# Chapter 3: roadmap

# TCP: overview RFCs: 793,1122, 2018, 5681, 7323

- **point-to-point:**
  - one sender, one receiver

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

- **full duplex data:**
  - bi-directional data flow in same connection
  - MSS: maximum segment size

- **cumulative ACKs**

- **pipelining:**
  - TCP congestion and flow control set window size

- **connection-oriented:**
  - handshaking (exchange of control messages) initializes sender, receiver state before data exchange

- **flow controlled:**
  - sender will not overwhelm receiver

# TCP: overview  RFCs: 793,1122, 2018, 5681, 7323

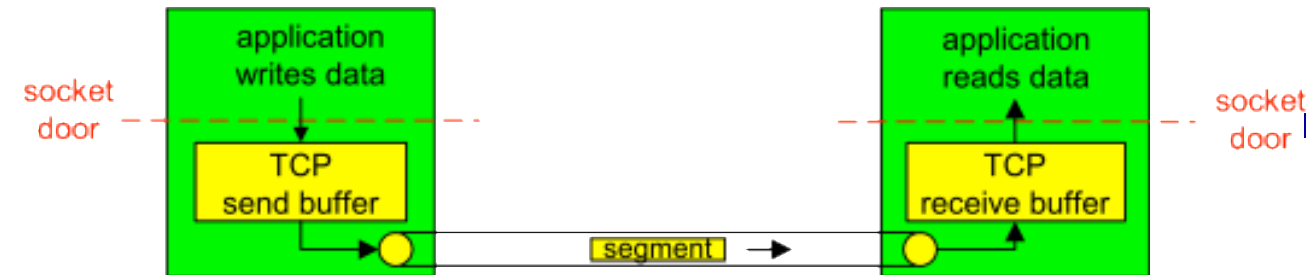- **cumulative ACKs**

- **pipelining:**
  - TCP congestion and flow control set window size

- **connection-oriented:**
  - handshaking (exchange of control messages) initializes sender, receiver state before data exchange

- **flow controlled:**
  - sender will not overwhelm receiver



socket door

application writes data

TCP send buffer

segment →

application reads data
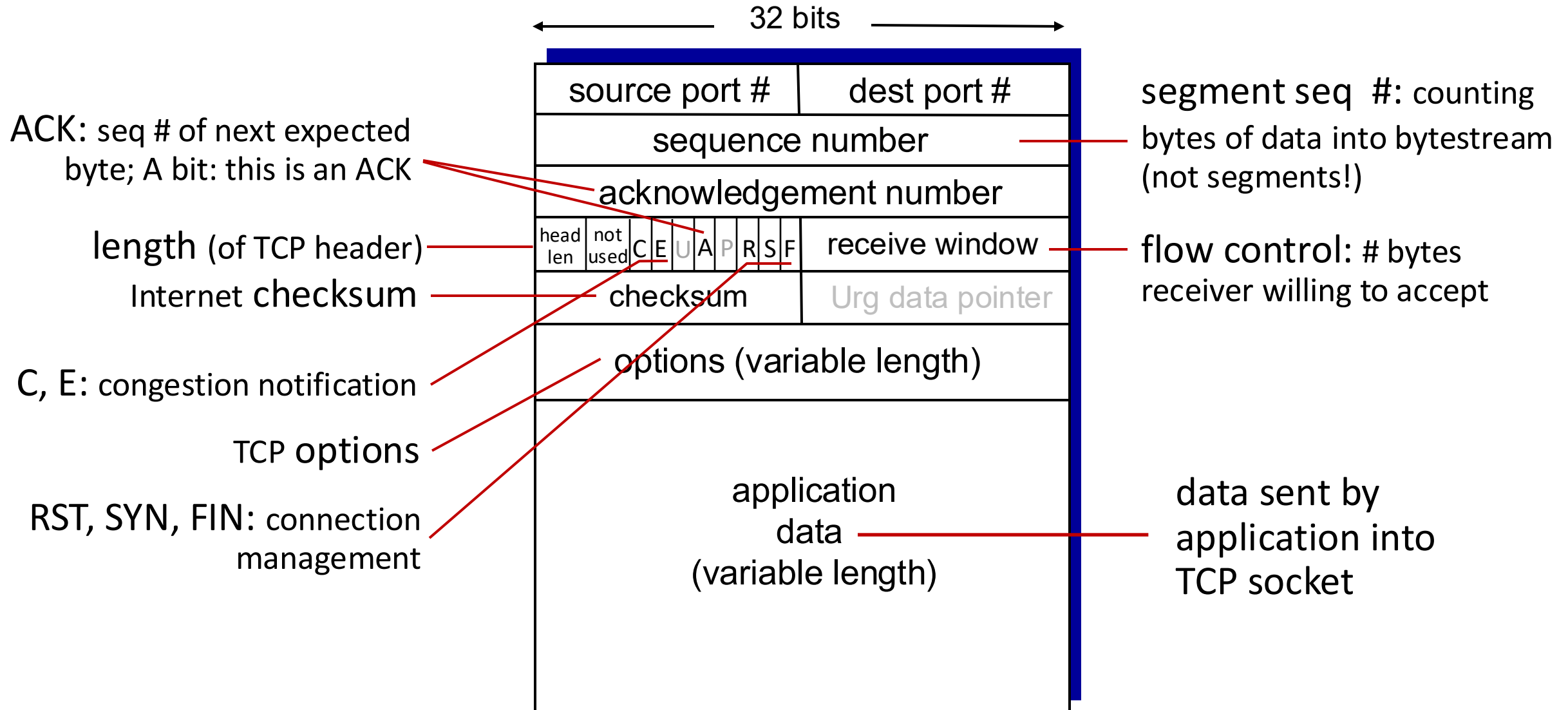
socket door

TCP receive buffer

# Question

- Given what we have discussed regarding protocols for a variety of lower level guarantees.

- And given that the network layer guarantees...nothing

  - ‣ Packets can be dropped

  - ‣ Packets can be duplicated

  - ‣ Packets can arrive out of order

  - ‣ Packets can arrive corrupted

- What do we need in the TCP header?

# 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

# TCP Options Examples

- **Maximum Segment Size (MSS):** Negotiates the largest amount of data, in bytes, that a TCP segment can carry, ensuring it does not exceed the network's Maximum Transmission Unit (MTU).

- **Window Scaling:** Allows for a larger receive window than the 16-bit field normally permits, which is crucial for high-bandwidth, high-latency networks.

- **Timestamps:** Helps compute an accurate Round Trip Time (RTT) and prevents issues with old duplicate packets on a network.

- **Selective Acknowledgements (SACK):** Allows the receiver to tell the sender which specific segments have been received, even if there are gaps, preventing the sender from having to retransmit all data after a single packet loss.

