# 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

# Guessing right?
## Karn's problem



**Scenario 1**

DATA

RTO

M

Retransmit DATA

ack

**Scenario 2**

RTO

M?

M?

DATA

retransmit

ack

How can we distinguish among an ACK to the original segment and to a duplicate?

# Solution to Karn's problem

r **Very simple: DO NOT update RTT when a segment has been retransmitted because of RTO expiration!**

r **Instead, use Exponential backoff**

    m *double RTO for every subsequent expiration of same segment*

        • When at 64 secs, stay

        • persist up to 9 minutes, then reset

# TCP reliable data transfer (more in detail)

r TCP creates rdt service on top of IP's unreliable service

r Pipelined segments

r Cumulative acks

r TCP uses single retransmission timer

r Retransmissions are triggered by:
  m timeout events
  m duplicate acks

r Initially consider simplified TCP sender:
  m  ignore duplicate acks
  m ignore flow control, congestion control

# TCP sender events:

## data rcvd from app:

r Create segment with seq #

r seq # is byte-stream number of first data byte in segment

r start timer if not already running (think of timer as for oldest unacked segment)

r expiration interval: `TimeOutInterval`

## timeout:

r retransmit segment that caused timeout

r restart timer

## Ack rcvd:

r If acknowledges previously unacked segments

m update what is known to be acked

m start timer if there are outstanding segments

```
NextSeqNum = InitialSeqNum
SendBase = InitialSeqNum
```

Purche' non si ecceda la finestra

# TCP sender (simplified)

```
loop (forever) {
   switch(event)

   event: data received from application above
        create TCP segment with sequence number NextSeqNum
        if (timer currently not running)
            start timer
        pass segment to IP
        NextSeqNum = NextSeqNum + length(data)

   event: timer timeout
        retransmit not-yet-acknowledged segment with
            smallest sequence number
        start timer

   event: ACK received, with ACK field value of y
        if (y > SendBase) {
            SendBase = y
            if (there are currently not-yet-acknowledged segments)
                start timer
        }

} /* end of loop forever */
```
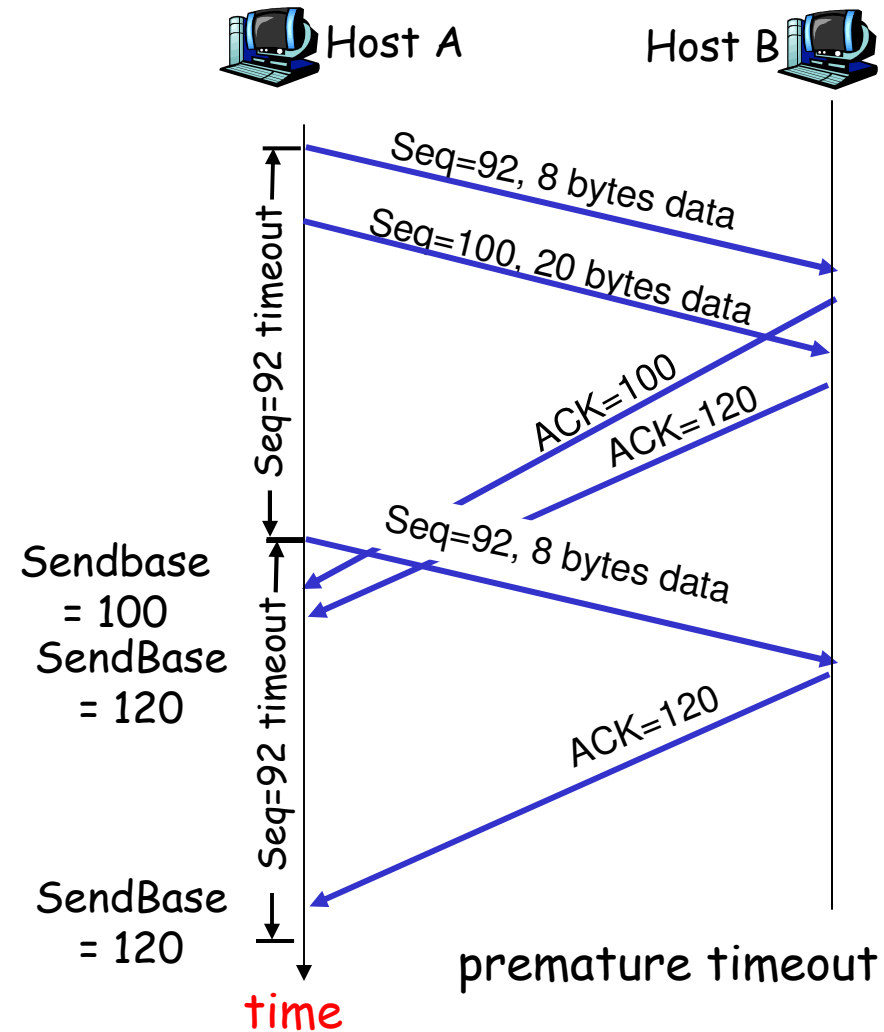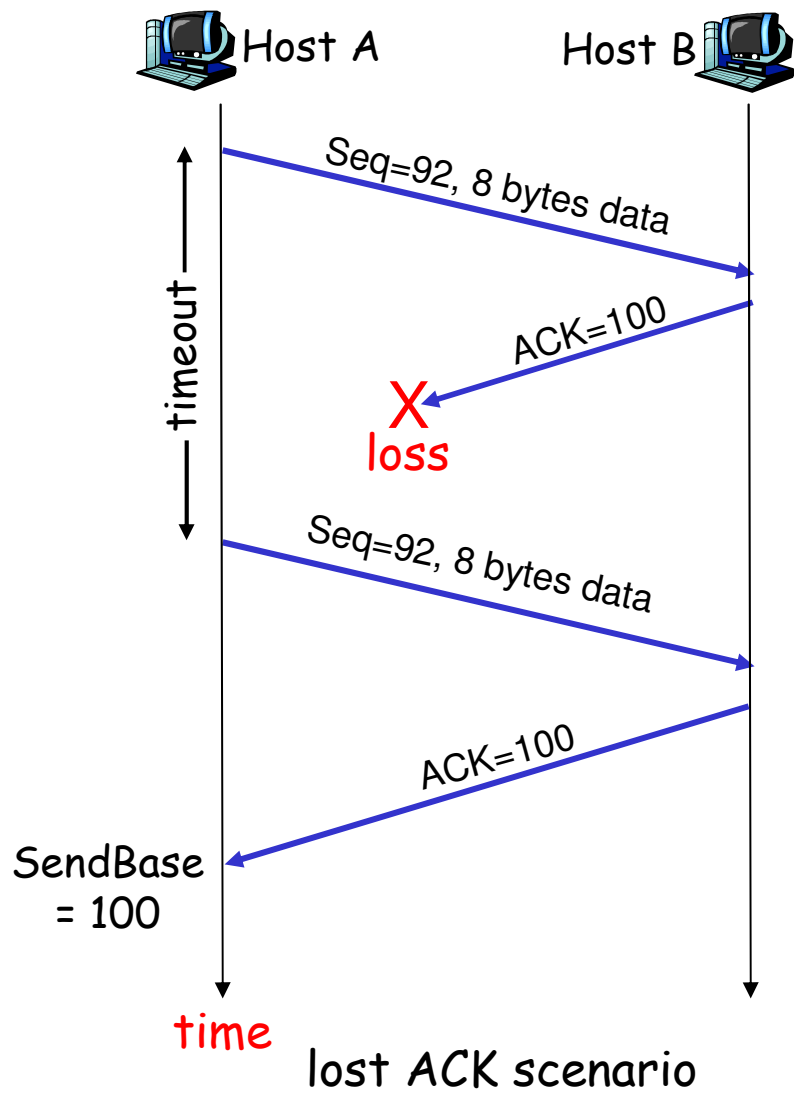
Comment:
• SendBase-1: last cumulatively ack'ed byte
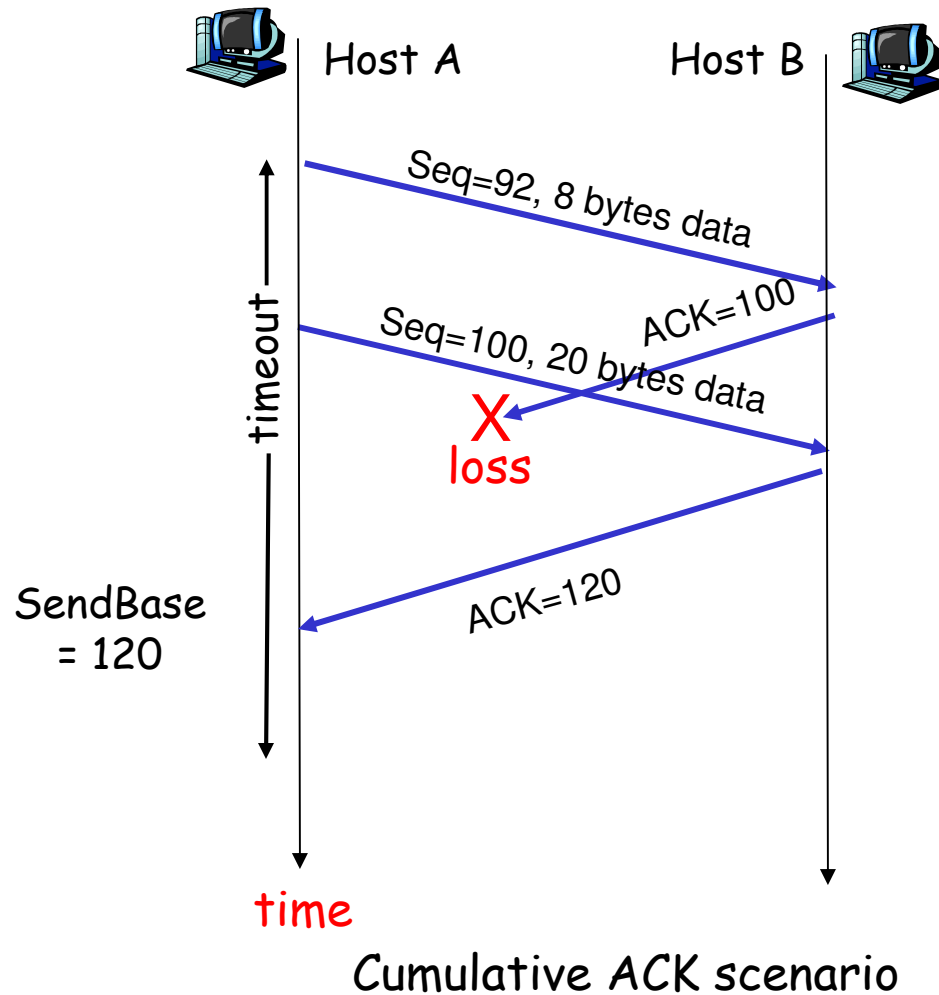Example:
• SendBase-1 = 71; y= 73, so the rcvr wants 73+ ;
y > SendBase, so that new data is acked

# TCP: retransmission scenarios



lost ACK scenario

premature timeout

# TCP retransmission scenarios (more)



Host A                    Host B

Seq=92, 8 bytes data

ACK=100

Seq=100, 20 bytes data

X
loss

timeout

SendBase
= 120

ACK=120

time

Cumulative ACK scenario

# TCP ACK generation [RFC 1122, RFC 2581]

Main motivation: performance

Favor piggybacking

| 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** |

Can advance source window

Duplicate ACK important feedback—more later
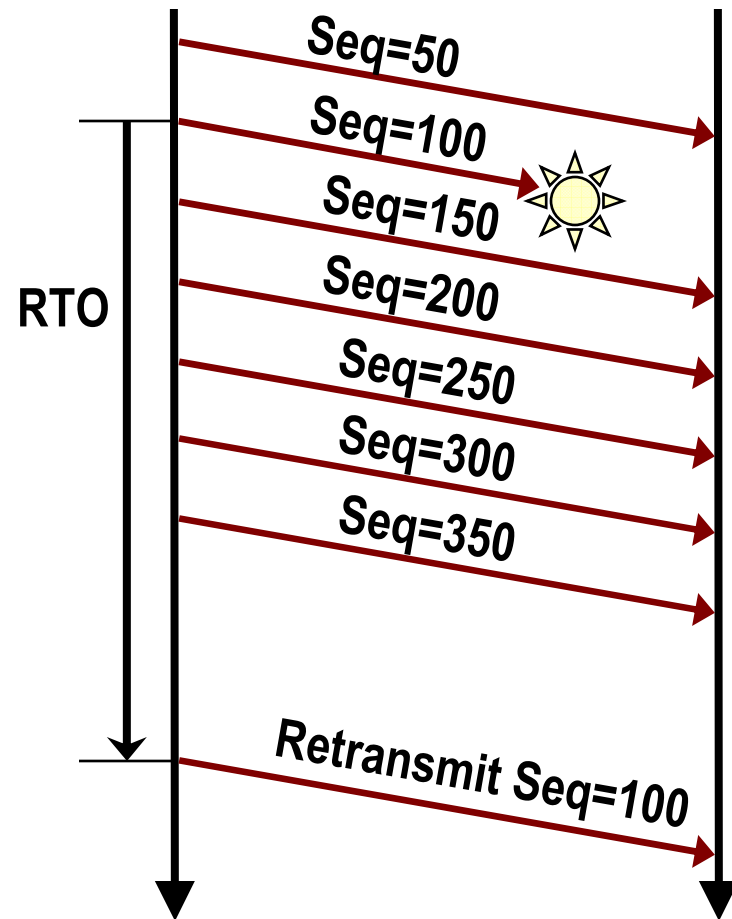
# So what is the TCP solution

r **Go-Back-N??**

r **Selective Repeat?**

r **A: An Hybrid solution.**

    m Possibility of buffering correctly received packets AND selective retransmission of packets, BUT NOT pure Selective Repeat, cumulative ACK, buffering not required (free implementation choice)

    m Shares some aspects with GBN BUT
- A single timer for the oldest unacked packet;
- when the timer experises ONLY that packet is retransmitted

# TCP: a reliable transport

r **TCP is a reliable protocol**

   m all data sent are guaranteed to be received

   m ***very important feature, as IP is unreliable network layer***

r **employs positive acknowledgement**

   m cumulative ack

   m selective ack may be activated when both peers implement it (use option) ⟶ **TCP SACKS**

r **does not employ negative ack**

   m error discovery via timeout (retransmission timer)

   m ...But "implicit NACK" is available (more later: fast retransmit)

# Need for implicit NACKs

➔ TCP does not support negative ACKs

➔ This can be a serious drawback

⇨ Especially in the case of single packet loss

➔ Necessary RTO expiration to start retransmit lost packet

⇨ As well as following ones!!

