# Chapter 3
# Transport Layer

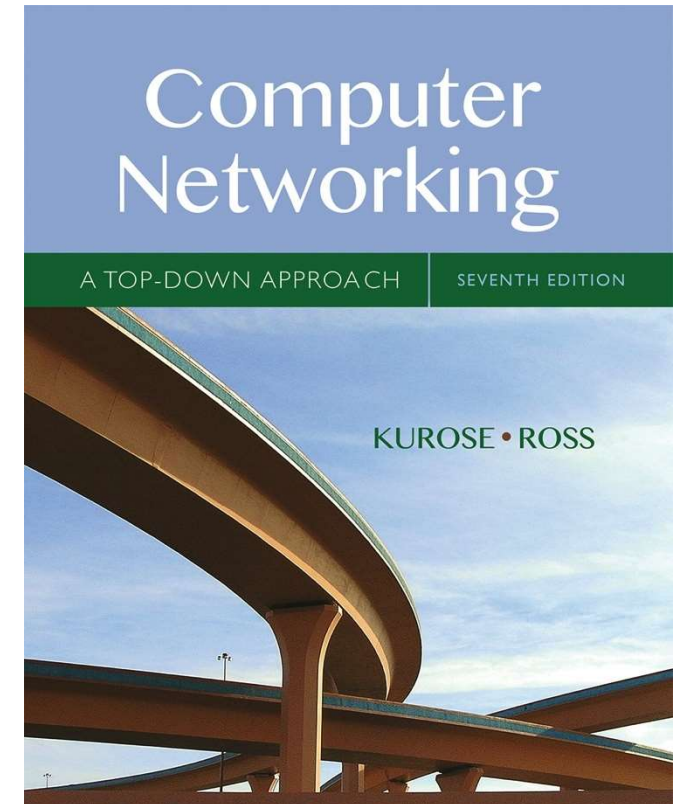## A note on the use of these Powerpoint slides:

We're making these slides freely available to all (faculty, students, readers). They're in PowerPoint form so you see the animations; and can add, modify, and delete slides  (including this one) and slide content to suit your needs. They obviously represent a *lot* of work on our part. In return for use, we only ask the following:

- If you use these slides (e.g., in a class) that you mention their source (after all, we'd like people to use our book!)
- If you post any slides on a www site, that you note that they are adapted from (or perhaps identical to) our slides, and note our copyright of this material.

Thanks and enjoy!  JFK/KWR

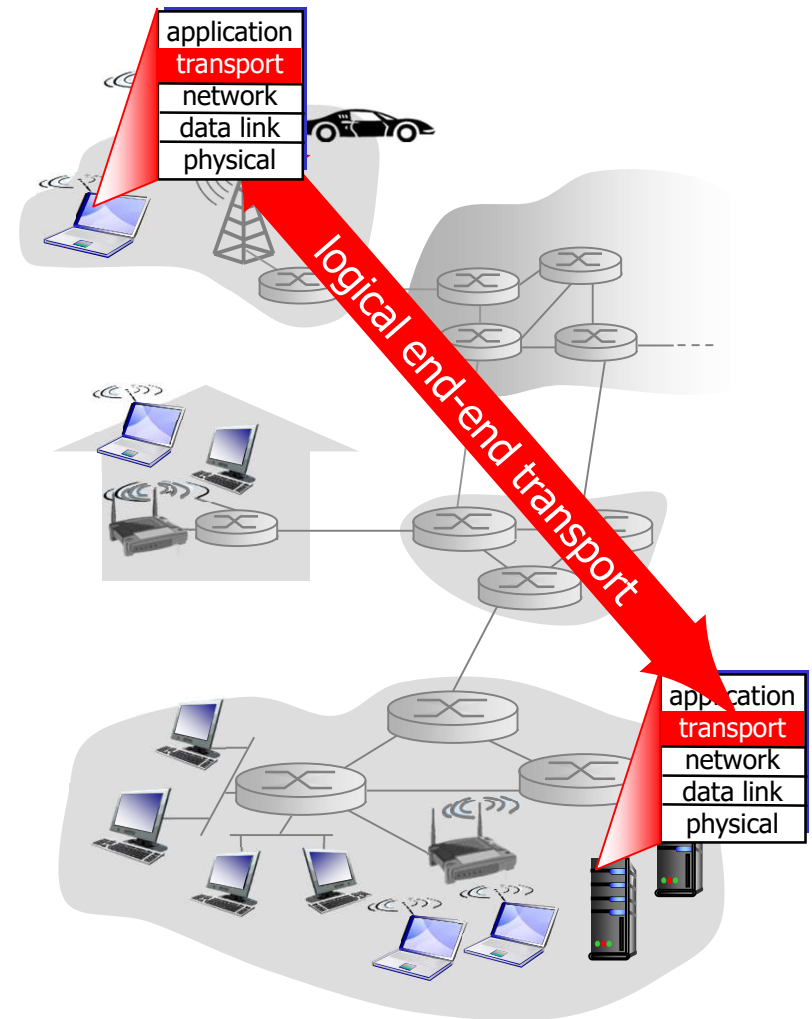*Computer Networking: A Top Down Approach*

7th edition
Jim Kurose, Keith Ross
Pearson/Addison Wesley
April 2016

# Transport services and protocols

- provide *logical communication* between app processes running on different hosts

- transport protocols run in end systems
  - send side: breaks app messages into *segments*, passes to network layer
  - rcv side: reassembles segments into messages, passes to app layer

- more than one transport protocol available to apps
  - Internet: TCP and UDP

application
transport
network
data link
physical

logical end-end transport

application
transport
network
data link
physical

# Transport vs. network layer

- *network layer:* logical communication between hosts

- *transport layer:* logical communication between processes
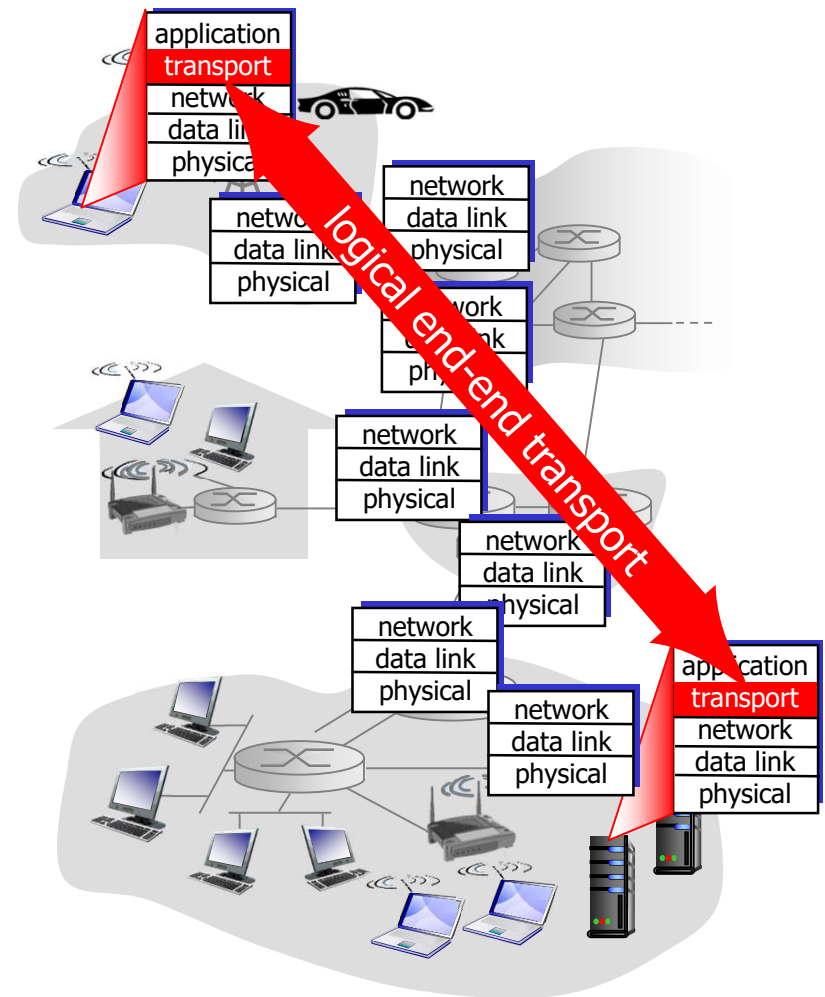  - relies on, enhances, network layer services

*household analogy:*

12 kids in Ann's house sending letters to 12 kids in Bill's house:

- hosts = houses
- processes = kids
- app messages = letters in envelopes
- transport protocol = Ann and Bill who demux to in-house siblings
- network-layer protocol = postal service

# Internet transport-layer protocols

- **reliable, in-order delivery (TCP)**
  - congestion control
  - flow control
  - connection setup
- **unreliable, unordered delivery: UDP**
  - no-frills extension of "best-effort" IP
- **services not available:**
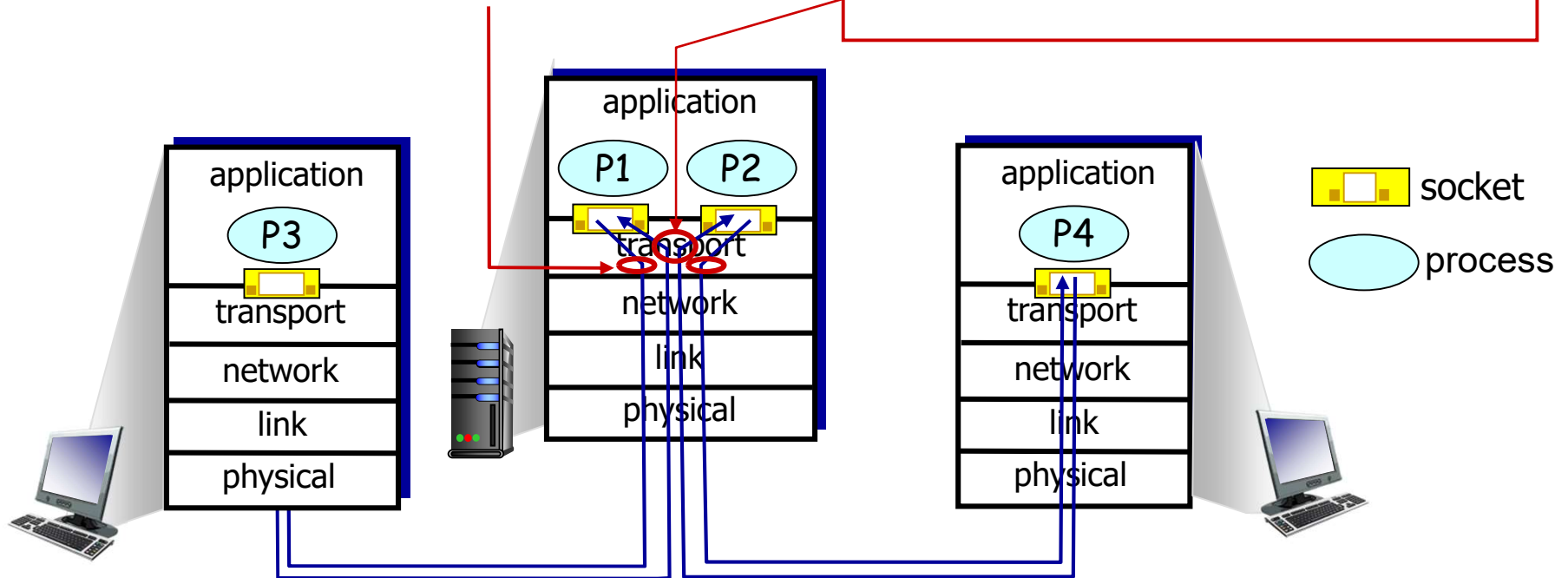  - delay guarantees
  - bandwidth guarantees

