# Connection-oriented transport: TCP

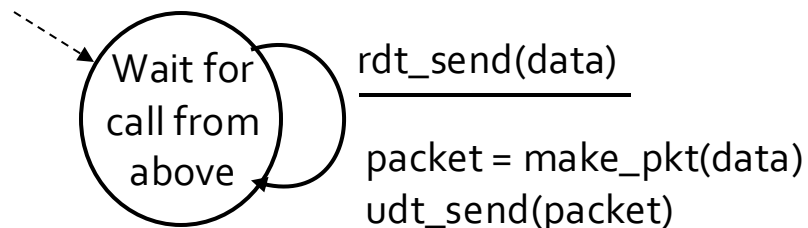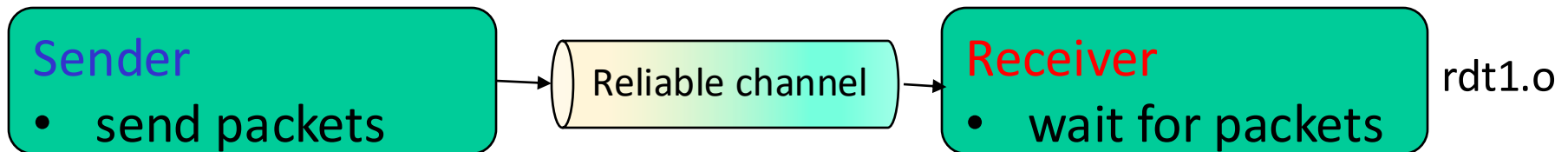## CE 352, Computer Networks

Salem Al-Agtash

Lecture 10

Slides are adapted from Computer Networking: A Top Down Approach, 7th Edition © J.F Kurose and K.W. Ross
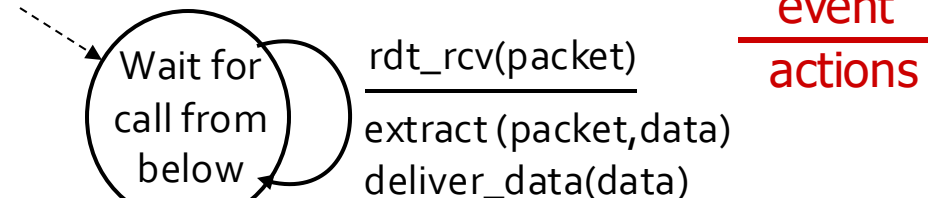
# Recap (transport layer)

- Transport layer provides communication between application processes (mux/demux using ports) so that reliable and efficient data delivery is achieved.

- Network layer can result in packets that are corrupted, delayed, dropped, reordered, or duplicated.
- Network layer gives no guidance on traffic volume to send and when

- TCP and UDP are the common transport protocols

- UDP is a minimalist lightweight communication between processes

- TCP offers a reliable, in-order, and byte stream communication with congestion control

# Recap (reliable transport channel)

- Reliable transport channel is easy in a perfect world
  - rdt1.0: Simple protocol
    - provides neither flow nor error control

| Sender | Reliable channel | Receiver | rdt1.0 |
|---|---|---|---|
| • send packets | | • wait for packets | |

Wait for call from above — rdt_send(data)

packet = make_pkt(data)
udt_send(packet)

sender

Wait for call from below — rdt_rcv(packet)
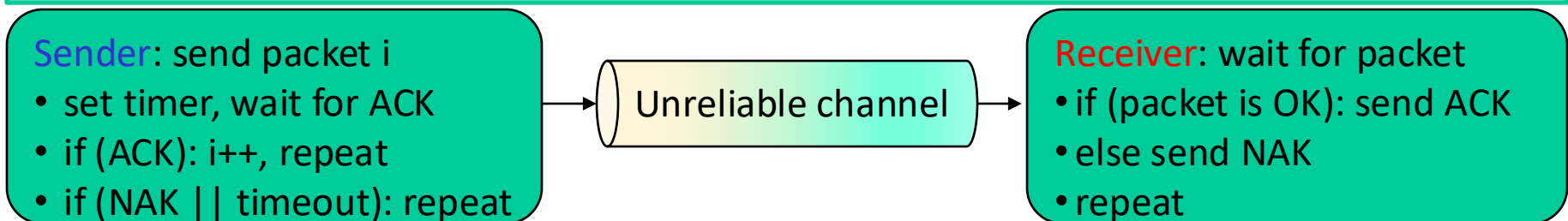
extract (packet,data)
deliver_data(data)