# TCP sequence numbers, ACKs

*Sequence numbers:*

- byte stream "number" of first byte in segment's data
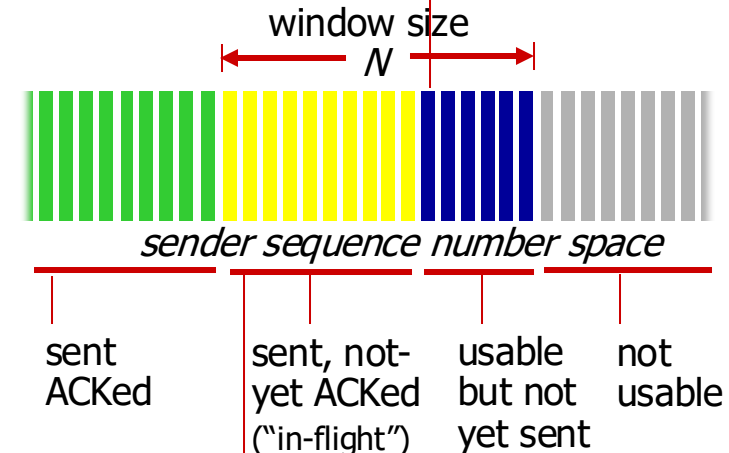
*Acknowledgements*:

- seq # of next byte expected from other side

- cumulative ACK

*Q*: how receiver handles out-of-order segments

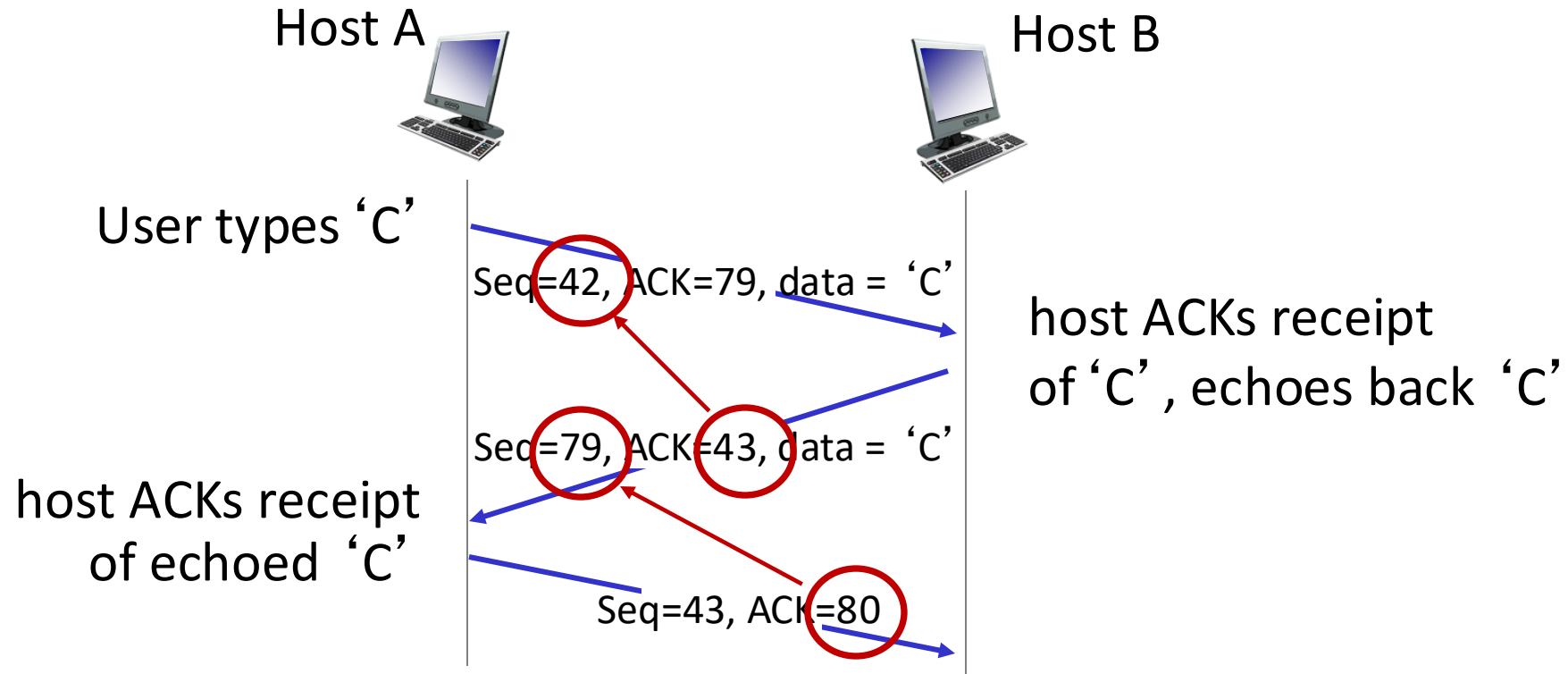- *A:* TCP spec doesn't say, - up to implementor

outgoing segment from sender

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

window size
N

*sender sequence number space*

sent
ACKed

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

usable but not yet sent

not usable

outgoing segment from receiver

| 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

# Another Question

- We know that when packets can be dropped, we need to set a timer.

- The question is, to what value?

- Think about that a bit.

  - And about what difficulties we might encounter

  - And think about the consequences if we get this wrong

# TCP round trip time, timeout

*Q:* how to set TCP timeout value?

