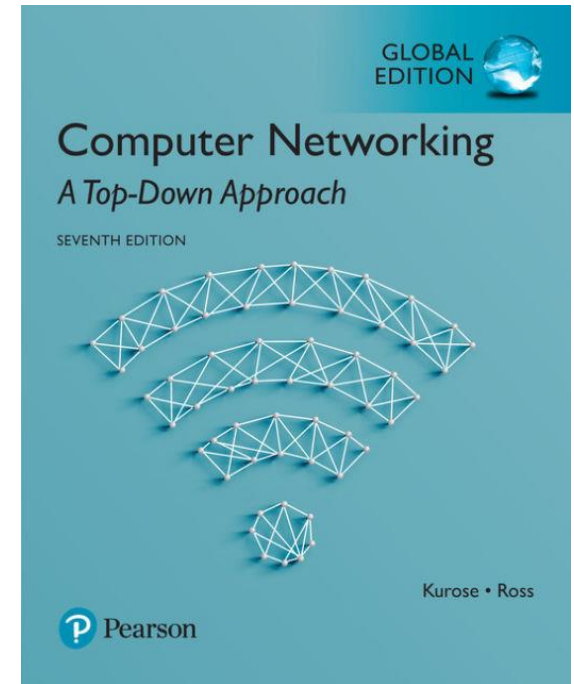# Transport Layer

*Computer Networking: A Top Down Approach*
7th edition
Jim Kurose, Keith Ross
Pearson, 2017

# Intended Learning Outcomes
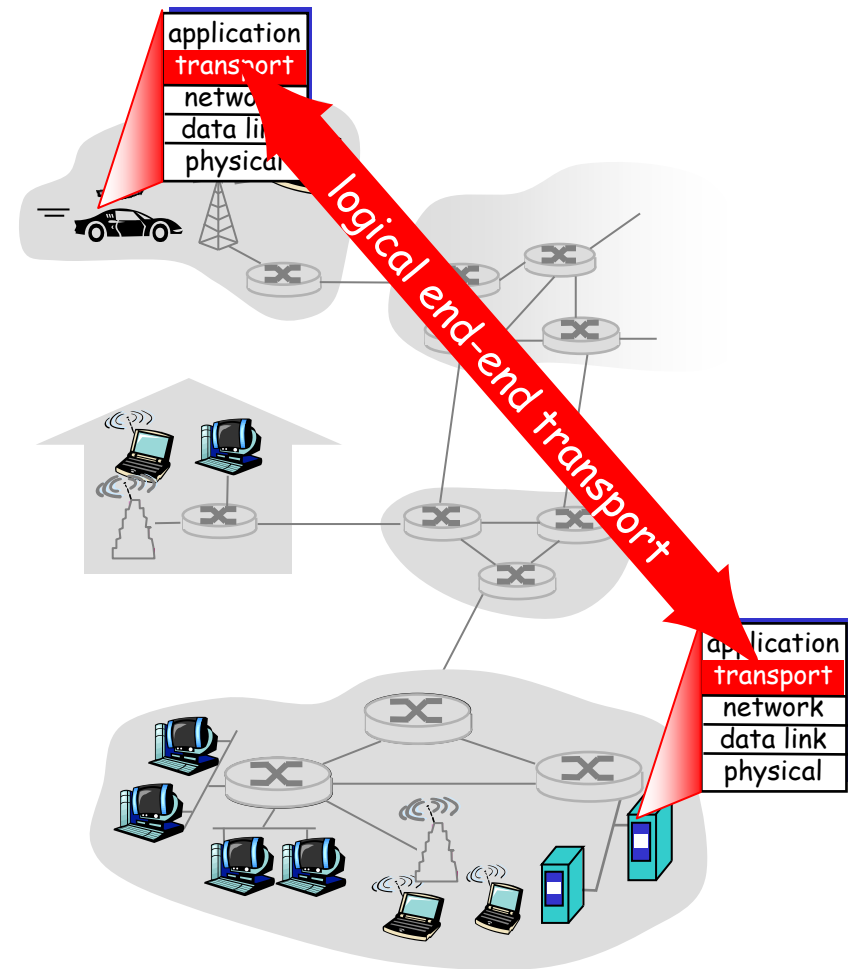
❑ understand principles behind transport layer services:

    ○ reliable data transfer

    ○ flow control

    ○ congestion control

❑ understand transport layer protocols in the Internet:

    ○ TCP: connection-oriented transport, connection management

    ○ TCP flow control

    ○ TCP congestion control

# Transport Layer

❑ <span style="color:red">Transport-layer services</span>
❑ Connectionless transport: UDP
❑ Connection-oriented transport: TCP
  o segment structure
  o reliable data transfer
  o flow control
  o connection management
❑ Principles of congestion control
❑ 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

# Transport Layer

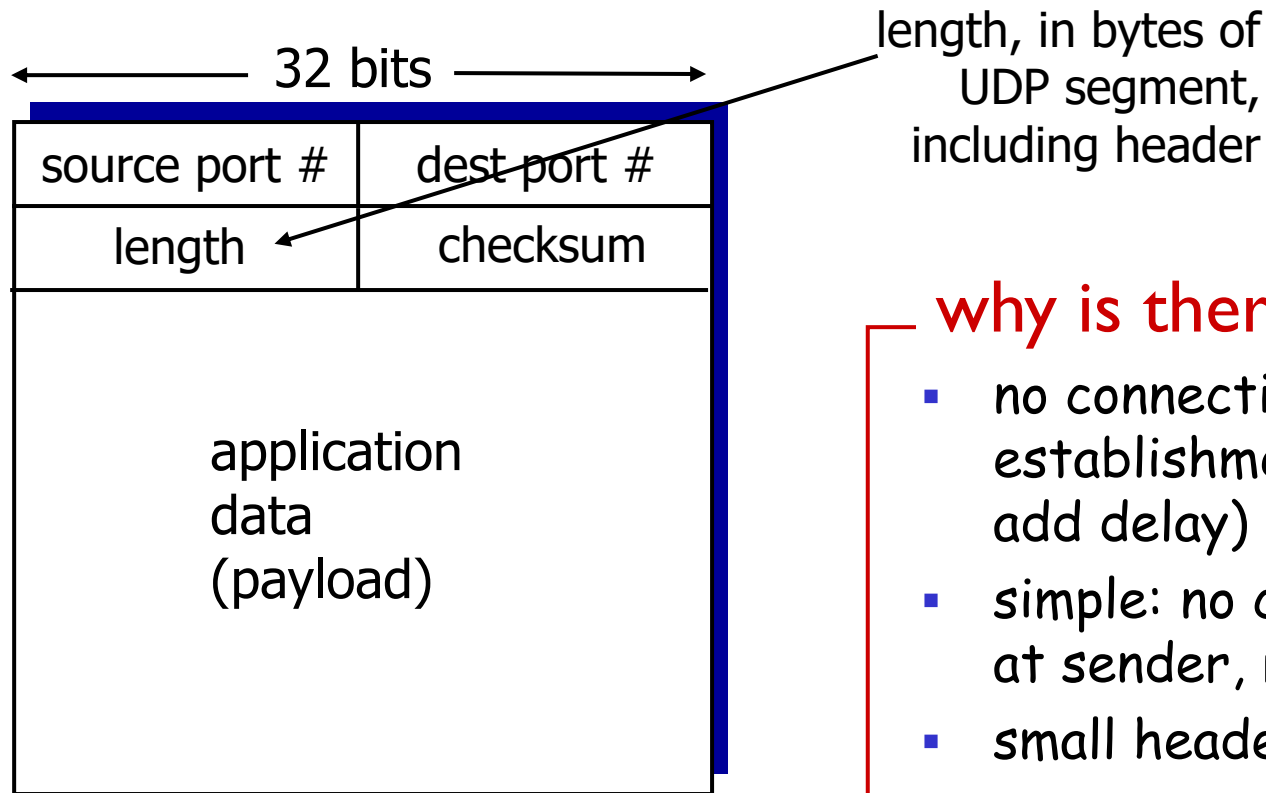❑ Transport-layer services
❑ Connectionless transport: UDP
❑ Connection-oriented transport: TCP
   o segment structure
   o reliable data transfer
   o flow control
   o connection management
❑ Principles of congestion control
❑ TCP congestion control

# UDP: User Datagram Protocol [RFC 768]

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

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

# UDP: segment header

32 bits

length, in bytes of
UDP segment,
including header

| source port # | dest port # |
|---------------|-------------|
| length | checksum |
| application data (payload) | |

UDP segment format

## 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 send out data as fast as desired
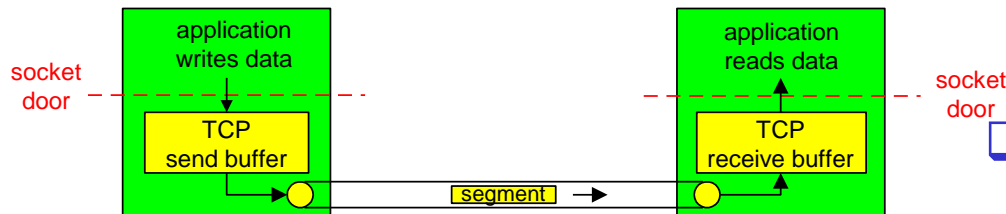
# Transport Layer

- ❑ Transport-layer services
- ❑ Connectionless transport: UDP
- ❑ Connection-oriented transport: TCP
  - o segment structure
  - o reliable data transfer
  - o flow control
  - o connection management
- ❑ Principles of congestion control
- ❑ TCP congestion control

# TCP: Review

RFCs: 793, 1122, 1323, 2018, 2581

❑ **point-to-point:**
  o one sender, one receiver

❑ **reliable, in-order *byte stream:***
  o no "message boundaries"

