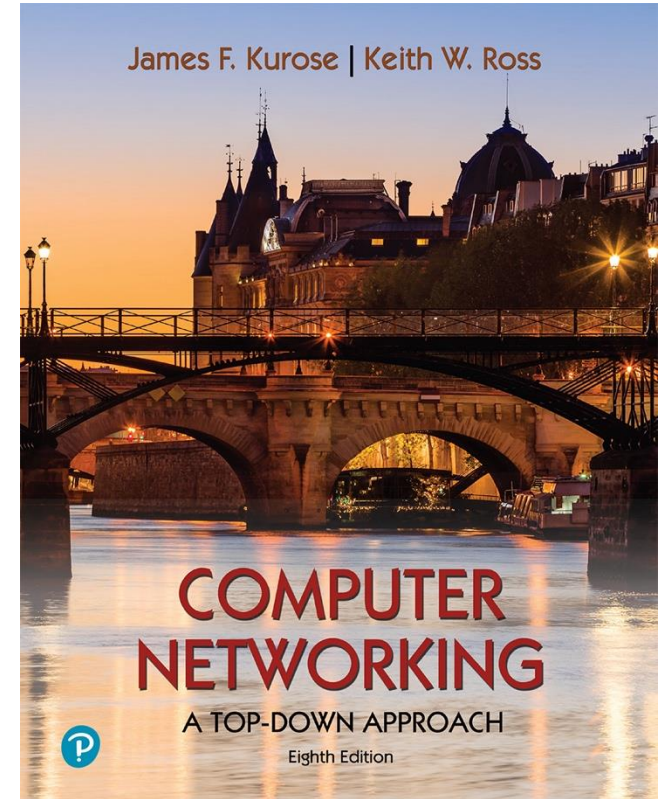# Chapter 3
# Transport Layer

## Yaxiong Xie

Department of Computer Science and Engineering
University at Buffalo, SUNY

Adapted from the slides of the book's authors
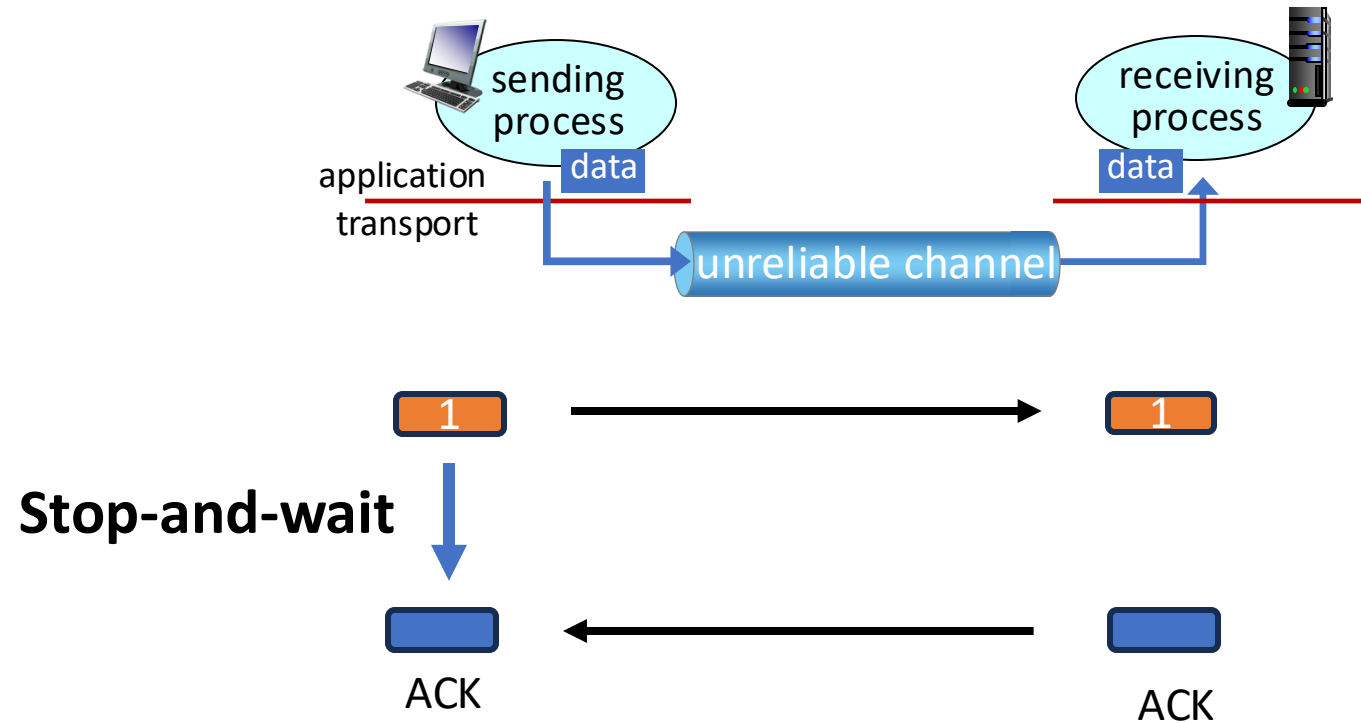
*Computer Networking: A Top-Down Approach*
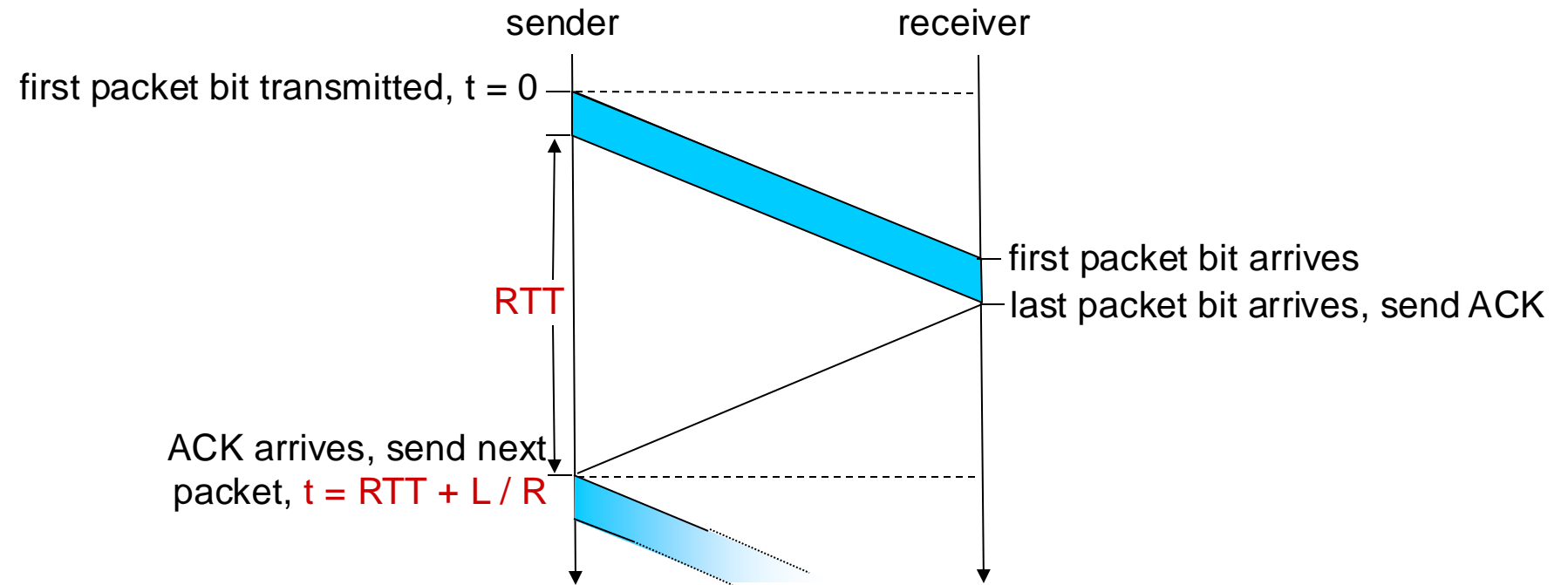
# Stop-and-wait operation

