

# Chapter 3

## Transport Layer

Reti di Elaboratori

Corso di Laurea in Informatica

Università degli Studi di Roma "La Sapienza"

Canale A-L

Prof.ssa Chiara Petrioli

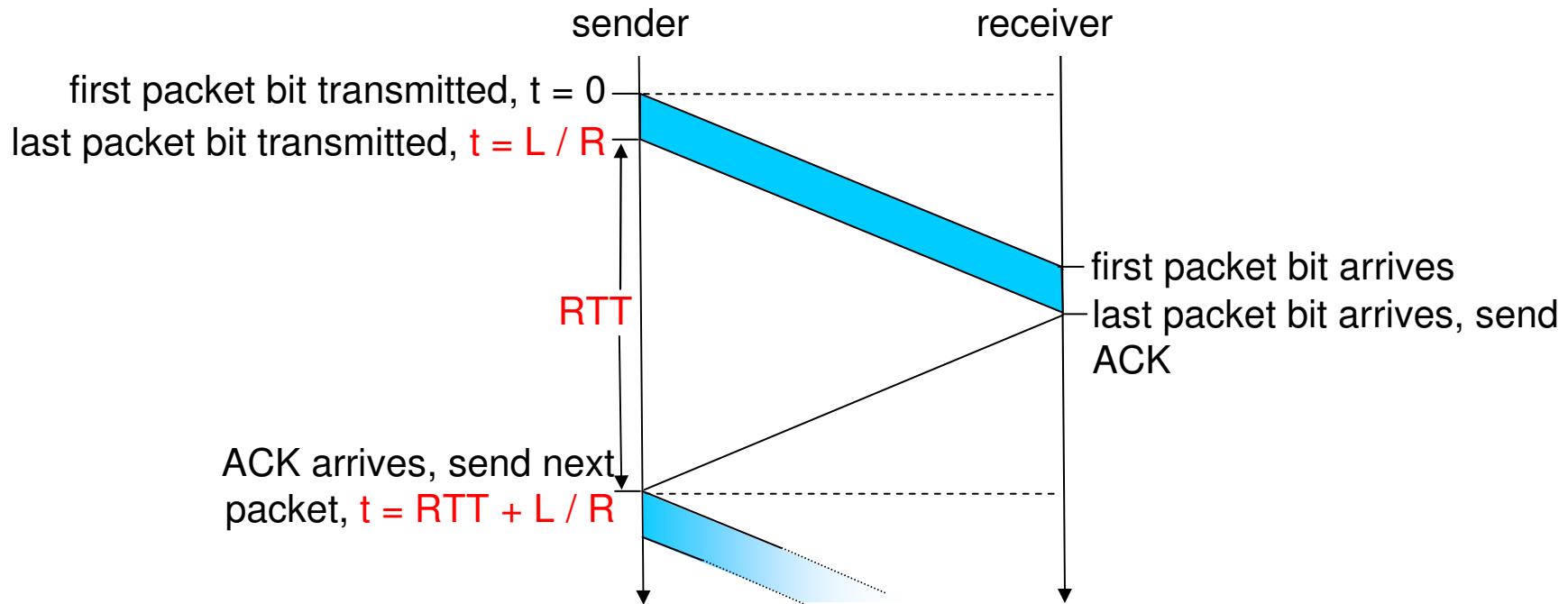
Parte di queste slide sono state prese dal materiale associato al libro  
*Computer Networking: A Top Down Approach*, 5th edition.

All material copyright 1996-2009

J.F Kurose and K.W. Ross, All Rights Reserved

Thanks also to Antonio Capone, Politecnico di Milano, Giuseppe Bianchi and  
Francesco LoPresti, Un. di Roma Tor Vergata

# rdt3.0: stop-and-wait operation

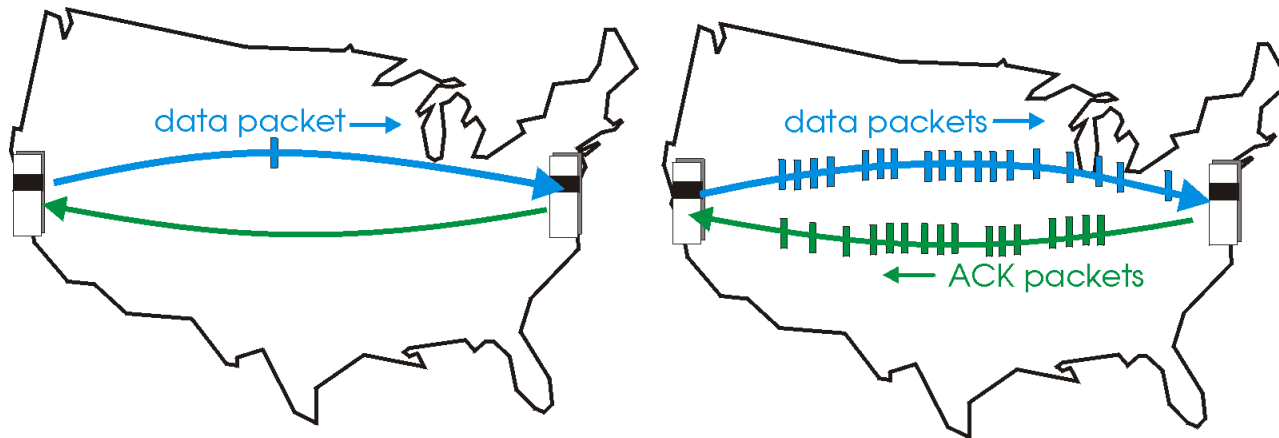


$$U_{\text{sender}} = \frac{L/R}{RTT + L/R} = \frac{.008}{30.008} = 0.00027$$

# Pipelined protocols

**Pipelining:** sender allows multiple, "in-flight", yet-to-be-acknowledged pkts

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

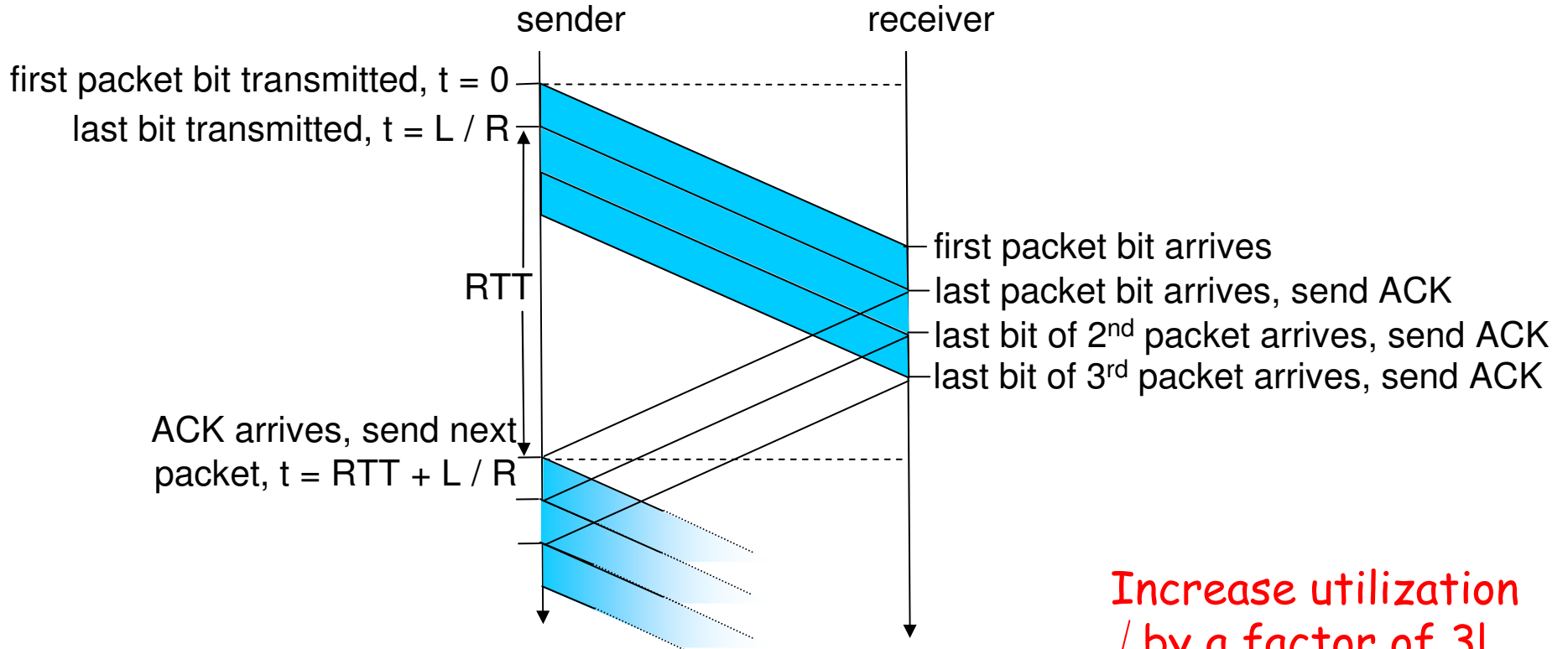


(a) a stop-and-wait protocol in operation

(b) a pipelined protocol in operation

- r Two generic forms of pipelined protocols: *go-Back-N*, *selective repeat*

# Pipelining: increased utilization



Increase utilization  
by a factor of 3!

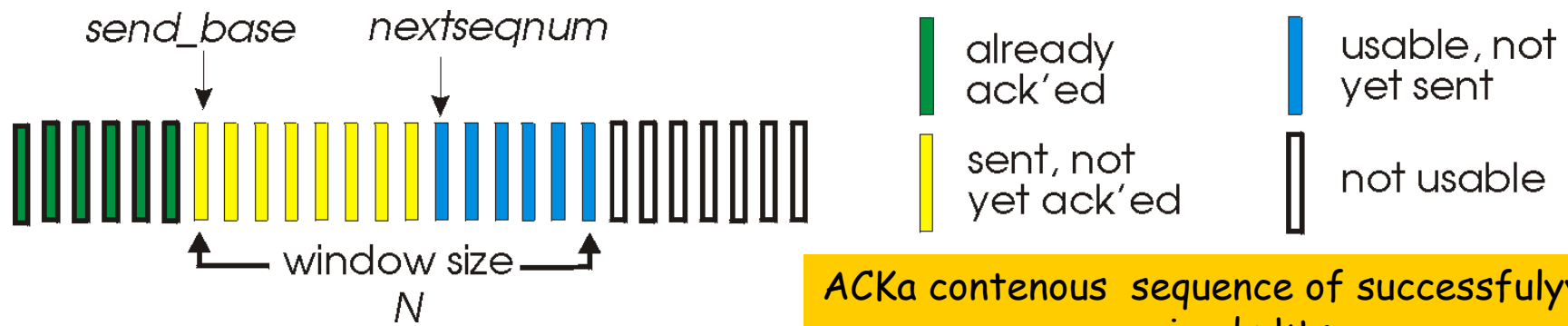
$$U_{\text{sender}} = \frac{3 * L / R}{RTT + L / R} = \frac{.024}{30.008} = 0.0008$$

# Go-Back-N

Diverso rispetto a Stop and Wait  
Q: Perché?

