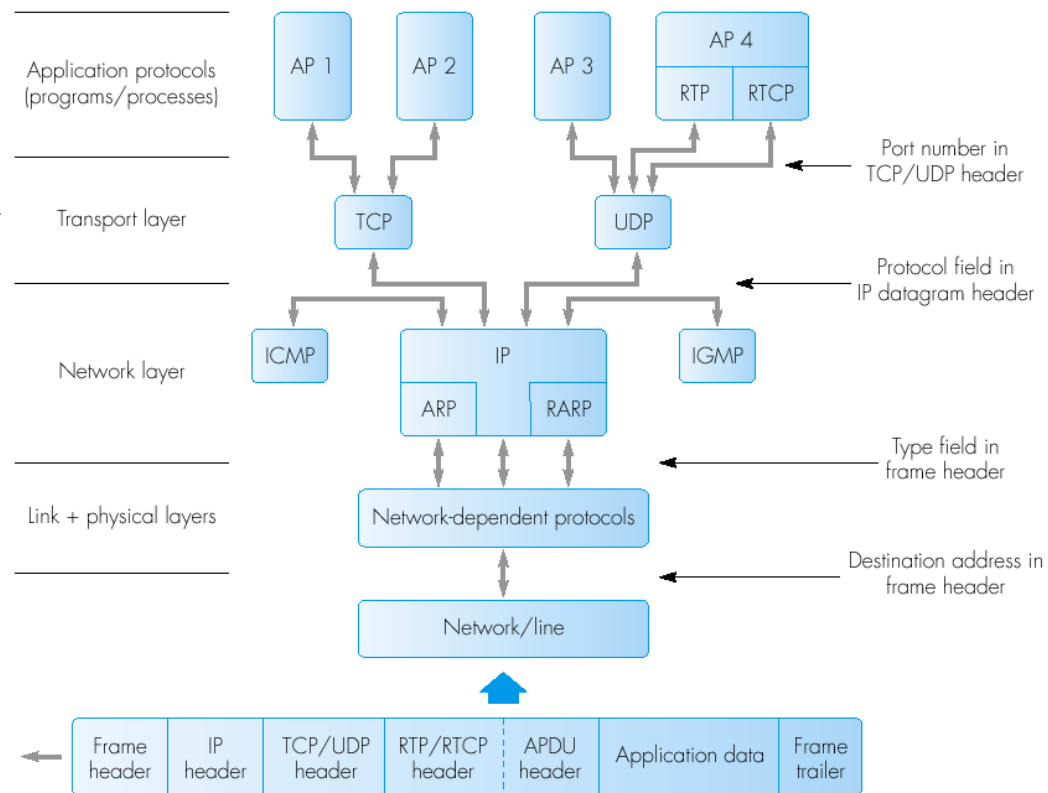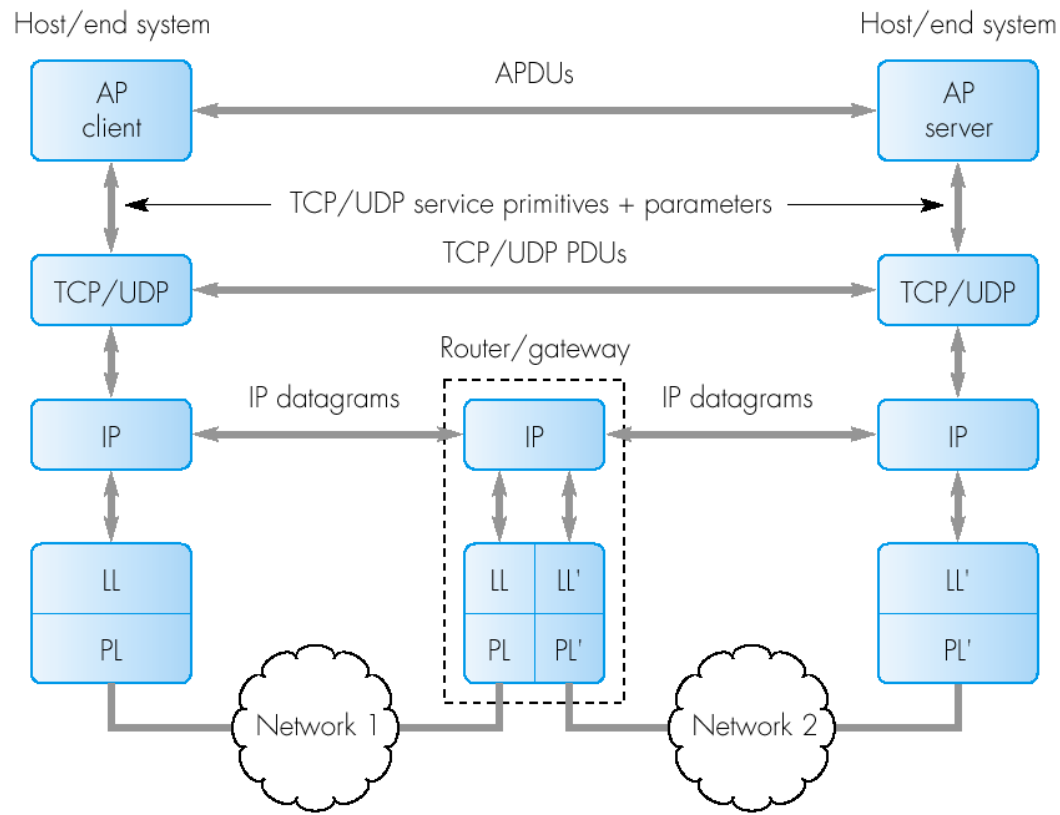# Transport Layer

## Figure 12.1 TCP/IP protocol suite and interlayer address selectors.

Layered architecture

Application protocols (programs/processes)

AP 1   AP 2   AP 3   AP 4
RTP   RTCP

Port number in TCP/UDP header

Transport layer

TCP   UDP

Protocol field in IP datagram header

Network layer

ICMP   IP   IGMP
ARP   RARP

Type field in frame header

Link + physical layers

Network-dependent protocols

Destination address in frame header

Network/line

| Frame header | IP header | TCP/UDP header | RTP/RTCP header | APDU header | Application data | Frame trailer |

# Figure 12.2  TCP/IP protocol suite interlayer communications.

Host/end system

Host/end system

APDUs

AP
client

AP
server

TCP/UDP service primitives + parameters

TCP/UDP PDUs

TCP/UDP

TCP/UDP

Router/gateway

IP datagrams

IP datagrams

IP

IP

IP

LL

LL

LL'

LL'

PL

PL

PL'

PL'

Network 1

Network 2

Transport Layer 3-3

# Chapter : Transport Layer

## our goals:

- understand principles behind transport layer services:
  - multiplexing, demultiplexing
  - flow control
  - congestion control

- learn about Internet transport layer protocols:
  - UDP: connectionless transport
  - TCP: connection-oriented reliable transport
  - TCP congestion control

