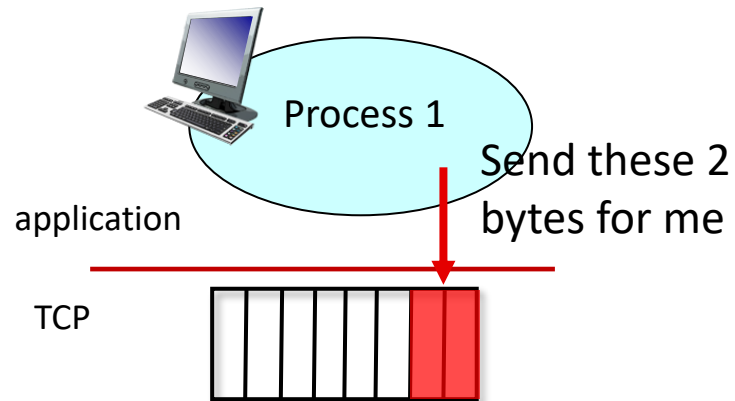- Guarantees reliable, in-order *byte steam:*
  - no "message boundaries"

# TCP: a widely-used reliable transport protocol

- Guarantees reliable, in-order *byte steam:*
  - no "message boundaries"

# TCP: a widely-used reliable transport protocol

- Guarantees reliable, in-order *byte steam:*
  - no "message boundaries"



application

Process 1

TCP

I'll fit them into 2 packets,
each with 3 bytes of the data

application

Process 2

TCP

# TCP: a widely-used reliable transport protocol

- Guarantees reliable, in-order *byte steam:*
  - no "message boundaries"

Process 1

application

TCP

Process 2

application

TCP

acks

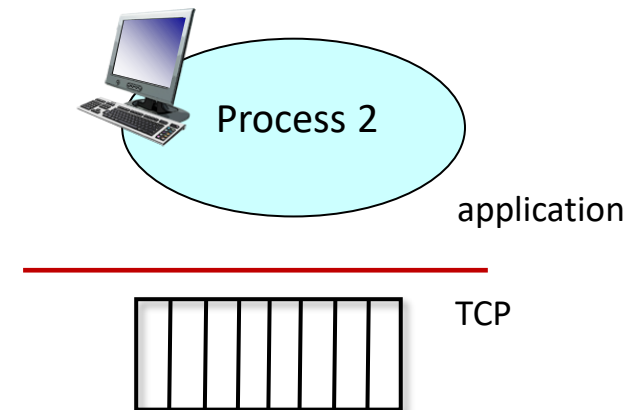I'll put these 6 bytes back together into a stream and send acks back

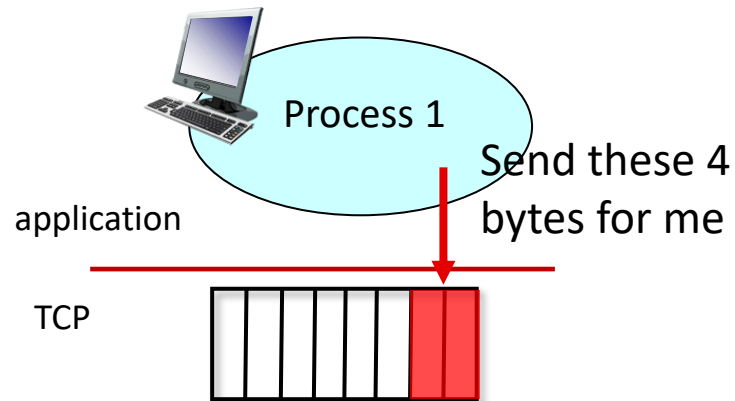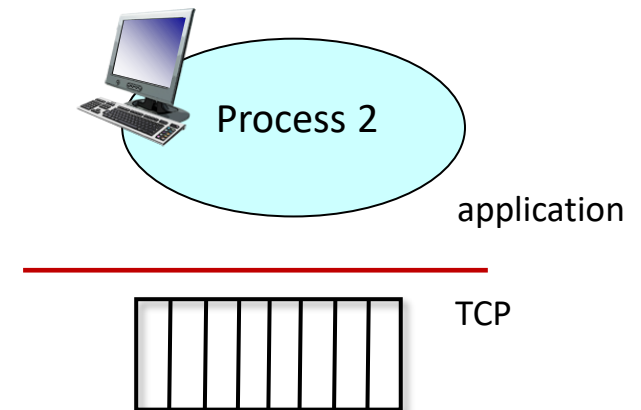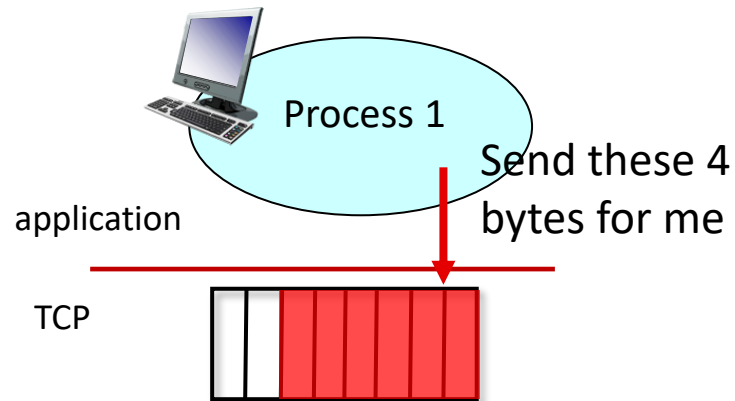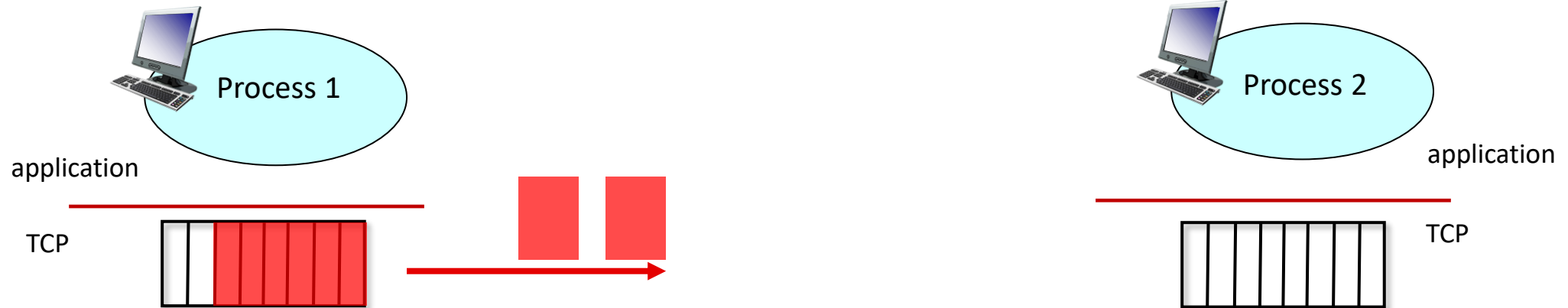# TCP: a widely-used reliable transport protocol

- Guarantees reliable, in-order *byte steam:*
  - no "message boundaries"

# TCP: a widely-used reliable transport protocol

- Guarantees reliable, in-order *byte steam:*
  - no "message boundaries"
- full duplex data: Possible to send data both ways once the two processes establish a connection

# TCP: a widely-used reliable transport protocol

- Guarantees reliable, in-order *byte steam:*
  - no "message boundaries"

- full duplex data:
  - Possible to send data both ways once the two processes establish a connection

- Uses the pipelining approach to reliable data transfer
  - A combination of techniques from Go-Back-N (cumulative acks) and Selective Repeat (only retransmitting presumably lost segment)
  - Performance optimizations like fast retransmit and delayed acks.