## Sender:

- r k-bit seq # in pkt header
- r "window" of up to N, consecutive unack'ed pkts allowed



- r ACK(n): ACKs all pkts up to, including seq # n - "cumulative ACK"
  - m may deceive duplicate ACKs (see receiver)
- r timer for each in-flight pkt
- r *timeout(n)*: retransmit pkt n and all higher seq # pkts in window

# GBN: sender extended FSM

rdt\_send(data)

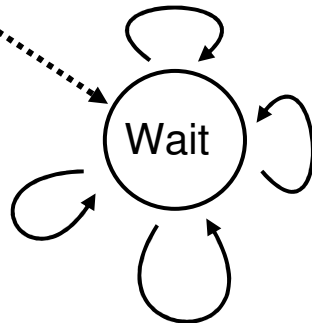
```

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)
    
```

Un solo timer in questa implementazione

$\Lambda$   
base=1  
nextseqnum=1

rdt\_rcv(rcvpkt)  
&& corrupt(rcvpkt)



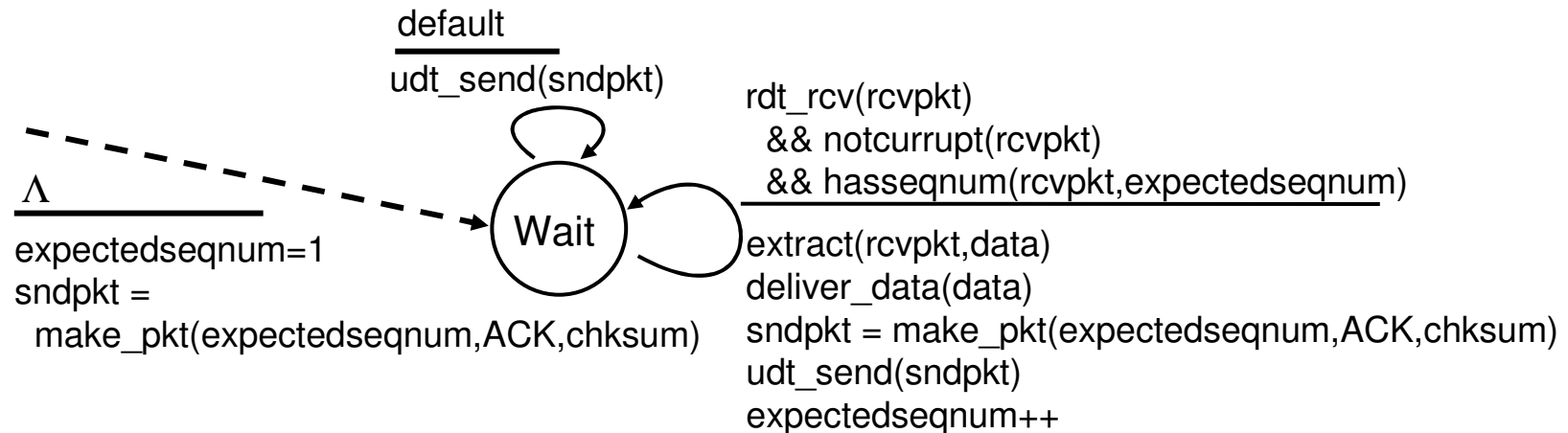
timeout  
start\_timer  
udt\_send(sndpkt[base])  
udt\_send(sndpkt[base+1])  
...  
udt\_send(sndpkt[nextseqnum-1])

All retransmitted

Gestione timer

rdt\_rcv(rcvpkt) && notcorrupt(rcvpkt)  
base = getacknum(rcvpkt)+1  
If (base == nextseqnum)  
stop\_timer  
else  
start\_timer

# GBN: receiver extended FSM



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

m may generate duplicate ACKs

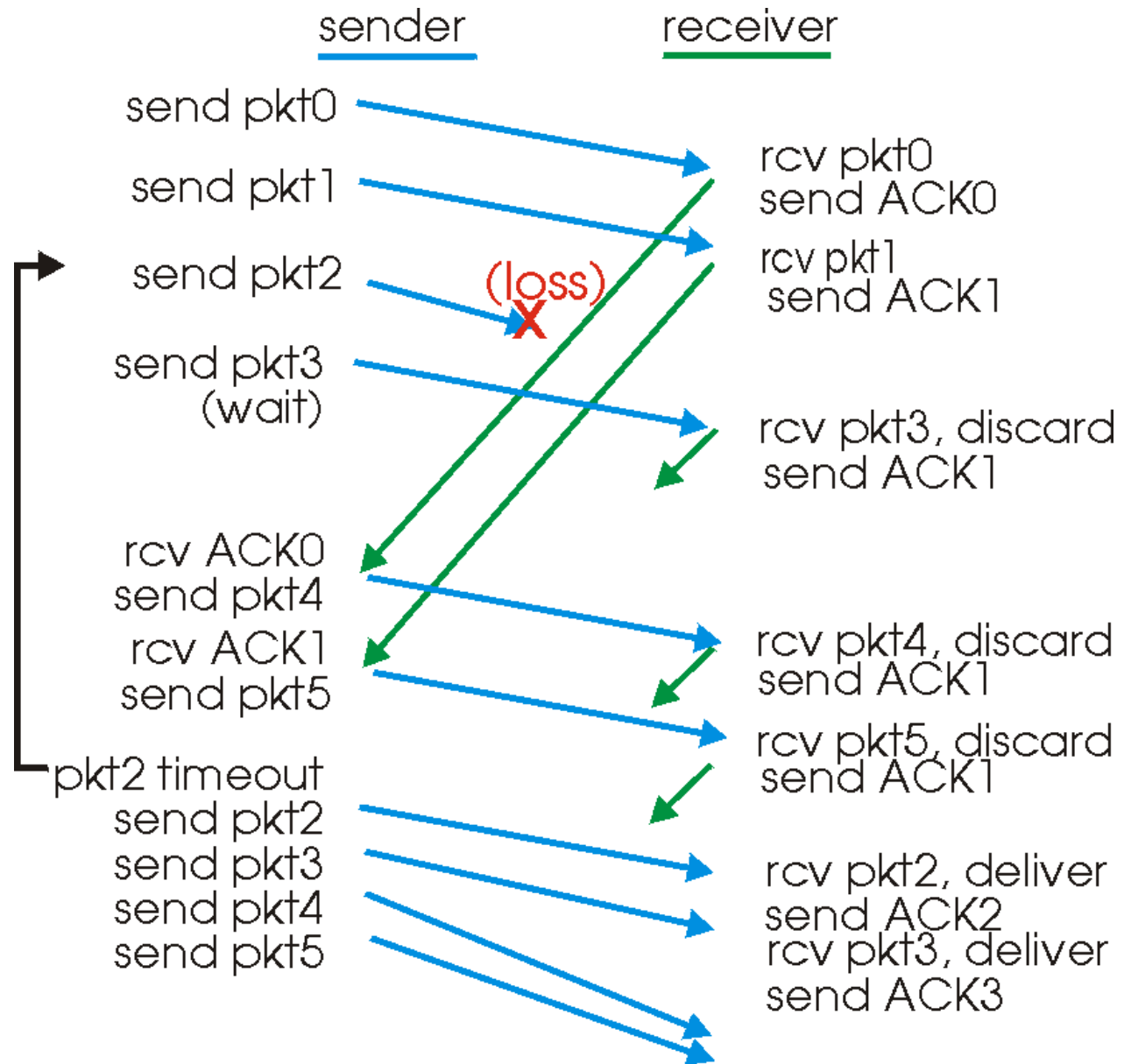
m need only remember **expectedseqnum**

r out-of-order pkt:

m discard (don't buffer) -> **no receiver buffering!**

m Re-ACK pkt with highest in-order seq #

# GBN in action



WINDOW SIZE ==4



# A few questions...

- r Why limiting the window size?
  - m max window size to improve performance related to RTT
  - m window size powerful tool to control data rate (important for flow control, congestion control)
  - m related to window size field length

# Selective Repeat

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

# Selective repeat

—sender—

**data from above :**

r if next available seq # in window, send pkt

**timeout(n):** Each packet has one Logical timer

r resend pkt n, restart timer

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

r mark pkt n as received

r if n smallest unACKed pkt, advance window base to next unACKed seq #

—receiver—

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

r send ACK(n)

r out-of-order: buffer

r in-order: deliver (also deliver buffered, in-order pkts), advance window to next not-yet-received pkt

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

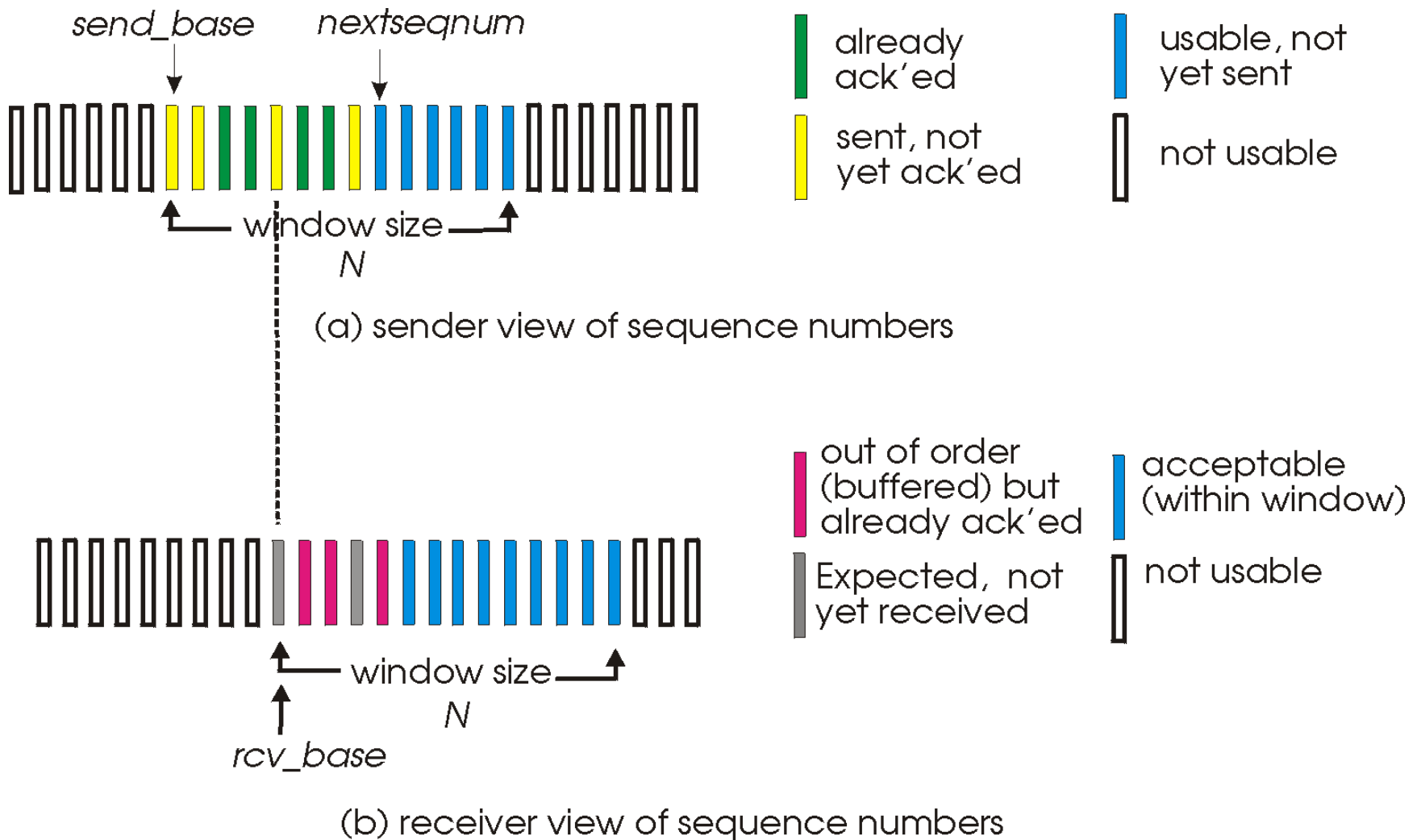
r ACK(n)

**otherwise:**

r ignore

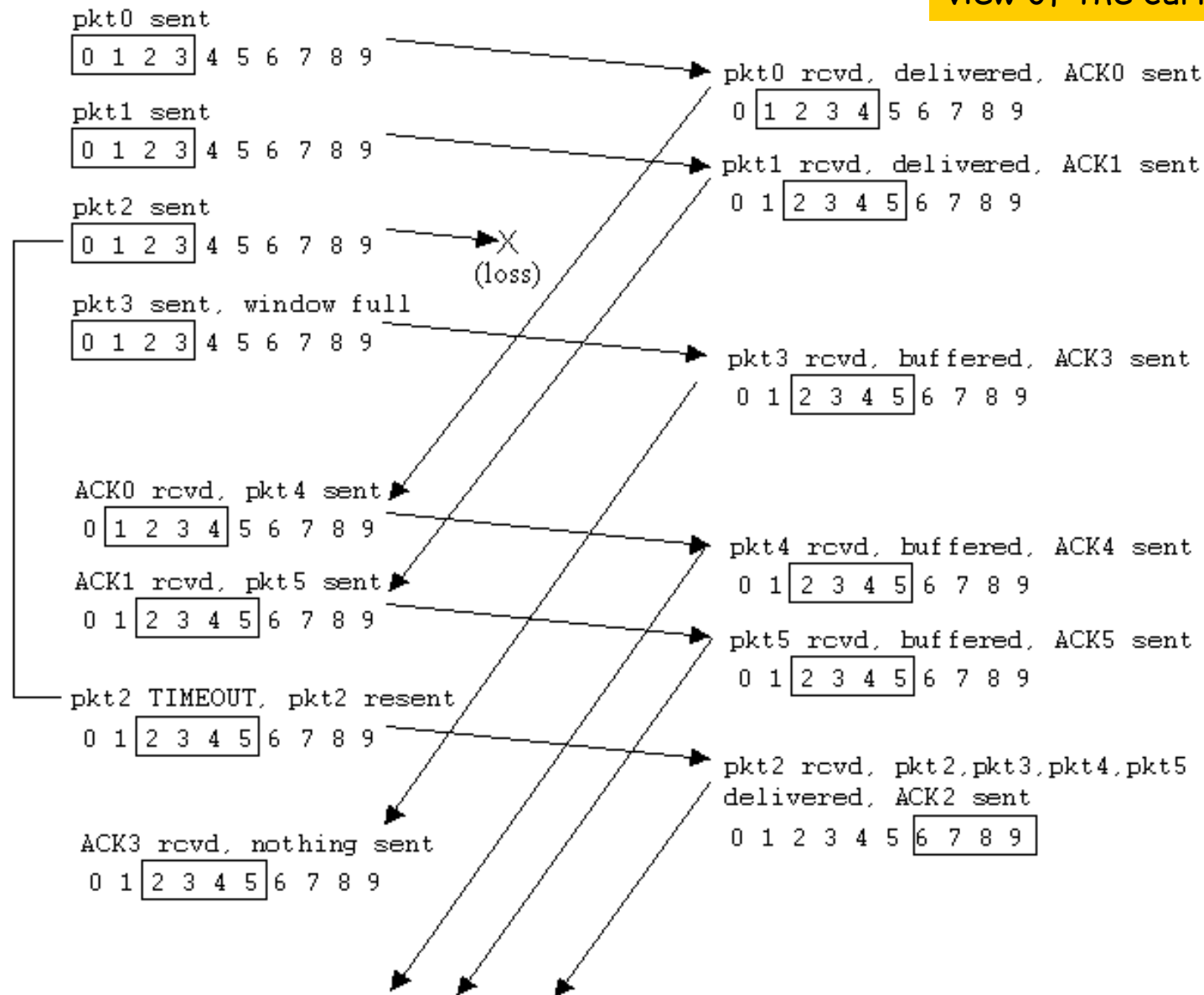
**Important!! Sender and receiver may have different views!!** Transport Layer 3-11