# Chapter outline

3.1 transport-layer services

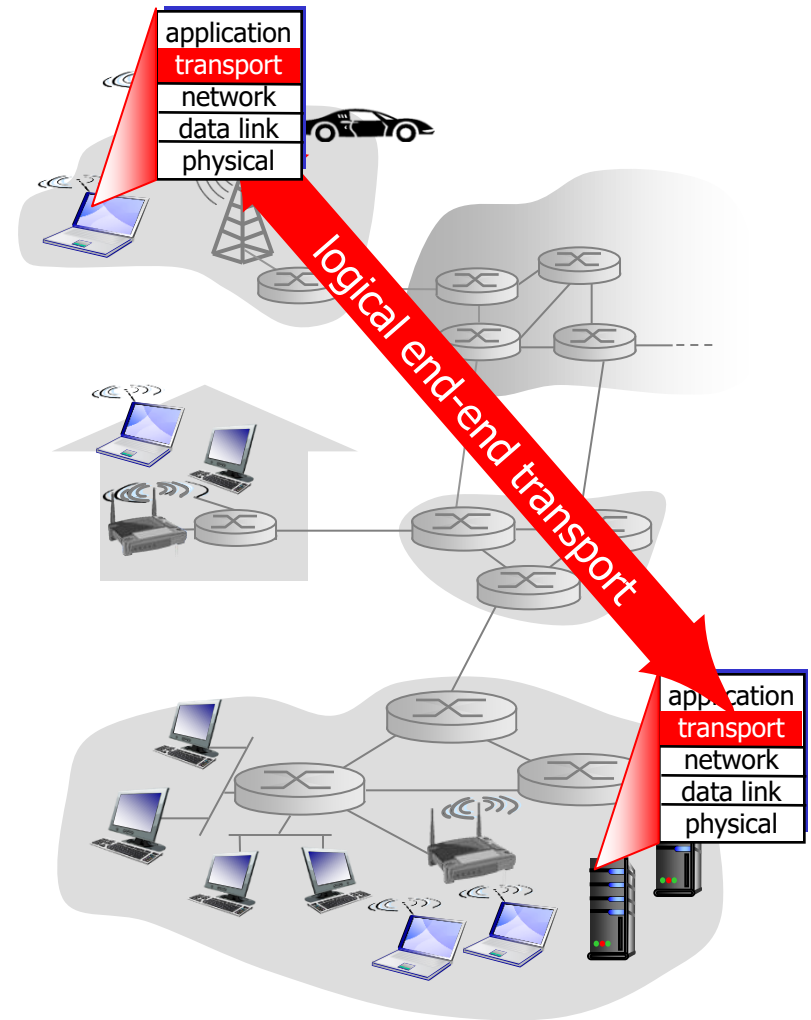3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management

3.6 TCP congestion control

# 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
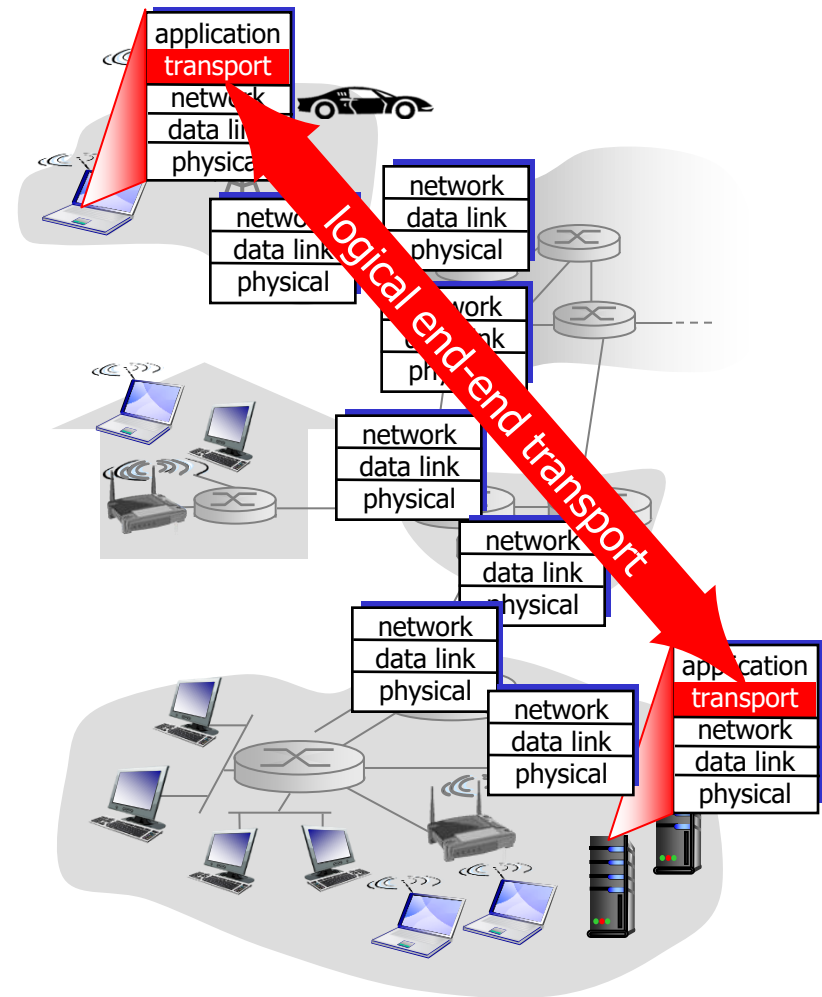  ▪ relies on, enhances, network layer services

*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



logical end-end transport

# Chapter 3 outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer

3.5 connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management

3.6 principles of congestion control
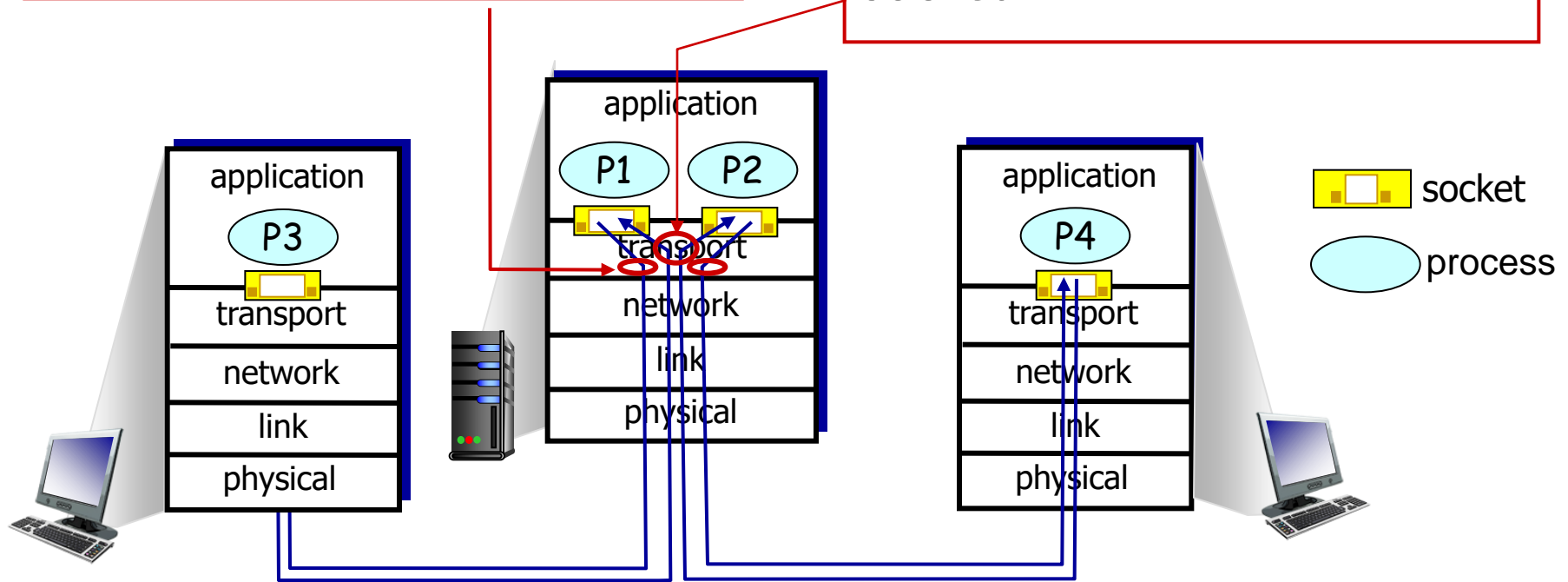
3.7 TCP congestion control