# Multiplexing/demultiplexing

_multiplexing at sender:_
handle data from multiple
sockets, add transport header
(later used for demultiplexing)

_demultiplexing at receiver:_
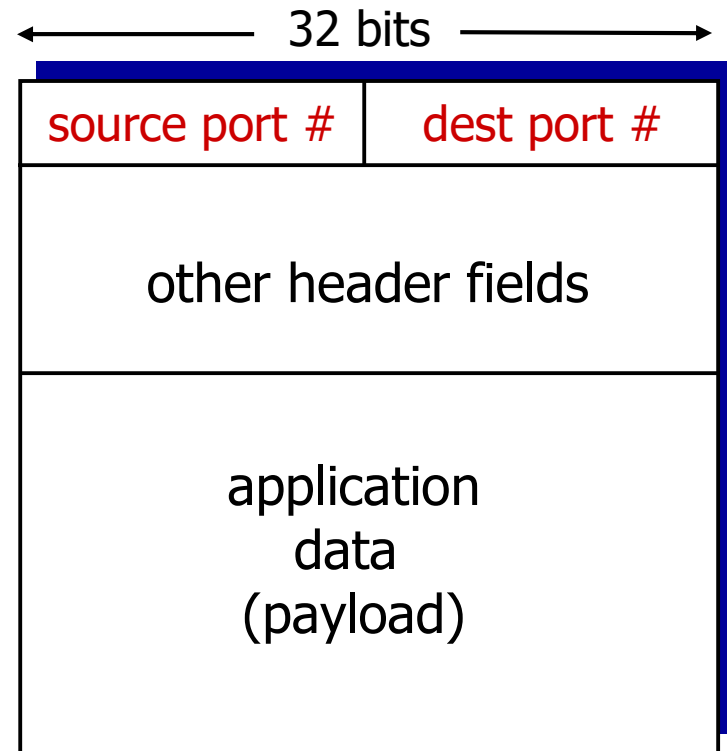use header info to deliver
received segments to correct
socket



application

P1    P2

application
P3

transport

application
P4

transport

network

network

network

link

link

link

physical

physical

physical

socket

process

# How demultiplexing works

- **host receives IP datagrams**
  - each datagram has source IP address, destination IP address
  - each datagram carries one transport-layer segment
  - each segment has source, destination port number
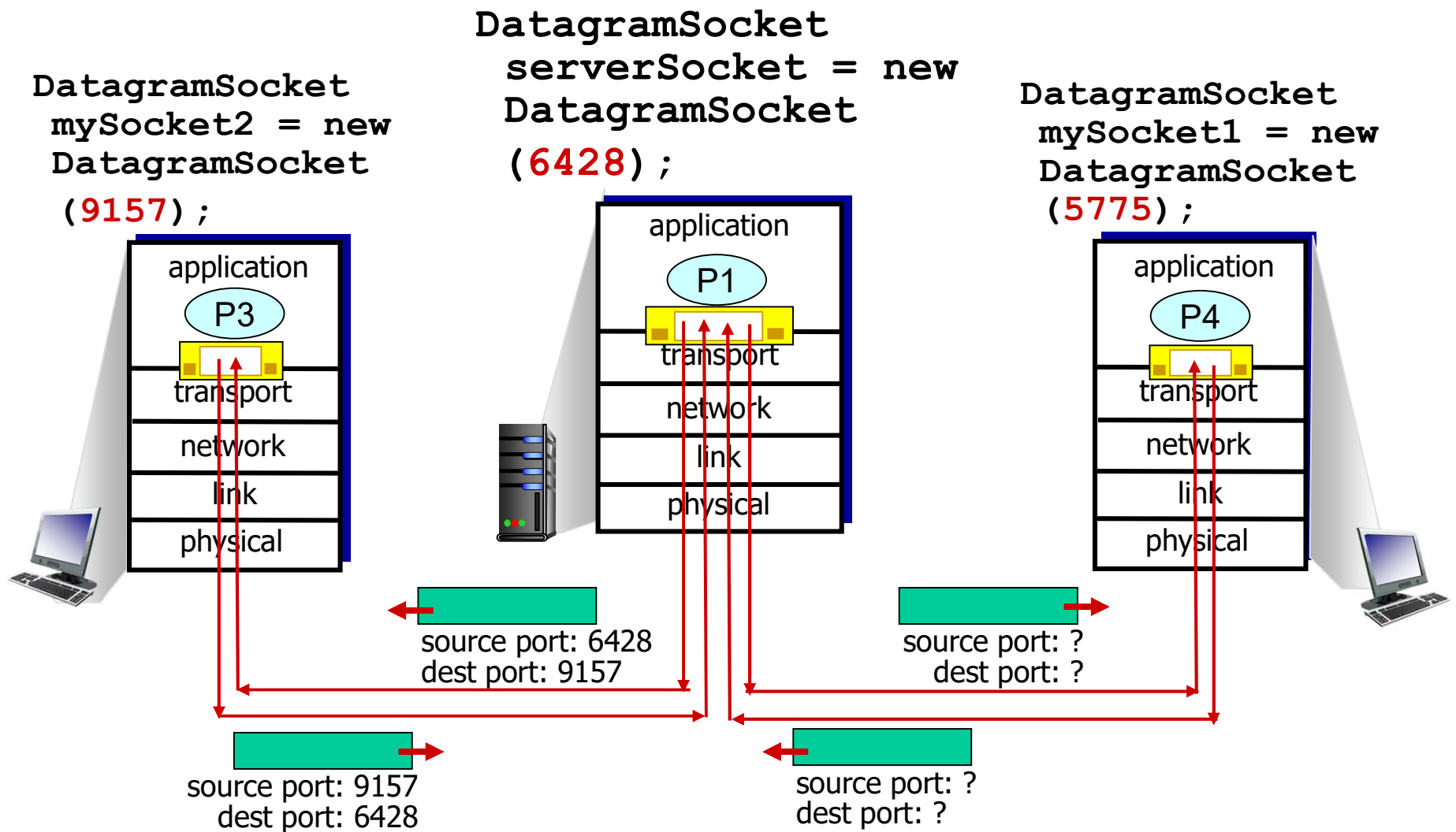- **host uses *IP addresses & port numbers* to direct segment to appropriate socket**

32 bits

| source port # | dest port # |
|---|---|
| other header fields | |

application
data
(payload)

TCP/UDP segment format

# Connectionless demultiplexing

- *recall:* created socket has host-local port #:

  ```
  DatagramSocket mySocket1
  = new DatagramSocket(12534);
  ```

- *recall:* when creating datagram to send into UDP socket, must specify
  - destination IP address
  - destination port #

---

- when host receives UDP segment:
  - checks destination port # in segment
  - directs UDP segment to socket with that port #

➡ IP datagrams with *same dest. port #,* but different source IP addresses and/or source port numbers will be directed to *same socket* at dest
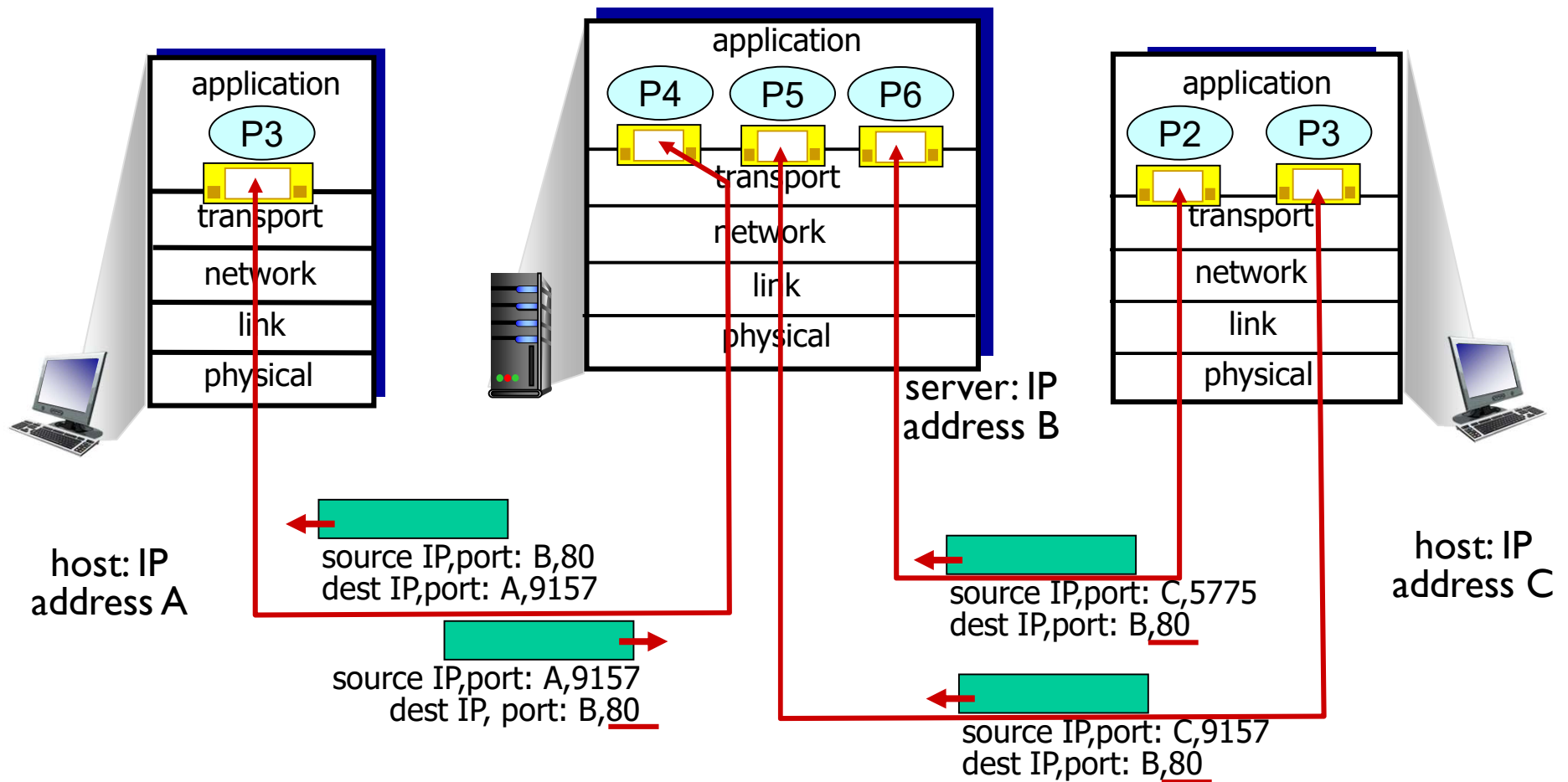
# Connectionless demux: example

DatagramSocket
serverSocket = new
DatagramSocket
(6428);

DatagramSocket
mySocket2 = new
DatagramSocket
(9157);

DatagramSocket
mySocket1 = new
DatagramSocket
(5775);

application

P3

transport

network

link

physical

application

P1

transport

network

link

physical

application

P4

transport

network

link

physical

source port: 6428
dest port: 9157

source port: ?
dest port: ?