# Selective repeat: sender, receiver windows



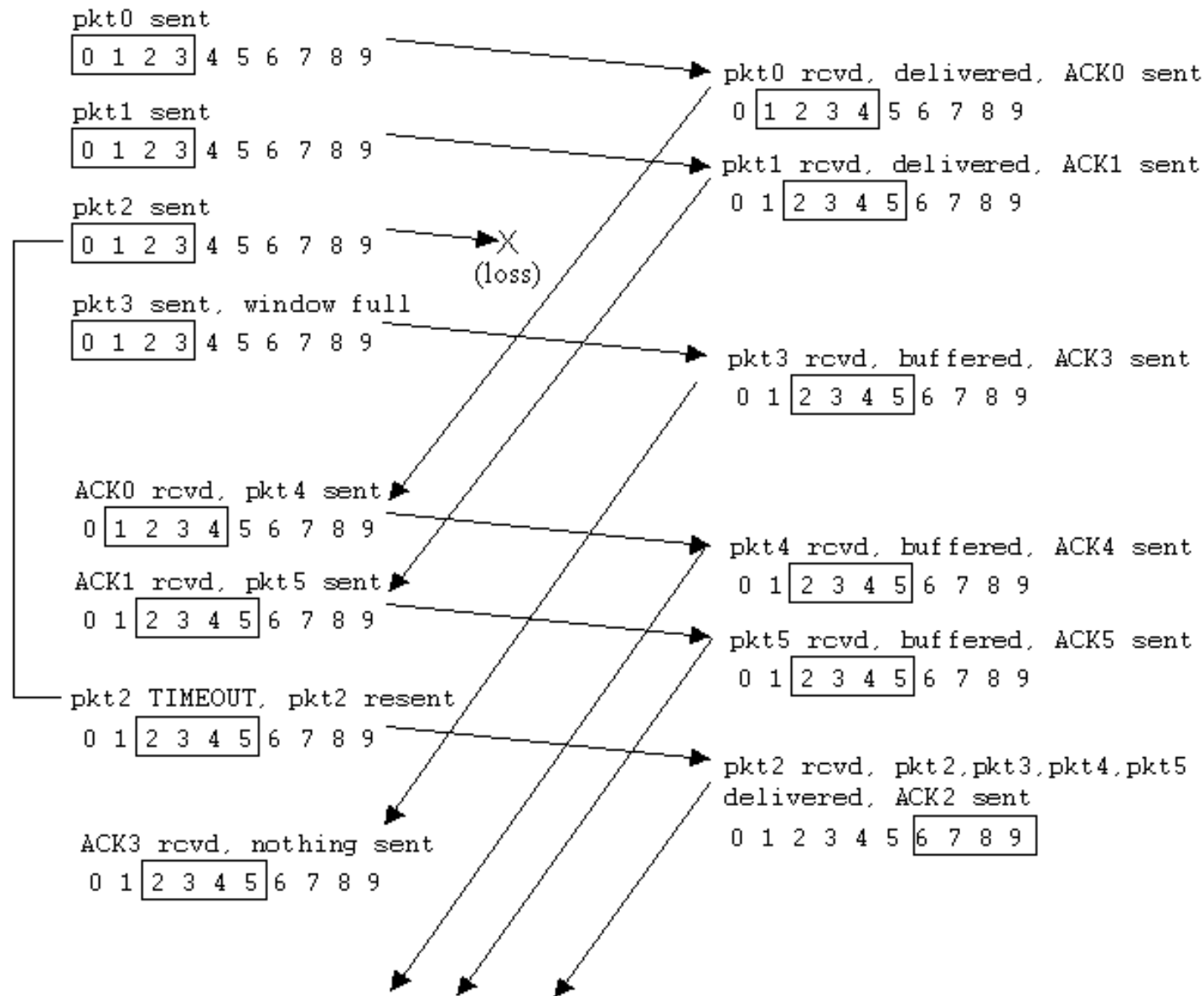
# Selective repeat in action

Transmitter and receiver can have different view of the current window



# Selective repeat in action

- 1) Delays in receiving ACKs
- 2) lost ACKs



# Selective repeat: dilemma

Example:

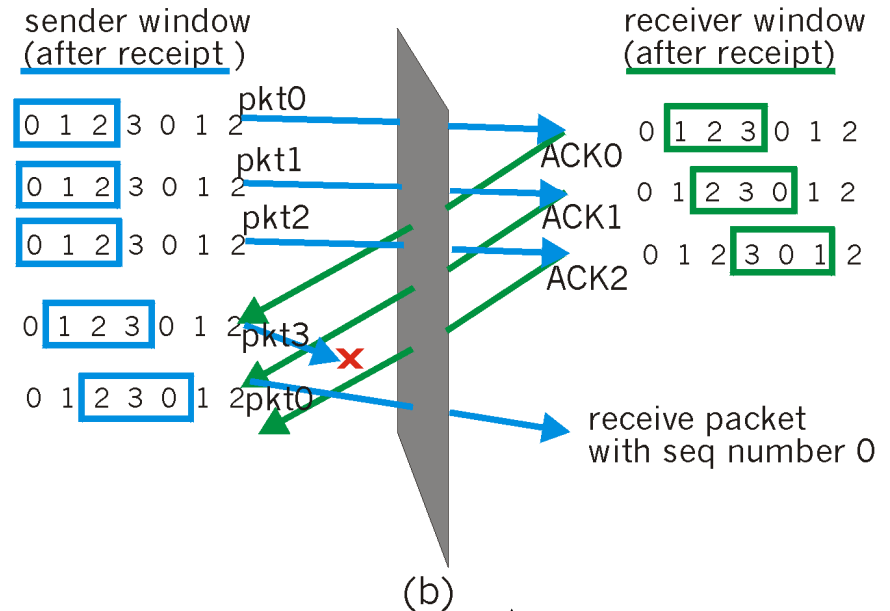
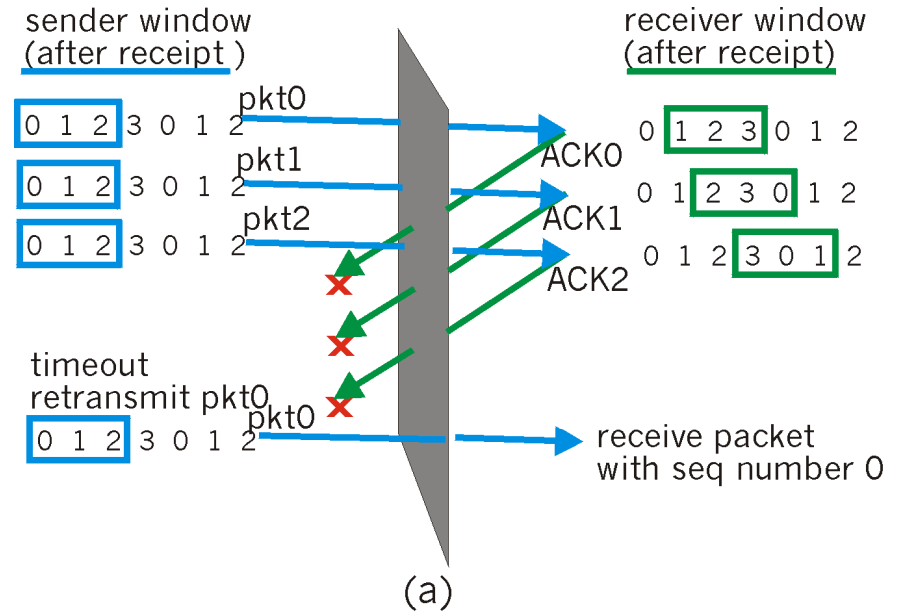
r seq #'s: 0, 1, 2, 3

r window size=3

r receiver sees no difference in two scenarios!

r incorrectly passes duplicate data as new in (a)

Q: what relationship between seq # size and window size?



Clearly at least the window must be small enough so that there is not ambiguity on sequence numbers!!! Is it enough in Selective Repeat??

## Answer to the dilemma

- r The window size must be less than or equal to half the size of the sequence number space for SR protocols

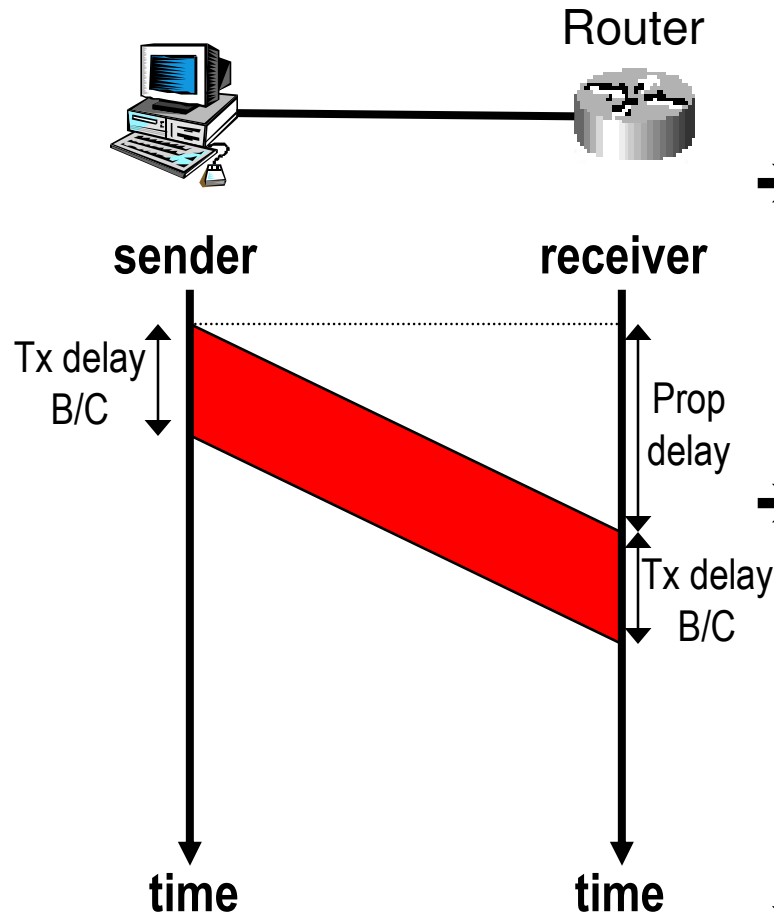
## Another issue

- r When ARQ solutions are applied at transport layer packets traverse not only one link but a path. Packets may arrive not in order, old packets may arrive with a long delay → in that case the answer is more involved. We cannot reuse a sequence number unless we are sure that old packets carrying that sequence number are out of the network (limit on the packet lifetime).