# Multiplexing/demultiplexing

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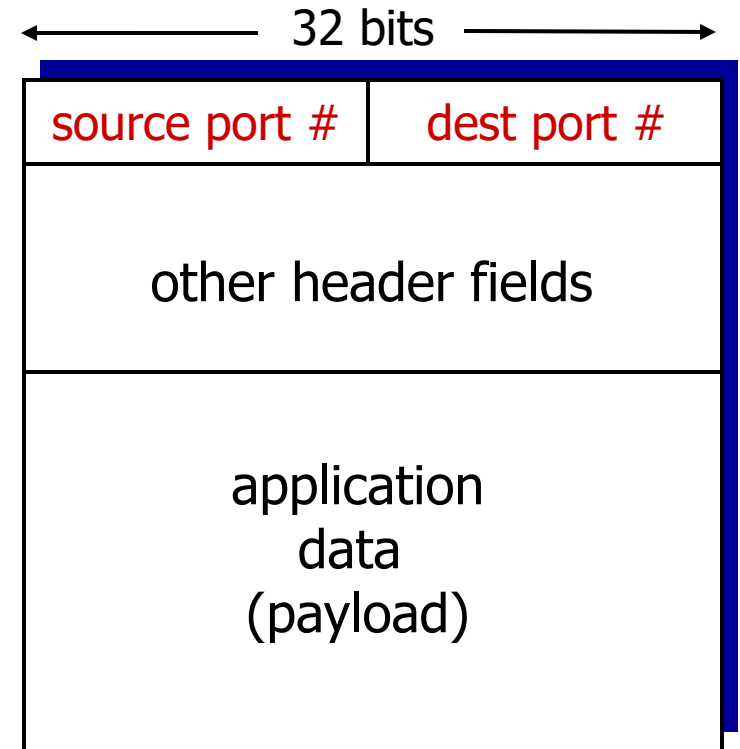handle data from multiple sockets, add transport header (later used for demultiplexing)

*demultiplexing at receiver:*
use header info to deliver received segments to correct socket

# 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

# Outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer

3.5 connection-oriented transport: TCP
- segment structure
- reliable data transfer
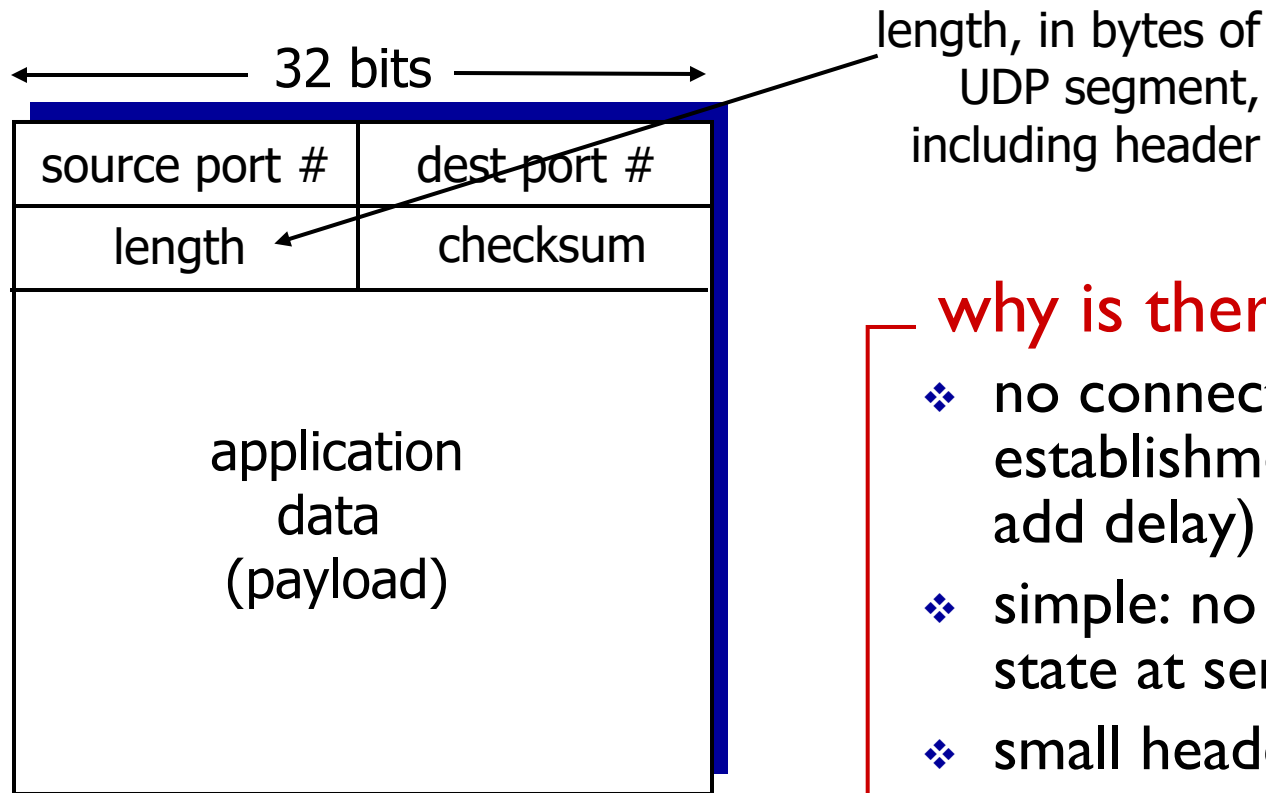- flow control
- connection management

3.6 principles of congestion control

3.7 TCP congestion control

# UDP: User Datagram Protocol [RFC 768]

- ❖ "no frills," 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!

# UDP: segment header

32 bits

| source port # | dest port # |
|---|---|
| length | checksum |

application
data
(payload)

UDP segment format

length, in bytes of
UDP segment,
including header

## why is there a UDP?

- ❖ no connection establishment (which can add delay)
- ❖ simple: no connection state at sender, receiver
- ❖ small header size
- ❖ no congestion control: UDP can blast away as fast as desired

# UDP checksum

*Goal:* detect "errors" (e.g., flipped bits) in transmitted segment

## sender:

❖ treat segment contents, including header fields, as sequence of 16-bit integers

❖ checksum: addition (one's complement sum) of segment contents

❖ sender puts checksum value into UDP checksum field

## receiver:

❖ compute checksum of received segment

❖ check if computed checksum equals checksum field value:

- NO - error detected
- YES - no error detected. *But maybe errors nonetheless?* More later ….

# Internet checksum: example

example: add two 16-bit integers

```
        1 1 1 0 0 1 1 0 0 1 1 0 0 1 1 0
        1 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1
```

wraparound  (1) 1 0 1 1 1 0 1 1 1 0 1 1 1 0 1 1

```
sum         1 0 1 1 1 0 1 1 1 0 1 1 1 1 0 0
checksum    0 1 0 0 0 1 0 0 0 1 0 0 0 0 1 1
```

*Note:* when adding numbers, a carryout from the most significant bit needs to be added to the result

