

Overview of last class

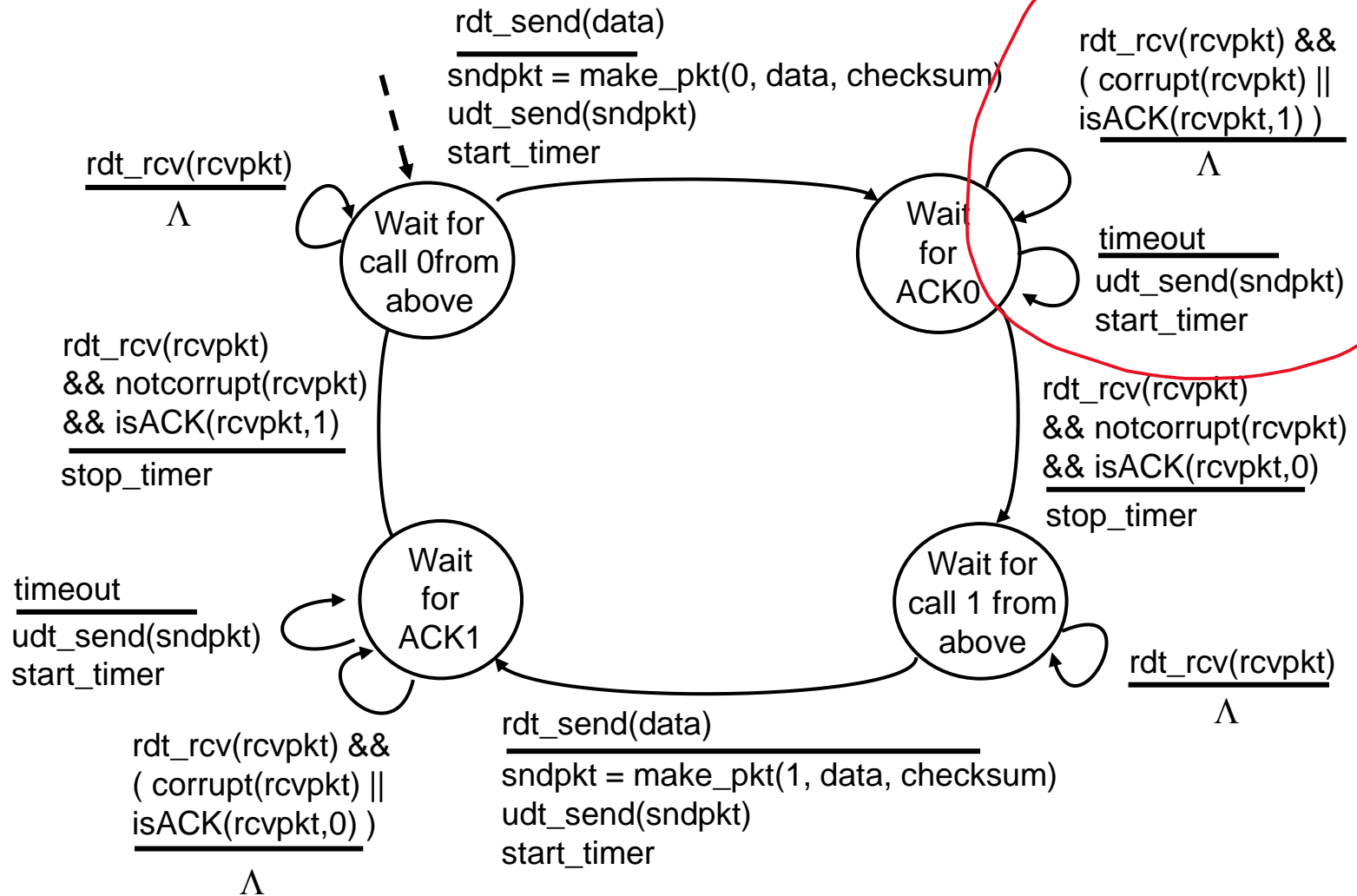
rdt2.2: a NAK-free protocol

- Receiver *explicitly* include seq # of pkt being ACKed
- Duplicate ACK at sender results in same action as NAK: *retransmit current pkt*
- No need to define two types of message; support general scenario
- ACK/NAK handles stop-and-wait is fine, difficult for pipeline case
 - Multiple new packets in transmission
 - Which one to be acknowledged?

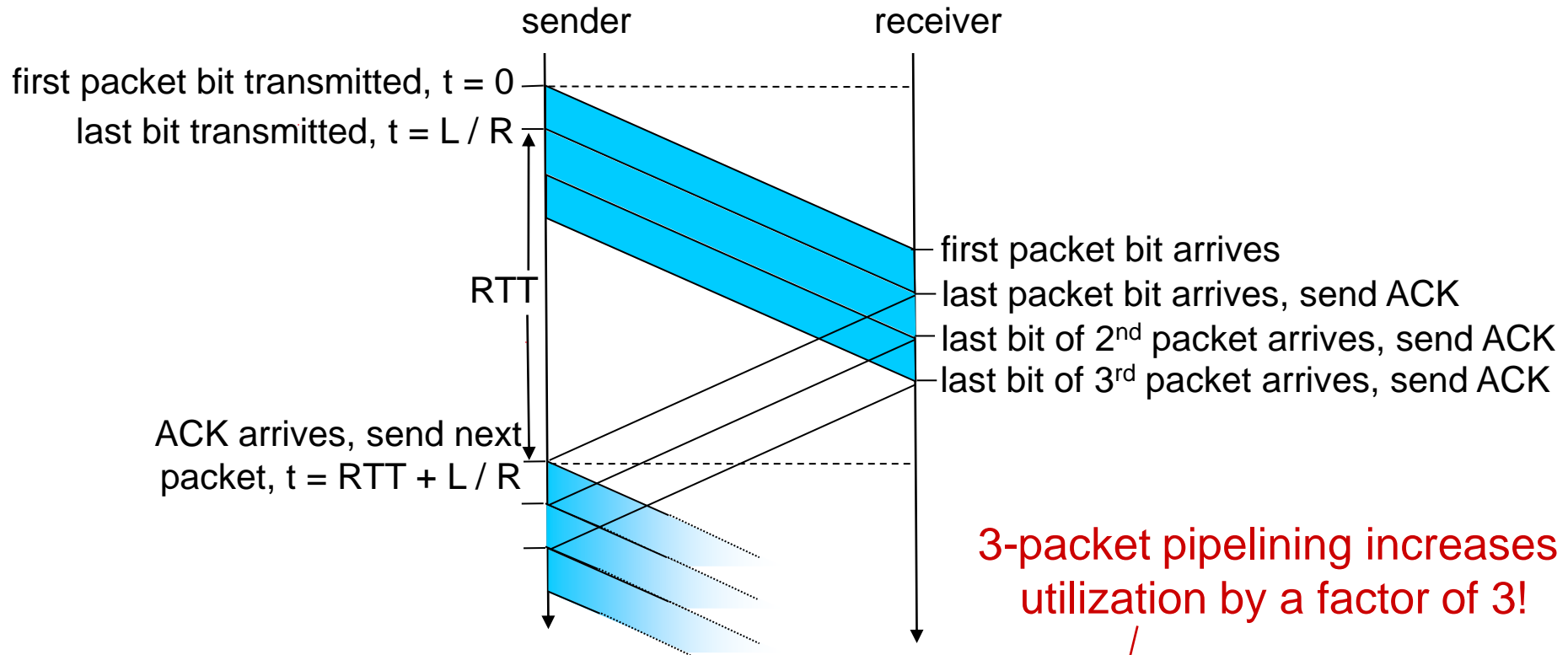
mechanisms for reliable commun. protocol design

- o checksum
- o feedback : (Ack + seq. number)
- o retransmission
- o seq. number
- o Timeout

rdt3.0 sender



Pipelining: increased utilization

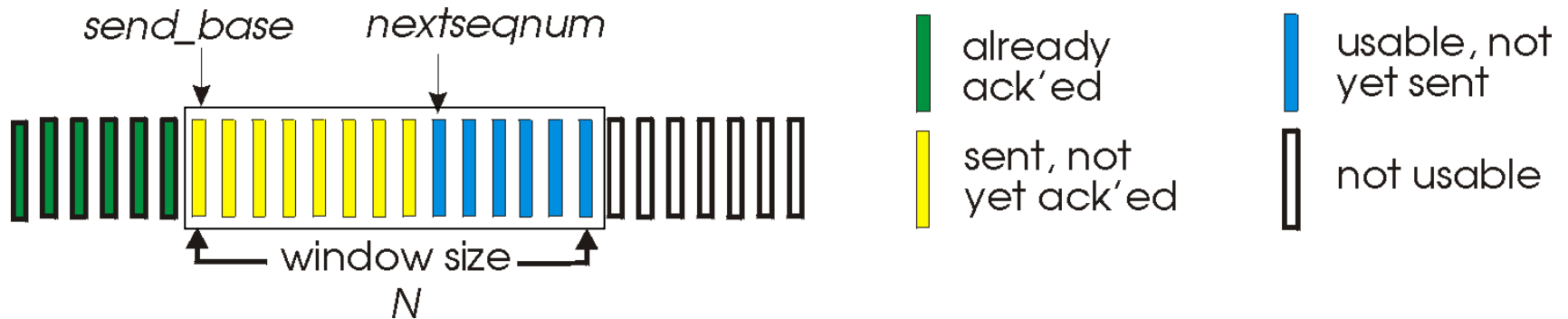


3-packet pipelining increases utilization by a factor of 3!

