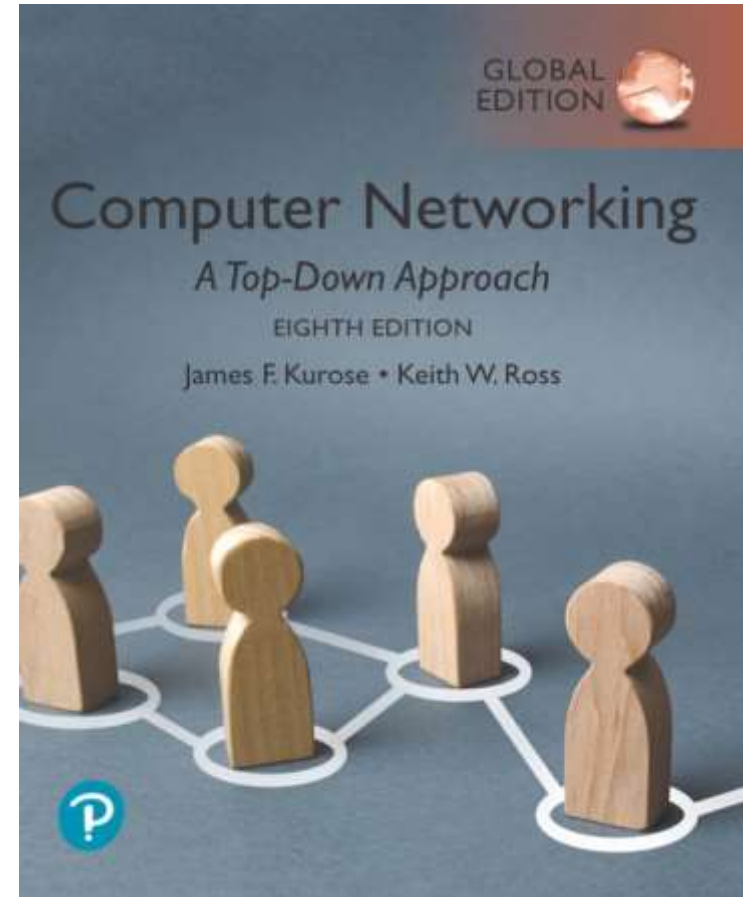# Chapter 3
# Transport Layer

Ren Ping Liu

renping.liu@uts.edu.au

adapted from textbook slides by JFK/KWR

25 March 2024

*Computer Networking: A Top-Down Approach*

8th Edition, Global Edition
Jim Kurose, Keith Ross
Copyright © 2022 Pearson Education Ltd

# Transport layer: roadmap

## 3.5 Connection-oriented transport: TCP

- segment structure
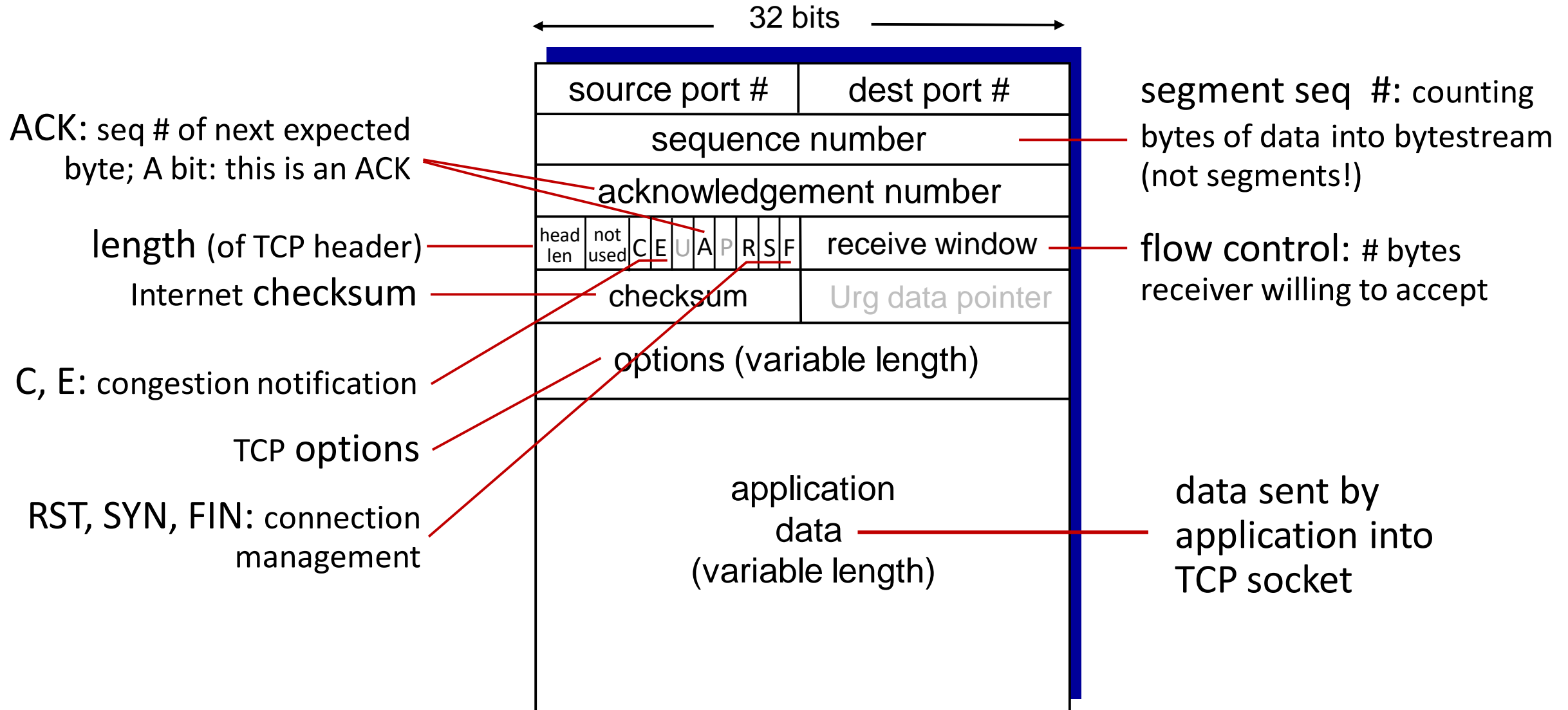- reliable data transfer
- flow control
- connection management

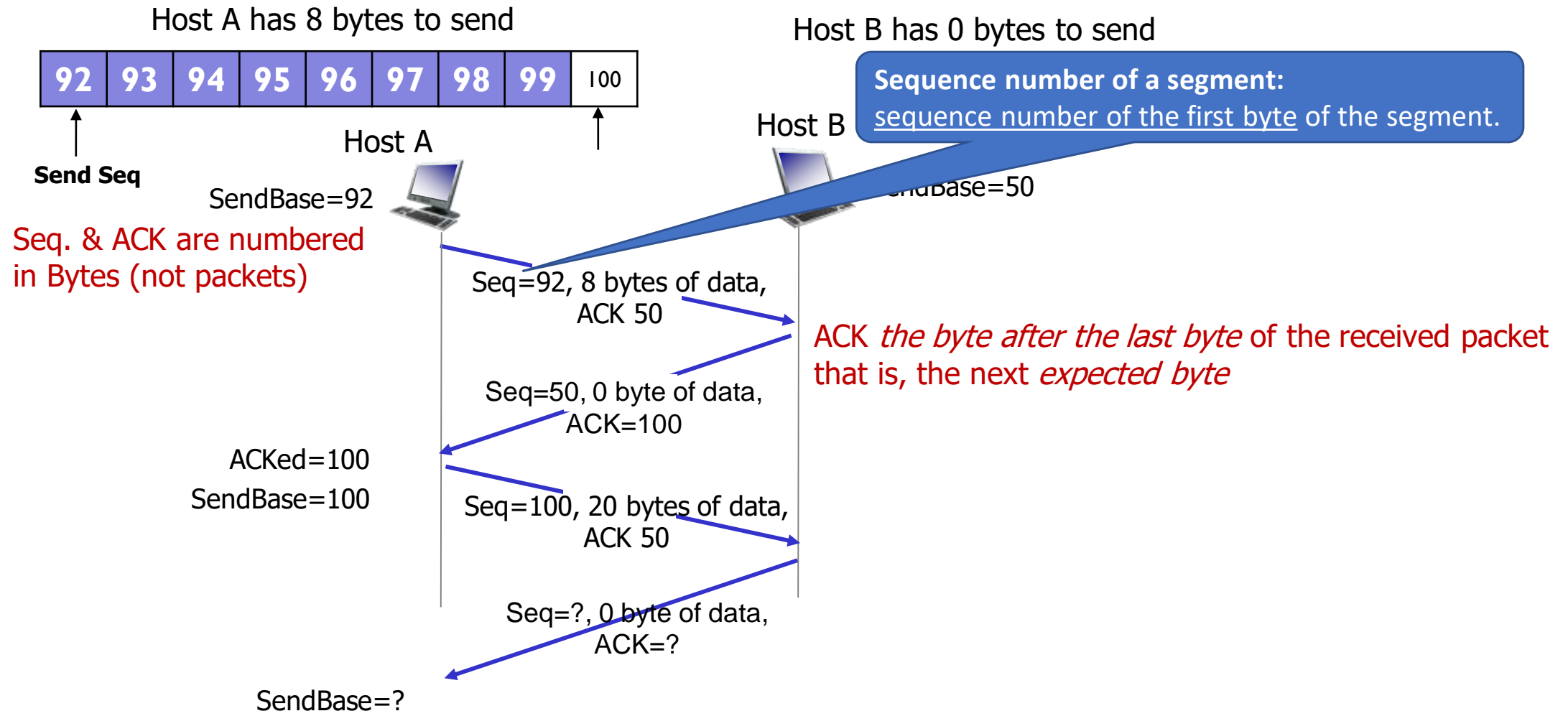# TCP: overview  RFCs: 793,1122, 2018, 5681, 7323