source port: 9157
dest port: 6428

source port: ?
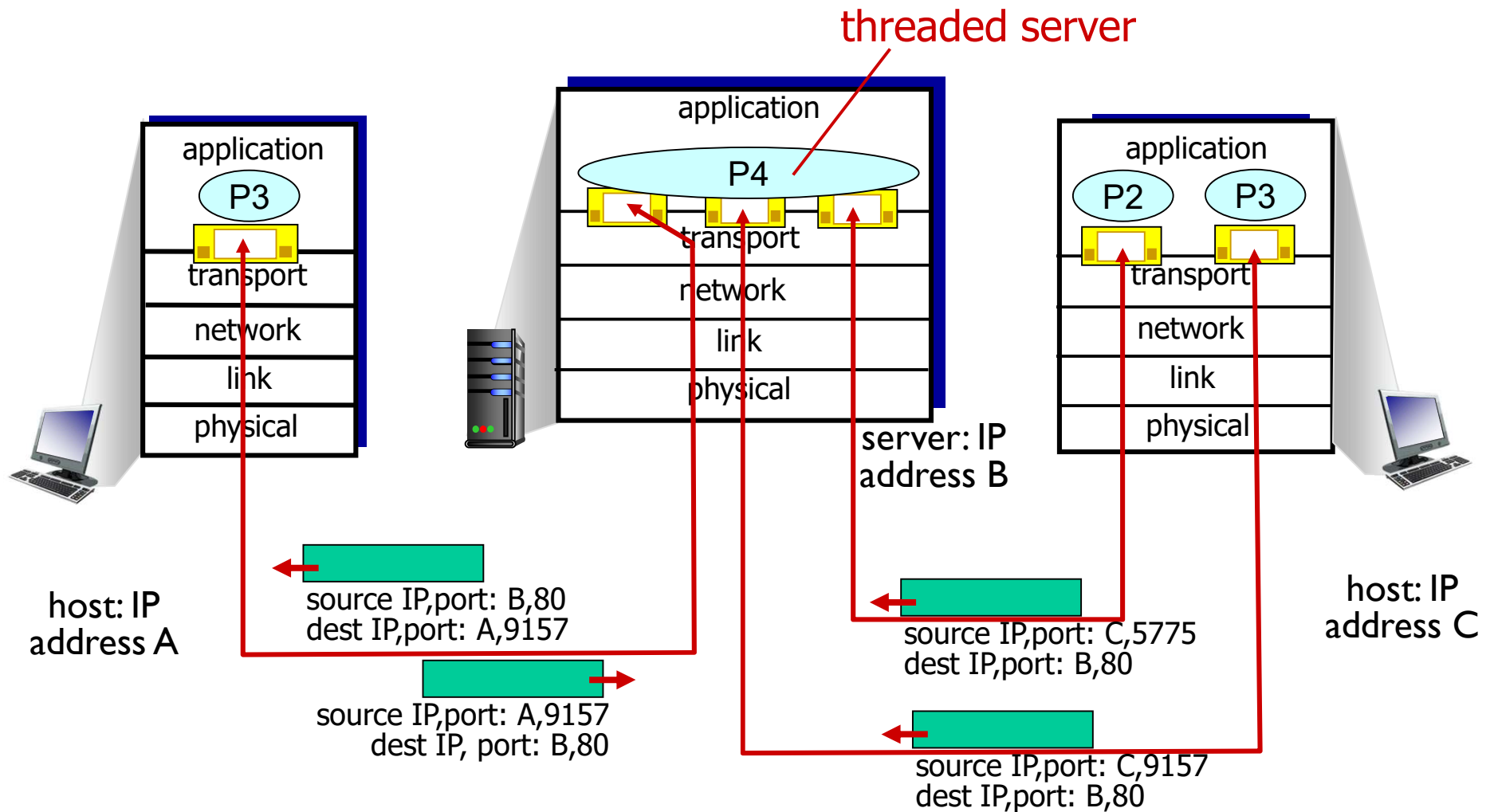dest port: ?

# Connection-oriented demux

- **TCP socket identified by 4-tuple:**
  - source IP address
  - source port number
  - dest IP address
  - dest port number
- **demux: receiver uses all four values to direct segment to appropriate socket**

- **server host may support many simultaneous TCP sockets:**
  - each socket identified by its own 4-tuple
- **web servers have different sockets for each connecting client**
  - non-persistent HTTP will have different socket for each request

# Connection-oriented demux: example



host: IP
address A

application

P3

transport

network

link

physical

application

P4   P5   P6

transport

network

link

physical

server: IP
address B

application

P2   P3

transport

network

link

physical

host: IP
address C

source IP,port: B,80
dest IP,port: A,9157

source IP,port: A,9157
dest IP, port: B,80

source IP,port: C,5775
dest IP,port: B,80

source IP,port: C,9157
dest IP,port: B,80

three segments, all destined to IP address: B,
dest port: 80 are demultiplexed to *different* sockets

# Connection-oriented demux: example
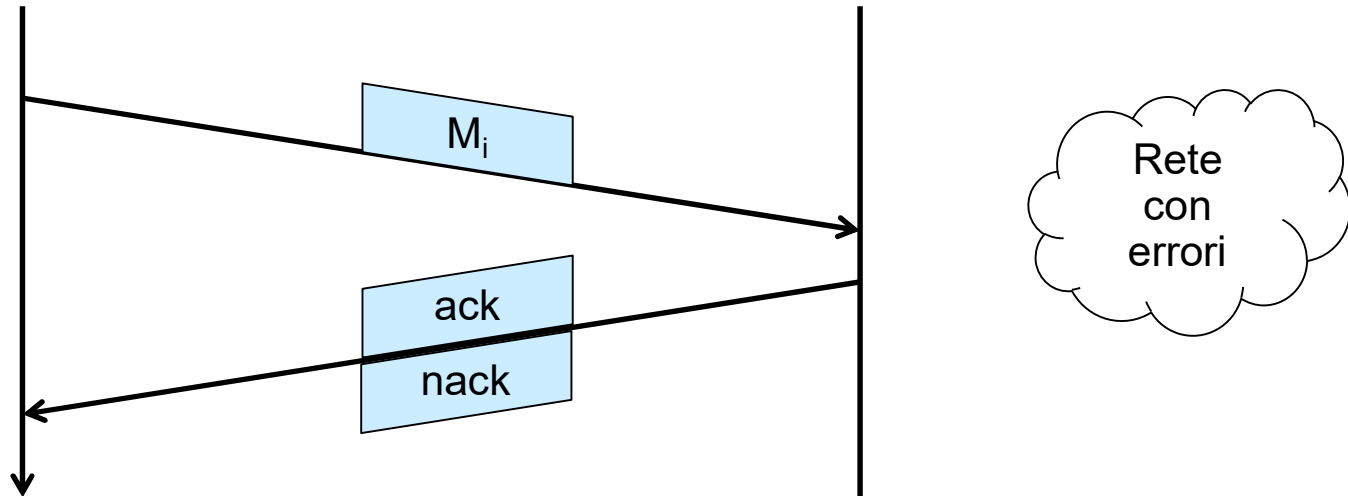
threaded server

application

P4

transport

network

link

physical

server: IP
address B

application

P3

transport

network

link

physical

host: IP
address A

application

P2    P3

transport

network

link

physical

host: IP
address C

source IP,port: B,80
dest IP,port: A,9157

source IP,port: A,9157
dest IP, port: B,80

source IP,port: C,5775
dest IP,port: B,80

source IP,port: C,9157
dest IP,port: B,80

# UDP: User Datagram Protocol [RFC 768]

- "no frills," "bare bones" Internet transport protocol
- "best effort" service, UDP segments may be:
  - lost
  - delivered out-of-order to app
- *connectionless:*
  - no handshaking between UDP sender, receiver
  - each UDP segment handled independently of others

- UDP use:
  - streaming multimedia apps (loss tolerant, rate sensitive)
  - DNS
  - SNMP
- reliable transfer over UDP:
  - add reliability at application layer
  - application-specific error recovery!

# Trasporto affidabile (1/10)



- Se la rete è «perfetta», ossia non introduce
    - errori sui bit
    - scarti
    - fuori sequenza
- Lo strato di trasporto non ha nulla da correggere e il protocollo è banale: il sender invia i messaggi uno dopo l'altro e il receiver li riceve tutti senza necessità di controlli
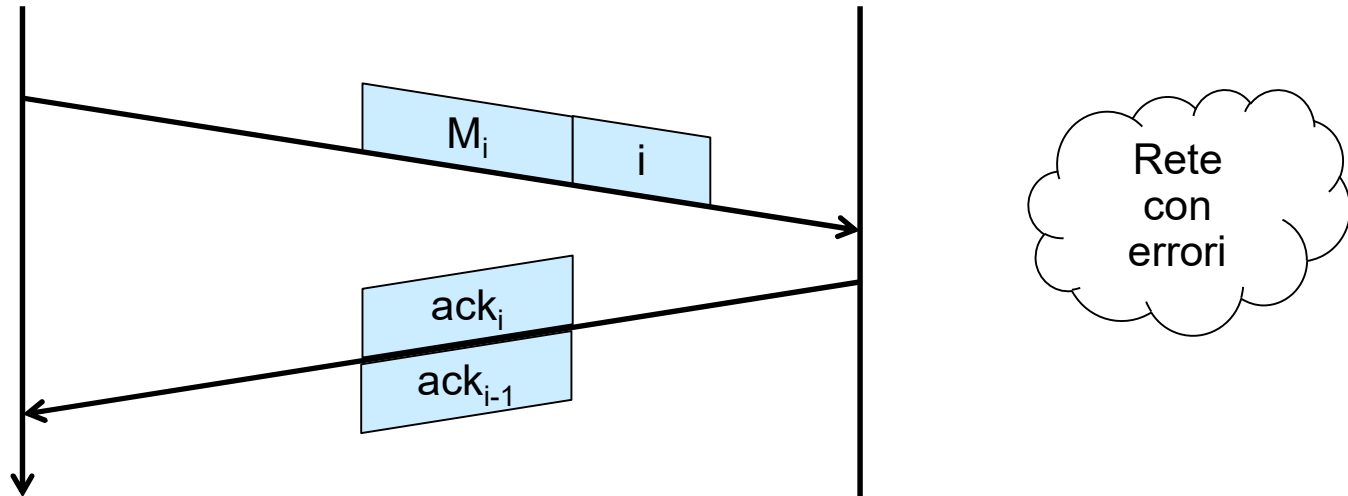
- In una rete con errori posso pensare di introdurre ack positivi e negativi. Ma il semplice algoritmo del sender:

```
IF ack
THEN M_{i+1}
ELSE IF nack
        THEN M_i
        ELSE ?
```
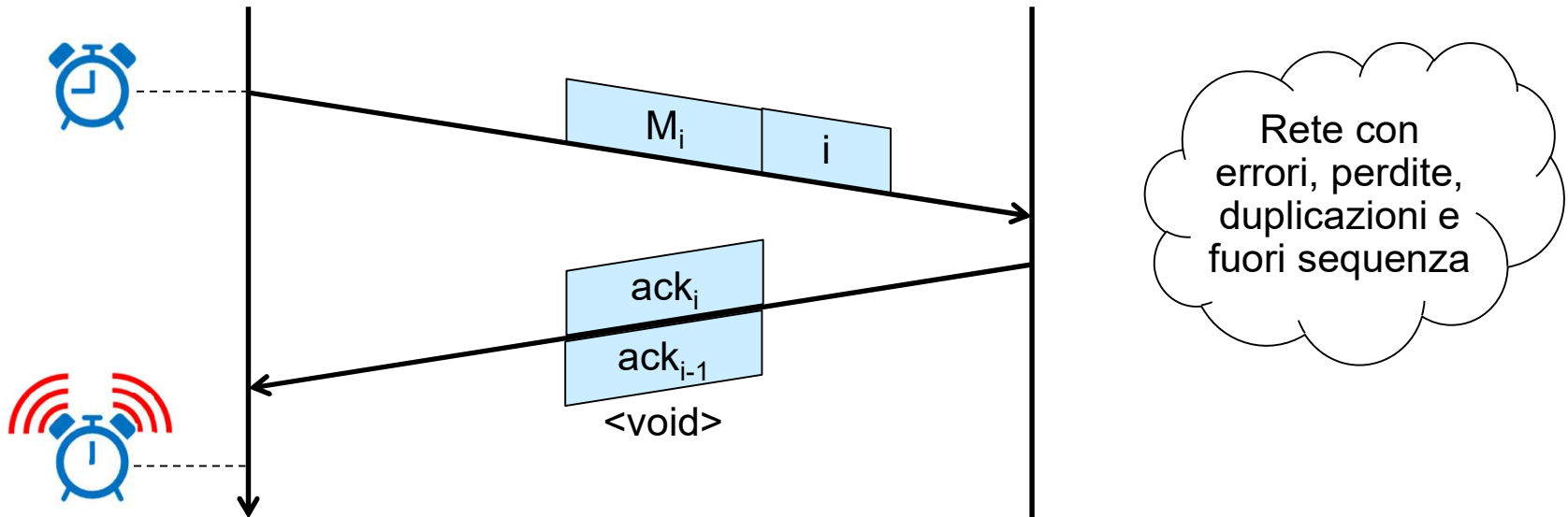
Non funziona! Perché anche ack/nack possono avere errori!