# Stop-and-wait operation



sender                    receiver

first packet bit transmitted, t = 0

RTT

first packet bit arrives

last packet bit arrives, send ACK

ACK arrives, send next packet, t = RTT + L / R
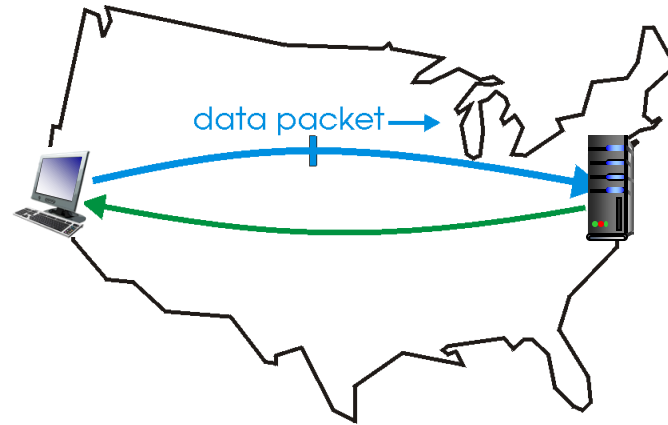
# Pipelined protocols operation
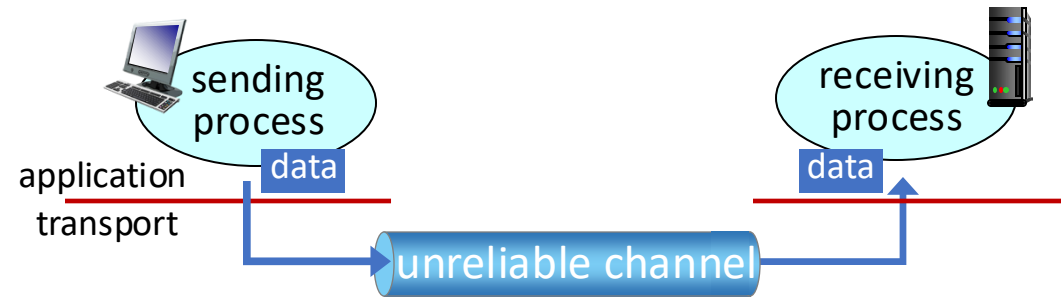
pipelining: sender allows multiple, "in-flight", yet-to-be-acknowledged packets

- range of sequence numbers must be increased
- buffering at sender and/or receiver



(a) a stop-and-wait protocol in operation

# Pipelining



How many packets shall we sent before receiving the ACK?

# Sliding Window



**Window Size N**

sent and ACKed

sent, not yet ACKed

Not yet sent

# Sliding Window

# Sliding Window



**Send one more packet**

**Window Size N**

ACK

sent and ACKed

sent, not yet ACKed

Not yet sent

# Sliding Window



**Window Size N**

ACK

sent and ACKed

sent, not yet ACKed

Not yet sent

# Sliding Window



Timeout : retransmission

sent and ACKed

Not yet sent

sent, not yet ACKed

Window Size N

ACK

- Timer for oldest in-flight packet
- *Timeout(n):* retransmit packet n and all higher seq # packets in window

# Sliding Window + cumulative ACK

# Sliding Window + cumulative ACK

**Window Size N**

sent and ACKed

Not yet sent

sent, not yet ACKed

**Sender**

**Receiver**

received and ACKed

Not yet received

received, not yet ACKed

# Sliding Window + cumulative ACK

Window Size N

... sent and ACKed    Not yet sent

sent, not yet ACKed

**Sender**

**Receiver**

❌

... received and ACKed    Not yet received

received, not yet ACKed

# Sliding Window + cumulative ACK



- *cumulative ACK:* ACK(*n*): ACKs all packets up to, including seq # *n*
  - on receiving ACK(*n*): move window forward to begin at *n+1*

# Sliding Window + cumulative ACK: out-of-order



**Two Question:**

1. do we need to send ACK when we receive an out-of-order packet? If yes, what kind of ACK?

# Sliding Window + cumulative ACK: out-of-order

**Window Size N**

Sender

Receiver

**?**

| | sent and ACKed | | Not yet sent |

| | sent, not yet ACKed |

| | received and ACKed | | Not yet received |

| | received, not yet ACKed |

**Two Question:**
1. do we need to send ACK when we receive an out-of-order packet? If yes, what kind of ACK?
2. how shall we deal with the out-of-order packets we received?

# Sliding Window + cumulative ACK: out-of-order



**Two Question:**
1. do we need to send ACK when we receive an out-of-order packet? If yes, what kind of ACK?
2. how shall we deal with the out-of-order packets we received?

# Sliding Window + cumulative ACK: out-of-order



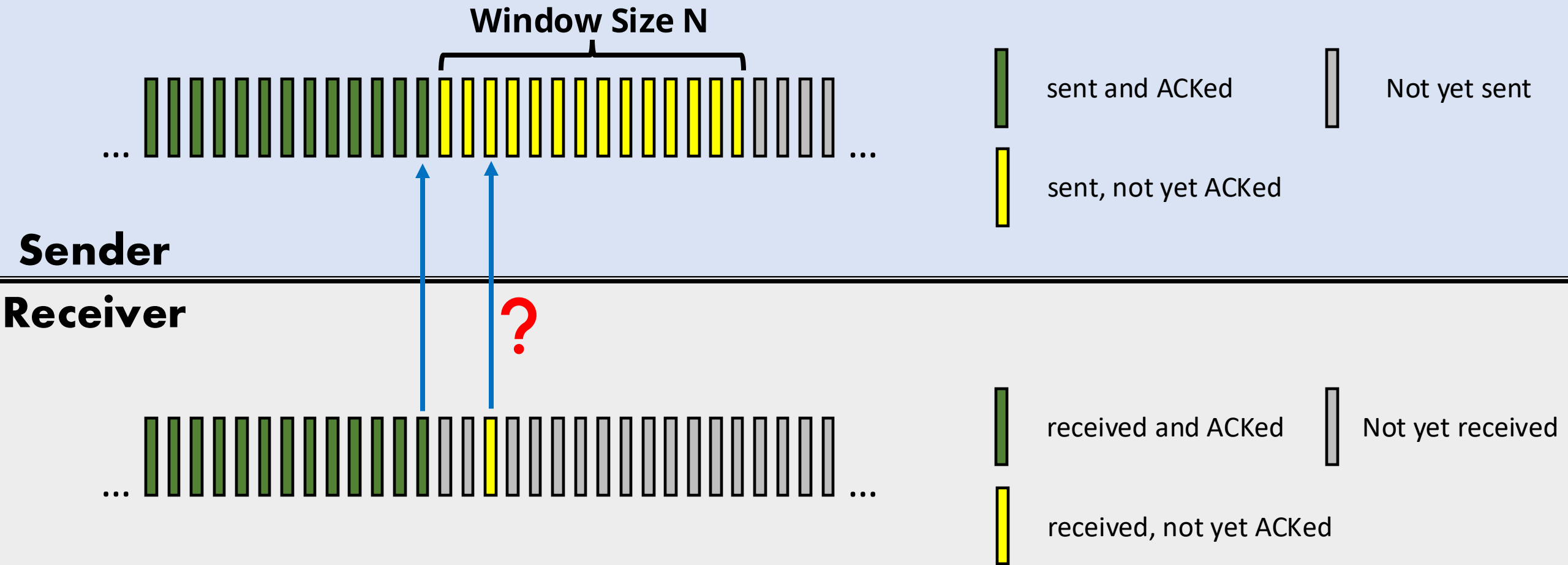- ACK-only: always send ACK for correctly-received packet so far, with highest *in-order* seq #
  - may generate duplicate ACKs

# Sliding Window + cumulative ACK: out-of-order

**Window Size N**

sent and ACKed

Not yet sent

sent, not yet ACKed

**Sender**

**Receiver**

*Retransmit*

received and ACKed

Not yet received

received, not yet ACKed

- on receipt of out-of-order packet:
  - can discard (don't buffer) or buffer: an implementation decision
  - re-ACK pkt with highest in-order seq #

# Go-Back-N: sender

▪ sender: "window" of up to N, consecutive transmitted but unACKed pkts

- • k-bit seq # in pkt header



send_base    nextseqnum

■ already ack'ed    ■ usable, not yet sent

■ sent, not yet ack'ed    ▯ not usable

window size N

▪ *cumulative ACK:* ACK($n$): ACKs all packets up to, including seq # $n$

- • on receiving ACK($n$): move window forward to begin at $n+1$

▪ timer for oldest in-flight packet

▪ *timeout(n):* retransmit packet n and all higher seq # packets in window

# Go-Back-N: receiver

- ACK-only: always send ACK for correctly-received packet so far, with highest *in-order* seq #
  - may generate duplicate ACKs
  - need only remember `rcv_base`

- on receipt of out-of-order packet:
  - can discard (don't buffer) or buffer: an implementation decision
  - re-ACK pkt with highest in-order seq #

Receiver view of sequence number space:



... rcv_base ...

received and ACKed

Out-of-order: received but not ACKed

Not received

# Go-Back-N in action

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                      receive pkt0, send ack0
`0 1 2 3` 4 5 6 7 8          send  pkt3    **X** *loss*      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