- 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
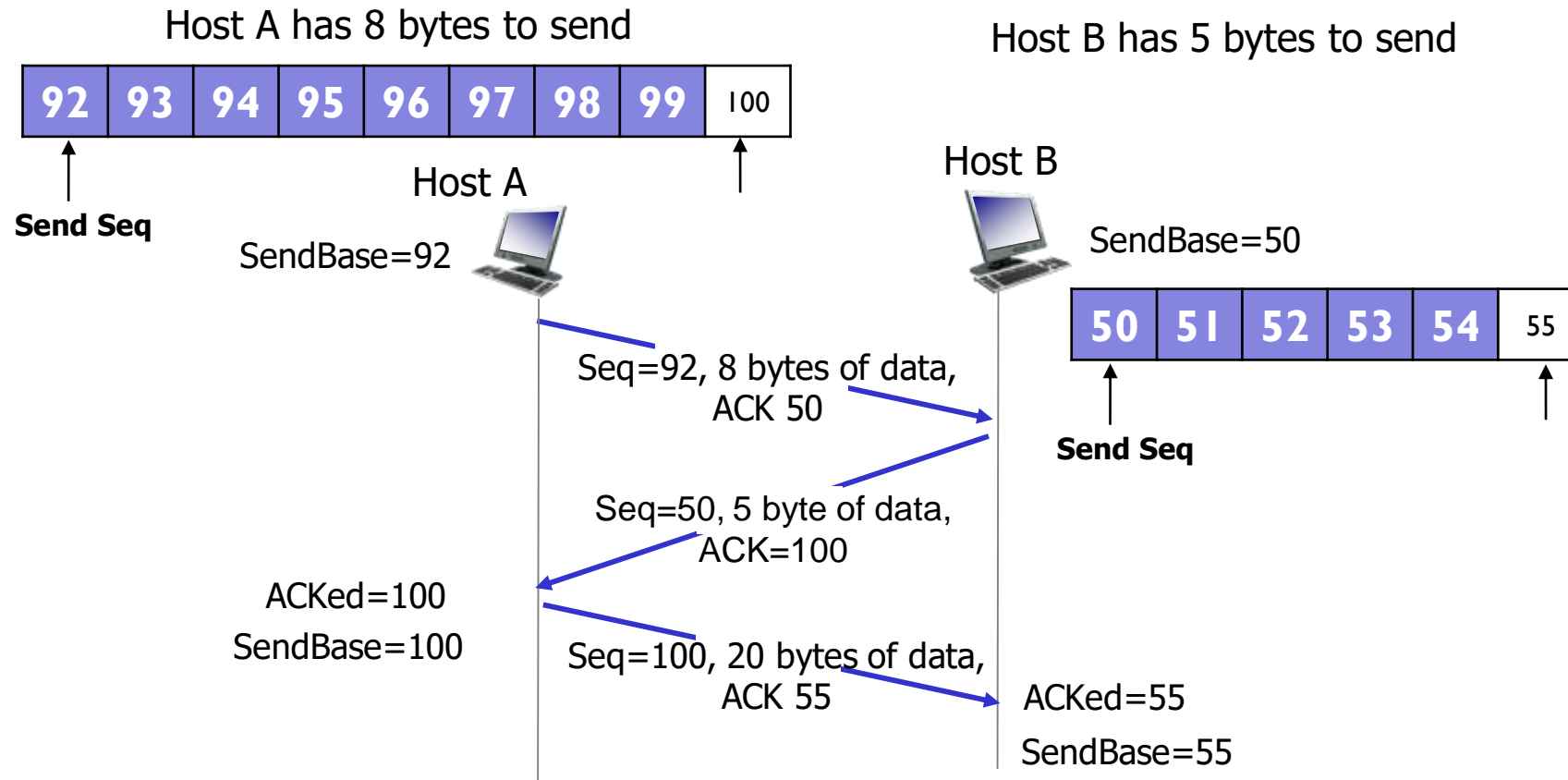- flow controlled:
  - sender will not overwhelm receiver

# TCP 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)

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

**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

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

data sent by application into TCP socket

# TCP: Seq. ACKs - one way: A➔B

Host A has 8 bytes to send

| 92 | 93 | 94 | 95 | 96 | 97 | 98 | 99 | 100 |

**Send Seq**

Host A

SendBase=92

Seq. & ACK are numbered
in Bytes (not packets)

Host B has 0 bytes to send

**Sequence number of a segment:**
<u>sequence number of the first byte</u> of the segment.

Host B

SendBase=50

Seq=92, 8 bytes of data,
ACK 50

ACK *the byte after the last byte* of the received packet
that is, the next *expected byte*

Seq=50, 0 byte of data,
ACK=100

ACKed=100

SendBase=100

Seq=100, 20 bytes of data,
ACK 50

Seq=?, 0 byte of data,
ACK=?

SendBase=?

# TCP: Seq. ACK - bidirectional

Host A has 8 bytes to send

| 92 | 93 | 94 | 95 | 96 | 97 | 98 | 99 | 100 |

**Send Seq**

Host A

SendBase=92

Host B has 5 bytes to send

Host B

SendBase=50

| 50 | 51 | 52 | 53 | 54 | 55 |

**Send Seq**

Seq=92, 8 bytes of data,
ACK 50

Seq=50, 5 byte of data,
ACK=100

ACKed=100
SendBase=100

Seq=100, 20 bytes of data,
ACK 55

ACKed=55
SendBase=55

# TCP Seq. ACKs - socket prog example

Client send 'test' (5Bytes) to Server, Server echo 'TEST' back

Client             Server

SendBase=42             SendBase=79

Client
types 'test'

Seq= ? , ACK= ? , data = 'test'

Server ACKs
receipt of 'test' ;
Echoes back  'TEST'

Seq= ? , ACK= ? , data = 'TEST'

Client ACKs receipt
of echoed 'TEST'

Seq= ? , ACK= ?

simple TCP client / server Python Socket Programming

# Transport layer: roadmap

# TCP round trip time, timeout

*Q:* how to set TCP timeout value?

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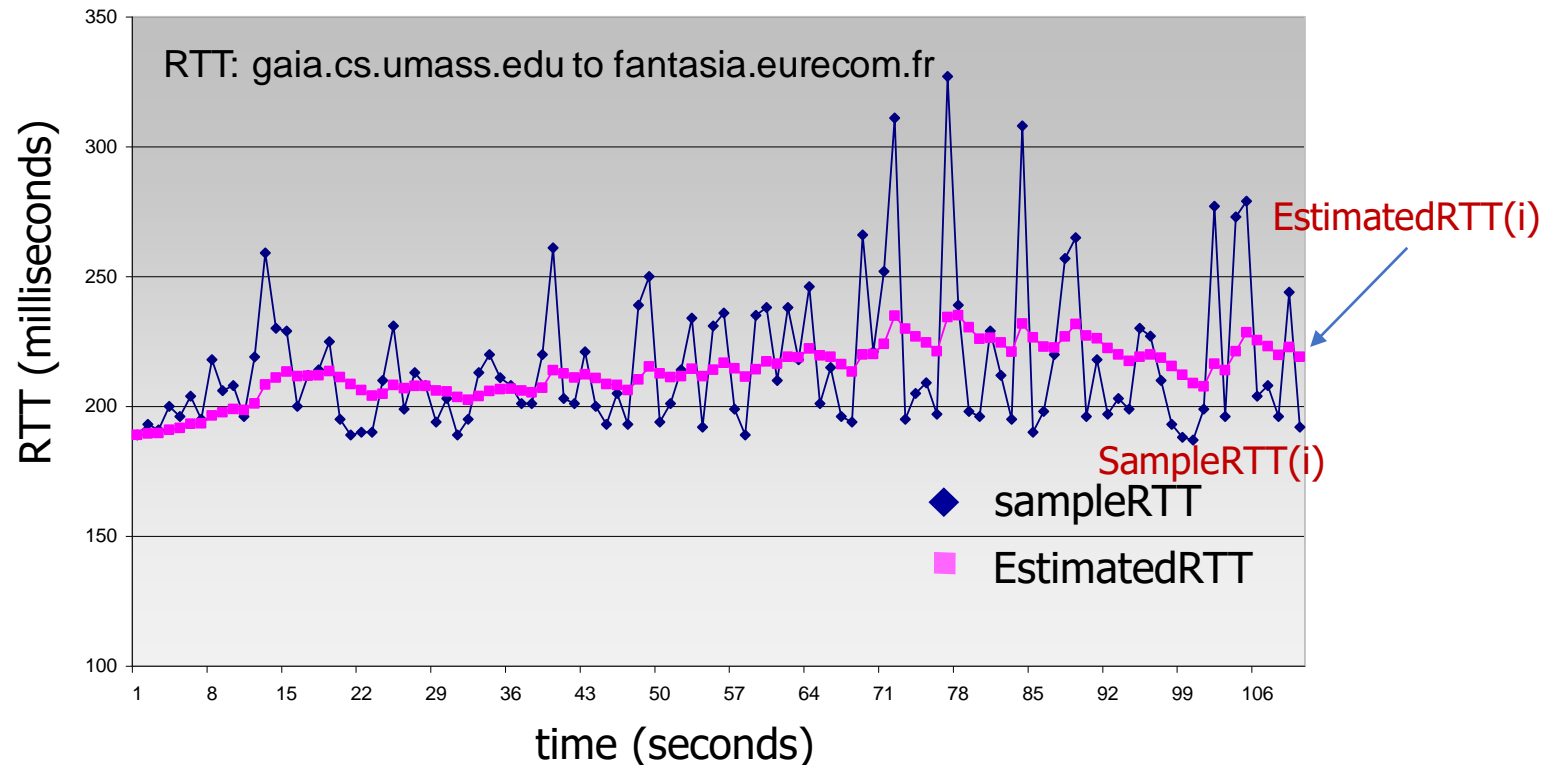- *too long:* slow reaction to segment loss