**May take too much time before retransmitting!!!**

➔ ISSUE: is there a way to have NACKs in an implicit manner????

Seq=50
Seq=100
Seq=150
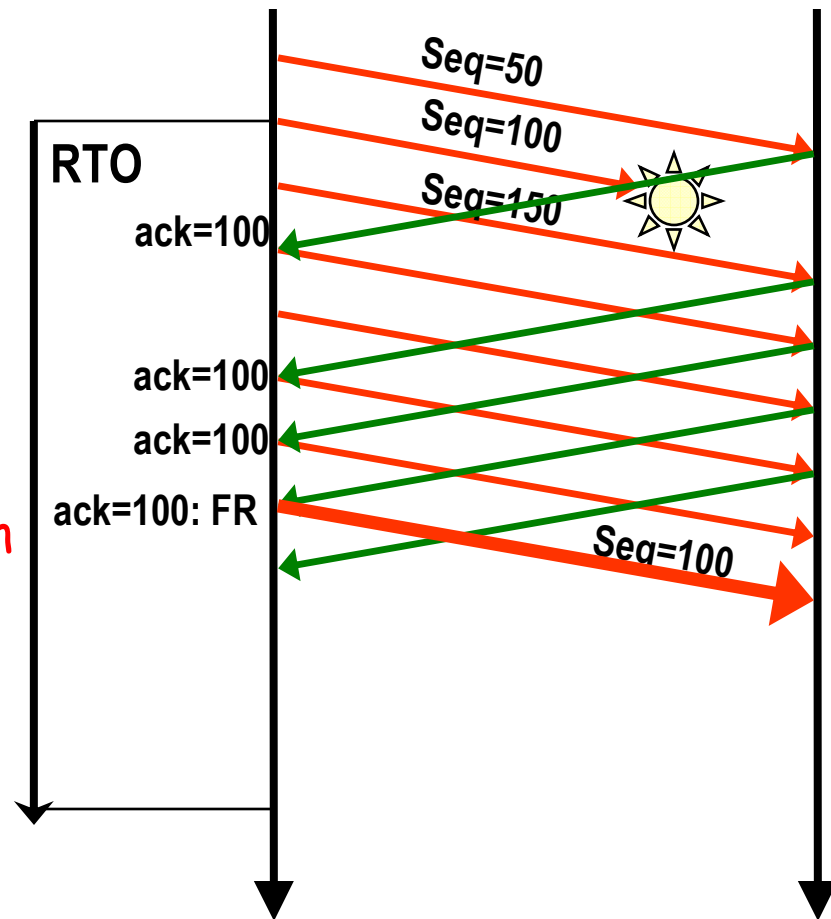Seq=200
Seq=250
Seq=300
Seq=350

RTO

Retransmit Seq=100

# The Fast Retransmit Algorithm

➔ Idea: use duplicate ACKs!

⇨ Receiver responds with an ACK every time it receives an out-of-order segment

⇨ ACK value = last correctly received segment

➔ FAST RETRANSMIT algorithm:

⇨ if 3 duplicate acks are received for the same segment, assume that the next segment has been lost. Retransmit it right away.

⇨ Helps if single packet lost. Not very effective with multiple losses

RTO

Seq=50

Seq=100

Seq=150

ack=100

ack=100

ack=100

ack=100: FR

Seq=100

# Fast retransmit algorithm:

```
event: ACK received, with ACK field value of y
        if (y > SendBase) {
             SendBase = y
             if (there are currently not-yet-acknowledged segments)
                  start timer
        }
        else {
             increment count of dup ACKs received for y
             if (count of dup ACKs received for y = 3) {
                  resend segment with sequence number y
             }
        }
```
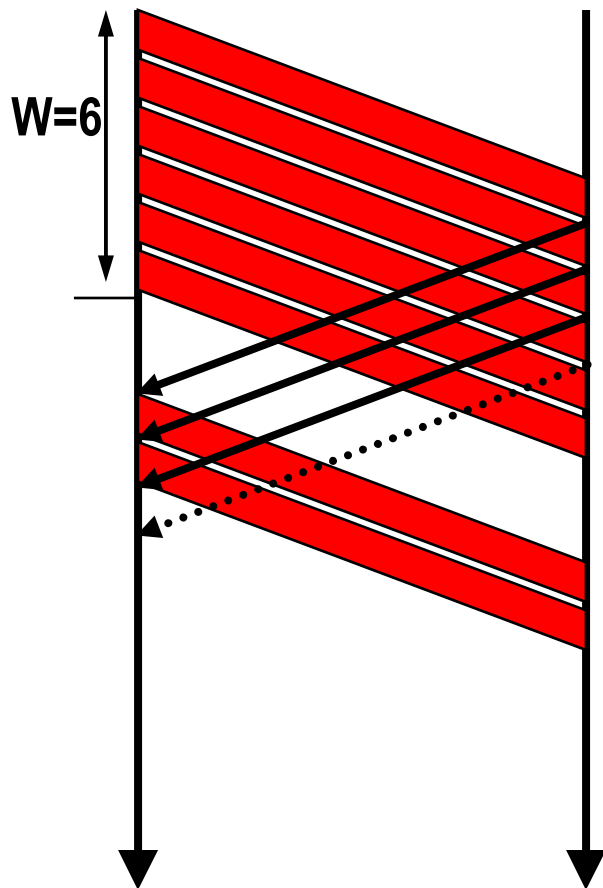
a duplicate ACK for
already ACKed segment

fast retransmit

# TCP mechanisms for:

r   flow control

r   congestion control

*Graphical examples (applet java) of several algorithms at:*
*http://www.ce.chalmers.se/~fcela/tcp-tour.html*
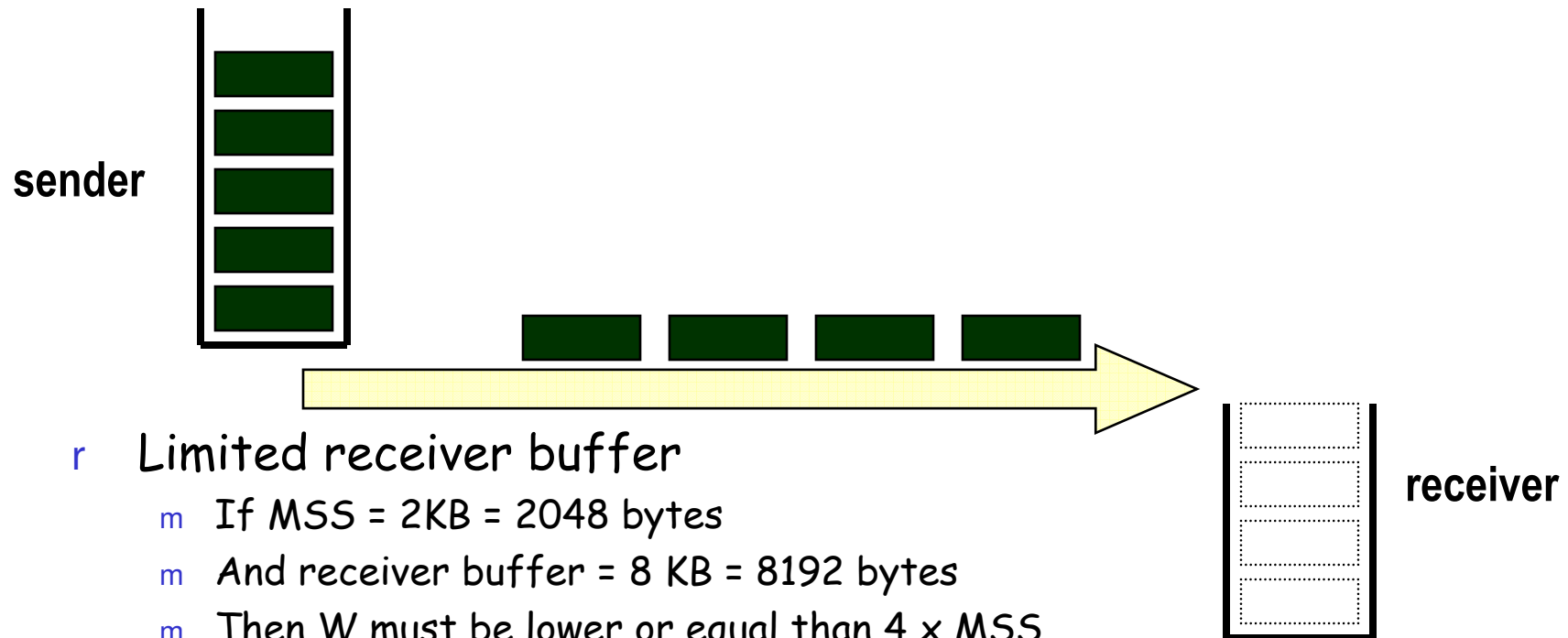
# TCP pipelining

W=6

r   **More than 1 segment "flying" in the network**

r   **Transfer efficiency increases with W**

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

r   **So, why an upper limit on W?**

m   **Esempio: flow control**

# Why flow control?

**sender**
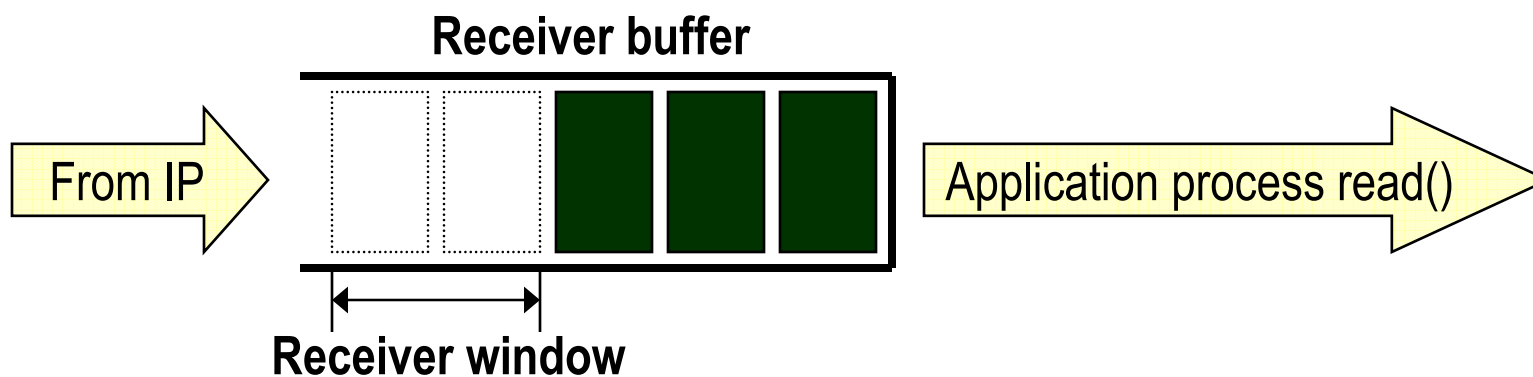
r  Limited receiver buffer

   m  If MSS = 2KB = 2048 bytes

   m  And receiver buffer = 8 KB = 8192 bytes

   m  Then W must be lower or equal than 4 x MSS

**receiver**

r  A possible implementation:

   m  During connection setup, exchange W value.

   m  *DOES NOT WORK. WHY?*

# Window-based flow control

➔ **receiver buffer capacity varies with time!**

  ⇨ Upon application process read()
  [asynchronous, not depending on OS, not predictable]



➔ MSS = 2KB = 2048 bytes
➔ Receiver Buffer capacity = 10 KB = 10240 bytes
➔ TCP data stored in buffer: 3 segments
➔ Receiver window = Spare room: 10-6 = 4KB = 4096 bytes
  ⇨ Then, at this time, W must be lower or equal than 2 x MSS

| 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 Window size field: used to advertise receiver's remaining storage capabilities

- m 16 bit field, on <u>every</u> packet
- m Measure unit: bytes, from 0 (included) to 65535
- m Sender rule:
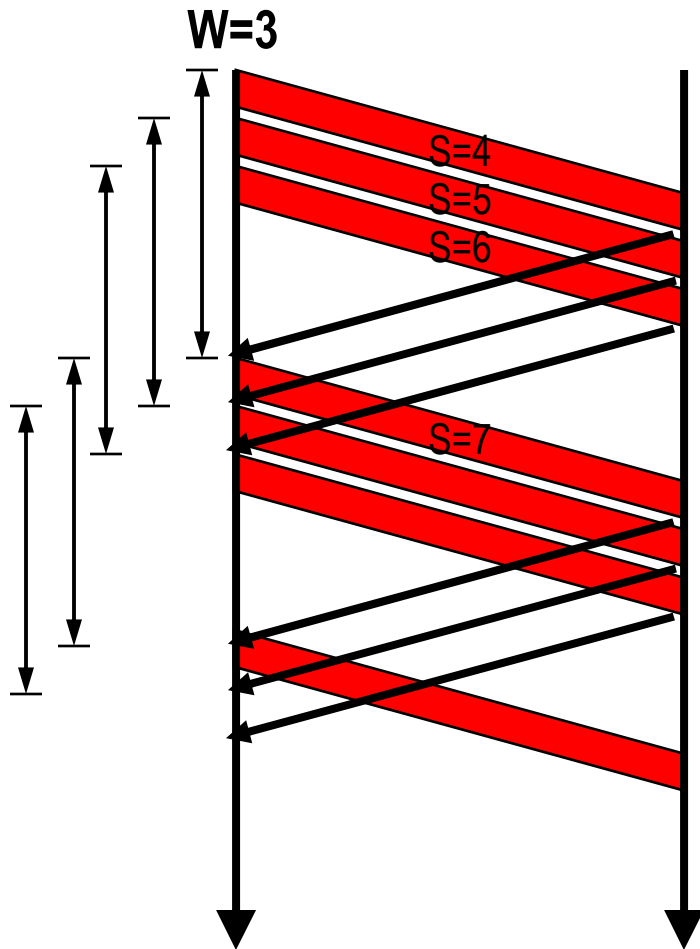
  **LastByteSent – LastByteAcked <= RcvWindow.**

- m W=2048 means:
  - I can accept other 2048 bytes since ack, i.e. bytes [ack, ack+W-1]
  - also means: sender may have 2048 bytes outstanding (in multiple segments)