# Performance issues with/without pipelining

# Link delay computation



→ Transmission delay:

→  $C$  [bit/s] = link rate

→  $B$  [bit] = packet size

→ transmission delay =  $B/C$  [sec]

→ Example:

→ 512 bytes packet

→ 64 kbps link

→ transmission delay =  $512 \cdot 8 / 64000 = 64\text{ms}$

→ Propagation delay - constant depending on

→ Link length

→ Electromagnetic waves propagation speed in considered media

→ 200 km/s for copper links

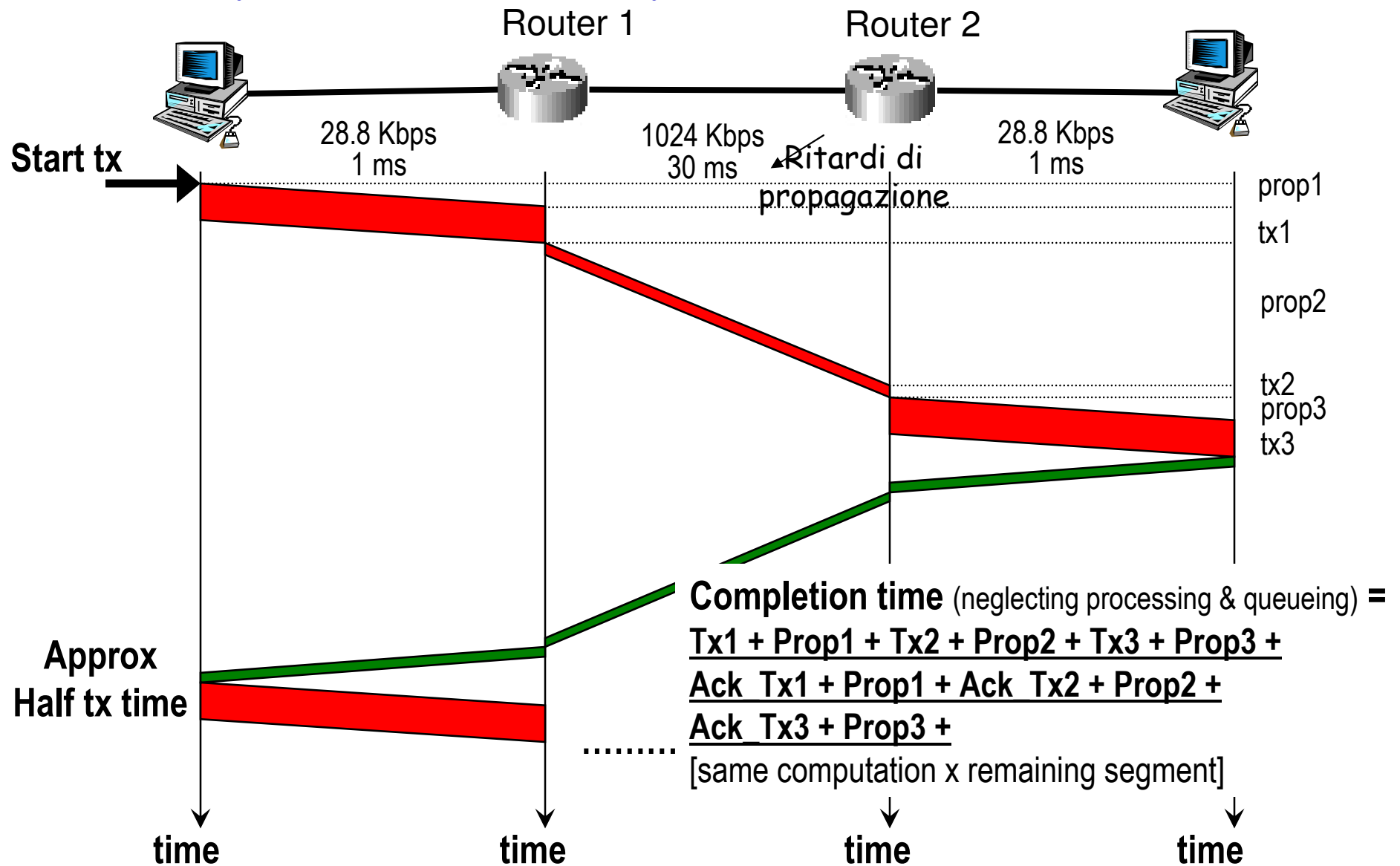
→ 300 km/s in air

→ other delays neglected

→ Queueing

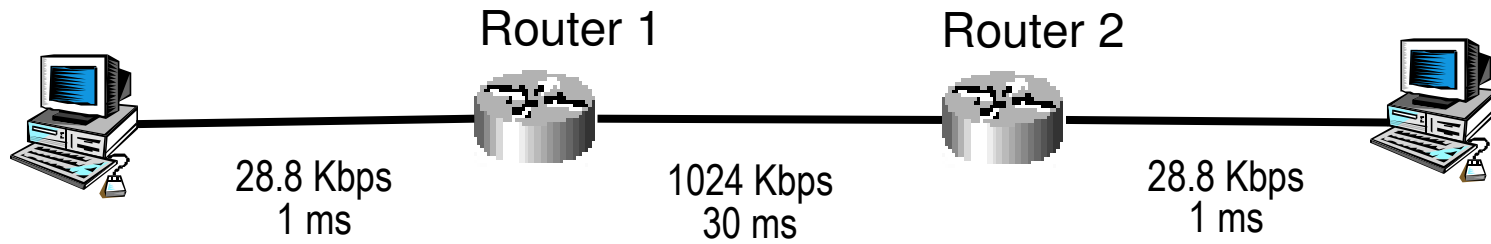
→ processing

# Stop-and-wait performance



# Stop-and-wait performance

## Numerical example



r Message:

- m 1024 bytes;
- m 2 segments:  
536+488 bytes
- m Overhead: 20 bytes  
TCP + 20 bytes IP
- m ACK = 40 bytes  
(header only)

Lower layer headers not considered

→ Segment 1:

- ⇒ Tx1 =  $576 \cdot 8 / 28,8 = 160 \text{ ms}$
- ⇒ Tx3 = Tx1
- ⇒ Tx2 =  $576 \cdot 8 / 1024 = 4,5 \text{ ms}$

→ Segment 2:

- ⇒ Tx1 =  $528 \cdot 8 / 28,8 = 146,7 \text{ ms}$
- ⇒ Tx3 = Tx1
- ⇒ Tx2 =  $528 \cdot 8 / 1024 = 4,1 \text{ ms}$

→ Acks:

- ⇒ Tx1 = Tx3 =  $40 \cdot 8 / 28,8 = 11,1 \text{ ms}$
- ⇒ Tx2 =  $40 \cdot 8 / 1024 = 0,3 \text{ ms}$

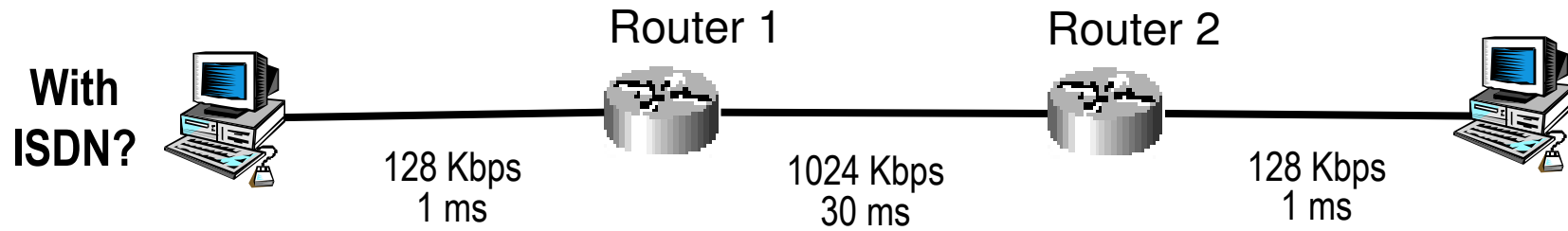
**RESULT:**

$$D = 667 \text{ (tx total)} + 2 \cdot \text{RTT} = 795 \text{ ms}$$

$$\text{THR} = 1024 \cdot 8 / 795 = 10,3 \text{ kbps}$$

# Stop-and-wait performance

## Numerical example



### → Segment 1:

$$\Rightarrow Tx1 = Tx3 = \frac{576 \cdot 8}{128} = 36 \text{ ms}$$

$$\Rightarrow Tx2 = \frac{576 \cdot 8}{1024} = 4,5 \text{ ms}$$

### → Acks:

$$\Rightarrow Tx1 = Tx3 = \frac{40 \cdot 8}{128} = 2,5 \text{ ms}$$

$$\Rightarrow Tx2 = \frac{40 \cdot 8}{1024} = 0,3 \text{ ms}$$

### → Segment 2:

$$\Rightarrow Tx1 = Tx3 = \frac{528 \cdot 8}{128} = 33 \text{ ms}$$

$$\Rightarrow Tx2 = \frac{528 \cdot 8}{1024} = 4,1 \text{ ms}$$

### RESULT:

$$D = 151,9 \text{ (tx total)} + 2 \cdot \text{RTT} = 279,9 \text{ ms}$$

$$\text{THR} = \frac{1024 \cdot 8}{279,9} = 29,3 \text{ kbps}$$

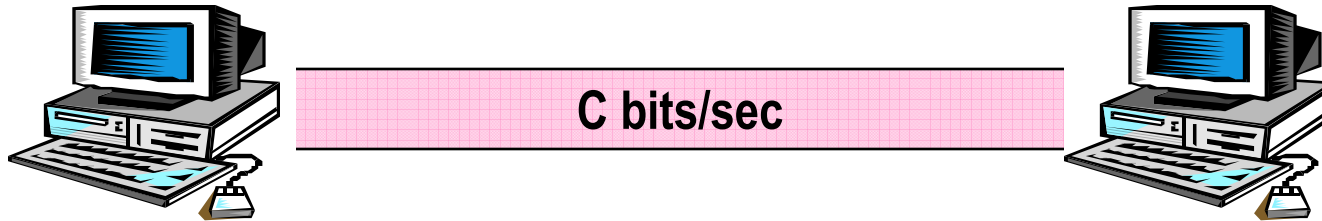
on Gbps fiber optics?

$$D = \text{negligible} + 2 \cdot \text{RTT} = 128 \text{ ms}$$

$$\text{THR} = \frac{1024 \cdot 8}{128} = 64 \text{ kbps}$$

MA È VERAMENTE MEGLIO  
AD ALTO DATA RATE? NO  
—DI QUI A POCO...

# Simplified performance model



Approximate analysis, much simpler than multi-hop  
Typically,  $C$  = bottleneck link rate