*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

**EstimatedRTT = (1- α)\*EstimatedRTT + α\*SampleRTT**
*(i)*                                            *(i-1)*                  *(i)*

- 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

# 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
  - expiration interval: **TimeOutInterval**

- transmit the segment

*event: ACK received*

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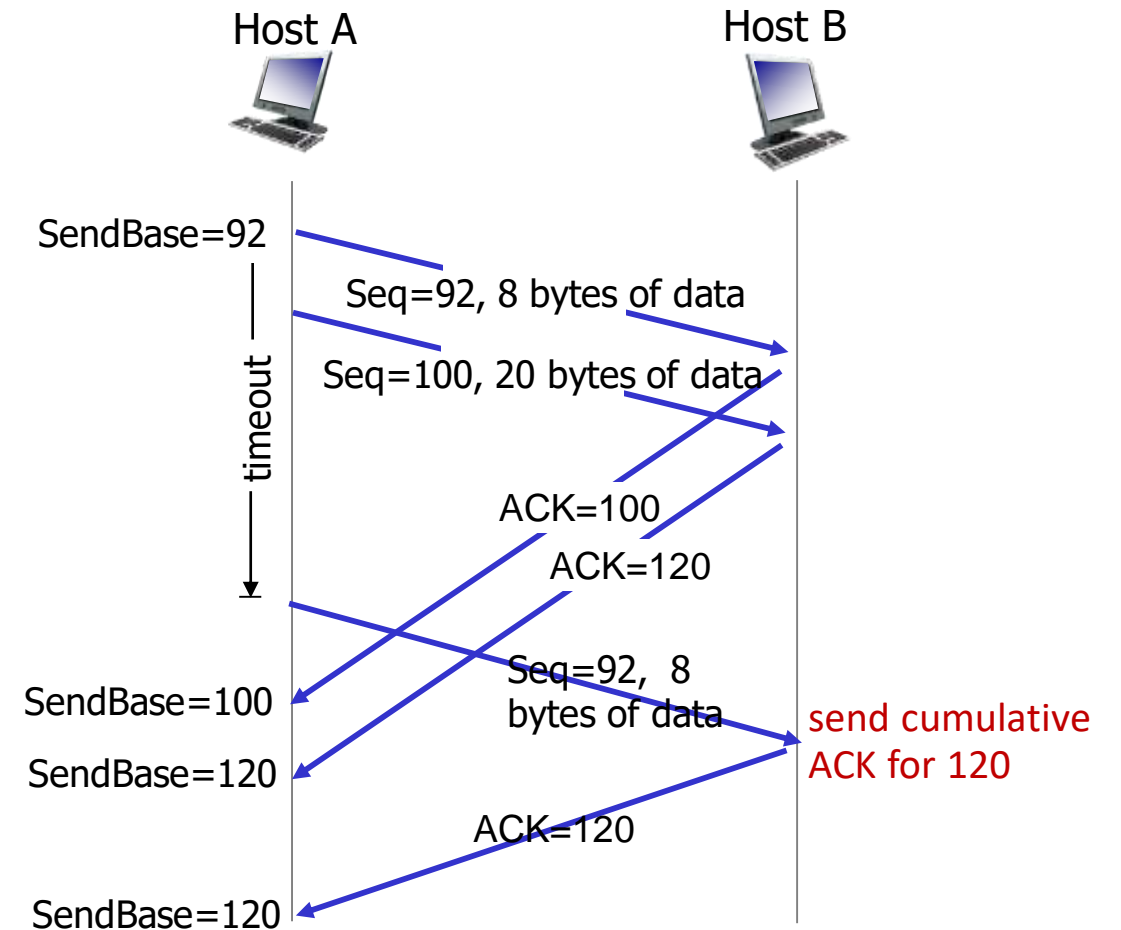  - start timer if there are  still unACKed segments

*event: timeout*

- retransmit segment that caused timeout
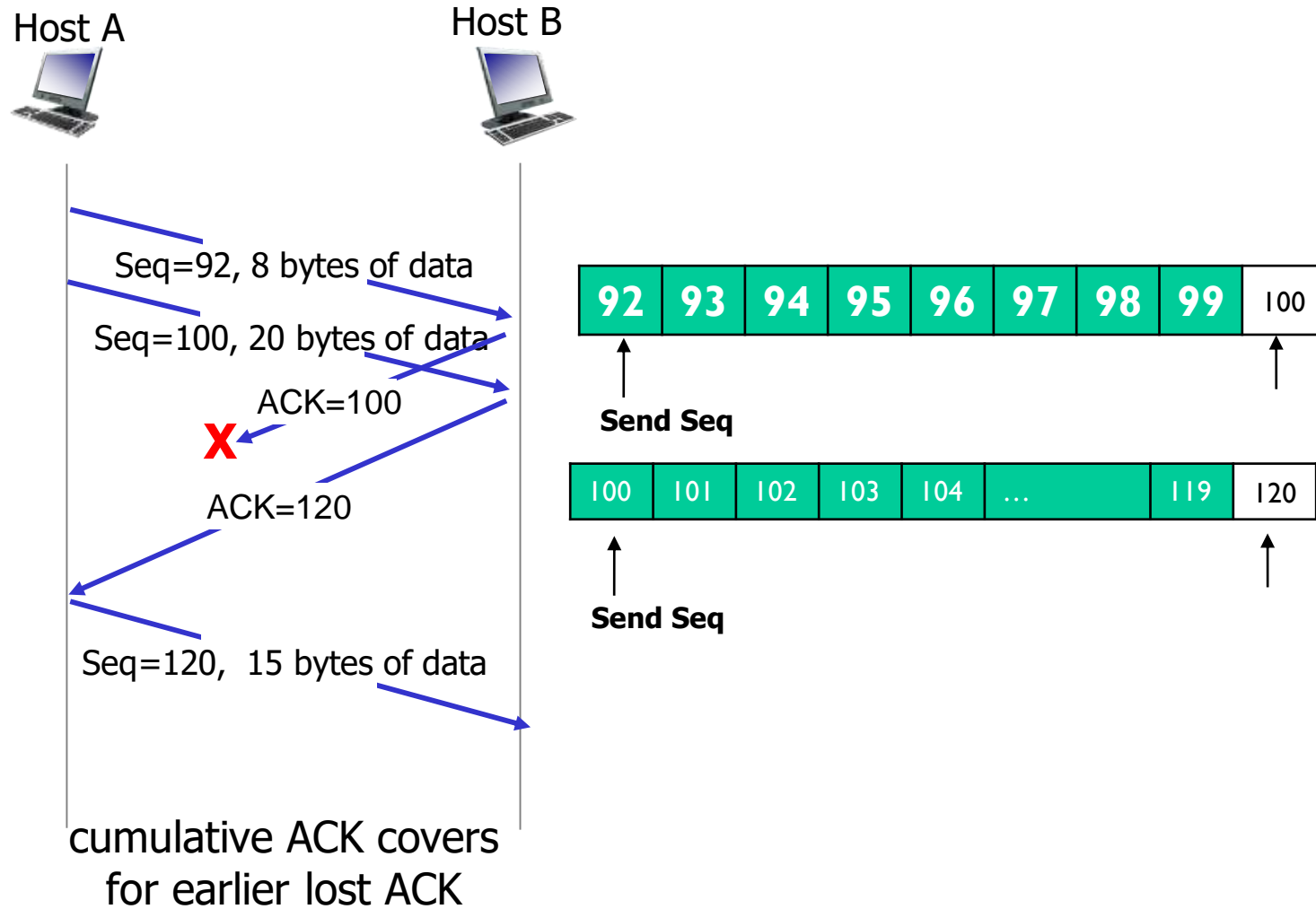- restart timer

# TCP retransmission - Timeout

| 92 | 93 | 94 | 95 | 96 | 97 | 98 | 99 | 100 |
|----|----|----|----|----|----|----|----|-----|