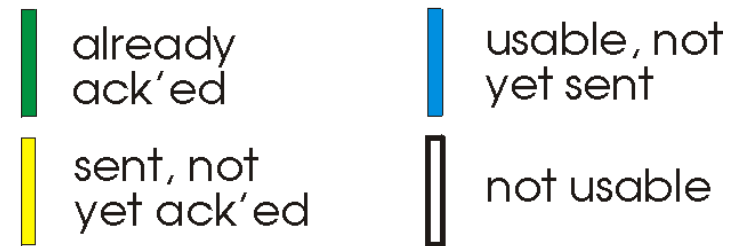
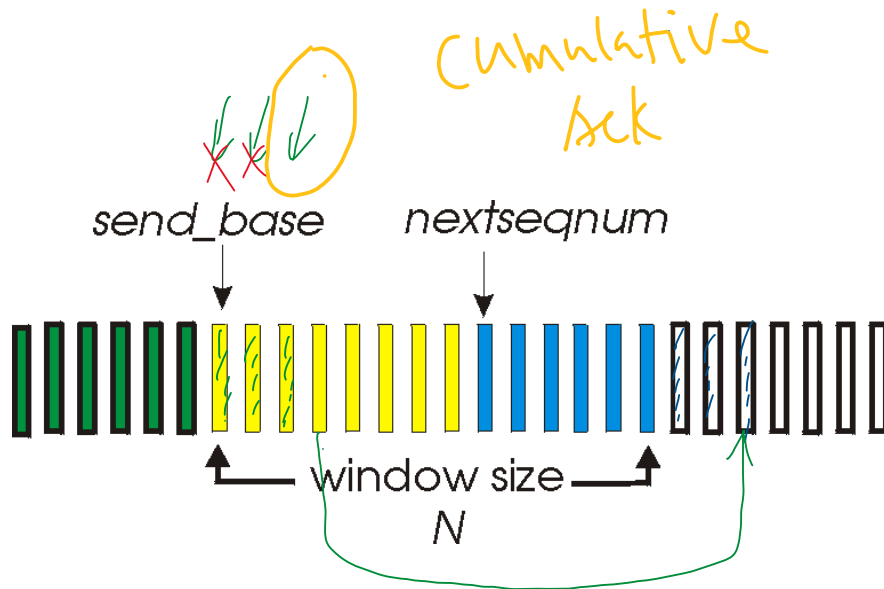
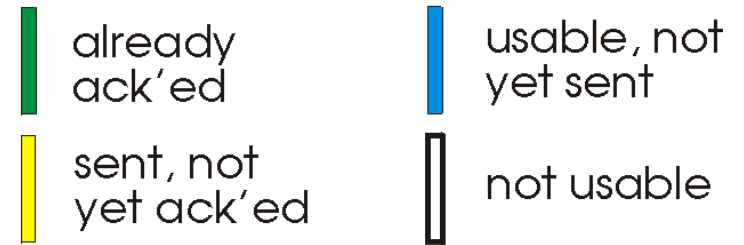
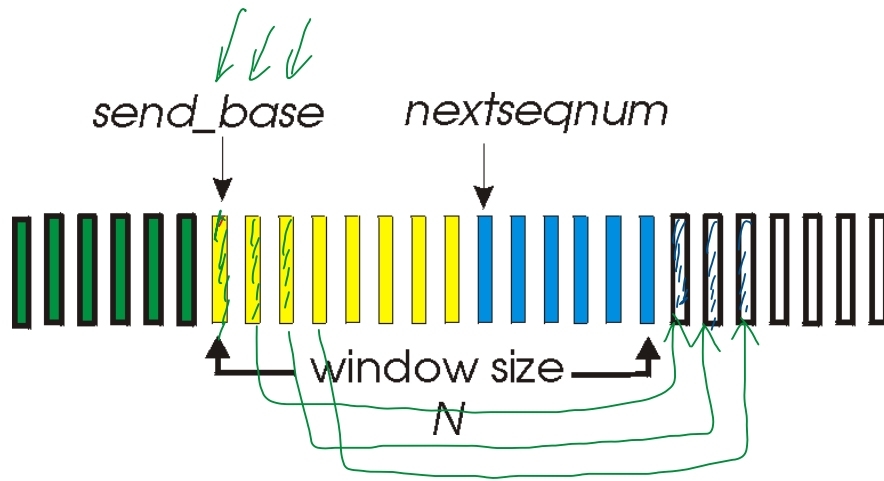
$$U_{\text{sender}} = \frac{3L / R}{RTT + L / R} = \frac{.0024}{30.008} = 0.00081$$

Go-Back-N: sender

- k-bit seq # in pkt header
- “window” of up to N, consecutive unack’ed pkts allowed



- ACK(n): ACKs all pkts up to, including seq # n - “*cumulative ACK*”
 - may receive duplicate ACKs (see receiver)
- timer for oldest in-flight pkt
- *timeout(n)*: retransmit packet n and all higher seq # pkts in window



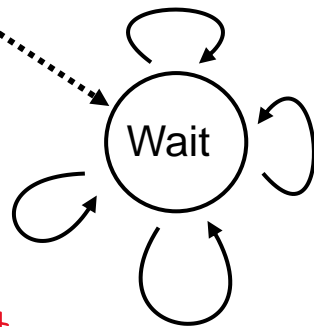
GBN: sender extended FSM

rdt_send(data)

```

if (nextseqnum < base+N) {
    sndpkt[nextseqnum] = make_pkt(nextseqnum,data,chksum)
    udt_send(sndpkt[nextseqnum])
    if (base == nextseqnum)
        start_timer
    nextseqnum++
}
else
    refuse_data(data)
    
```

Λ
base=1
nextseqnum=1



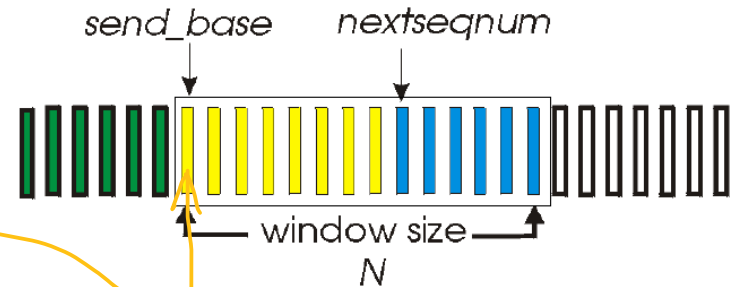
rdt_rcv(rcvpkt)
&& corrupt(rcvpkt)

Retransmission by timeout

*if base ≤
 getacknum(rcvpkt)*

rdt_rcv(rcvpkt) &&
notcorrupt(rcvpkt)
base = getacknum(rcvpkt)+1
 If (base == nextseqnum)
 stop_timer
 else
 start_timer