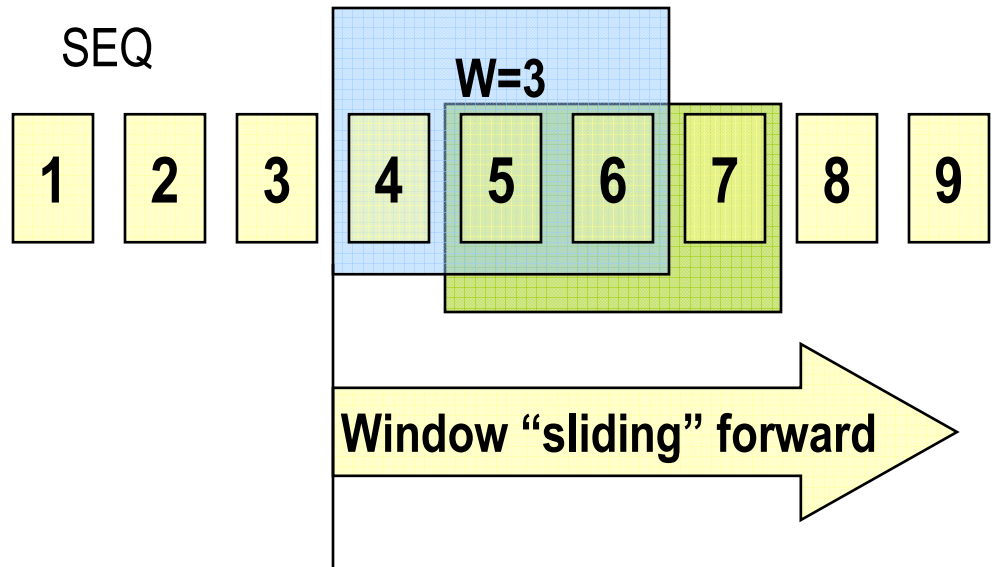# What is flow control needed for?

r   Window flow control guarantees receiver buffer to be able to accept outstanding segments.

r   When receiver buffer full, just send back win=0

r   in essence, flow control guarantees that transmission bit rate never exceed receiver rate
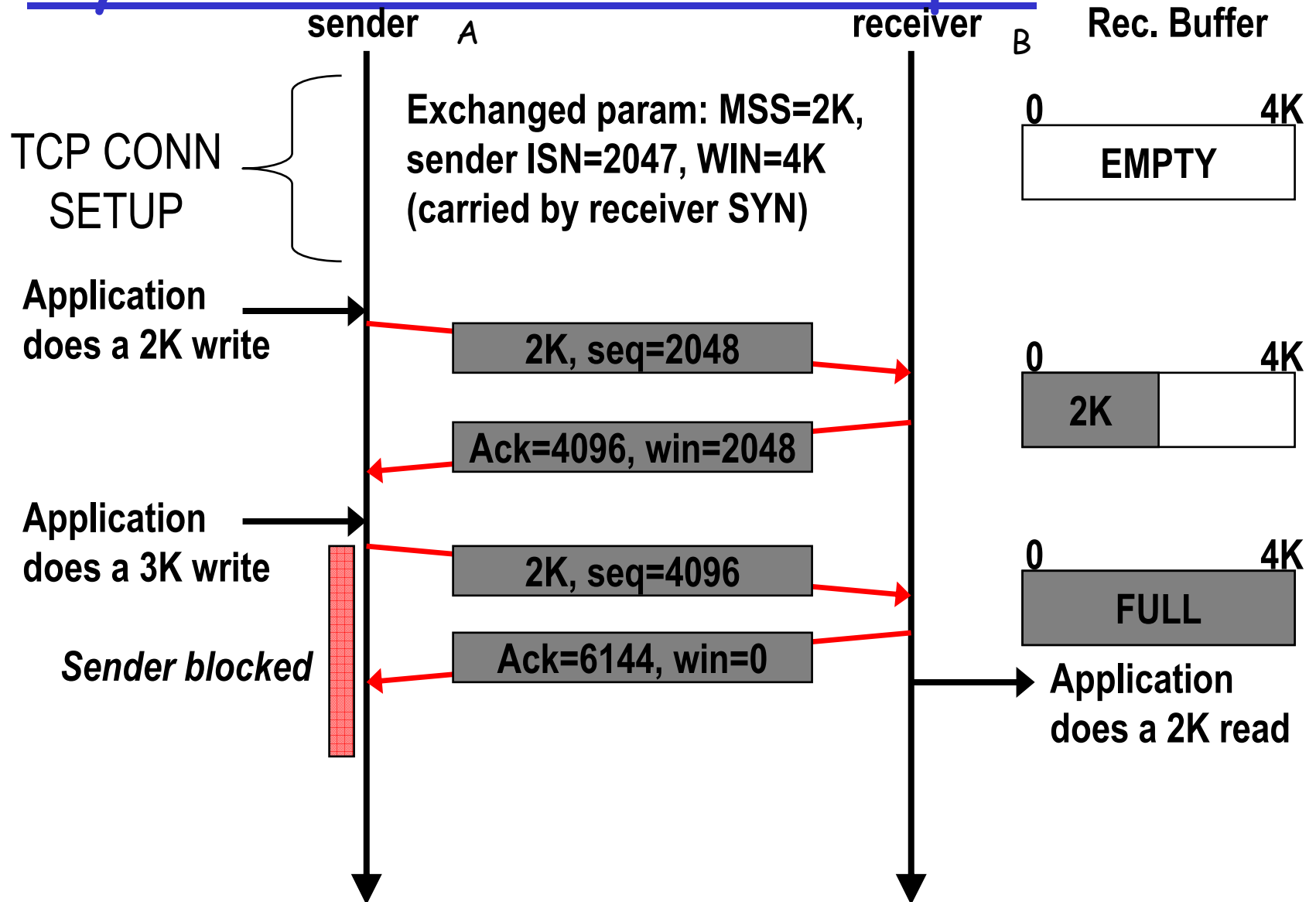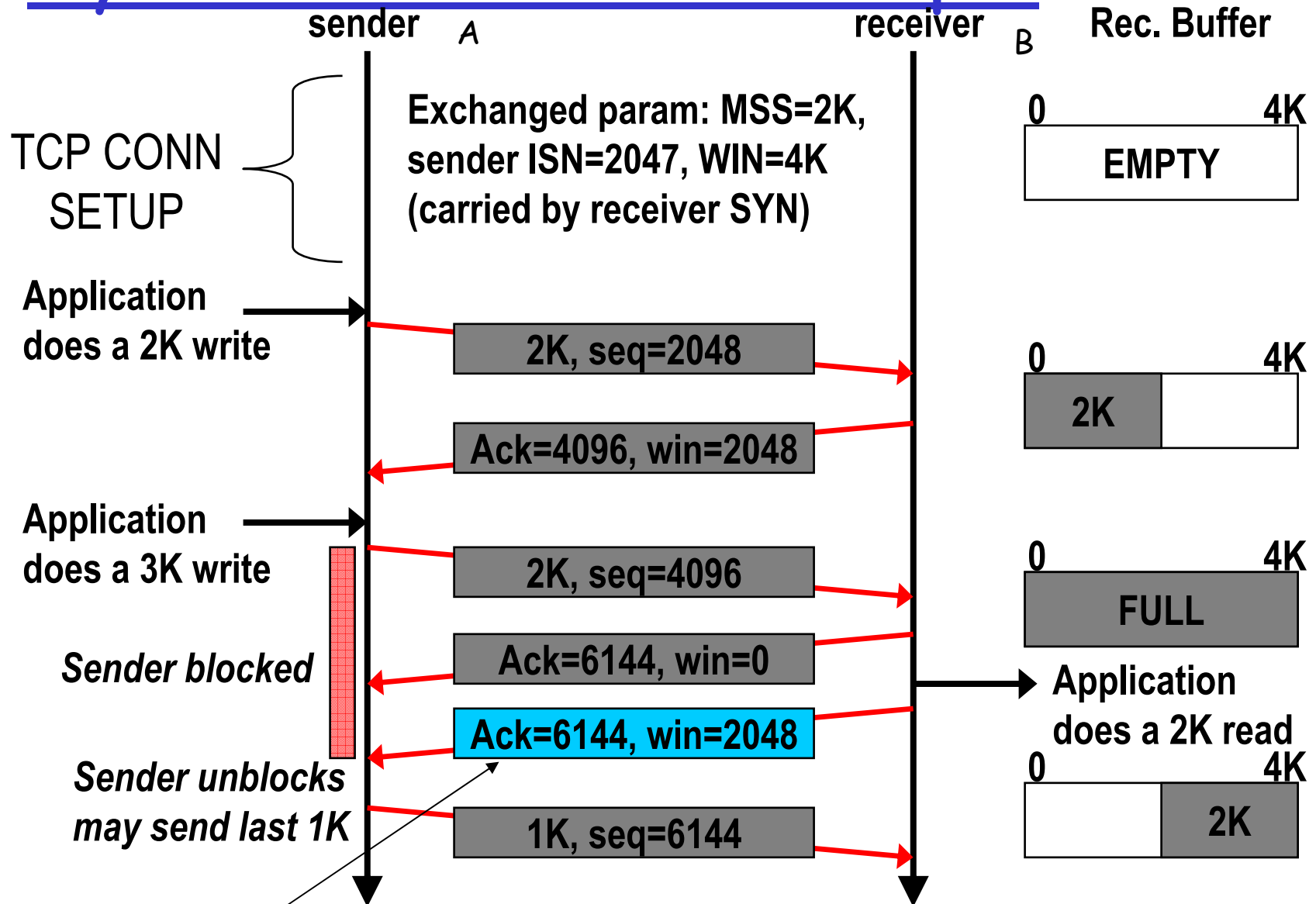
# Sliding window

W=3

S=4
S=5
S=6

S=7

**Dynamic window based reduces to pure sliding window when receiver app is very fast in reading data…**

SEQ

W=3

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |

**Window "sliding" forward**

# Dynamic window - example

sender A            receiver B      Rec. Buffer

TCP CONN SETUP
Exchanged param: MSS=2K, sender ISN=2047, WIN=4K (carried by receiver SYN)

0                                4K
| EMPTY |

Application does a 2K write

2K, seq=2048

0                                4K
| 2K |        |

Ack=4096, win=2048

Application does a 3K write

2K, seq=4096

0                                4K
| FULL |

Sender blocked

Ack=6144, win=0

Application does a 2K read

# Dynamic window - example

sender A                                    receiver B        Rec. Buffer

TCP CONN SETUP
Exchanged param: MSS=2K, sender ISN=2047, WIN=4K (carried by receiver SYN)

0                    4K
EMPTY

Application does a 2K write

2K, seq=2048

0                    4K
2K

Ack=4096, win=2048

Application does a 3K write

2K, seq=4096

0                    4K
FULL

Sender blocked

Ack=6144, win=0

Application does a 2K read

Ack=6144, win=2048
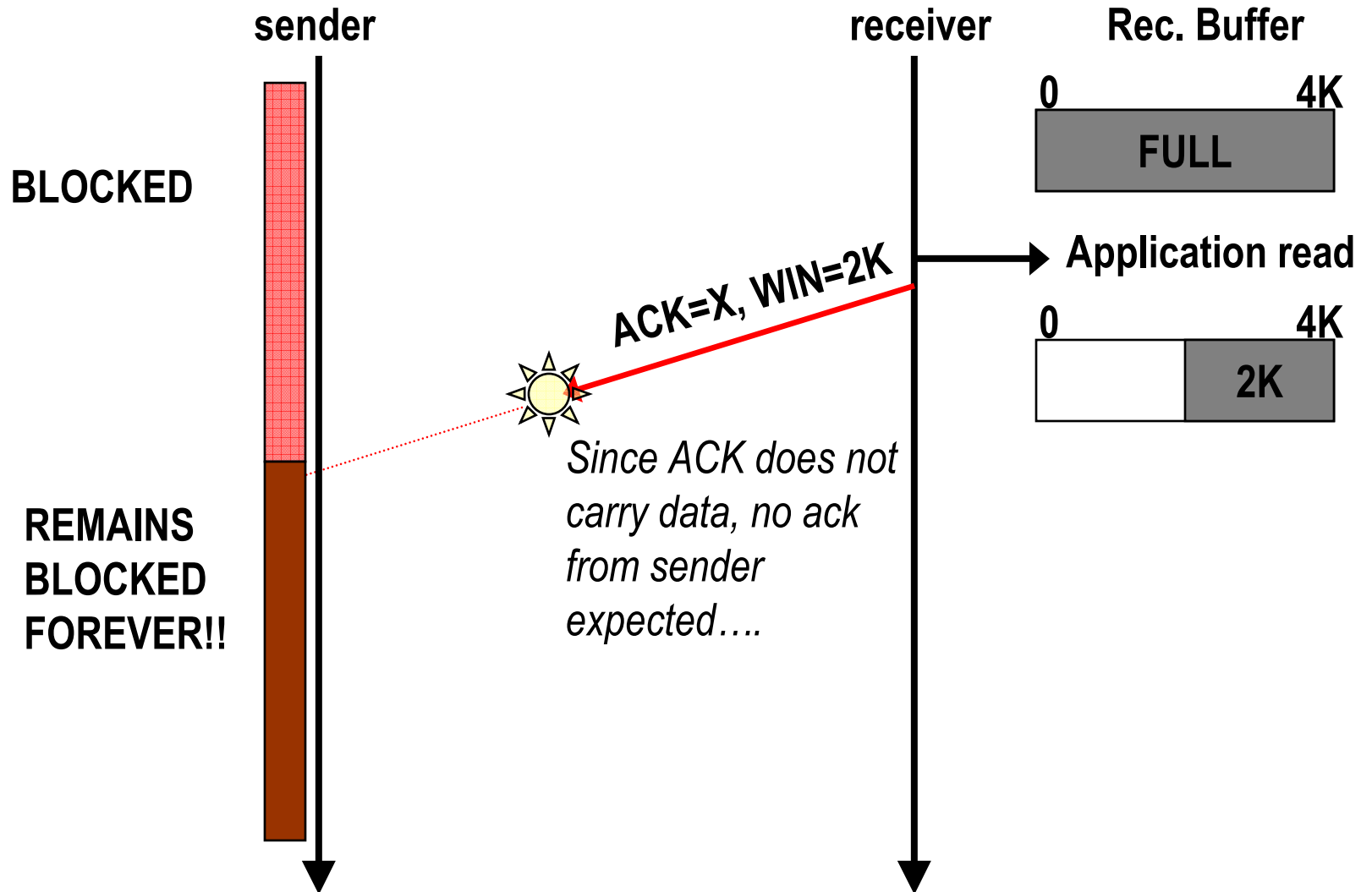
Sender unblocks may send last 1K