- Se numero i pacchetti non corro rischi anche se invio duplicati:

```
IF ack_i
THEN M_{i+1}
ELSE M_i
```

Inoltre numerando gli ack posso eliminare la necessità dei nack grazie alla regola: un secondo $ack_{i-1}$ equivale ad un $nack_i$

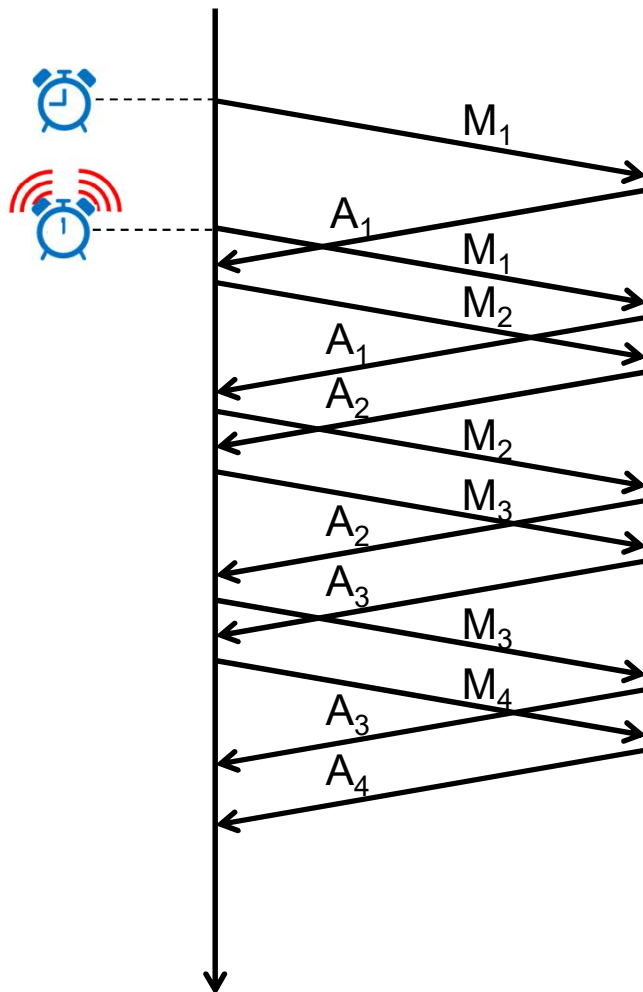- Purtroppo le reti con errori e senza perdite non esistono

- Se inseriamo un timer possiamo gestire le perdite (del pacchetto o degli ack)

```
IF T-off
THEN Mi,set T
ELSE IF acki
    THEN Mi+1,set T
    ELSE Mi,set T
```
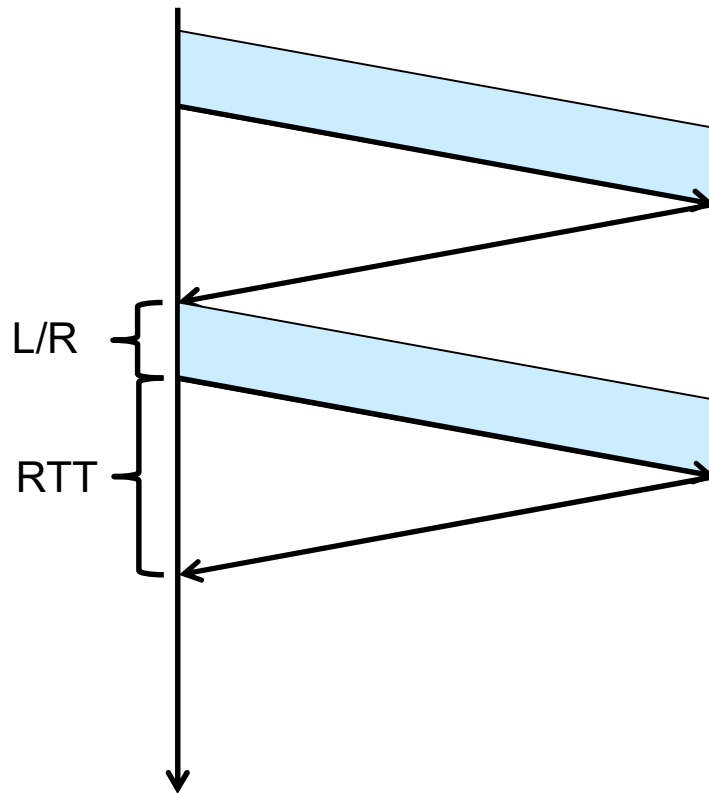
- Fissare la deadline di un timer è un problema complesso

- Un timer troppo corto può generare ritrasmissioni inutili

- Nel caso visto finora (trasmissione e riscontro di un pacchetto per volta) si possono addirittura generare sequenze molto lunghe di ritrasmissioni inutili

- Un timer troppo lungo ferma la trasmissione per troppo tempo nel caso di una perdita
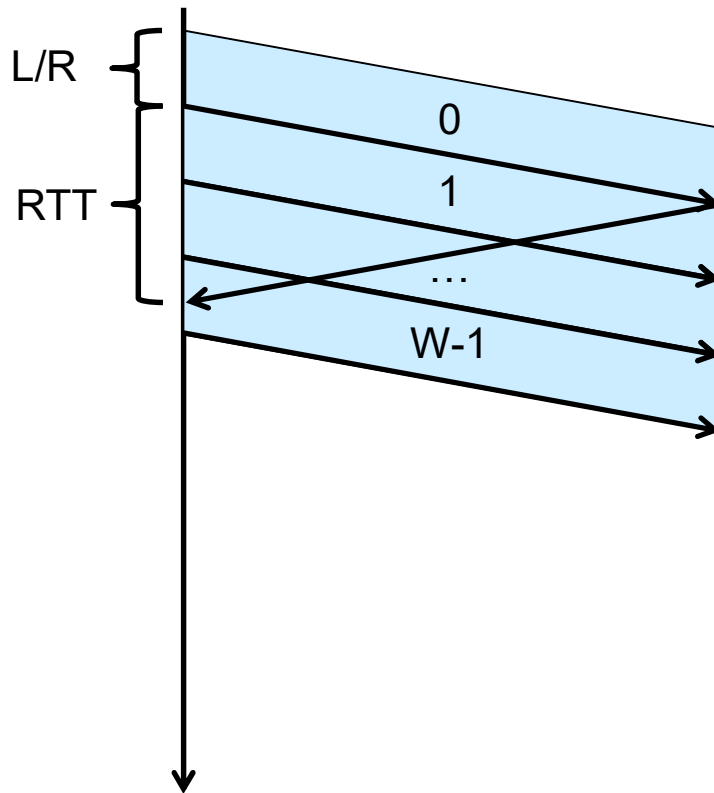
# Trasporto affidabile (6/10)



- Trasmettere un pacchetto per volta ed aspettare il riscontro prima del successivo porta a limitare in maniera significativa le prestazioni

$$U = \frac{L/R}{RTT + L/R}$$

- Esempio: RTT=100ms, L=1kbyte, R=100Mbit/s

U=0,0008

- I protocolli a finestra trasmettono fino a W pacchetti in attesa di ricevere il riscontro del primo

- La condizione per una trasmissione continua è che la finestra non si chiuda prima dell'arrivo del primo ack
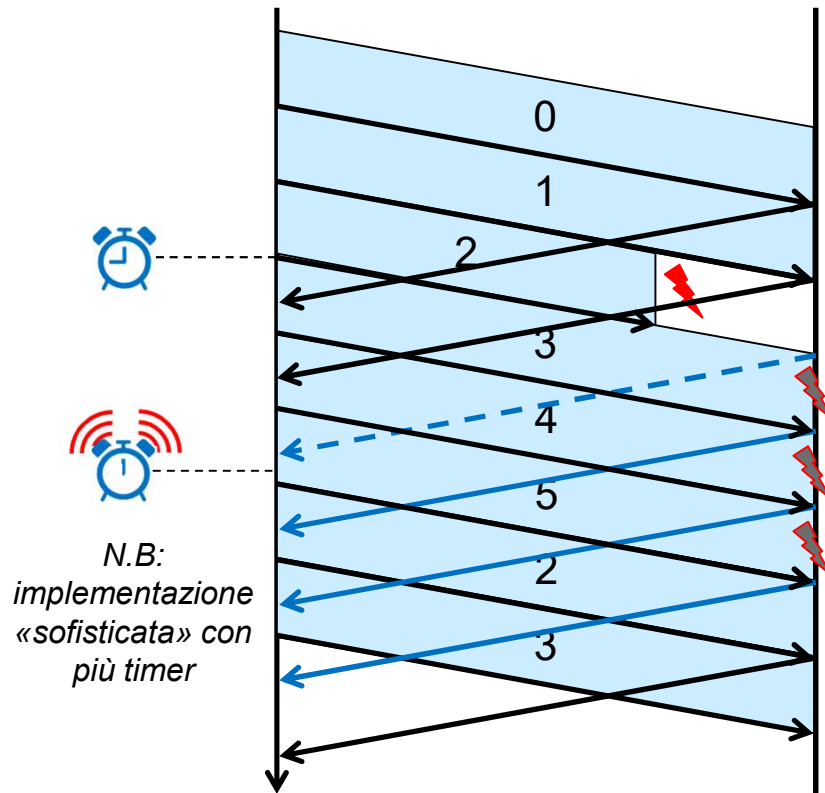
$$W \cdot {}^L\!/_R \geq RTT + {}^L\!/_R$$

$$W \geq \frac{RTT \cdot R}{L} + 1$$

- Nell'esempio precedente

$$W \geq 1251$$

- Cosa ritrasmetto quando perdo un pacchetto?