receiver

event
actions
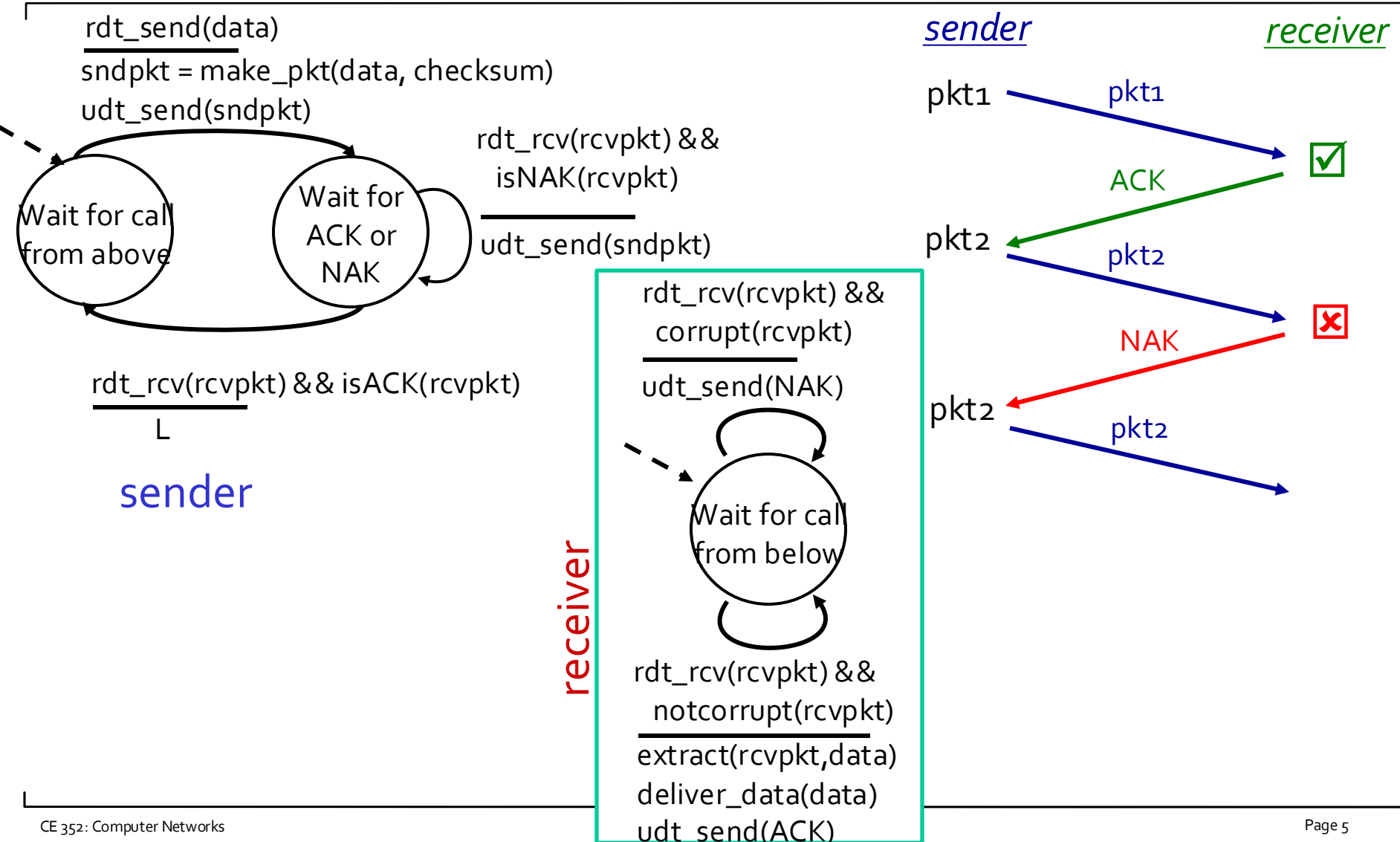
# Recap (unreliable transport channel)
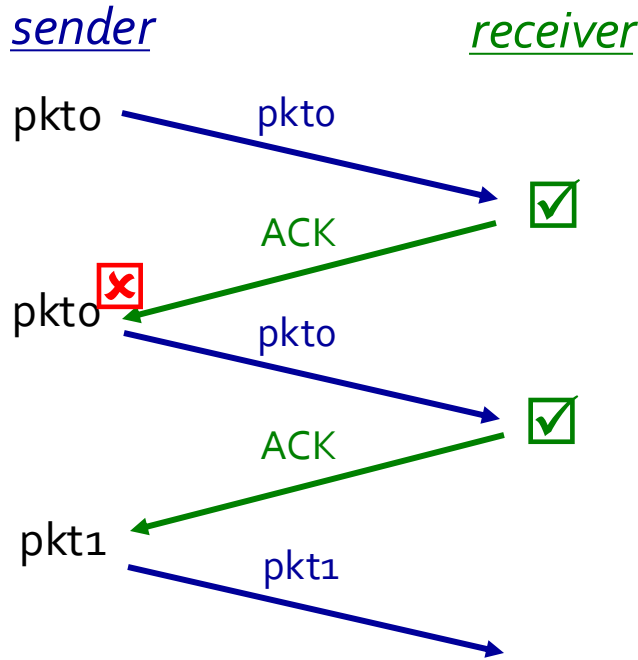
- Components of solution to send packets over unreliable channel:
    - checksums (to detect bit errors)
    - timers (to detect loss)
    - acknowledgements (Ack, NAK)
    - sequence numbers (to deal with duplicates)
- Stop-and-wait
    - rdt2.0: deals with packet corruption
    - rdt2.1: deals with garbled Ack/NAKs
    - rdt2.2: NAK-free by the use of sequence numbers
    - rdt3.0: deals with packet loss – timer driven loss detection
- Pipelined protocol (sliding window)
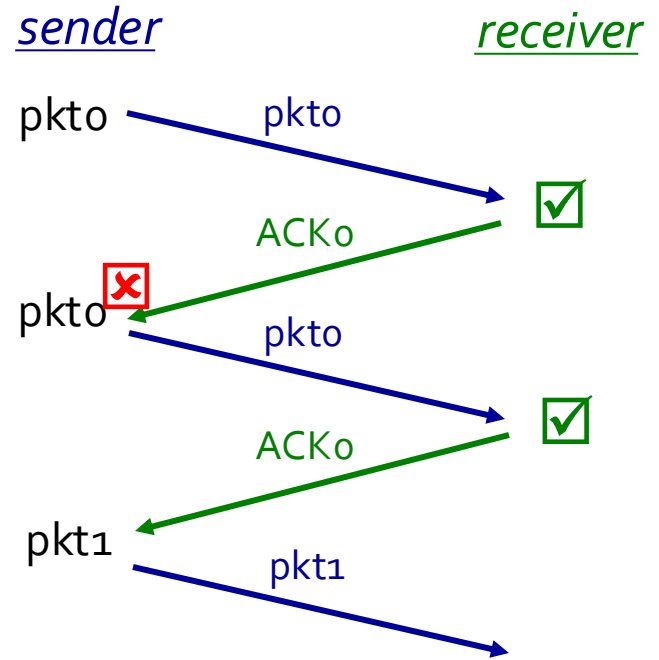    - Go-Back-N
    - Selective repeat

**Sender**: send packet i
- set timer, wait for ACK
- if (ACK): i++, repeat
- if (NAK || timeout): repeat

Unreliable channel

**Receiver**: wait for packet
- if (packet is OK): send ACK
- else send NAK
- repeat

# Recap (rdt2.0: deals with packet corruption)

rdt_send(data)

sndpkt = make_pkt(data, checksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) &&
isNAK(rcvpkt)

Wait for call from above    Wait for ACK or NAK

udt_send(sndpkt)

rdt_rcv(rcvpkt) && isACK(rcvpkt)

L

sender

rdt_rcv(rcvpkt) &&
corrupt(rcvpkt)

udt_send(NAK)

Wait for call from below

receiver

rdt_rcv(rcvpkt) &&
notcorrupt(rcvpkt)

extract(rcvpkt,data)
deliver_data(data)
udt_send(ACK)

