

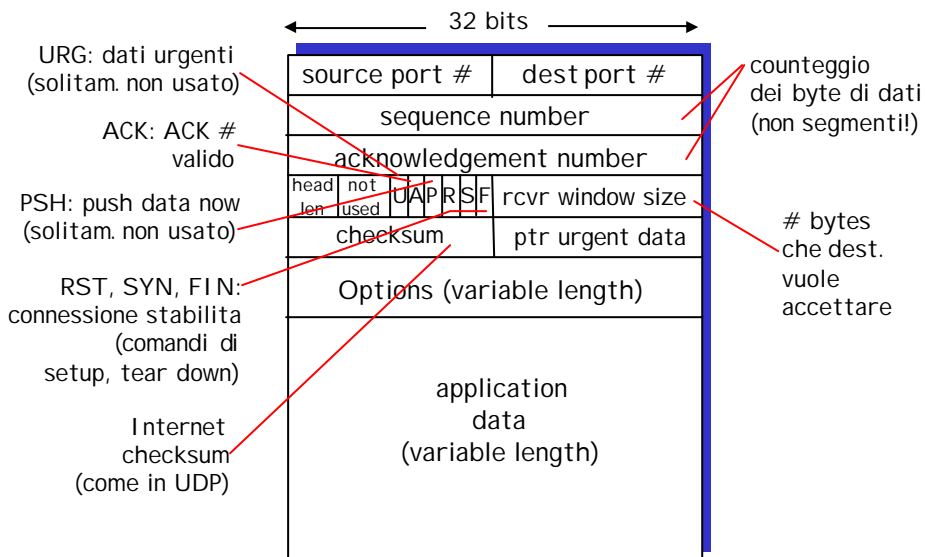
# TCP: Panoramica RFC: 793, 1122, 1323, 2018, 2581

- r **punto-punto:**
  - m un mittente, un destinatario
- r **flusso di byte affidabile e ordinato**
- r **protocollo pipeline:**
  - m il controllo di flusso e di congestione definisce la dimensione della window
- r **buffer send & receive**
- r **dati full duplex :**
  - m Flusso bi-direzionale nella stessa connessione
  - m MSS: maximum segment size
- r **connection-oriented:**
  - m handshaking per inizializzare lo stato del mittente e destinatario
- r **flusso controllato:**
  - m Il mittente non sovraccarica il destinatario



II Livello Trasporto 3b-1

# TCP: struttura del segmento



II Livello Trasporto 3b-2

## TCP: Connessione

Ricordate: nel TCP si stabilisce una "connessione" prima di scambiare segmenti dati

- r inizializzare variabili TCP:
  - m seq. #
  - m info buffers, controllo
  - m flusso (es., **RcvWindow**)

r *client*: avvia connessione  
`Socket clientSocket = new Socket("hostname", "port number");`

r *server*: contattato da client

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

### Three way handshake:

Passo 1: il client invia un SYN al server

- m SYN=1, specifica il seq # iniziale

Passo 2: il server riceve SYN, risponde SYNACK

- m alloca buffers
- m ACK del SYN, specifica server-> seq. # iniziale

Passo 3: client ric. SYNACK

- m alloca buffers
- m invia riscontro (SYN = 0, seq# =iniziale+1, ACK del seq# server+1)

## TCP: Connessione (cont.)

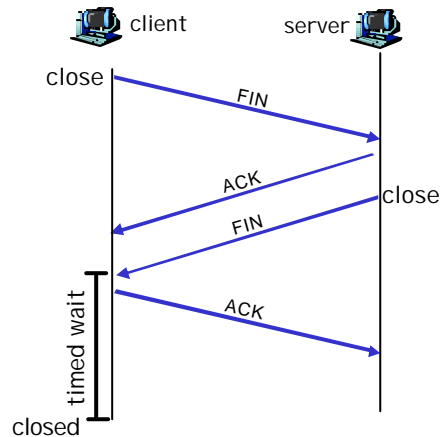
### Chiusura connessione:

client closes socket:

```
clientSocket.close();
```

Passo 1: il *client* invia FIN al server

Passo 2: il *server* riceve FIN, replica con ACK. Chiude la connessione, invia FIN.