## Go-Back-N

- o trasmetto ed avanzo la finestra ad ogni ack
- o se T-off o ack ripetuto ritrasmetto tutto dall'ultimo pacchetto riscontrato
- o nessun buffer sul ricevitore: i pacchetti ricevuti dopo una perdita vengono scartati fino a quando il pacchetto in sequenza viene ritrasmesso
- o gli ack sono cumulativi (si recuperano molti eventi di perdita di ack)

*N.B: implementazione «sofisticata» con più timer*

## Selective Repeat

- gli ack sono individuali
- ogni pacchetto ha un suo timer
- se T-off ritrasmetto solo il pacchetto perso
- il ricevitore mantiene un buffer in cui i pacchetti vengono risequenziati e inviati in ordine all'applicazione

0

1

2

3

4

5

2

6

# Trasporto affidabile (10/10)
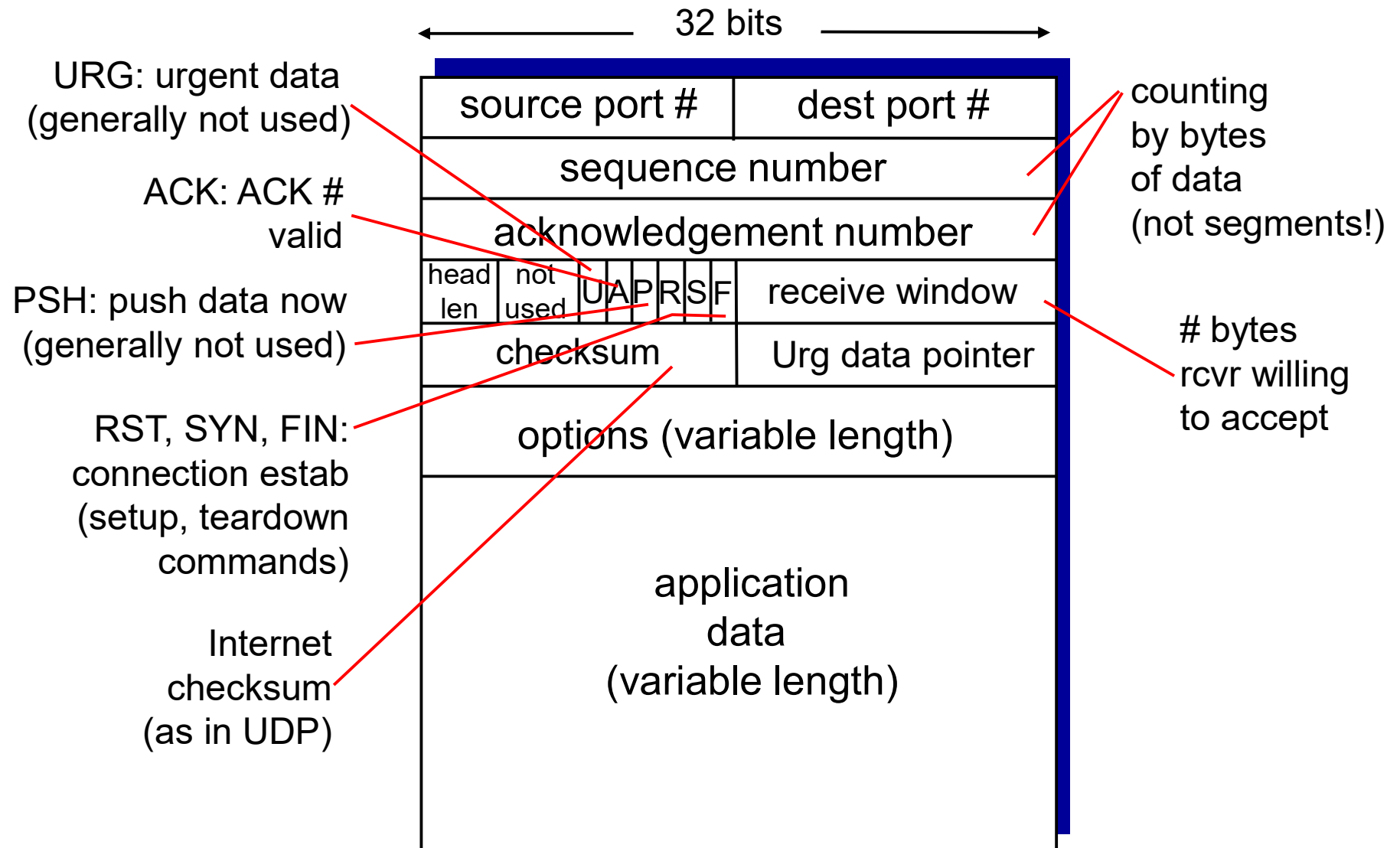
## Osservazioni a margine

- Non esistono solo due possibilità (Go-Back-N puro o Selective Repeat puro): i protocolli reali (come TCP) sono spesso un ibrido di soluzioni

- I numeri di sequenza sono rappresentati con un numero finito di bit. Ci sono delle regole che legano il valore massimo del numero di sequenza alla dimensione della finestra e impediscono che dopo un ritorno a zero si creino confusioni

  - Go-Back-N: $N \geq W+1$
  - Selective Repeat: $N \geq 2W$

# TCP: Overview  RFCs: 793,1122,1323, 2018, 2581

- **point-to-point:**
  - one sender, one receiver
- **reliable, in-order *byte stream:***
  - no "message boundaries"
- **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) inits sender, receiver state before data exchange
- **flow controlled:**
  - sender will not overwhelm receiver

# TCP segment structure

URG: urgent data
(generally not used)

ACK: ACK #
valid

PSH: push data now
(generally not used)

RST, SYN, FIN:
connection estab
(setup, teardown
commands)

Internet
checksum
(as in UDP)

32 bits

| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |
| head len | not used | U A P R S F | receive window |
| checksum | Urg data pointer |
| options (variable length) | |
| application data (variable length) | |

counting
by bytes
of data
(not segments!)

# bytes
rcvr willing
to accept

# TCP seq. numbers, ACKs

## sequence numbers:
- byte stream "number" of first byte in segment's data

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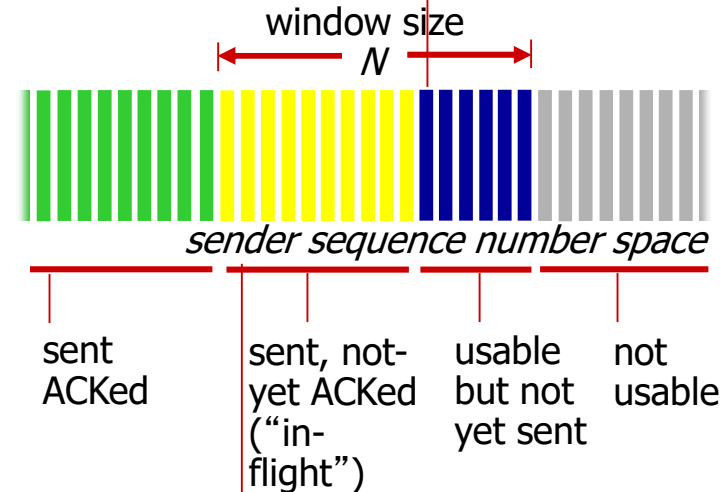- seq # of next byte expected from other side
- cumulative ACK

Q: how receiver handles out-of-order segments
- A: TCP spec doesn't say, - up to implementor

outgoing segment from sender

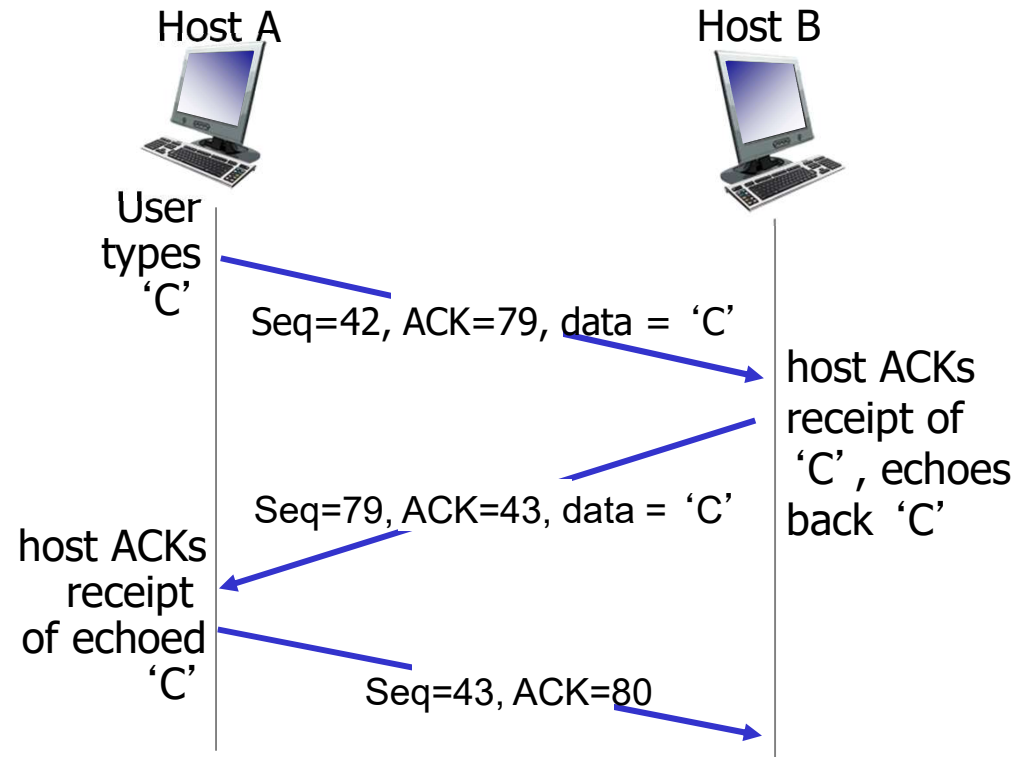| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |
| | rwnd |
| checksum | urg pointer |

window size
N



sender sequence number space

sent ACKed

sent, not-yet ACKed ("in-flight")

usable but not yet sent

not usable

incoming segment to sender

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

# TCP seq. numbers, ACKs

Host A                          Host B

User
types
'C'
    Seq=42, ACK=79, data = 'C'

                         host ACKs
                         receipt of
                         'C', echoes

    Seq=79, ACK=43, data = 'C'   back 'C'

host ACKs
receipt
of echoed
'C'
    Seq=43, ACK=80

simple telnet scenario

# TCP round trip time, timeout

Q: how to set TCP timeout value?

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

Q: how to estimate RTT?