**MSS** = segment size (ev. ignore overhead)

**MSIZE** = message size

Ignore ACK transmission time

No loss of segments

**W** = number of outstanding segments

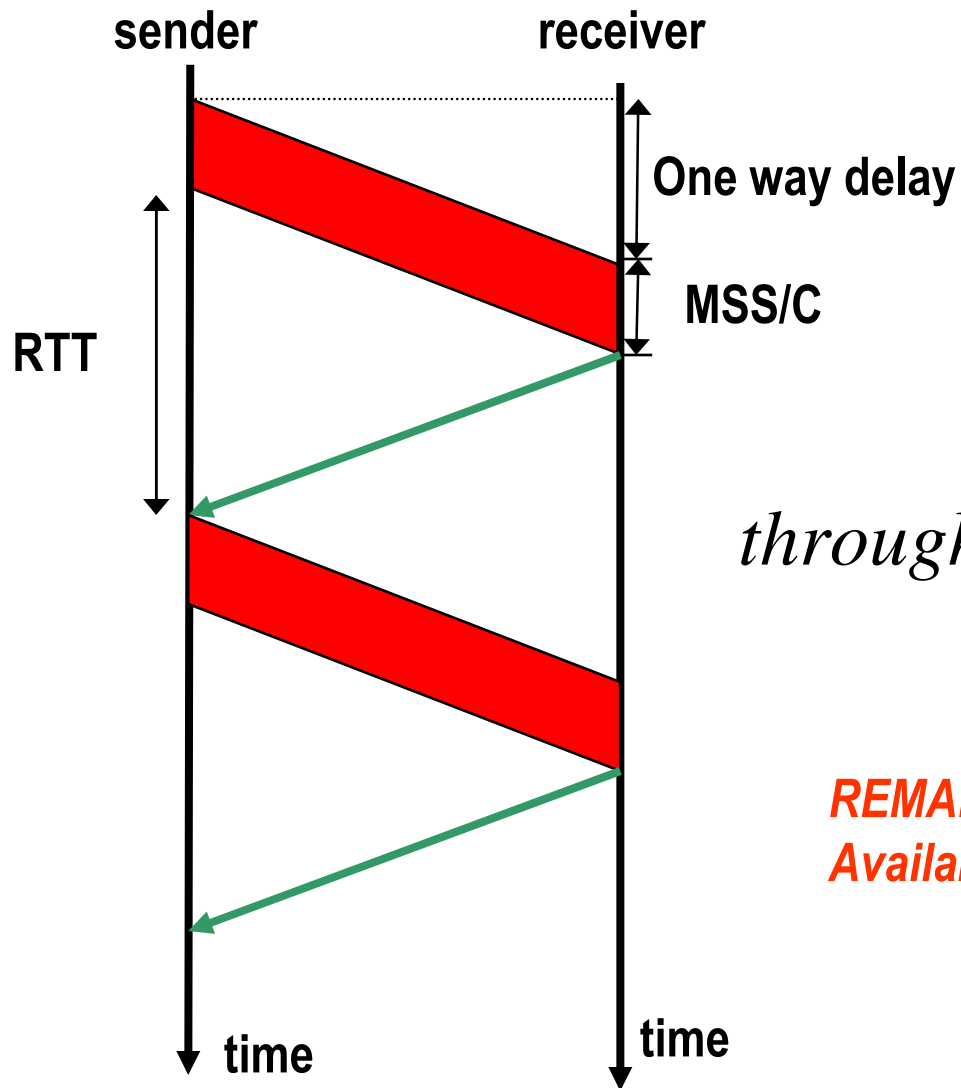
**W=1**: stop-and-wait

**W>1**: go-back-N (sliding window)

*This is a highly dynamic parameter in TCP!!*

*For now, consider W fixed*

# W=1 case (stop-and-wait)

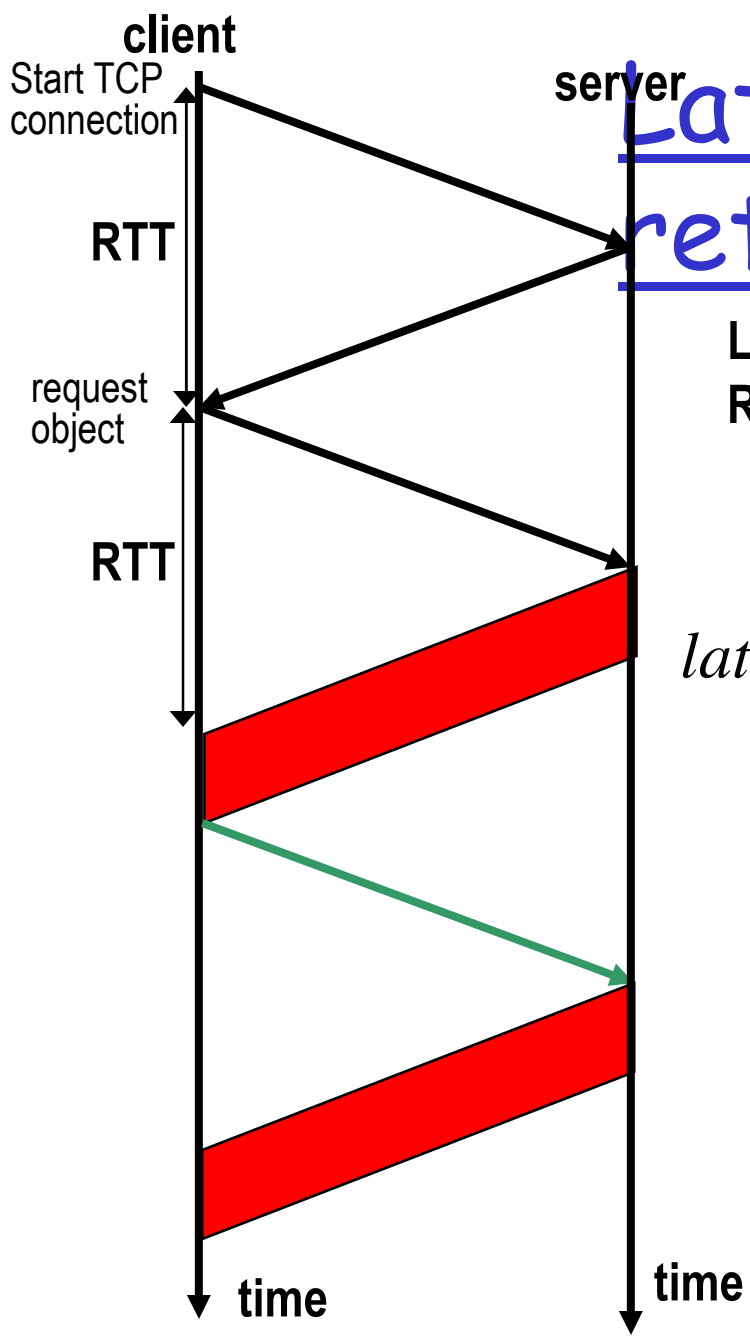


$$\text{throughput} = \frac{MSS}{RTT + MSS / C}$$

**REMARK: throughput always lower than Available link rate!**

# Latency in TCP retrieval model

Latency: time elapsing between TCP connection Request, and last bit received at client



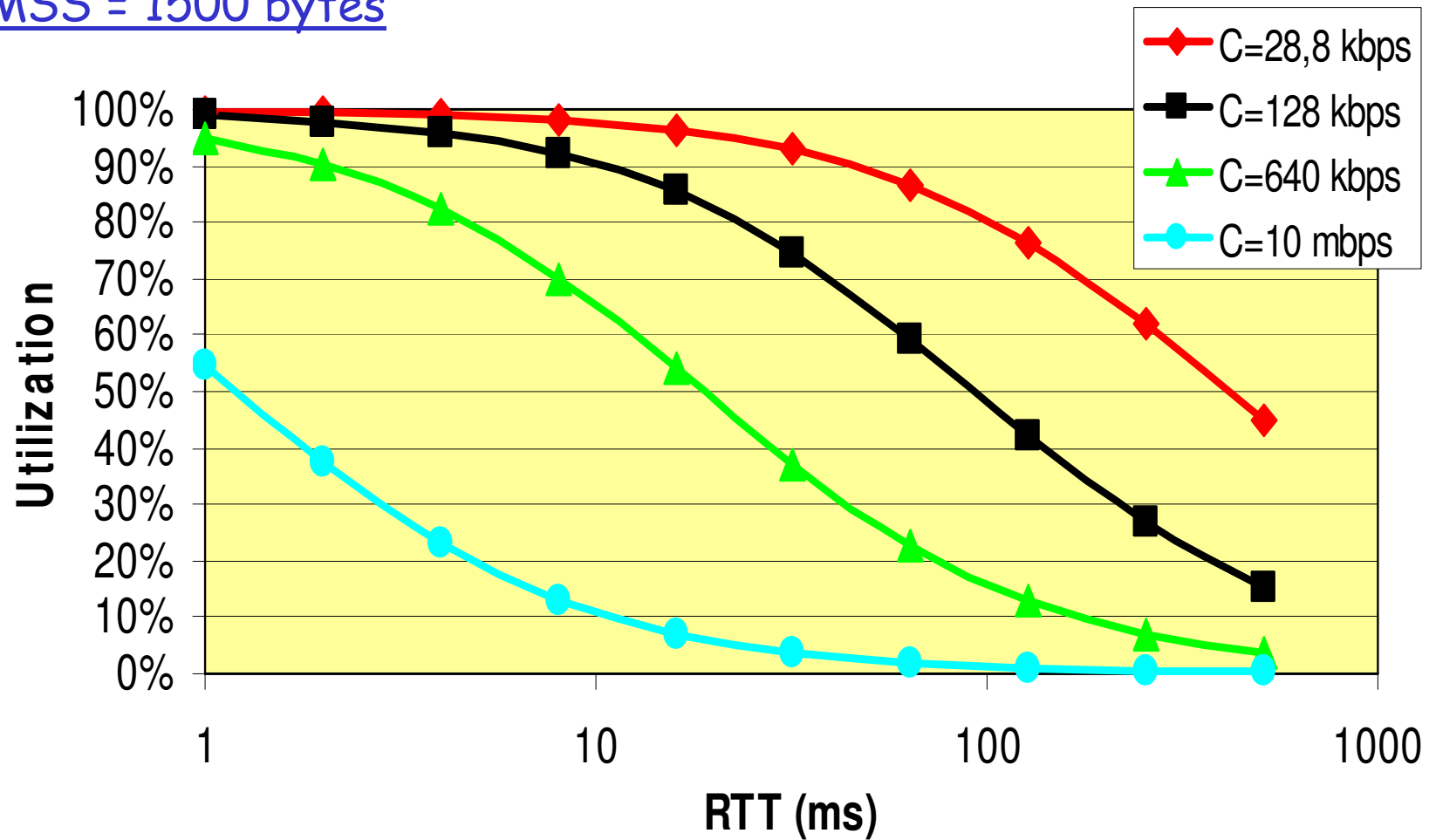
$$latency = 2RTT + \frac{MSIZE}{C} + \left[ \frac{MSIZE}{MSS} - 1 \right] RTT$$

↑  
Number of segments  
In which message  
is split



# W=1 case (stop-and-wait)

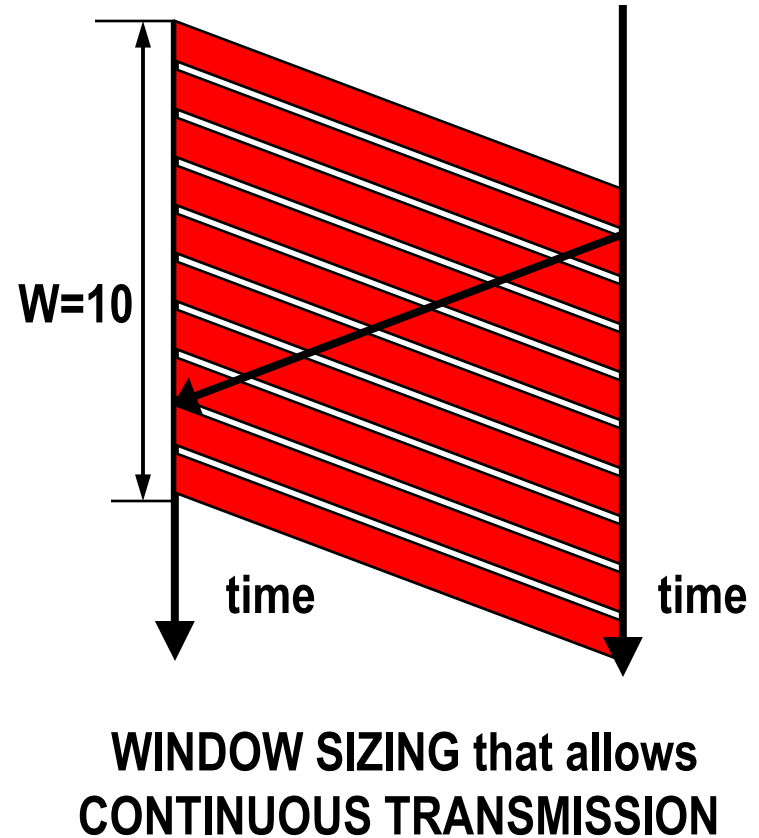
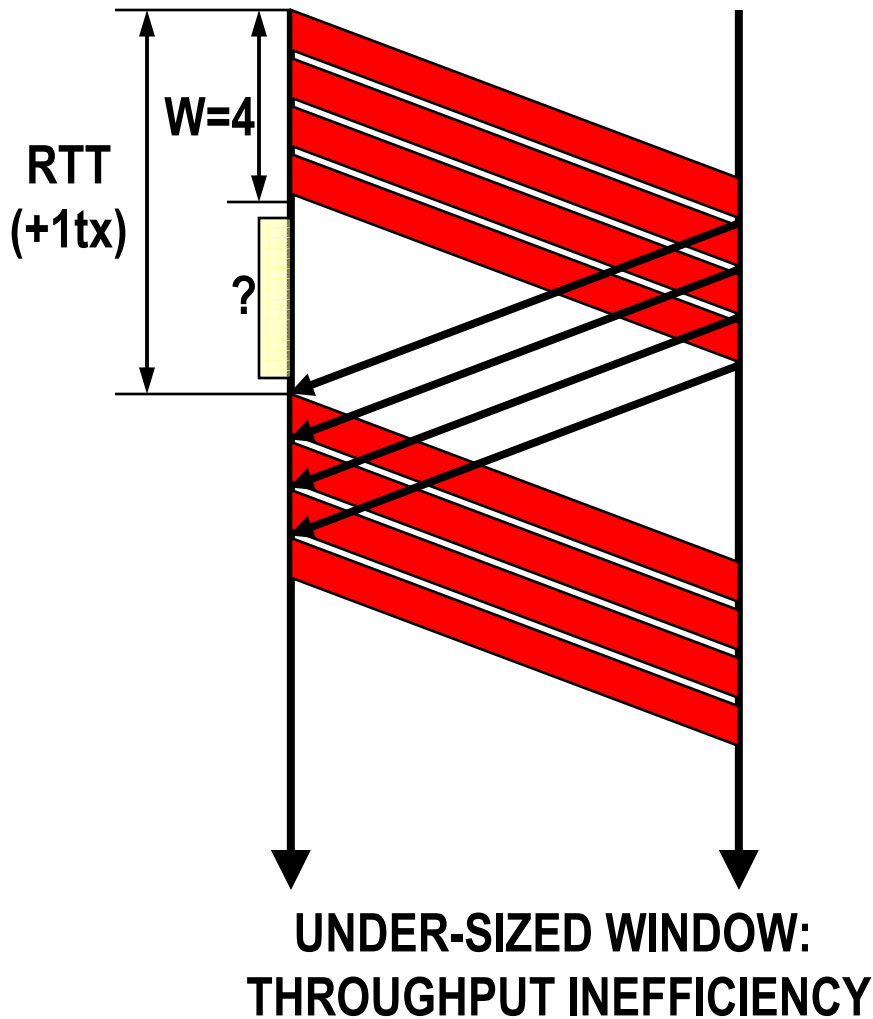
MSS = 1500 bytes



**Under-utilization with: 1) high capacity links, 2) large RTT links**

# Pipelining ( $W > 1$ ) analysis

## two cases



# Continuous transmission

*Condition in which link rate is fully utilized*

$$\underbrace{W \cdot \frac{MSS}{C}}_{\text{Time to transmit } W \text{ segments}} > \underbrace{RTT + \frac{MSS}{C}}_{\text{Time to receive Ack of first segment}}$$

We may elaborate:

$$W \cdot MSS > RTT \cdot C + MSS \approx RTT \cdot C$$

This means that full link utilization is possible when window size (in bits) is Greater than the bandwidth (C bit/s) delay (RTT s) product!