# outline

3.1  transport-layer services

3.2  multiplexing and demultiplexing

3.3  connectionless transport: UDP

3.5  connection-oriented transport: TCP
- segment structure
- reliable data transfer
- flow control
- connection management
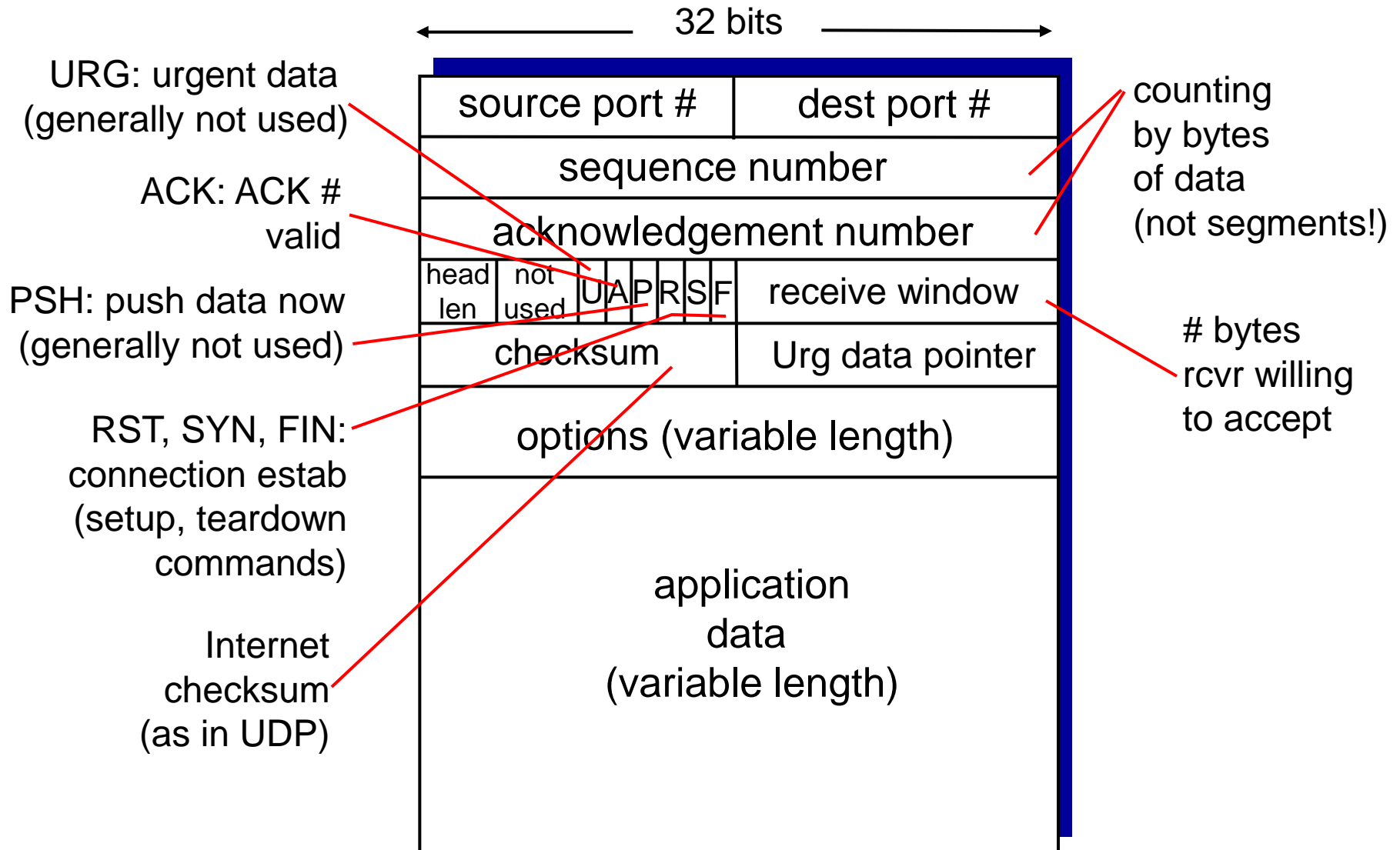
3.6  principles of congestion control

3.7 TCP congestion control

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

- ❖ point-to-point:
  - one sender, one receiver

- ❖ reliable, in-order *byte steam:*
  - 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

32 bits

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)

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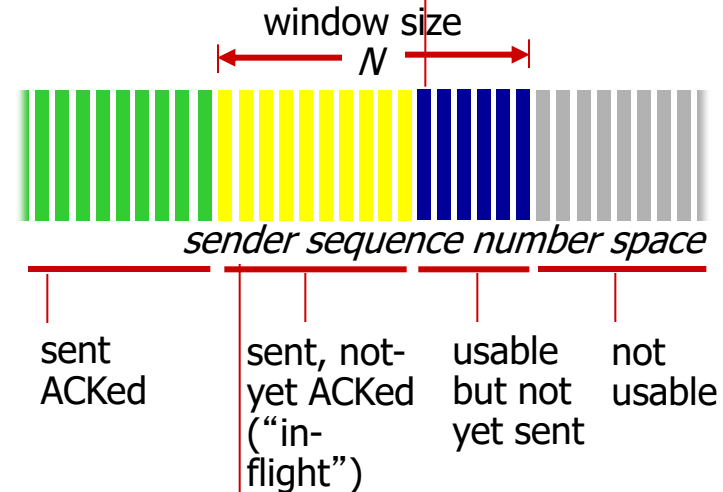- byte stream "number" of first byte in segment's data

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

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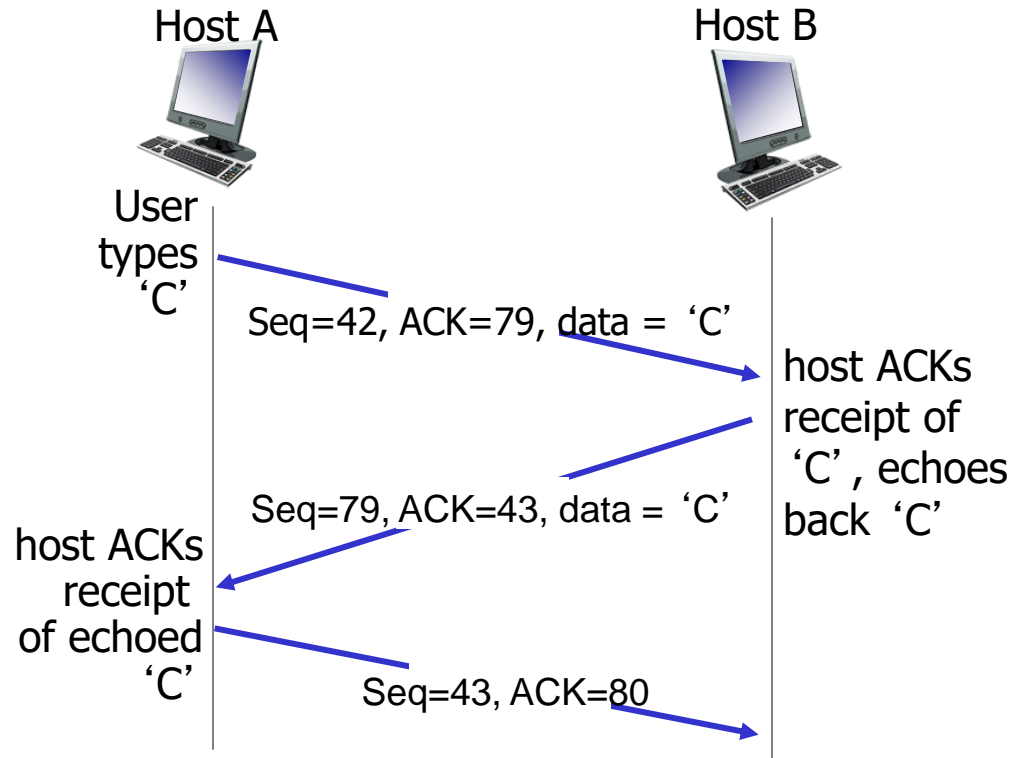- A: TCP spec doesn't say, - up to implementor

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 round trip time, timeout

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

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



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

RTT (milliseconds)

time (seconds)

◆ sampleRTT
■ EstimatedRTT

# TCP round trip time, timeout

❖ **timeout interval: `EstimatedRTT` plus "safety margin"**
  - large variation in `EstimatedRTT` -> larger safety margin

❖ estimate SampleRTT deviation from EstimatedRTT:

```
DevRTT = (1-β)*DevRTT +
              β*|SampleRTT-EstimatedRTT|

         (typically, β = 0.25)
```

**`TimeoutInterval = EstimatedRTT + 4*DevRTT`**

estimated RTT          "safety margin"

# TCP flow control

application may
remove data from
TCP socket buffers ….

application
process

application
- - - - - - - -
OS

TCP socket
receiver buffers

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

TCP
code

IP
code

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

from sender

receiver protocol stack
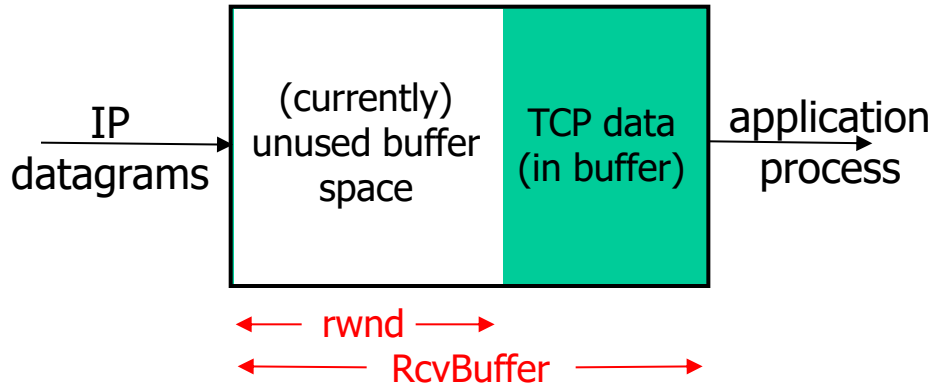
# TCP flow control