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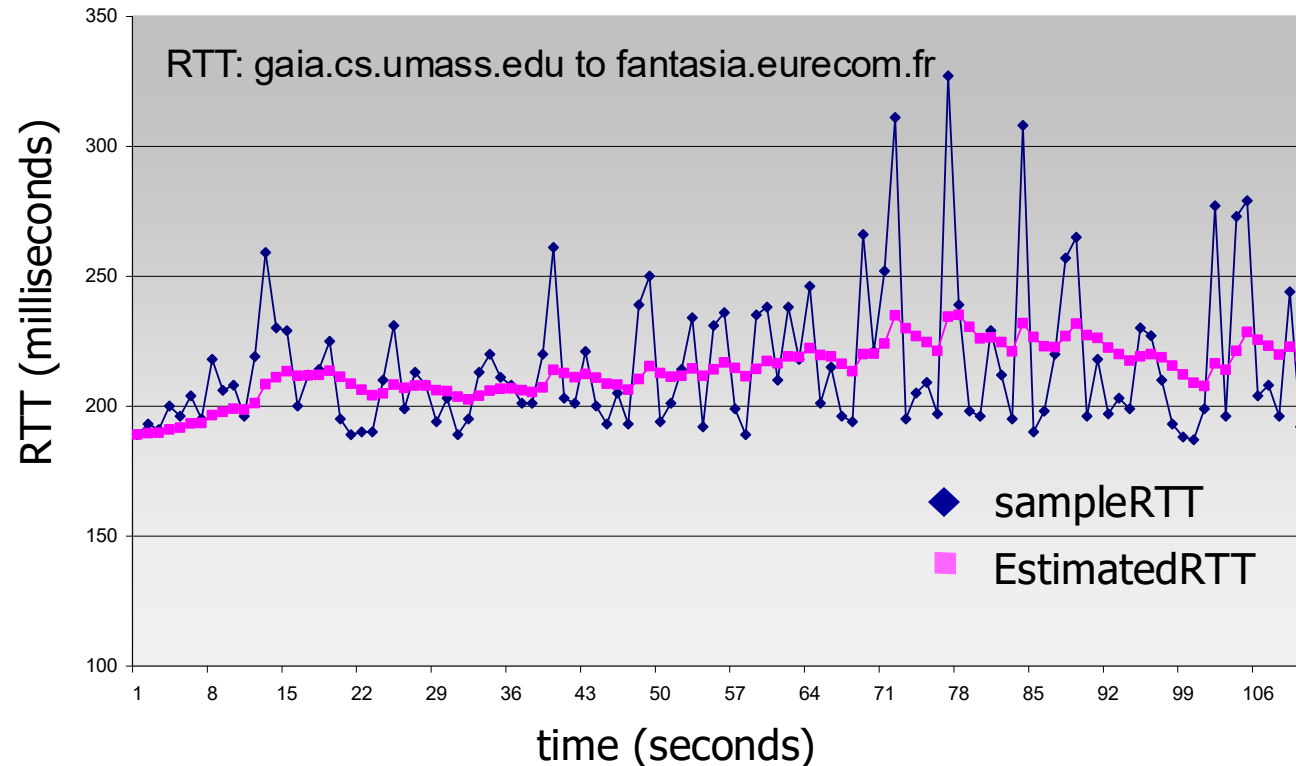- *too long:* slow reaction to segment loss

*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

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

- exponential <u>w</u>eighted <u>m</u>oving <u>a</u>verage (EWMA)
- influence of past sample decreases exponentially fast
- typical value: $\alpha$ = 0.125



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

◆ sampleRTT
■ EstimatedRTT

RTT (milliseconds)

time (seconds)

# TCP round trip time, timeout

▪ timeout interval: **EstimatedRTT** plus "safety margin"

  • large variation in **EstimatedRTT**:  want a larger safety margin

**TimeoutInterval = EstimatedRTT + 4\*DevRTT**

estimated RTT        "safety margin"

▪ **DevRTT**: EWMA of **SampleRTT**  deviation from **EstimatedRTT**:

**DevRTT = (1-β)\*DevRTT + β\*|SampleRTT-EstimatedRTT|**

(typically, β = 0.25)

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

# TCP Sender (simplified)

**event: data received from application**

- 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`
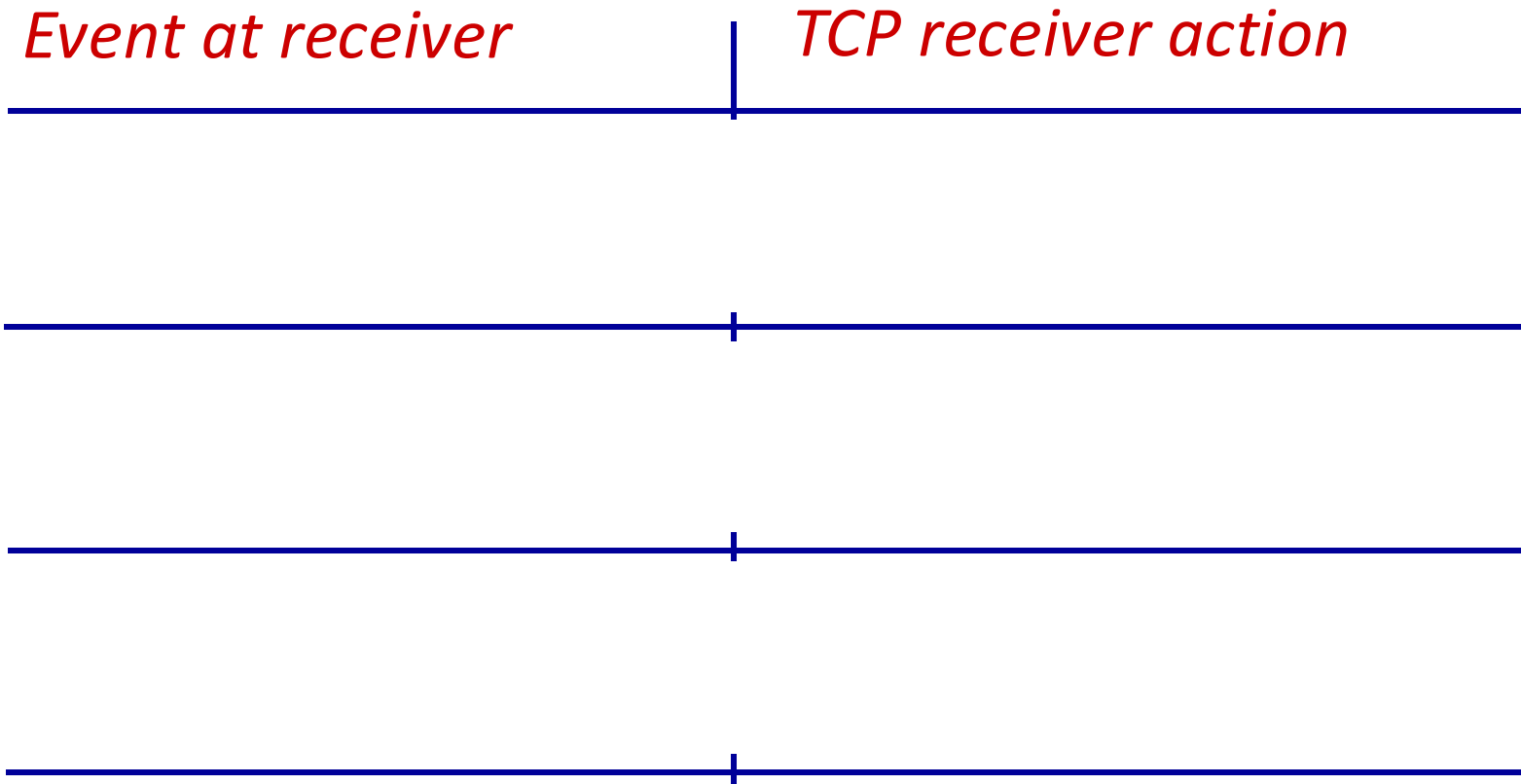
*event: timeout*