❑ **pipelined:**
  o TCP congestion and flow control set window size

❑ ***send & receive buffers***
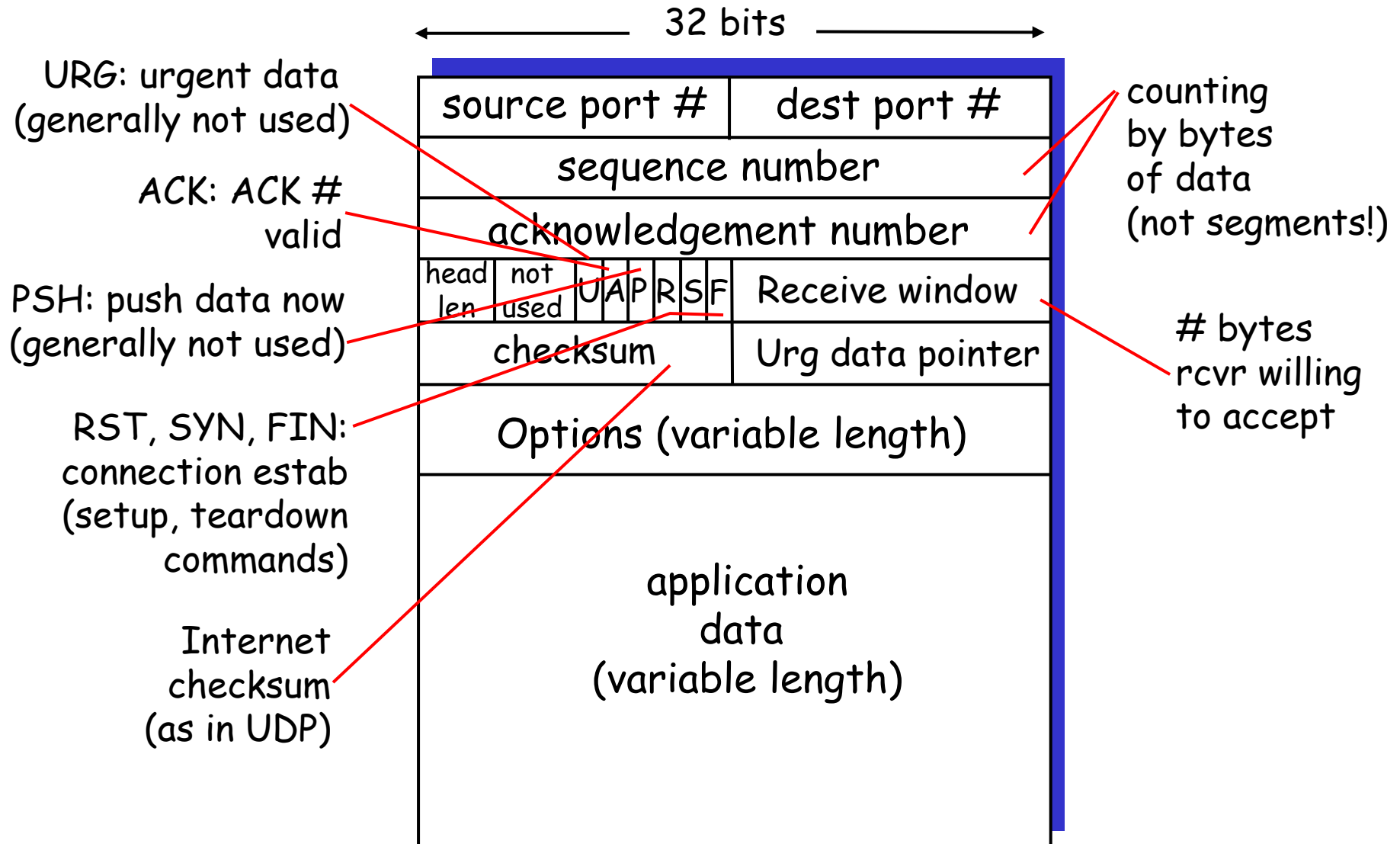


❑ **full duplex data:**
  o bi-directional data flow in same connection
  o MSS: maximum segment size

❑ **connection-oriented:**
  o handshaking (exchange of control msgs) init's sender & receiver states before data exchange

❑ **flow controlled:**
  o sender will not overload receiver

❑ **congestion controlled:**
  o sender will not overload network

Transport Layer     9

# Transport Layer

❑ Transport-layer services
❑ Connectionless transport: UDP
❑ Connection-oriented transport: TCP
  o <span style="color:red">segment structure</span>
  o reliable data transfer
  o flow control
  o connection management
❑ Principles of congestion control
❑ TCP congestion control

# 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

# Transport Layer

- ❑ Transport-layer services
- ❑ Connectionless transport: UDP
- ❑ Connection-oriented transport: TCP
  - o segment structure
  - o reliable data transfer
  - o flow control
  - o connection management
- ❑ Principles of congestion control
- ❑ TCP congestion control

# TCP seq. numbers, ACKs

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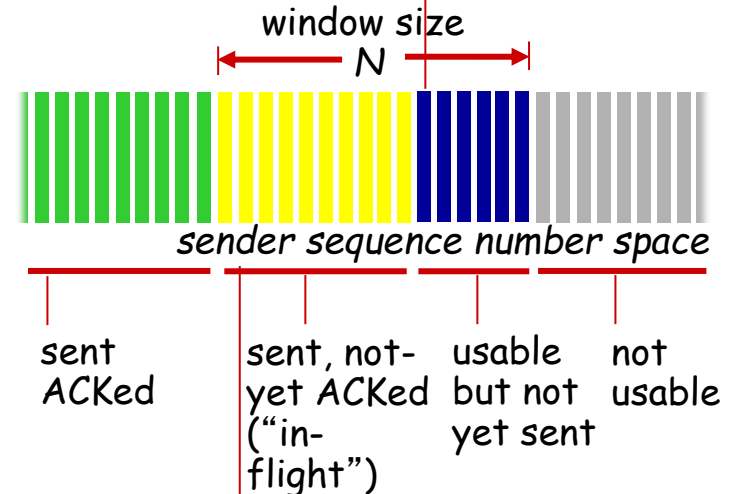- o byte stream "number" of first byte in segment's data

acknowledgement number:
- o seq # of next byte expected from other side
- o cumulative ACK

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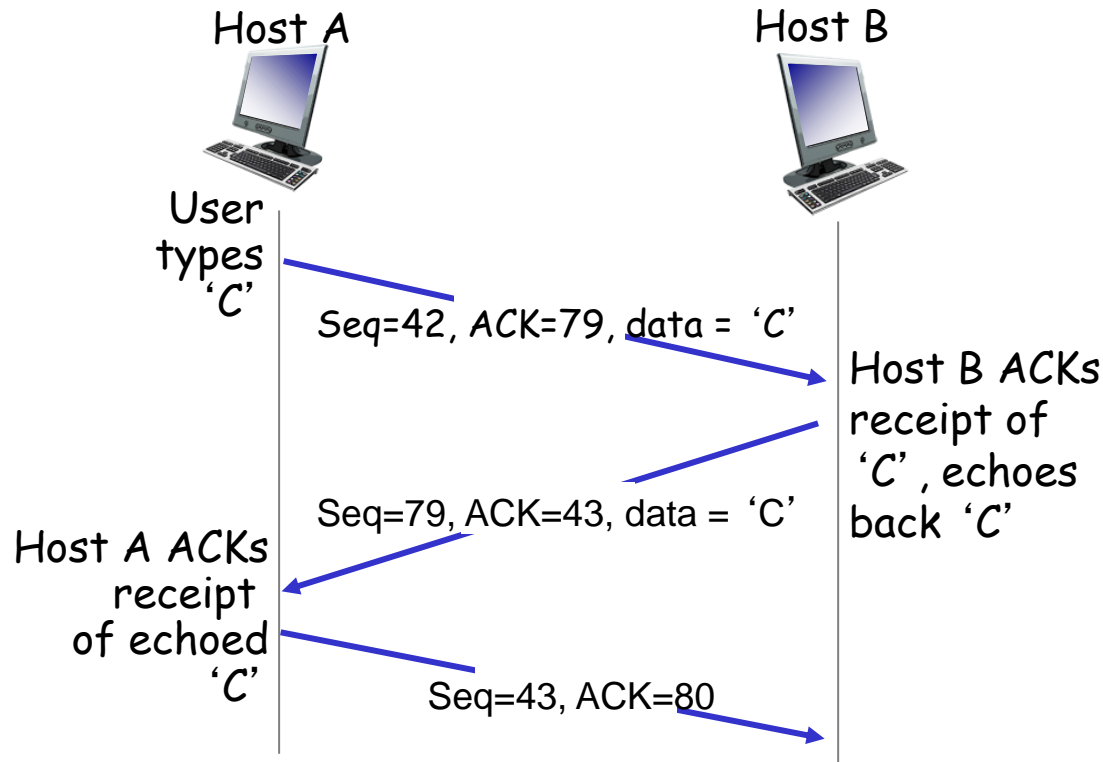- o Ans: TCP spec doesn't say, - up to implementor

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, ACK numbers

Host A                          Host B

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

Host B ACKs
receipt of
'C', echoes
back 'C'

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

Host A ACKs
receipt
of echoed
'C'