1K, seq=6144

0                    4K
2K

Piggybacked in a packet sent from B to A

Window thus source rate limited by reading speed and buffer size at the receiver

23

# Blocked sender deadlock problem

sender  receiver  Rec. Buffer

**BLOCKED**

0  4K
**FULL**

Application read

ACK=X, WIN=2K

0  4K
**2K**

*Since ACK does not carry data, no ack from sender expected….*
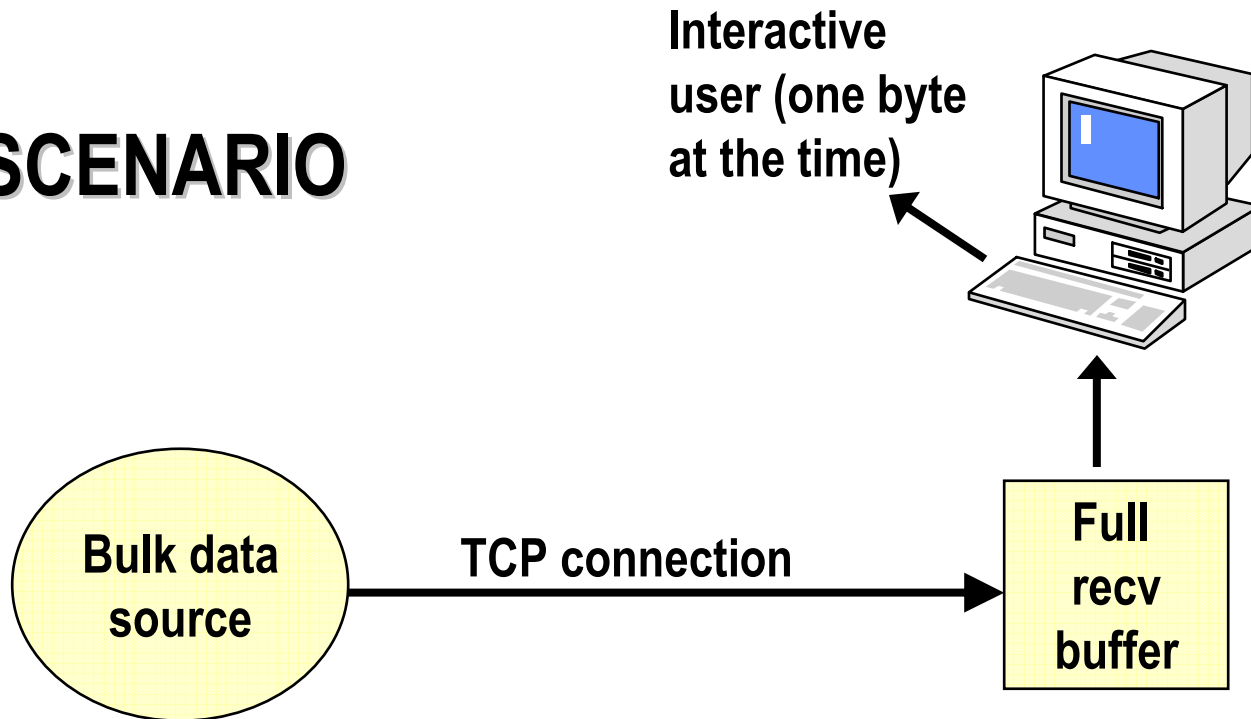
**REMAINS BLOCKED FOREVER!!**

# Solution: Persist timer
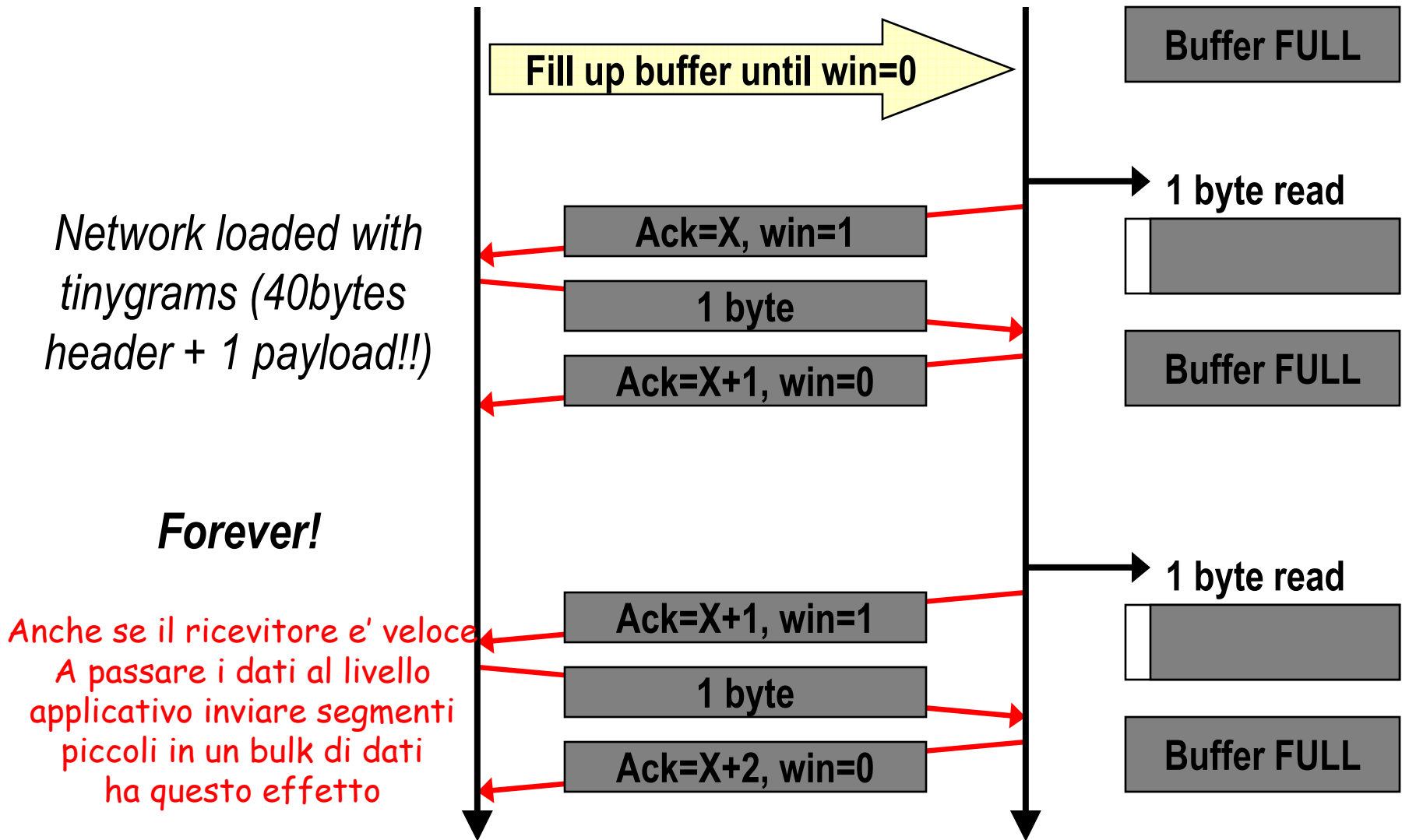
r   When win=0 (blocked sender), sender starts a "persist" timer
  - Initially 500ms (but depends on implementation)

r   When persist timer elapses AND no segment received during this time, sender transmits "probe"

  m   Probe = 1byte segment; makes receiver reannounce next byte expected and window size
  - this feature necessary to break deadlock
  - if receiver was still full, rejects byte
  - otherwise acks byte and sends back actual win

r   Persist time management (exponential backoff):

  m   Doubles every time no response is received
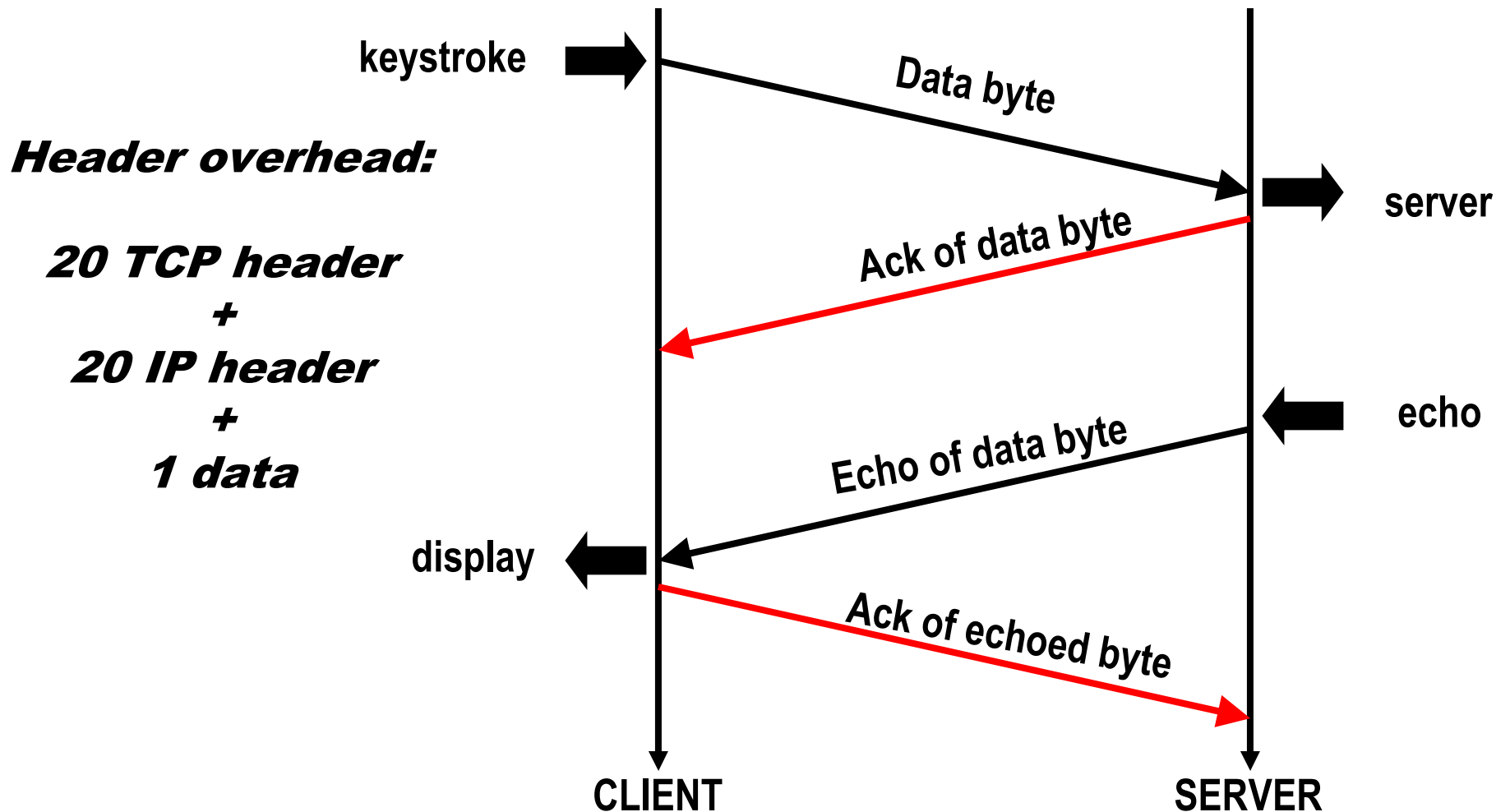  m   Maximum = 60s

# The silly window syndrome

**SCENARIO**

Interactive user (one byte at the time)

Bulk data source →(TCP connection)→ Full recv buffer → Interactive user

# The silly window syndrome

Fill up buffer until win=0

Buffer FULL

*Network loaded with tinygrams (40bytes header + 1 payload!!)*

1 byte read

Ack=X, win=1

1 byte

Buffer FULL

Ack=X+1, win=0

**Forever!**

Anche se il ricevitore e' veloce A passare i dati al livello applicativo inviare segmenti piccoli in un bulk di dati ha questo effetto

1 byte read

Ack=X+1, win=1

1 byte

Buffer FULL

Ack=X+2, win=0

# Silly window solution

r Problem discovered by David Clark (MIT), 1982

r easily solved, by preventing receiver to send a window update for 1 byte

r rule: send window update when:

- receiver buffer can handle a whole MSS

  or

- half received buffer has emptied (if smaller than MSS)

r sender also may apply rule

- by waiting for sending data when win low

# Interactive applications

keystroke ➡️

Header overhead:

20 TCP header
+
20 IP header
+
1 data

Data byte

→ server

Ack of data byte

echo ⬅️

Echo of data byte

display ⬅️

Ack of echoed byte

CLIENT                                          SERVER

*Interactive apps: create some tricky situations….*

# Nagle's algorithm
## (RFC 896, 1984)

UNLESS MSS data (or at least half the window size bytes) are ready to be transmitted

**NAGLE RULE: inhibit sending new segments if any previously transmitted data unacked**

*self-clocking algorithm:*

*on LANs, plenty of tynigrams*

*on slow WANs, data aggregation*

1 byte

WAIT

1 byte

ack

WAIT

2 byte

ack

L'idea e' che si possoa trasmettere segmenti piu' lunghi avendo bufferizzato dati nel frattempo

Transport Layer   3-30

# PUSH flag

| 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 **Used to notify**

    m TCP sender to send data

        • but for this an header flag NOT needed! Sufficient a "push" type indication in the TCP sender API

    m TCP receiver to pass received data to the application

# Urgent data

| Source port | | | | | | | Destination port | |
|---|---|---|---|---|---|---|---|---|
| 32 bit Sequence number | | | | | | | | |
| 32 bit acknowledgement number | | | | | | | | |
| Header length | 6 bit Reserved | URG | ACK | PSH | RST | SYN | FIN | Window size |
| checksum | | | | | | | Urgent pointer | |

r   URG on: notifies rx that "urgent" data placed in segment.

r   When URG on, *urgent pointer* contains position of *the last octet* of urgent data