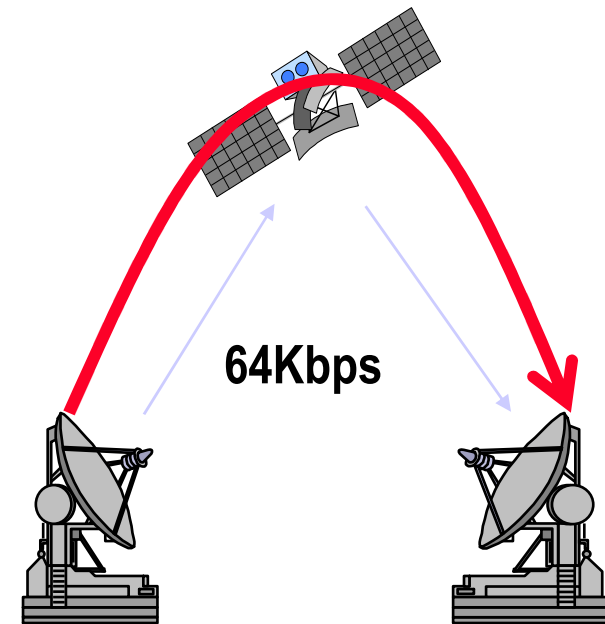
# Bandwidth-delay product



→ **Network: like a pipe**

→  **$C$  [bit/s]  $\times$   $D$  [s]**

- ⇒ number of bits “flying” in the network
- ⇒ number of bits injected in the network by the tx, before that the first bit is rxd



A 15360 (64000x0.240) bits  
“worm” in the air!!

bandwidth-delay product = no of bytes that saturate network pipe

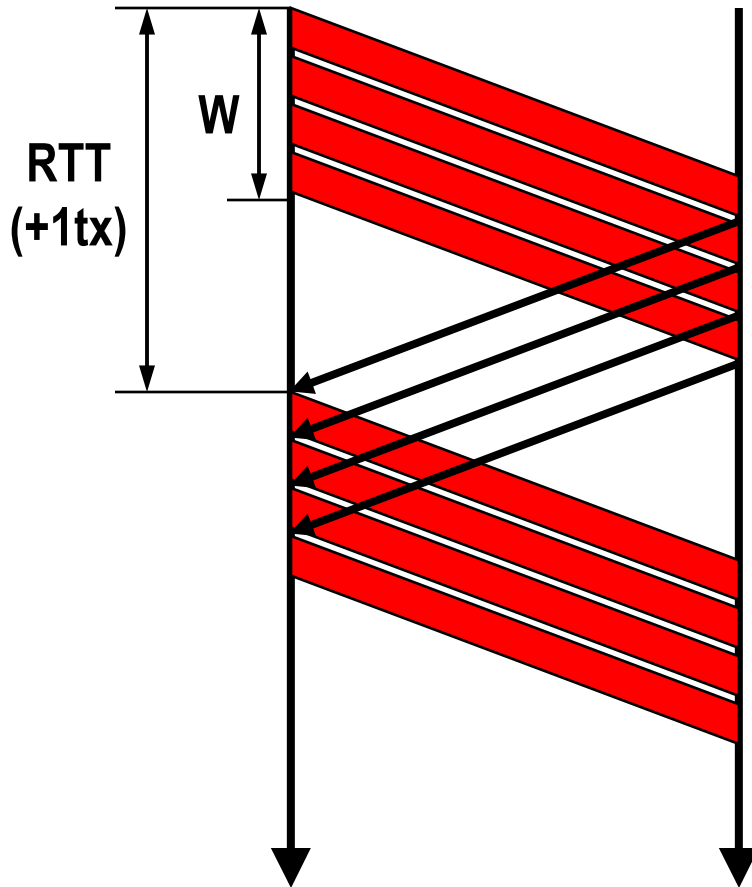
# Long Fat Networks

**LFNs (el-ef-an(t)s): large bandwidth-delay product**

NETWORK	RTT (ms)	rate (kbps)	BxD (bytes)
Ethernet	3	10.000	3.750
T1, transUS	60	1.544	11.580
T1 satellite	480	1.544	92.640
T3 transUS	60	45.000	337.500
Gigabit transUS	60	1.000.000	7.500.000

**The 65535 (16 bit field in TCP header) maximum window size  $W$  may be a limiting factor!**

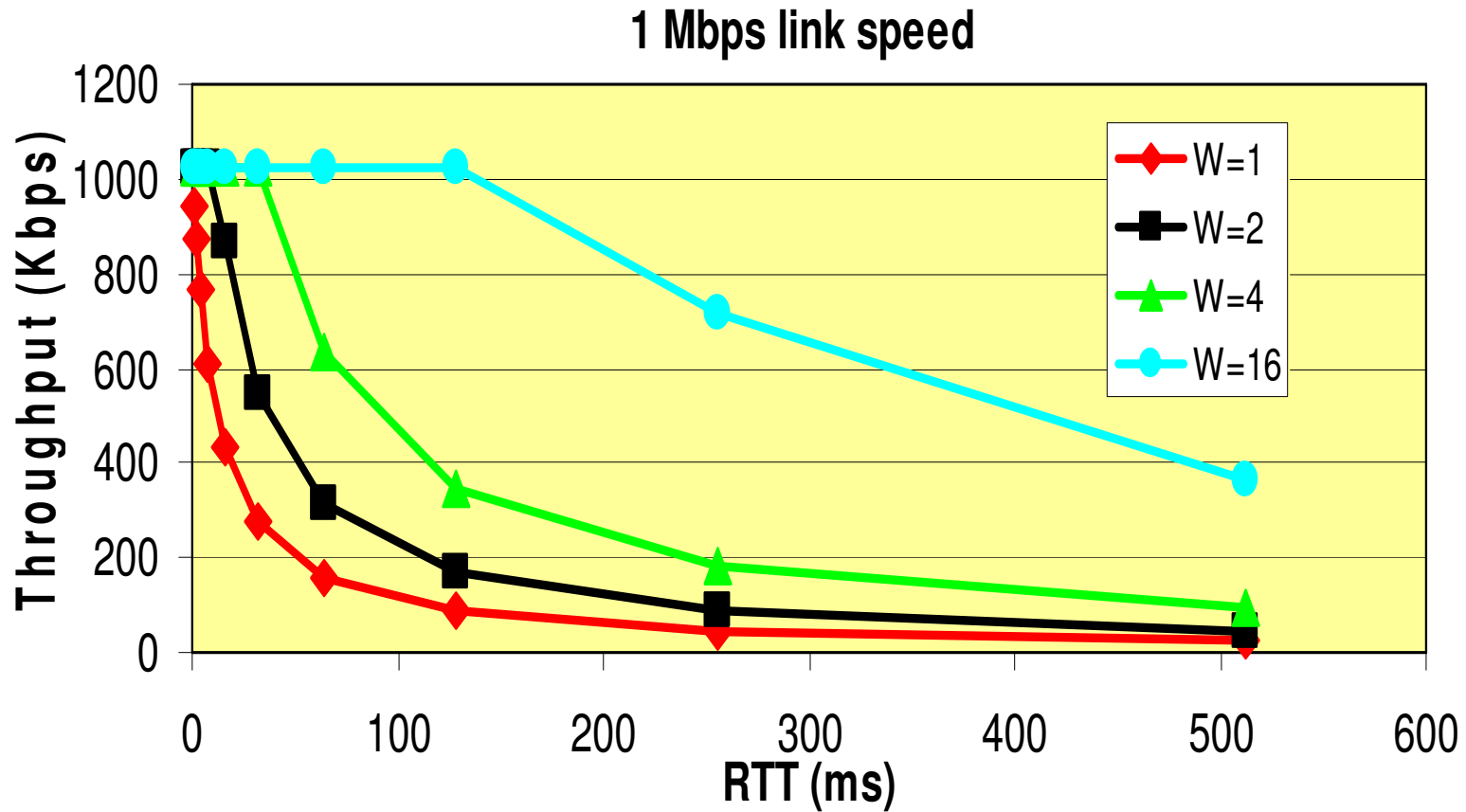
# Pipelining ( $W > 1$ ) analysis



$$thr = \min \left( C, \frac{W \cdot MSS}{RTT + MSS / C} \right)$$

# Throughput for pipelining

MSS = 1500 bytes



# Chapter 3 outline

- r 3.1 Transport-layer services
- r 3.2 Multiplexing and demultiplexing
- r 3.3 Connectionless transport: UDP
- r 3.4 Principles of reliable data transfer
- r 3.5 Connection-oriented transport: TCP
  - m segment structure
  - m reliable data transfer
  - m flow control
  - m connection management
- r 3.6 Principles of congestion control
- r 3.7 TCP congestion control



# TCP: Overview

RFCs: 793, 1122, 1323, 2018, 2581

r **point-to-point:**

m one sender, one receiver

r **reliable, in-order *byte stream*:**

m no "message boundaries"