Seq=43, ACK=80

simple telnet scenario

# TCP reliable data transfer (rdt)

❖ 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
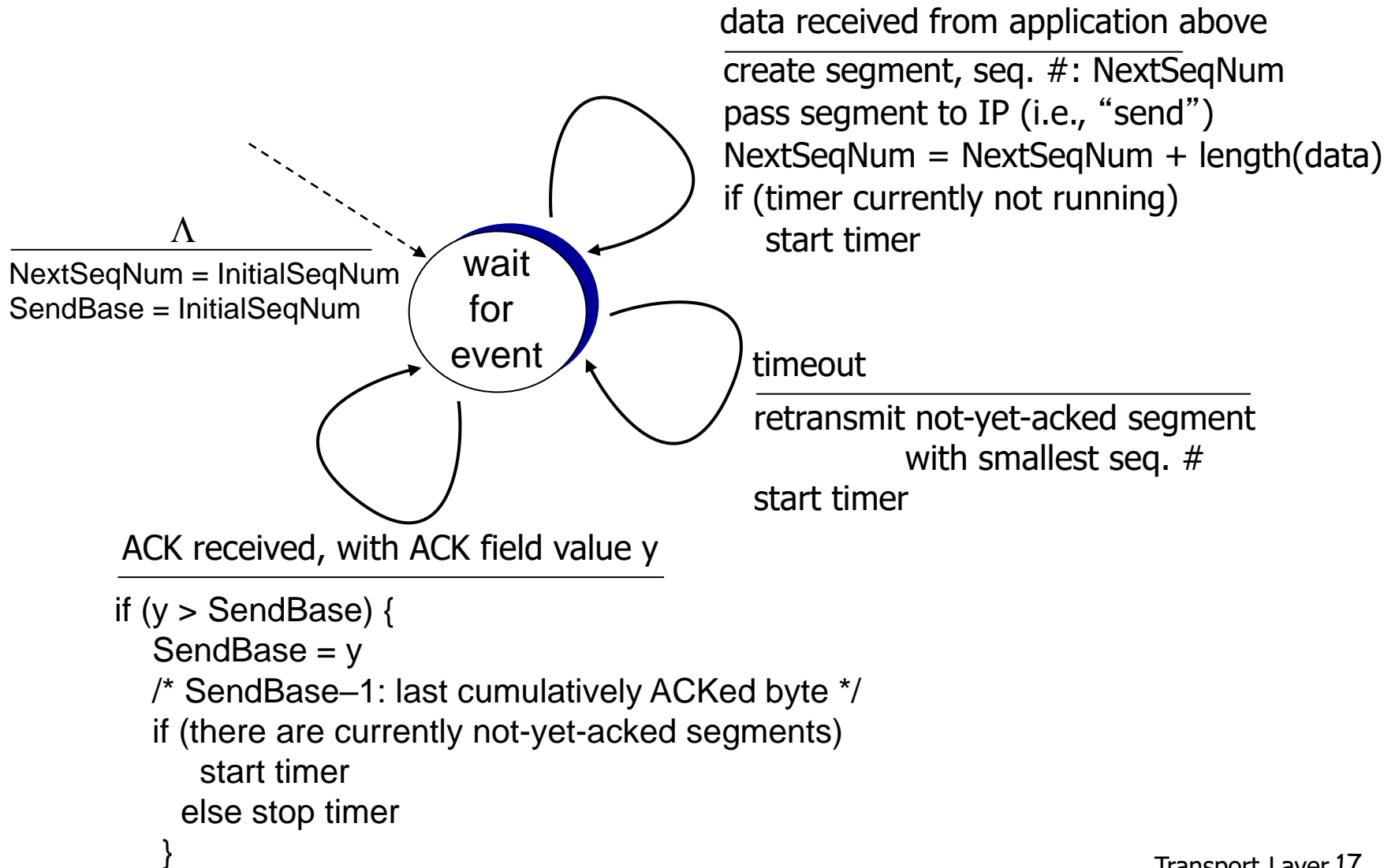  ▪ expiration interval: `TimeOutInterval`

*timeout:*

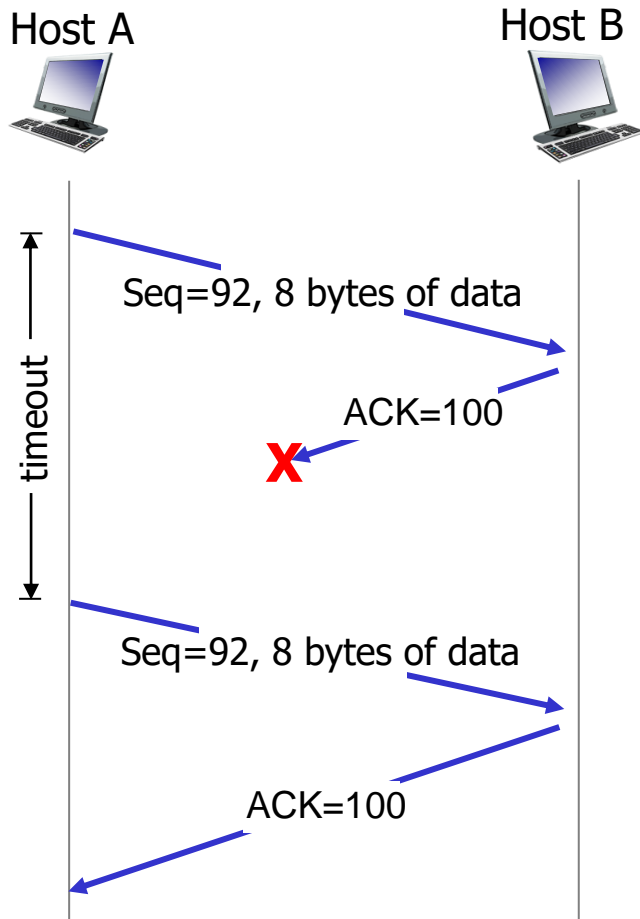❖ retransmit segment that caused timeout

❖ restart timer

*ack rcvd:*

❖ if ack acknowledges previously unacked segments
  ▪ update what is known to be ACKed
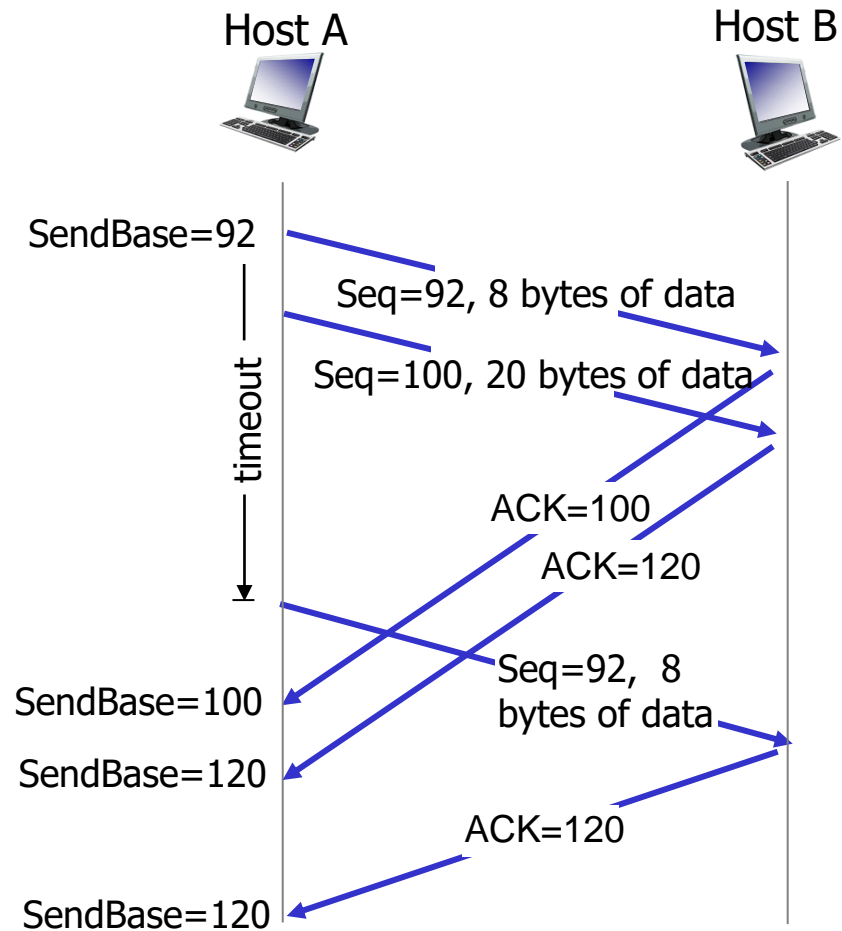  ▪ start timer if there are still unacked segments

# TCP sender (simplified)

data received from application above
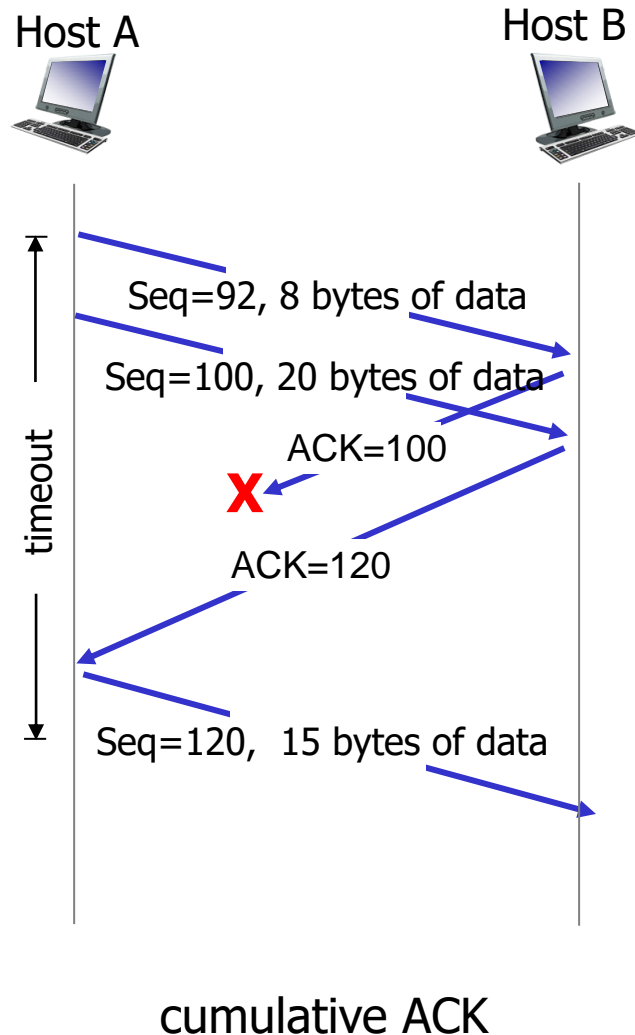
create segment, seq. #: NextSeqNum
pass segment to IP (i.e., "send")
NextSeqNum = NextSeqNum + length(data)
if (timer currently not running)
    start timer

$\Lambda$

NextSeqNum = InitialSeqNum
SendBase = InitialSeqNum

wait
for
event

timeout

retransmit not-yet-acked segment
            with smallest seq. #
start timer

ACK received, with ACK field value y

if (y > SendBase) {
    SendBase = y
    /* SendBase–1: last cumulatively ACKed byte */
    if (there are currently not-yet-acked segments)
        start timer
      else stop timer
    }