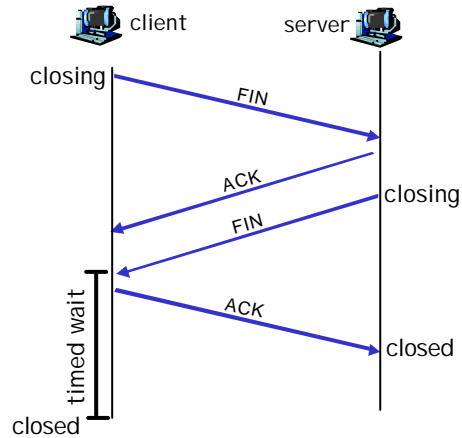
## TCP: Connessione (cont.)

**Passo 3:** il client riceve FIN, replica con ACK.

- entra in "attesa" risponde con ACK ai FIN ricevuti

**Passo 4:** server, riceve ACK, chiude la connessione

**Nota:** con poche modifiche, può gestire FIN simultanei

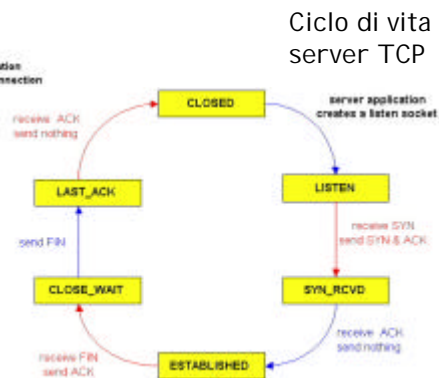


II Livello Trasporto 3b-5

## TCP: Connessione (cont)



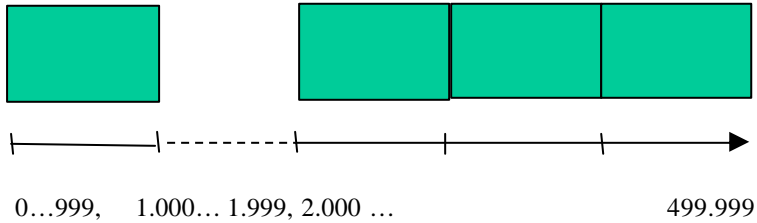
Ciclo di vita client TCP



Ciclo di vita server TCP

II Livello Trasporto 3b-6

## TCP seq. # e ACK



TCP: fornisce riscontri cumulativi

TCP: non ci sono regole fisse nello standard per i segmenti fuori ordine.

## TCP seq. # e ACK

### Seq. #:

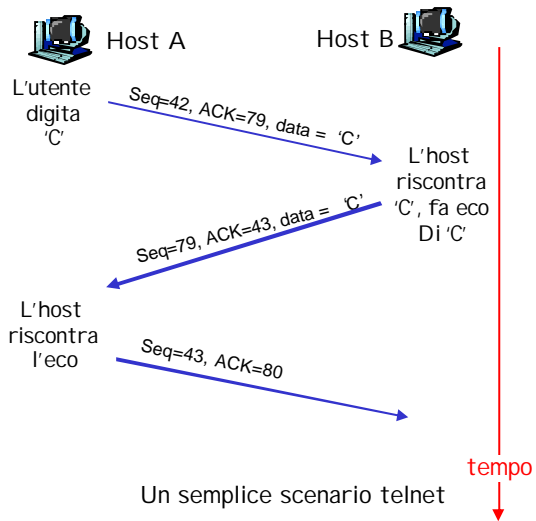
- m numero del primo byte nei dati del segmento (primo # scelto a caso)

### ACK:

- m seq # del prossimo byte atteso
- m ACK cumulativi

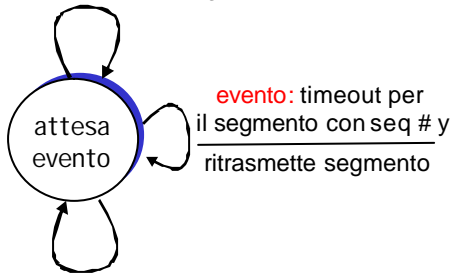
**D:** come il destinatario tratta i segmenti fuori ordine ?