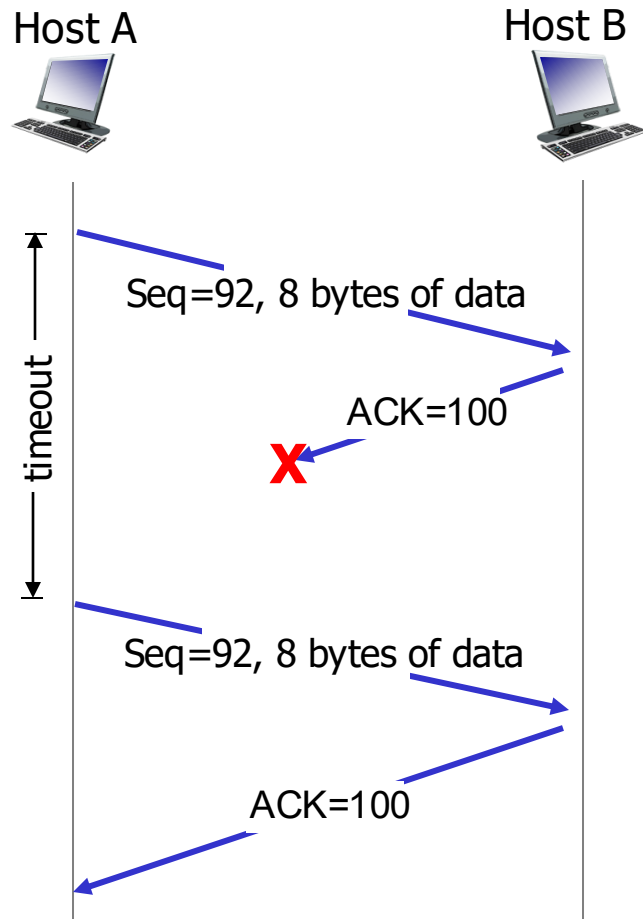- retransmit segment that caused timeout
- restart timer

*event: ACK received*

- if ACK acknowledges previously unACKed segments
  - update what is known to be ACKed
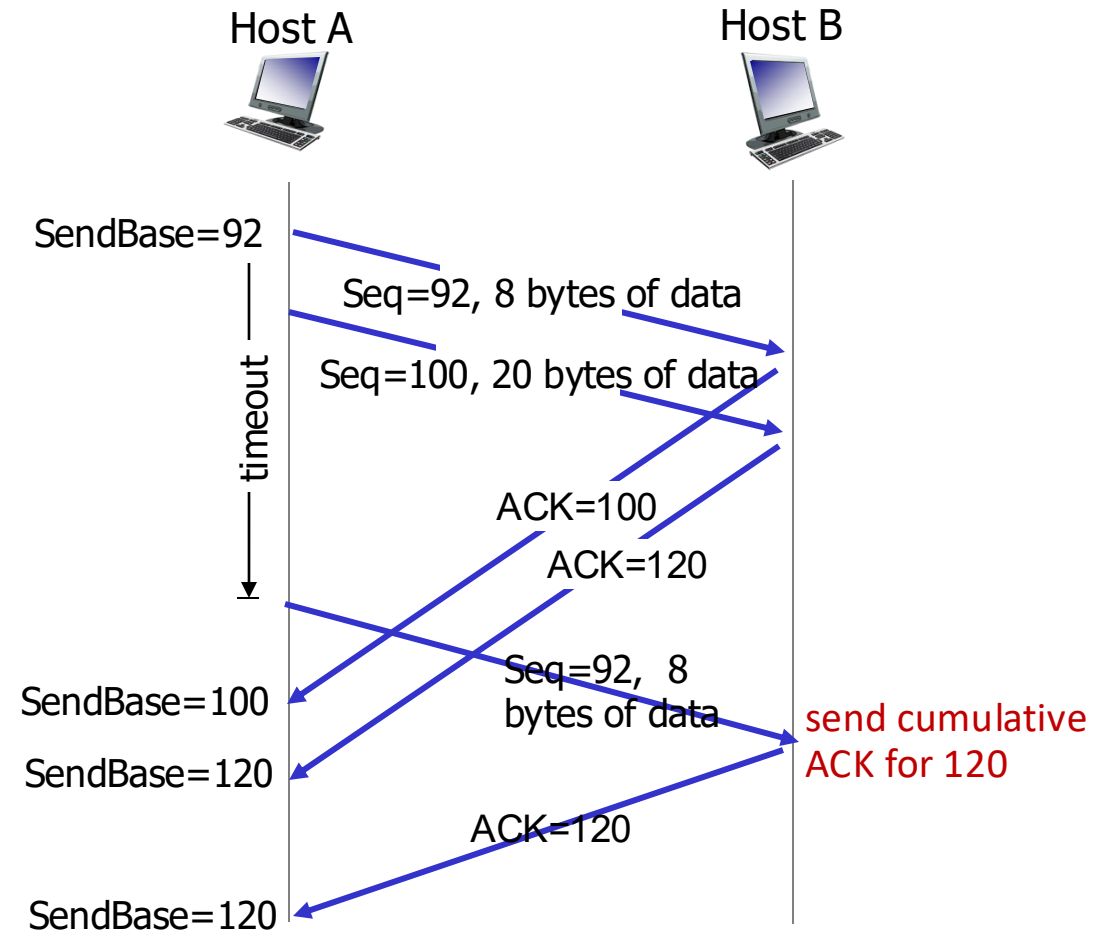  - start timer if there are still unACKed segments
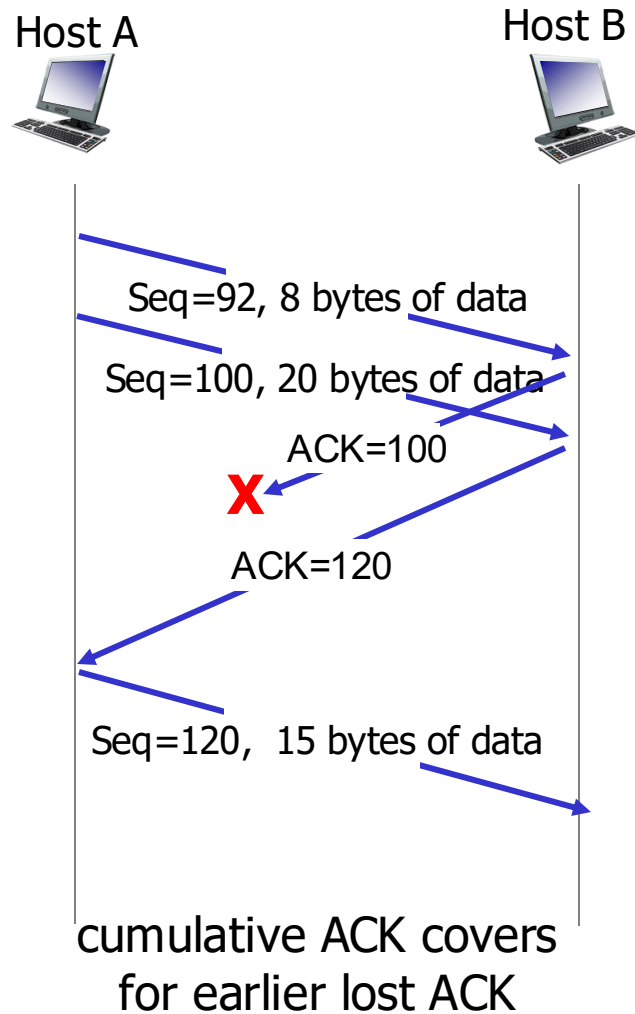
# TCP Receiver: ACK generation [RFC 5681]