# TCP: a widely-used reliable transport protocol

- ## Connection-oriented
  - Connection establishment: Control messages prior to data exchange to initialize the proper state in the communication endpoints
  - Connection tear-down: Control messages after data exchange to end connection

- ## Flow controlled
  - sender will not overwhelm receiver

# Transport layer: roadmap

- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- **Connection-oriented transport: TCP**
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- Principles of congestion control
- TCP congestion control

# TCP segment structure

32 bits

| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |

| head len | not used | C | E | U | A | P | R | S | F | receive window |
|---|---|---|---|---|---|---|---|---|---|---|
| checksum | | | | | | | | | | Urg data pointer |

options (variable length)

application
data
(variable length)

**ACK:** seq # of next expected byte; A bit: this is an ACK

**length** (of TCP header)

Internet **checksum**

**C, E:** congestion notification

TCP **options**

**RST, SYN, FIN:** connection management

**segment seq #:** counting bytes of data into bytestream (not segments!)

**flow control:** # bytes receiver willing to accept

data sent by application into TCP socket

# Transport layer: roadmap

# TCP reliable data transfer

*TCP uses all the reliable data transfer tools we have discussed!*

- Checksum
- Sequence number
- Receiver feedback (ACK)
- Timer
- Sliding window/pipelining

# TCP sequence numbers – one for every byte

- The interface between a sending process and TCP is a byte stream.

- TCP assigns a sequence number to *every byte*
  - As opposed to every segment, as we discussed in the last lecture

- It keeps track of the "status" of every byte
  - Is it sent yet? Is it acknowledged yet?

The Nth byte has sequence number init_seq + N - 1

First byte has sequence number init_seq (initial sequence number)

Next byte has sequence number init_seq + 1

*Sender's view of sequence number space*

(Colors represent segment status -- see next slide)

# TCP sequence numbers

outgoing segment from sender

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

window size
*N*

*Sender's view of sequence number space*

sent &
ACKed

sent, not-yet
ACKed
("in-flight")

usable
but not
yet sent

not
Usable not data associated
with it, application hasn't sent
any data beyond this point)

# TCP ACKs

Sequence number = init_seq

received

Received out of order (optional to track)

Grey ones are not received

*Receiver's view of the sequence number space*

outgoing segment from receiver

- **Cumulative ACK**
  - Has seq number of next expected in-order byte
- **ACK(n) means:**
  - All bytes in [init_seq, n − 1] are received.
  - The receiver is expecting byte n next
- **Note the difference from Go-Back-N ack**

*Q*: What about out-of-order segments?
- *A:* TCP spec doesn't specify, - up to implementor

| source port # | dest port # |
|---|---|
| sequence number ||
| acknowledgement number ||
| | A | rwnd |
| checksum | urg pointer |

# TCP sequence numbers, ACKs

Host A                    Host B

User types 'C'

Seq=42, ACK=79, data = 'C'

host ACKs receipt
of 'C', echoes back 'C'

Seq=79, ACK=43, data = 'C'

host ACKs receipt
of echoed 'C'

Seq=43, ACK=80

simple telnet scenario

# TCP sequence numbers, ACKs

Host A                    Host B

User types 'C'

Seq=42, ACK=79, data = 'C'

host ACKs receipt
of 'C', echoes back 'C'

Seq=79, ACK=43, data = 'C'

host ACKs receipt
of echoed 'C'

Seq=43, ACK=80

simple telnet scenario

- Note the bi-directional communication!
- There are two data streams:
  - one in each direction
  - each with its own sequence number space

# TCP Sender (simplified)

event: data received from application

- create segment with seq #

- seq # is byte-stream offset of first data byte in segment

- start timer if not already running
  - think of timer as for oldest unACKed segment
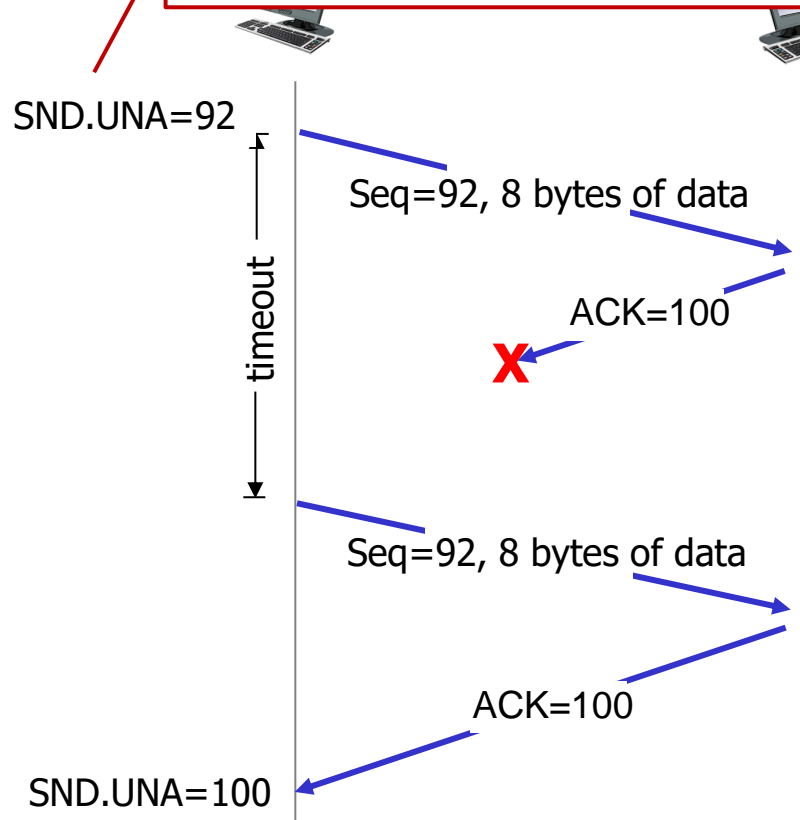  - expiration interval: `TimeOutInterval`

*event: timeout*

- retransmit segment that caused timeout
- restart timer

*event: ACK received*

- if ACK acknowledges previously unACKed segments
  - update what is known to be ACKed
  - restart timer if there are still unACKed segments

# TCP: retransmission scenarios

First Sent but unacknowledged sequence number
The sequence number at the beginning of the window

Host A                    Host B

**lost ACK scenario**

SND.UNA=92

timeout

Seq=92, 8 bytes of data

ACK=100
X

Seq=92, 8 bytes of data

ACK=100

SND.UNA=100

**premature timeout**

SND.UNA=92

timeout

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

ACK=100

ACK=120

SND.UNA=100

SND.UNA=120

Seq=92, 8 bytes of data

send cumulative ACK for 120

ACK=120

SND.UNA=120

# TCP: retransmission scenarios



Host A                    Host B

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

ACK=100

X

ACK=120

Seq=120,  15 bytes of data

cumulative ACK covers
for earlier lost ACK

# TCP: retransmission scenarios



Host A        Host B

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

ACK=100

**X**

ACK=120

Seq=120, 15 bytes of data

cumulative ACK covers
for earlier lost ACK

- (short) in class exercise:
  - What is the value of SND.UNA after sending and receiving each packet?

# TCP: retransmission scenarios

Host A                          Host B

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

ACK=100

**X**

ACK=120

Seq=120,  15 bytes of data

cumulative ACK covers
for earlier lost ACK

- Q: How is TCP similar to Go-Back-N? How is it different? How about Selective Repeat?

# Knowledge Check

- Make sure you understand and can complete a TCP send and receive timeline.

- This includes, but is not limited to
  - sequence and acknowledgement numbers on packets going back and forth
  - how the sender and receiver view of the sequence number space changes as a result of packets being sent and received (e.g., status of the bytes, position of the sliding window, etc.)

# TCP round trip time, timeout

*Q:* how to set TCP timeout value?