- *indeed it contains the positive offset from the segment sequence number*
- *and the position of the first octet of urgent data? No way to specify it!*
- *Changed wrt RFC 793*

r   receiver is expected to pass all data up to urgent ptr to app

- *interpretation of urgent data is left to the app*

r   typical usage: ctrlC (interrupt) in rlogin & telnet; abort in FTP

r   urgent data is a <u>second</u> exception to blocked sender

# Chapter 3 outline

# TCP connection

Application (client)

Socket

TCP software

**State variables:**
- **conn status**
- **MSS**
- **windows**
- **...**

**buffer space**

normally 4 to 16 Kbytes

64+ Kbytes possible

*"Logical" connection*
*only end hosts are aware!*

**TCP**

INTERNET

Application (server)

Socket

TCP software

*Connection described by client&server status*
*Connection SET-UP duty:*
*1) initializes state variables* → Transmission control block
*2) reserves buffer space*

Contains also info on: sockets, pointers to the users' send and receive buffers, to the retransmit queue and to the current segment

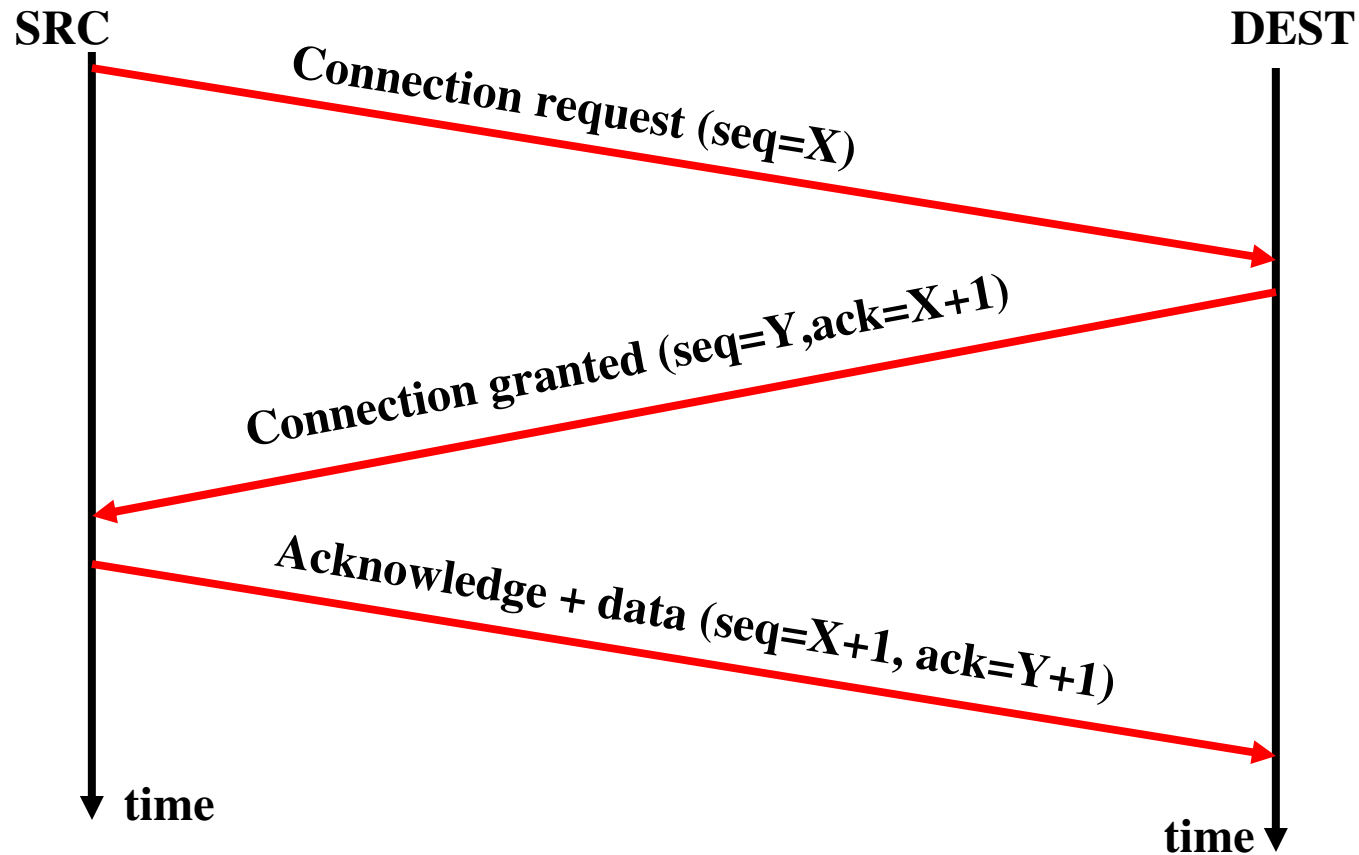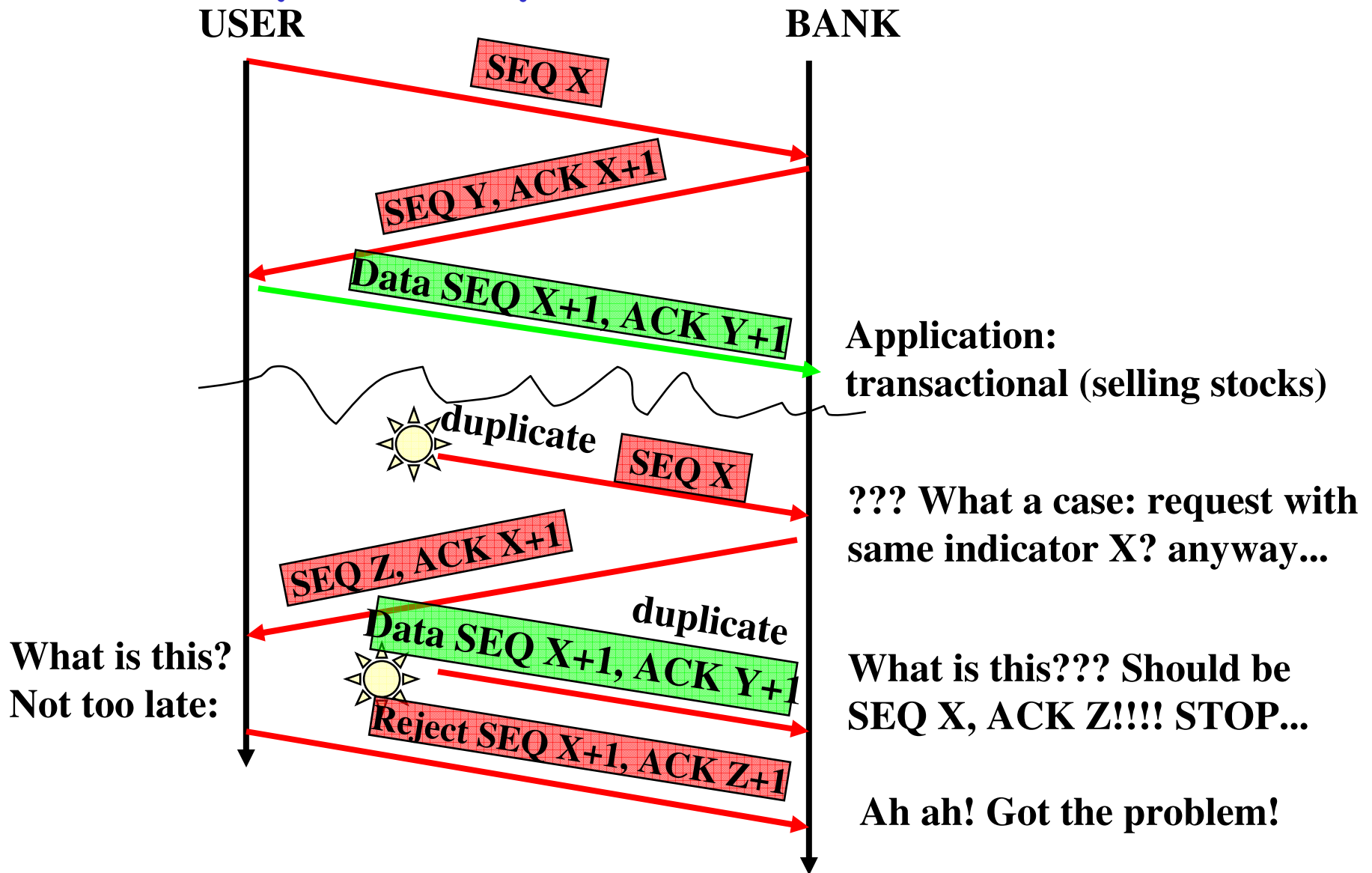# Connection establishment: simplest approach (non TCP)

Connection request

Connection granted

Transmit data

time

time

# Delayed duplicate problem

USER                                    BANK

REQ

Data

Application:
transactional (sell
100000$ stocks)

duplicate

REQ

ACK

What is this?
Oh my God!
Too late!!!

duplicate
Data

Selling other 100000$
stocks!!!!!

# Solution: three way handshake
## Tomlinson 1975



SRC                     DEST

Connection request (seq=X)

Connection granted (seq=Y,ack=X+1)

Acknowledge + data (seq=X+1, ack=Y+1)

time                     time

# Delayed duplicate detection

USER            BANK

SEQ X

SEQ Y, ACK X+1

Data SEQ X+1, ACK Y+1

Application:
transactional (selling stocks)

duplicate

SEQ X

??? What a case: request with
same indicator X? anyway...

SEQ Z, ACK X+1

duplicate

Data SEQ X+1, ACK Y+1

What is this?
Not too late:

What is this??? Should be
SEQ X, ACK Z!!!! STOP...

Reject SEQ X+1, ACK Z+1

Ah ah! Got the problem!

Disaster could not be avoided with a two-way handshake

| 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 SYN (synchronize sequence numbers): used to open connection

  m SYN present: this host is setting up a connection

  m SEQ with SYN: means initial sequence number (ISN)

  m data bytes numbered from ISN+1.

r FIN: no more data to send

  m used to close connection

  *...more later about connection closing...*

# Three way handshake in TCP



Full *duplex connection: opened in both ways*
*SRC: performs ACTIVE OPEN*
*DEST: Performs PASSIVE OPEN*

# Initial Sequence Number

r Should change in time

   m RFC 793 (but not all implementations are conforming) suggests to generate ISN as a sample of a 32 bit counter incrementing at 4μs rate (4.55 hour to wrap around—Maximum Segment Lifetime much shorter)

r transmitted whenever SYN (Synchronize sequence numbers) flag active

   m note that both src and dest transmit THEIR initial sequence number (remember: full duplex)

r Data Bytes numbered from ISN+1

   m necessary to allow SYN segment ack

# Forbidden Region

r Obiettivo: due sequence number identici non devono trovarsi in rete allo stesso tempo



Forbidden region

r Aging dei pacchetti→ dopo un certo tempo MSL (Maximum Segment Lifetime) i pacchetti eliminati dalla rete

r Sequence numbers basati sul clock

r Un ciclo del clock circa 4 ore; MSL circa 2 minuti.

r → Se non ci sono crash che fanno perdere il valore dell'ultimo sequence number usato NON ci sono problemi (si riusa lo stesso sequence number ogni 4 ore circa, quando il segmento precedentemente trasmesso con quel sequence number non è più in rete)

r → Cosa succede nel caso di crash? RFC suggerisce l'uso di un 'periodo di silenzio' in cui non vengono inviati segmenti dopo il riavvio pari all'MSL (per evitare che pacchetti precedenti connessioni siano in giro).

# Maximum Segment Size - MSS

r   Announced at setup by both ends.

r   Lower value selected (indeed min of lower value and largest size permitted by IP layer).

r   MSS sent in the Options header of the SYN segment

   m   clearly cannot (=ignored if happens) send MSS in a non SYN segment, as connection has been already setup

   m   when SYN has no MSS, default value 536 used

r   goal: the larger the MSS, the better...

   m   until fragmentation occurs

   m   e.g. if host is on ethernet, sets MSS=1460

      • 1500 max ethernet size - 20 IP header - 20 TCP header

# MSS advertise

**CLIENT (C_MSS)**                                    **SERVER (S_MSS)**

Conn request (C_MSS, SYN, seq=C_ISN)

*If (S_MSS<C_MSS)*
*MSS = S_MSS;*
*else MSS = C_MSS;*

Conn granted (MSS, SYN, seq=S_ISN, ack=C_ISN+1)

*Use recv MSS*

Acknowledge (seq=C_ISN+1,ack=S_ISN+1)

**time**                                              **time**

***Does not avoid fragmentation to occur WITHIN the network!!***

# TCP Connection Management:Summary

## Recall: TCP sender, receiver establish "connection" before exchanging data segments

r  initialize TCP variables:

  m  seq. #s

  m  buffers, flow control info (e.g. `RcvWindow`)

  m  MSS

r  *client:* connection initiator

```
Socket clientSocket = new
Socket("hostname","port

number");
```

r  *server:* contacted by client

```
Socket connectionSocket =
welcomeSocket.accept();
```

## Three way handshake:

Step 1: client host sends TCP SYN segment to server

  m  specifies initial seq #

  m  no data

Step 2: server host receives SYN, replies with SYNACK segment

  m  server allocates buffers

  m  specifies server initial seq. #

Step 3: client receives SYNACK, allocates buffer and variables,replies with ACK segment, which may contain data