- m R: le specifiche del TCP non lo dicono, dipende dalla implementazione



## TCP: trasferimento affidabile

**evento:** ricezione dati  
dall'applicazione  
crea e invia un segmento



**evento:** ACK ricevuto  
con ACK # y  
Esamina ACK

mittente semplice, assumendo

- trasferimento one way
- nessun controllo di flusso o congestione

II Livello Trasporto 3b-9

## TCP: trasferim affidabile

mittente  
TCP  
semplificato

```
00 sendbase = initial_sequence number
01 nextseqnum = initial_sequence number
02
03 loop (forever) {
04   switch(event)
05     event: data received from application above
06       create TCP segment with sequence number nextseqnum
07       start timer for segment nextseqnum
08       pass segment to IP
09       nextseqnum = nextseqnum + length(data)
10     event: timer timeout for segment with sequence number y
11       retransmit segment with sequence number y
12       compute new timeout interval for segment y
13       restart timer for sequence number y
14     event: ACK received, with ACK field value of y
15       if (y > sendbase) { /* cumulative ACK of all data up to y */
16         cancel all timers for segments with sequence numbers < y
17         sendbase = y
18       }
19     else { /* a duplicate ACK for already ACKed segment */
20       increment number of duplicate ACKs received for y
21       if (number of duplicate ACKs received for y == 3) {
22         /* TCP fast retransmit */
23         resend segment with sequence number y
24         restart timer for segment y
25       }
26   } /* end of loop forever */
```

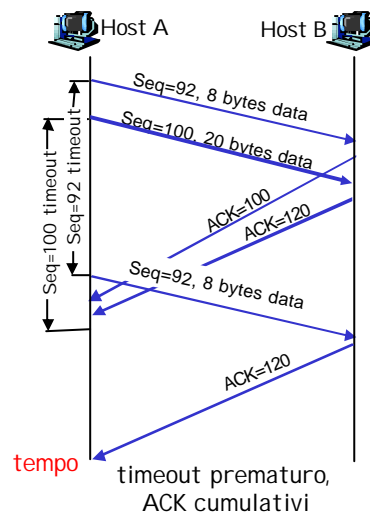
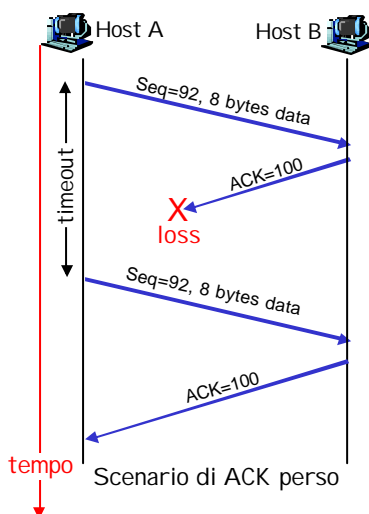
II Livello Trasporto 3b-10

## TCP: generaz. ACK [RFC 1122, RFC 2581]

Evento	azioni del destinatario TCP
arrivo ordinato segmenti, non ci sono "buchi", tutto il resto già riscontrato	ACK ritardato. Attendi 500ms per prossimo segmento. Se non arriva, Invia ACK
arrivo ordinato segmenti, non ci sono "buchi", un ACK ritardato sospeso	Invia immediatamente un singolo ACK cumulativo
arrivo segmento fuori ordine seq. # maggiore di quello atteso rilevato un "buco" (gap)	invia ACK duplicato che indica il seq. # del prossimo byte atteso
arrivo di un segmento che colma parzialmente o completamente il gap	ACK immediato se il segmento inizia all'estremo inferiore del gap

II Livello Trasporto 3b-11

## TCP: ritrasmissione



II Livello Trasporto 3b-12

## TCP: Controllo Flusso

### flow control

il mittente non sovraccarica il buffer del destinatario trasmettendo troppo, troppo in fretta

**Destinatario:** informa esplicitamente il mittente dello spazio disponibile (cambia dinamicamente)

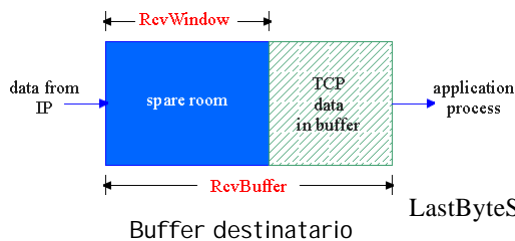
m **campo RcvWindow** nel segmento TCP

**RcvBuffer** = dimensione TCP Receive Buffer

**RcvWindow** = quantità di spazio rimasto nel Buffer

$\text{LastByteRcvd} - \text{LastByteRead} \leq \text{RcvBuffer}$

$\text{RcvWindow} = \text{RcvBuffer} - [\text{LastByteRcvd} - \text{LastByteRead}]$