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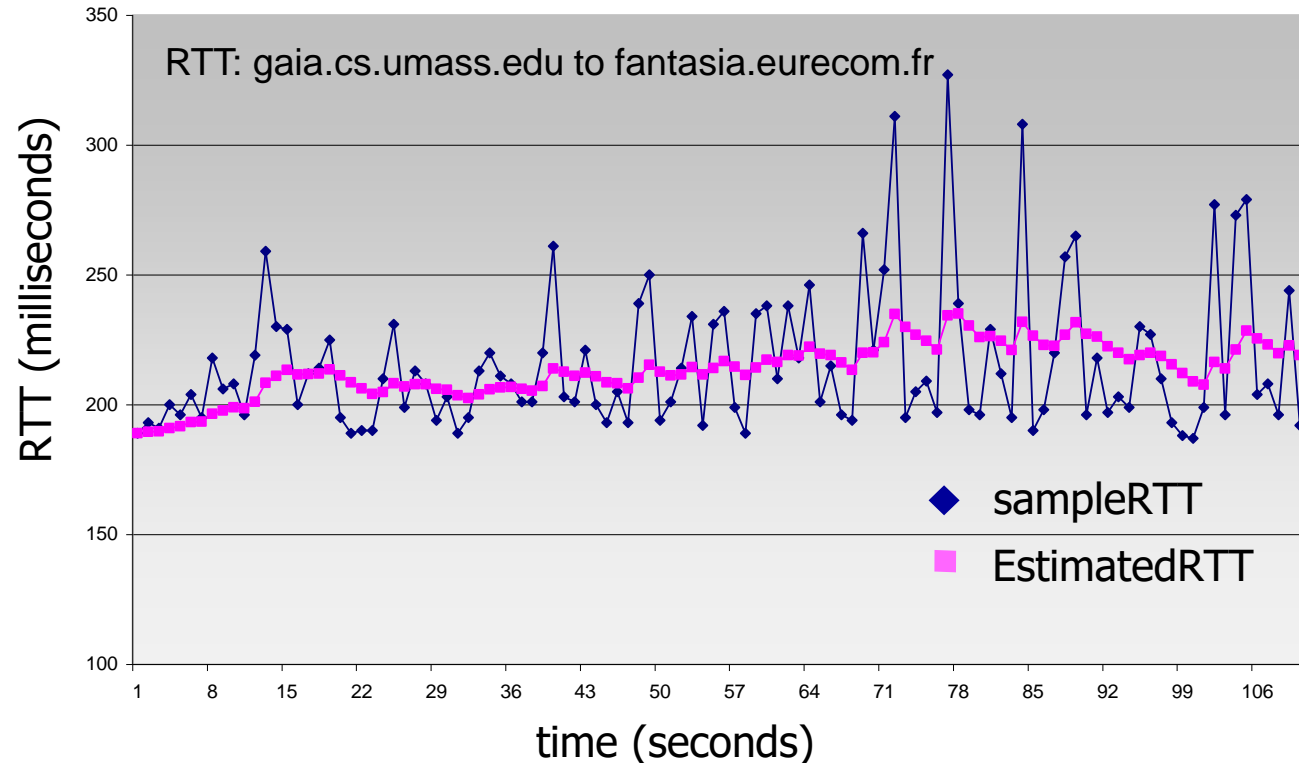- *too long:* slow reaction to segment loss

*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 (EWMA)
- influence of past sample decreases exponentially fast
- typical value: $\alpha$ = 0.125



RTT: gaia.cs.umass.edu to fantasia.eurecom.fr

# TCP round trip time, timeout

- timeout interval: `EstimatedRTT` plus "safety margin"

  - large variation in `EstimatedRTT`: want a larger safety margin

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

estimated RTT          "safety margin"

- `DevRTT`: EWMA of `SampleRTT` deviation from `EstimatedRTT`:

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

(typically, $\beta$ = 0.25)

* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

# Performance optimizations for TCP

- So far, we have covered "the basics" of TCP's rdt
  - Sequence number
  - Cumulative ACKs
  - Pipelined segments
  - Retransmission timer
- Next, we will discuss some optimizations
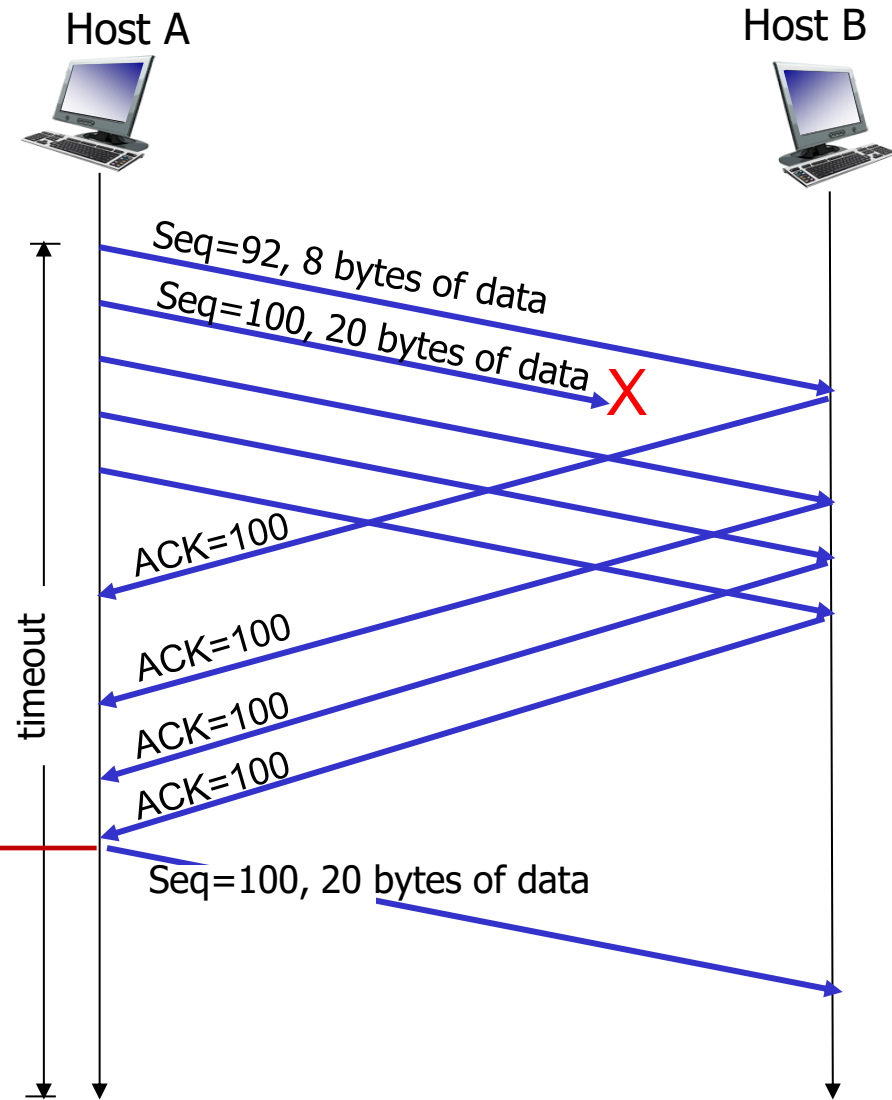
# Optimization 1: Fast Retransmit

TCP fast retransmit

if sender receives 3 additional ACKs for same data ("triple duplicate ACKs"), resend unACKed segment with smallest seq #

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

💡 Receipt of three duplicate ACKs indicates 3 segments received after a missing segment – lost segment is likely. So retransmit!