| Event at receiver | TCP receiver action |
|---|---|
| | |
| | |
| | |

# TCP: retransmission scenarios



lost ACK scenario

premature timeout

# TCP: retransmission scenarios



cumulative ACK covers
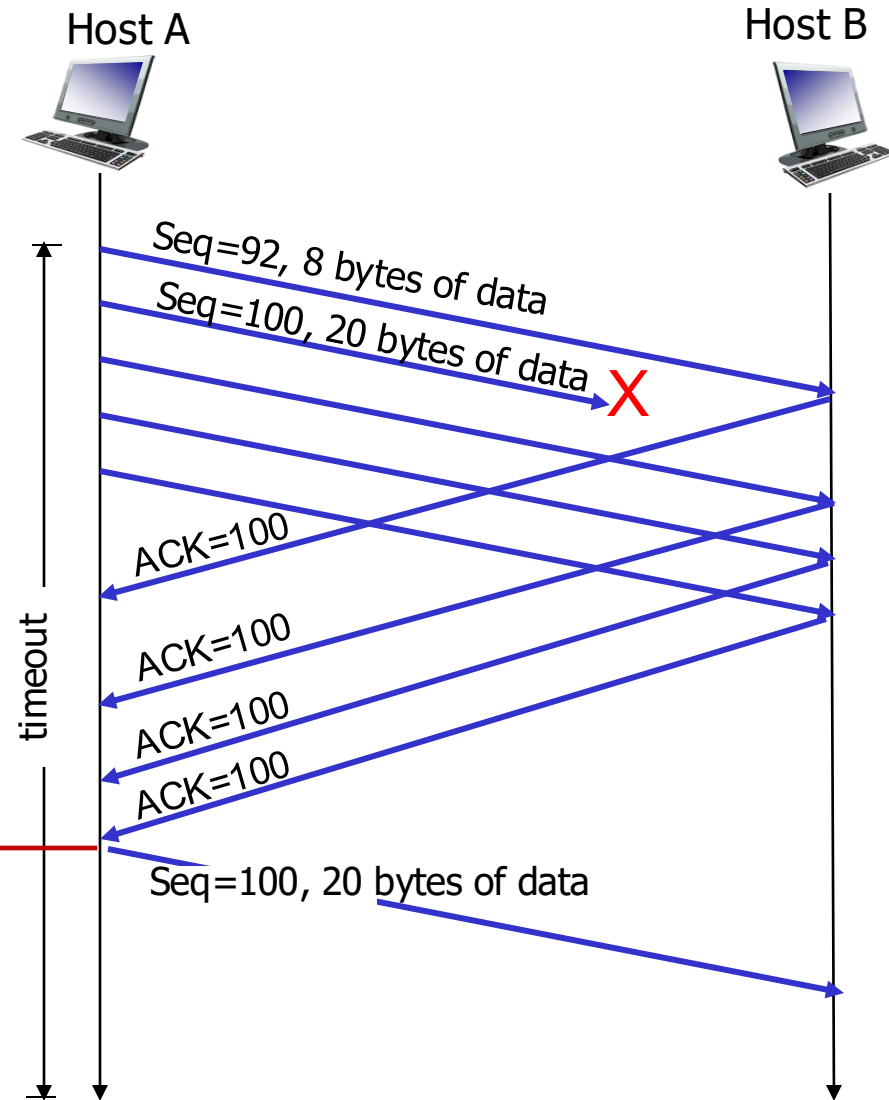for earlier lost ACK