# TCP: retransmission scenarios

Host A                                          Host B

Seq=92, 8 bytes of data

timeout

ACK=100

**X**

Seq=92, 8 bytes of data

ACK=100

lost ACK scenario

Host A                                          Host B

SendBase=92

Seq=92, 8 bytes of data

timeout

Seq=100, 20 bytes of data

ACK=100

ACK=120

SendBase=100

Seq=92, 8 bytes of data

SendBase=120

ACK=120

SendBase=120

premature timeout

# TCP: retransmission scenarios: Poll 6



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 ACK generation [RFC 1122, RFC 2581]

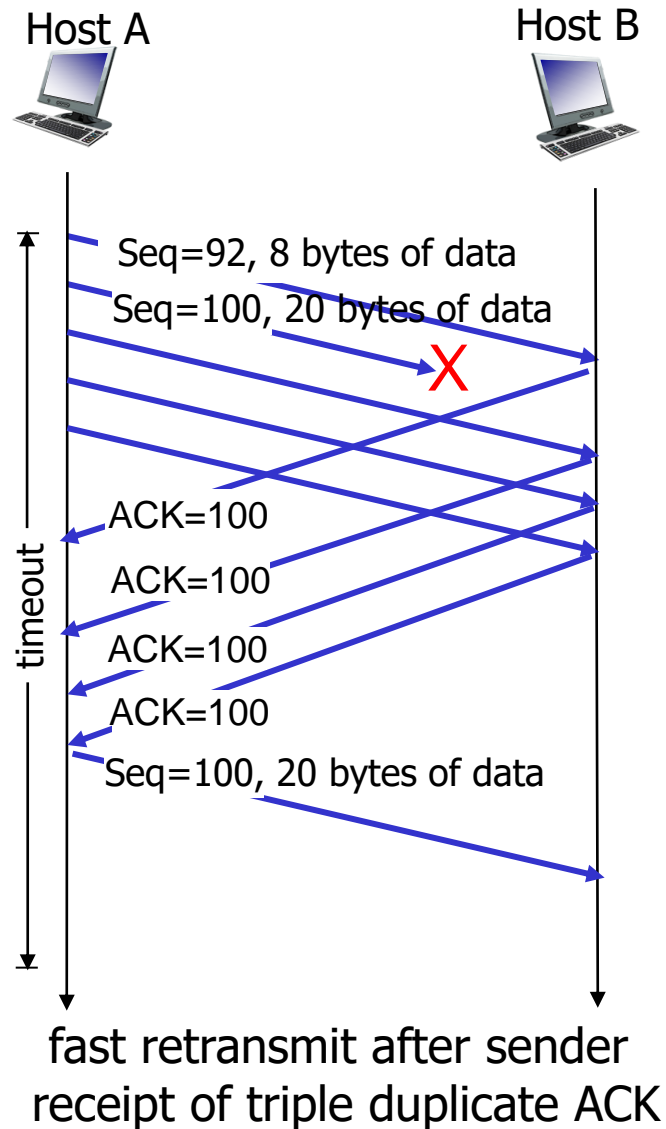| *event at receiver* | *TCP receiver action* |
|---|---|
| arrival of in-order segment with expected seq #. All data up to expected seq # already ACKed | delayed ACK. Wait up to 500ms for next segment. If no next segment, send ACK |
| arrival of in-order segment with expected seq #. One other segment 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 |

# TCP fast retransmit

❖ **time-out period often relatively long:**
  - long delay before resending lost packet
❖ **detect lost segments via duplicate ACKs.**
  - sender often sends many segments back-to-back
  - if segment is lost, there will likely be many duplicate ACKs.

*TCP fast retransmit*

if sender receives 4 ACKs for same data ("triple duplicate ACKs"), resend unacked segment with smallest seq #
  - likely that unacked segment lost, so don't wait for timeout

# TCP fast retransmit

Host A                                        Host B

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

X

ACK=100

ACK=100

timeout

ACK=100

ACK=100

Seq=100, 20 bytes of data

fast retransmit after sender
receipt of triple duplicate ACK

# Round trip time (RTT) vs timeout

Q: how to set TCP timeout value?

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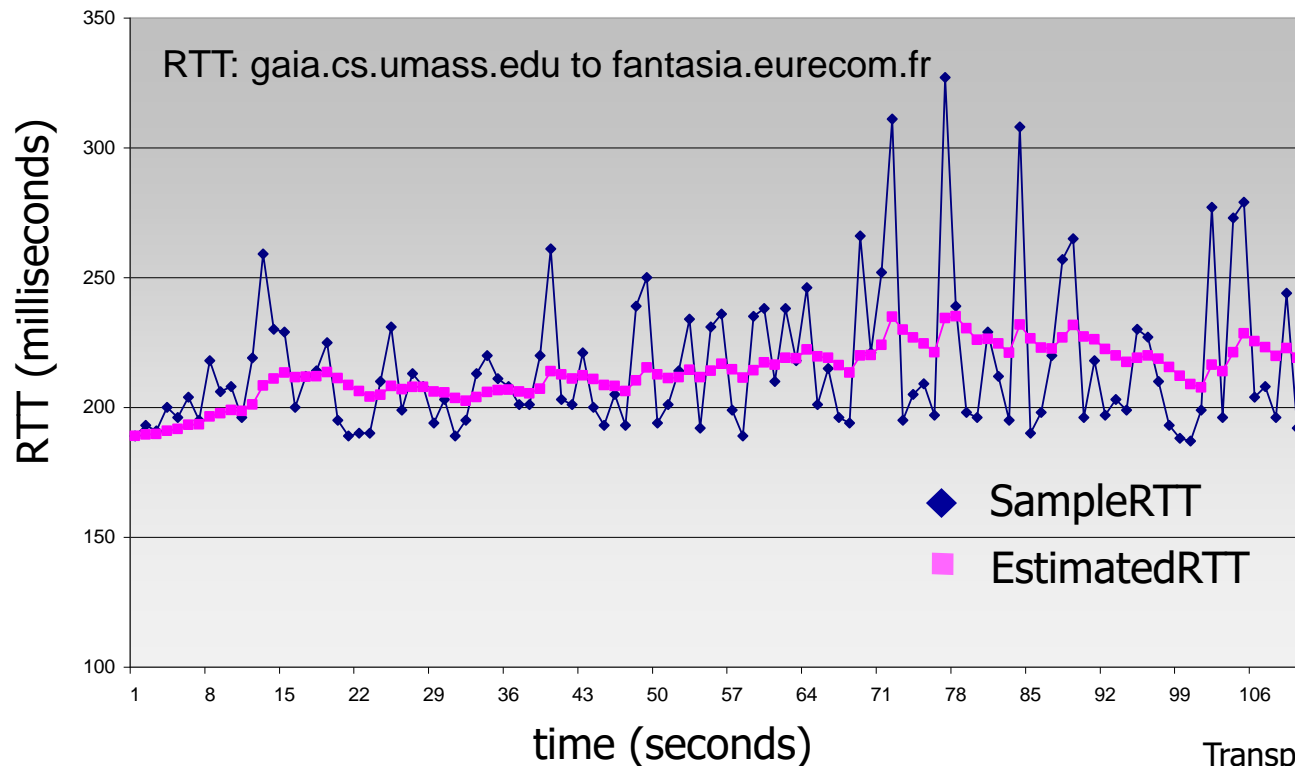❖ *too long:* slow reaction to segment loss

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**

# Round trip time (RTT) vs timeout

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

- ❖ 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

time (seconds)

# Round trip time (RTT) vs timeout

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

❖ estimate SampleRTT deviation from EstimatedRTT:

$$DevRTT_{i+1} = (1-\beta)*DevRTT_i + $$
$$\beta*|SampleRTT_i - EstimatedRTT_i|$$
$$(typically, \beta = 0.25)$$

$$TimeoutInterval_i = EstimatedRTT_i + 4*DevRTT_i$$

estimated RTT          "safety margin"

# Summary for TCP reliable data transfer

❑ **How to make reliable data?**
  - o Sequence number, retransmission timer, cumulative ACK

❑ **How to shorten retransmission delay?**
  - o "Fast Retransmit": lost segments detection via duplicate ACKs

❑ **How to set time-out value?**
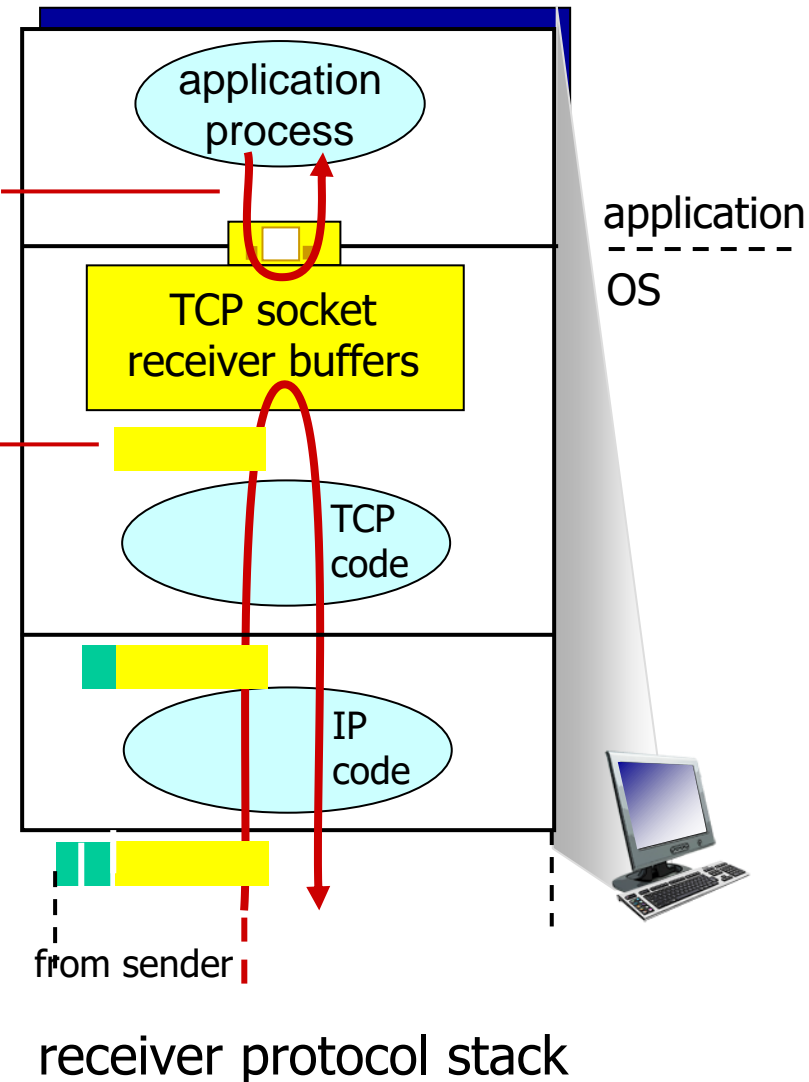  - o Exponential weighted moving average

# Transport Layer

❑ Transport-layer services
❑ Connectionless transport: UDP
❑ Connection-oriented transport: TCP
  o segment structure
  o reliable data transfer
  o <span style="color:red">flow control</span>
  o connection management
❑ Principles of congestion control
❑ TCP congestion control

# TCP flow control

application may
remove data from
TCP socket buffers ….