❖ 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
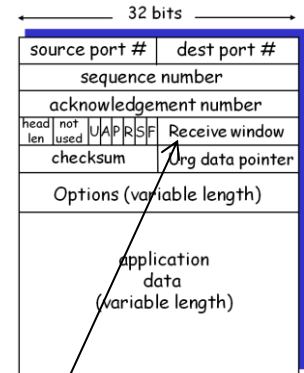
❖ guarantees receive buffer will not overflow

*to application process*

**RcvBuffer**

buffered data

**rwnd**

free buffer space

*TCP segment payloads*

*receiver-side buffering*

# TCP Flow control: how it works



(suppose TCP receiver discards out-of-order segments)

- ❖ unused buffer space:

```
= rwnd
= RcvBuffer-[LastByteRcvd -
     LastByteRead]
```

- ❖ receiver: advertises unused buffer space by including rwnd value in segment header

- ❖ sender: limits # of unACKed bytes to rwnd
  - ■ guarantees receiver's buffer doesn't overflow

# TCP Connection Management

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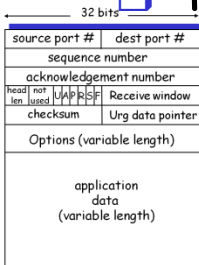


- ❑ initialize TCP variables:
  - ○ seq. #s
  - ○ buffers, flow control info (e.g. `RcvWindow`)
- ❑ *client:* connection initiator

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

- ❑ *server:* contacted by client

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

## Three way handshake:

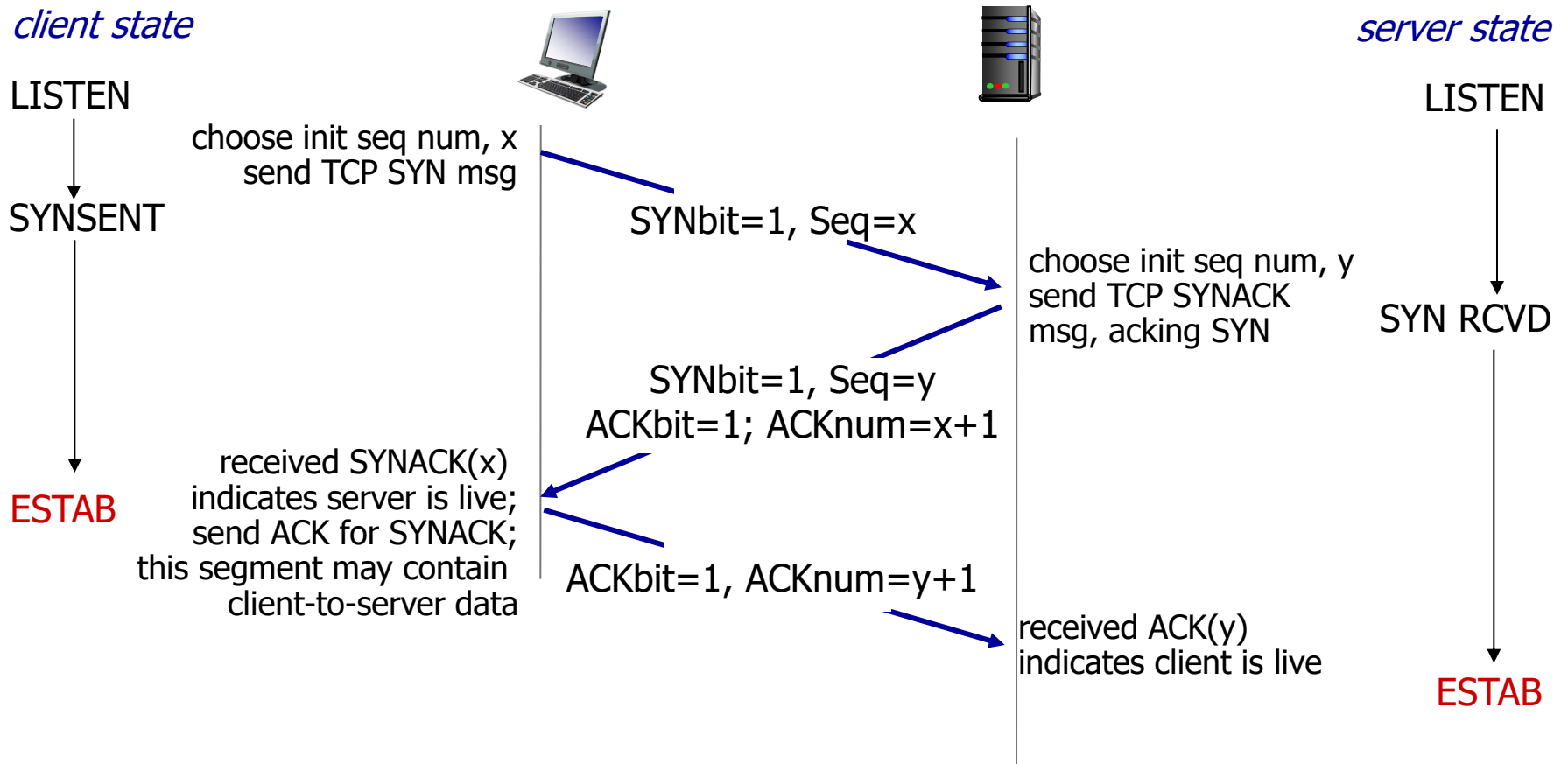<u>Step 1</u>: client host sends TCP SYN segment to server
  - ○ specifies initial seq #
  - ○ no data

<u>Step 2</u>: server host receives SYN, replies with SYNACK segment

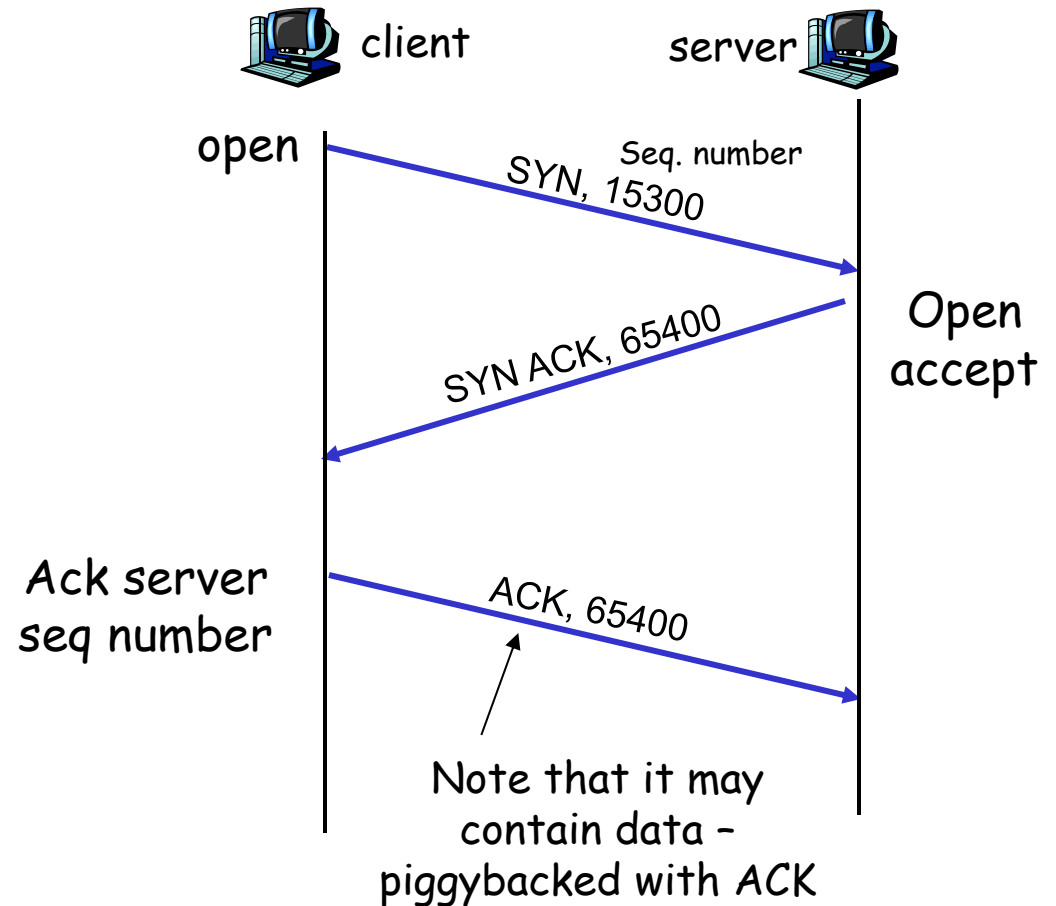  - ○ server allocates buffers
  - ○ specifies server initial seq #