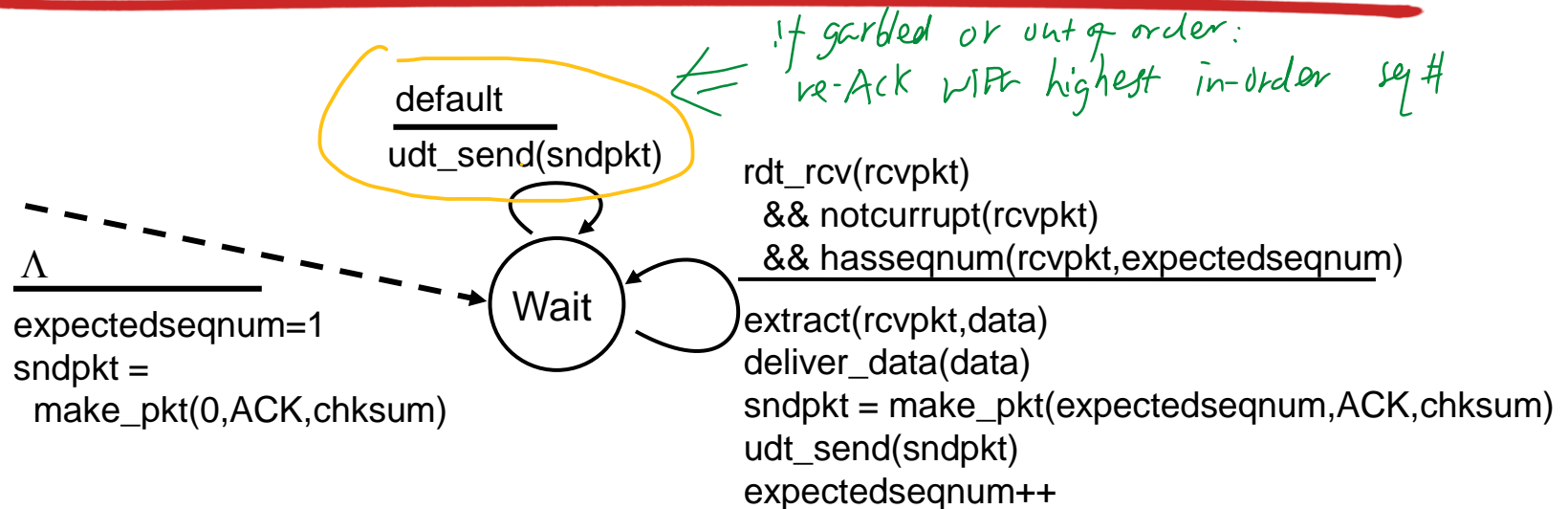
timeout
start_timer
udt_send(sndpkt[base])
udt_send(sndpkt[base+1])
 ...
udt_send(sndpkt[nextseqnum-1])



Timeout timer

All yellow ones

GBN: receiver extended FSM



ACK-only: always send ACK for correctly-received pkt with highest *in-order* seq #

- may generate duplicate ACKs
- need only remember **expectedseqnum**
- out-of-order pkt:
 - discard (don't buffer): *no receiver buffering!*
 - re-ACK pkt with highest in-order seq #

GBN in action

sender window (N=4)

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

↓ sliding

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

sender

send pkt0

send pkt1

send pkt2

send pkt3

(wait)