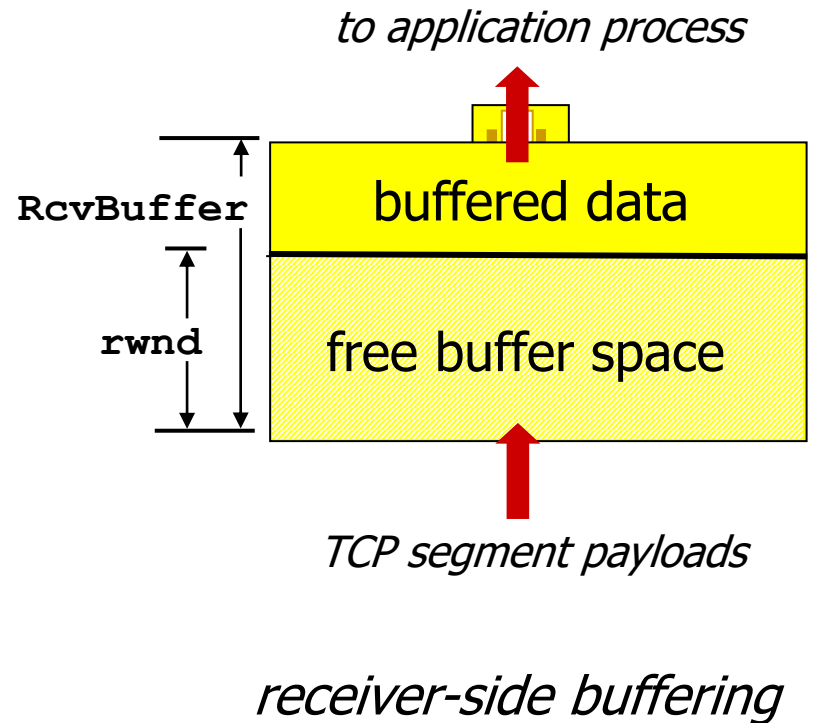application
process

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
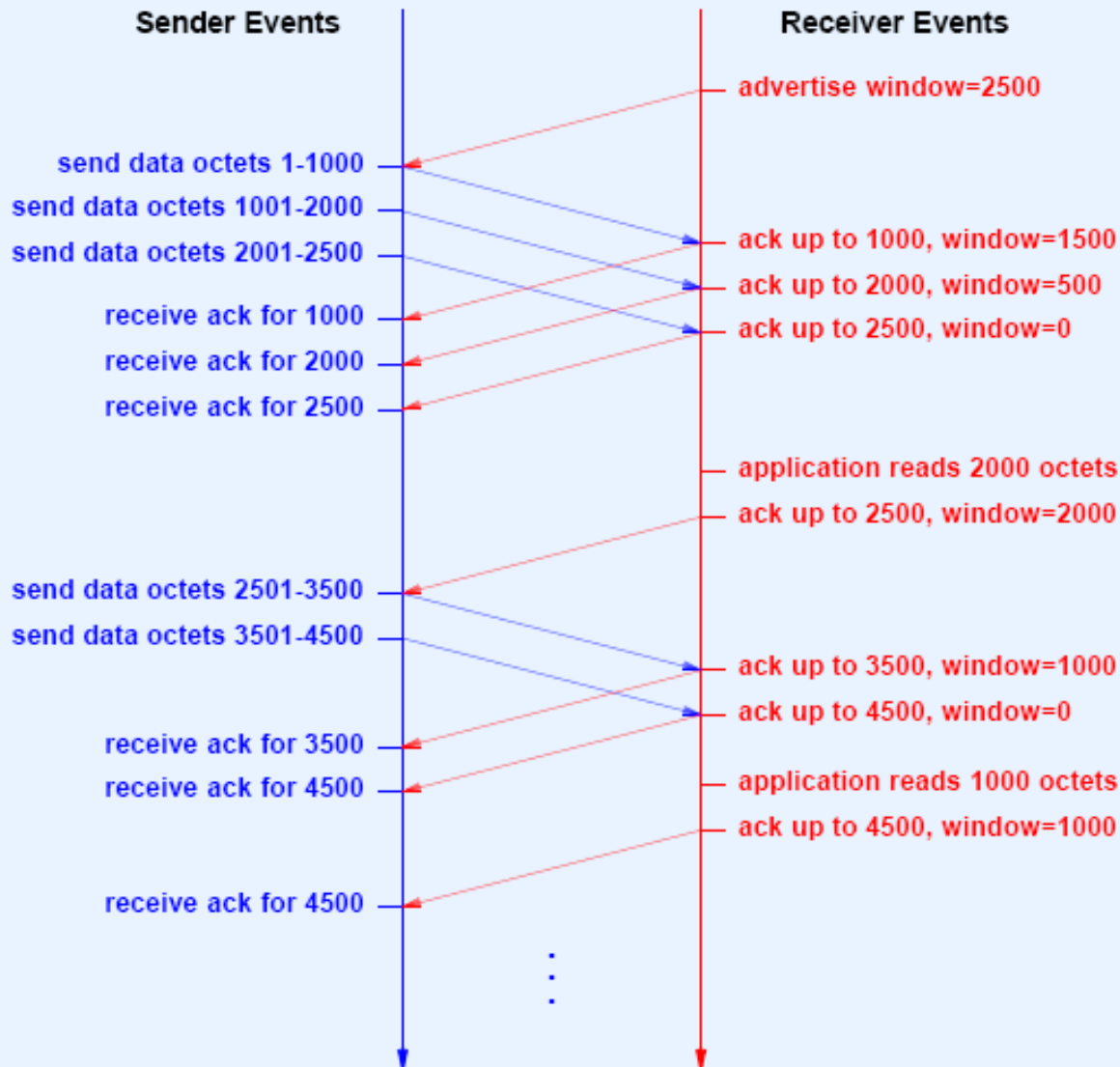
receiver protocol stack

# TCP flow control (cont'ed)

❖ 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 receiver buffer will not overflow

*to application process*

RcvBuffer

buffered data

rwnd

free buffer space

*TCP segment payloads*

*receiver-side buffering*

# Window Advertisement

❑ Each ACK carries new window information:
   o Acknowledgement number (*AN*)
   o Window size (*W*)

❑ ACK contains *AN* = *i*, *W* = *j*:
   o Bytes through *SN* = *i* − 1 acknowledged
      • Cumulative ACK
      • Byte *i* has not been received (It is the next byte expected)
   o Permission is granted to send *W* = *j* more bytes
      • i.e. bytes *i* through *i* + *j* − 1
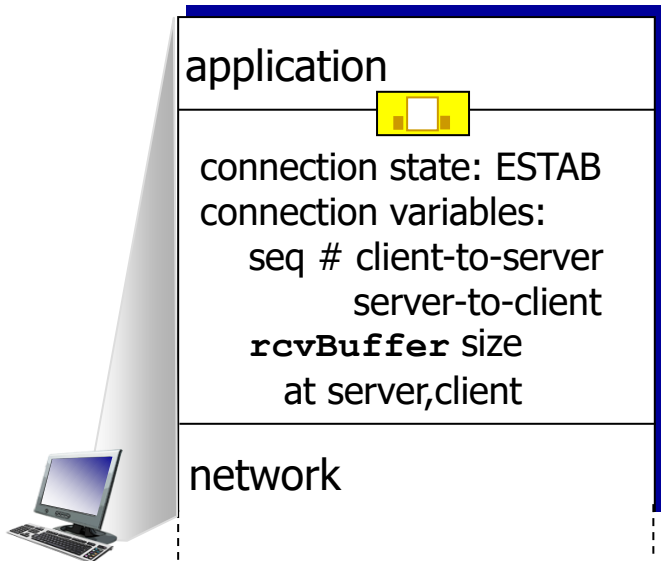
# Illustration: Window Advertisement

# Transport Layer

❑ Transport-layer services
❑ Connectionless transport: UDP
❑ Connection-oriented transport: TCP
  o segment structure
  o reliable data transfer
  o flow control
  o connection management
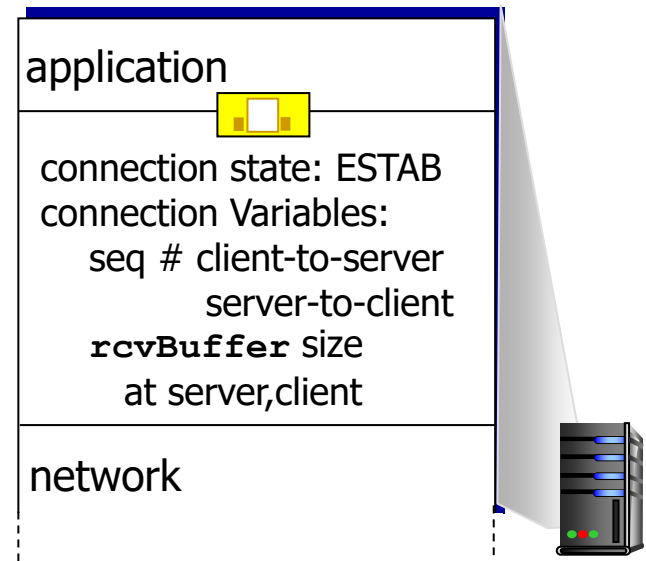❑ Principles of congestion control
❑ TCP congestion control