Per chiudere la connessione uno dei due estremi invia un messaggio con FIN flag a 1 a cui l'altro estremo della connessione risponde con ACK
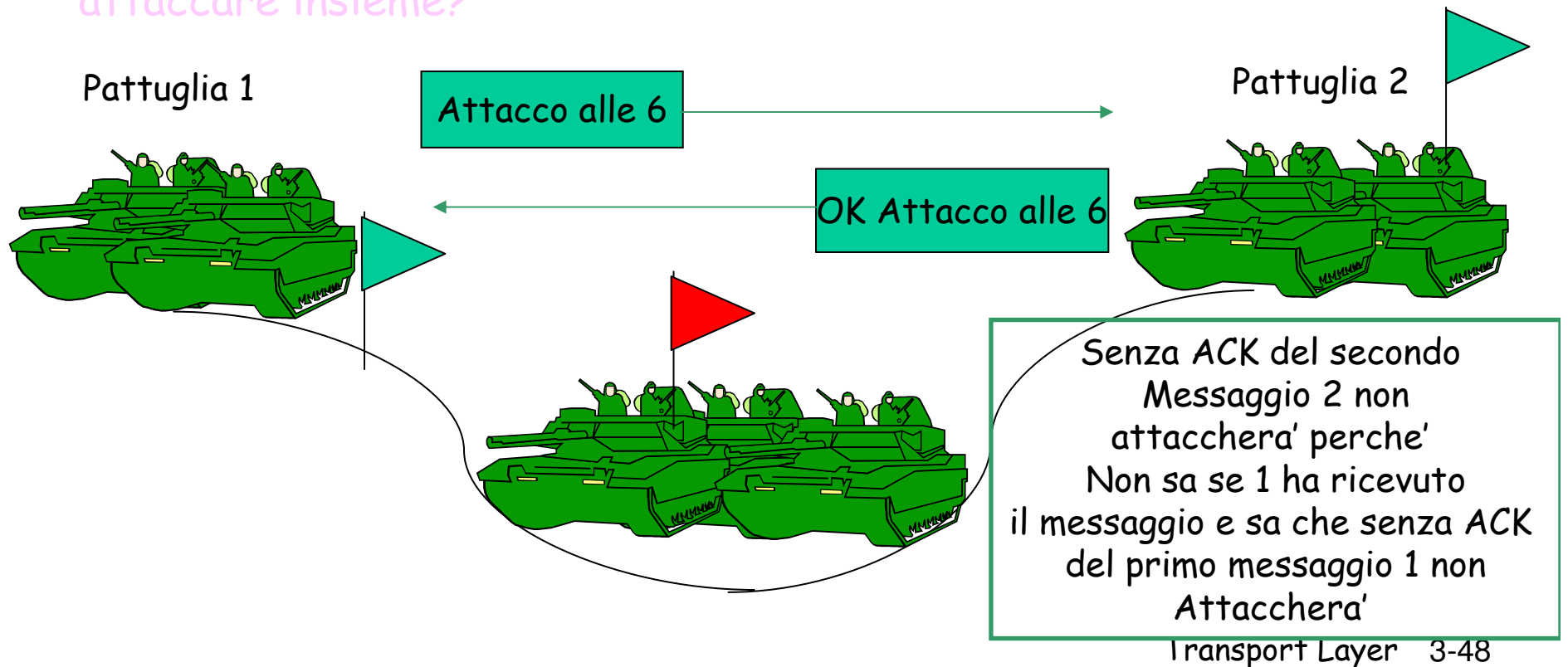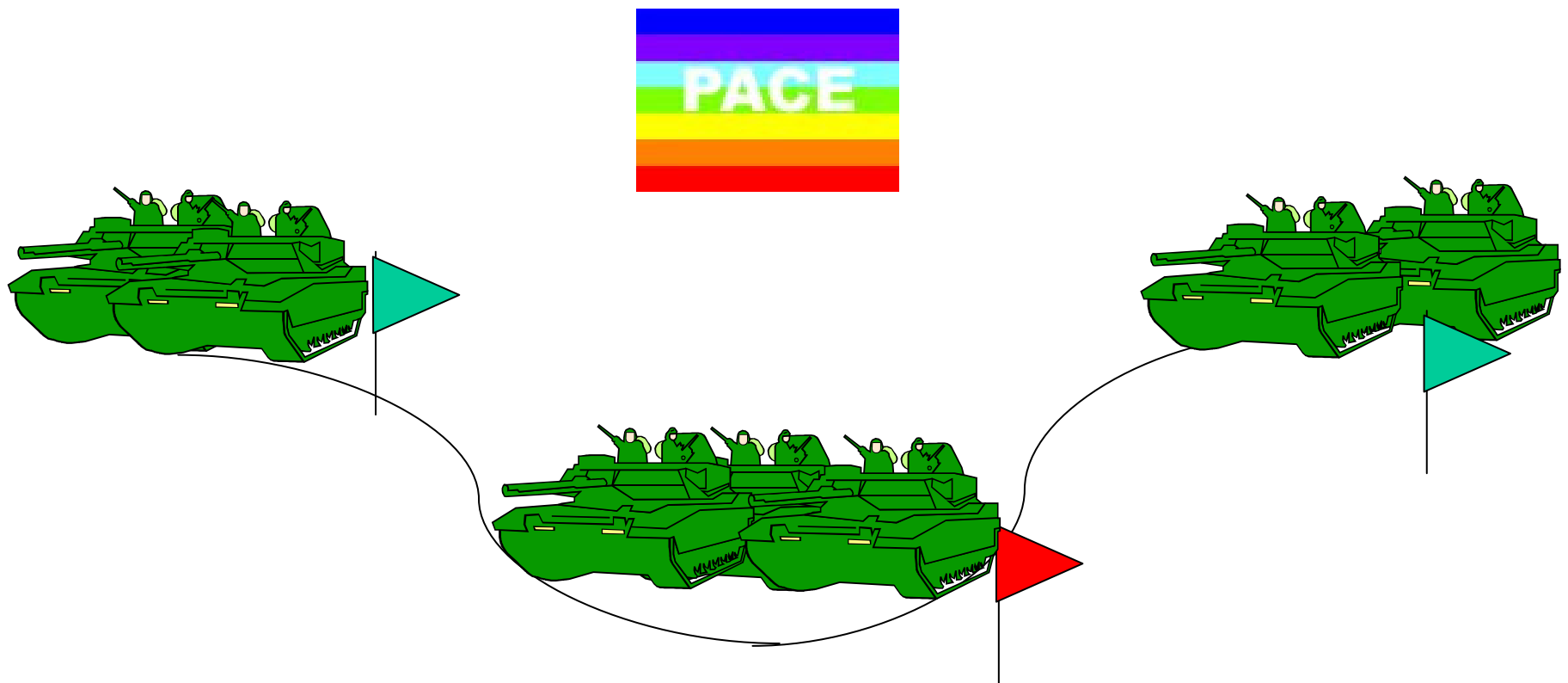
# Problema dei due eserciti

r  L'esercito rosso e' globalmente più debole. Se le due pattuglie verdi attaccano insieme lo sconfiggono, altrimenti perdono. Possono scambiarsi messaggi relativi all'orario in cui attaccheranno e di ACK di un messaggio ricevuto. I messaggeri che li portano possono pero' essere catturati e quindi il messaggio può non arrivare correttamente a destinazione. Come fanno a mettersi d'accordo per attaccare insieme?
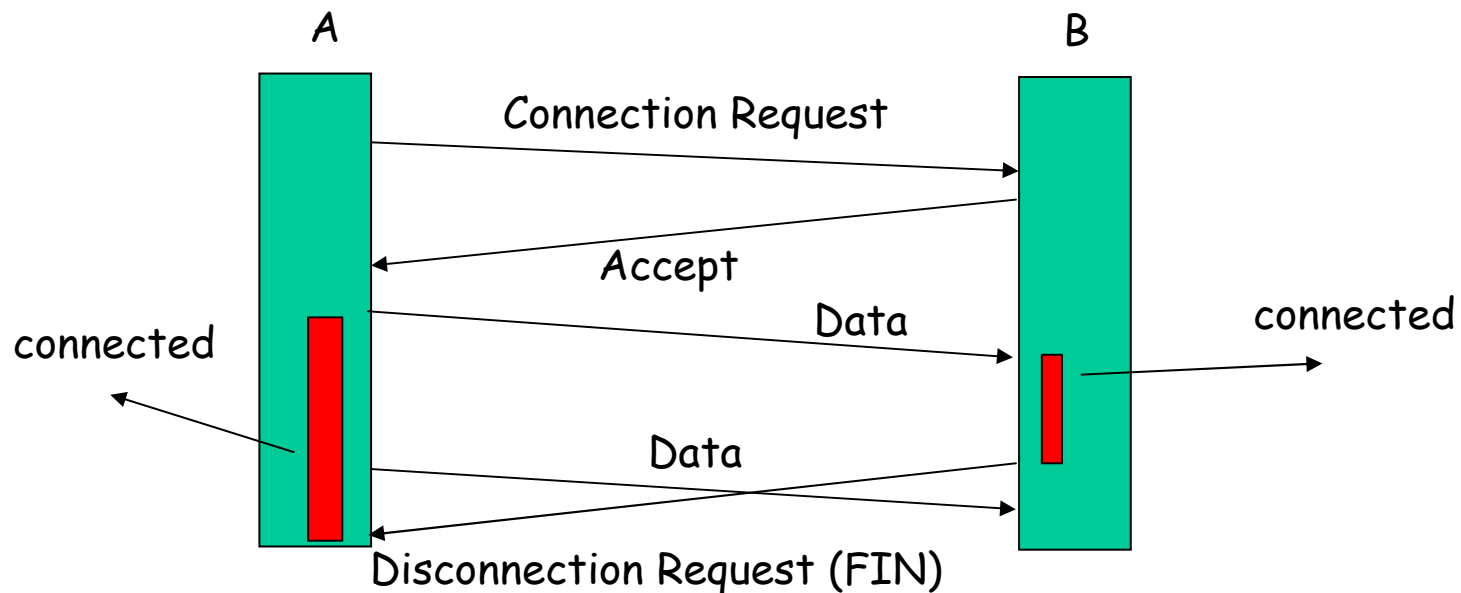
# Problema dei due eserciti

r  L'esercito rosso e' globalmente più debole. Se le due pattuglie verdi attaccano insieme lo sconfiggono, altrimenti perdono. Possono scambiarsi messaggi relativi all'orario in cui attaccheranno e di ACK di un messaggio ricevuto. I messaggeri che li portano possono pero' essere catturati e quindi il messaggio può non arrivare correttamente a destinazione. Come fanno a mettersi d'accordo per attaccare insieme?

Pattuglia 1

Pattuglia 2

Attacco alle 6

Senza ACK 1 non
Attacchera' perche'
Non sa se 2 ha ricevuto
Il messaggio

# Problema dei due eserciti

r L'esercito rosso e' globalmente più debole. Se le due pattuglie verdi attaccano insieme lo sconfiggono, altrimenti perdono. Possono scambiarsi messaggi relativi all'orario in cui attaccheranno e di ACK di un messaggio ricevuto. I messaggeri che li portano possono pero' essere catturati e quindi il messaggio può non arrivare correttamente a destinazione. Come fanno a mettersi d'accordo per attaccare insieme?

Pattuglia 1

Pattuglia 2

Attacco alle 6

OK Attacco alle 6

Senza ACK del secondo
Messaggio 2 non
attacchera' perche'
Non sa se 1 ha ricevuto
il messaggio e sa che senza ACK
del primo messaggio 1 non
Attacchera'

# Problema dei due eserciti

r   In generale: se N scambi di messaggi /Ack etc. necessari a raggiungere la certezza dell'accordo per attaccare allora cosa succede se l'ultimo messaggio 'necessario' va perso?

r   →E' impossibile raggiungere questa certezza. Le due pattuglie non attaccheranno mai!!

# Problema dei due eserciti: cosa ha a che fare con le reti e TCP??

r Chiusura di una connessione. Vorremmo un accordo tra le due peer entity o rischiamo di perdere dati.



A pensa che il secondo pacchetto sia stato ricevuto. La connessione e'
Stata chiusa da B prima che ciò avvenisse→ secondo pacchetto perso!!!

# Quando si può dire che le due peer entity abbiano raggiunto un accordo???

r  Problema dei due eserciti!!!



**Ma se l'ACK va perso????**

**Soluzione: si e' disposti a correre piu' rischi quando si butta giu' una connessione di quando si attacca un esercito nemico. Possibili malfunzionamenti. Soluzioni 'di recovery' in questi casi**

# TCP Connection Management (cont.)

**Since it is impossible to solvethe proble use simple solution: two way handshake**

## Closing a connection:

client closes socket:
`clientSocket.close();`

Step 1: client end system
sends TCP FIN control
segment to server

Step 2: server receives
FIN, replies with ACK.
Closes connection, sends
FIN.

client          server

close                      FIN

                    ACK

                    FIN                close

timed wait          ACK

closed

# TCP Connection Management (cont.)

**Step 3:** client receives FIN, replies with ACK.

  m Enters "timed wait" - will respond with ACK to received FINs

**Step 4:** server, receives ACK. Connection closed.

client      server

closing

FIN

closing

ACK

FIN

timed wait

ACK

closed

closed

# Connection states - Client

# Connection States - Server

# Why TIME_WAIT?

r MSL (Maximum Segment Lifetime): **maximum time a segment can live in the Internet**

- no timers on IP packets! Only hop counter
- RFC 793 specifies MSL=2min, but each implementation has its own value (from 30s to 2min)

r TIME_WAIT state: **2 x MSL**

m **allows to "clean" the network of delayed packets belonging to the connection**

m **2xMSL because a lost FIN_ACK implies a new FIN from server**

r during TIME_WAIT conn sock pair reserved

m **many implementations even more restictive (local port non reusable)**

m **clearly this may be a serious problem when restarting server daemon (must pause from 1 to 4 minutes...)**

| 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 **RST (Reset)**

   m sent whenever a segment arrives and does not apparently belong to the connection

   m typical RST case: connection request arriving to port not in use

r **Sending RST within an active connection:**

   m allows **aborting release** of connection (versus **orderly release**)

      · any queued data thrown away

      · receiver of RST can notify app that abort was performed at other end

# Chapter 3 outline

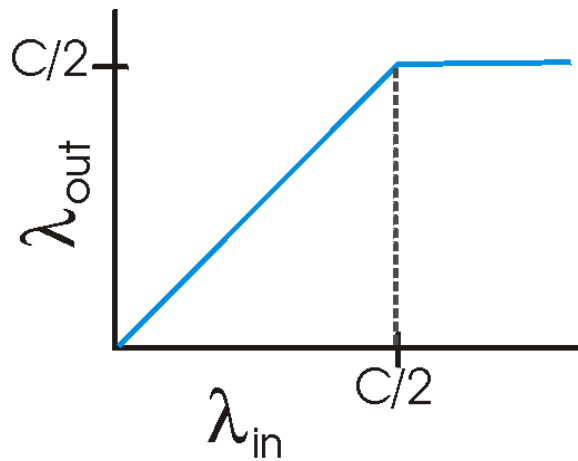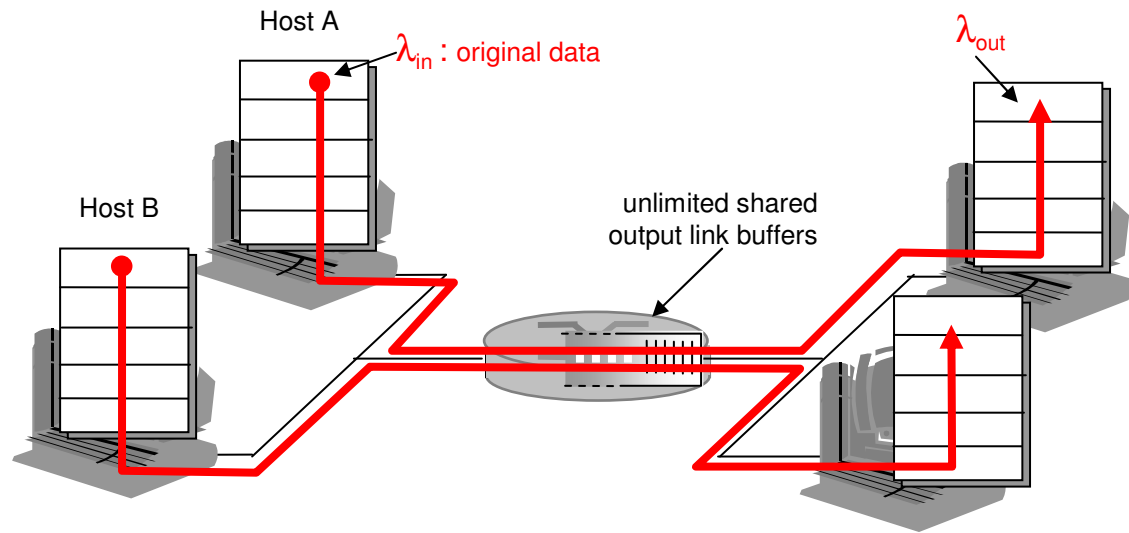# Principles of Congestion Control

## Congestion:

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

r   different from flow control!

r   manifestations:

   m lost packets (buffer overflow at routers)

   m long delays (queueing in router buffers)

r   a top-10 problem!