Host A

Host B

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

X

ACK=100

ACK=100

ACK=100

ACK=100

timeout

Seq=100, 20 bytes of data

# Optimization 2: Delayed ACKs

- Instead of generating an ACK in response to every segment the moment it arrives
  - Wait for some time to see if there is another segment right afterwards
  - Create one ACK for both.
- Benefits?
  - Saves bandwidth
- Disadvantages?
  - Increases delay in responding to the sender.
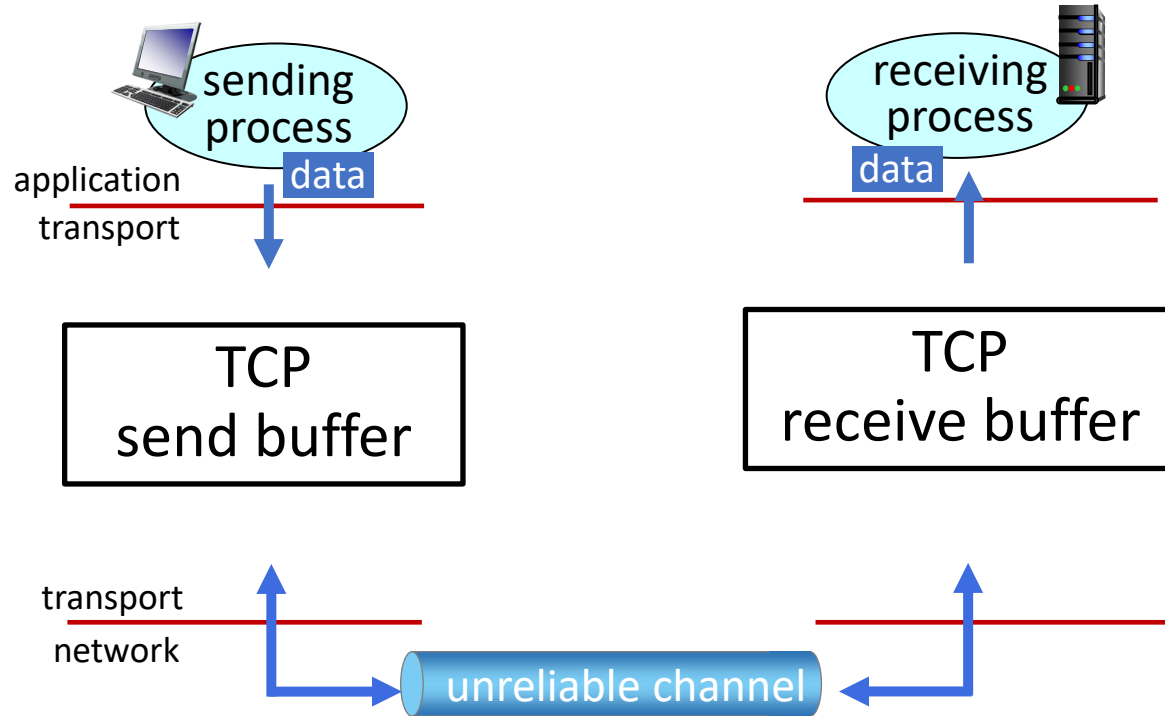
# Optimizations 2: Delays ACKs (cont.)

| Event at receiver | TCP receiver action |
|---|---|
| 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 *duplicate ACK,* indicating seq. # of next expected byte |

# Transport layer: roadmap

- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- **Connection-oriented transport: TCP**
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- Principles of congestion control
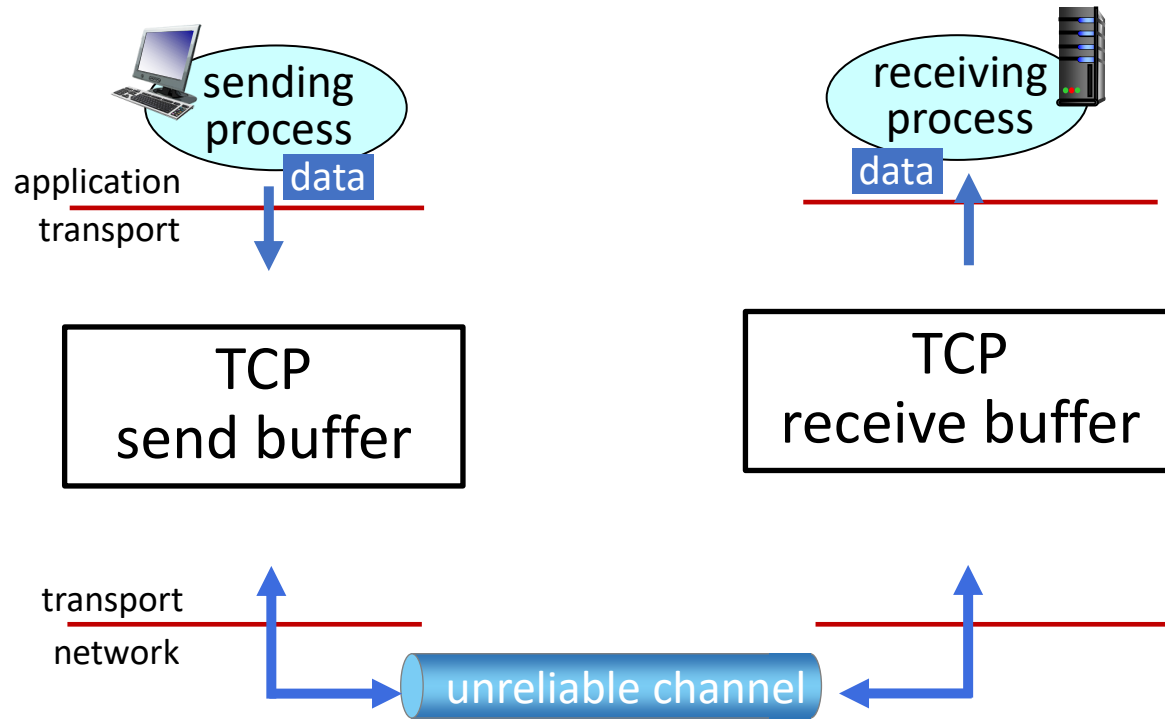- TCP congestion control
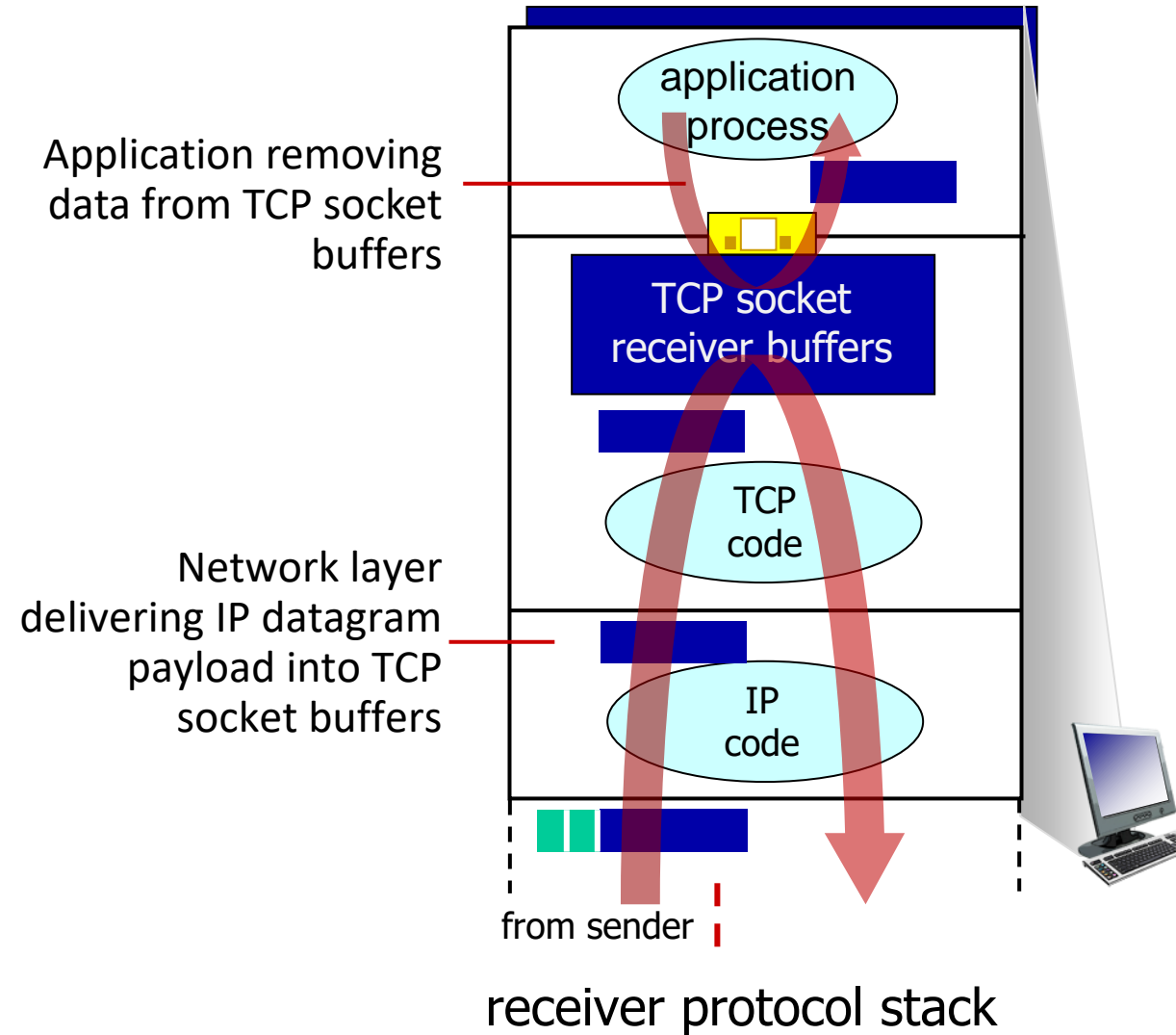
# TCP flow control