# Sliding Window + selective repeat
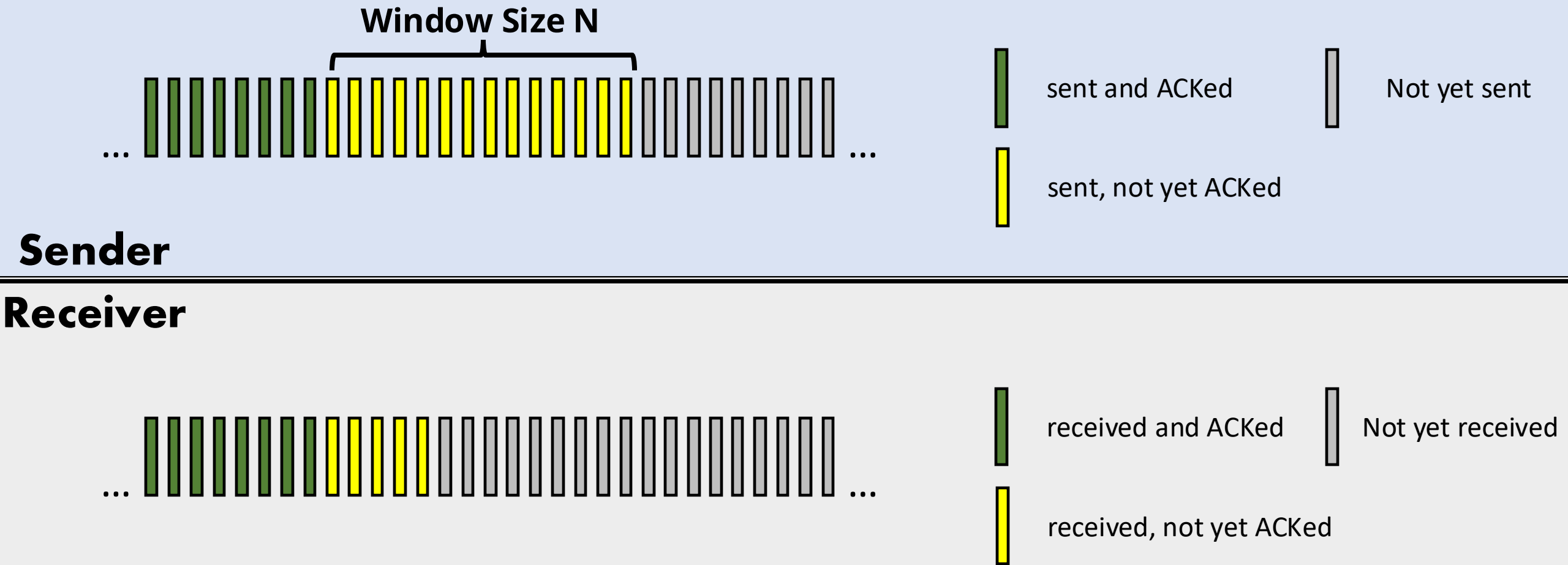
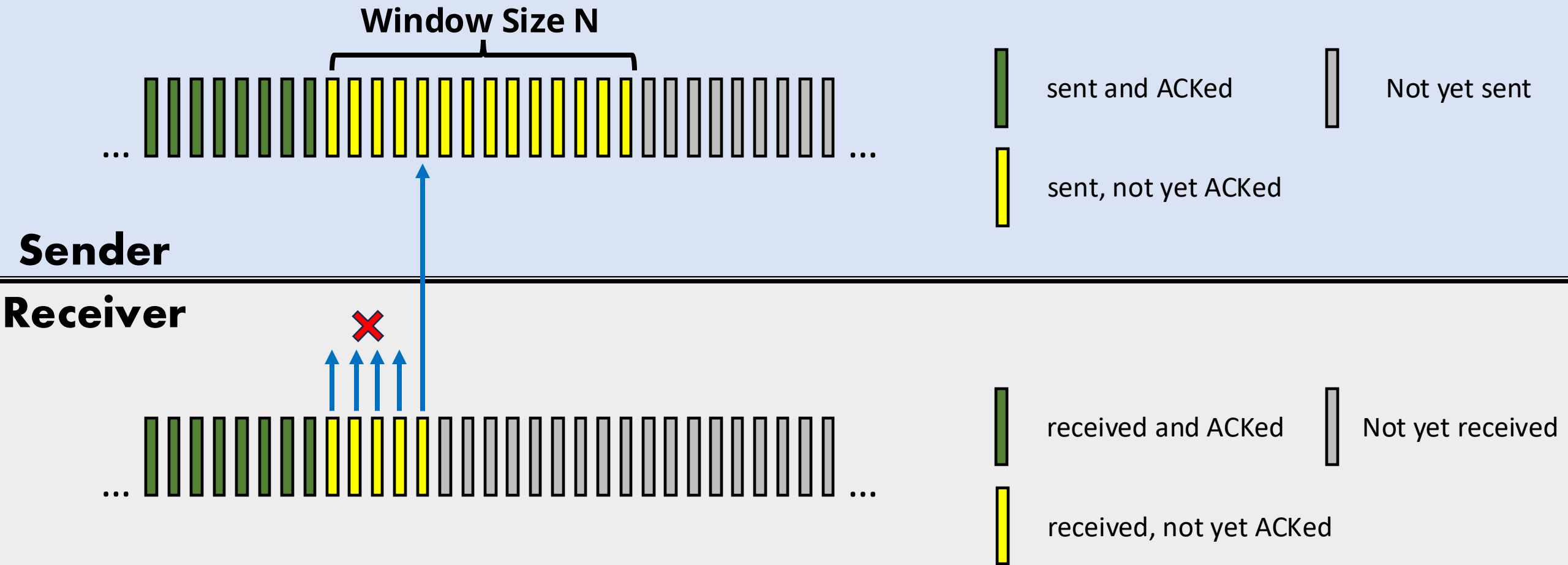# Sliding Window + selective repeat

# Sliding Window + selective repeat

# Sliding Window + selective repeat



**Window Size N**

**Sender**

- **sent and ACKed**
- **sent, not yet ACKed**
- **Not yet sent**

**Receiver**

- **received and ACKed**
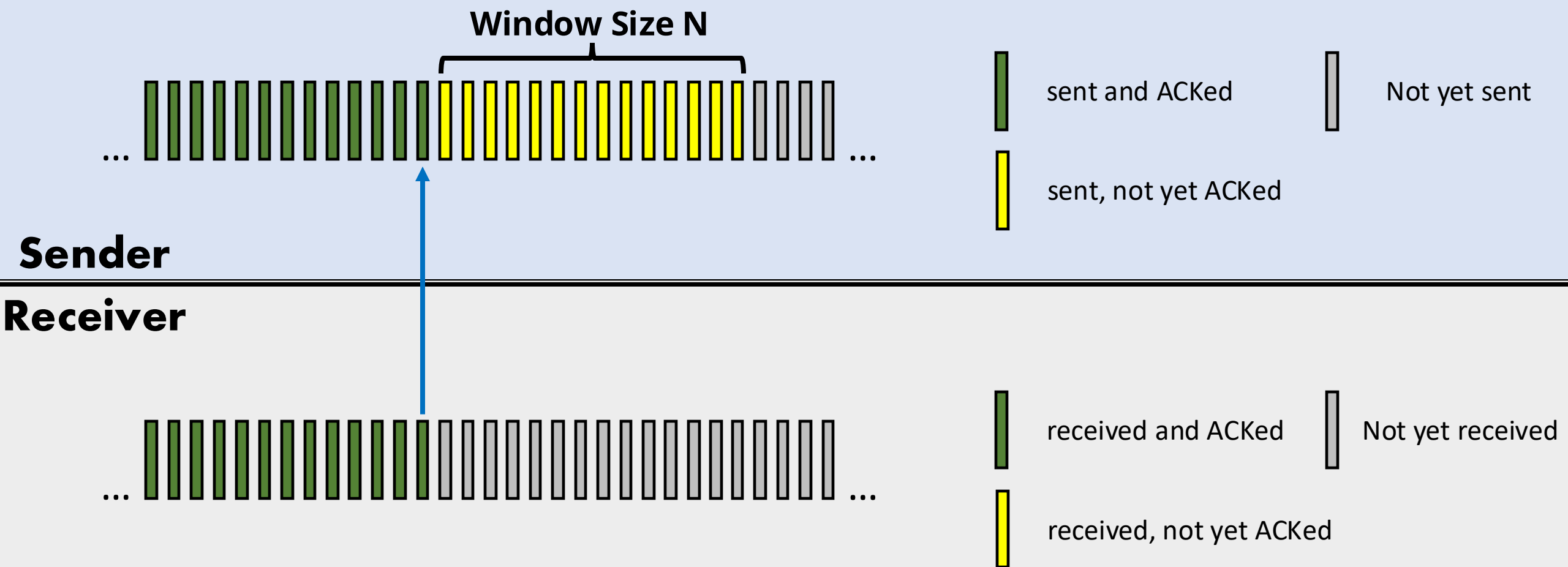- **received, not yet ACKed**
- **Not yet received**

# Sliding Window + selective repeat: out-of-order

# Sliding Window + selective repeat: out-of-order

# Sliding Window + selective repeat: out-of-order



- *receiver individually ACKs* all correctly received packets
  - buffers packets, as needed, for in-order delivery to upper layer