# Connection Management

before exchanging data, sender/receiver "handshake":

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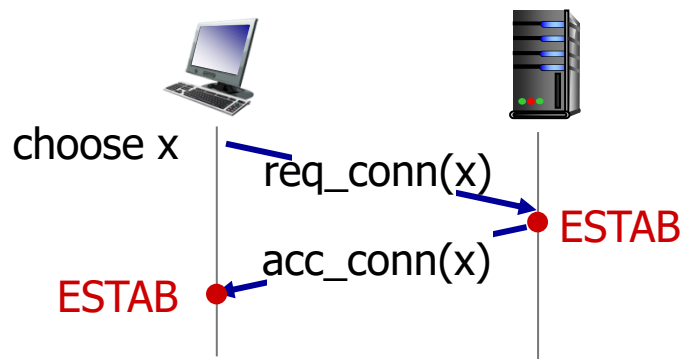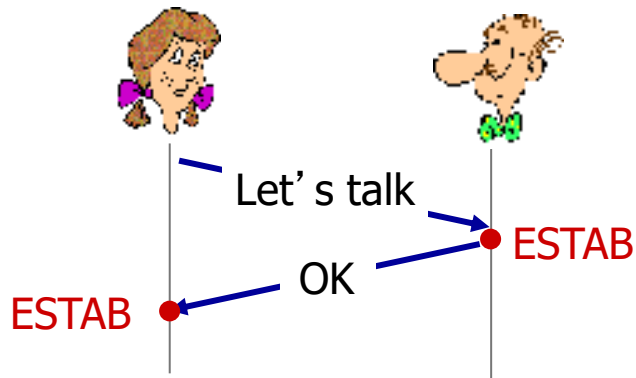
❖ agree on connection parameters

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:
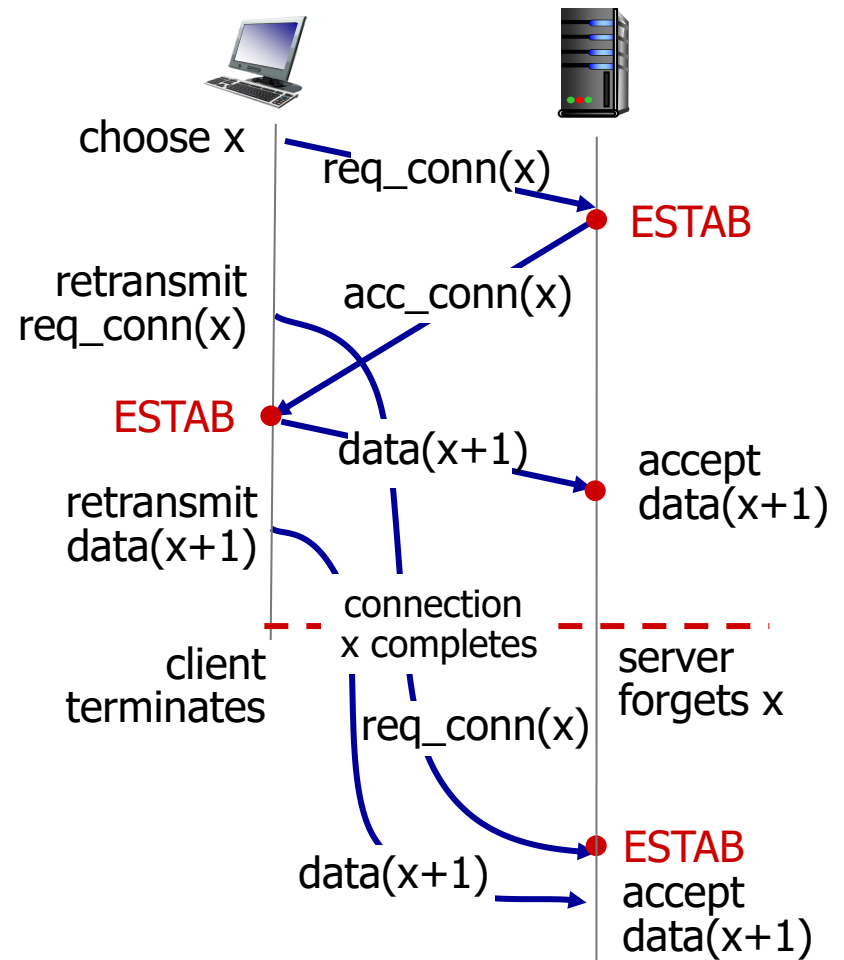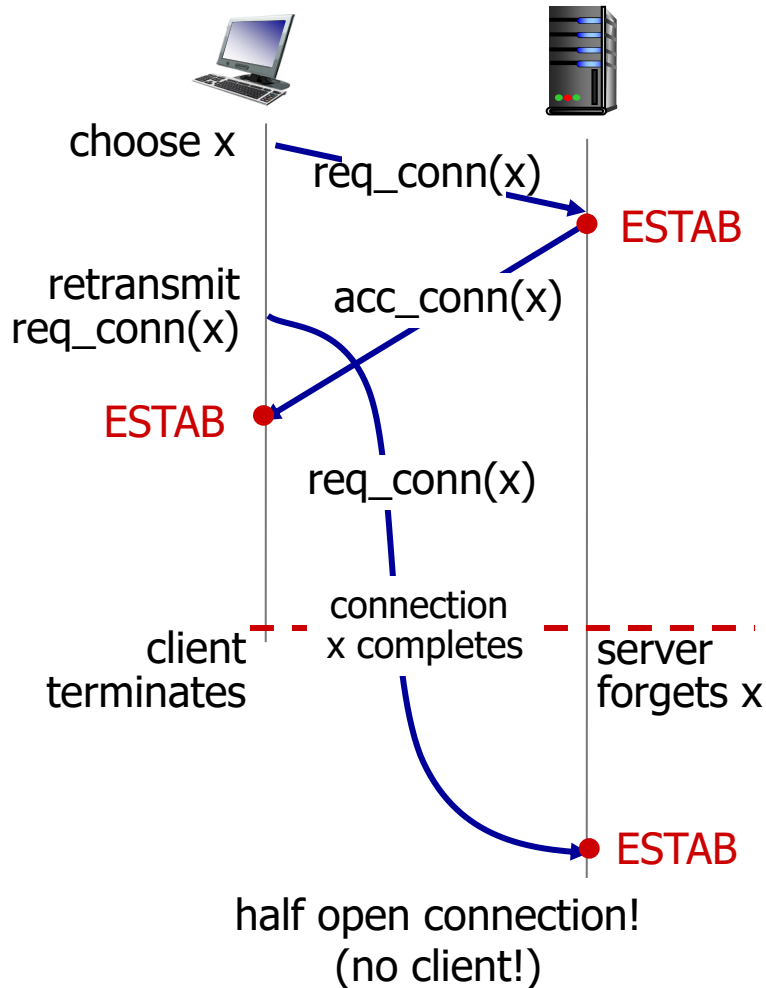


*Q:* will 2-way handshake always work in network?

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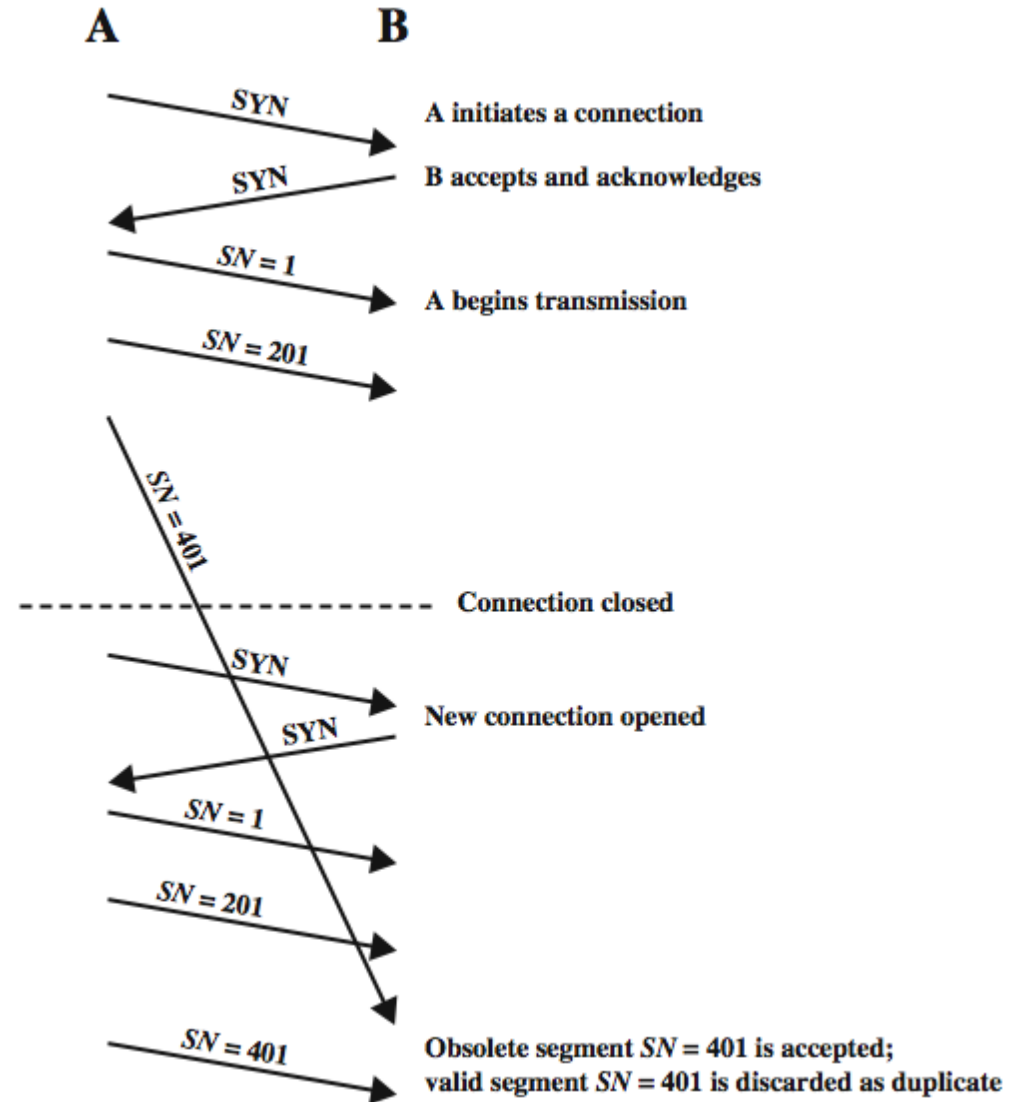❖ can't "see" other side