**Send Seq**

Host A           Host B

timeout

Seq=92, 8 bytes of data

ACK=100

X

Seq=92, 8 bytes of data

ACK=100

lost ACK scenario

Host A           Host B

SendBase=92

timeout

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

ACK=100

ACK=120

SendBase=100

SendBase=120

Seq=92, 8 bytes of data

send cumulative ACK for 120

ACK=120

SendBase=120

premature timeout

# TCP retransmission - Cumulative ACK



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

| 92 | 93 | 94 | 95 | 96 | 97 | 98 | 99 | 100 |

Send Seq

| 100 | 101 | 102 | 103 | 104 | ... | 119 | 120 |

Send Seq

# TCP retransmission - 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

Seq=120

Seq=140

Seq=160

Seq=180

# Transport layer: roadmap

# TCP flow control

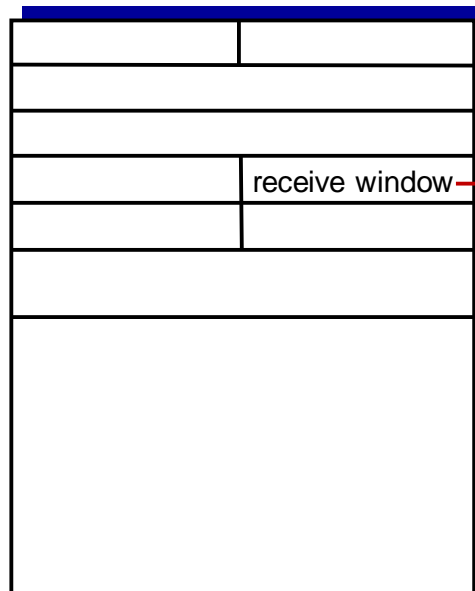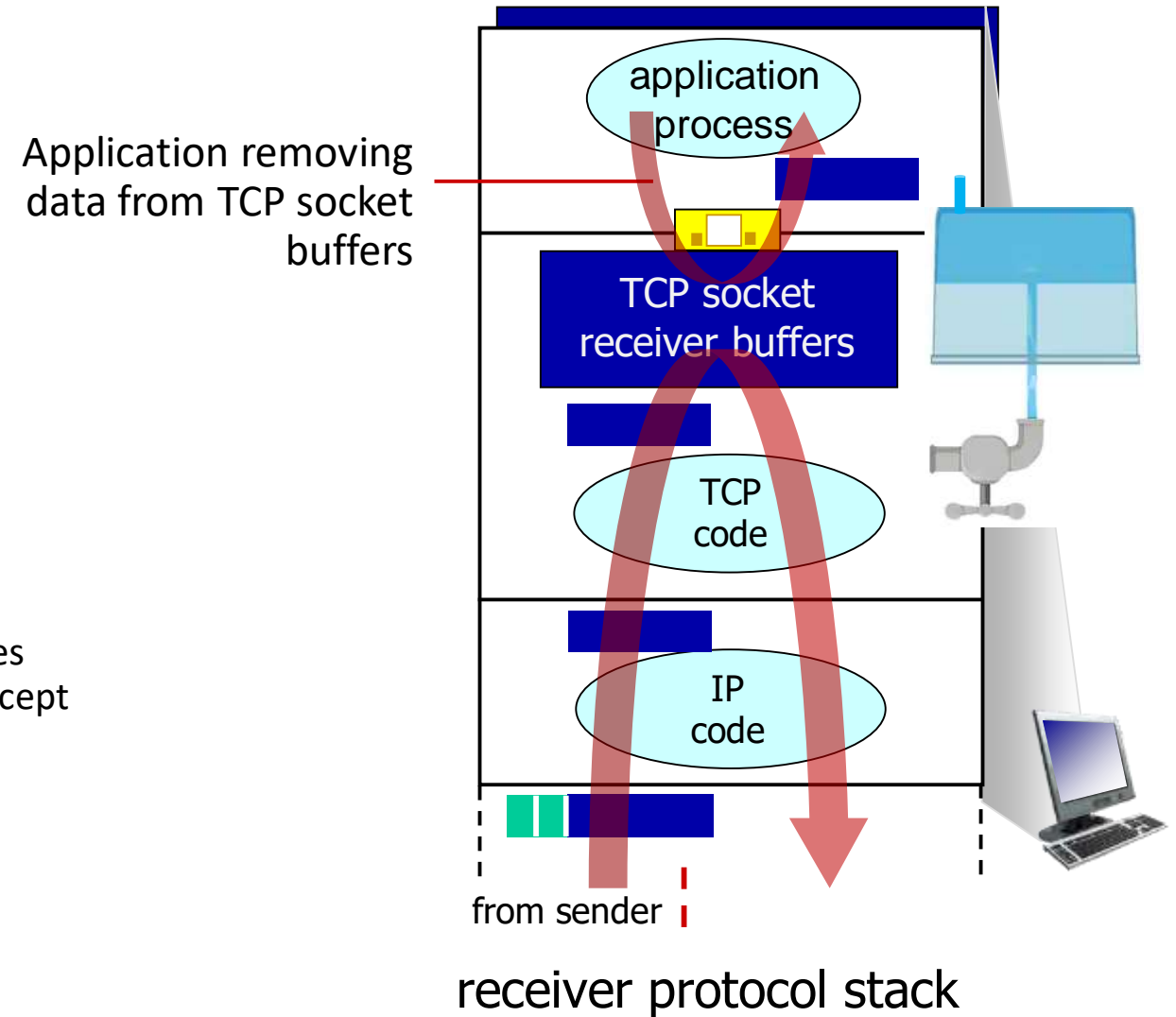Q: What happens if network layer delivers data faster than application layer removes data from socket buffers?

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
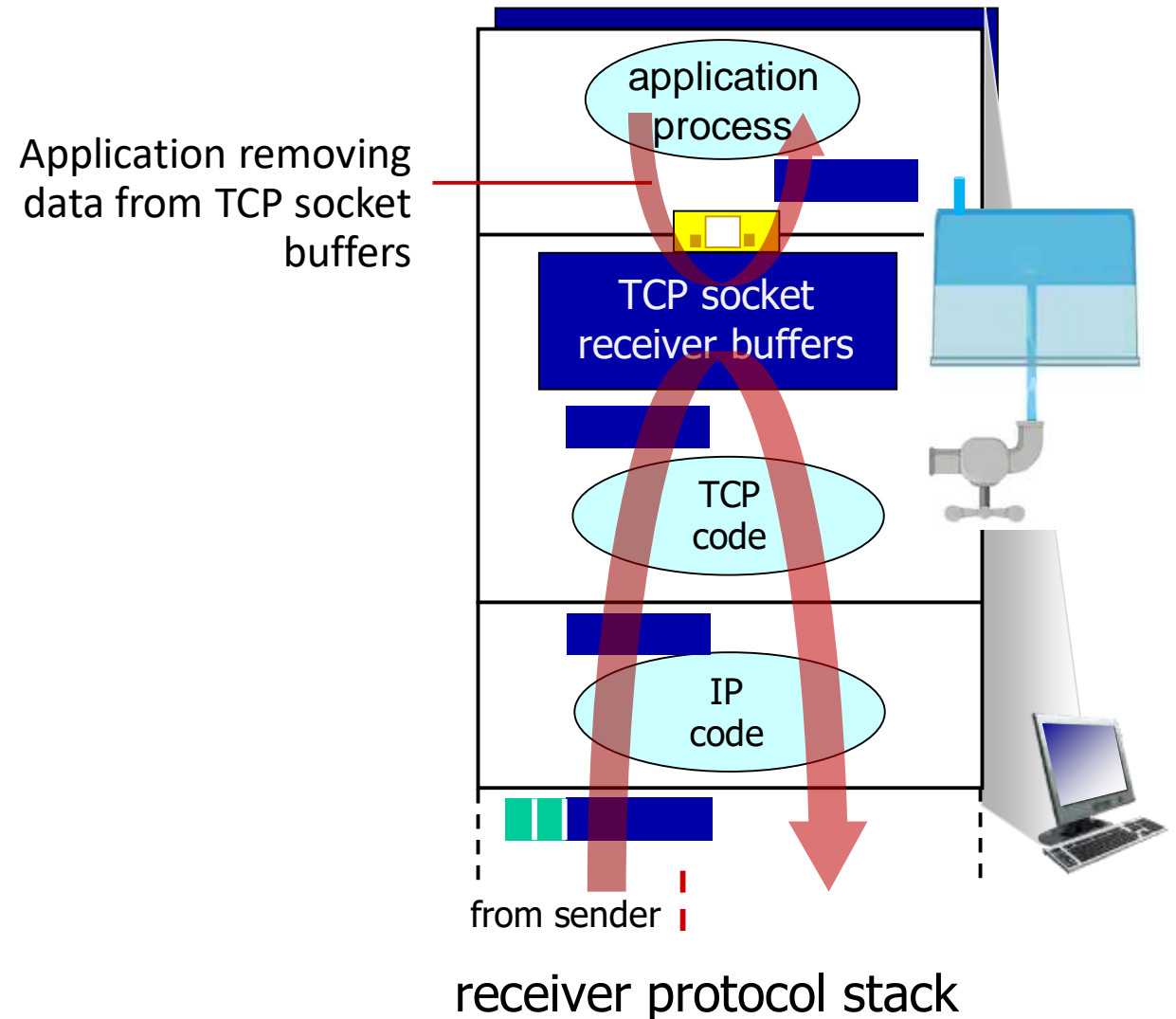
IP code

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?