# Sliding Window + selective repeat: out-of-order

**Window Size N**



**Sender**

**Receiver**

| | sent and ACKed | | Not yet sent |
| | sent, not yet ACKed | | |

| | received and ACKed | | Not yet received |
| | received, not yet ACKed | | |

- ▪ *receiver individually ACKs* all correctly received packets
  - buffers packets, as needed, for in-order delivery to upper layer

# Sliding Window + selective repeat: out-of-order



**receiver individually ACKs** all correctly received packets
- buffers packets, as needed, for in-order delivery to upper layer

# Sliding Window + selective repeat: out-of-order



- *receiver individually ACKs* all correctly received packets
  - buffers packets, as needed, for in-order delivery to upper layer

# Sliding Window + selective repeat: out-of-order

**Window Size N**



**Sender**

**Receiver**

| | sent and ACKed | | Not yet sent |
| | sent, not yet ACKed | | |

| | received and ACKed | | Not yet received |
| | received, not yet ACKed | | |

- *receiver individually ACKs* all correctly received packets
  - buffers packets, as needed, for in-order delivery to upper layer

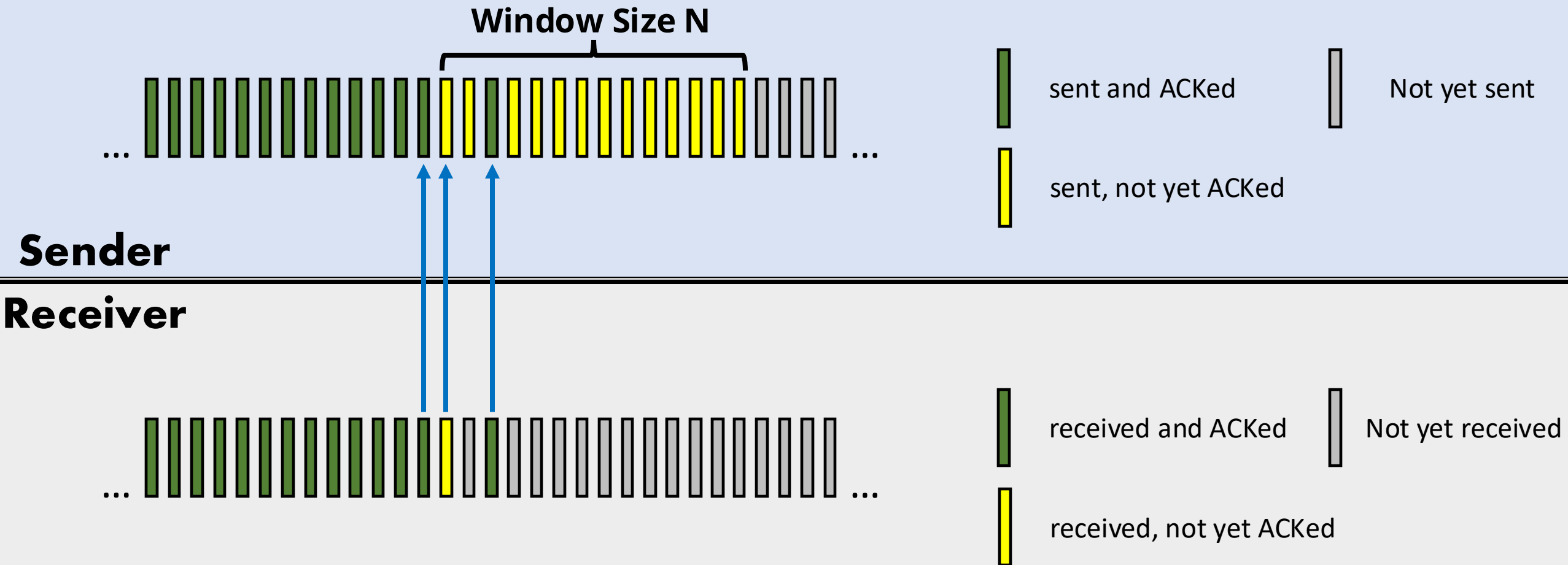# Sliding Window + selective repeat: out-of-order

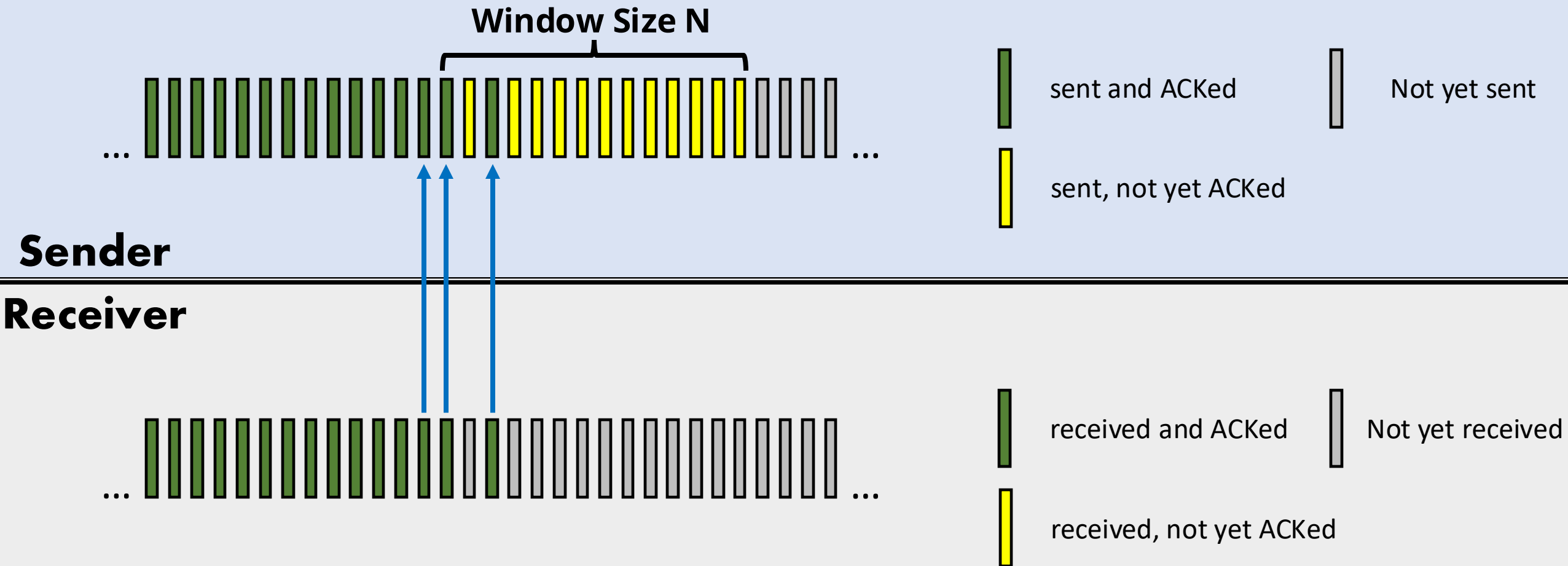

- *receiver individually ACKs* all correctly received packets
  - buffers packets, as needed, for in-order delivery to upper layer