*sender*                    *receiver*

pkt1 ——— pkt1 ☑

ACK

pkt2 ——— pkt2 ☒

NAK

pkt2 ——— pkt2

# Recap (rdt2.1/2.2: deals with garbled ACK-NAK/ data and ack packets carry sequence numbers)

sender    receiver        sender    receiver

pkt0 ———— pkt0 ————→ ☑         pkt0 ———— pkt0 ————→ ☑
         ←—— ACK ————                   ←—— ACK0 ————
pkt0 ✖                           pkt0 ✖
     ———— pkt0 ————→ ☑               ———— pkt0 ————→ ☑
     ←—— ACK ————                        ←—— ACK0 ————
pkt1 ———— pkt1 ————→             pkt1 ———— pkt1 ————→

       rdt2.1                          rdt2.2

# Recap (rdt3.0: deals with packet loss Timer-driven loss detection)

*sender*  *receiver*

pkt0 → pkt0 → ❌

Timeout

pkt0 → pkt0 → ☑

ACK0

pkt1 → pkt1

*sender*  *receiver*

pkt0 → pkt0 → ☑

Timeout

❌ ← pkt0

pkt0 → pkt0 → duplicate

ACK0

pkt1 ← pkt1

*sender*  *receiver*

pkt0 → pkt0 → ☑

Timeout

ACK0

pkt0 → pkt0 → duplicate

ACK0

# Recap (stop-and-wait is inefficient)

TRANS

Inefficient
TRANS << RTT

sender                                    receiver

first packet bit transmitted, t = 0

last packet bit transmitted, $t = L / R$

RTT

first packet bit arrives

last packet bit arrives, send ACK

ACK

ACK arrives, send next
packet, $t = RTT + L / R$

- R= 1 Gbps link, 15 ms prop. delay, L= 8000 bit packet → total t = 30.008 msec

→ 1,000 bytes in 30.008 milliseconds, gives throughput of only 267 kbps

# Recap (Pipelined protocol: Sliding window)

- Window: set of adjacent sequence numbers (window size *n*)
- Send up to *n* packets at a time
  - Sender can send packets in its window
  - Receiver can accept packets in its window
  - Window slides on successful reception/ acknowledgement

*sender*

n

A

sequence number →

*receiver*

n

B

Already ACK'd
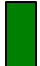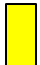
Sent but not ACK'd

Cannot be sent

Received and ACK'd

Acceptable but not yet received

Cannot be received
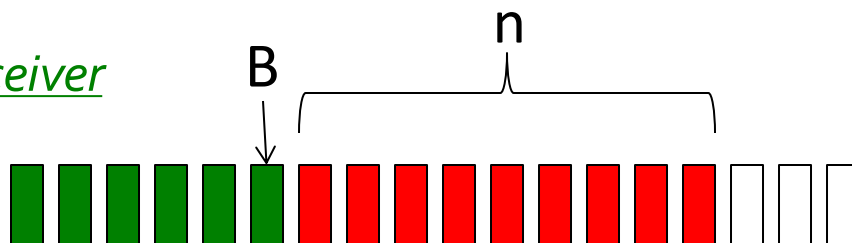
# Recap (Go-Back-N)

- Sender transmits up to *n* unacknowledged packets
- Receiver only accepts packets in order
  - Receiver discards out-of-order packets
  - Receiver uses cumulative acknowledgements
  - Sender sets timer for 1st outstanding ack (A+1), if timeout, retransmit A+1, ... A+n
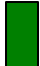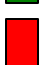
*sender*

n

A

sequence number →

- Already ACK'd
- Sent but not ACK'd
- Cannot be sent

*receiver*

n

B

- Received and ACK'd
- Acceptable but not yet received
- Cannot be received

# Recap (GBN example)

*sender window (N=4)*      *sender*                   *receiver*

`0 1 2 3` 4 5 6 7 8      send pkt0

`0 1 2 3` 4 5 6 7 8      send pkt1

`0 1 2 3` 4 5 6 7 8      send pkt2      **X** *loss*     receive pkt0, send ack0

`0 1 2 3` 4 5 6 7 8      send pkt3               receive pkt1, send ack1

                    (wait)

                                       receive pkt3, discard, (re)send ack1

0 `1 2 3 4` 5 6 7 8      rcv ack0, send pkt4

0 1 `2 3 4 5` 6 7 8      rcv ack1, send pkt5

                                       receive pkt4, discard, (re)send ack1

                   ignore duplicate ACK      receive pkt5, discard, (re)send ack1

                     *pkt 2 timeout*

0 1 `2 3 4 5` 6 7 8      send pkt2

0 1 `2 3 4 5` 6 7 8      send pkt3

0 1 `2 3 4 5` 6 7 8      send pkt4               rcv pkt2, deliver, send ack2

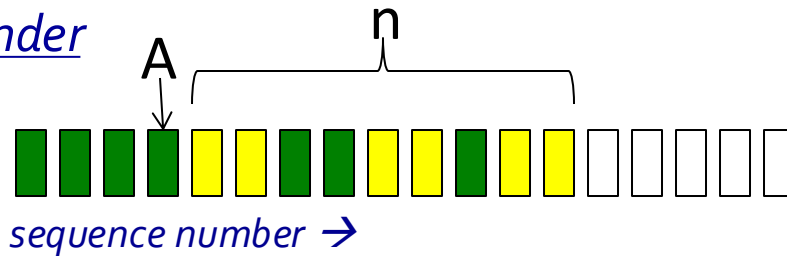0 1 `2 3 4 5` 6 7 8      send pkt5               rcv pkt3, deliver, send ack3

                                       rcv pkt4, deliver, send ack4

                                       rcv pkt5, deliver, send ack5

# Recap (Selective repeat)

- Sender transmits up to *n* unacknowledged packets
- Assume packet m is lost, m+1 is not
- Receiver indicates packet m+1 is correctly received
- Sender retransmits only packet on m timeout
- Efficient in retransmission, but requires book-keeping

*sender*

A    n

*sequence number* →

■ Already ACK'd

■ Sent but not ACK'd

□ Cannot be sent

*receiver*

B    n

■ Received and ACK'd

■ Acceptable but not yet received

□ Cannot be received

# Recap (Selective repeat example)

*sender window (N=4)*  *sender*  *receiver*

| | |
|---|---|
| `0 1 2 3`4 5 6 7 8 | send pkt0 |
| `0 1 2 3`4 5 6 7 8 | send pkt1 |
| `0 1 2 3`4 5 6 7 8 | send pkt2 |
| `0 1 2 3`4 5 6 7 8 | send pkt3 |
| | (wait) |