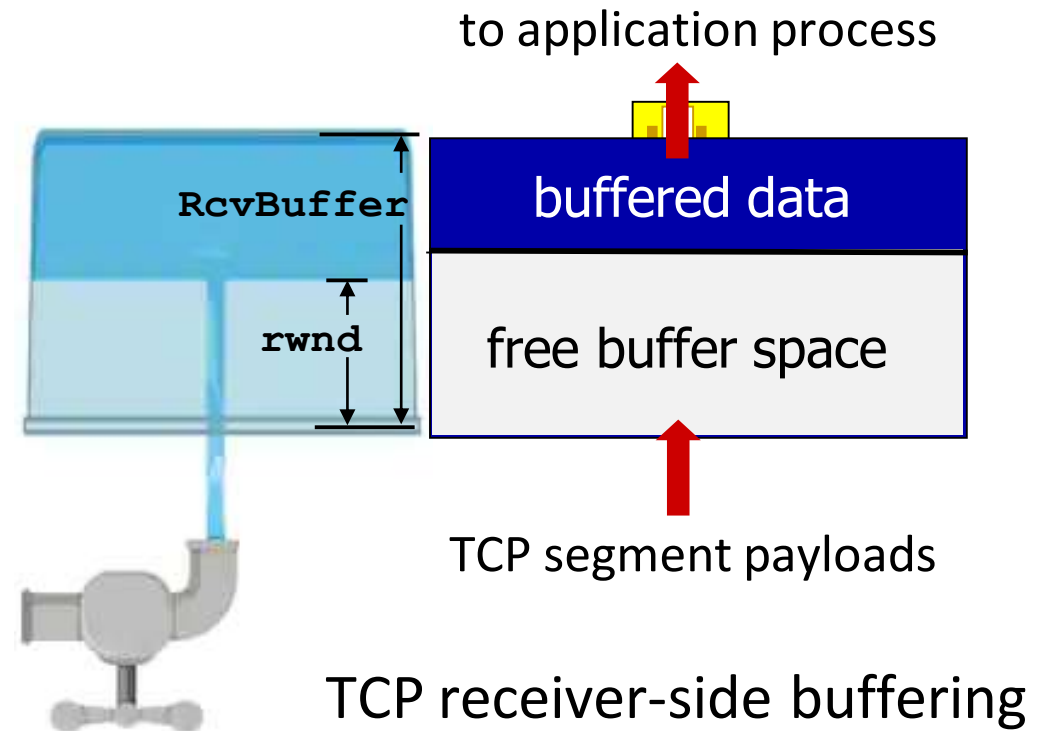Application removing data from TCP socket buffers

TCP socket receiver buffers

TCP code

Network layer delivering IP datagram payload into TCP socket buffers

IP code

application process

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

receive window ——— flow control: # bytes receiver willing to accept

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

buffered data

**RcvBuffer**

**rwnd**

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 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

incoming segment from receiver

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

# TCP window flow control



fast retransmit after sender
receipt of triple duplicate ACK

# TCP window flow control

SendBase=0000
MSS=1024
Window=5120

Seq=0000, segment=0-1023

Seq=1024

Seq=2048

Seq=3072

Seq=4096

ACKed=1024

ACKed=2048

Seq=5120

Seq=6144

Seq=?

Seq=?

fast retransmit after sender
receipt of triple duplicate ACK

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") | not usable |
|---|---|---|

incoming segment to sender

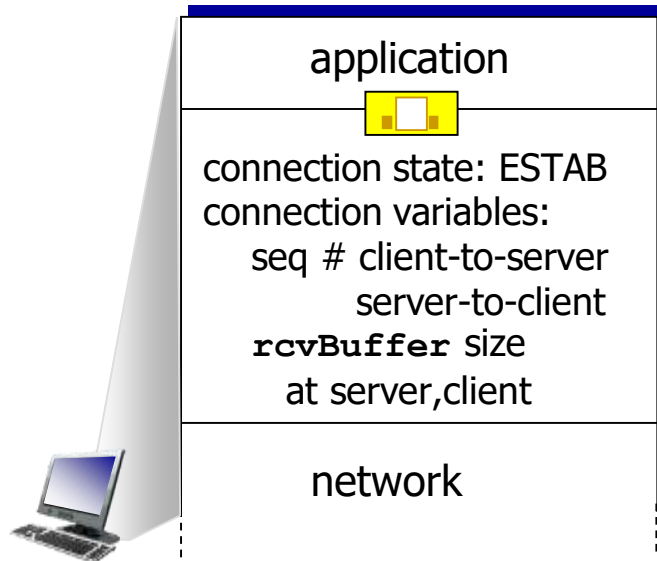| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |
| A | rwnd |
| checksum | urg pointer |

# Mid-break
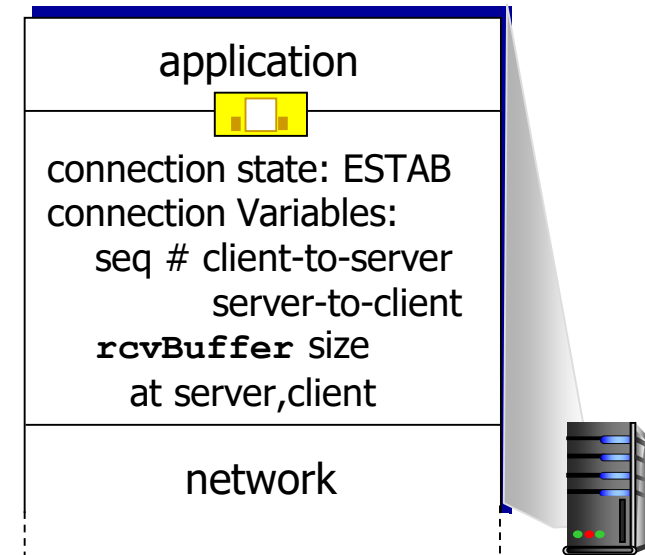
- Q & A

# Transport layer: roadmap

# TCP connection management

before exchanging data, sender/receiver "handshake":

- agree to establish connection (each knowing the other willing to establish connection)
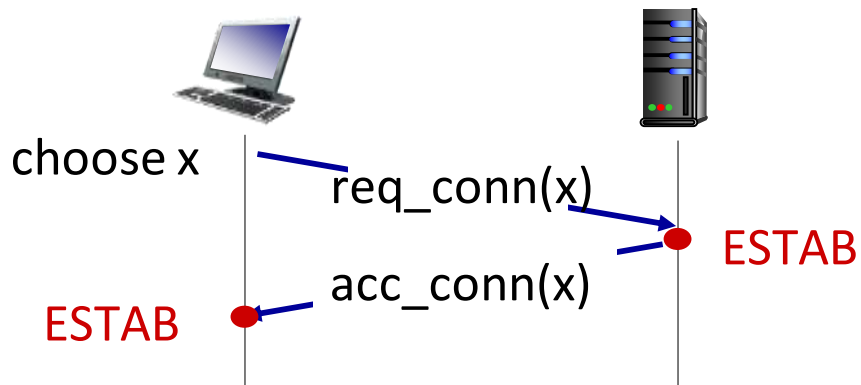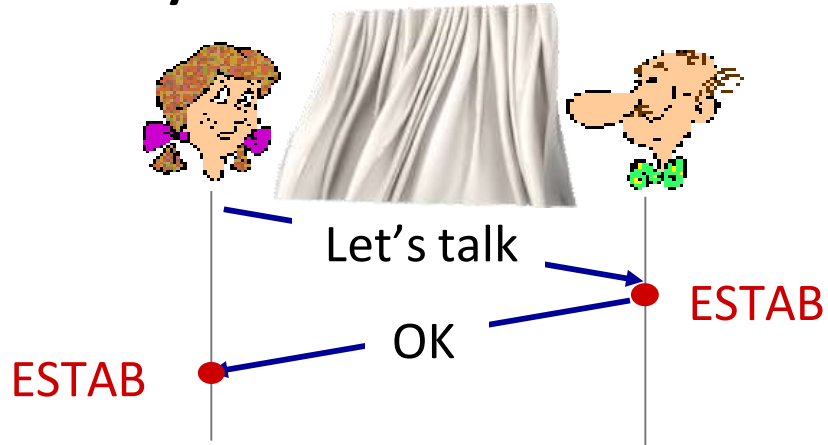- agree on connection parameters (e.g., starting seq #s)

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

```
clientSocket.connect("server","port");
...
clientSocket.send(…)
```

```
connSocket,addr = serverSocket.accept();
...
connSocket.recv(…)
```
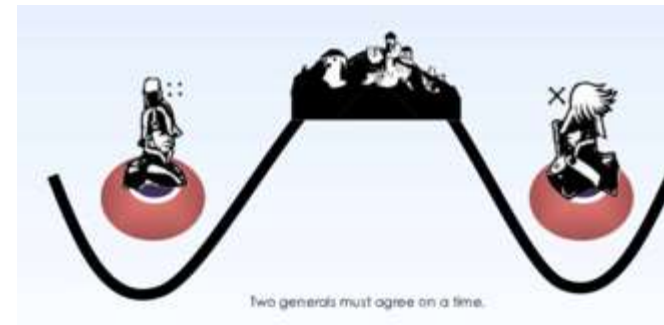
# Agreeing to establish a connection