# TCP 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

timeout

ACK=100
ACK=100
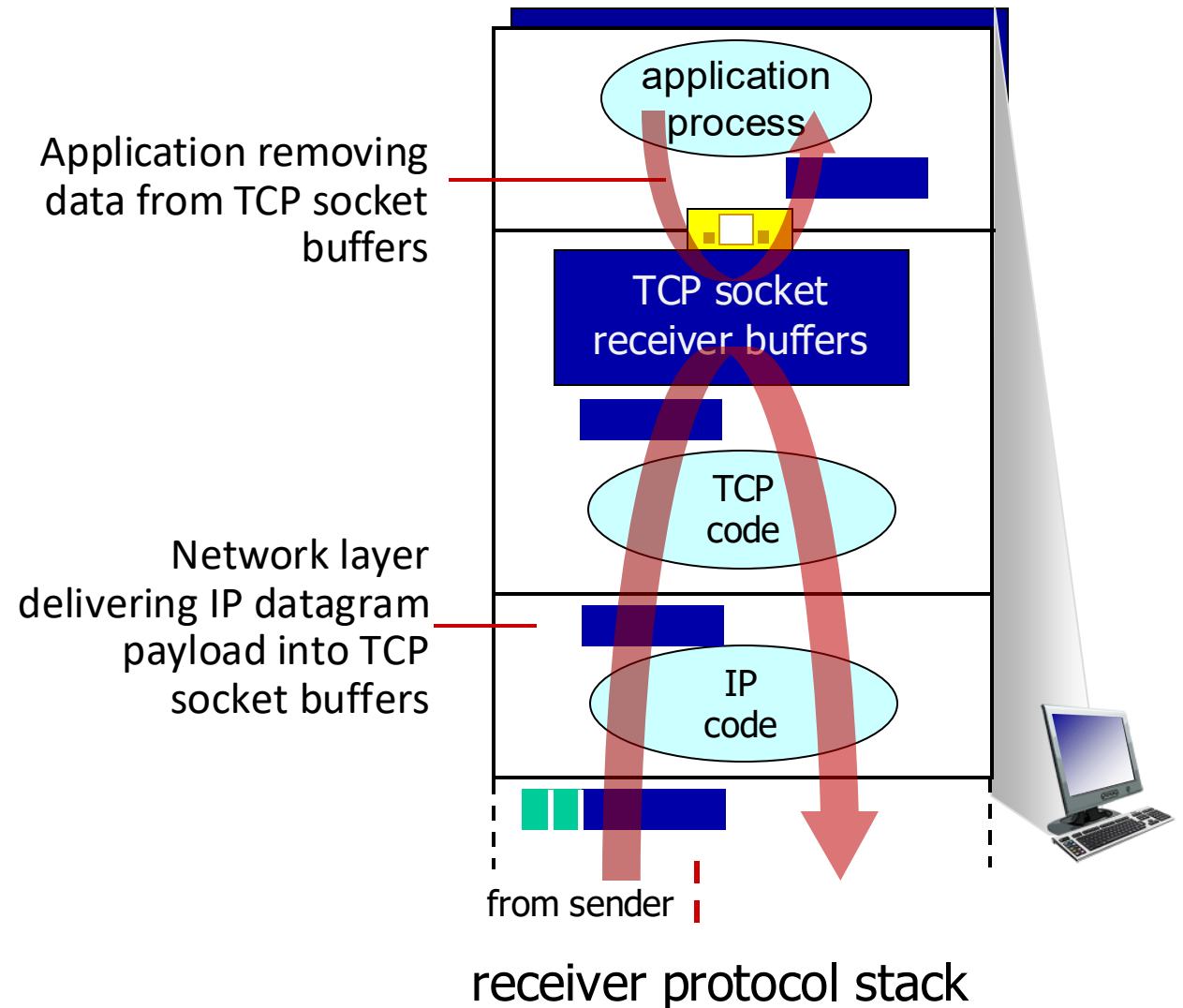ACK=100
ACK=100

Seq=100, 20 bytes of data

# Chapter 3: 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**
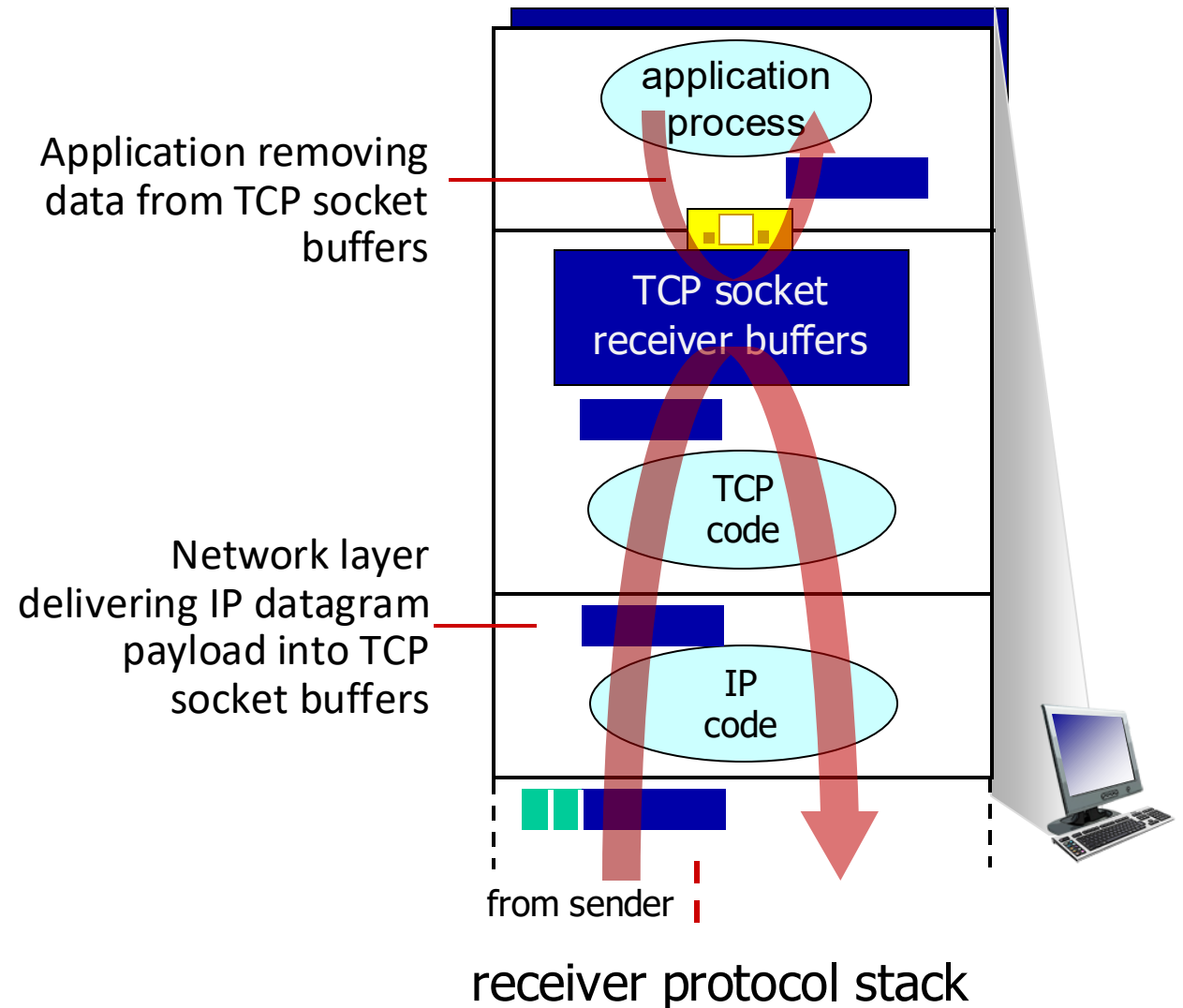- Principles of congestion control
- TCP congestion control

# TCP flow control

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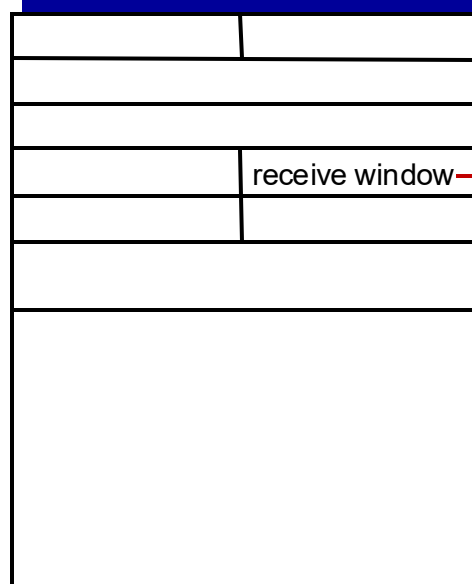
Application removing data from TCP socket buffers

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?