- `SampleRTT`: measured time from segment transmission until ACK receipt
  - ignore retransmissions
- `SampleRTT` will vary, want estimated RTT "smoother"
  - average several *recent* measurements, not just current `SampleRTT`

# TCP reliable data transfer

- **TCP creates rdt service on top of IP's unreliable service**
  - pipelined segments
  - cumulative acks
  - single retransmission timer
- **retransmissions triggered by:**
  - timeout events
  - duplicate acks

let's initially consider simplified TCP sender:
  - ignore duplicate acks
  - ignore flow control, congestion control

# TCP sender events:

**data rcvd from app:**

- create segment with seq #
- seq # is byte-stream number of first data byte in segment
- start timer if not already running
  - think of timer as for oldest unacked segment
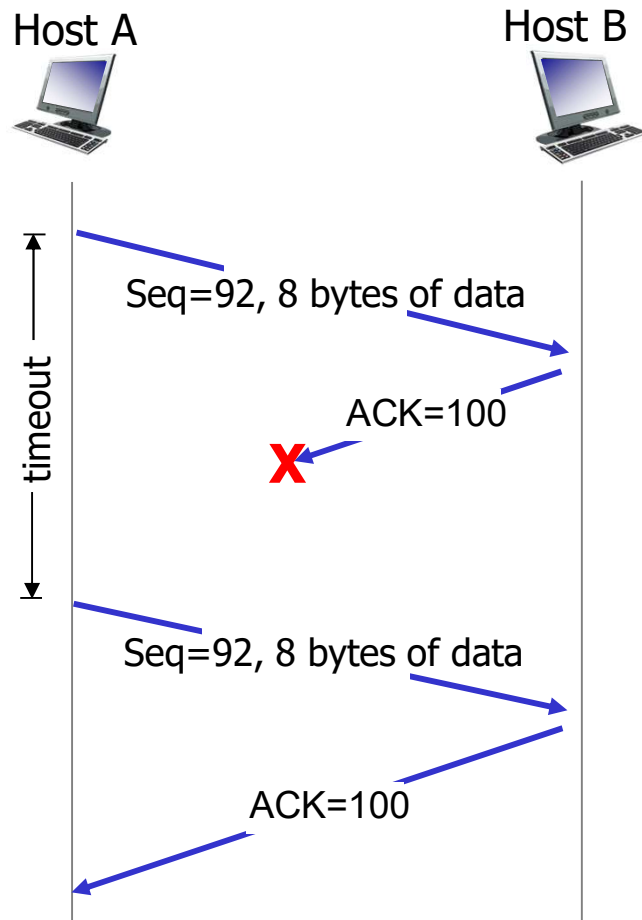  - expiration interval: `TimeOutInterval`

**timeout:**
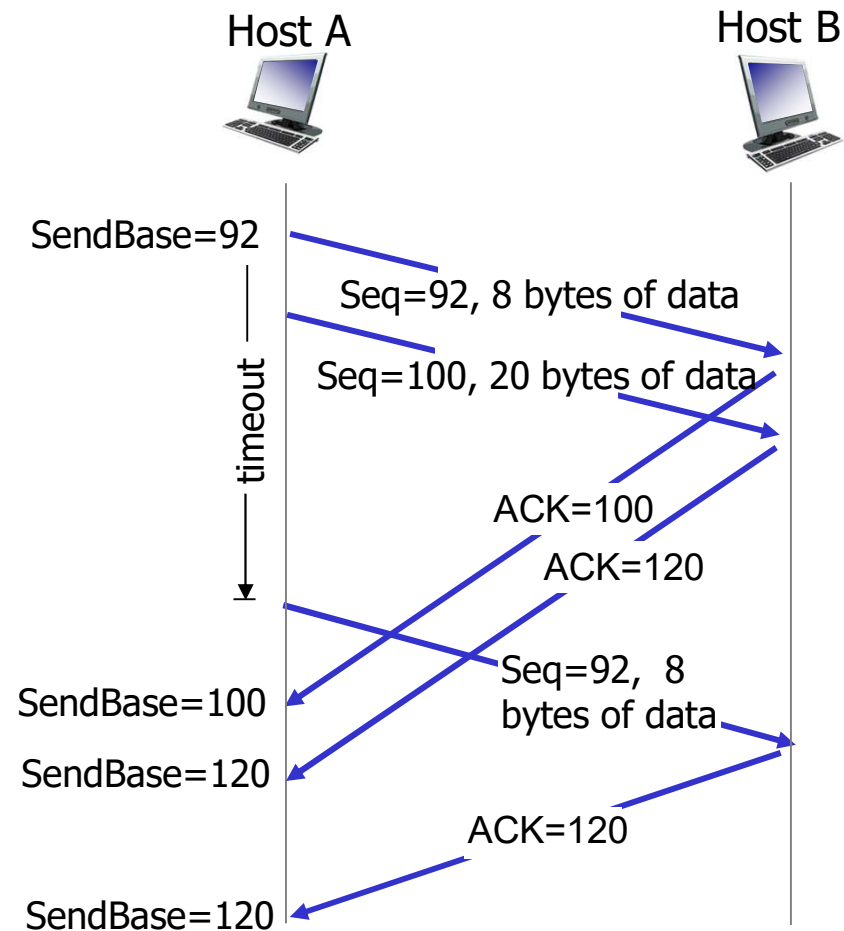
- retransmit segment that caused timeout
- restart timer

**ack rcvd:**

- if ack acknowledges previously unacked segments
  - update what is known to be ACKed
  - start timer if there are still unacked segments
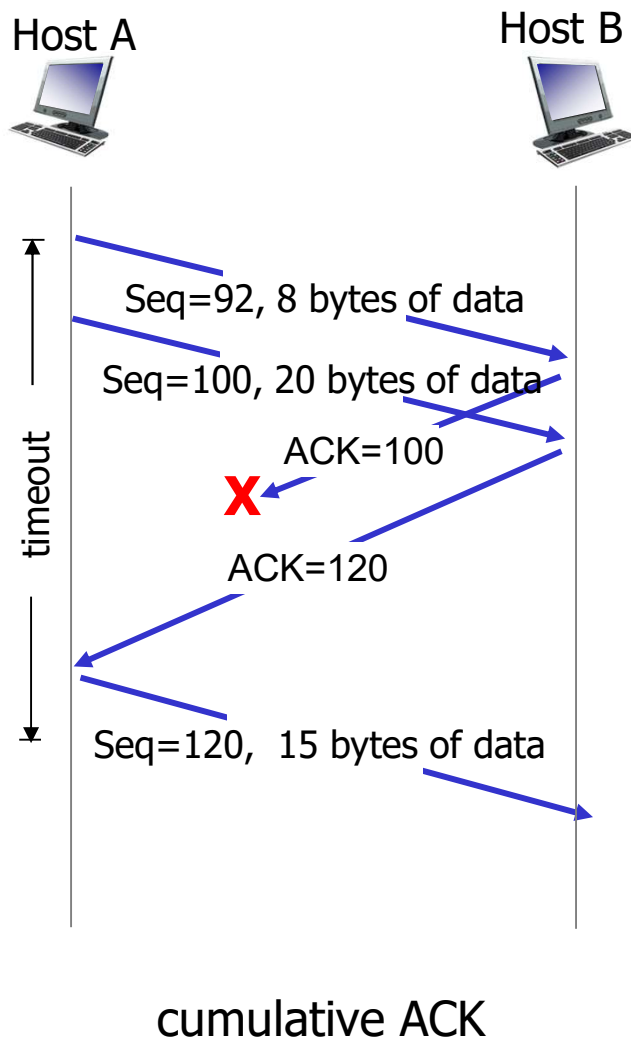
# TCP: retransmission scenarios



lost ACK scenario

premature timeout
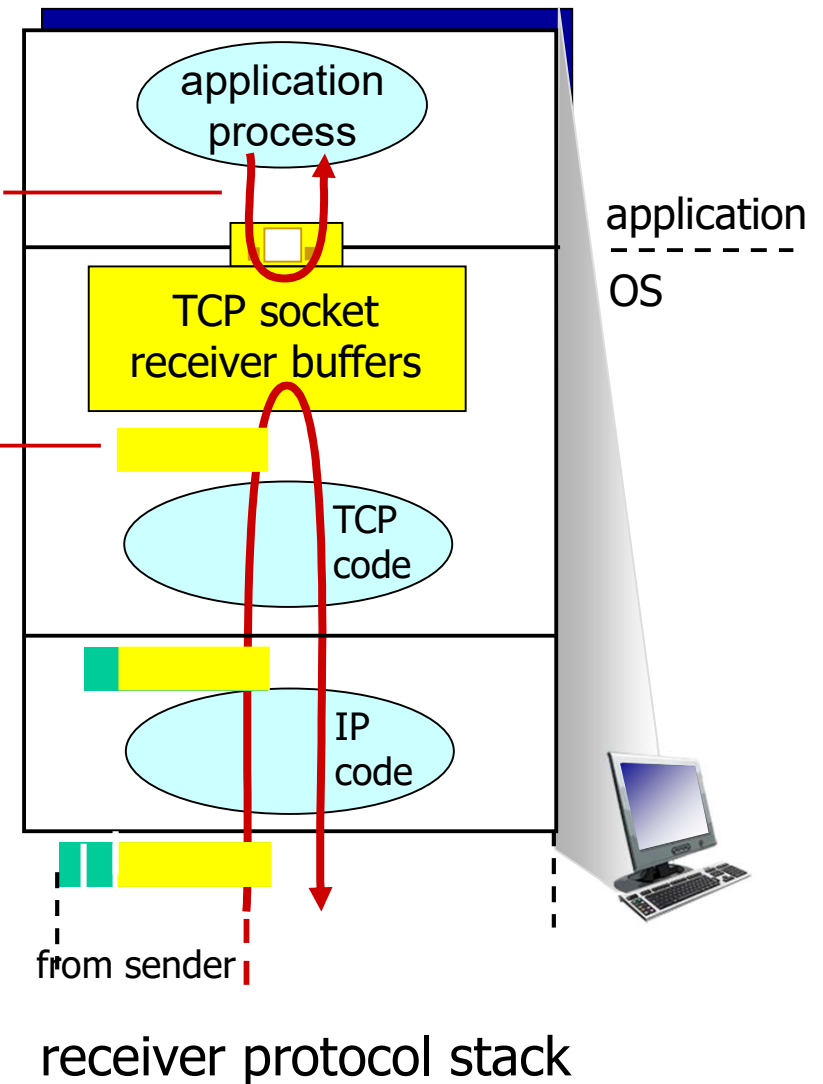
# TCP: retransmission scenarios

Host A

Host B

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

ACK=100

X

ACK=120

timeout

Seq=120, 15 bytes of data

cumulative ACK

# TCP flow control

application may remove data from TCP socket buffers ....