**Mittente:** fa sì che i dati trasmessi ma non riscontrati < del più recente **RcvWindow** ricevuto

$\text{LastByteSent} - \text{LastByteAcked} \leq \text{RcvWindow}$

II Livello Trasporto 3b-13

## TCP Round Trip Time e Timeout

**D:** come si stabilisce il valore del timeout?

- r Maggiore di RTT
  - m nota: RTT varia
- r Troppo breve: timeout prematuro
  - m Ritrasmissioni inutili
- r Troppo lungo: reazione lenta allo smarrimento di segmenti

**D:** Come stimare RTT?

- r **SampleRTT**: misura del tempo dalla trasmissione del segmento fino alla ricezione dell'ACK
  - m ignora ritrasmissioni, ACK cumulativi
- r **SampleRTT** varia, occorre stimare RTT in maniera opportuna
  - m media di molte misure recenti e non solo del **SampleRTT** corrente

II Livello Trasporto 3b-14

## TCP Round Trip Time e Timeout

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

- r media variabile pesata esponenzialmente
- r influenza di un dato campione diminuisce esponenzialm.
- r Valore tipico di x: 0.1

### Impostazione del timeout

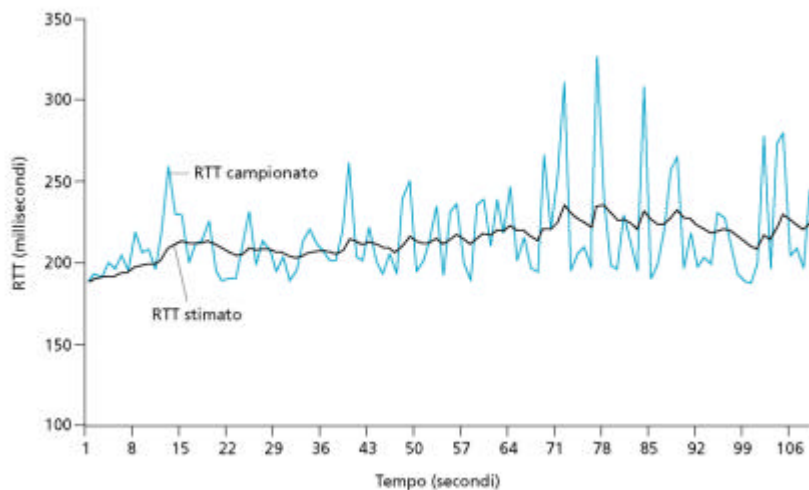
- r **EstimatedRTT** più "margine di sicurezza"
- r Ampie variazioni **EstimatedRTT** -> margine di sicurezza maggiore

$$\text{Timeout} = \text{EstimatedRTT} + 4 * \text{Deviation}$$

$$\text{Deviation} = (1-x) * \text{Deviation} + x * |\text{SampleRTT} - \text{EstimatedRTT}|$$

II Livello Trasporto 3b-15

## RTT campionati e RTT stimati



II Livello Trasporto 3b-16



## Principi di Controllo Congestione

### Congestione:

- r informalmente: "troppe sorgenti che inviano troppi dati troppo in fretta perché la *rete* sia in grado di gestirli"
- r È diverso dal controllo di flusso!
- r effetti:
  - m Pacchetti persi (a causa dell' overflow del buffer ai router)
  - m Ritardi lunghi (a causa delle code nei buffer dei router)
- r Un problema nella top-10!

II Livello Trasporto 3b-17

## Approcci per il controllo congestione

Due approcci principali:

### Controllo

#### end-to-end :

- r Non c'è feedback esplicito dalla rete
- r stato congestione ricavato dai livelli di perdita e ritardo osservati agli end-system
- r L'approccio del TCP

### Controllo

#### network-assisted:

- r I router forniscono feedback agli end system
  - m Un bit indica la congestione (SNA, DECbit, TCP/IP ECN, ATM)
  - m Viene specificato esplicitamente a quale velocità il mittente deve trasmettere

II Livello Trasporto 3b-18

# Controllo Congestione del TCP

- r controllo end-to-end (niente feedback da network)
- r La velocità trasmissiva è limitata dalla dimensione della finestra di congestione, **Congwin**, sui segmenti:



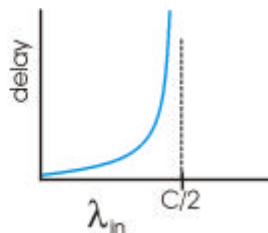
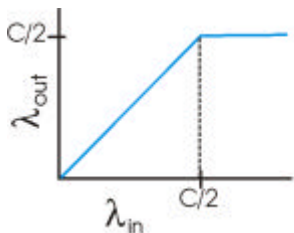
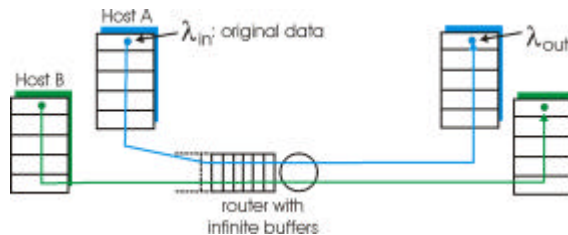
- r  $w$  segmenti, ciascuno invia MSS bytes in un RTT:

$$\text{throughput} = \frac{w * \text{MSS}}{\text{RTT}} \text{ Bytes/sec}$$

II Livello Trasporto 3b-19

## Cause/costi congestione: scenario 1

- r Due mittenti, due destinatari
- r un router, buffer infinito
- r No ritrasmissioni

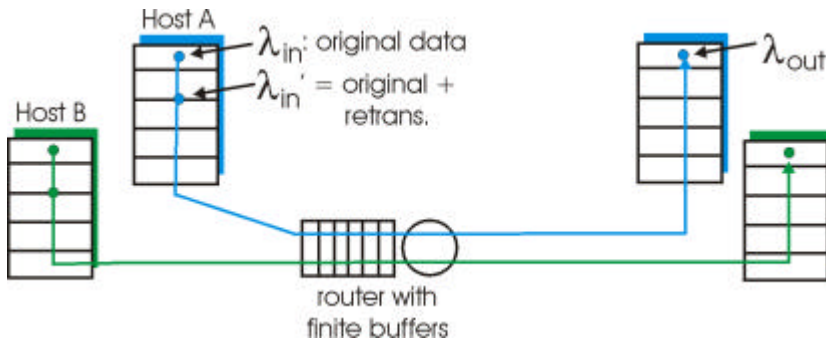


- r grossi ritardi in caso congestione
- r massimo throughput ottenibile

II Livello Trasporto 3b-20

## Cause/costi congestione: scenario 2

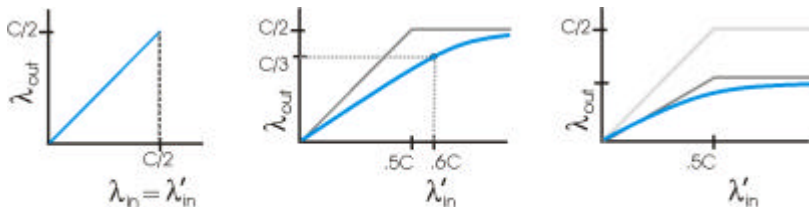
- r un router, buffer *finito*
- r Ritrasmissione dei pacchetti persi



II Livello Trasporto 3b-21

## Cause/costi congestione: scenario 2

- r sempre:  $\lambda_{in} = \lambda_{out}$  (goodput)
- r ritrasmissione "perfetta" solo quando:  $\lambda_{in}' > \lambda_{out}$
- r ritrasmissione dei pacchetti ritardati (non persi) rende  $\lambda_{in}'$  più lungo (del caso perfetto) per alcuni  $\lambda_{out}$



"costi" della congestione:

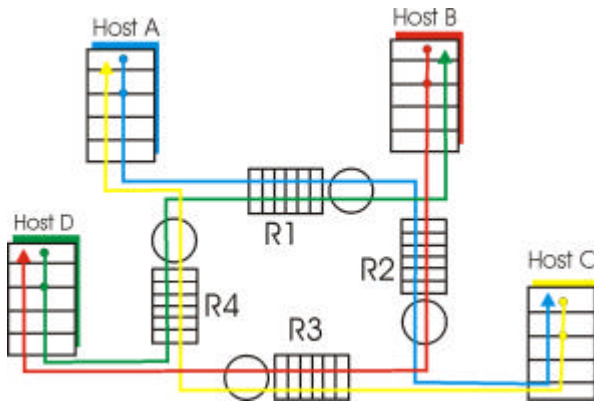
- r più lavoro (ritrasm) per un dato "goodput"
- r ritrasmissioni non necessarie: più copie del pkt sul link