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

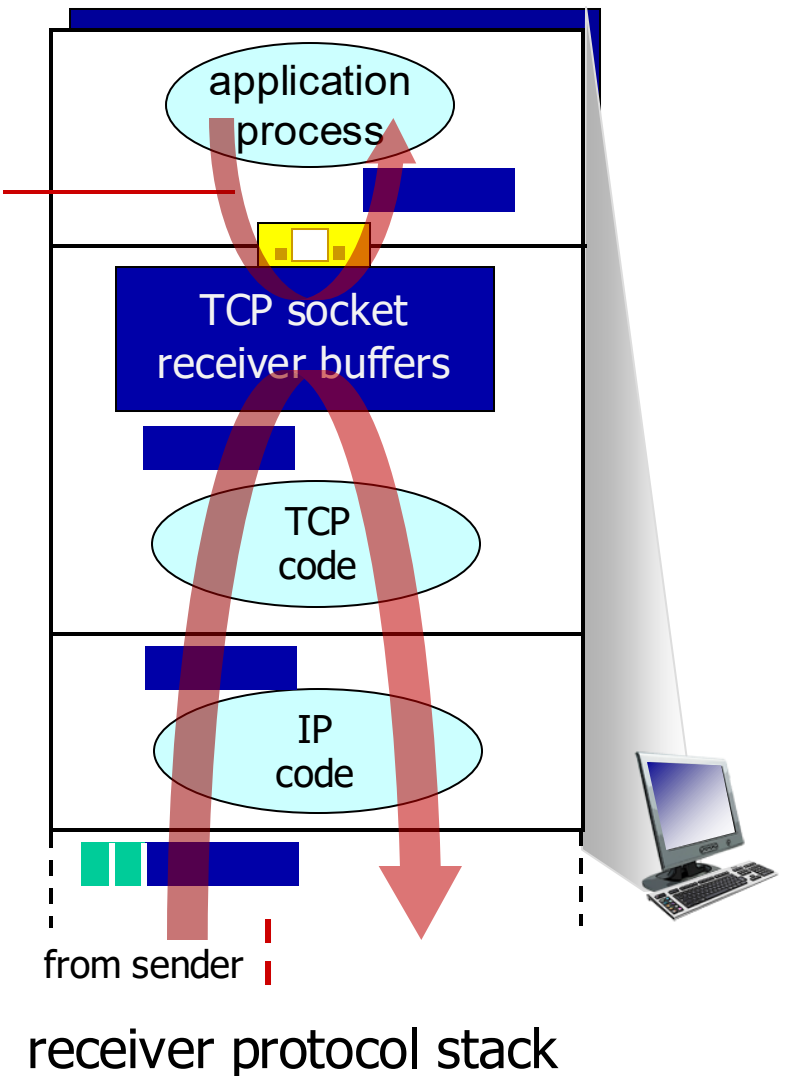*Q:* What happens if network layer delivers data faster than application layer removes data from socket buffers?