rcv ack0, send pkt4

rcv ack1, send pkt5

ignore duplicate ACK



pkt 2 timeout

send pkt2

send pkt3

send pkt4

send pkt5

receiver

receive pkt0, send ack0

receive pkt1, send ack1

expect 2

receive pkt3, discard,
(re)send ack1

receive pkt4, discard,
(re)send ack1

receive pkt5, discard,
(re)send ack1

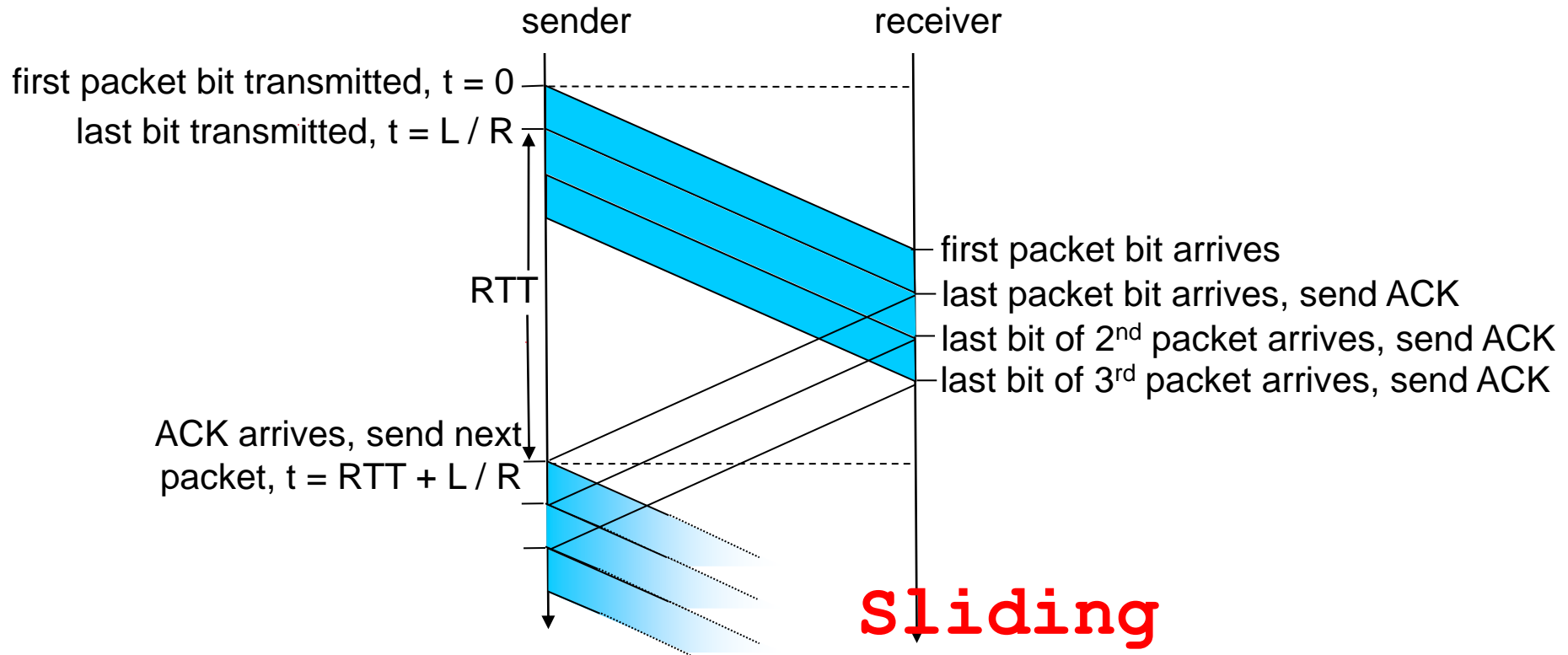
rcv pkt2, deliver, send ack2

rcv pkt3, deliver, send ack3

rcv pkt4, deliver, send ack4

rcv pkt5, deliver, send ack5

Window size: sending rate control



$$U_{\text{sender}} = \frac{NL / R}{RTT + L / R}$$

Sliding

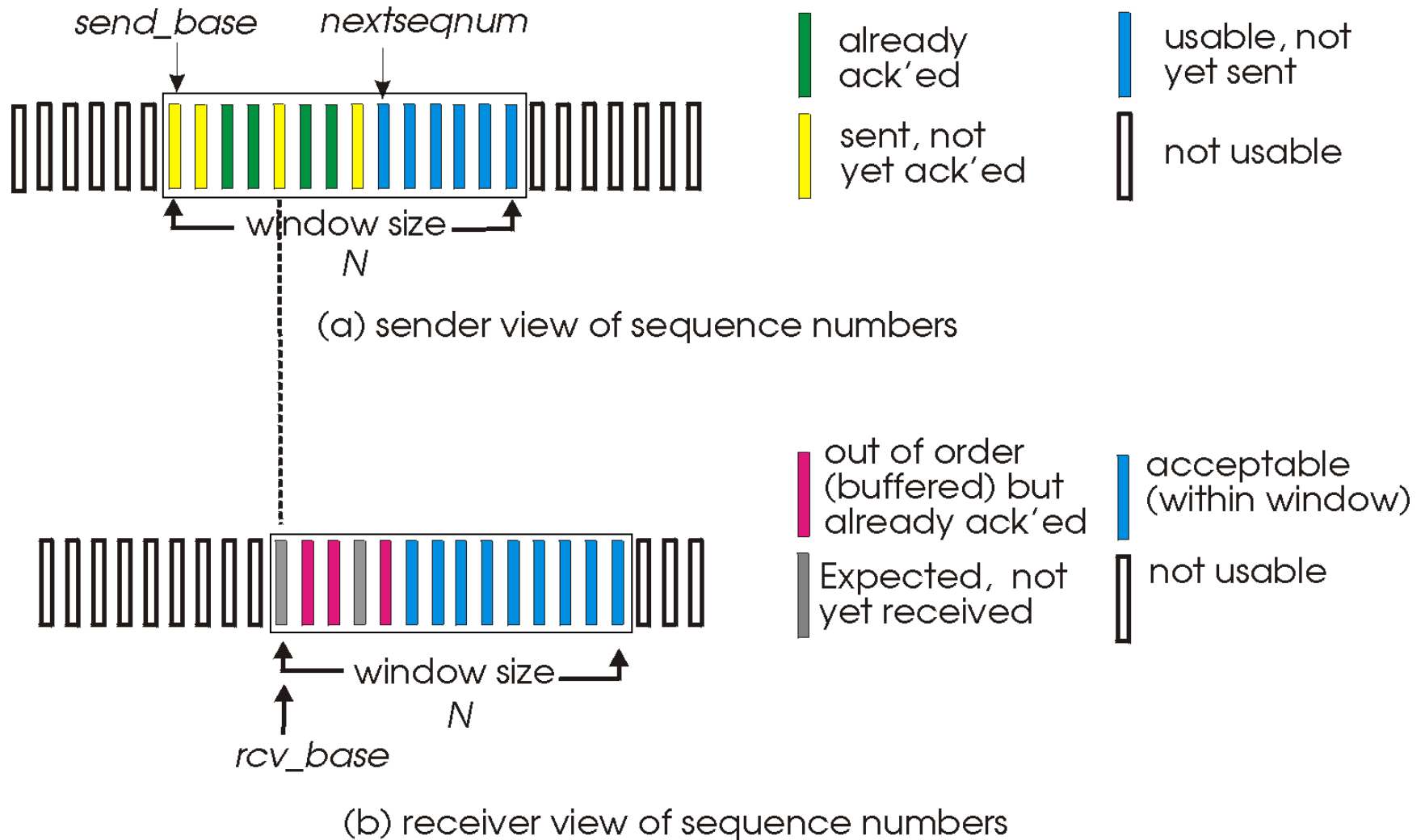
- Window of size N : up to N packets in air; states to be determined
- For the data well received in sequence, considered as processed, so move on

Class Today

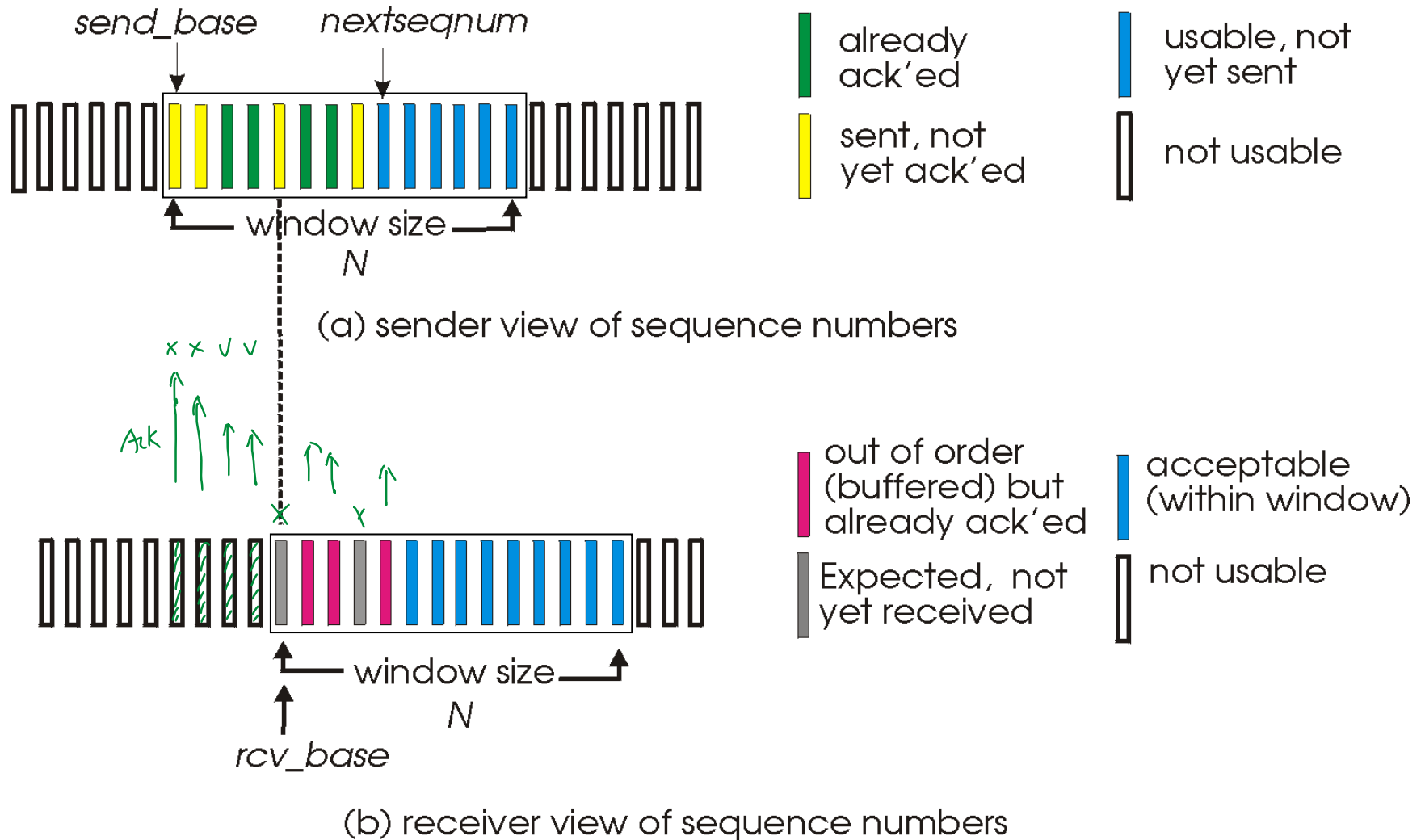
Selective repeat

- receiver *individually* acknowledges all correctly received pkts
 - buffers pkts, as needed, for eventual in-order delivery to upper layer
- sender only resends pkts for which ACK not received
 - sender timer for each unACKed pkt
- sender window
 - N consecutive seq #'s
 - limits seq #s of sent, unACKed pkts

Selective repeat: sender, receiver windows



Selective repeat: sender, receiver windows



Selective repeat

sender

data from above:

- if next available seq # in window, send pkt

timeout(n):

- resend pkt n, restart timer

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

- mark pkt n as received
- if n smallest unACKed pkt, advance window base to next unACKed seq #

receiver

pkt n in [rcvbase, rcvbase+N-1]

- send ACK(n)
- out-of-order: buffer
- in-order: deliver (also deliver buffered, in-order pkts), advance window to next not-yet-received pkt

pkt n in [rcvbase-N, rcvbase-1]

- ACK(n)

otherwise:

- ignore

Selective repeat in action

sender window (N=4)

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

sender

send pkt0

send pkt1

send pkt2

send pkt3

(wait)

rcv ack0, send pkt4

rcv ack1, send pkt5

record ack3 arrived
& buffered



pkt 2 timeout

send pkt2

record ack4 arrived

record ack5 arrived

receiver

receive pkt0, send ack0

receive pkt1, send ack1

receive pkt3, buffer,
send ack3

receive pkt4, buffer,
send ack4

receive pkt5, buffer,
send ack5

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

Q: what happens when ack2 arrives?