# Sliding Window + selective repeat: out-of-order

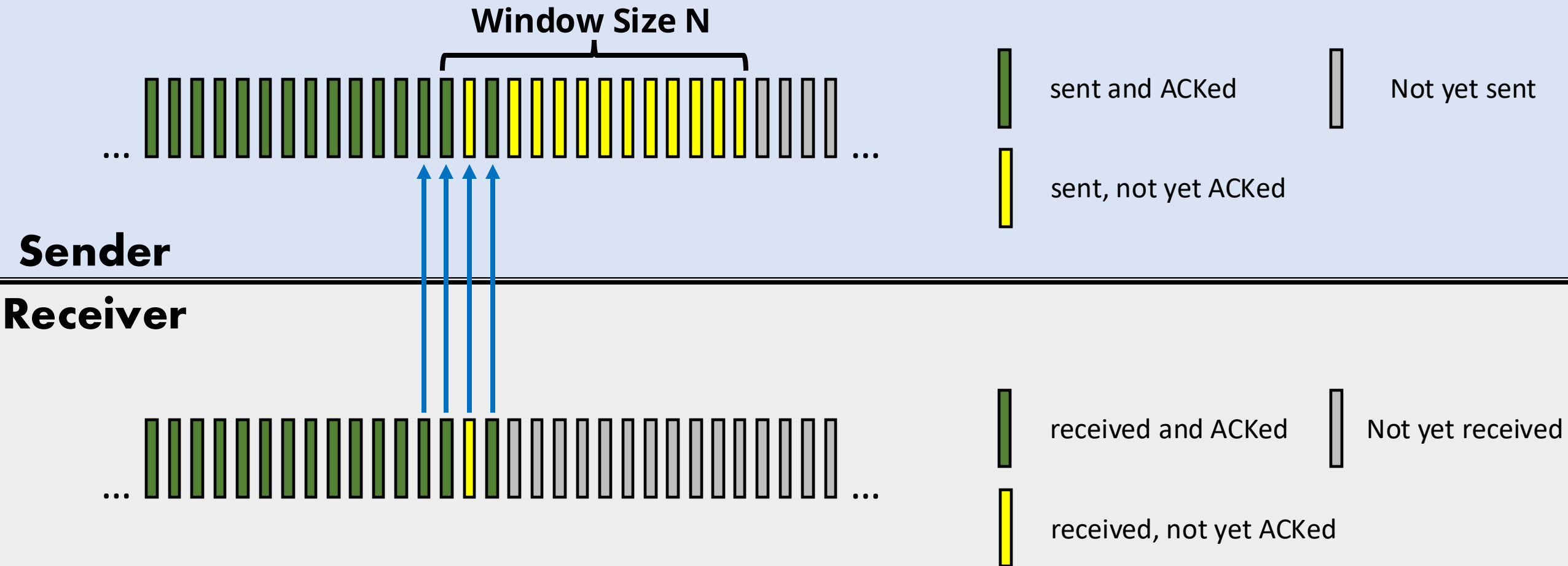**Window Size N**



**Sender**

**Receiver**

■ *receiver individually ACKs* all correctly received packets
  • buffers packets, as needed, for in-order delivery to upper layer

# Sliding Window + selective repeat

**Window Size N**



**Sender**

- sent and ACKed
- Not yet sent
- sent, not yet ACKed

**Receiver**

- received and ACKed
- Not yet received
- received, not yet ACKed

- Maintains (conceptually) a timer for each unACKed pkt
  - timeout: retransmits single unACKed packet associated with timeout

# Selective repeat: the approach

- *pipelining*: *multiple* packets in flight

- *receiver individually ACKs* all correctly received packets
  - buffers packets, as needed, for in-order delivery to upper layer

- sender:
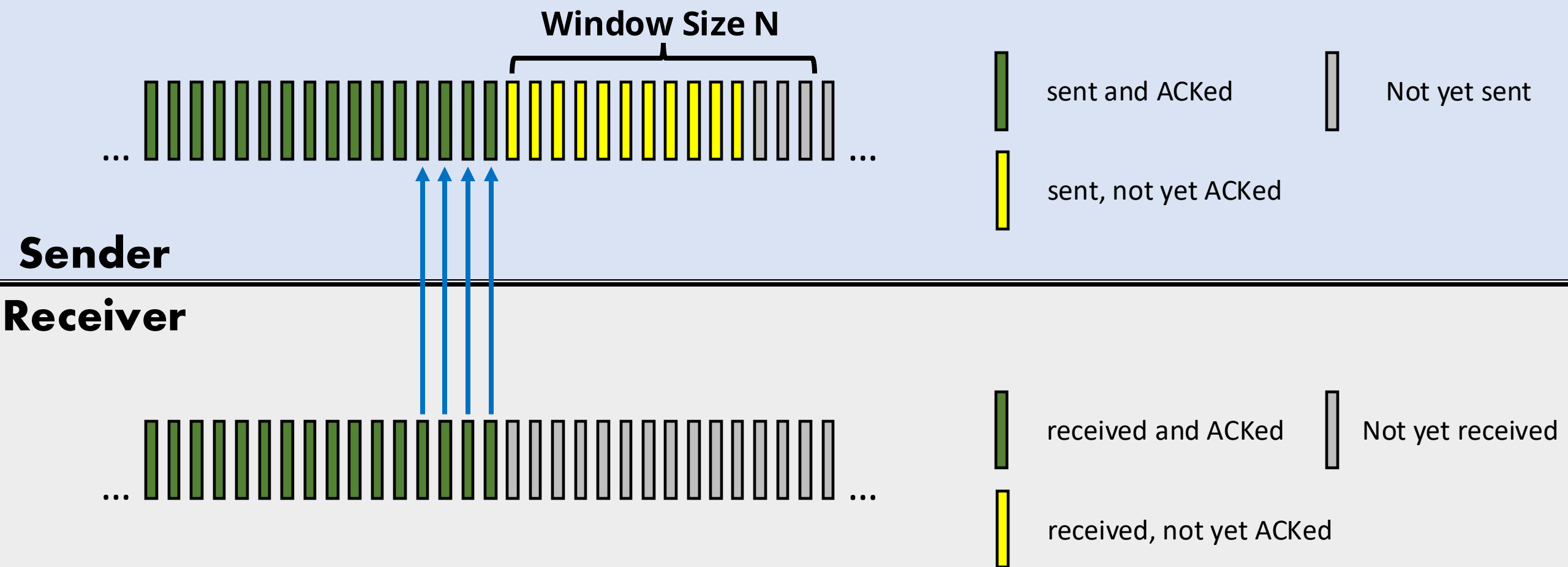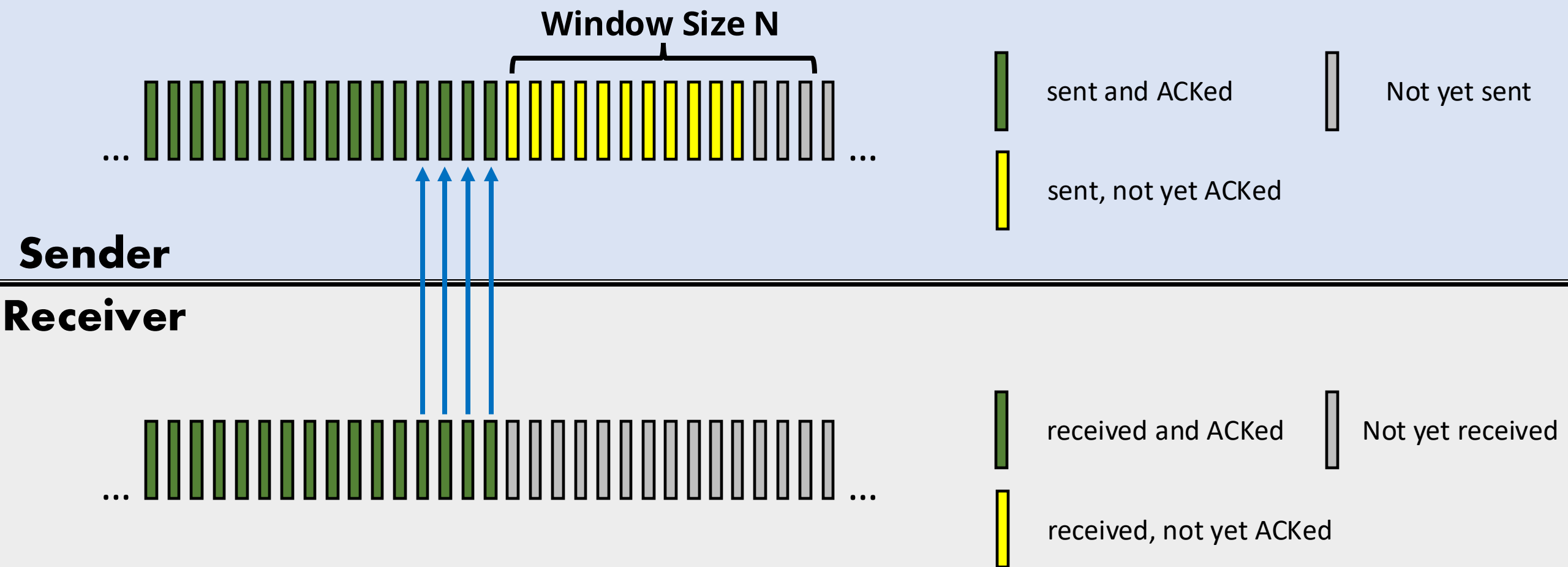  - maintains (conceptually) a timer for each unACKed pkt
    - timeout: retransmits single unACKed packet associated with timeout
  - maintains (conceptually) "window" over  *N* consecutive seq #s
    - limits pipelined, "in flight" packets to be within this window

# Selective repeat: sender, receiver windows



(a) sender view of sequence numbers

# Selective repeat: sender and receiver

## sender

### data from above:

- if next available seq # in window, send packet

### timeout(*n*):

- resend packet *n*, restart timer

### ACK(*n*) in [sendbase,sendbase+N-1]:

- mark packet *n* as received

- if n smallest unACKed packet, advance window base to next unACKed seq #

## receiver

### packet *n* in [rcvbase, rcvbase+N-1]

- send ACK(*n*)

- out-of-order: buffer

- in-order: deliver (also deliver buffered, in-order packets), advance window to next not-yet-received packet

### packet *n* in [rcvbase-N,rcvbase-1]

- ACK(*n*)