<u>Step 3</u>: client receives SYNACK, replies with ACK segment, which may contain data
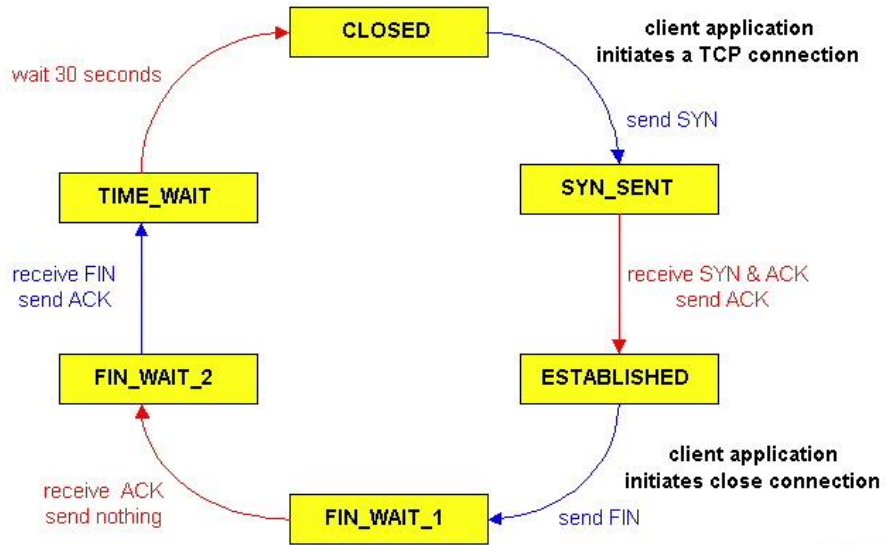
# TCP 3-way handshake

**client state**

LISTEN

SYNSENT

ESTAB

**server state**

LISTEN

SYN RCVD

ESTAB

choose init seq num, x
send TCP SYN msg

SYNbit=1, Seq=x

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

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

ACKbit=1, ACKnum=y+1

received ACK(y)
indicates client is live

# TCP Connection Management

□ 3 way handshake

client       server

open

SYN, 15300   Seq. number

Open
accept

SYN ACK, 65400

Ack server
seq number

ACK, 65400

Note that it may
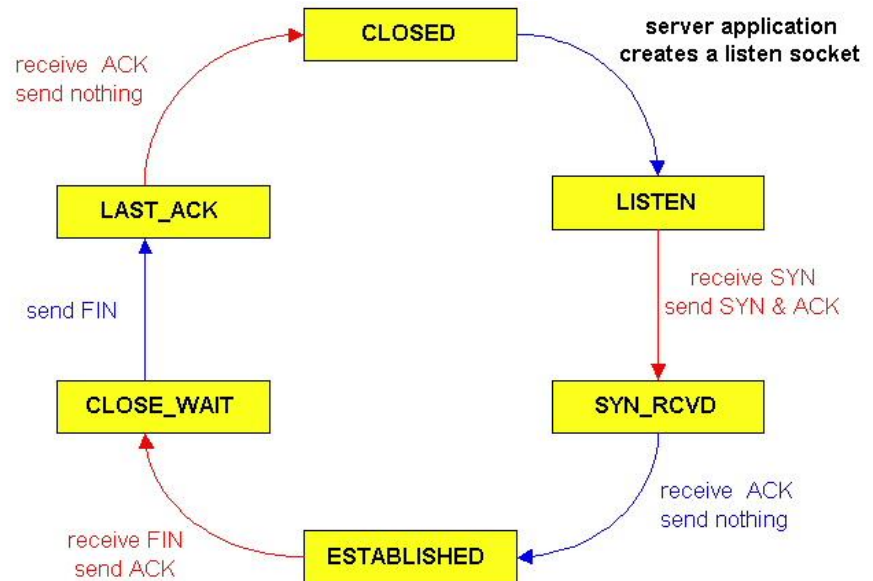contain data –
piggybacked with ACK

# TCP Connection Management (cont)



TCP client lifecycle

TCP server lifecycle

# TCP: closing a connection

- ❖ client, server each close their side of connection
  - ▪ send TCP segment with FIN bit = 1
- ❖ respond to received FIN with ACK
  - ▪ on receiving FIN, ACK can be combined with own FIN
- ❖ simultaneous FIN exchanges can be handled

# TCP: closing a connection

client state

ESTAB

clientSocket.close()

FIN_WAIT_1    can no longer
              send but can
              receive data