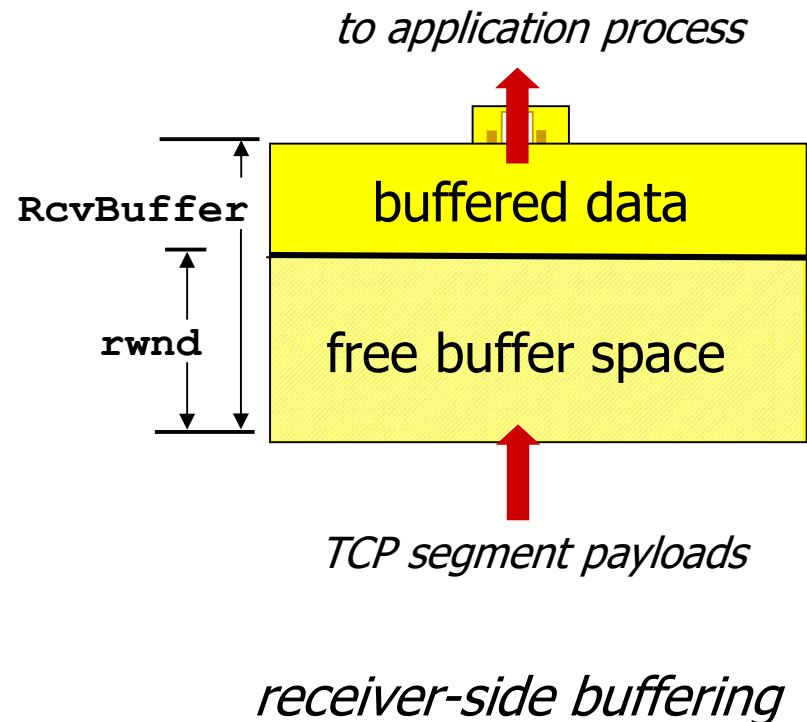... slower than TCP receiver is delivering (sender is sending)

application process

TCP socket receiver buffers

TCP code

IP code

application
OS

*flow control*
receiver controls sender, so sender won't overflow receiver's buffer by transmitting too much, too fast

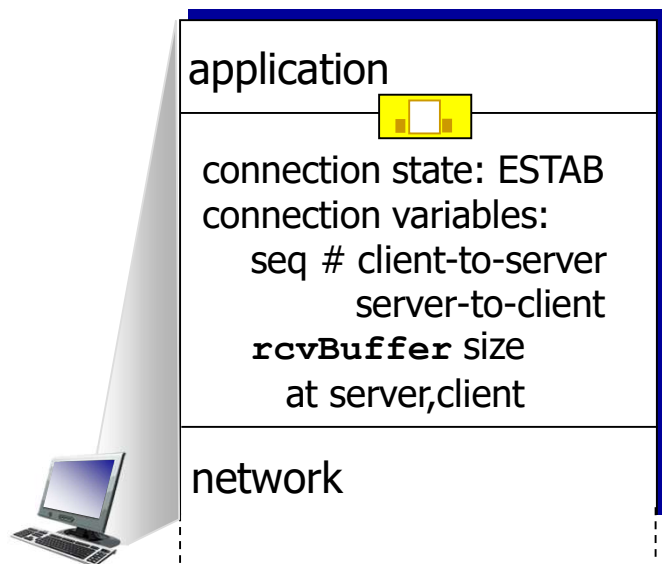from sender

receiver protocol stack

# TCP flow control

- receiver "advertises" free buffer space by including **rwnd** value in TCP header of receiver-to-sender segments

  - **RcvBuffer** size set via socket options (typical default is 4096 bytes)
  - many operating systems autoadjust **RcvBuffer**

- sender limits amount of unacked ("in-flight") data to receiver's **rwnd** value

- guarantees receive buffer will not overflow

*to application process*

RcvBuffer

buffered data

rwnd

free buffer space
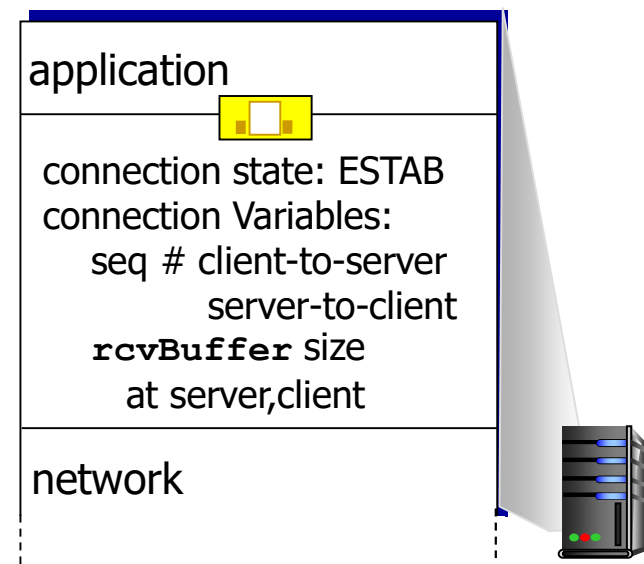
*TCP segment payloads*

*receiver-side buffering*

# Connection Management

before exchanging data, sender/receiver "handshake":

- agree to establish connection (each knowing the other willing to establish connection)
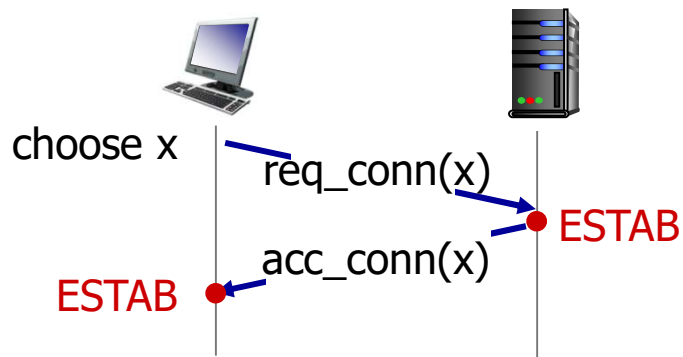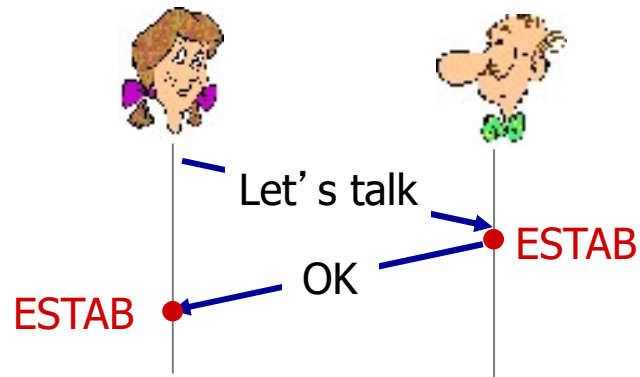- agree on connection parameters



application

connection state: ESTAB
connection variables:
    seq # client-to-server
        server-to-client
  **rcvBuffer** size
    at server,client

network

```
Socket clientSocket =
  newSocket("hostname","port
  number");
```

application

connection state: ESTAB
connection Variables:
    seq # client-to-server
        server-to-client
  **rcvBuffer** size
    at server,client

network

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

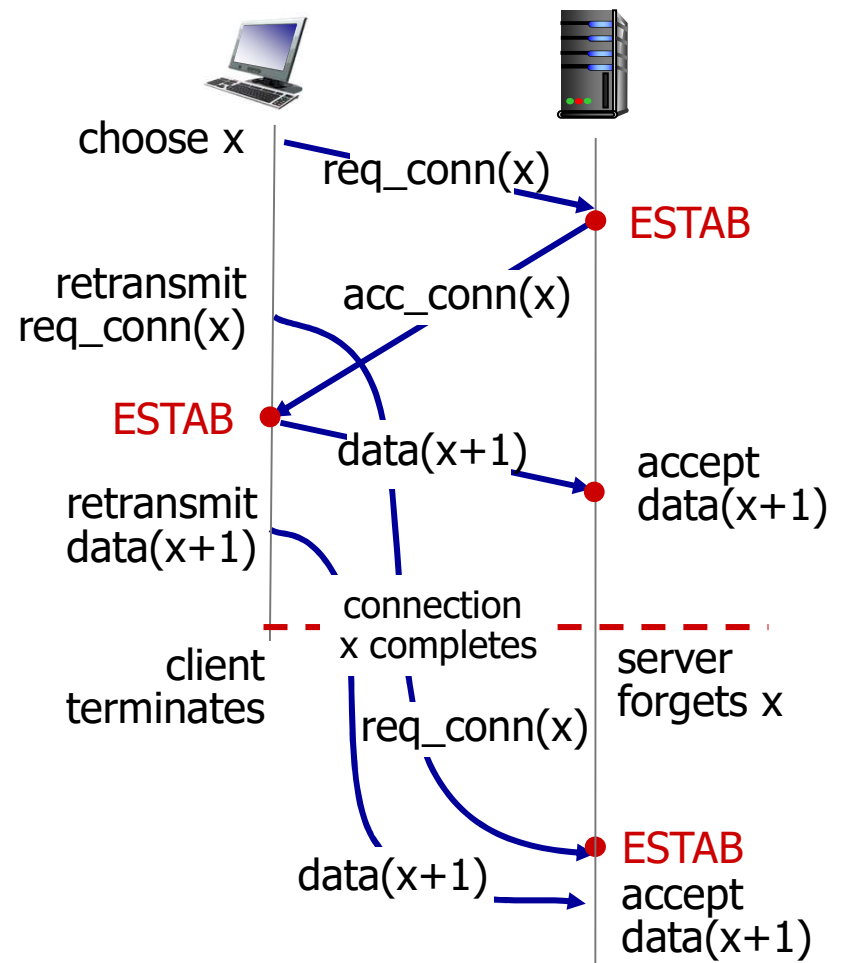# Agreeing to establish a connection
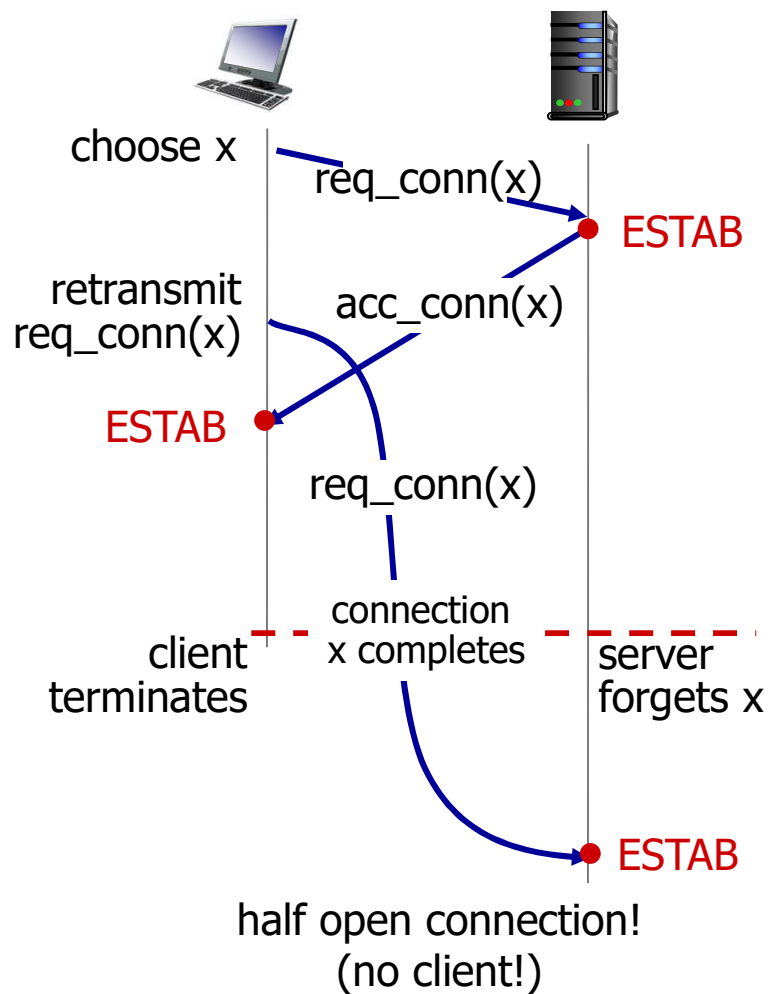
2-way handshake:
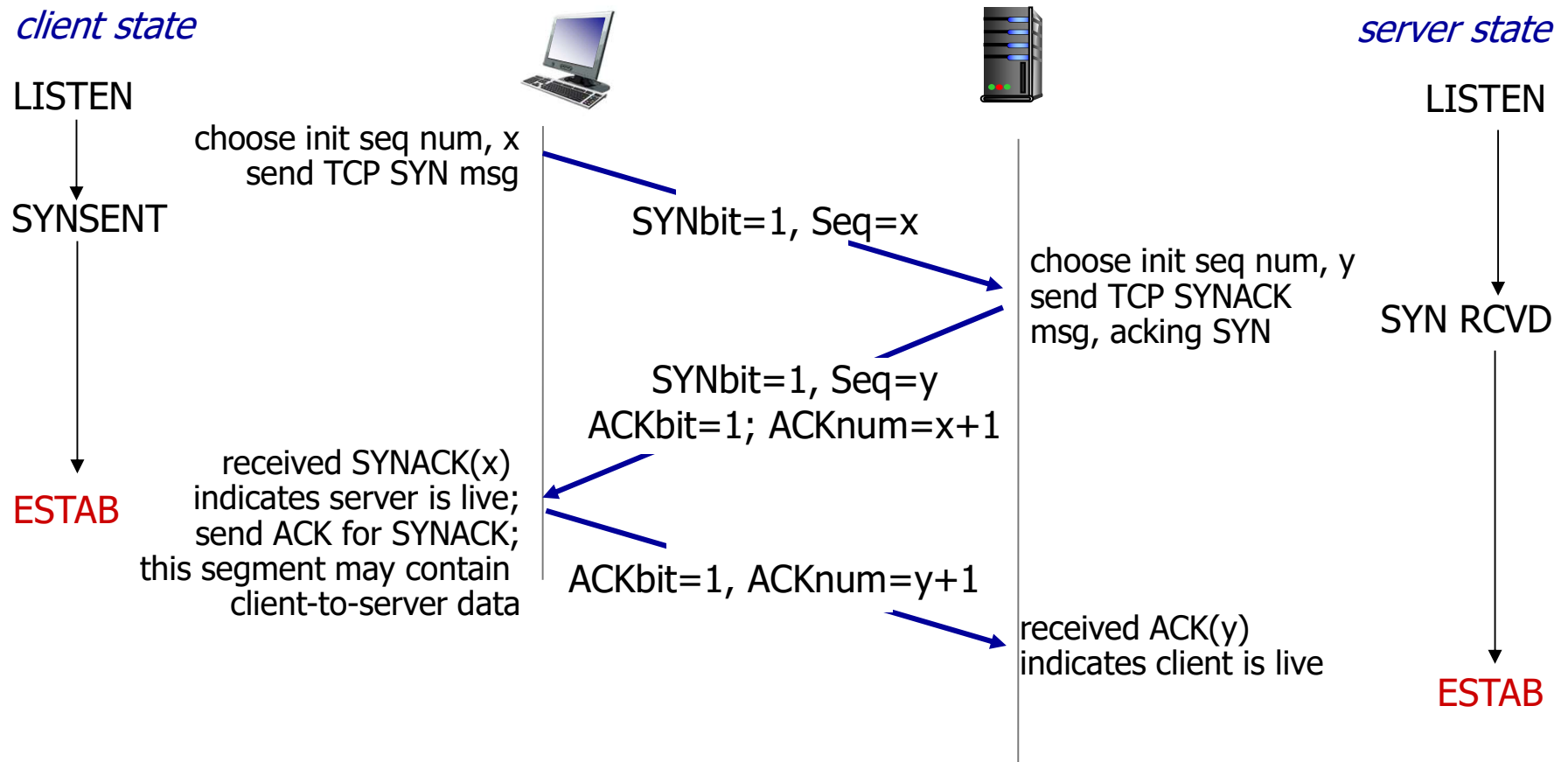