II Livello Trasporto 3b-22

## Cause/costi congestione: scenario 3

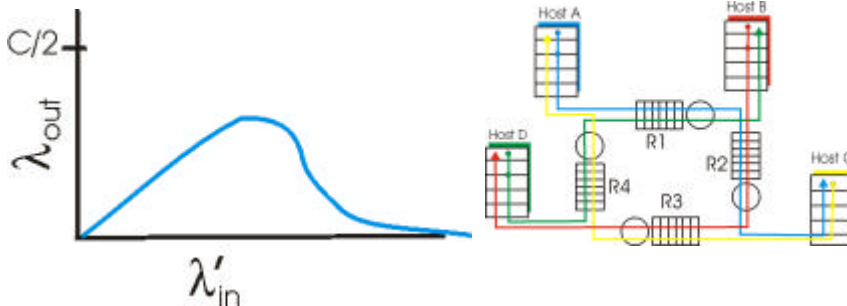
- r quattro mittenti
- r Percorsi multihop
- r timeout/retransmit

D: cosa succede se  $\lambda_{in}$  e  $\lambda'_{in}$  aumentano ?



II Livello Trasporto 3b-23

## Cause/costi congestione: scenario 3



Un altro "costo" della congestione:

- r Quando un pacchetto viene scartato, ogni capacità di trasmissione "upstream" usata per il pacchetto viene sprecata!

II Livello Trasporto 3b-24

## Studio di un Caso: controllo congestione dell' ATM (ABR)

### ABR: available bit rate:

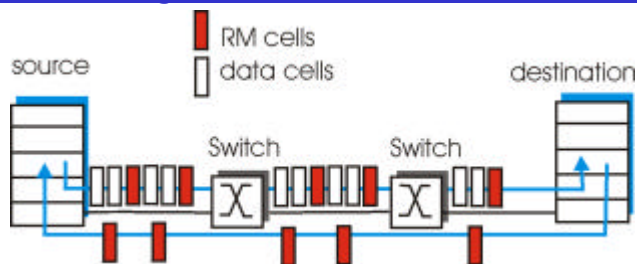
- r "servizio elastico"
- r se il percorso del mittente è "scarico":
  - m Il mittente deve usare la banda disponibile
- r se il percorso mittente è "congestionato":
  - m Il mittente viene riportata ad una velocità minima garantita

### celle di RM (resource management) :

- r Il mittente le "mischia" con le celle dati
- r bits nella cella RM scritti dagli switch ("*network-assisted*")
  - m NI bit: non aumentare la velocità (mild congestion)
  - m CI bit: congestion indication
- r Le celle RM vengono restituite dal destinatario con i bit inalterati

II Livello Trasporto 3b-25

## Studio di un Caso: controllo congestione dell' ATM (ABR)



- r campo di due-byte ER (explicit rate) nella cella RM
  - m uno switch in congestione può diminuire il valore di ER
- r nelle celle dati c'è il bit EFCI : posto a 1 negli switch congestionati
  - m se la cella dati che precede la cella RM ha EFCI , il mittente pone il bit CI a 1 nella cella RM di ritorno

II Livello Trasporto 3b-26

# Controllo Congestione del TCP:

- r "sondando" la disponibilità di banda:
  - m **idealmente**: trasmetti il più velocemente possibile (**Congwin** il più ampia possibile) senza perdite
  - m **incrementa Congwin** finché iniziano le perdite (congestione)
  - m congestione: **decrementa Congwin**, poi inizia a sondare di nuovo
- r due "fasi"
  - m **slow start**
  - m **congestion avoidance**
- r variabili importanti:
  - m **Congwin**
  - m **threshold**: definisce la soglia tra le due fasi di slow start e congestion control

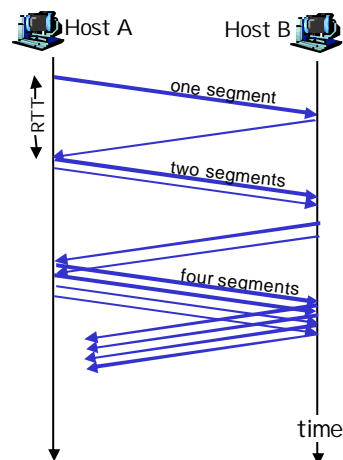
II Livello Trasporto 3b-27

# TCP Slowstart

**algoritmo Slowstart**

```
initialize: Congwin = 1
for (each segment ACKed)
  Congwin++
until (loss event OR
      CongWin > threshold)
```

- r Incremento esponenziale (per RTT) nel window size (non è poi così lento!)
- r evento loss : timeout (Tahoe TCP) and/or tre ACK duplicati (Reno TCP)