# Problems with Two-way Handshake

❑ In an *unreliable* network (e.g. the Internet), lost or delayed segments can cause problems in connection establishment, data transfer and connection termination

# Agreeing to establish a connection

2-way handshake failure scenarios:

# Two Way Handshake: Obsolete Data Segment



A initiates a connection — SYN

B accepts and acknowledges — SYN

SN = 1 — A begins transmission

SN = 201

SN = 401

Connection closed

SYN — New connection opened

SYN

SN = 1

SN = 201

SN = 401 — Obsolete segment SN = 401 is accepted; valid segment SN = 401 is discarded as duplicate

# Two Way Handshake: Obsolete SYN Segment



A           B

SYN $i$

- - - - - - - - - - - - - - - - - - - -    Connection closed

Obsolete SYN $i$ arrives

SYN $k$      SYN $j$      B responds; A sends new SYN

B discards duplicate SYN

$SN = k+1$

B rejects segment as out of sequence

# TCP 3-way handshake

client state

LISTEN

SYNSENT

ESTAB

server state

LISTEN

choose init seq num, x
send TCP SYN msg

SYNbit=1, Seq=x

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

SYN RCVD

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

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

ACKbit=1, ACKnum=y+1

received ACK(y)
indicates client is live

ESTAB

# TCP 3-way handshake: FSM

closed

$\dfrac{\text{Socket connectionSocket =}}{\text{welcomeSocket.accept();}}$

$\Lambda$

$\dfrac{\text{Socket clientSocket =}}{\text{newSocket("hostname","port number");}}$

SYN(seq=x)

$\dfrac{\text{SYN(x)}}{\begin{array}{c}\text{SYNACK(seq=y,ACKnum=x+1)}\\\text{create new socket for}\\\text{communication back to client}\end{array}}$

listen

SYN rcvd

SYN sent

ESTAB

$\dfrac{\text{ACK(ACKnum=y+1)}}{\Lambda}$

$\dfrac{\text{SYNACK(seq=y,ACKnum=x+1)}}{\text{ACK(ACKnum=y+1)}}$

# Three Way Handshake: Examples



A initiates a connection — SYN $i$

B accepts and acknowledges — SYN $j$, $AN = i + 1$

A acknowledges and begins transmission — $SN = i + 1$, $AN = j + 1$

(a) Normal operation

Obsolete SYN arrives — SYN $i$

B accepts and acknowledges — SYN $j$, $AN = i + 1$

A rejects B's connection — RST, $AN = j$

(b) Delayed SYN

SYN $i$    SYN $k$, $AN = p$

A initiates a connection
Old SYN arrives at A; A rejects

B accepts and acknowledges — RST, $AN = k$

SYN $j$, $AN = i + 1$

A acknowledges and begins transmission — $SN$ $i + 1$, $AN = j + 1$

(c) Delayed SYN, ACK

# Closing a connection

## Question:

Do we have a perfect solution for synchronizing the disconnection on both end systems if data can be lost in the network?

See: the two-army problem or
the Romeo and Juliet problem

# Closing a connection (cont'd)

## Answer:

No, but "three-way handshake" is an acceptable solution.

# TCP: closing a Connection

- For better understanding, think of a TCP connection as a pair of simplex connections.
  - "Simplex" means uni-directional data flow
  - Note: A TCP connection is full duplex (i.e., bi-directional.)
- Each simplex connection is released independently using these two steps:
  - Send a TCP segment with the FIN bit set to one.
  - When the FIN is acknowledged, that direction is shut down.
- Timers are used for graceful disconnection to avoid the two-army problem.
  - Not a perfect solution, i.e. graceful disconnection cannot be guaranteed
  - In fact, there is no perfect solution at all!

# 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

# Summary for connection management

❑ Problem: In an *unreliable* network (e.g. the Internet), lost or delayed segments can cause problems in connection establishment, data transfer and connection termination.

❑ Acceptable Solution: three way handshake

❑ Three way handshake is much better than two way handshake.

❑ Timers are used for graceful disconnection to avoid the two-army problem.

# Transport Layer

❑ Transport-layer services
❑ Connectionless transport: UDP
❑ Connection-oriented transport: TCP
  o segment structure
  o reliable data transfer
  o flow control
  o connection management
❑ Principles of congestion control
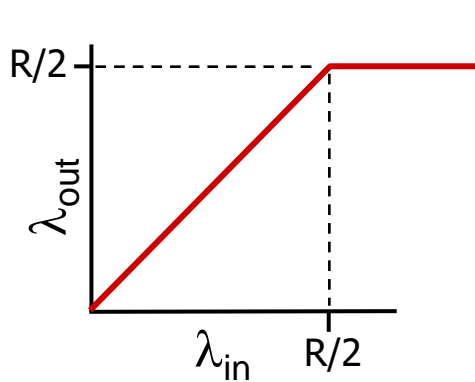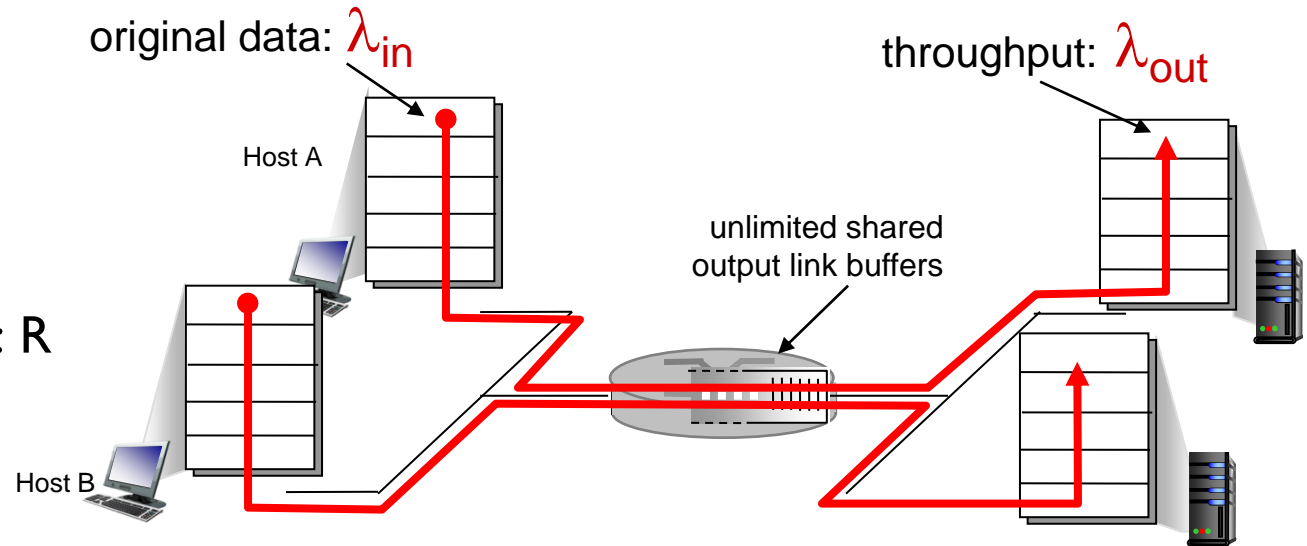❑ TCP congestion control

# Principles of Congestion Control

Congestion:

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

❑ different from flow control!

❑ Signs indicating congestion:

○ lost packets (buffer overflow at routers)

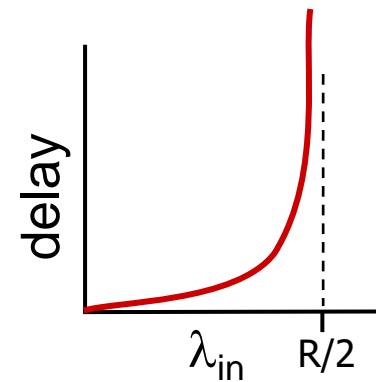○ long delays (queueing in router buffers)

❑ a top-10 problem!

# Causes/costs of congestion: scenario 1

original data: $\lambda_{in}$

throughput: $\lambda_{out}$

- ❖ two senders, two receivers
- ❖ one router, infinite buffers
- ❖ output link capacity: R
- ❖ no retransmission