Selective repeat: dilemma

00, 01, 10, 11 | 00, 01, 10, 11

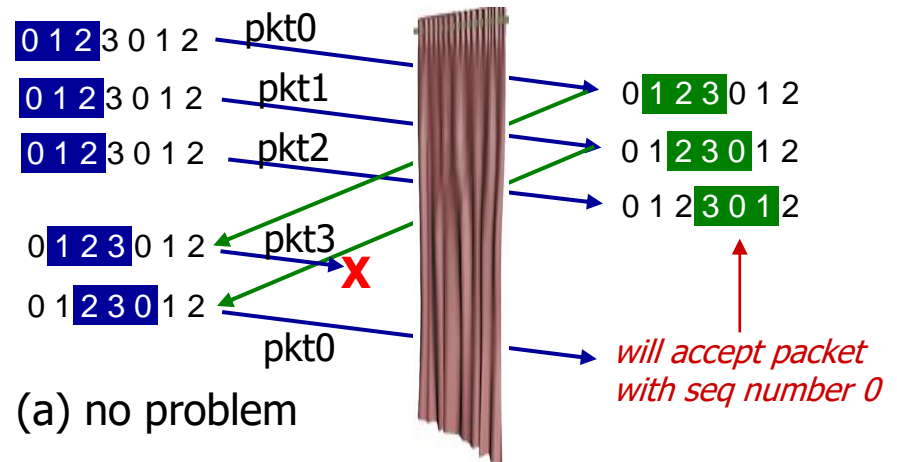
example:

- seq #'s: 0, 1, 2, 3
- window size=3
- receiver sees no difference in two scenarios!
- duplicate data accepted as new in (b)

Q: what relationship between seq # size and window size to avoid problem in (b)?

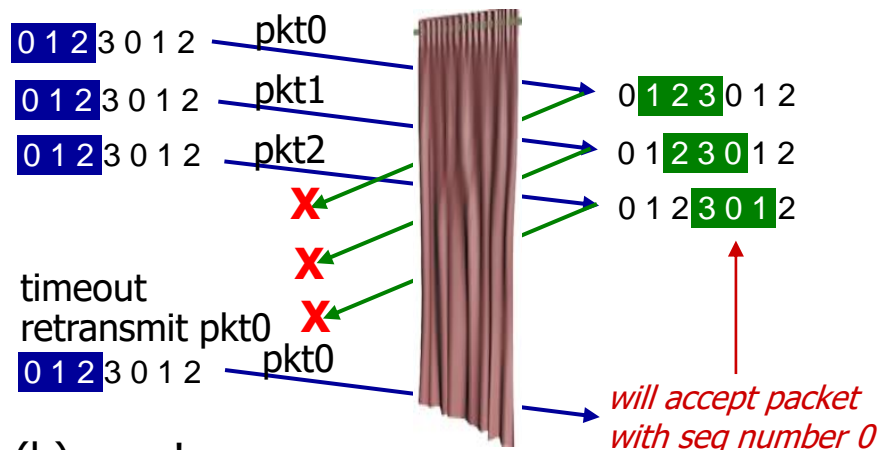
sender window
(after receipt)

receiver window
(after receipt)



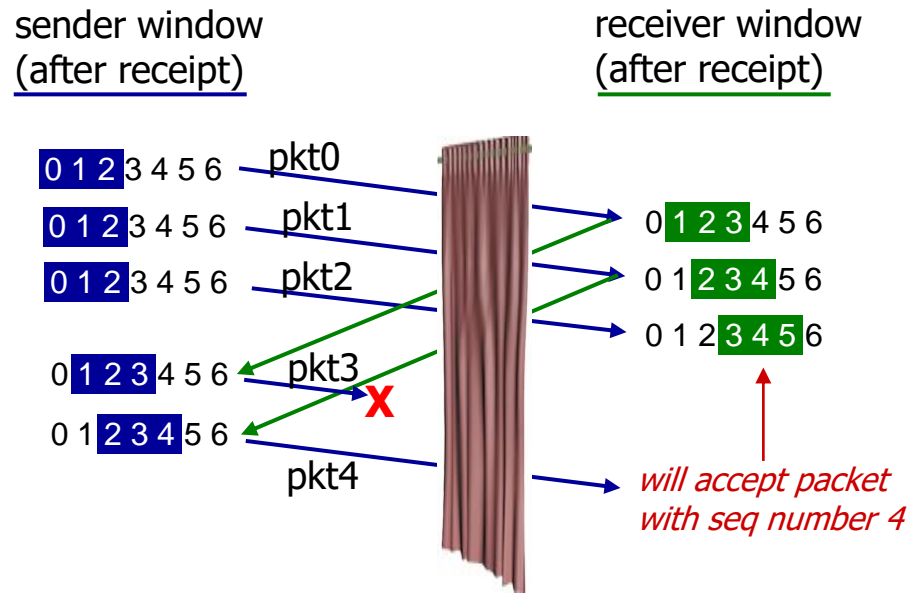
(a) no problem

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

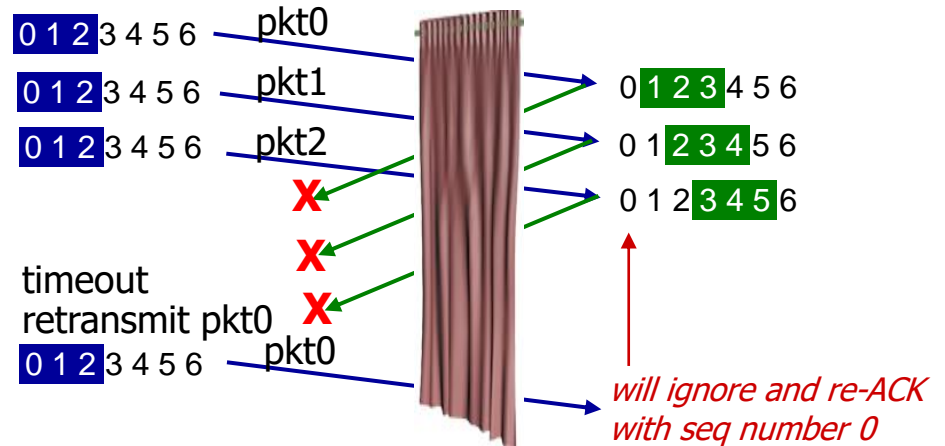


(b) oops!

Address the dilemma: with sequence number



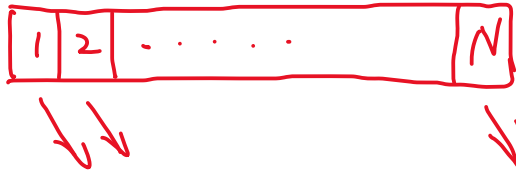
(a) Receive the new pkt4



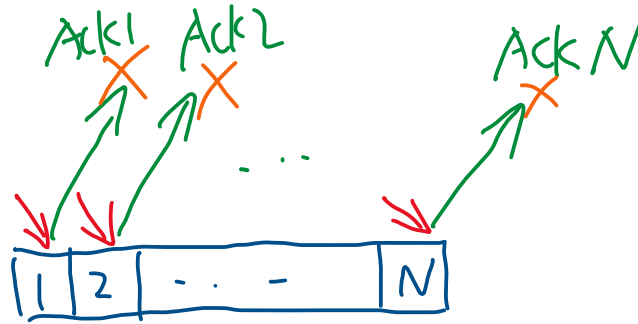
(b) Receive the retransmitted pkt 0

Sequence number size Analysis

sender

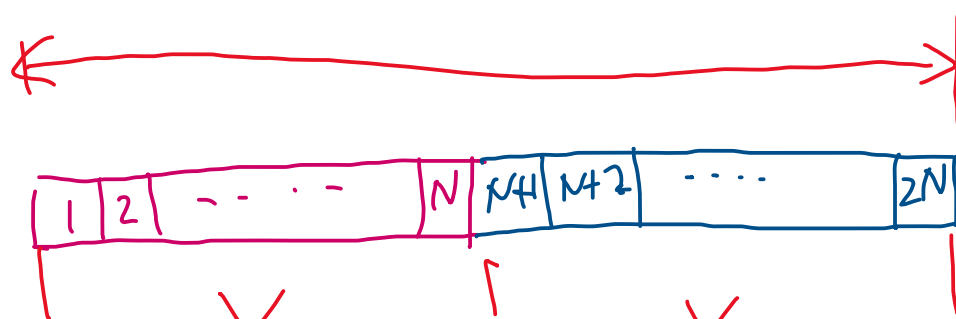


Receiver



seq # cover
range of

$\log_2(2N)$



\Rightarrow # of bits for seq #.