**X** *loss*

receive pkt0, send ack0
receive pkt1, send ack1

receive pkt3, buffer,
   send ack3

| | |
|---|---|
| 0 `1 2 3 4`5 6 7 8 | rcv ack0, send pkt4 |
| 0 1 `2 3 4 5`6 7 8 | rcv ack1, send pkt5 |

receive pkt4, buffer,
   send ack4
receive pkt5, buffer,
   send ack5

record ack3 arrived

*pkt 2 timeout*

| | |
|---|---|
| 0 1 `2 3 4 5`6 7 8 | send pkt2 |
| 0 1 `2 3 4 5`6 7 8 | record ack4 arrived |
| 0 1 `2 3 4 5`6 7 8 | record ack5 arrived |
| 0 1 `2 3 4 5`6 7 8 | |

rcv pkt2; deliver pkt2,
pkt3, pkt4, pkt5; send ack2

0 1 2 3 4 `6 7 8 9`  *what happens when ack2 arrives?*

# TCP header

URG: urgent data
(generally not used)

ACK: ACK #
valid

PSH: push data now
(generally not used)

RST, SYN, FIN:
connection estab
(setup, teardown
commands)

Internet
checksum
(as in UDP)

| 32 bits | |
|---|---|
| source port # | dest port # |
| sequence number | |
| acknowledgement number | |
| head len / not used / UAPRSF | receive window |
| checksum | Urg data pointer |
| options (variable length) | |
| application data (variable length) | |

counting
by bytes
of data
(not segments!)
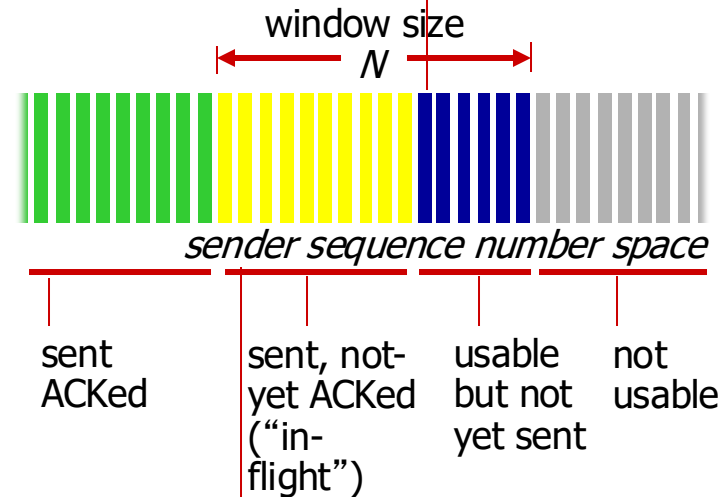
# bytes
rcvr willing
to accept

# TCP seq. numbers, ACKs

- <u>Sequence numbers:</u>
  - byte stream "number" of first byte in segment's data

- <u>Acknowledgements:</u>
  - seq # of next byte expected from other side
  - cumulative 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") | usable but not yet sent | not usable |
|---|---|---|---|

incoming segment to sender

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

# TCP seq. numbers, ACKs

- Starting seq. no. 42 for client and 79 for server
- After TCP established, before data sent, client waits for byte 79 and server waits for byte 42
- Server replies with ack 43, seg 79 and echo back 'C'
- Acknowledgment is said to be piggybacked

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 reliable data transfer

TCP creates rdt service on top of IP's unreliable service
- sender/receiver agree to establish connection "handshake'
- pipelined segments for efficiency
- cumulative acks for acknowledgements
- single retransmission timer (for loss detection)
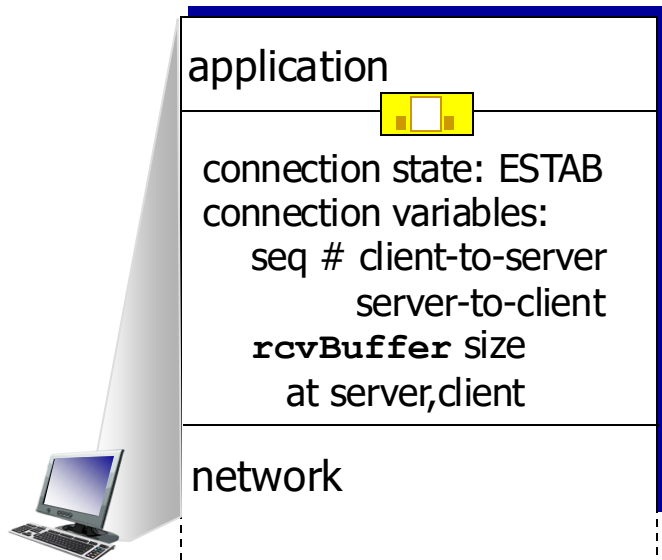- checksums (for error detection)

retransmissions  triggered by:
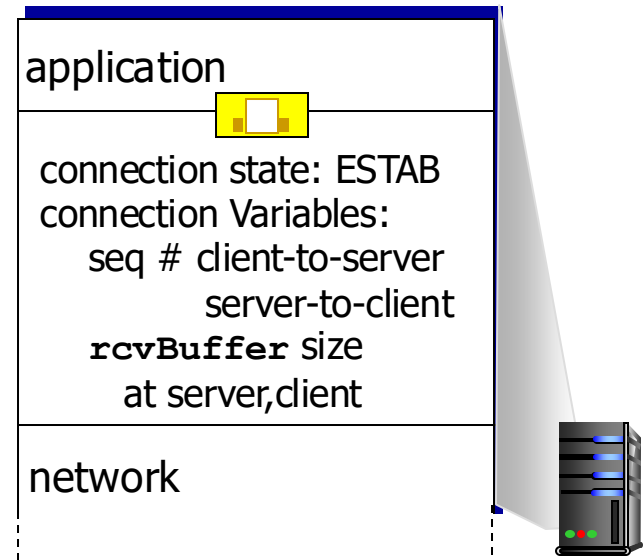- timeout events
- duplicate acks

simplified TCP sender:
- ignore duplicate acks
- ignore flow control, congestion control

# Connection Management

before exchanging data, sender/receiver "handshake":