r **pipelined:**

m TCP congestion and flow control set window size

r ***send & receive buffers***

r **full duplex data:**

m **bi-directional data flow** in same connection

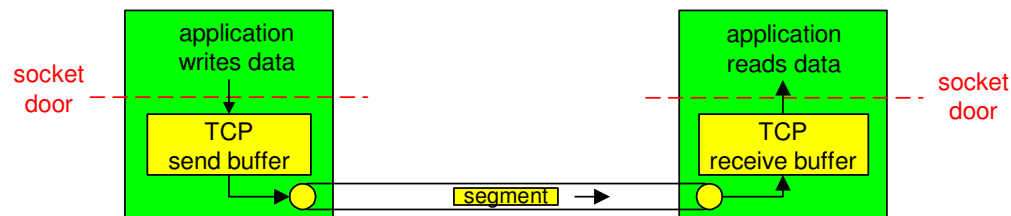
m **MSS: maximum segment size**

r **connection-oriented:**

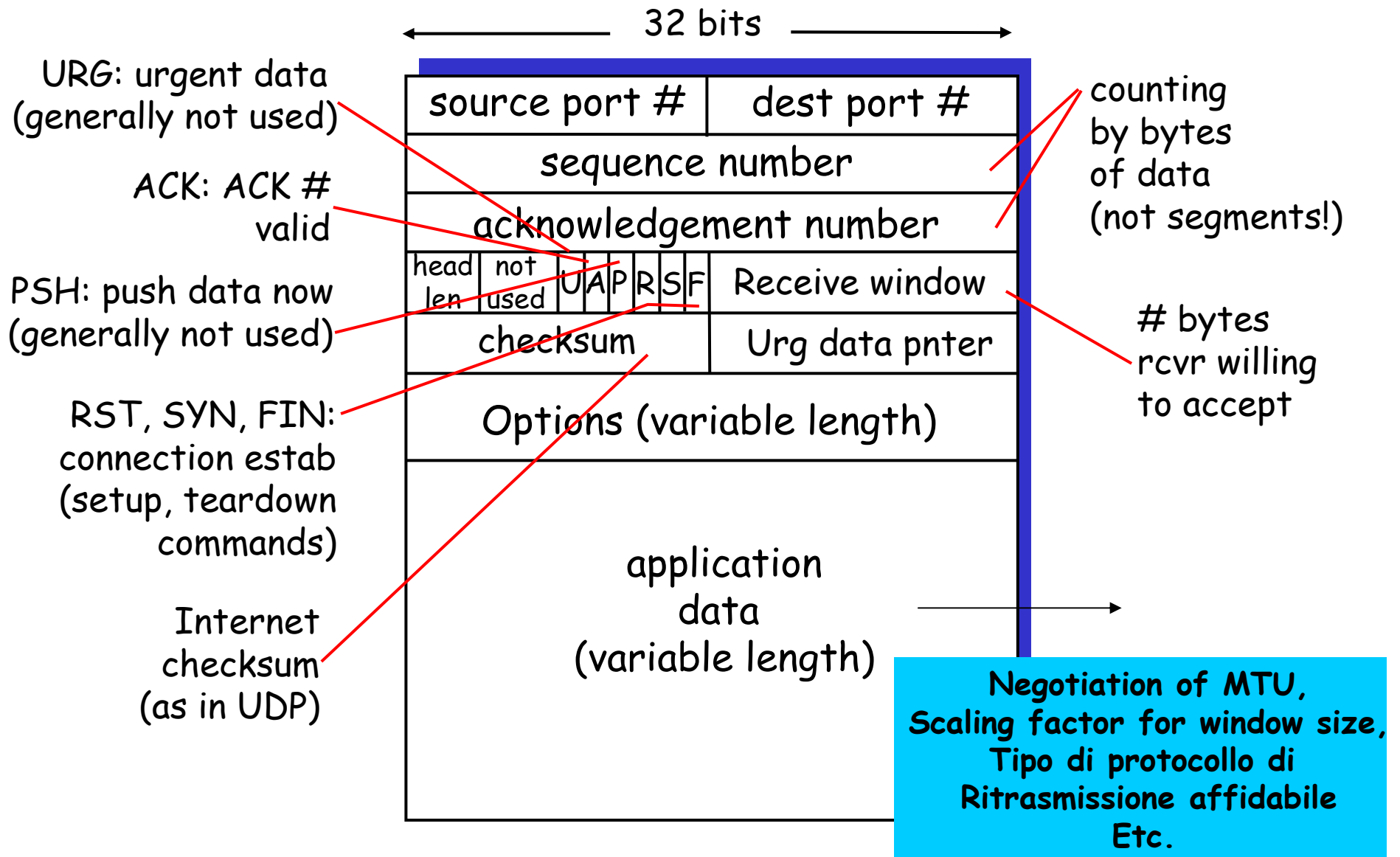
m handshaking (exchange of control msgs) init's sender, receiver state before data exchange

r **flow controlled:**

m sender will not overwhelm receiver



# TCP segment structure



# Window Scale Option

- r Appears in SYN segment
  - m operates only if both peers understand option
- r allows client & server to agree on a different W scale
  - m specified in terms of bit shift (from 1 to 14)
  - m maximum window:  $65535 * 2^b$
  - m  $b=14$  means max  $W = 1.073.725.440$  bytes!!

Source port				Destination port				
32 bit Sequence number								
32 bit acknowledgement number								
Header length	6 bit Reserved	U R G	A C K	P S H	R S T	S Y N	F I N	Window size
checksum				Urgent pointer				

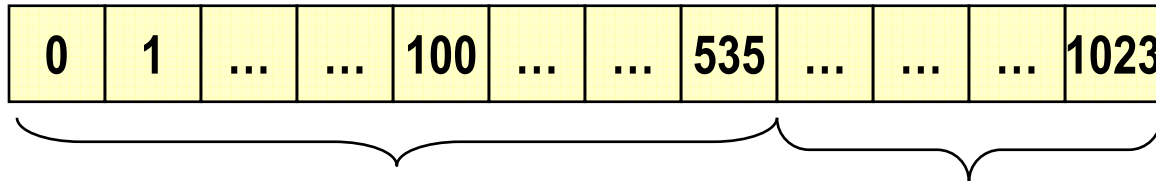
- r Sequence number:
  - m Sequence number of the *first* byte in the segment.
  - m When reaches  $2^{32}-1$ , next wraps back to 0
- r Acknowledgement number:
  - m valid only when ACK flag on
  - m Contains the *next* byte sequence number that the host *expects* to receive (= last successfully received byte of data + 1)
  - m grants successful reception for all bytes up to ack# - 1 (cumulative)
- r When seq/ack reach  $2^{32}-1$ , next wrap back to 0

# TCP data transfer management

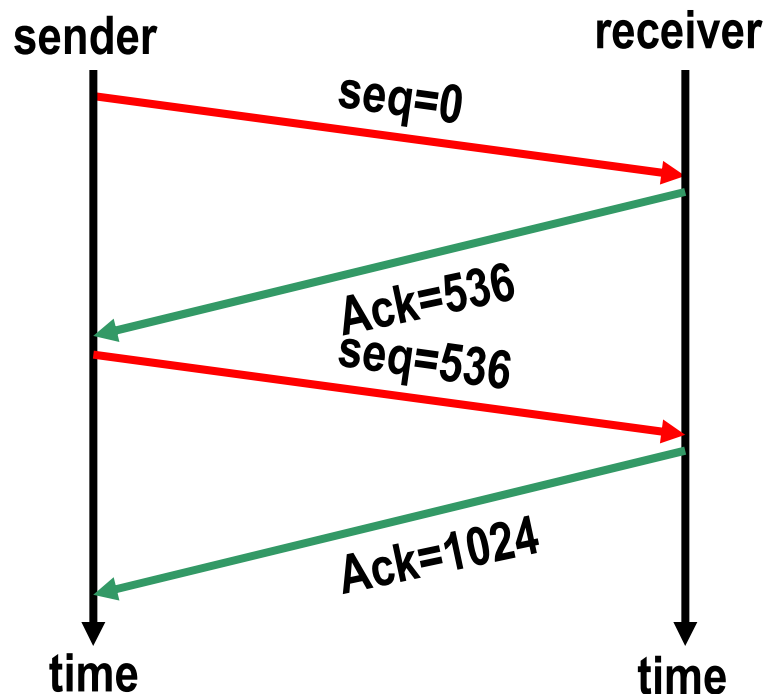
- r Full duplex connection
  - m data flows in both directions, independently
  - m To the application program these appear as two unrelated data streams
- r each end point maintains a sequence number
  - m Independent sequence numbers at both ends
  - m Measured in bytes
- r acks often carried on top of reverse flow data segments (piggybacking)
  - m But ack packets alone are possible

# Byte-oriented

Example: 1 Kbyte message – 1024 bytes



Example: segment size = 536 bytes → 2 segments: 0-535; 536-1023

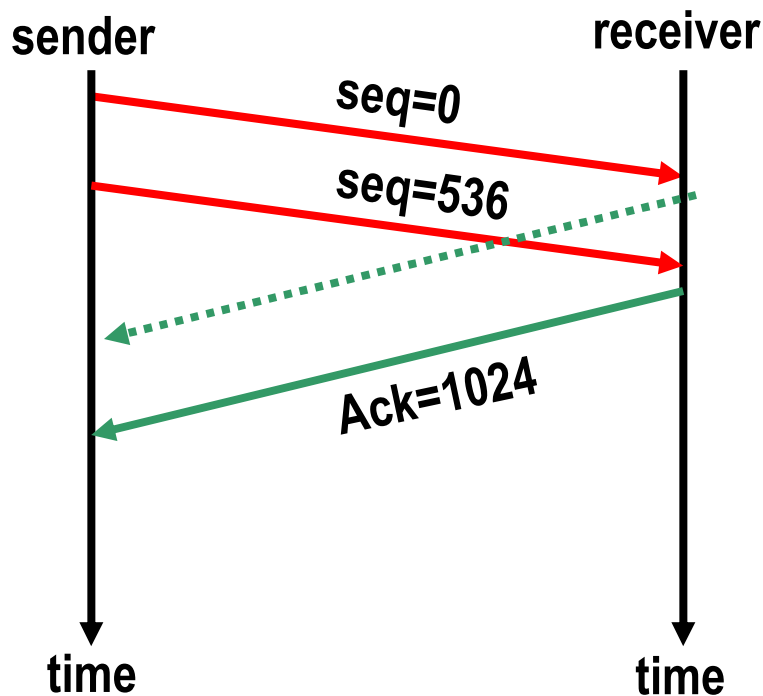
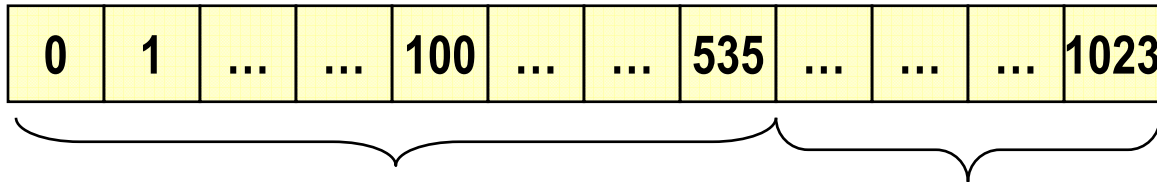


→ **No explicit segment size indication**

- ⇒ Seq = first byte number
- ⇒ Returning Ack = last byte number + 1
- ⇒ Segment size = Ack-seq#

# Pipelining - cumulative ack

Example: 1024 bytes msg; seg\_size = 536 bytes → 2 segments: 0-535; 536-1023



## → Cumulative ack

⇒ ACK = all previous bytes correctly received!

⇒ E.g. ACK=1024: all bytes 0-1023 received

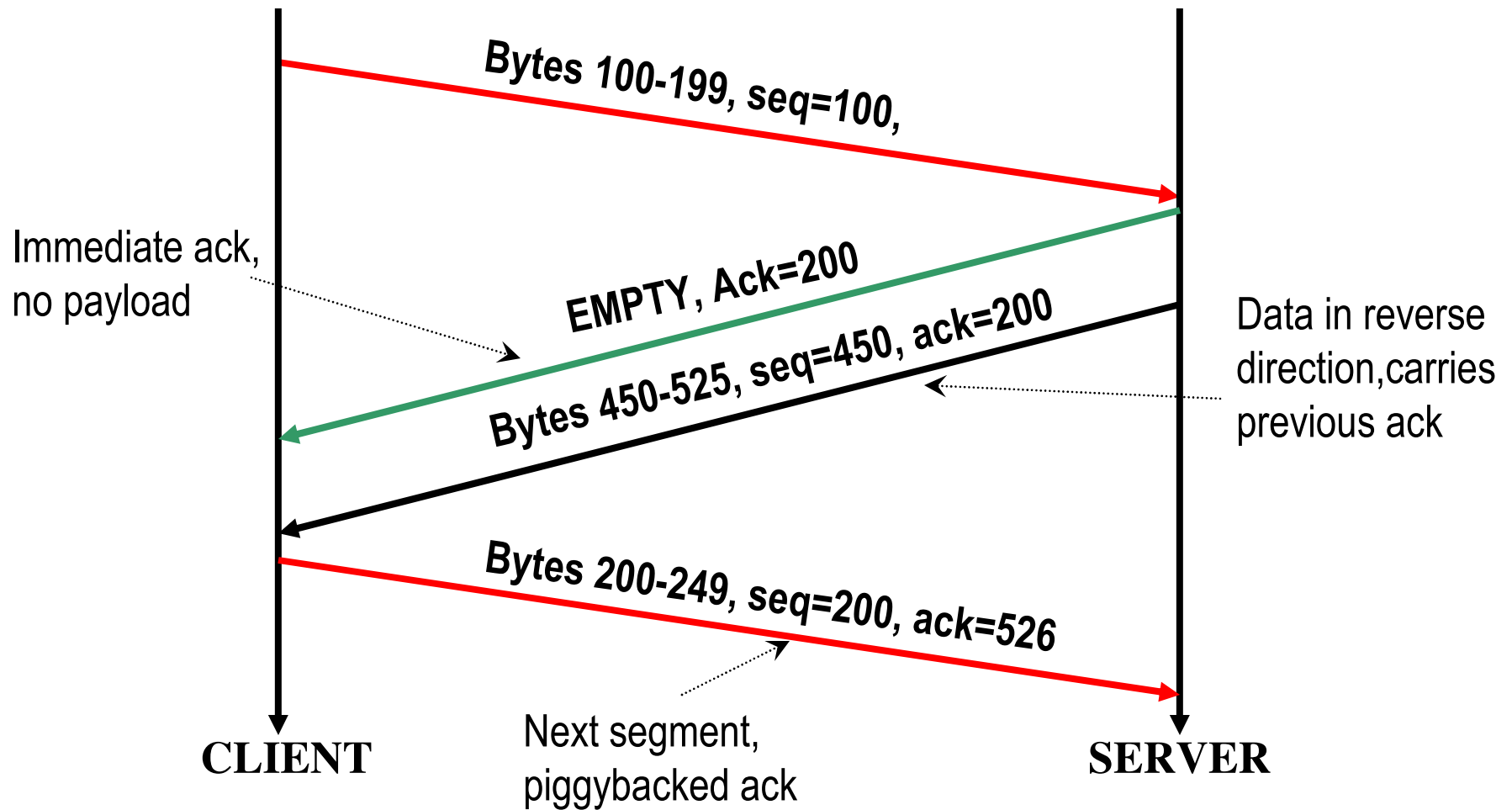
## ⇒ Other names of pipelining:

⇒ Go-Back-N ARQ mechanisms

⇒ Sliding window mechanisms

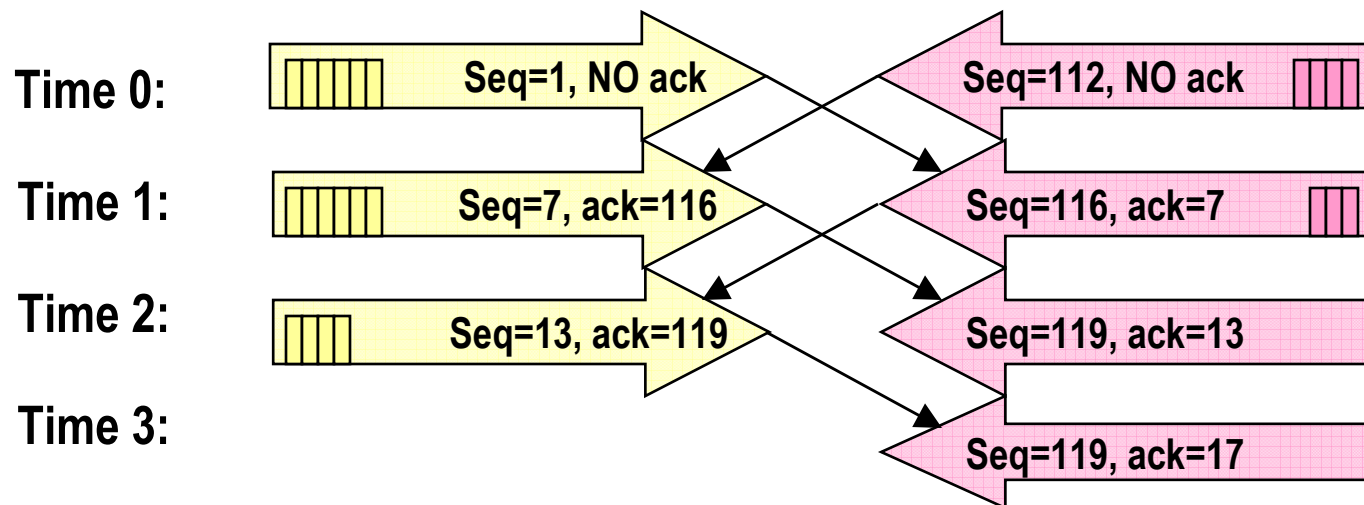
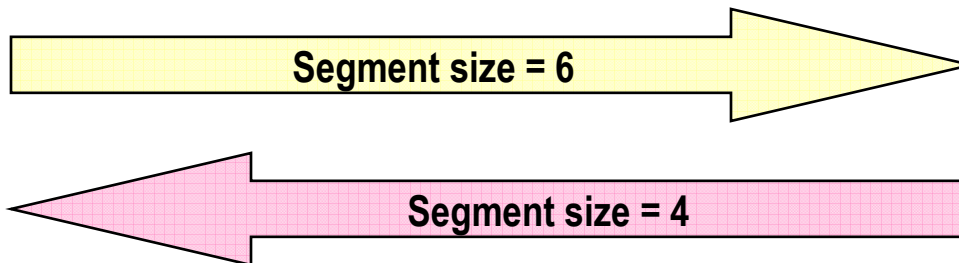
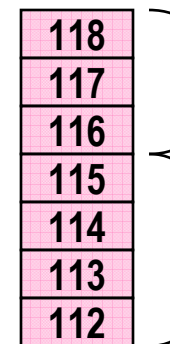
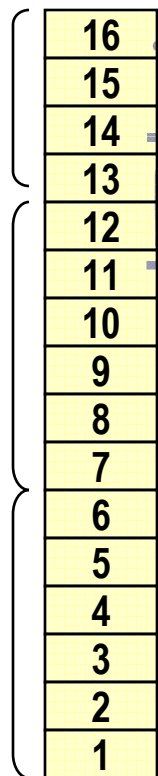
***Why pipelining? Dramatic improvement in efficiency!***

# Multiple acks; Piggybacking





# TCP data transfer bidirectional example



# TCP seq. #'s and ACKs

TCP Solution: Go Back N like

## Seq. #'s:

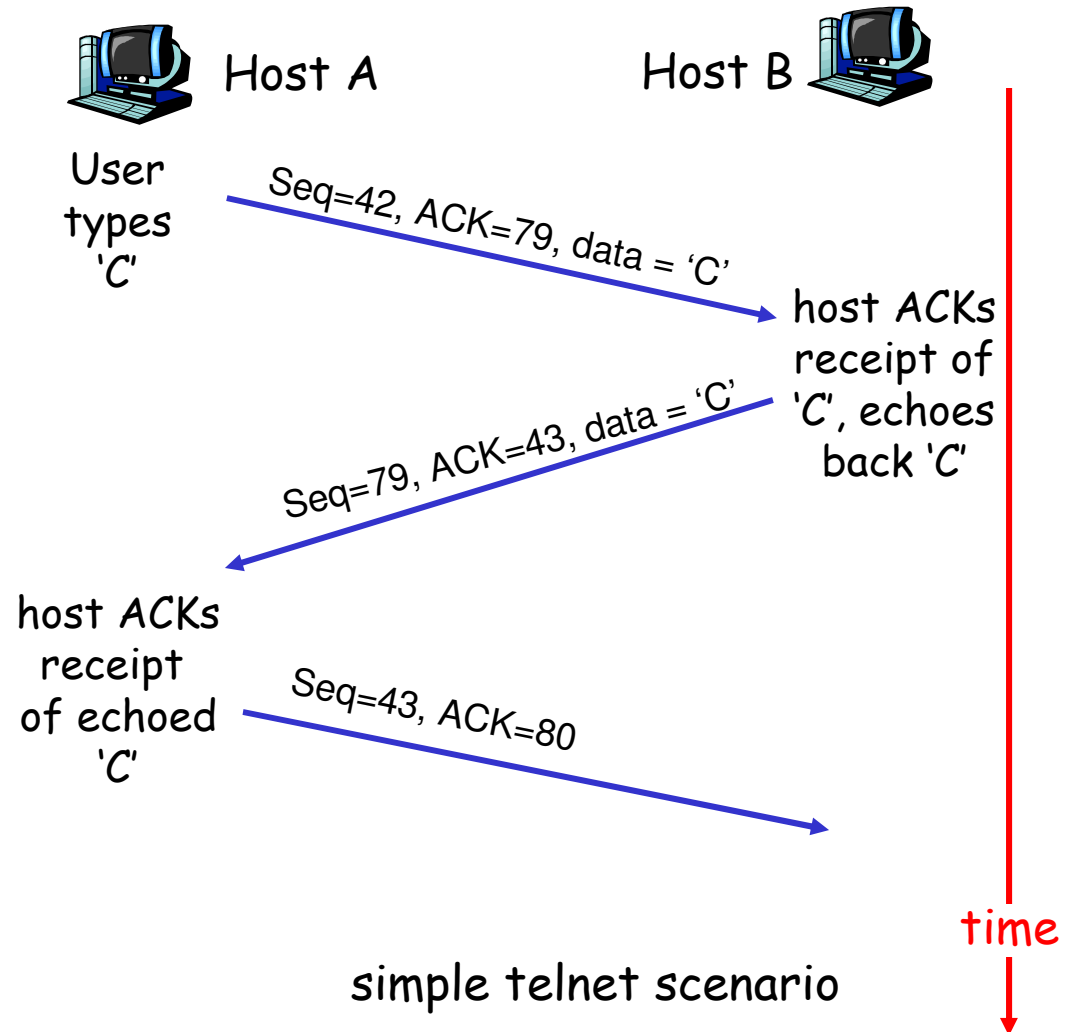
- m byte stream
- "number" of first byte in segment's data

## ACKs:

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

Q: how receiver handles out-of-order segments

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



# TCP Round Trip Time and Timeout

**Q:** how to set TCP timeout value? (not trivial, highly varying, it is a RTT over a network path)

- r longer than RTT
  - m but RTT varies
- r too short: premature timeout
  - m unnecessary retransmissions
- r too long: slow reaction to segment loss

**Q:** how to estimate RTT?

r **SampleRTT:** measured time from segment transmission until ACK receipt

m ignore retransmissions

**Why??**

- r **SampleRTT** will vary, want estimated RTT "smoother"
  - m average several recent measurements, not just current **SampleRTT**

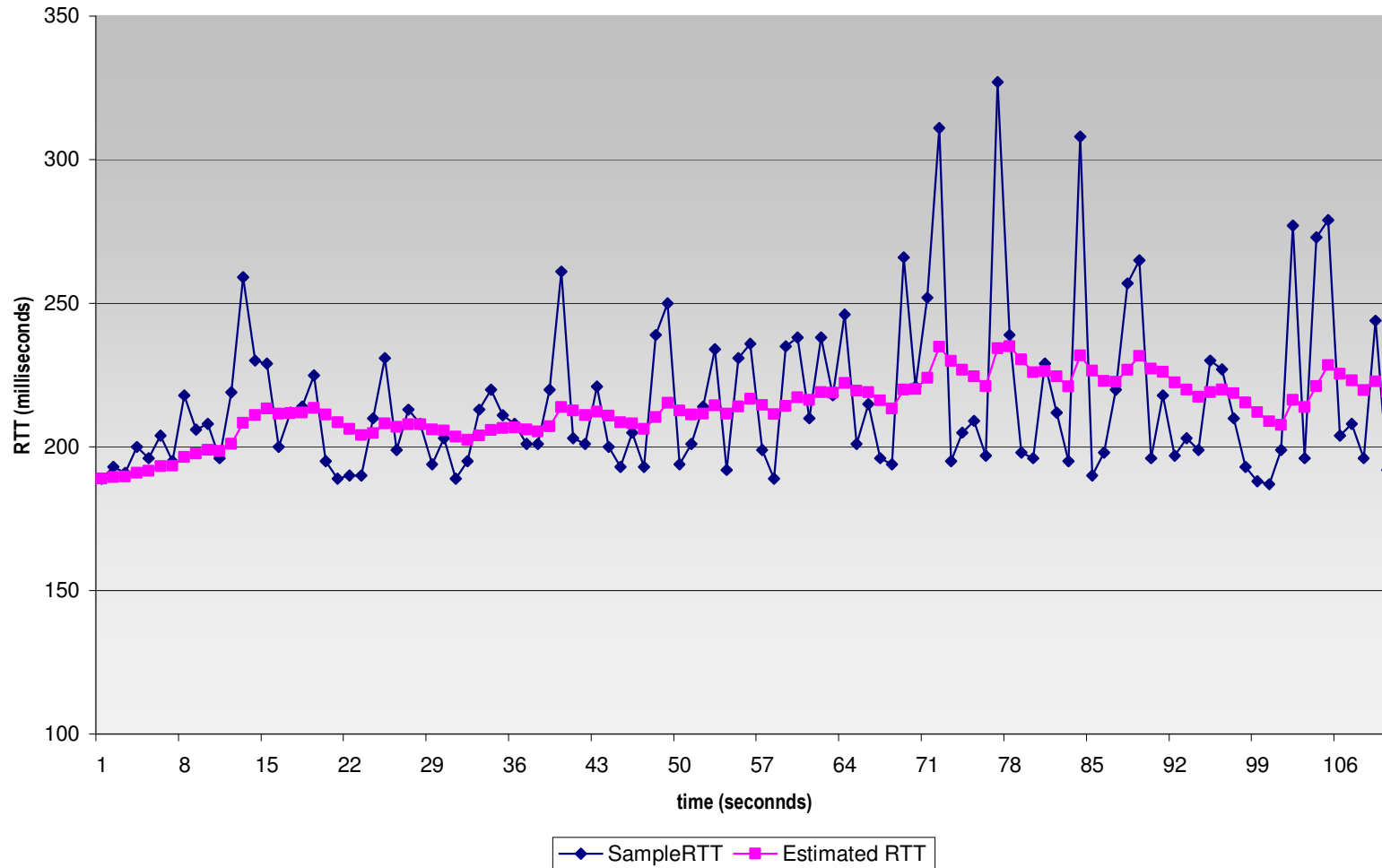
# TCP Round Trip Time and Timeout

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

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

# Example RTT estimation:

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



# TCP Round Trip Time and Timeout

## Setting the timeout

- r EstimatedRTT plus "safety margin"
  - m large variation in EstimatedRTT → larger safety margin
- r first estimate of how much SampleRTT deviates from EstimatedRTT:

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

(typically,  $\beta = 0.25$ )

Then set timeout interval:

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