- The send buffer holds the data the application sends to TCP until it is delivered
- The receive buffer holds the data TCP receives from the network until it is delivered to the application
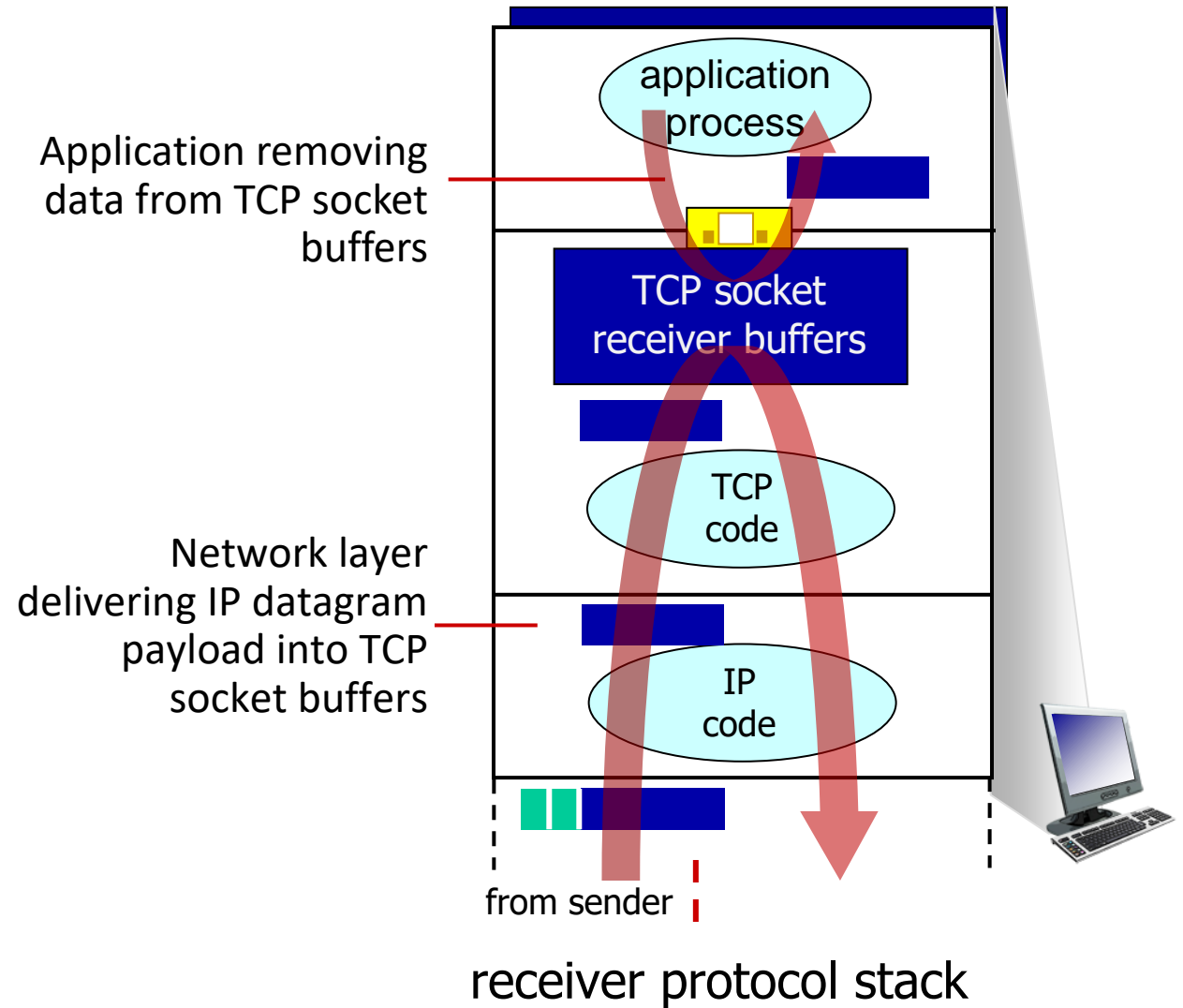
# TCP flow control

sending process

data

application
transport

TCP send buffer

transport
network

unreliable channel

receiving process

data

TCP receive buffer

application
transport

transport
network

application process

Application removing data from TCP socket buffers

TCP socket receiver buffers

TCP code

Network layer delivering IP datagram payload into TCP socket buffers
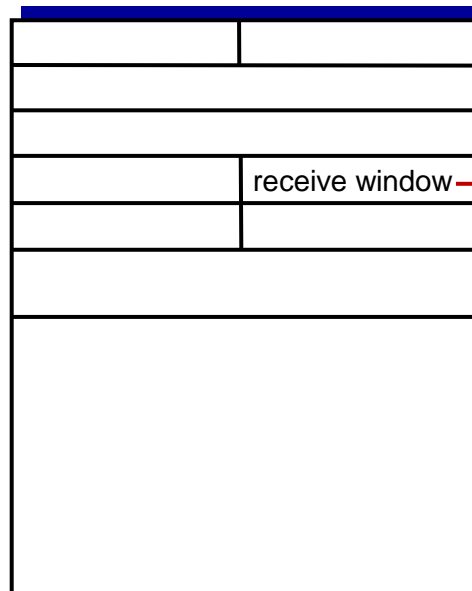
IP code

from sender

receiver protocol stack

*Q:* What happens if network layer delivers data faster than application layer removes data from socket buffers?