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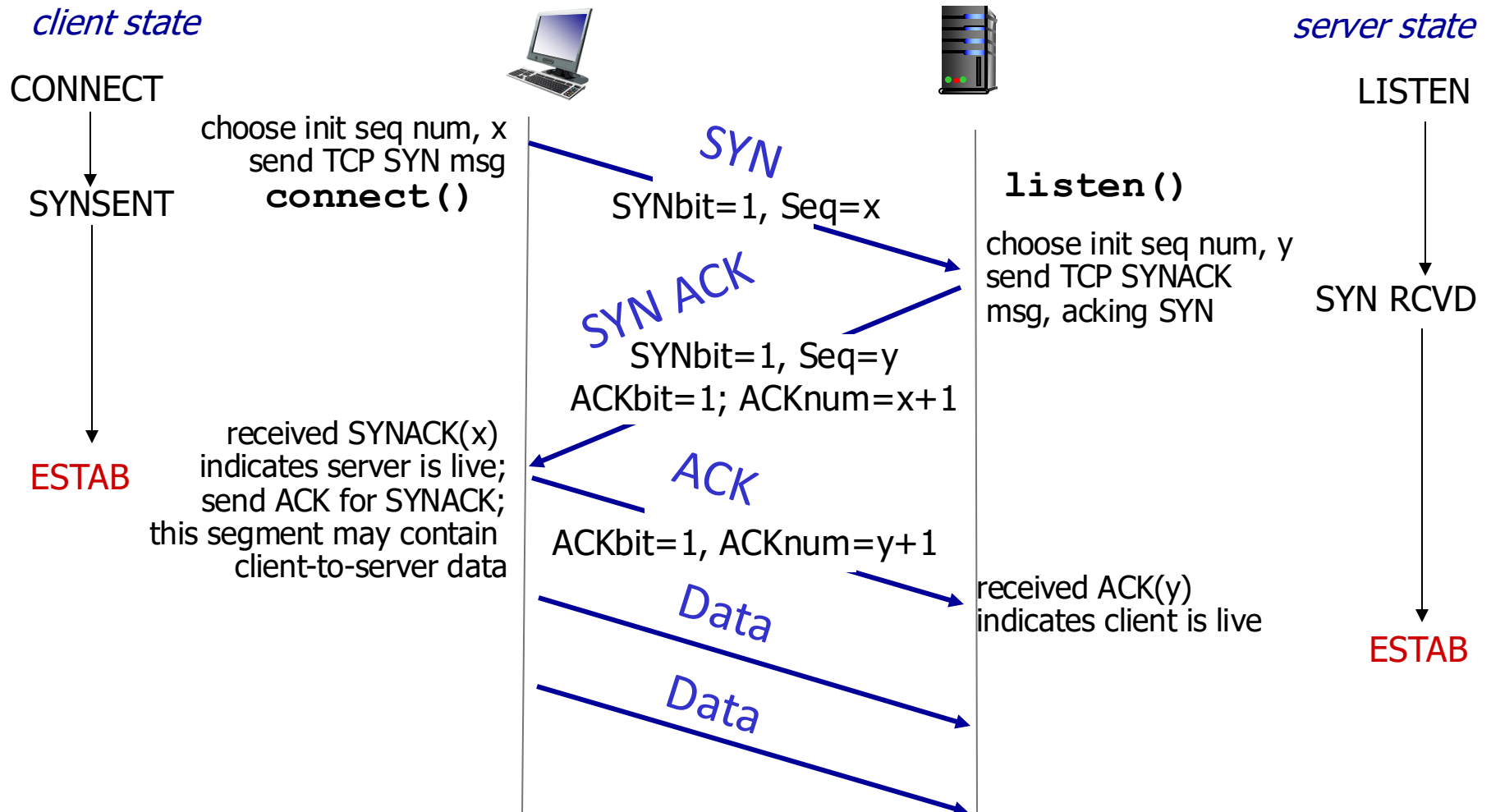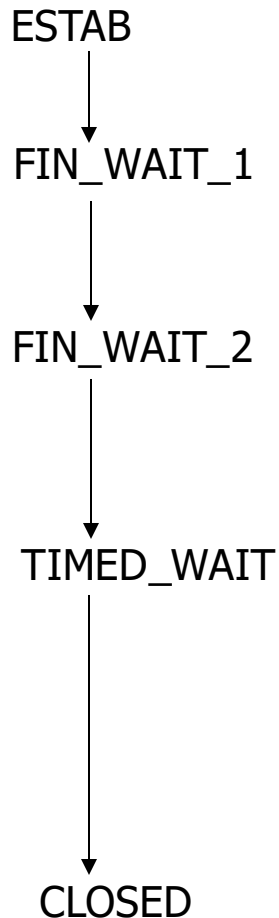
agree on connection parameters

application

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

network

`connect (sockfd,..)`

application

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

network

`listen (sockfd,n);`
`connfd = accept();`

# Establishing a TCP Connection: 3-way handshake

*client state*

**CONNECT**

SYNSENT

ESTAB

*server state*

LISTEN

**listen()**

SYN RCVD

ESTAB

choose init seq num, x
send TCP SYN msg
**connect()**

*SYN*

SYNbit=1, Seq=x

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

*SYN ACK*

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

*ACK*

ACKbit=1, ACKnum=y+1

received ACK(y)
indicates client is live

*Data*

*Data*

# Closing a TCP connection

*client state*

*server state*

ESTAB

ESTAB

**close()**

FIN_WAIT_1

can no longer
send but can
receive data

FINbit=1, seq=x

CLOSE_WAIT

ACKbit=1; ACKnum=x+1

can still
send data

FIN_WAIT_2

wait for server
close

FINbit=1, seq=y

LAST_ACK

TIMED_WAIT

can no longer
send data

ACKbit=1; ACKnum=y+1

timed wait
for 2*max
segment lifetime

CLOSED

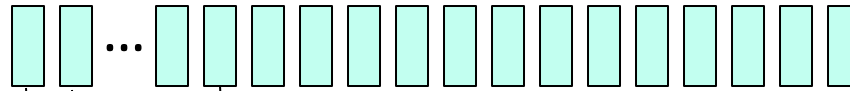CLOSED

# TCP Segments

*Sender*

… 

Byte

TCP Data
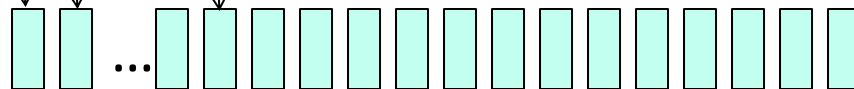
Sender sends packet and starts timer
  Data starts with sequence number X
  Packet contains B bytes [X, X+1, X+2, ….X+B-1]
  Expiration interval: `TimeOutInterval`