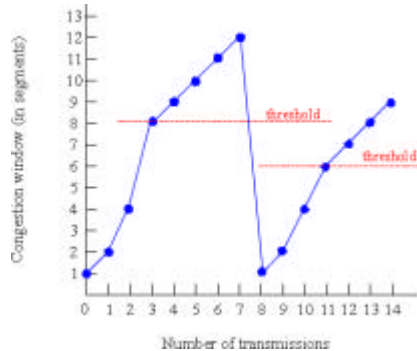
II Livello Trasporto 3b-28

# TCP Congestion Avoidance

## Congestion avoidance

```

/* slowstart is over */
/* Congwin > threshold */
Until (loss event) {
  every w segments ACKed:
    Congwin++
}
threshold = Congwin/2
Congwin = 1
perform slowstart1
    
```



1: TCP Reno salta lo slowstart (fast recovery) dopo tre ACK duplicati

II Livello Trasporto 3b-29

## AIMD

TCP congestion avoidance:

r **AIMD**: *additive increase, multiplicative decrease*

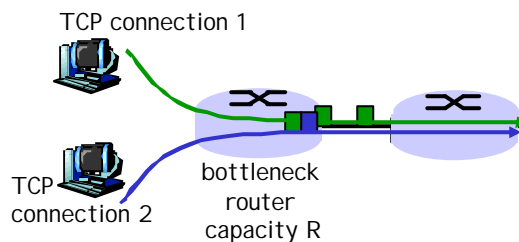
m Aumenta la window di 1 per RTT

m Decrementa la window di un fattore 2 se c'è congestione

## TCP: Equità (Fairness)

obiettivi della Fairness:

se N sessioni TCP condividono lo stesso link, ciascuna deve ottenere 1/N della capacità del link

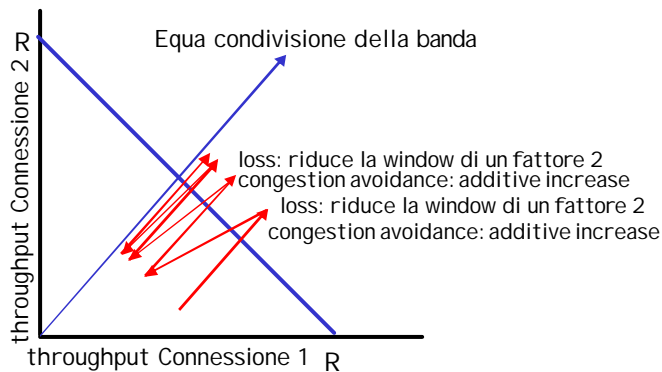


II Livello Trasporto 3b-30

## Perchè il TCP è equo?

Due sessioni in competizione:

- r L'**additive increase** da pendenza 1, se il throughput cresce
- r Il **multiplicative decrease** riduce il throughput proporzionalmente



II Livello Trasporto 3b-31

## TCP: modello di latenza

**D:** Quanto tempo occorre per ricevere un oggetto da un Web server dopo aver inviato una richiesta?

- r Stabilire una connessione TCP
- r Ritardo trasferimento dati

**Notazioni, assunzioni:**

- r Un solo link tra client e server di velocità  $R$
- r Window di congestione fissa,  $W$  segmenti
- r  $S$ : MSS (bits)
- r  $O$ : object size (bits)
- r no ritrasmissioni (no loss, no corruption)

**due casi da considerare:**

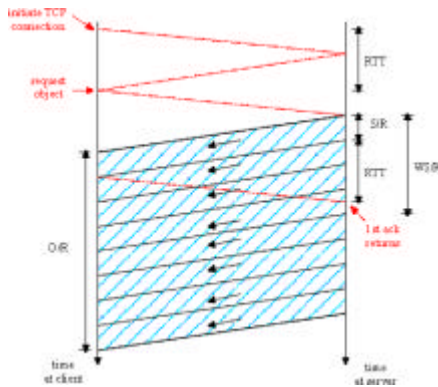
- r  $WS/R > RTT + S/R$ : ACK del primo segmento nella window torna prima
- r  $WS/R < RTT + S/R$ : aspetta ACK dopo aver spedito il valore della finestra dati inviata