### otherwise:

- ignore

# Selective Repeat in action

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, 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
                             record ack3 arrived                  receive pkt5, buffer,
                                                                       send ack5
                             *pkt 2 timeout*
`0 1 2 3 4 5 6 7 8`          send  pkt2
`0 1 2 3 4 5 6 7 8`          (but not 3,4,5)
`0 1 2 3 4 5 6 7 8`                                               rcv pkt2; deliver pkt2,
`0 1 2 3 4 5 6 7 8`                                               pkt3, pkt4, pkt5; send ack2

                             *Q: what happens when ack2 arrives?*

# Selective repeat: a dilemma!

example:

- seq #s: 0, 1, 2, 3 (base 4 counting)
- window size=3

pkt0
pkt1
pkt2
pkt3
pkt0

will accept packet
with seq number 0

(a) no problem

pkt0
pkt1
pkt2

timeout
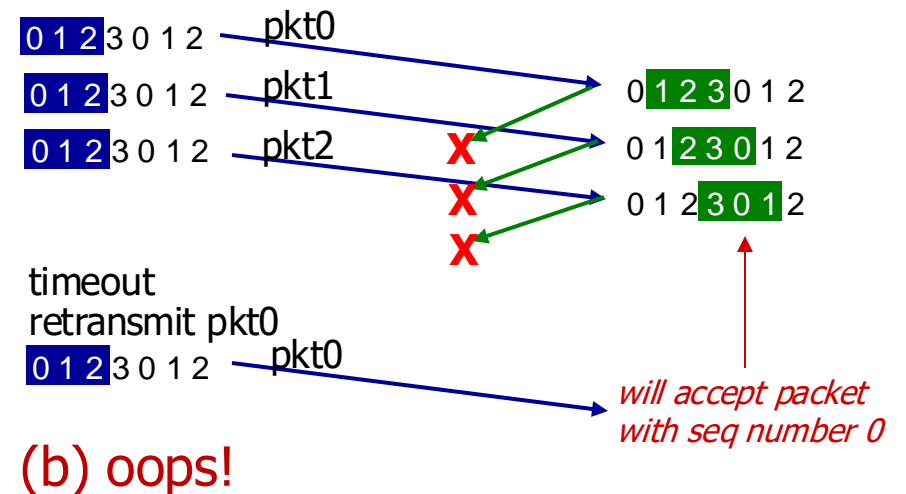retransmit pkt0
pkt0

will accept packet
with seq number 0

(b) oops!

# Selective repeat: a dilemma!

example:
- seq #s: 0, 1, 2, 3 (base 4 counting)
- window size=3

- *receiver can't see sender side*
- *receiver behavior identical in both cases!*
- *something's (very) wrong!*

pkt0
pkt1
pkt2

pkt3

0 1 2 3 0 1 2
0 1 2 3 0 1 2
0 1 2 3 0 1 2

*will accept packet with seq number 0*

pkt2

0 1 2 3 0 1 2
0 1 2 3 0 1 2
0 1 2 3 0 1 2

timeout
retransmit pkt0
pkt0

*will accept packet with seq number 0*

(b) oops!

# 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
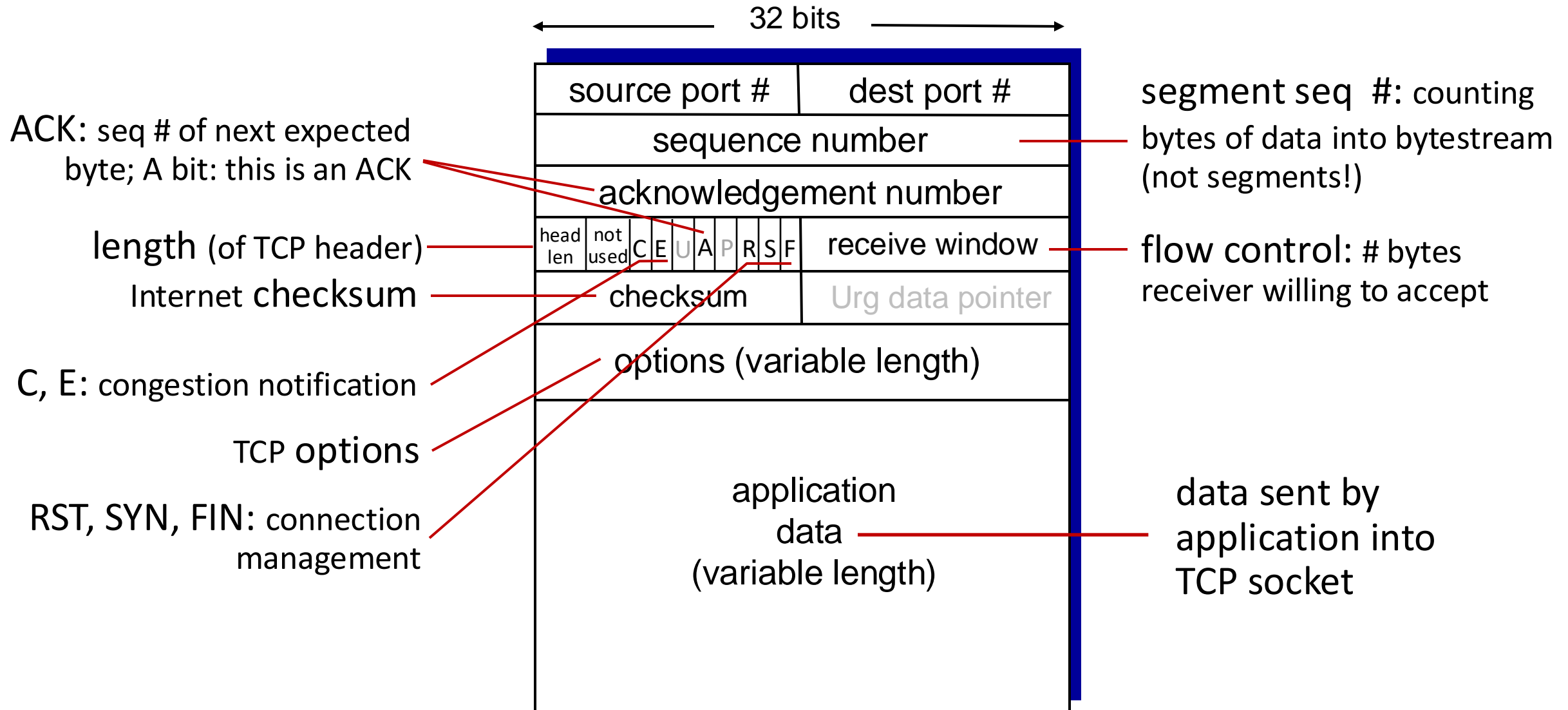- TCP congestion control

# TCP: 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

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

# TCP sequence numbers, ACKs

*Sequence numbers:*

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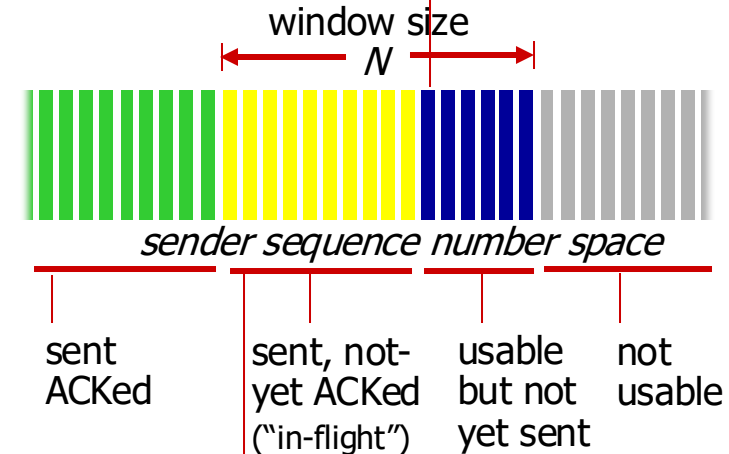
*Acknowledgements:*