FIN_WAIT_2    wait for server
              close

TIMED_WAIT

              timed wait
              for 2*max
              segment lifetime

CLOSED

FINbit=1, seq=x

ACKbit=1; ACKnum=x+1

FINbit=1, seq=y

ACKbit=1; ACKnum=y+1

server state

ESTAB

CLOSE_WAIT    can still
              send data

LAST_ACK      can no longer
              send data

CLOSED

# TCP Connection Management (cont.)

## Closing a connection:

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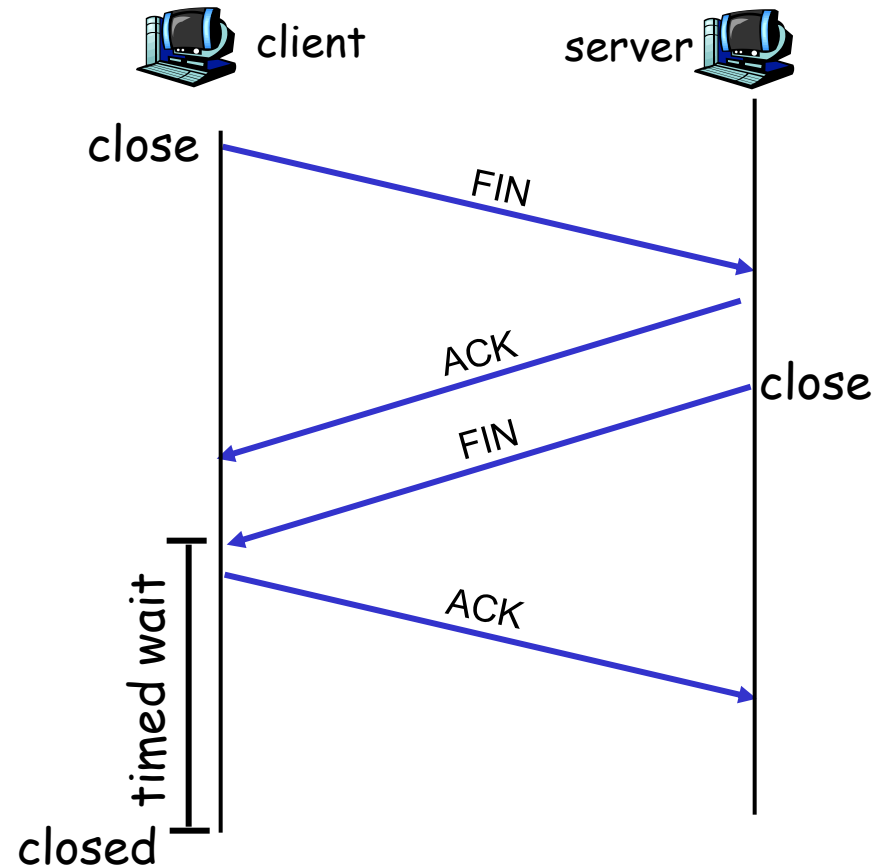```
clientSocket.close();
```

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

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

Step 3: client receives FIN, replies with ACK.

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

Step 4: server, receives ACK. Connection closed.
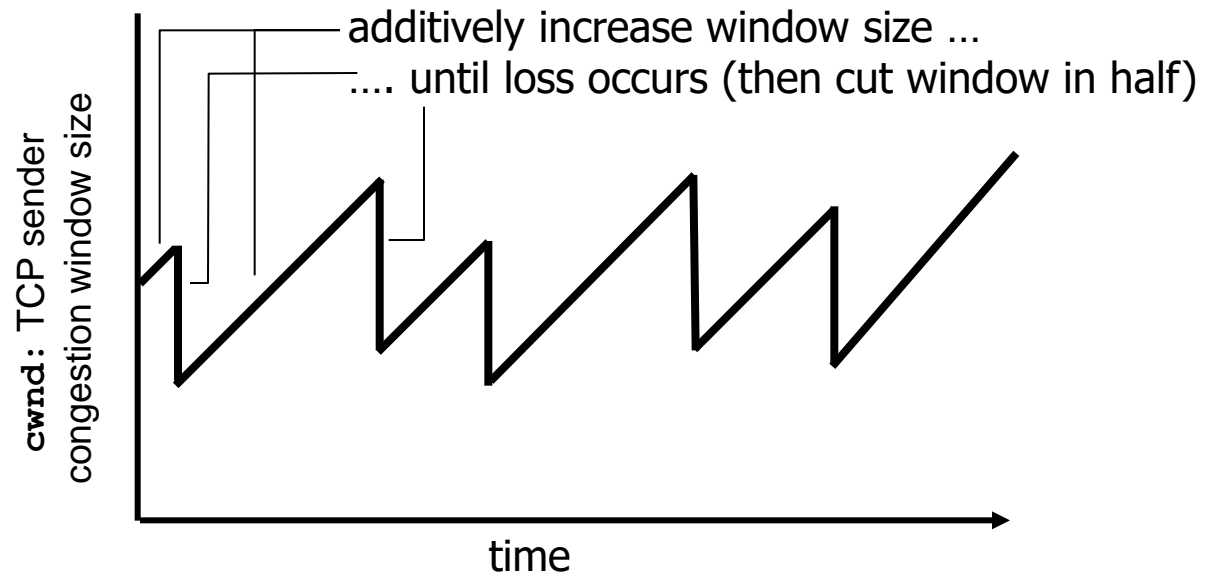
# Principles of congestion control

*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)
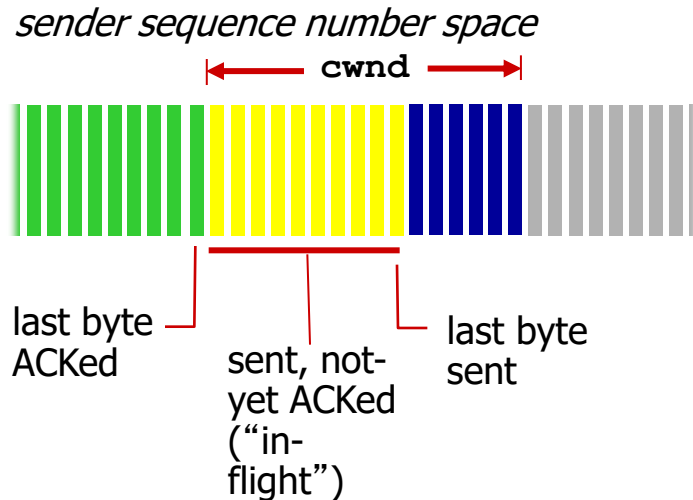  - ▪ long delays (queueing in router buffers)
- ❖ a top-10 problem!

# 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:

$$LastByteSent - LastByteAcked \leq cwnd$$

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

*TCP sending rate:*

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

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

# TCP Congestion Control

- TCP has a mechanism for congestion control. The mechanism is implemented at the sender

- The window size at the sender is set as follows:

  - **Send Window = MIN (flow control window, congestion window)**

where

- flow control window is advertised by the receiver
- congestion window is adjusted based on feedback from the network

# TCP Congestion Control

- The sender has two additional parameters:
  - **Congestion Window** (**cwnd**)
    Initial value is 1 MSS (=maximum segment size) counted as bytes
  - **Slow-start threshold Value** (**ssthresh)**
    Initial value is the advertised window size)

- Congestion control works in <u>two modes</u>:
  - **slow start** (cwnd < ssthresh)
  - **congestion avoidance** (cwnd >= ssthresh)
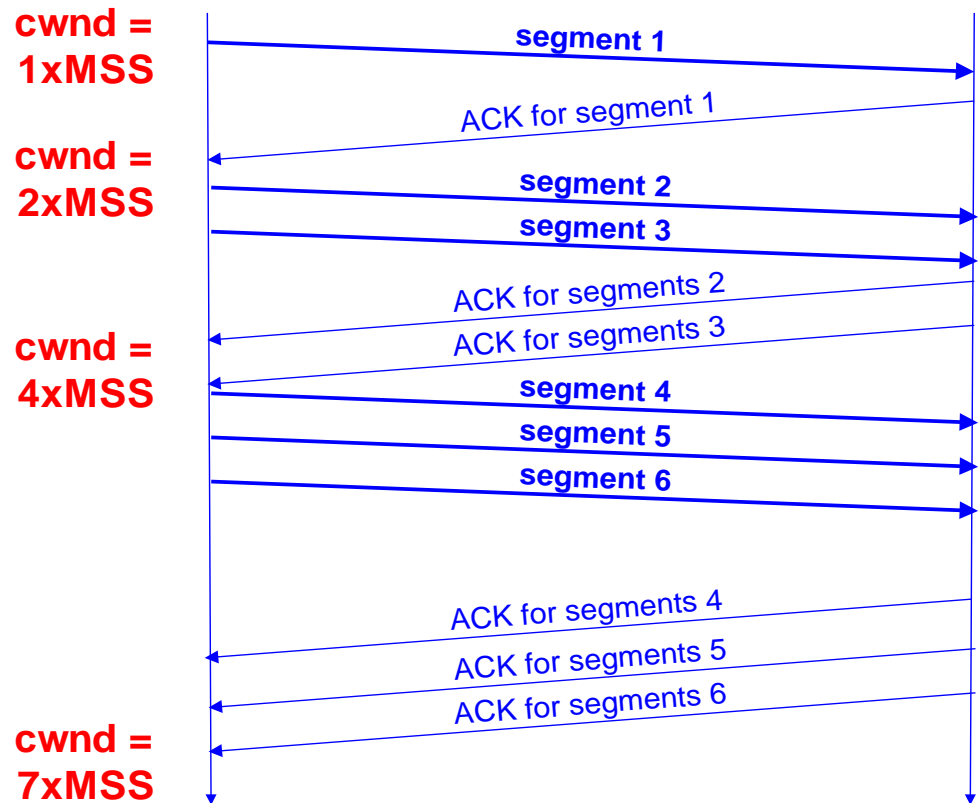
# Slow Start

- Initial value:

  - **cwnd = 1 segment**

- **Note: cwnd is actually measured in bytes:
  1 segment = MSS bytes**

- Each time an ACK is received, the congestion window is increased by MSS bytes.

  - **cwnd = cwnd + MSS**

  - If an ACK acknowledges two segments, cwnd is still increased by only 1 segment.

  - Even if ACK acknowledges a segment that is smaller than MSS bytes long, cwnd is increased by MSS.

- Does Slow Start increment slowly? Not really.