Segment number
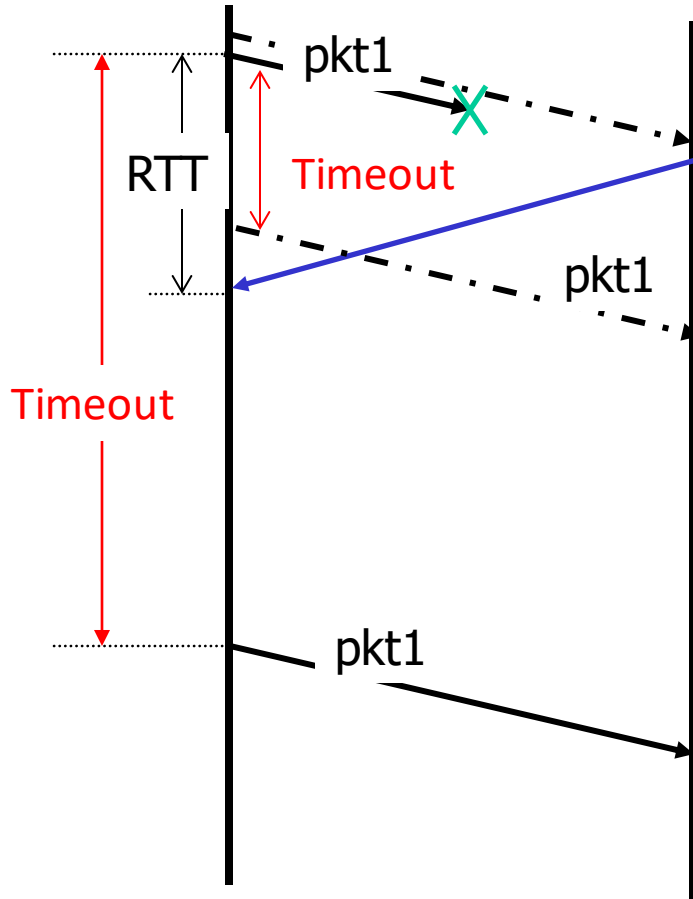  = 1st byte in
    segment