# TCP flow control

*Q:* What happens if network layer delivers data faster than application layer removes data from socket buffers?
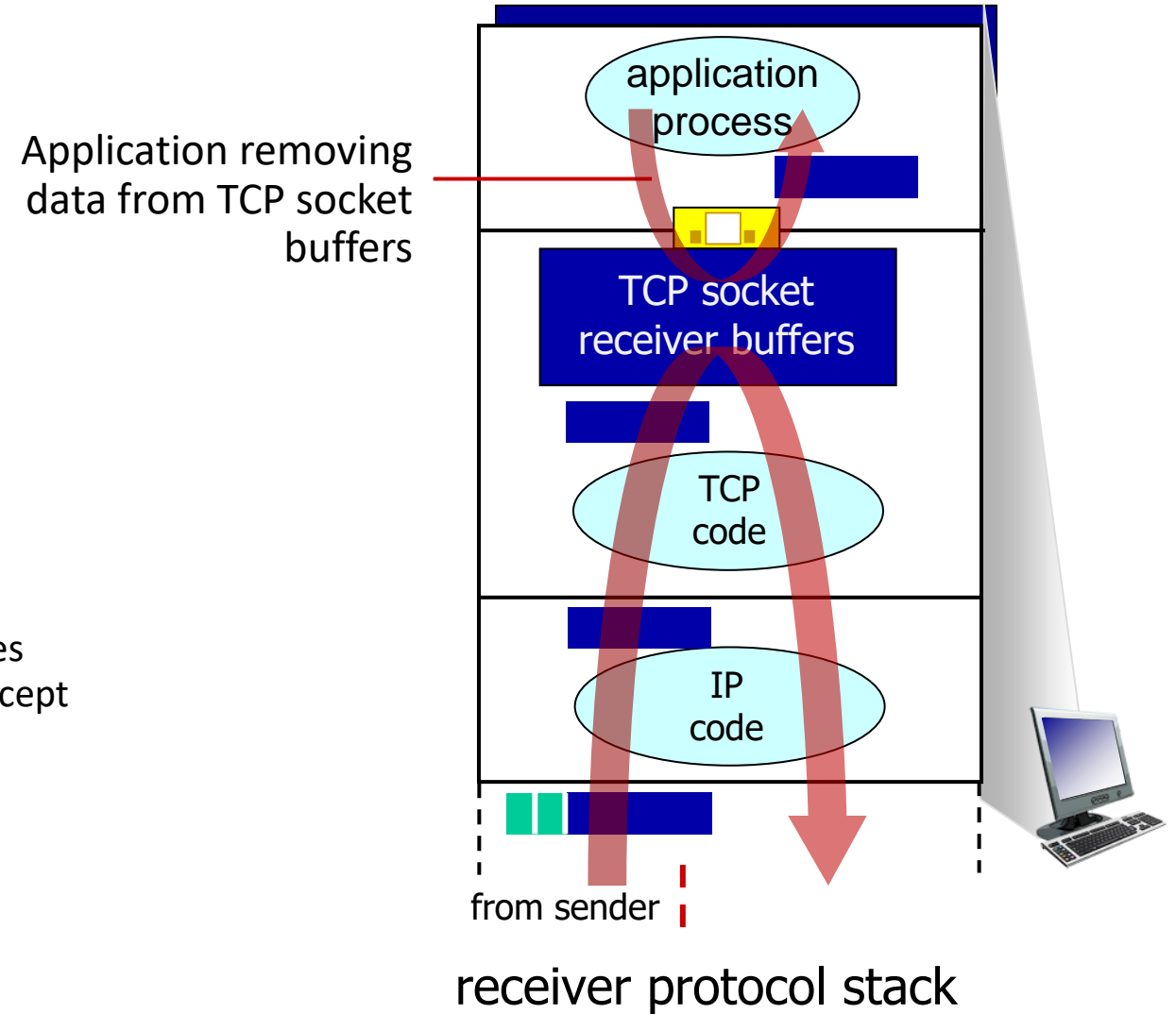


Application removing data from TCP socket buffers

application process

TCP socket receiver buffers

TCP code

Network layer delivering IP datagram payload into TCP socket buffers

IP code

from sender

receiver protocol stack

# TCP flow control

*Q:* What happens if network layer delivers data faster than application layer removes data from socket buffers?

receive window

flow control: # bytes receiver willing to accept

Application removing data from TCP socket buffers

application process

TCP socket receiver buffers

TCP code

IP code

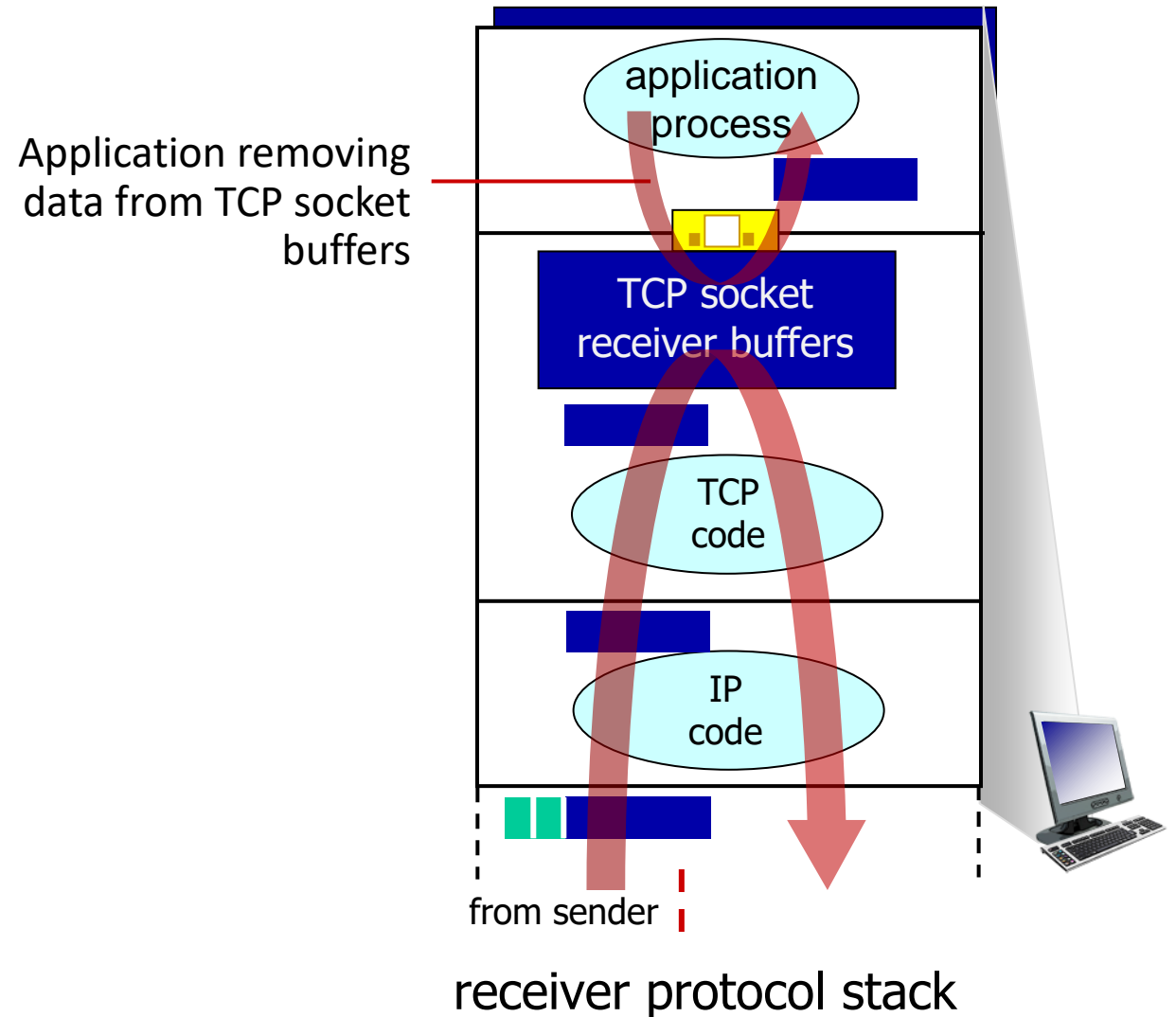from sender

receiver protocol stack

# TCP flow control

*Q:* What happens if network layer delivers data faster than application layer removes data from socket buffers?