## 2-way handshake:



Let's talk → ESTAB

ESTAB ← OK

choose x → req_conn(x) → ESTAB
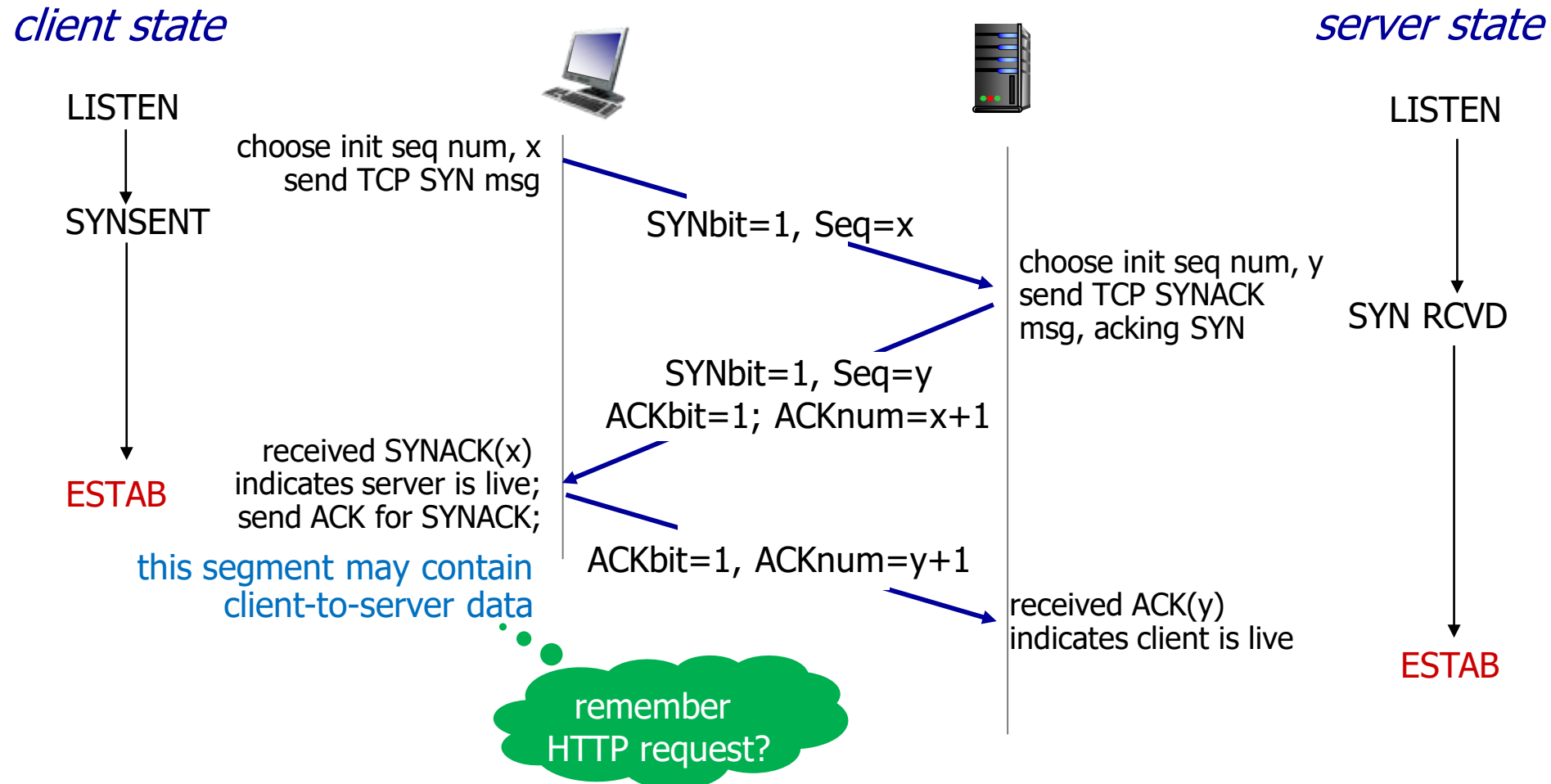
ESTAB ← acc_conn(x)

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

- variable delays
- unreliable channel,
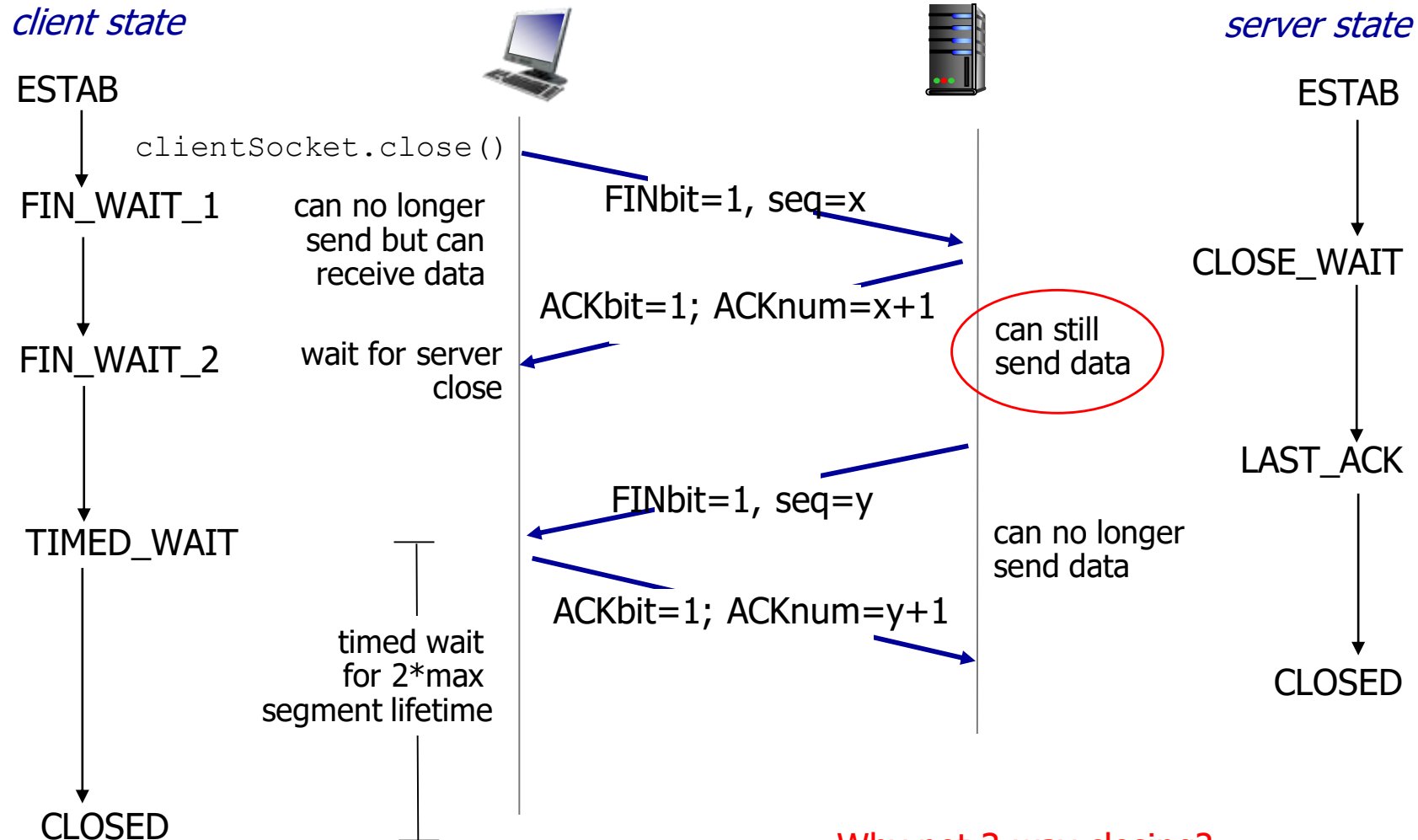- retransmitted messages due to message delay, loss

**Byzantine Generals Problem**



Two generals must agree on a time.

# TCP 3-way handshake

*client state*　　　　　　　　　　　　　　　　　　　　　　　　　*server state*

LISTEN　　　　　　　　　　　　　　　　　　　　　　　　　　　　　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

received SYNACK(x)
indicates server is live;
send ACK for SYNACK;

ESTAB

this segment may contain
client-to-server data

ACKbit=1, ACKnum=y+1

received ACK(y)
indicates client is live

ESTAB

remember
HTTP request?
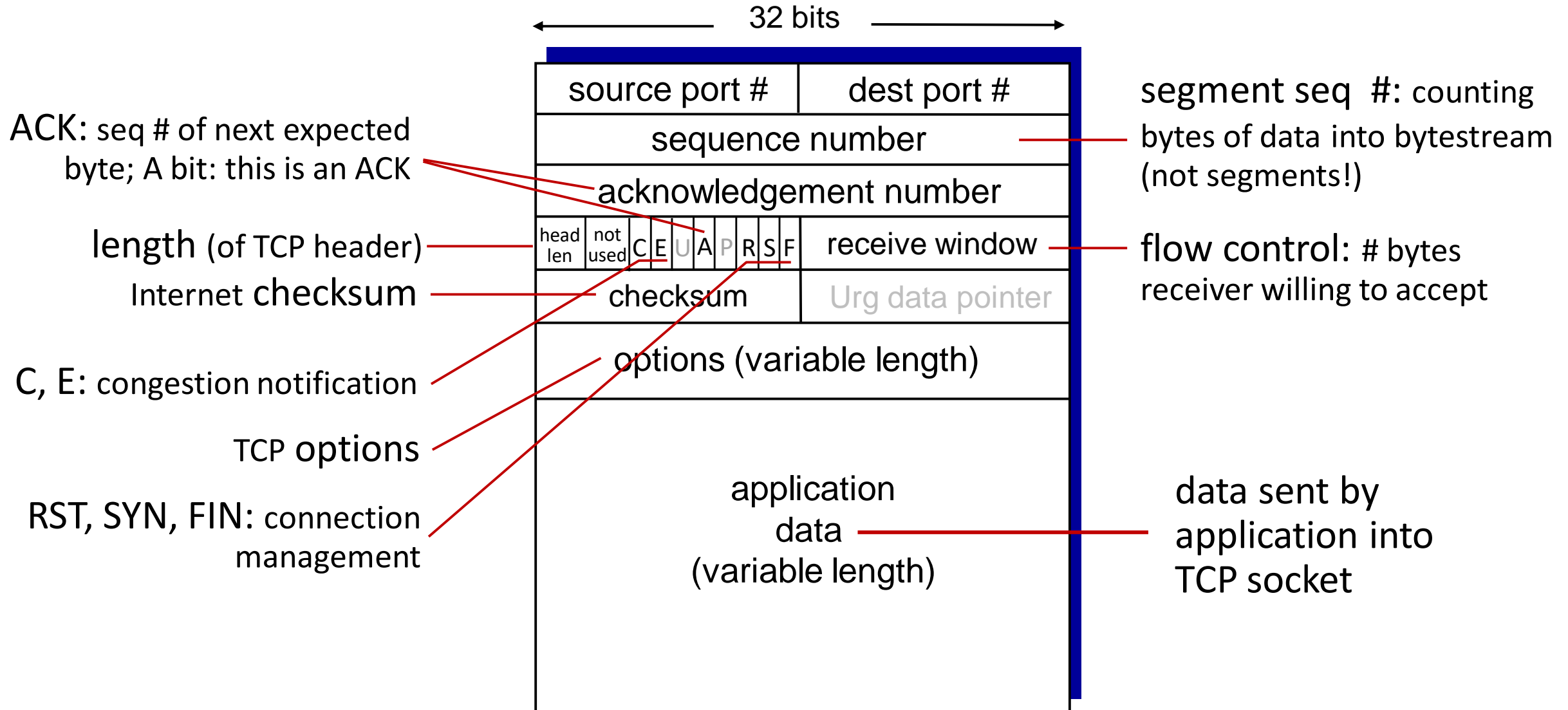
# Closing a TCP connection

- 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

# Closing a TCP connection

# TCP segment structure - review

32 bits

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

| 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) |
|---|

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

# Principles of congestion control

## Congestion:

- informally: "too many sources sending too much data too fast for *network* to handle"

- manifestations:
  - long delays (queueing in router buffers)
  - packet loss (buffer overflow at routers)

- different from flow control (one-to-one)!