II Livello Trasporto 3b-32

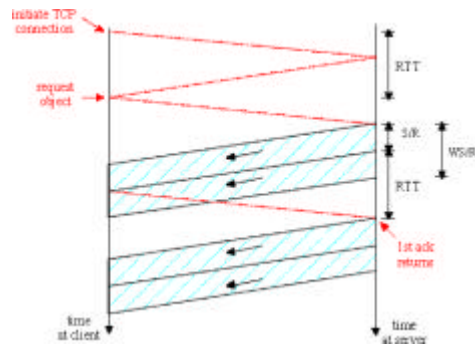


# TCP Latency Modeling

$$K := O/WS$$



Case 1: latency =  $2RTT + O/R$



Case 2: latency =  $2RTT + O/R + (K-1)[S/R + RTT - WS/R]$

## TCP Latency Modeling: Slow Start

- r Now suppose window grows according to slow start.
- r Will show that the latency of one object of size  $O$  is:

$$Latency = 2RTT + \frac{O}{R} + P \left[ RTT + \frac{S}{R} \right] - (2^P - 1) \frac{S}{R}$$

where  $P$  is the number of times TCP stalls at server:

$$P = \min\{Q, K-1\}$$

- where  $Q$  is the number of times the server would stall if the object were of infinite size.
- and  $K$  is the number of windows that cover the object.

## TCP Latency Modeling: Slow Start (cont.)

Example:

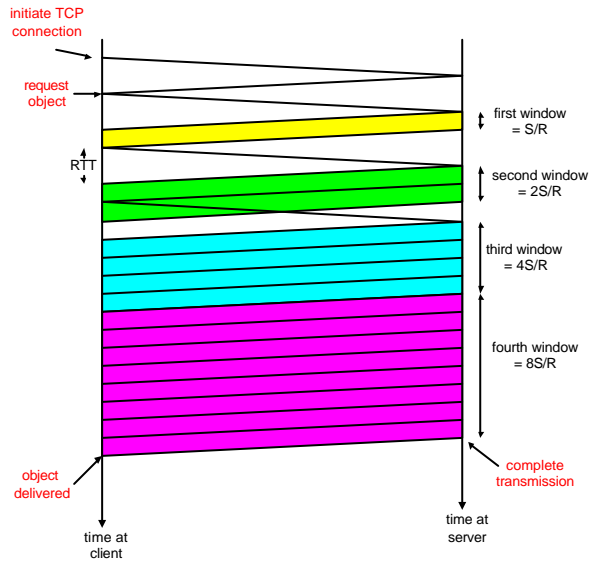
$O/S = 15$  segments

$K = 4$  windows

$Q = 2$

$P = \min\{K-1, Q\} = 2$

Server stalls  $P=2$  times.



II Livello Trasporto 3b-35

## TCP Latency Modeling: Slow Start (cont.)

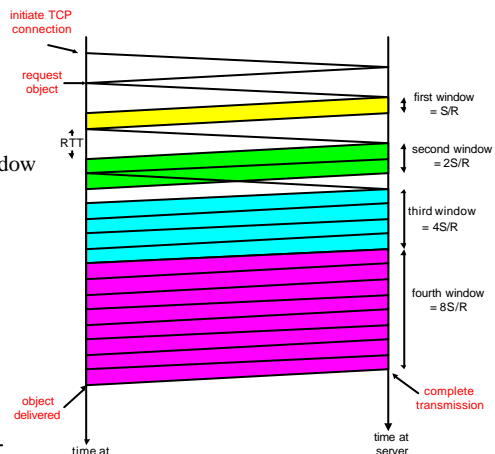
$\frac{S}{R} + RTT =$  time from when server starts to send segment

until server receives acknowledgement

$2^{k-1} \frac{S}{R} =$  time to transmit the  $k$ th window

$\left[ \frac{S}{R} + RTT - 2^{k-1} \frac{S}{R} \right]^+ =$  stall time after the  $k$ th window

$$\begin{aligned} \text{latency} &= \frac{O}{R} + 2RTT + \sum_{p=1}^P \text{stallTime}_p \\ &= \frac{O}{R} + 2RTT + \sum_{k=1}^P \left[ \frac{S}{R} + RTT - 2^{k-1} \frac{S}{R} \right] \\ &= \frac{O}{R} + 2RTT + P \left[ RTT + \frac{S}{R} \right] - (2^P - 1) \frac{S}{R} \end{aligned}$$



II Livello Trasporto 3b-36

## Chapter 3: Summary

- r principles behind transport layer services:
  - m multiplexing/demultiplexing
  - m reliable data transfer
  - m flow control
  - m congestion control
- r instantiation and implementation in the Internet
  - m UDP
  - m TCP

### Next:

- r leaving the network “edge” (application transport layer)
- r into the network “core”