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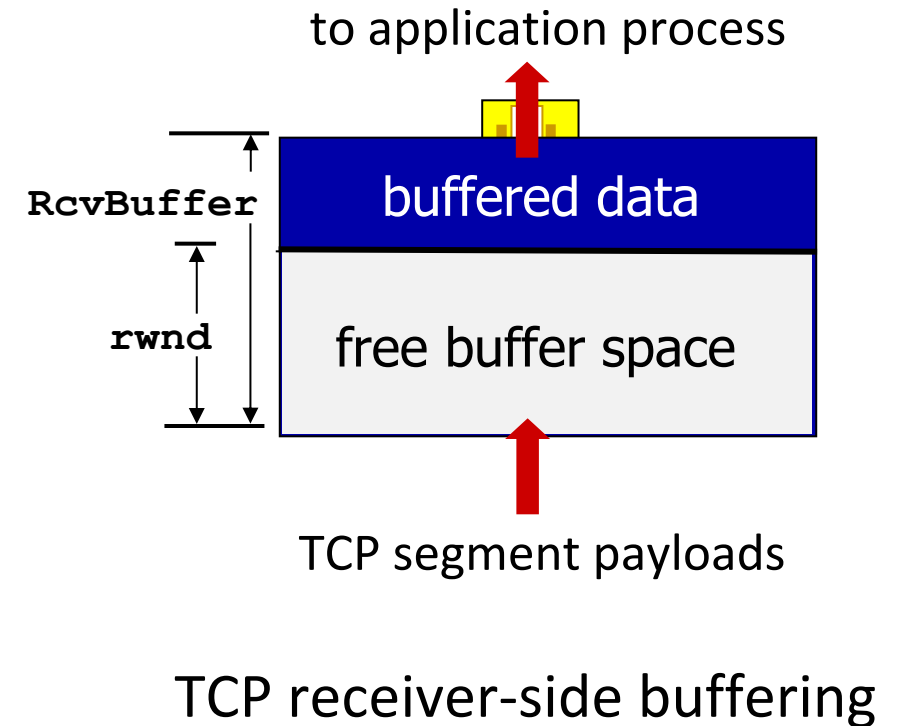
receiver controls sender, so sender won't overflow receiver's buffer by transmitting too much, too fast

Application removing data from TCP socket buffers

application process

TCP socket receiver buffers

TCP code

IP code
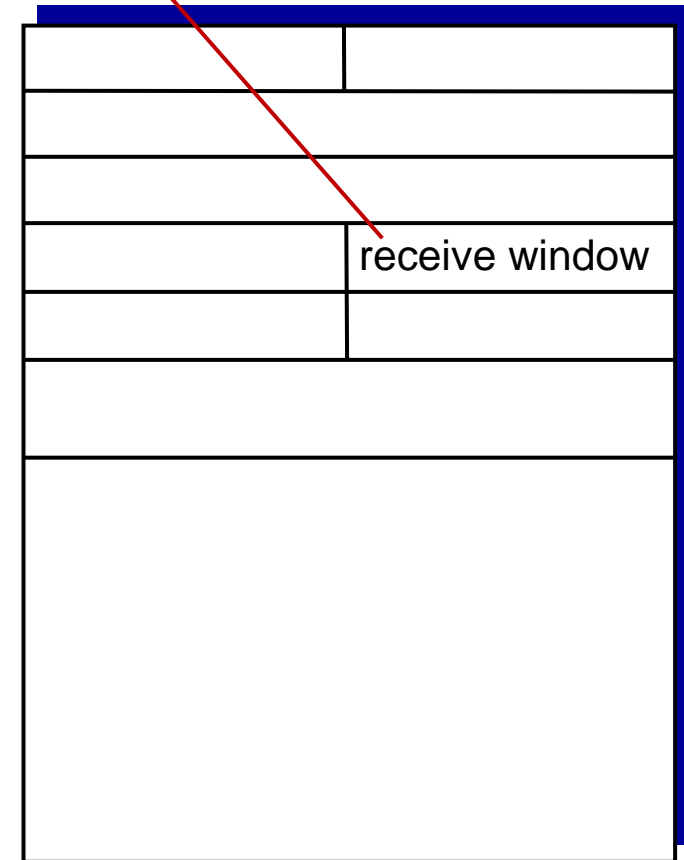
from sender

receiver protocol stack

# TCP flow control

- TCP receiver "advertises" free buffer space in `rwnd` field in TCP header
  - `RcvBuffer` size set via socket options
  - many operating systems auto-adjust `RcvBuffer`

- sender limits amount of unACKed ("in-flight") data to received `rwnd`

- guarantees receive buffer will not overflow



TCP receiver-side buffering

# TCP flow control

- TCP receiver "advertises" free buffer space in **rwnd** field in TCP header
  - **RcvBuffer** size set via socket options
  - many operating systems auto-adjust **RcvBuffer**

- sender limits amount of unACKed ("in-flight") data to received **rwnd**

- guarantees receive buffer will not overflow

flow control: # bytes receiver willing to accept

receive window
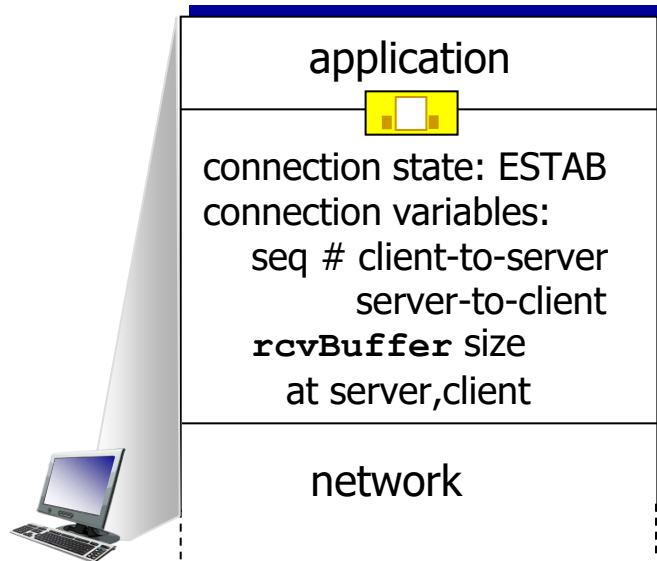
TCP segment format

# Transport layer: roadmap

- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- **Connection-oriented transport: TCP**
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- Principles of congestion control
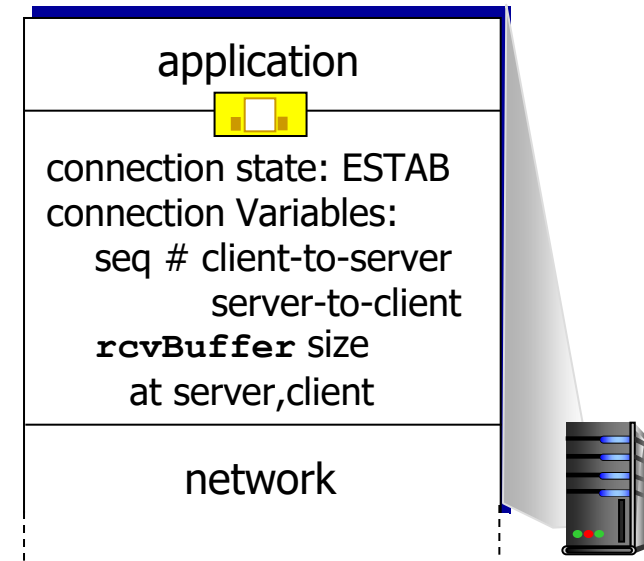- TCP congestion control

# TCP connection management

before exchanging data, sender/receiver "handshake":
- agree to establish connection (each knowing the other willing to establish connection)
- agree on connection parameters (e.g., starting seq #s)



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

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

# TCP 3-way handshake

## Client state

```
clientSocket = socket(AF_INET, SOCK_STREAM)
```

LISTEN

```
clientSocket.connect((serverName,serverPort))
```

Pay close attention to sequence and ack numbers during handshake

## Server state

```
serverSocket = socket(AF_INET,SOCK_STREAM)
serverSocket.bind(('',serverPort))
serverSocket.listen(1)
connectionSocket, addr = serverSocket.accept()
```