TCP Data

Receiver sends an ACK
  cumulative (GBN) or selective acknowledges

*receiver*

…

# TCP retransmission timer



Timeout too long – inefficient
Timeout too short – duplicate packets
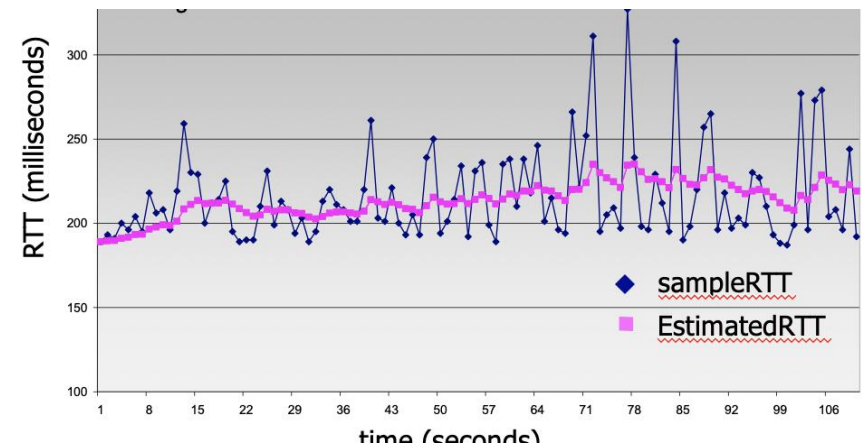Make timeout proportional to RTT
**EstimatedRTT = (1- α)*EstimatedRTT +
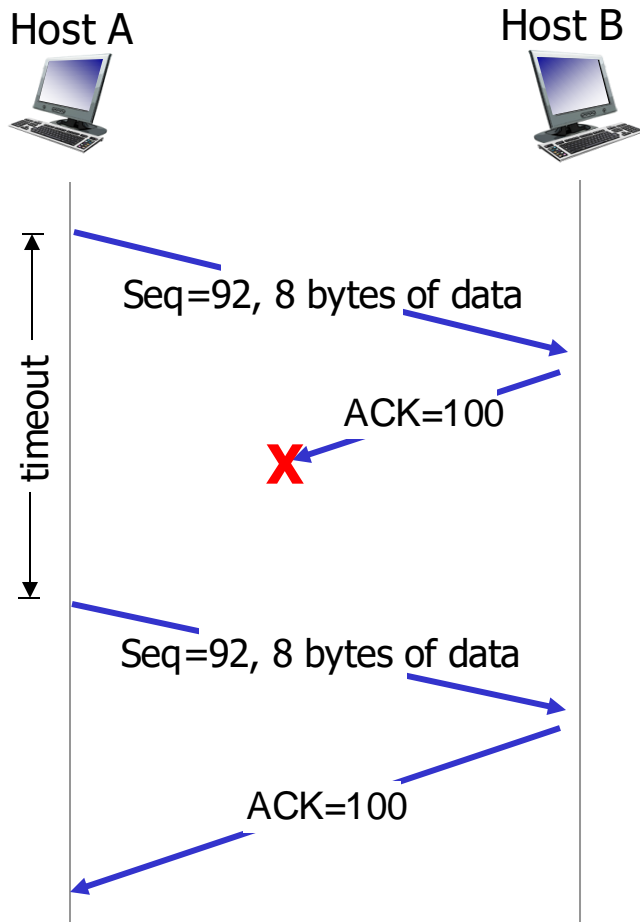α*SampleRTT**
(exponential weighted moving average)
SampleRTT = t segment sent – t segment acknowledged
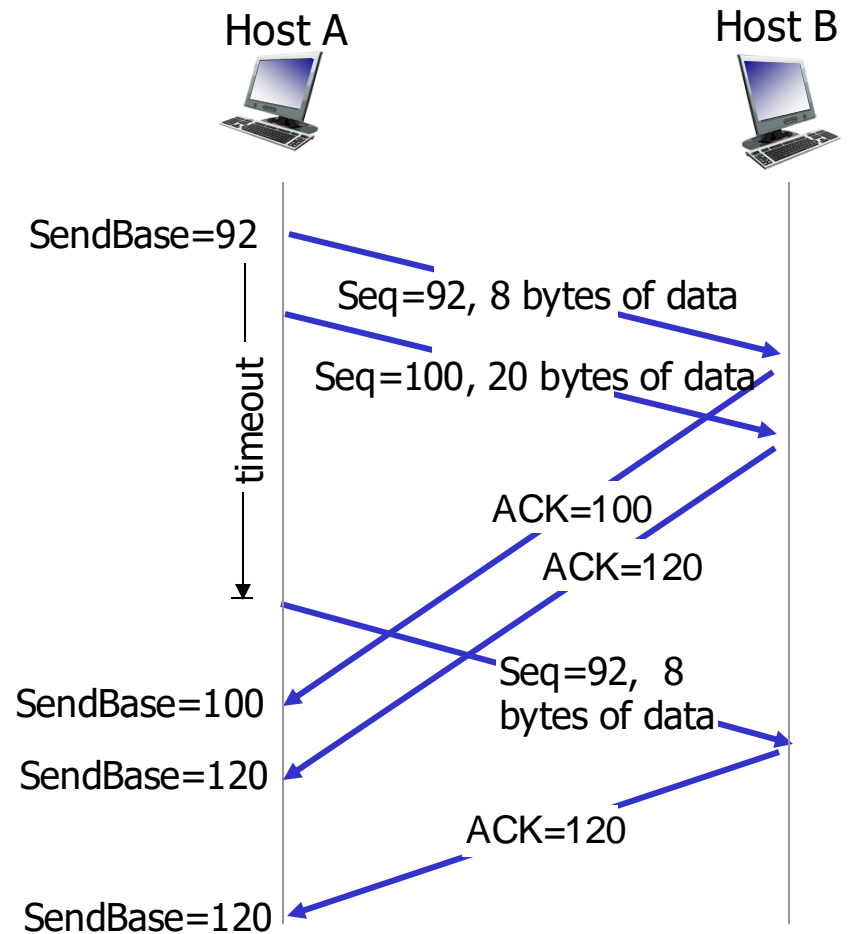```
This time keeps changing with network
conditions
```
- influence of past sample decreases exponentially fast
- typical value: $\alpha$ = 0.125

# TCP: retransmission scenarios



**Host A**         **Host B**         **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

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

ACK=120

SendBase=120

premature timeout

# TCP: retransmission scenarios



Host A

Host B

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

ACK=100

**X**

ACK=120

timeout

Seq=120, 15 bytes of data

# TCP fast retransmit

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

fast retransmit after sender
receipt of triple duplicate ACK
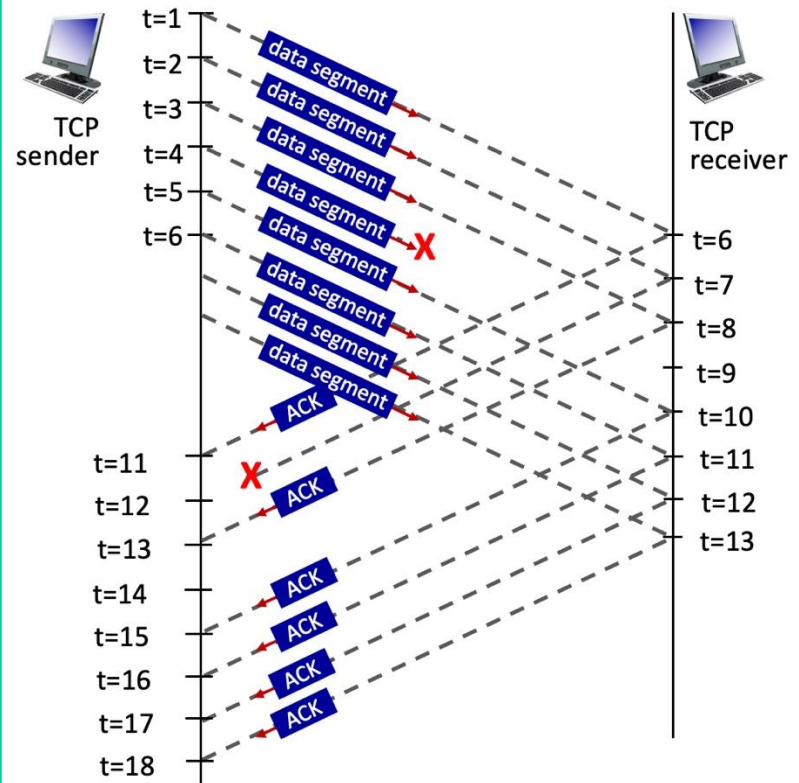
# TCP example

TCP sender sends 8 TCP segments at t = 1, 2, 3, 4, 5, 6, 7, 8. Suppose the initial value of the sequence number is 0 and every segment sent to the receiver each contains 100 bytes. The delay between the sender and receiver is 5 time units, and so the first segment arrives at the receiver at t = 6.

- Sender's sequence numbers: 0, 100, 200, 300,….
- Receiver's acknowledgement numbers: 100, 200, 300, 300,..

# TCP ACK generation [RFC 1122, RFC 2581]

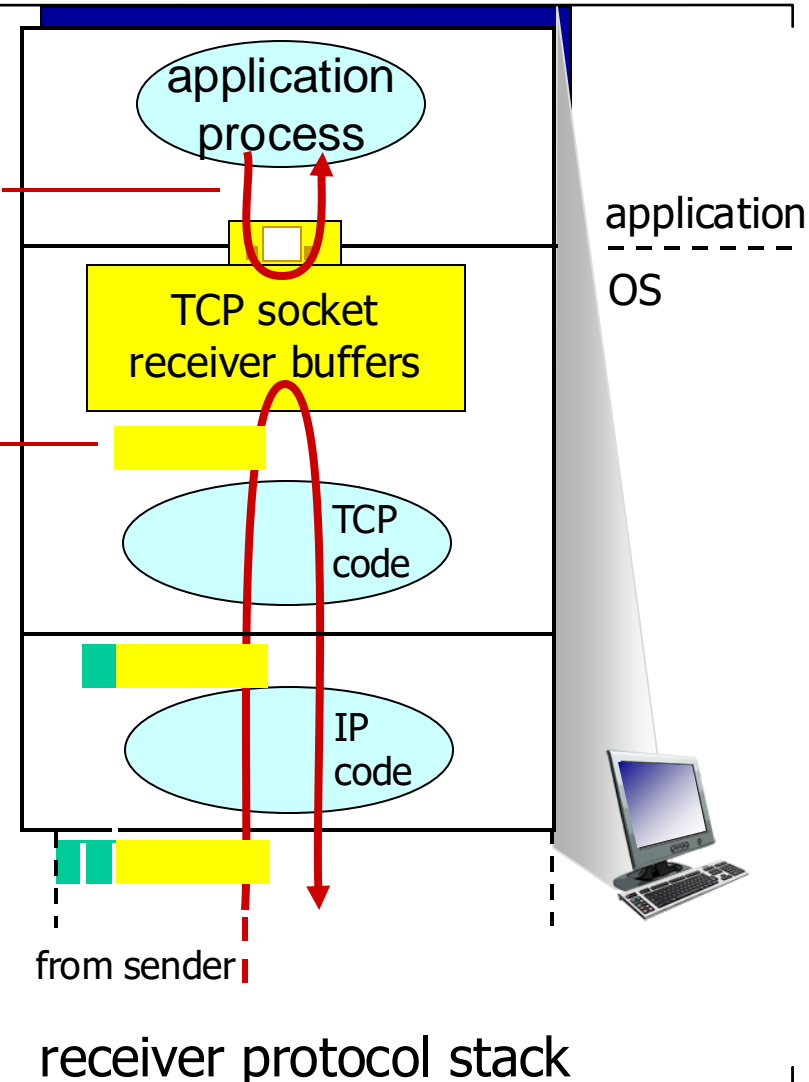| *event at receiver* | *TCP receiver action* |
|---|---|