  In fact, the increase of cwnd can be exponential

# Slow Start Example

- The congestion window size grows very rapidly
  - For every ACK, we increase cwnd by 1 irrespective of the number of segments ACK'ed
- TCP slows down the increase of *cwnd* when **cwnd > ssthresh**

**cwnd = 1xMSS** → segment 1

ACK for segment 1 ←

**cwnd = 2xMSS** → segment 2, segment 3

ACK for segments 2 ←
ACK for segments 3 ←

**cwnd = 4xMSS** → segment 4, segment 5, segment 6

ACK for segments 4 ←
ACK for segments 5 ←
ACK for segments 6 ←

**cwnd = 7xMSS**

# Congestion Avoidance

- Congestion avoidance phase is started if cwnd has reached the slow-start threshold value

- If cwnd >= ssthresh then each time an ACK is received, increment cwnd  as follows:
  - cwnd = cwnd + MSS(MSS/ cwnd)

- So *cwnd* is increased by one segment (=MSS bytes) only if all segments have been acknowledged.

# Slow Start / Congestion Avoidance

- Here we give a more accurate version than in our earlier discussion of Slow Start:

**If** $cwnd < ssthresh$ **then**

      Each time an Ack is received:

      cwnd = cwnd $+ MSS$

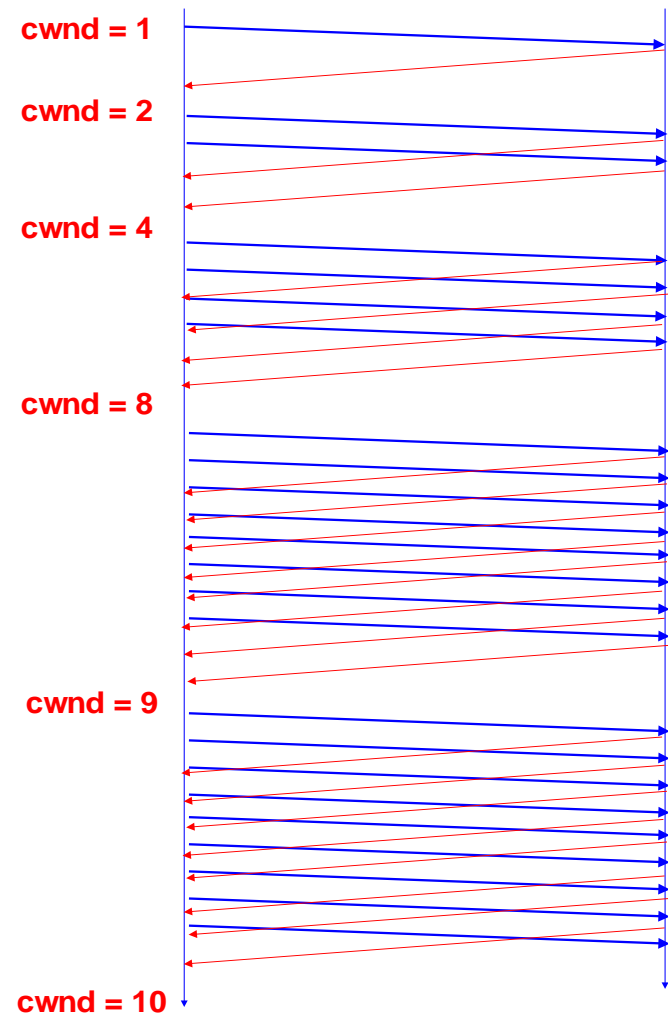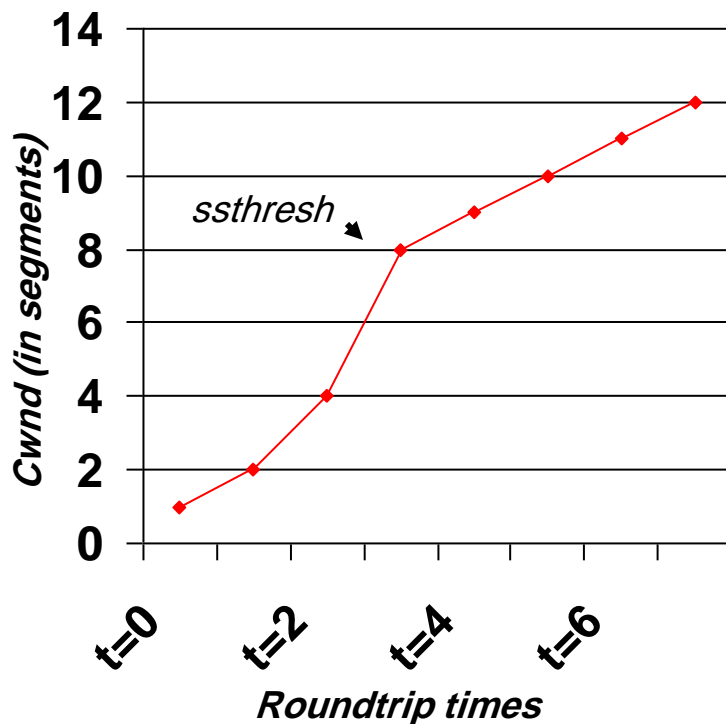**else**    **/\*** $cwnd >= ssthresh$ **\*/**

      Each time an Ack is received :

      $cwnd = cwnd + MSS. MSS / cwnd$

**endif**

# Example of
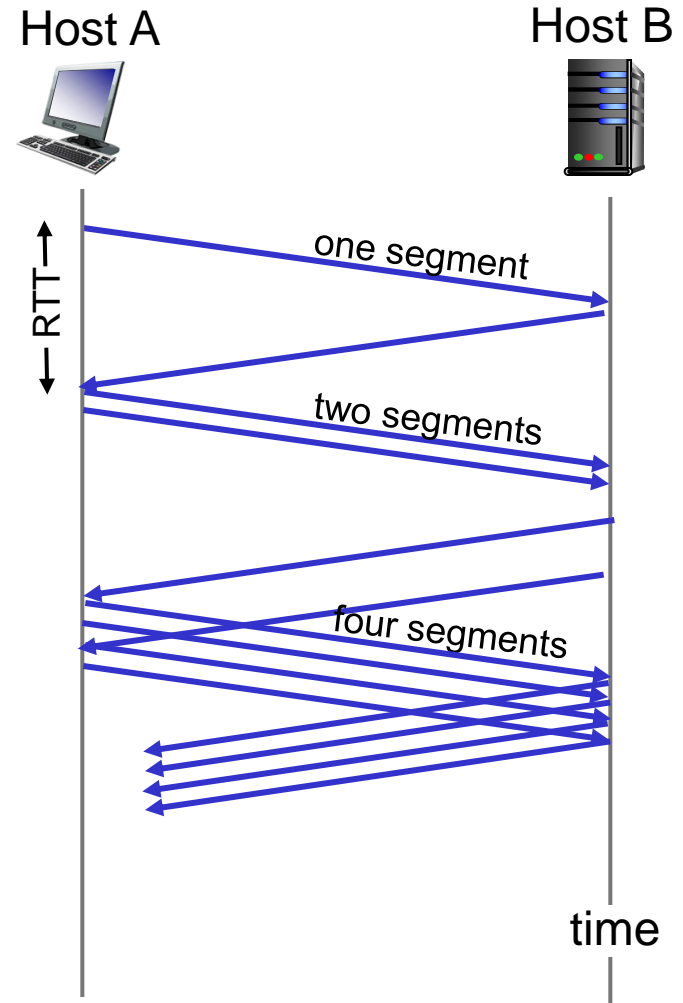# Slow Start/Congestion Avoidance

Assume that *ssthresh = 8*

# Responses to Congestion

- Most often, a packet loss in a network is due to an overflow at a congested router (rather than due to a transmission error)

- So, TCP assumes there is congestion if it detects a packet loss

- A TCP sender can detect lost packets via:
  - Timeout of a retransmission timer
  - Receipt of a duplicate ACK

- When TCP assumes that a packet loss is caused by congestion it reduces the size of the sending window

# 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

# TCP: detecting, reacting to loss

❖ loss indicated by timeout:
  - ▪ `cwnd` set to 1 MSS;
  - ▪ window then grows exponentially (as in slow start) to threshold, then grows linearly

❖ loss indicated by 3 duplicate ACKs: TCP RENO
  - ▪ dup ACKs indicate network capable of delivering some segments
  - ▪ `cwnd` is cut in half window then grows linearly

❖ TCP Tahoe always sets `cwnd` to 1 (timeout or 3 duplicate acks)

# TCP: switching from slow start to CA

Q: when should the exponential increase switch to linear?

A: when `cwnd` gets to 1/2 of its value before timeout.

## Implementation:

❖ variable `ssthresh`

❖ on loss event, `ssthresh` is set to 1/2 of `cwnd` just before loss event