LISTEN

choose init seq num, x
send TCP SYN msg

SYNSENT

SYNbit=1, Seq=x

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

SYN RCVD

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

ESTAB

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

ACKbit=1, ACKnum=y+1

received ACK(y)
indicates client is live

ESTAB

# A human 3-way handshake protocol

# Closing a TCP 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
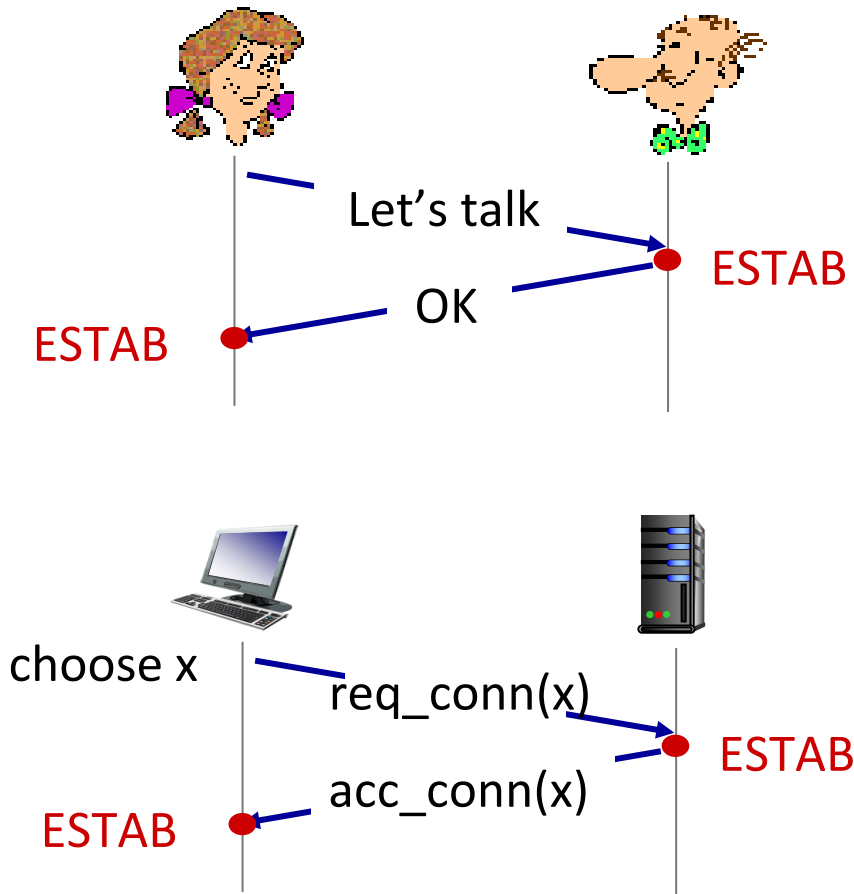- simultaneous FIN exchanges can be handled

# Knowledge Check

- Make sure you understand and can complete a TCP connection timeline

  - From connection establishment, through reliable data transfer (with optimizations and flow control), to connection tear-down

- This includes, but is not limited to

  - sequence and acknowledgement numbers on packets going back and forth
  - how the sender and receiver view of the sequence number space changes as a result of packets being sent and received (e.g., status of the bytes, position of the sliding window, etc.)

# Additional Slides
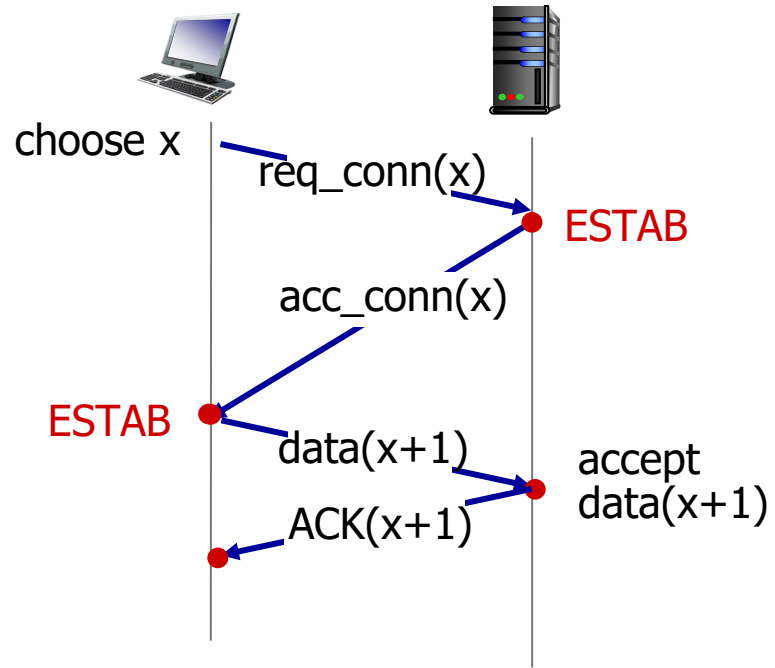
# Agreeing to establish a connection

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
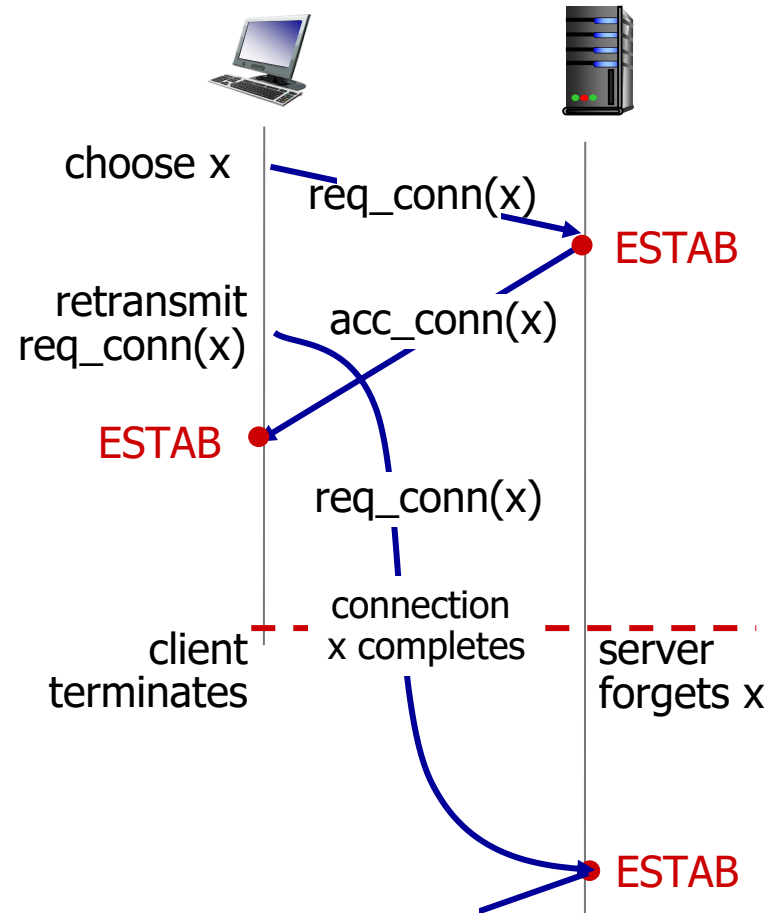- message reordering
- can't "see" other side

# 2-way handshake scenarios

choose x

req_conn(x)

ESTAB

acc_conn(x)

ESTAB

data(x+1)

accept
data(x+1)

ACK(x+1)

No problem!

✅

# 2-way handshake scenarios

choose x

req_conn(x)

ESTAB

retransmit
req_conn(x)

acc_conn(x)

ESTAB

req_conn(x)

connection
x completes

client
terminates

server
forgets x

ESTAB

❌ Problem: half open
connection! (no client)

# 2-way handshake scenarios



choose x
req_conn(x)
ESTAB
retransmit
req_conn(x)
acc_conn(x)
ESTAB
data(x+1)
accept
data(x+1)
retransmit
data(x+1)
connection
x completes
client
terminates
server
forgets x

req_conn(x)

ESTAB

data(x+1)

accept
data(x+1)

**Problem: dup data accepted!**