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

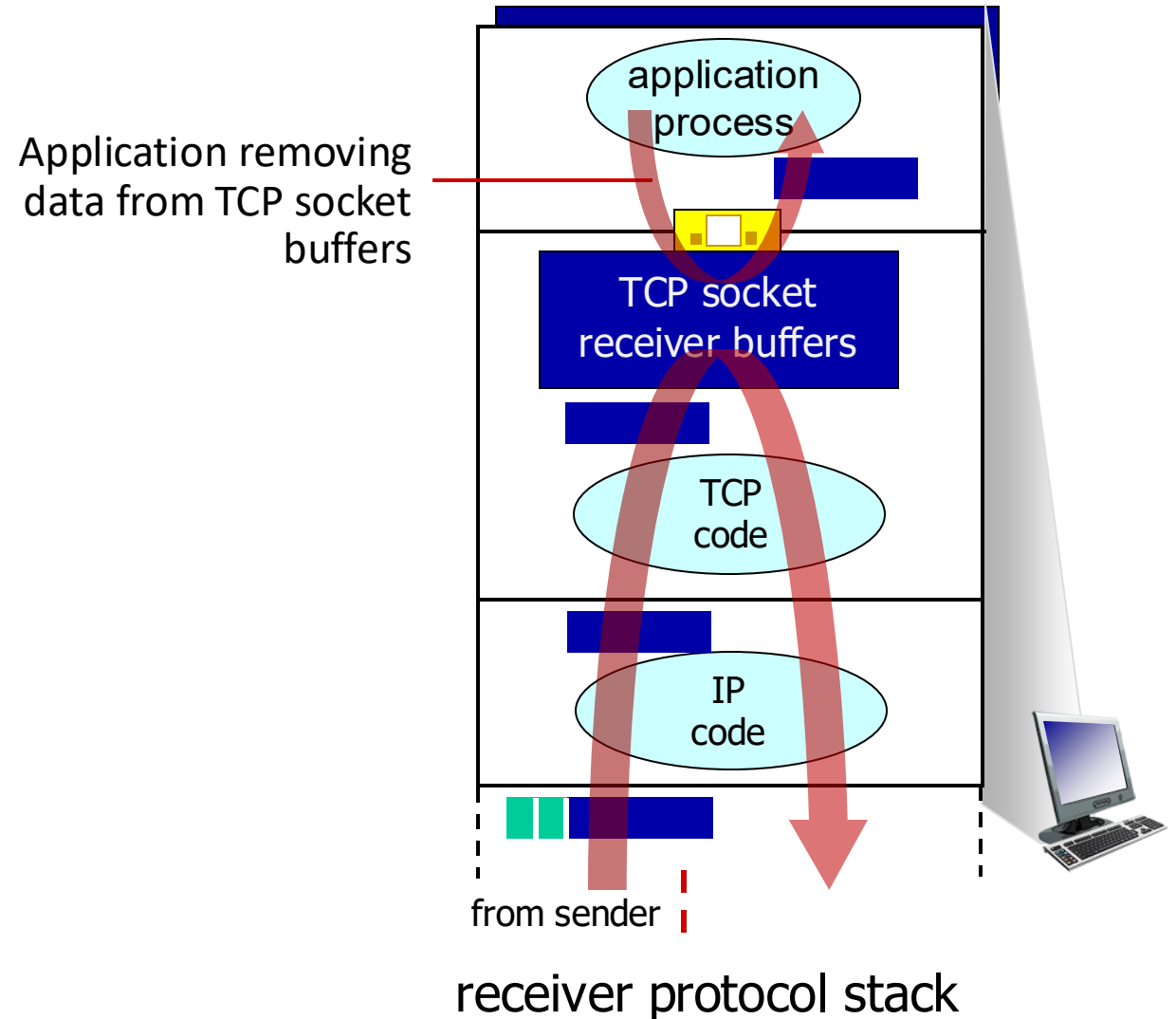from sender

receiver protocol stack

# TCP flow control

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

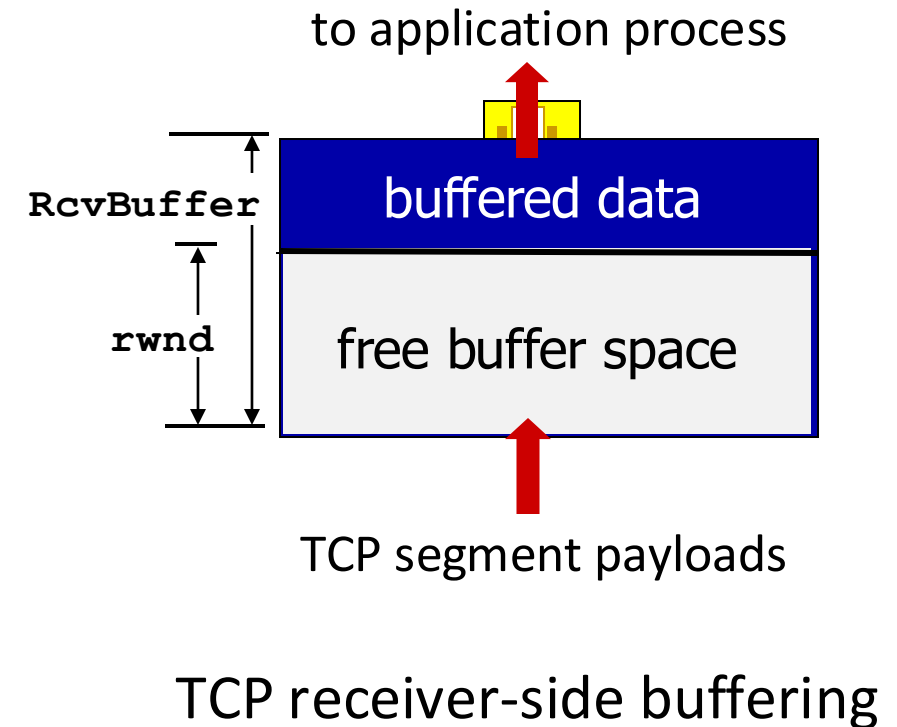**flow control**
receiver controls sender, so sender won't overflow receiver's buffer by transmitting too much, too fast

Application removing data from TCP socket buffers

application process

TCP socket receiver buffers

TCP code

IP code

from sender

**receiver protocol stack**
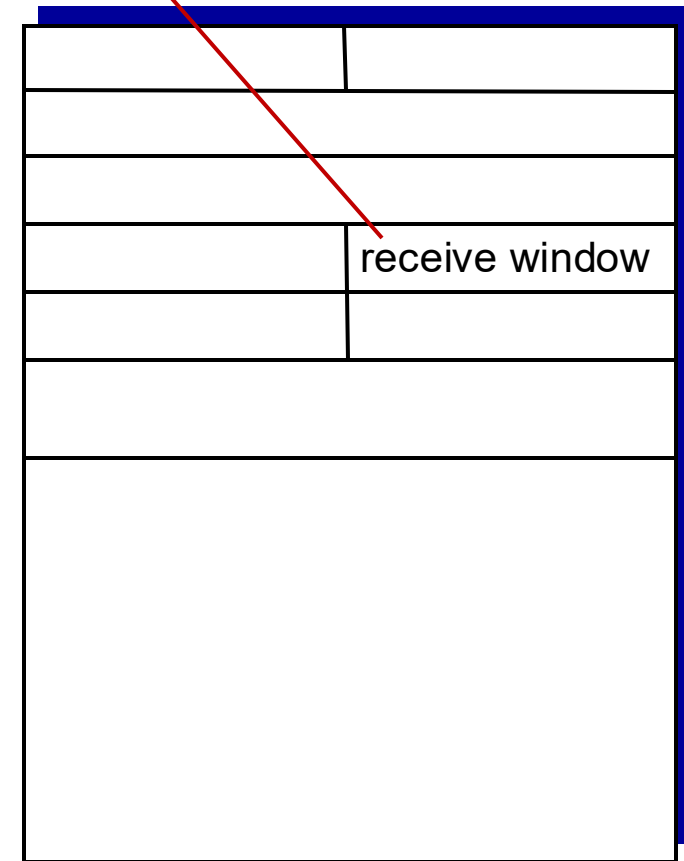
# TCP flow control

- TCP receiver "advertises" free buffer space in **rwnd** field in TCP header

    - **RcvBuffer** size set via socket options (typical default is 4096 bytes)

    - many operating systems autoadjust **RcvBuffer**

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

- guarantees receive buffer will not overflow

to application process

RcvBuffer

rwnd

buffered data

free buffer space

TCP segment payloads

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 (typical default is 4096 bytes)

  - many operating systems autoadjust **RcvBuffer**

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

- guarantees receive buffer will not overflow

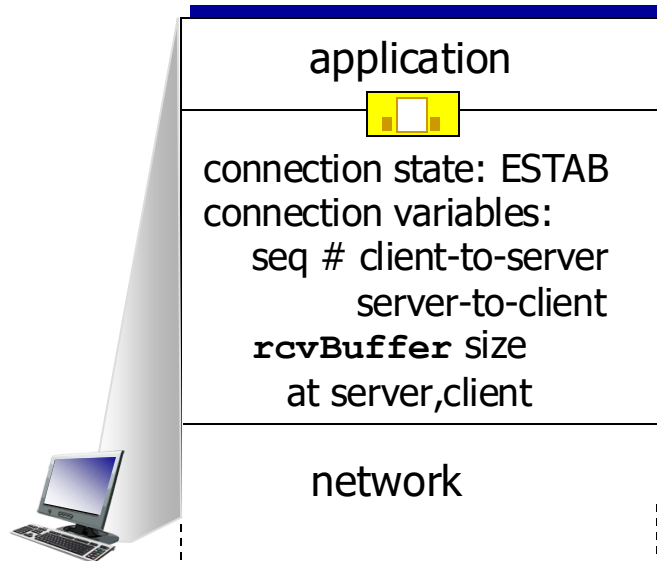flow control: # bytes receiver willing to accept
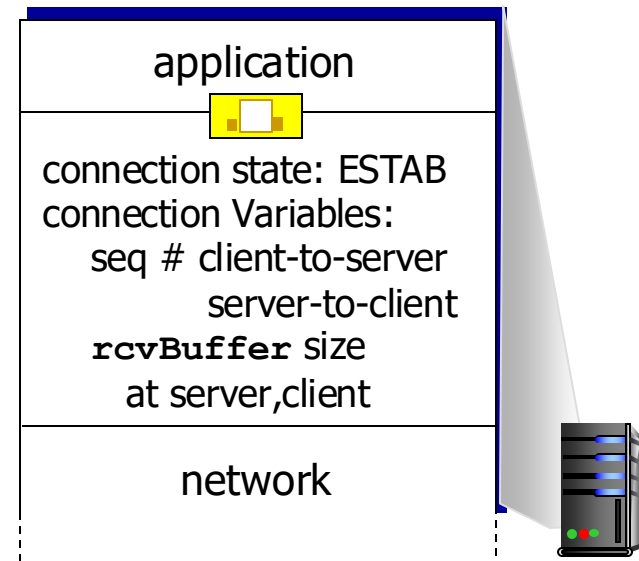
receive window

TCP segment format

# 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)

application

connection state: ESTAB
connection variables:
   seq # client-to-server
      server-to-client
  **rcvBuffer** size
   at server,client

network

application

connection state: ESTAB
connection Variables:
   seq # client-to-server
      server-to-client
  **rcvBuffer** size
   at server,client

network

```
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))
```

SYNSENT

choose init seq num, x
send TCP SYN msg

SYNbit=1, Seq=x

ESTAB

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

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

ACKbit=1, ACKnum=y+1

## Server state

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

LISTEN

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

SYN RCVD

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