| arrival of in-order segment with expected seq #. All data up to expected seq # already ACKed | delayed ACK. Wait up to 500ms for next segment. If no next segment, send ACK |
| arrival of in-order segment with expected seq #. One other segment 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 |
| arrival of segment that partially or completely fills gap | immediate send ACK, provided that segment starts at lower end of gap |

# TCP flow control

"no one can drink from a firehose"

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
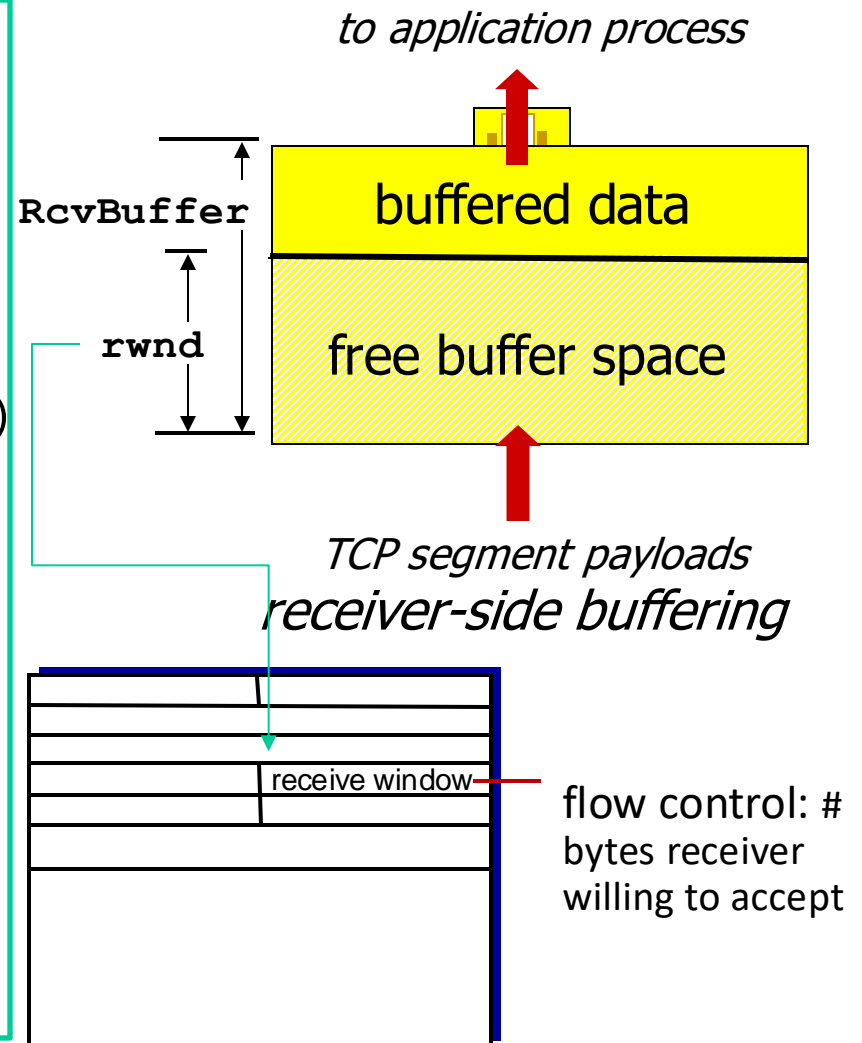
application process

TCP socket receiver buffers

TCP code

IP code

from sender

application

OS

receiver protocol stack

# Receiver window - rwnd

receiver "advertises" free buffer space by including **rwnd** (*receiver window*) 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

*to application process*

**RcvBuffer**

buffered data

**rwnd**

free buffer space

*TCP segment payloads*

*receiver-side buffering*

receive window

flow control: # bytes receiver willing to accept

# TCP flow control

*to application process*

**RcvBuffer**

buffered data

**rwnd**

free buffer space

*TCP segment payloads*

rwnd=RcvBuffer-[LastByteRcvd-LastByteRead]
Sender keeps: LastByteSent - LastByteAcked ≤ rwnd

**Buffered data**

last byte
Read (application

last byte Rcvd
(network)

application
process

application

OS

TCP socket
receiver buffers

TCP
code

IP
code

from sender

receiver protocol stack

# Summary

Today:
- TCP reliable data transfer
- TCP connection
- TCP segments
- TCP sender and receiver ACK generation
- TCP flow control

Canvas discussion:
- Reflection
- Exit ticket

Next time:
- read 3.6 and 3.7 of K&R (TCP congestion control)
- follow on Canvas! Material and announcements

# Any questions?