when out of order
data needs
to be buffered

Chapter 3 outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer

3.5 connection-oriented transport: TCP

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

3.6 principles of congestion control

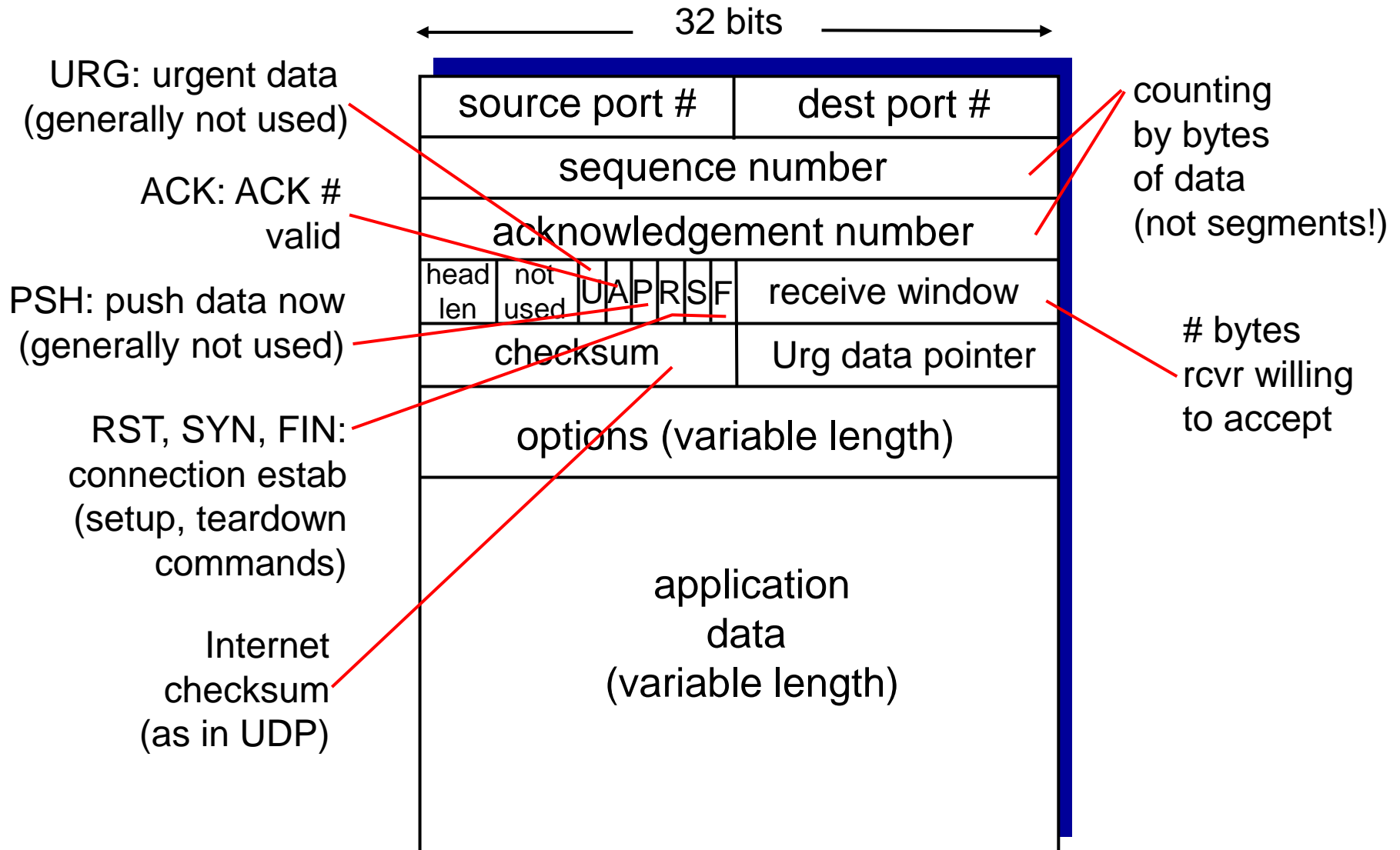
3.7 TCP congestion control

TCP: Overview

RFCs: 793, 1122, 1323, 2018, 2581

- **point-to-point:**
 - one sender, one receiver
- **reliable, in-order *byte stream*:**
 - no “message boundaries”
 - Stream of bytes instead of number of segments facilitates quickly putting data into the right position.
- **pipelined:**
 - TCP congestion and flow control set window size
- **full duplex data:**
 - bi-directional data flow in same connection
 - MSS: maximum segment size
- **connection-oriented:**
 - handshaking (exchange of control msgs) initializes sender, receiver state before data exchange
- **flow controlled:**
 - sender will not overwhelm receiver

TCP segment structure



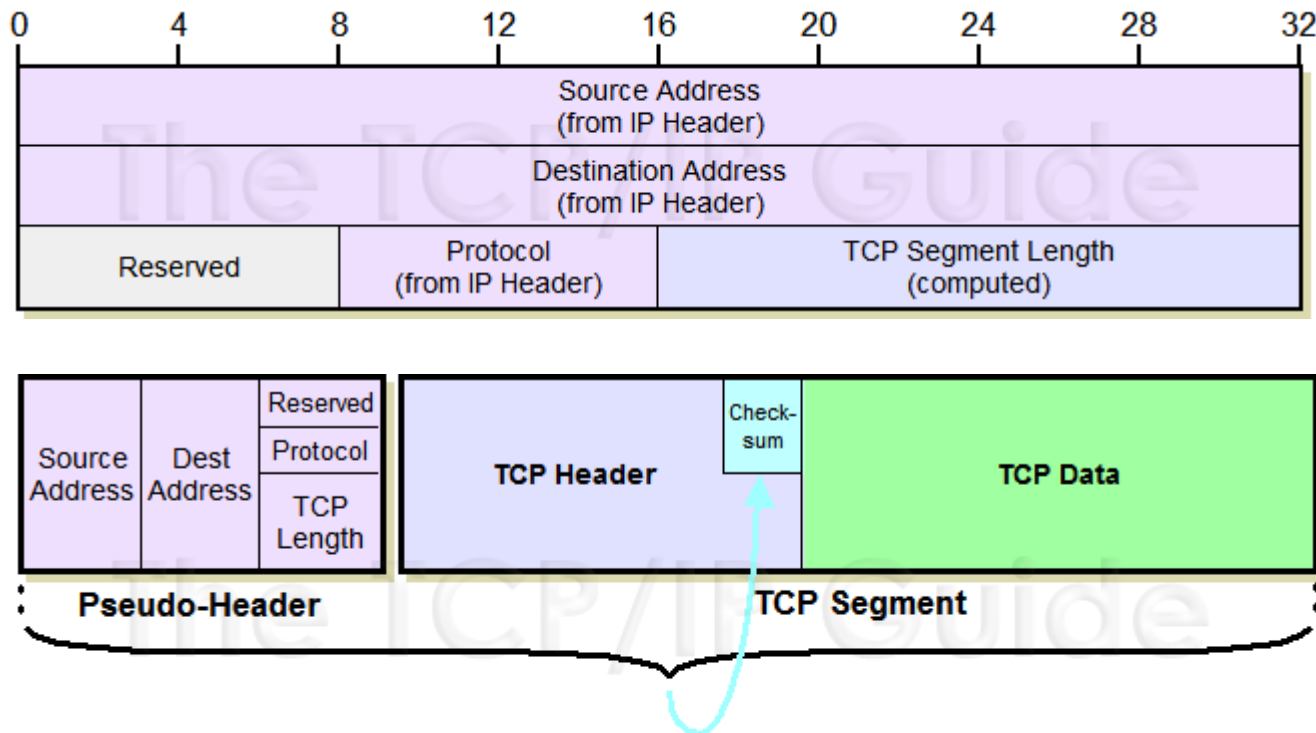
More on TCP segment

- “Options” field:
 - negotiate maximum segment size; scaling factor for high speed network; timestamp
- Flag bits:
 - PSH: indicate receiver to pass data to upper layer immediately
 - URG: Sender upper layer indicate data urgent
 - Urg data pointer: pointer for the urgent data

TCP Checksum

- To calculate the TCP segment header's Checksum field, the TCP pseudo header is first constructed and placed, logically, before the TCP segment. The checksum is then calculated over both the pseudo header and the TCP segment. The pseudo header is then discarded.