- seq # of next byte expected from other side

- cumulative ACK

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

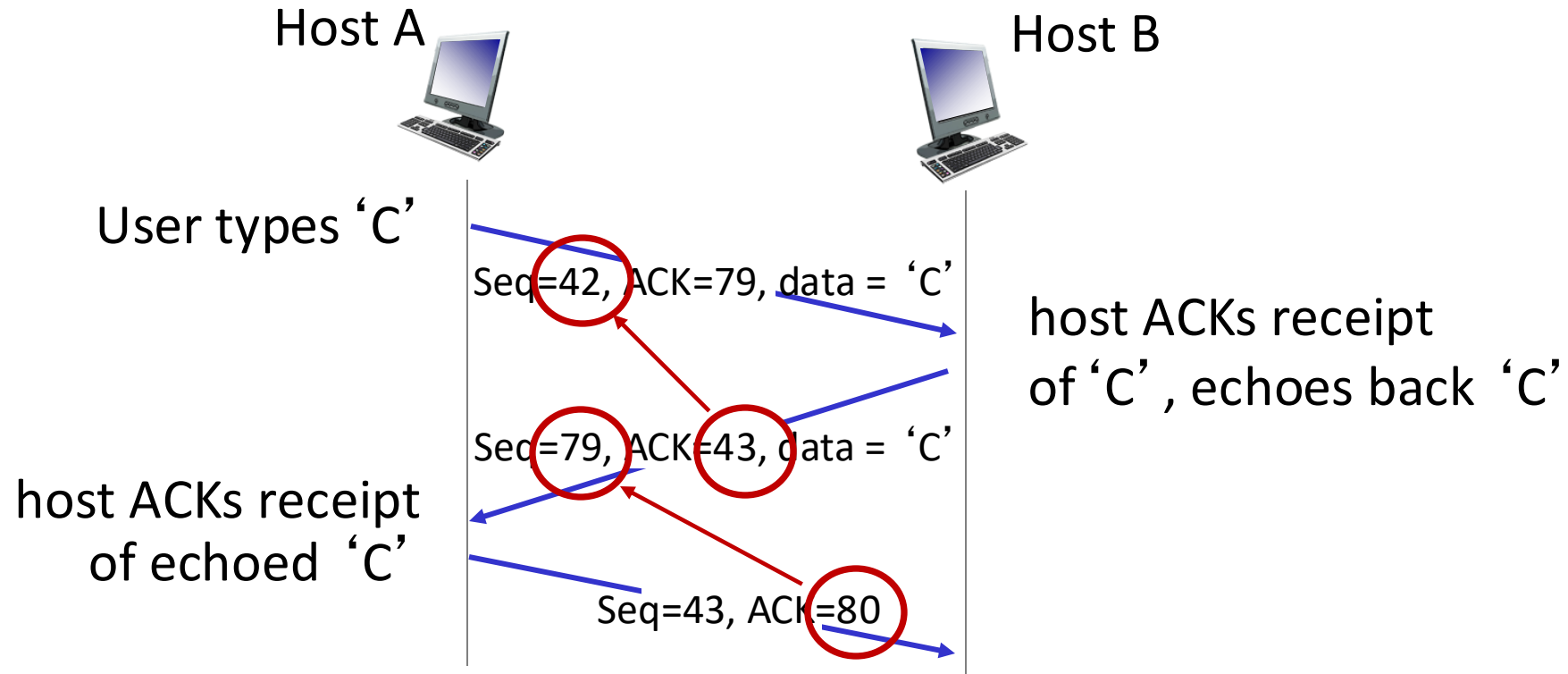- *A:* TCP spec doesn't say, - up to implementor

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

# 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
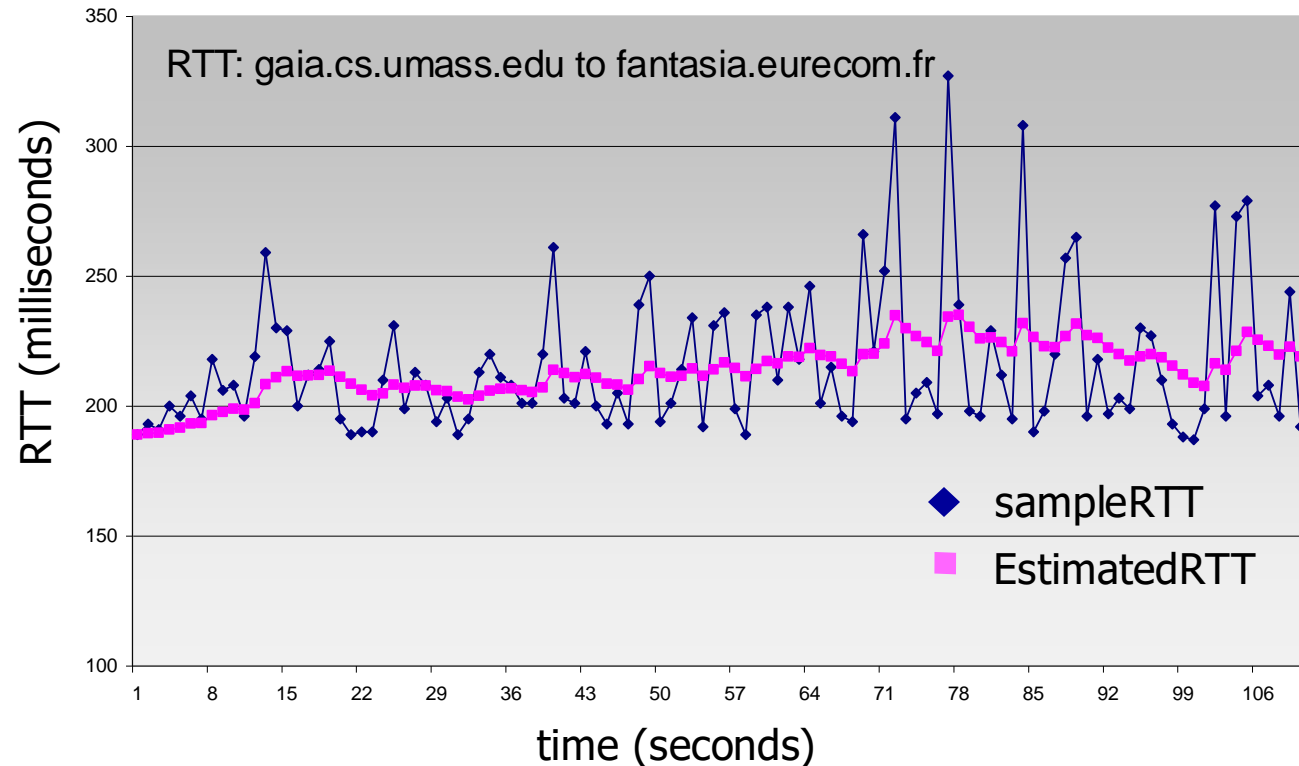
*Q:* how to estimate RTT?

- `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 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

◆ 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, $\beta$ = 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
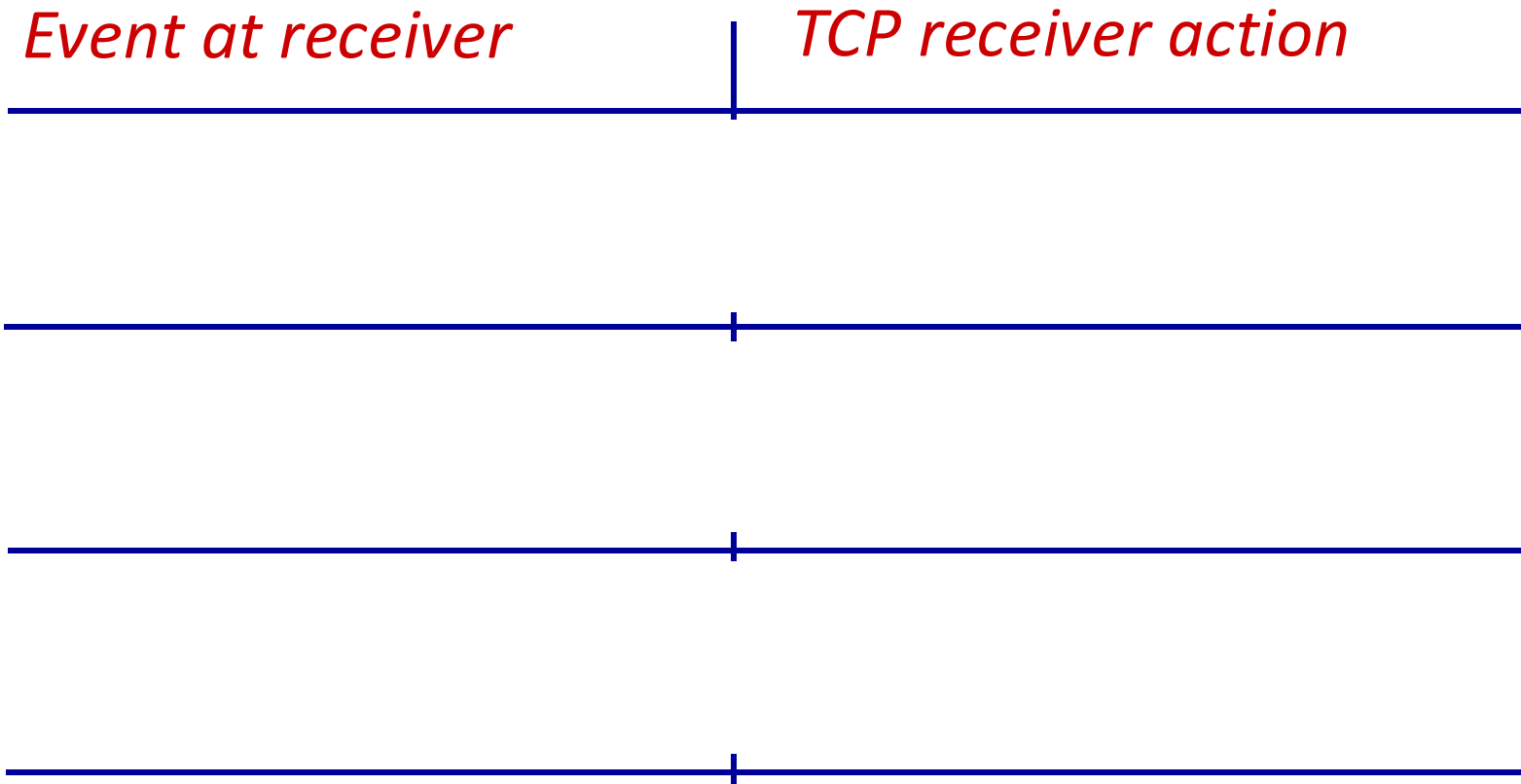  - expiration interval: **`TimeOutInterval`**