# Causes/costs of congestion: scenario 1
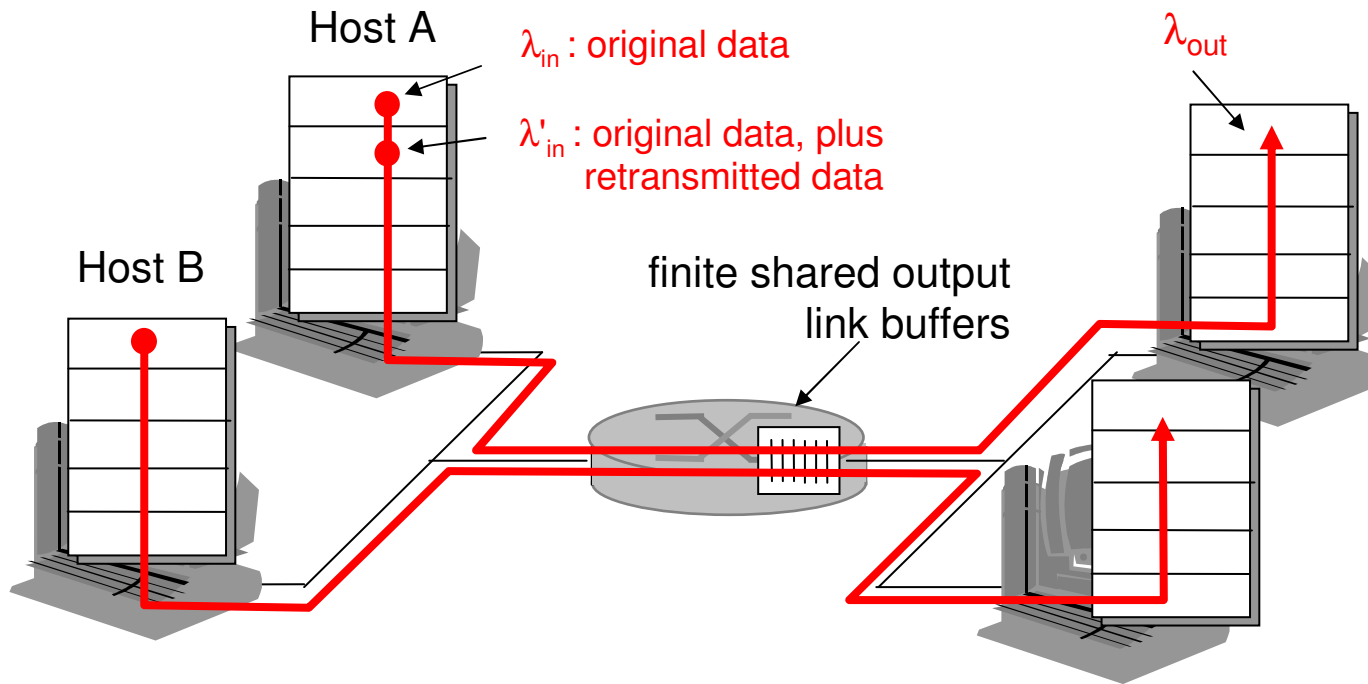
r two senders, two receivers

r one router, infinite buffers

r no retransmission

Host A    $\lambda_{in}$ : original data                                    $\lambda_{out}$

Host B    unlimited shared output link buffers
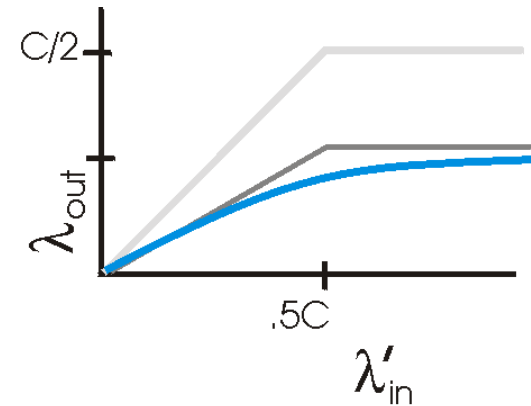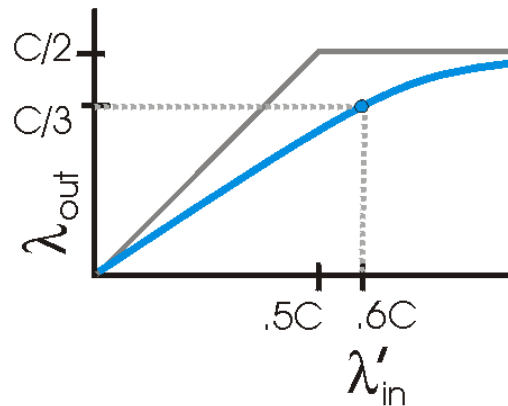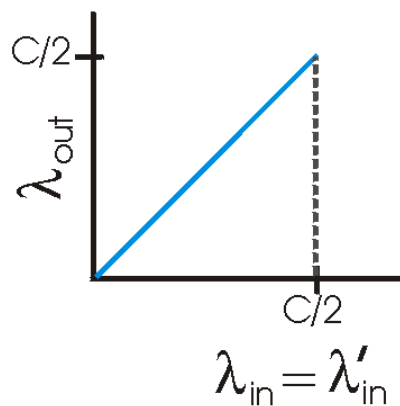
r large delays when congested

r maximum achievable throughput

# Causes/costs of congestion: scenario 2

r  one router, *finite* buffers

r  sender retransmission of lost packet

Host A

$\lambda_{in}$ : original data

$\lambda'_{in}$ : original data, plus retransmitted data

$\lambda_{out}$

Host B

finite shared output link buffers

# Causes/costs of congestion: scenario 2

r   always we want: $\lambda_{in} = \lambda_{out}$ (goodput)

r   Second step ...retransmission only when loss: $\lambda'_{in} > \lambda_{out}$

r   retransmission of delayed (not lost) packet makes $\lambda'_{in}$ larger (than second case) for same $\lambda_{out}$



Caso in cui ciascun pacchetto instradato
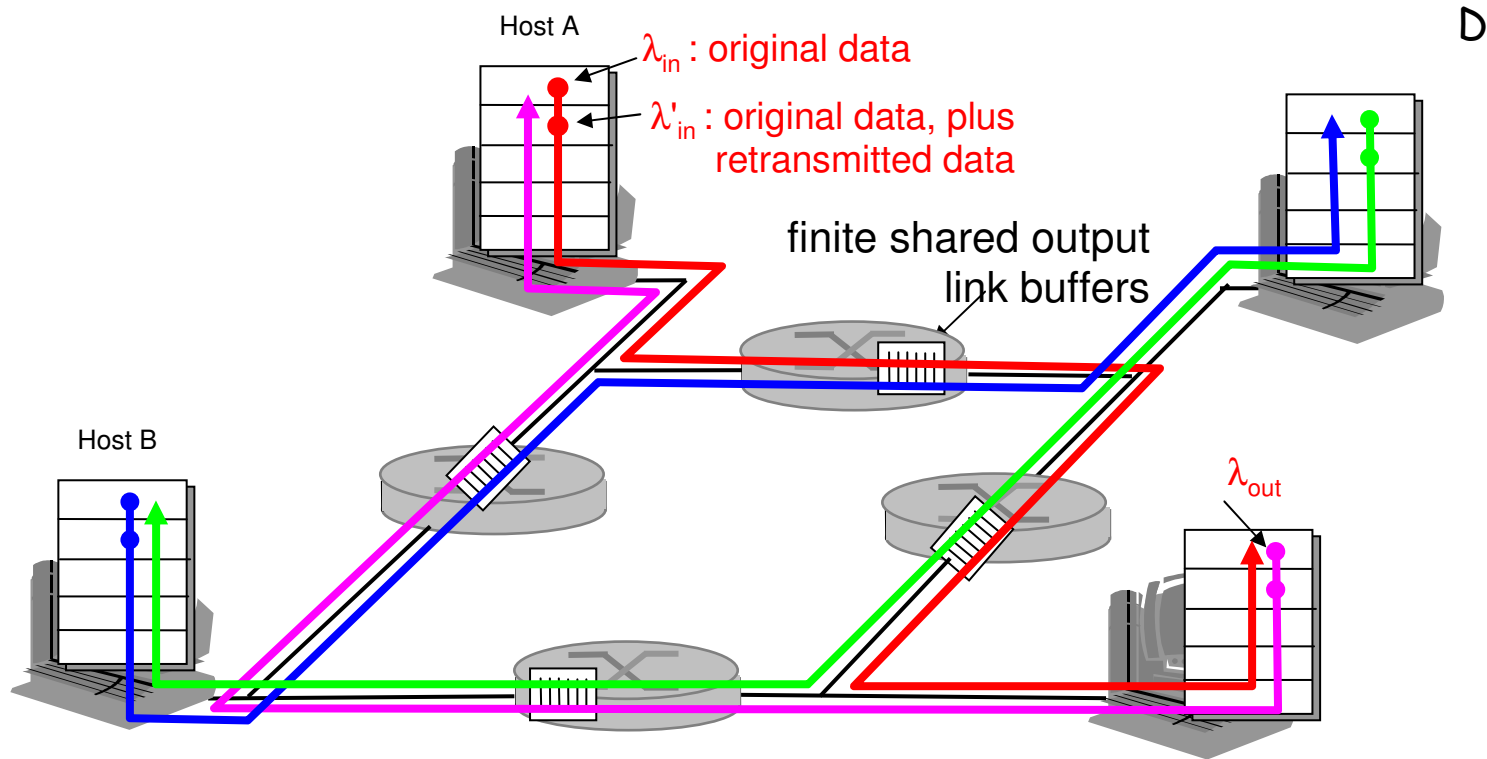Sia trasmesso mediamente due volte dal router

**"costs" of congestion:**

r   more work (retrans) for given "goodput"

r   unneeded retransmissions: link carries multiple copies of pkt
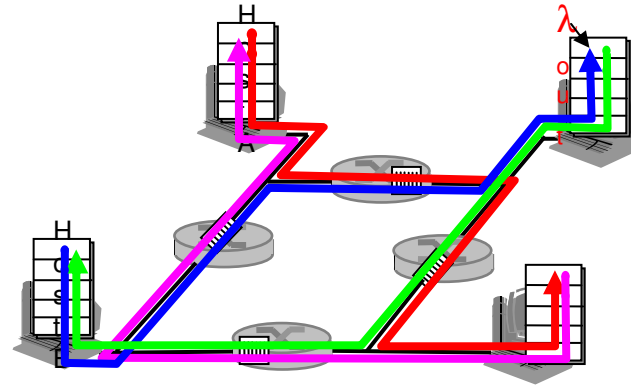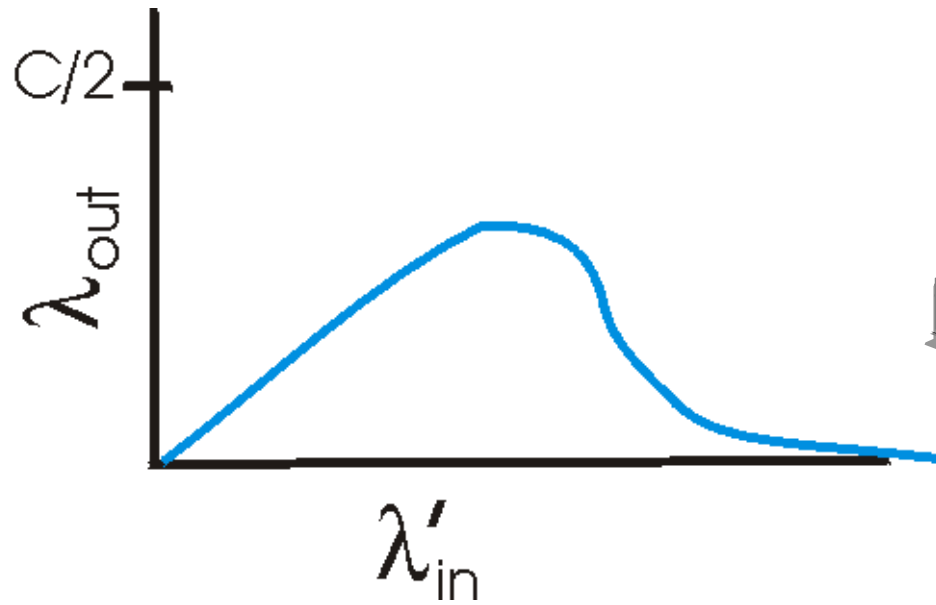
2

# Causes/costs of congestion: scenario 3

r   four senders

r   multihop paths

r   timeout/retransmit

$Q$: what happens as $\lambda_{in}$ and $\lambda'_{in}$ increase ?