(http://www.tcpipguide.com/free/t_TCPChecksumCalculationandtheTCPPseudoHeader-2.htm)



Checksum Calculated Over Pseudo Header and TCP Segment

TCP seq. numbers, ACKs

sequence numbers:

- byte stream “number” of first byte in segment’s data

acknowledgements:

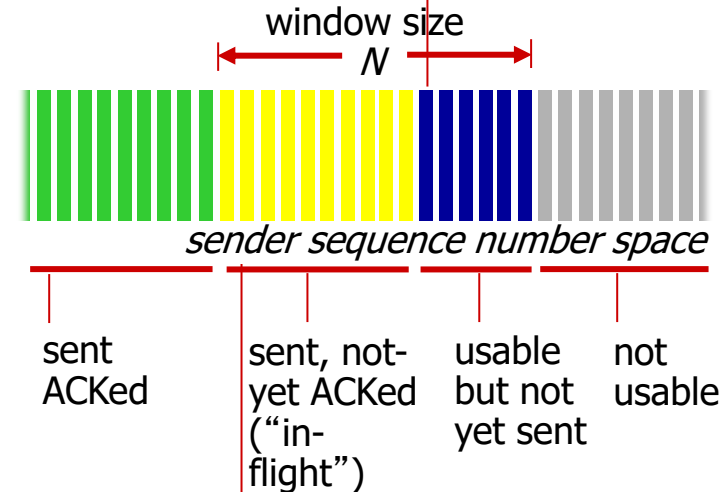
- seq # of next byte expected from other side
- cumulative ACK

Q: how receiver handles out-of-order segments

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

outgoing segment from sender

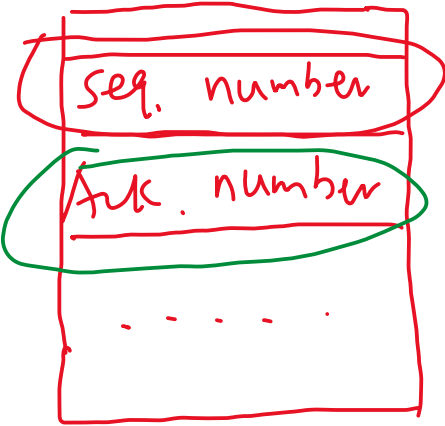
source port #		dest port #	
sequence number			
acknowledgement number			
			rwnd
checksum		urg pointer	



incoming segment to sender

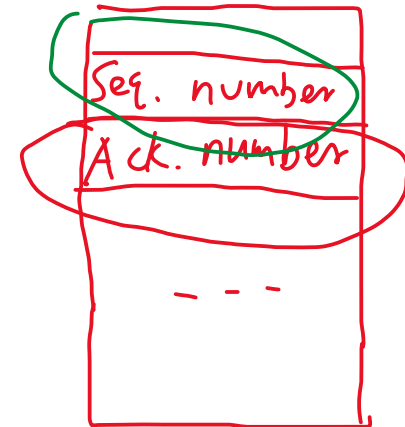
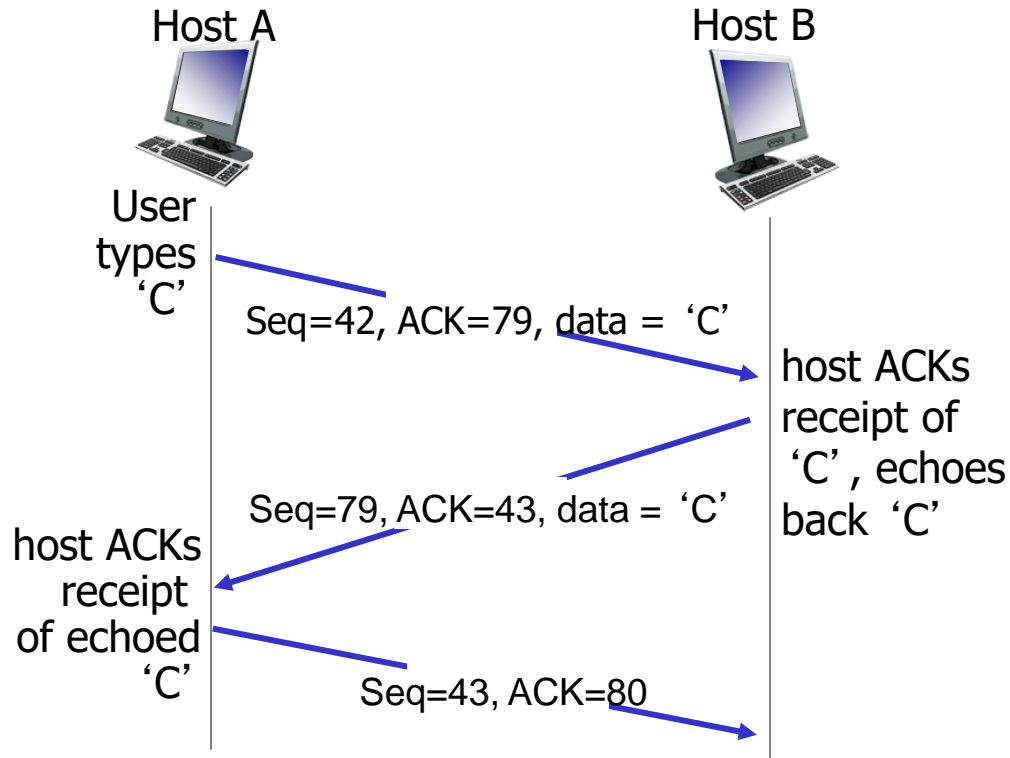
source port #		dest port #	
sequence number			
acknowledgement number			
		A	rwnd
checksum		urg pointer	

TCP seq. numbers, ACKs



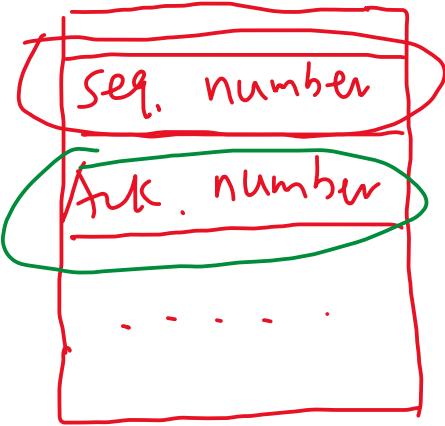
A → B

B → A



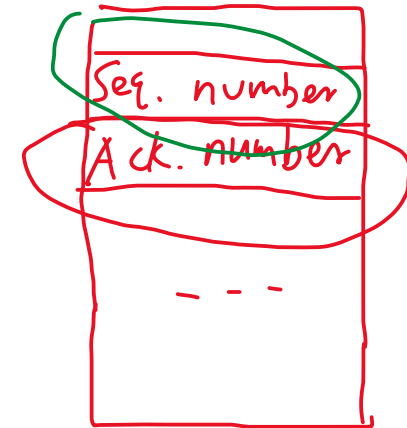
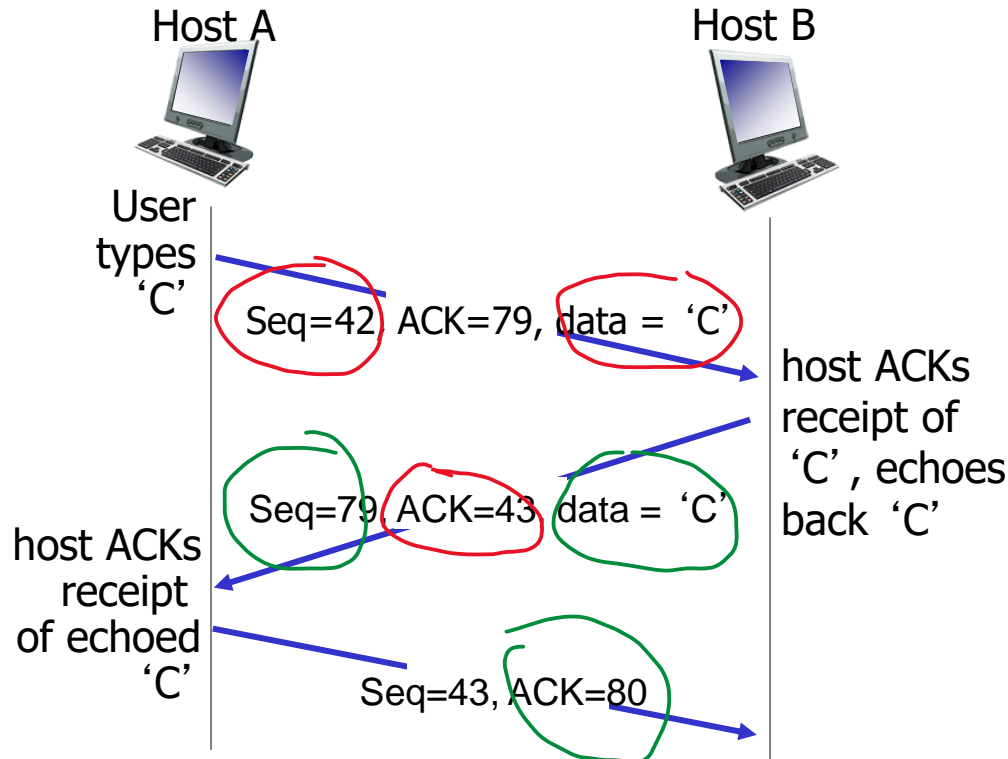
simple telnet scenario