*Q:* will 2-way handshake always work in network?

- variable delays
- retransmitted messages (e.g. req_conn(x)) due to message loss
- message reordering
- can't "see" other side

# Agreeing to establish a connection

2-way handshake failure scenarios:

# TCP 3-way handshake

**client state**

LISTEN

choose init seq num, x
send TCP SYN msg

SYNSENT

SYNbit=1, Seq=x

SYNbit=1, Seq=y
ACKbit=1; ACKnum=x+1

received SYNACK(x)
indicates server is live;
send ACK for SYNACK;
this segment may contain
client-to-server data

ESTAB

ACKbit=1, ACKnum=y+1

**server state**

LISTEN

choose init seq num, y
send TCP SYNACK
msg, acking SYN

SYN RCVD

received ACK(y)
indicates client is live

ESTAB

# TCP: closing a connection

client state

server state

ESTAB

ESTAB

`clientSocket.close()`

FIN_WAIT_1    can no longer
send but can
receive data

FINbit=1, seq=x

CLOSE_WAIT

ACKbit=1; ACKnum=x+1

FIN_WAIT_2    wait for server
close

can still
send data

FINbit=1, seq=y

LAST_ACK

TIMED_WAIT

can no longer
send data

ACKbit=1; ACKnum=y+1
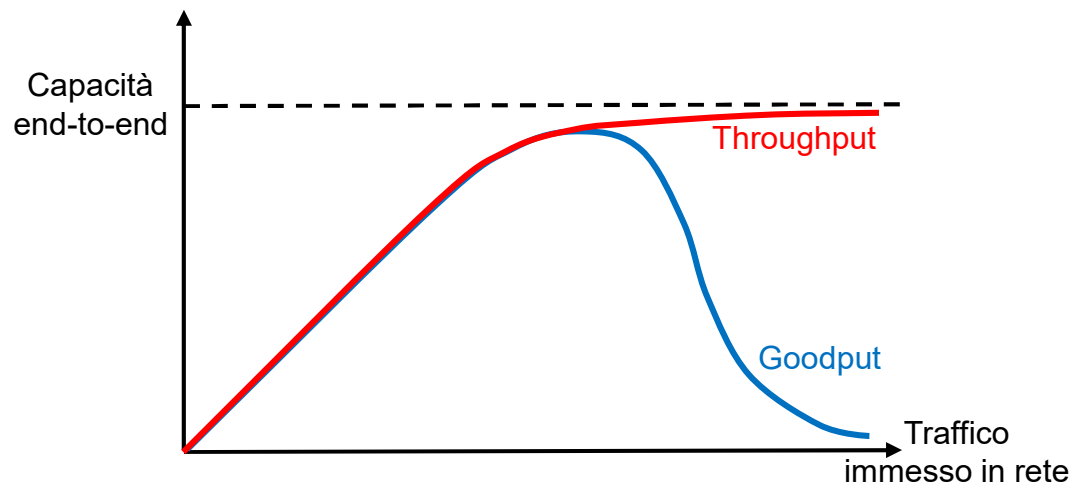
timed wait
for 2*max
segment lifetime

CLOSED

CLOSED

# Principles of congestion control

*congestion:*

- informally: "too many sources sending too much data too fast for *network* to handle"
- different from flow control!
- manifestations:
  - lost packets (buffer overflow at routers)
  - long delays (queueing in router buffers)
- a top-10 problem!

# Congestione

- Se la rete si avvicina alla saturazione delle risorse disponibili (capacità), il ritardo e la percentuale di perdite cresce

- Se il trasporto ritrasmette, aumenta il numero medio di ritrasmissioni di ogni pacchetto

- Mentre il throughput (pacchetti che attraversano la rete) si avvicina al 100% della capacità, il «goodput» visto dall'applicazione decresce !
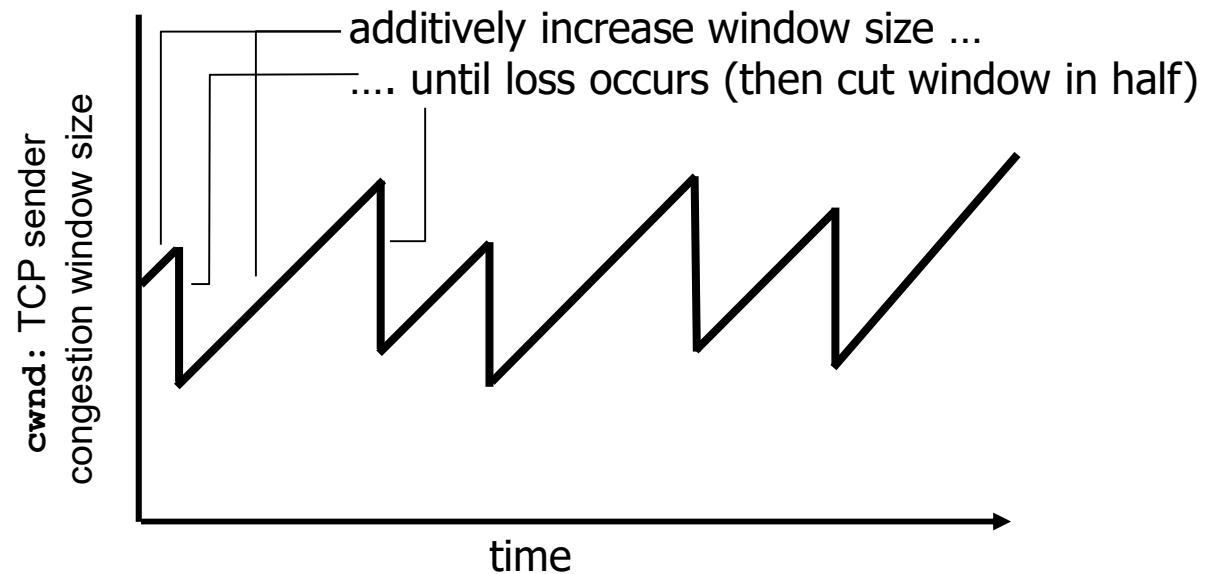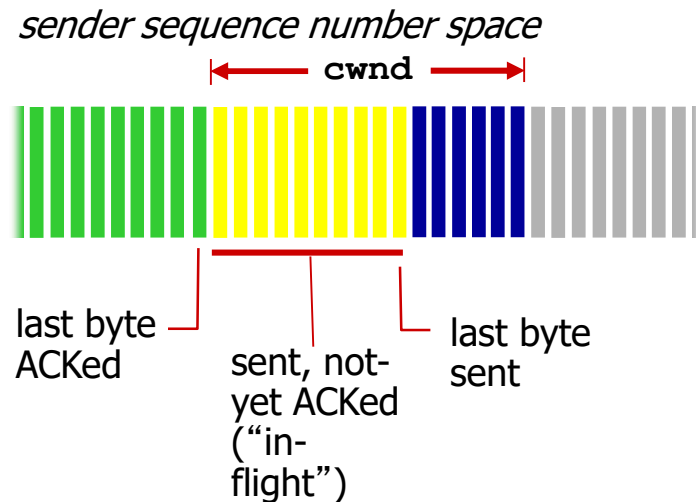
# TCP congestion control: additive increase multiplicative decrease

- *approach:* sender increases transmission rate (window size), probing for usable bandwidth, until loss occurs
  - *additive increase:* increase `cwnd` by 1 MSS every RTT until loss detected
  - *multiplicative decrease:* cut `cwnd` in half after loss

AIMD saw tooth behavior: probing for bandwidth

additively increase window size …

…. until loss occurs (then cut window in half)

**cwnd:** TCP sender congestion window size

time

# TCP Congestion Control: details

*sender sequence number space*



last byte ACKed

sent, not-yet ACKed ("in-flight")

last byte sent

- **sender limits transmission:**

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

- **cwnd** is dynamic, function of perceived network congestion

*TCP sending rate:*

- *roughly:* send cwnd bytes, wait RTT for ACKS, then send more bytes

$$\text{rate} \approx \frac{\text{cwnd}}{\text{RTT}} \text{ bytes/sec}$$

# TCP Slow Start

- when connection begins, increase rate exponentially until first loss event:
  - initially `cwnd` = 1 MSS
  - double `cwnd` every RTT
  - done by incrementing `cwnd` for every ACK received
- *summary:* initial rate is slow but ramps up exponentially fast

Host A                                      Host B

RTT↑↓

one segment

two segments

four segments

time