- a top-10 problem!
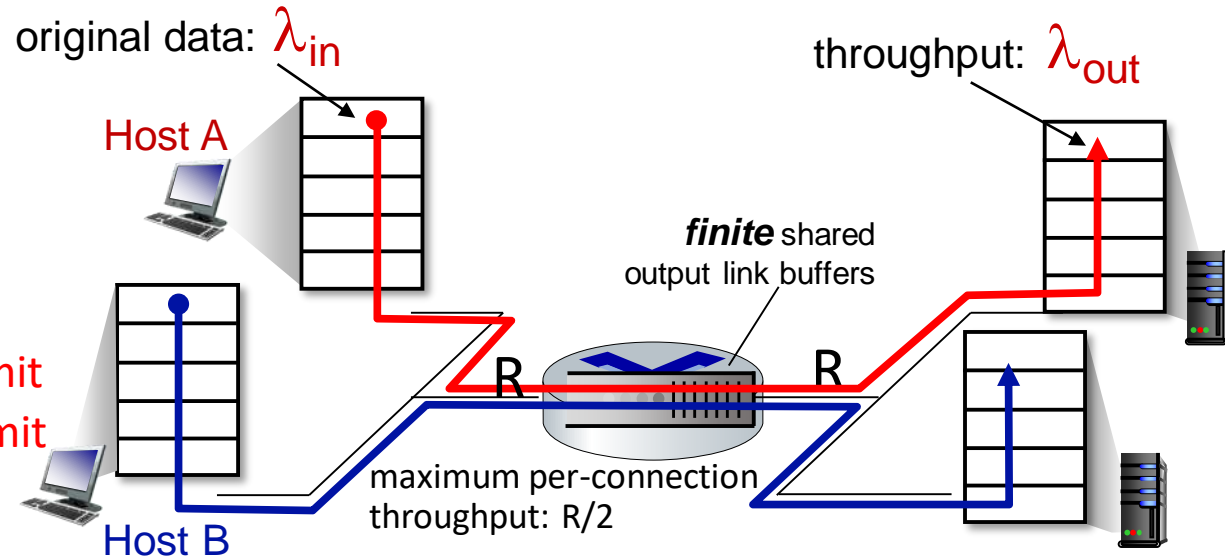
congestion control:
too many senders,
sending too fast

flow control: one sender
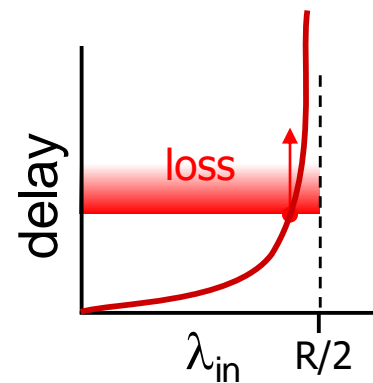too fast for one receiver

# Causes/costs of congestion: scenario 1

- two senser/receiver pairs
- one router, two flows,
- link capacity: R
- As traffic increases → R/2
  - ✗ Delay increase - timeout → retransmit
  - ✗ Packet loss → retransmit

original data: $\lambda_{in}$

Host A

throughput: $\lambda_{out}$

*finite* shared output link buffers

R   R

maximum per-connection throughput: R/2

Host B

- **How can we?**
  - ✓ Avoid congestion
  - ✓ Resolve congestion

delay    loss

$\lambda_{in}$   R/2

✗ Delay increase
✗ Packet loss

R/2

throughput: $\lambda_{out}$

"wasted" capacity due to un-needed retransmissions

$\lambda'_{in}$   R/2
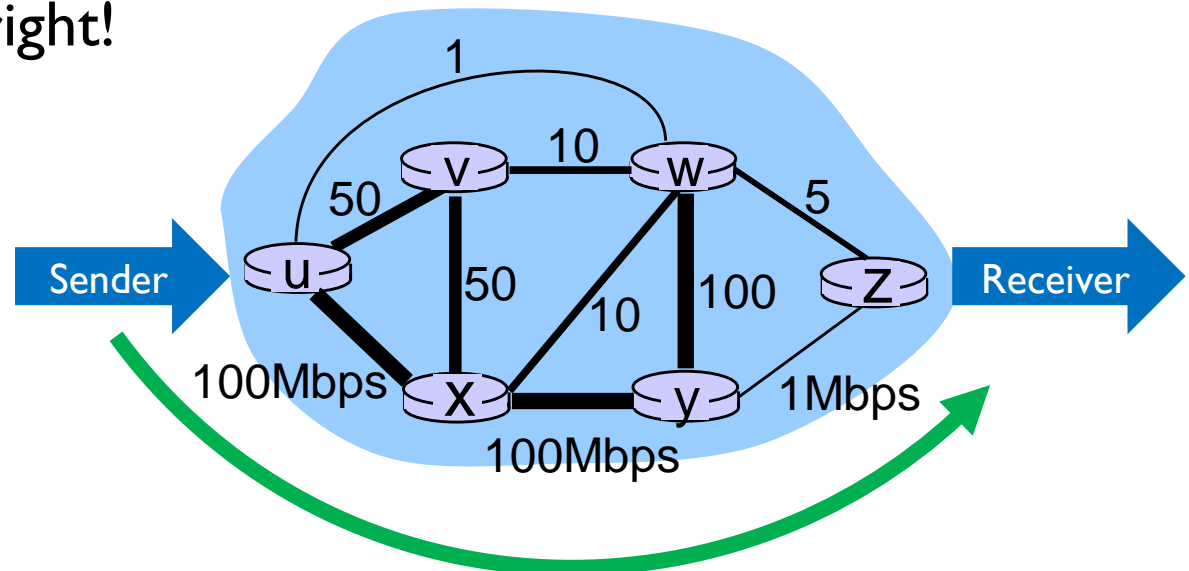
# Congestion: multi-hop scenario

- Sender ➜ Receiver
  - Download 100MB file, using FTP/TCP
  - Via multiple hops of various BW (unknown!)
- Question:
  - Decide sending data rate? (wnd/RTT)
  - not too low, not too high – just right!
- Strategy:
  - Start from low rate
  - Slowly ramp up
  - Stabilise at allowed

# Transport layer: roadmap

# 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

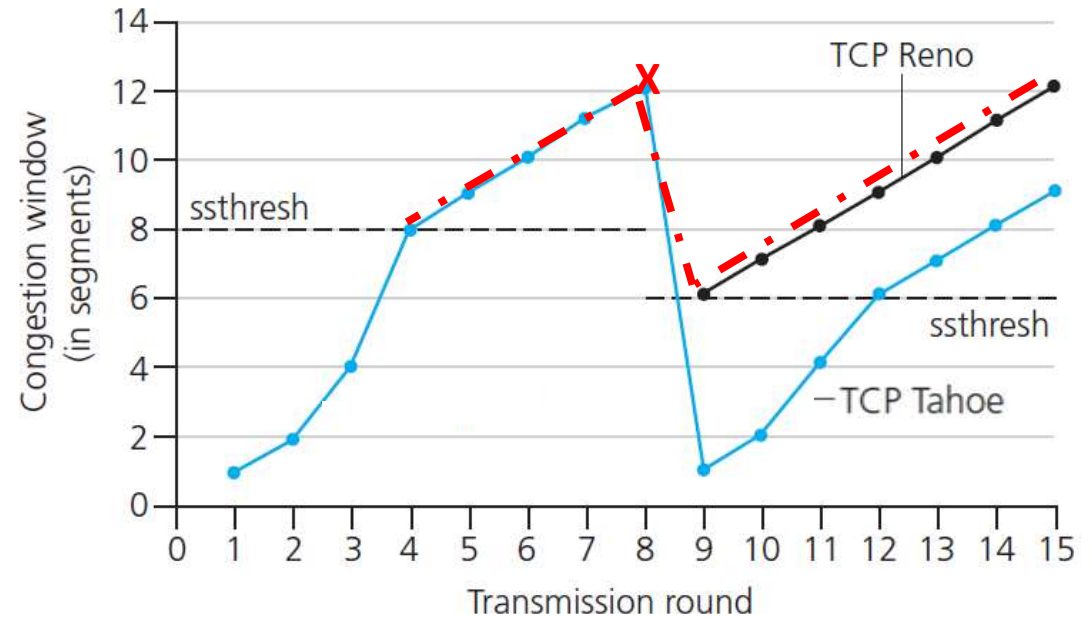- *summary:* initial rate is slow, but ramps up exponentially fast

# TCP: from slow start to congestion avoidance

*Q:* when should the exponential increase switch to linear?