Host A

$\lambda_{in}$ : original data

$\lambda'_{in}$ : original data, plus retransmitted data

D

finite shared output link buffers

Host B

$\lambda_{out}$

D-B traffic high

# Causes/costs of congestion: scenario 3



**Another "cost" of congestion:**

r   when packet dropped, any "upstream transmission capacity used for that packet was wasted!

# Approaches towards congestion control

Two broad approaches towards congestion control:

**End-end congestion control:**

- r no explicit feedback from network
- r congestion inferred from end-system observed loss, delay
- r approach taken by TCP

**Network-assisted congestion control:**

- r routers provide feedback to end systems
  - m single bit indicating congestion (SNA, DECbit, TCP/IP ECN, ATM)
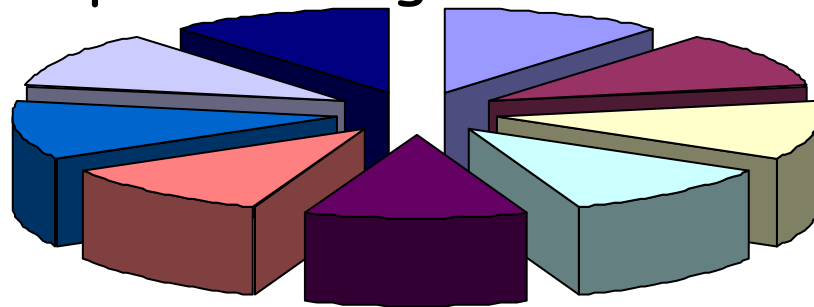  - m explicit rate sender should send at

# Chapter 3 outline

# TCP: controllo di congestione

r Il TCP ha dei meccanismi di controllo della congestione

   m il flusso dei dati in ingresso in rete è anche regolato dalla situazione di traffico in rete

   m se il traffico in rete porta a situazioni di congestione il TCP riduce velocemente il traffico in ingresso

   m in rete non vi è nessun meccanismo per notificare esplicitamente le situazioni di congestione

   m il TCP cerca di scoprire i problemi di congestione sulla base degli eventi di perdita dei pacchetti

# TCP: controllo di congestione

r il meccanismo si basa ancora sulla sliding window la cui larghezza viene dinamicamente regolata in base alle condizioni in rete

r in linea di principio scopo del controllo è far si che il flusso emesso da ciascuna sorgente venga regolato in modo tale che il flusso complessivo offerto a ciascun canale non superi la sua capacità

r tutti i flussi possono essere ridotti in modo tale che la capacità della rete venga condivisa da tutti in misura se possibile uguale

# The problem of congestion

SENDERs
(bulk flows)

RECEIVERs
(large capacity)

Advertise large win

Several outstanding segments

Internal
network
congestion:
- queues build up
- delay increases
- RTOs expire
-more segments transmitted, more
Segments retransmitted -> more congestion!

# The goal of congestion control

**SENDERs**
**(bulk flows)**

**RECEIVERs**
**(large capacity)**

**Bottleneck link rate C**

**N=4 TCP connections**
**Each should transmit at C/4 rate.**

**Since:**

$$thr \approx \frac{W \cdot MSS}{RTT}$$

**Each should adapt W accordingly...**
**How sources can be lead to know the RIGHT value of W??**

# TCP approach for detecting and controlling congestion

r IP protocol does not implement mechanisms to detect congestion in IP routers
- Unlike other networks, e.g. ATM

r necessary indirect means (TCP is an end-to-end protocol)

r TCP approach: congestion detected by lack of acks
- couldn't work efficiently in the 60s & 70s (error prone transmission lines)
- OK in the 80s & 90s (reliable transmission)
- what about wireless networks???

r Controlling congestion: use a SECOND window (congestion window)
- Locally computed at sender
- Outstanding segments: min(receiver_window, congestion_window)

# TCP Congestion Control

r   end-end control (no network assistance)

r   sender limits transmission:

**LastByteSent-LastByteAcked**

$\leq$ **CongWin**

r   Roughly,

$$\text{rate} = \frac{\text{CongWin}}{\text{RTT}} \text{ Bytes/sec}$$

r   **CongWin** is dynamic, function of perceived network congestion

**How does sender perceive congestion?**

r   loss event = timeout *or* 3 duplicate acks

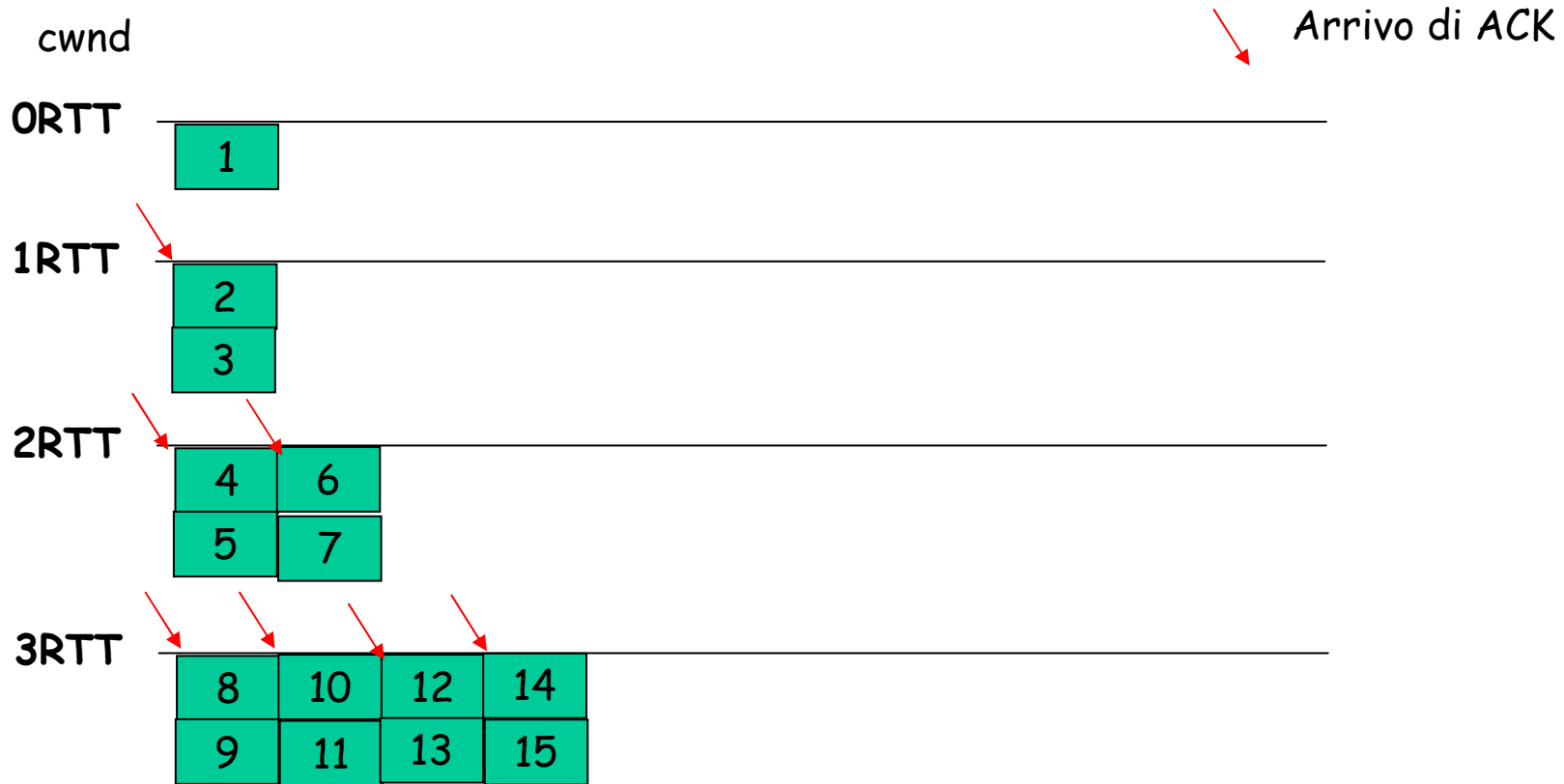r   TCP sender reduces rate (**CongWin**) after loss event

**three mechanisms:**

m   AIMD

m   slow start

m   conservative after timeout events

# Starting a TCP transmission

- r A new offered flow may suddenly overload network nodes

  - m receiver window is used to avoid recv buffer overflow
  - m But it may be a large value (16-64 KB)

- r Idea: slow start

  - m Start with small value of cwnd
  - m And increase it as soon as packets get through

    - Arrival of ACKs = no packet losts = no congestion

- r Initial cwnd size:

  - m Just 1 MSS!
  - m Recent (1998) proposals for more aggressive starts (up to 4 MSS) have been found to be dangerous
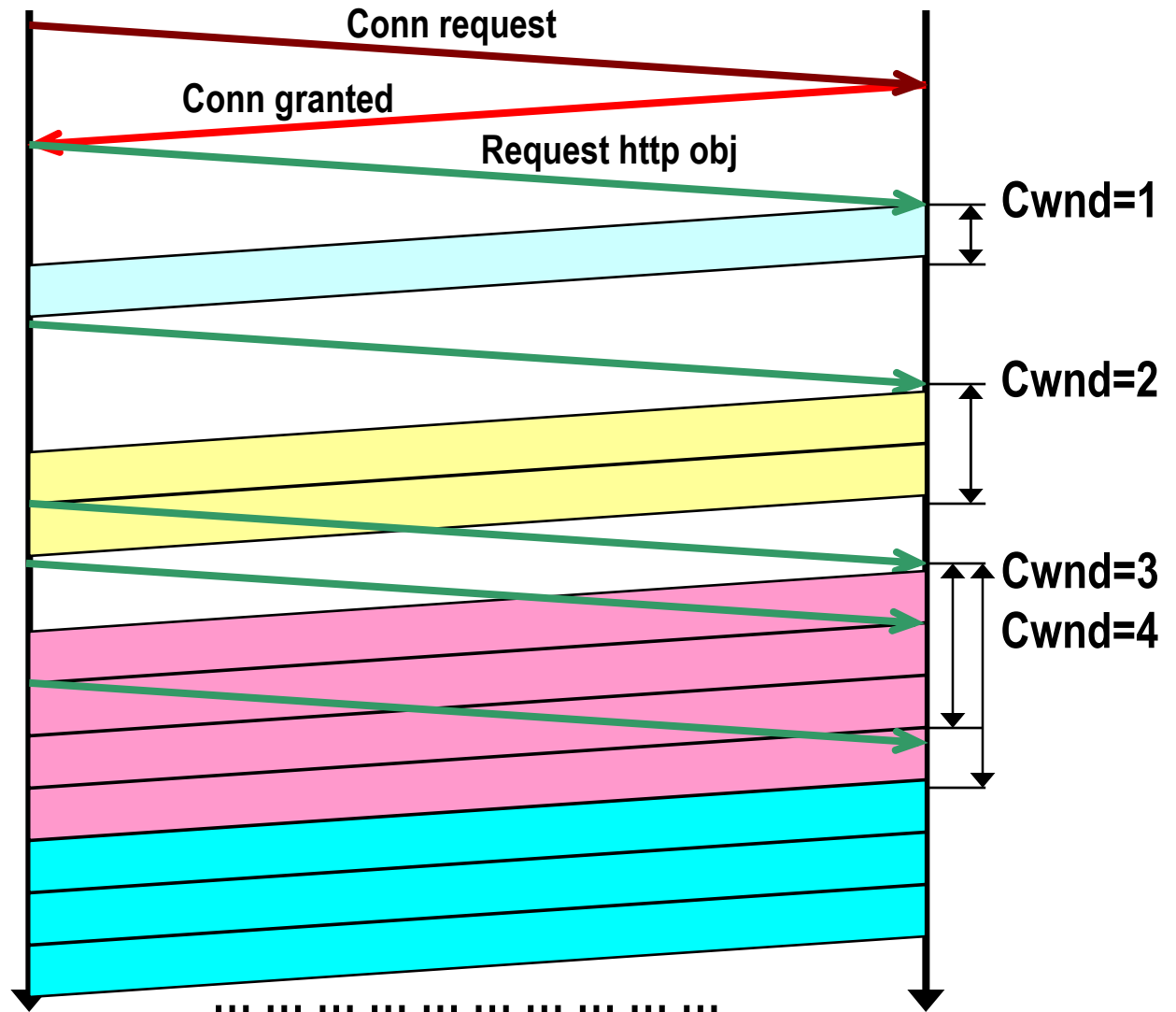
# Slow start: the idea



Si trasmette il minimo tra window e cwd pacchetti

# Slow start – exponential increase

➔ First start: set congestion window cwnd = 1MSS

➔ send cwnd segments
  ⇨ assume cwnd <= receiver win

➔ upon successful reception:
  ⇨ Cwnd +=1 MSS
  ⇨ i.e. double cwnd every RTT
  ⇨ until reaching receiver window advertisement
  ⇨ <u>OR a segment gets lost</u>

Conn request

Conn granted

Request http obj

Cwnd=1

Cwnd=2

Cwnd=3
Cwnd=4

… … … … … … … … … …

# Detecting congestion and restarting

r **Segment gets lost**

   m Detected via RTO expiration

   m Indirectly notifies that one of the network nodes along the path has lost segment

      – Because of full queue

r **Restart from cwnd=1 (slow start)**

r **But introduce a supplementary control: slow start threshold**

      • sstresh = max(min(cwnd,window)/2,2MSS)

   m The idea is that we now KNOW that there is congestion in the network, and we need to increase our rate in a more careful manner…

   m Ssthresh defines the "congestion avoidance" region

# Congestion avoidance

r **If cwnd < ssthresh**

   m Slow start region: Increase rate exponentially

r **If cwnd >= ssthresh**

   m Congestion avoidance region : Increase rate linearly

   m At rate 1 MSS per RTT     *Corrisponde ad un segmento per finestra*

     • Practical implementation:
       cwnd += MSS*MSS/cwnd

     • Good approximation for 1 MSS per RTT

     • Alternative (exact) implementations: count!!

r **Which initial ssthresh?**

      – ssthresh initially set to 65535: unreachable!

*In essence, congestion avoidance is flow control imposed by sender while advertised window is flow control imposed by receiver*

# Simplified example (overall)



Timeout:
cwnd = 1
ssthresh=8

Timeout:
cwnd = 1
ssthresh=6

Congestion window cwnd (in MSS)

16
14
12
10
8
6
4
3
2
1

1

Number of transmissions