*event: timeout*

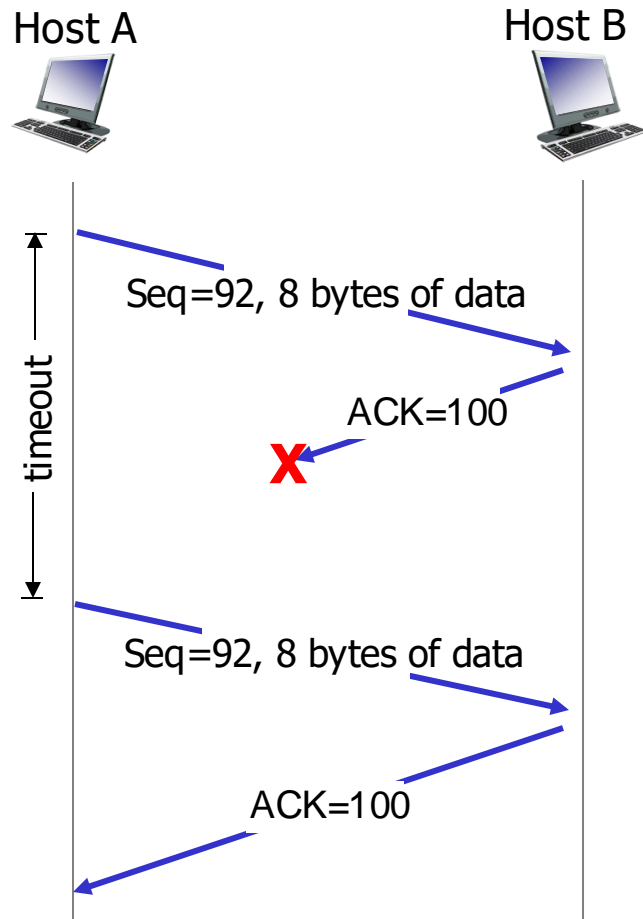- retransmit segment that caused timeout
- restart timer

*event: ACK received*

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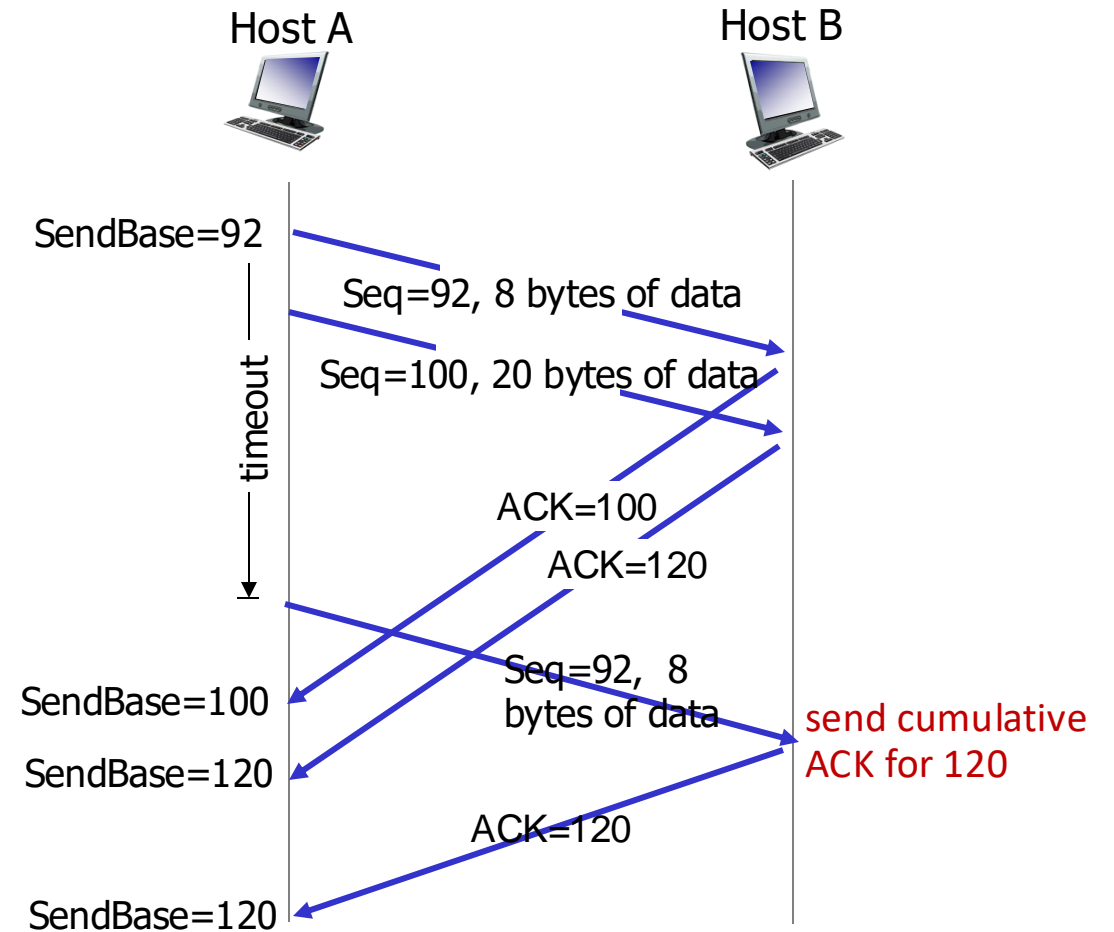  - start timer if there are still unACKed segments
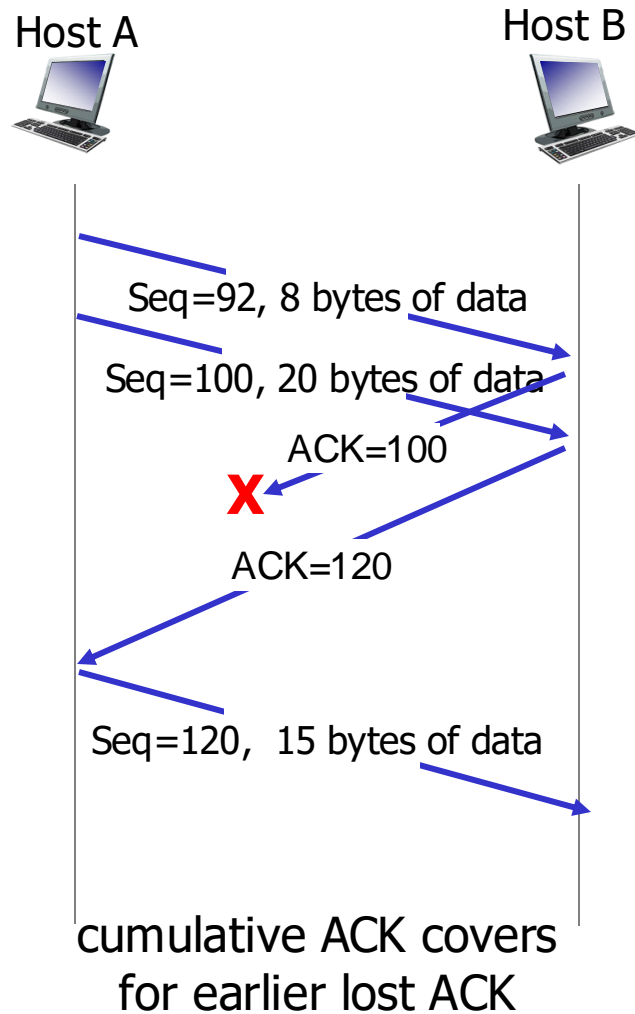
# TCP Receiver: ACK generation [RFC 5681]

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

# TCP: retransmission scenarios



lost ACK scenario

premature timeout

# TCP: retransmission scenarios



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 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

timeout

Seq=92, 8 bytes of data
Seq=100, 20 bytes of data    X

ACK=100

ACK=100

ACK=100

ACK=100

Seq=100, 20 bytes of data