*A:* when **cwnd** ≥ **ssthresh**

## Implementation:

- variable **ssthresh**

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



* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/
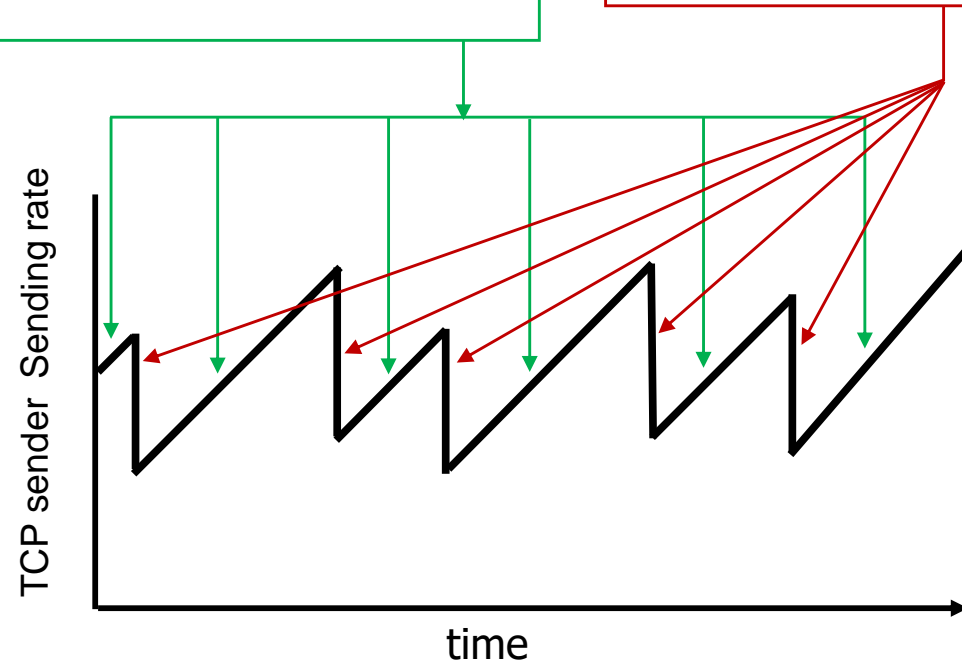
# TCP congestion control: AIMD

■ *approach:* senders can increase sending rate until packet loss (congestion) occurs, then decrease sending rate on loss event

*Additive Increase*

increase sending rate by 1 maximum segment size every RTT until loss detected

*Multiplicative Decrease*

cut sending rate in half at each loss event



TCP sender  Sending rate

time

**AIMD** sawtooth behavior: *probing* for bandwidth
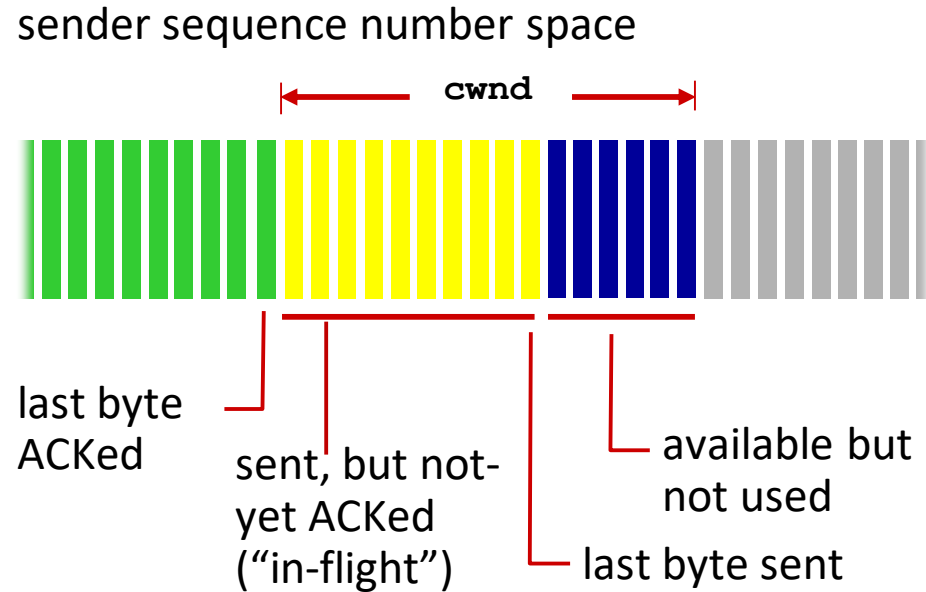
# TCP AIMD: more

*Multiplicative decrease* detail:  sending rate is

- Cut in half on loss detected by triple duplicate ACK (TCP Reno)
- Cut to 1 MSS (maximum segment size) when loss detected by timeout (TCP Tahoe)

Why AIMD?

- AIMD – a distributed, asynchronous algorithm – has been shown to:
  - optimize congested flow rates network wide!
  - have desirable stability properties

# TCP congestion control: details

sender sequence number space

cwnd

last byte ACKed

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

last byte sent
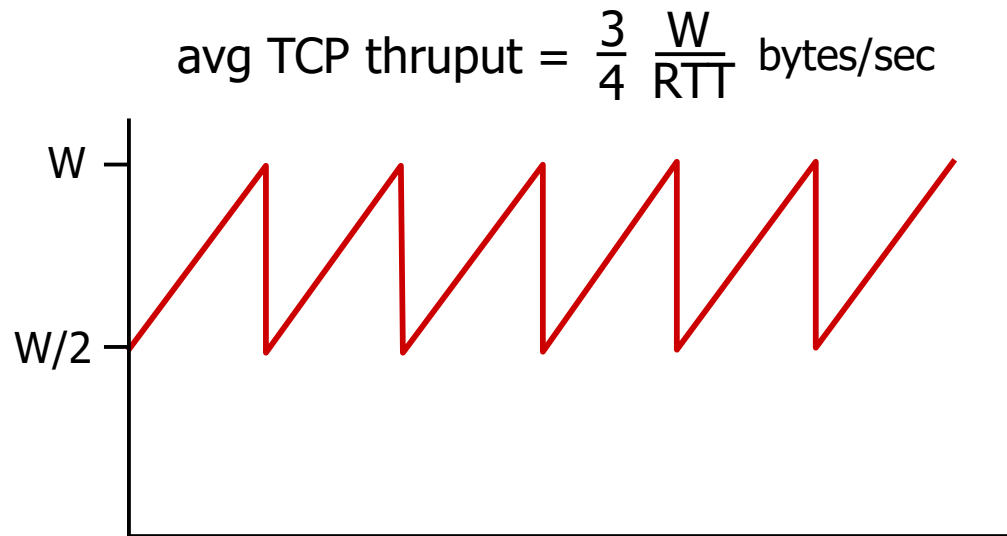
available but not used

TCP sending behavior:

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

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

- TCP sender limits transmission: `LastByteSent- LastByteAcked ` $\leq$ `cwnd`
- `cwnd` is dynamically adjusted in response to observed network congestion (implementing TCP congestion control)

# TCP throughput

- avg. TCP thruput as function of window size, RTT?
  - ignore slow start, assume there is always data to send
- W: window size (measured in bytes) where loss occurs
  - avg. window size (# in-flight bytes) is ¾ W
  - avg. thruput is 3/4W per RTT

avg TCP thruput = $\dfrac{3}{4}\dfrac{W}{RTT}$ bytes/sec

# Chapter 3: summary

- **principles behind transport layer services:**
  - multiplexing, demultiplexing
  - reliable data transfer
  - flow control
  - congestion control
- **instantiation, implementation in the Internet**
  - UDP
  - TCP

Up next:

- **leaving the network "edge"** (application, transport layers)
- **into the network "core"**
- **two network-layer chapters:**
  - data plane
  - control plane

# Assignment2 - 5%

- Assignment2 – 5%: Transport Layer
  - Based on weeks 5 & 6 learning materials
  - Answer in required format - **strictly!**
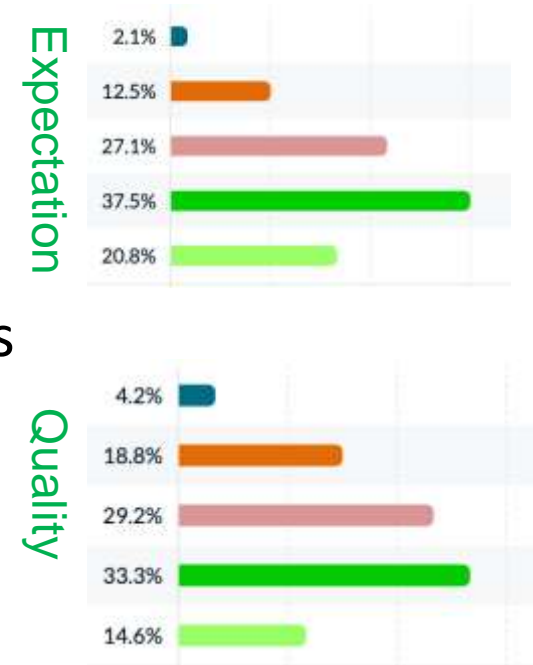  - Due by Sunday 31st March - <u>No late submission accepted</u>!

# Student Feedback Survey (SFS) Results

- **Some Positive**
  - Detailed explanation – that's the way!
  - Engaging Q&A during lecture – making the most!
  - Tutorial/Lab – generally positive, some not - passed to the tutors
- **Some Issues:**
  - Content heavy – reducing load, keeping fundamental
  - overwhelmed by IT concepts – recursive teaching in CS



Expectation

| | |
|---|---|
| 2.1% | |
| 12.5% | |
| 27.1% | |
| 37.5% | |
| 20.8% | |

Quality

| | |
|---|---|
| 4.2% | |
| 18.8% | |
| 29.2% | |
| 33.3% | |
| 14.6% | |

# Lecture done ✓

- Q & A