TCP seq. numbers, ACKs



A → B

B → A



simple telnet scenario

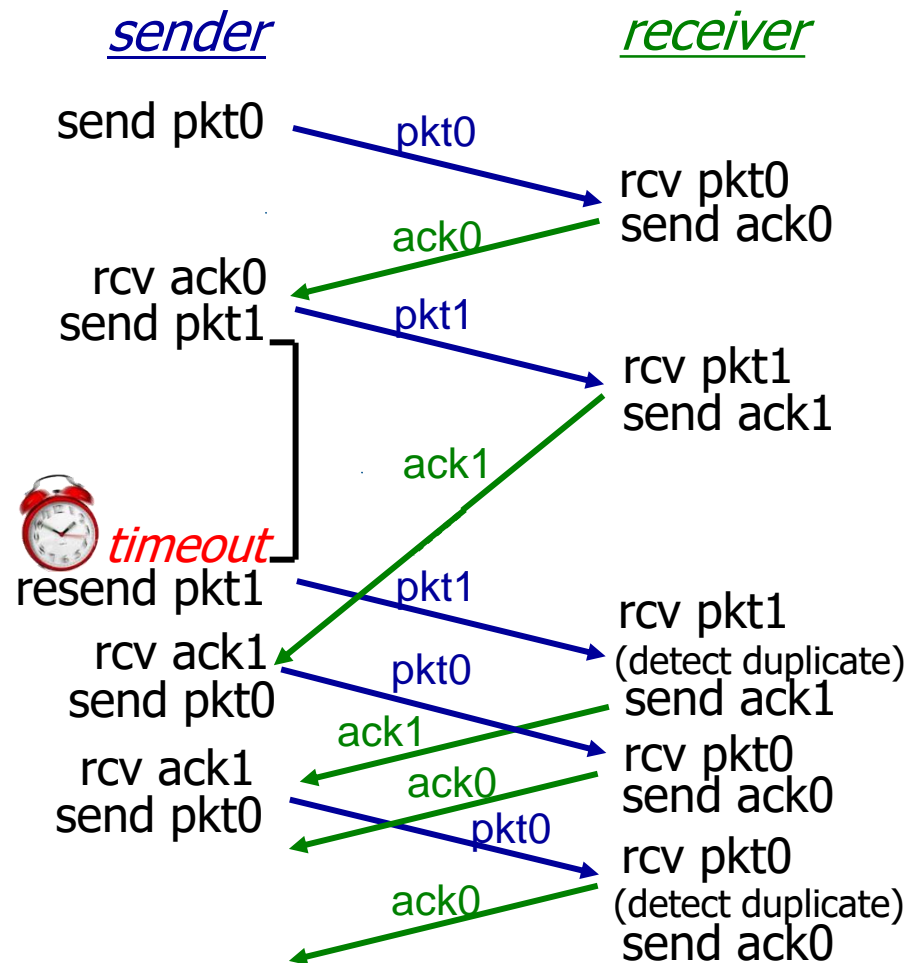
TCP round trip time, timeout

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**

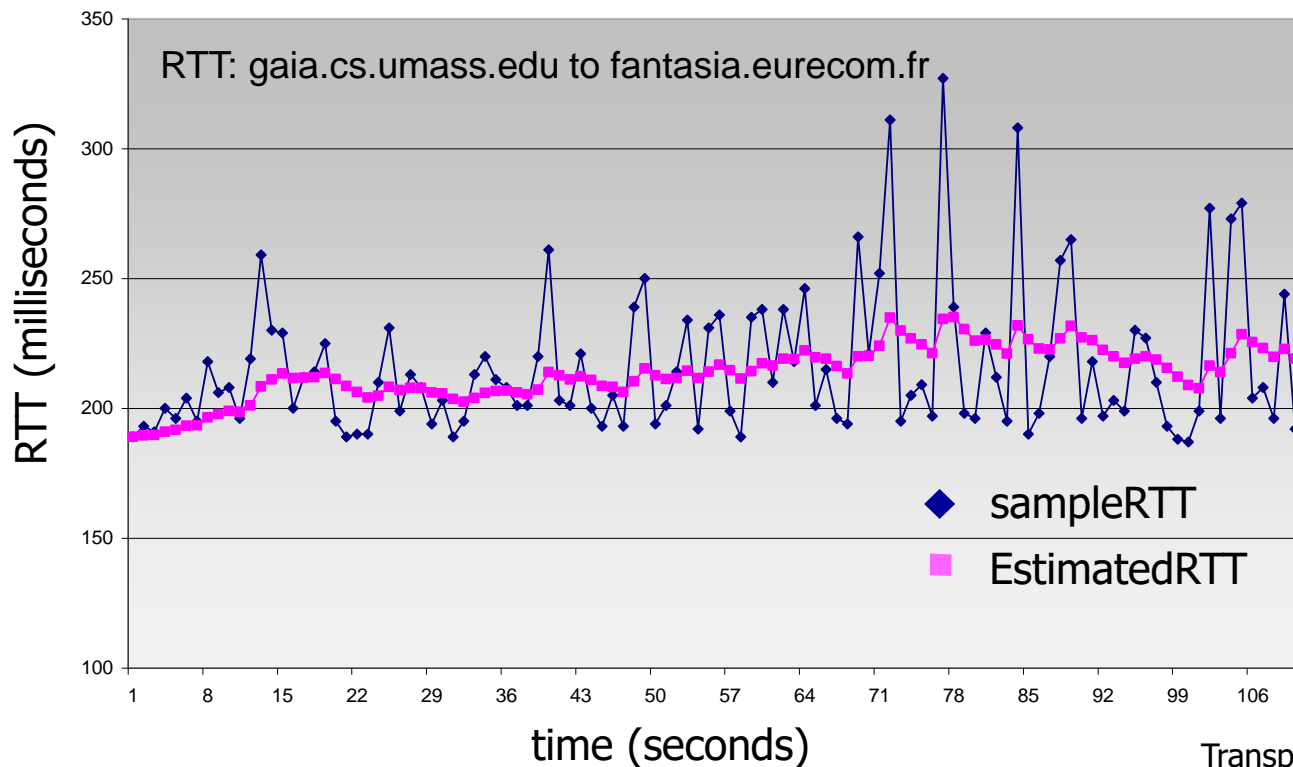


(d) premature timeout/ delayed ACK

TCP round trip time, timeout

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

- exponential weighted moving average
- influence of past sample decreases exponentially fast
- typical value: $\alpha = 0.125$



Exponential Weighted Moving Average

$$\text{EstimateRTT} = \text{SampleRTT}_1$$

$i=1$
while (1)

$i++$;

$$\text{EstimateRTT} = (1-\alpha) \text{EstimateRTT} + \alpha \text{SampleRTT}_i$$

end

$$\text{EstimateRTT} = \text{SampleRTT}_1$$

$$\text{EstimateRTT} = (1-\alpha) \text{SampleRTT}_1 + \alpha \text{SampleRTT}_2$$

$$\begin{aligned} \text{EstimateRTT} &= (1-\alpha) \left[(1-\alpha) \text{SampleRTT}_1 + \alpha \text{SampleRTT}_2 \right] + \alpha \text{SampleRTT}_3 \\ &= (1-\alpha)^2 \text{SampleRTT}_1 + (1-\alpha)\alpha \text{SampleRTT}_2 + \alpha \text{SampleRTT}_3 \end{aligned}$$

- . . . -

$$\begin{aligned} &= (1-\alpha)^{n-1} \text{SampleRTT}_1 + (1-\alpha)^{n-2} \alpha \text{SampleRTT}_2 + \dots + (1-\alpha) \alpha \text{SampleRTT}_{n-1} \\ &\quad + \alpha \text{SampleRTT}_n \end{aligned}$$

TCP round trip time, timeout

- **timeout interval:** **EstimatedRTT** plus “safety margin”
 - large variation in **EstimatedRTT** -> larger safety margin
- estimate SampleRTT deviation from EstimatedRTT:

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

(typically, $\beta = 0.25$)

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



↑
estimated RTT

↑
“safety margin”