Host A

Host B

unlimited shared output link buffers
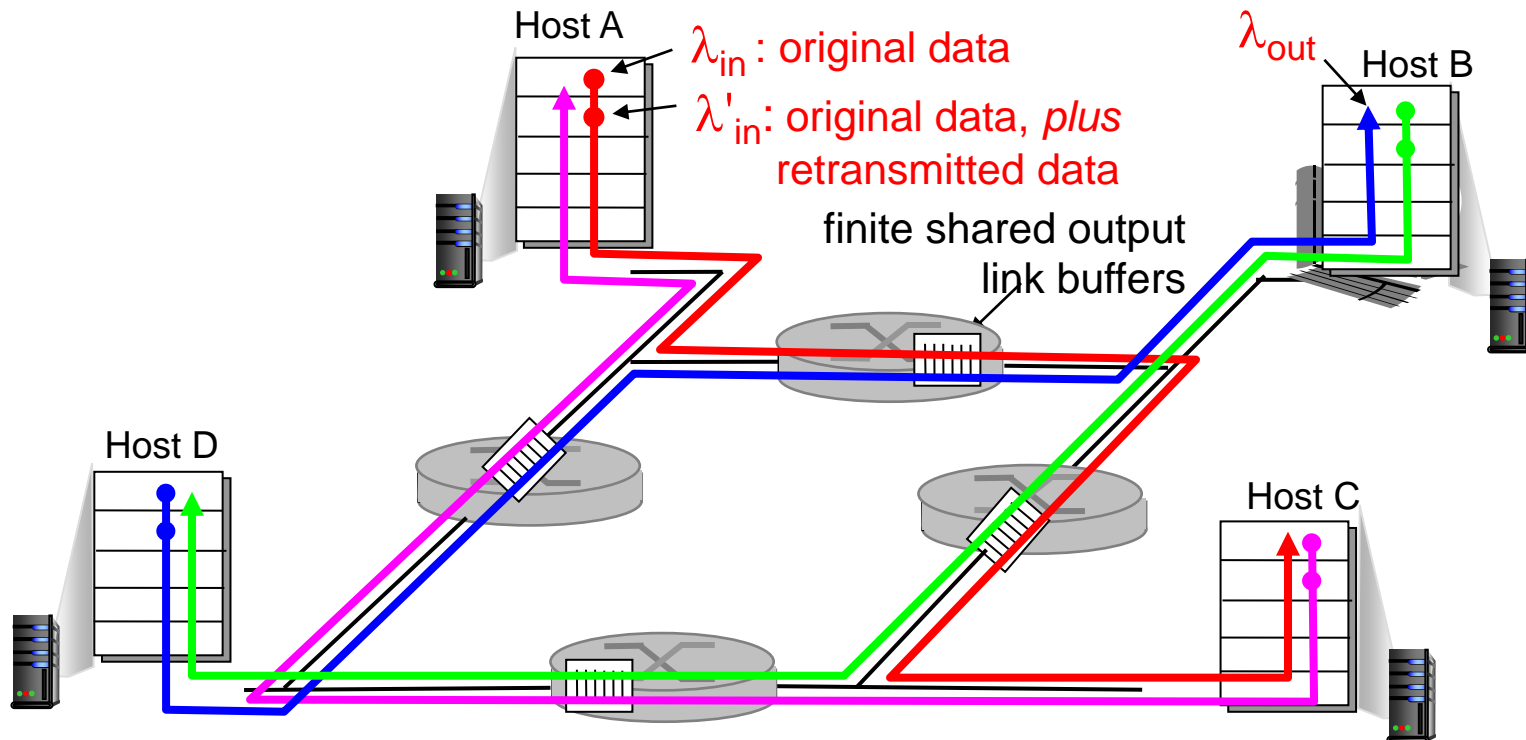


- ❖ maximum per-connection throughput: R/2
- ❖ large delays as arrival rate, $\lambda_{in}$, approaches capacity
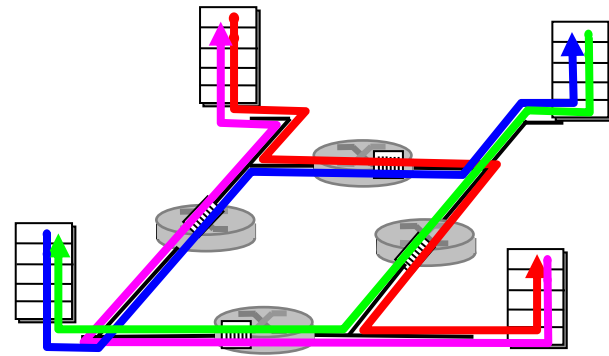
# Causes/costs of congestion: scenario 2

- ❖ four senders/receivers
- ❖ multihop paths
- ❖ timeout/retransmission

<u>Q</u>: what happens as $\lambda_{in}$ and $\lambda_{in}'$ increase ?

<u>A</u>: as red $\lambda_{in}'$ increases, all arriving blue pkts at upper queue are dropped, blue throughput → 0



Host A

$\lambda_{in}$ : original data

$\lambda_{in}'$: original data, *plus* retransmitted data

$\lambda_{out}$

Host B

finite shared output link buffers

Host D

Host C

# Causes/costs of congestion: scenario 2



**another "cost" of congestion:**

❖ when a 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:

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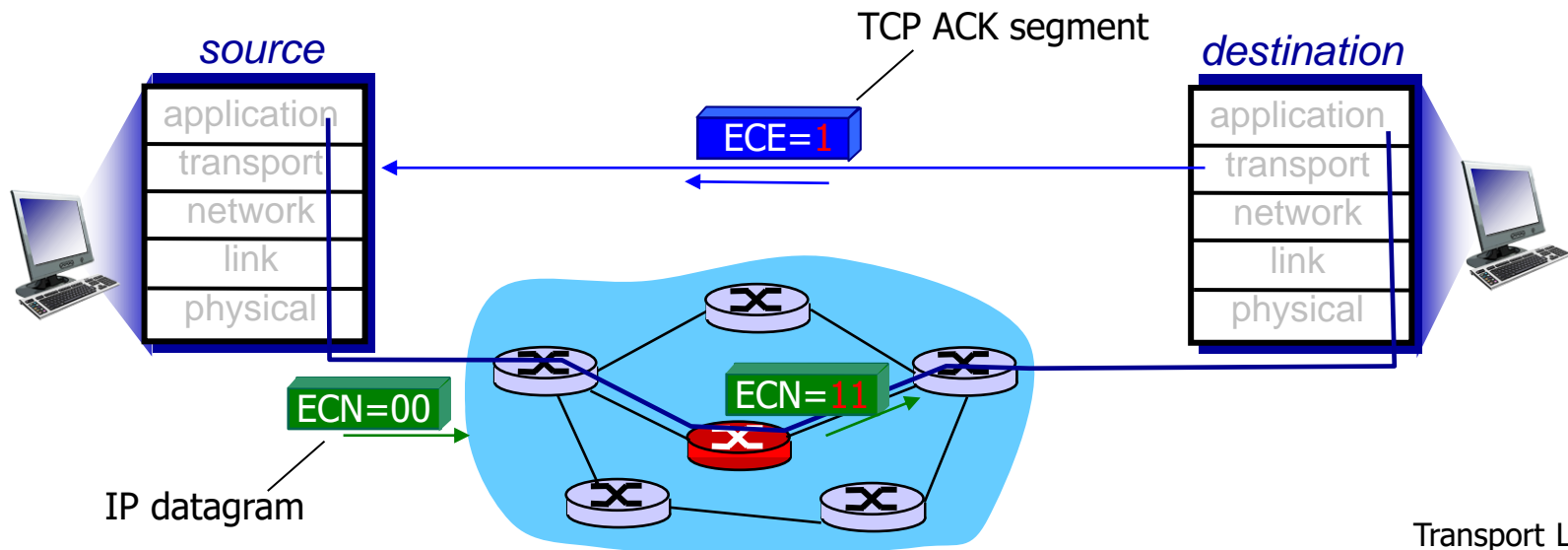- ❖ approach taken by TCP

### network-assisted congestion control:

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

# Explicit Congestion Notification (ECN)

*network-assisted congestion control:*

- two bits in ToS (Type of Service) field of IP header marked *by network router* to indicate congestion
- congestion indication carried to receiving host
- receiver (seeing congestion indication in IP datagram) ) sets ECN-Echo (ECE) bit on receiver-to-sender ACK segment to notify sender of congestion

# Transport Layer

- ❑ Transport-layer services
- ❑ Connectionless transport: UDP
- ❑ Connection-oriented transport: TCP
  - o segment structure
  - o reliable data transfer
  - o connection management
- ❑ Principles of congestion control
- ❑ TCP congestion control

# TCP congestion control:

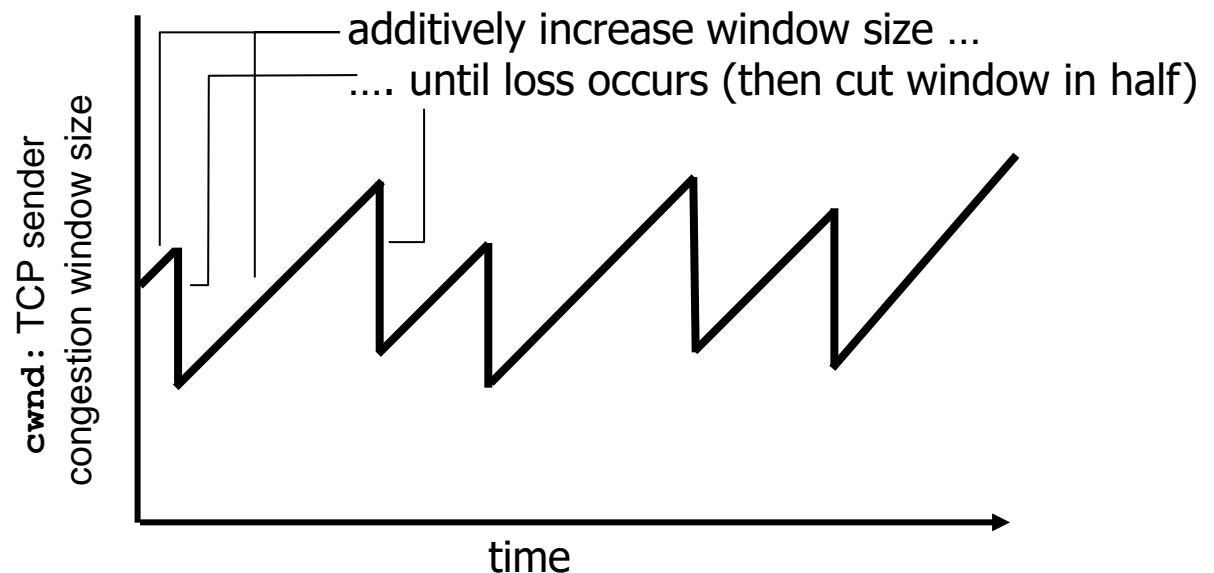❑ *goal:* TCP sender should transmit as fast as possible, but without congesting network

o <u>Q:</u> how to find rate *just* below congestion level?

❑ decentralized: each TCP sender sets its own rate, based on *implicit* feedback:

o *ACK:* segment received (a good thing!), network not congested, so increase sending rate

o *lost segment:* assume loss due to congested network, so decrease sending rate

# 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 Maximum Segment Size (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:

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

❖ **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

# Transitioning into/out of slowstart

**ssthresh: cwnd** threshold maintained by TCP

❑ on loss event: **set ssthresh to cwnd/2**

    o  remember (half of) TCP rate when congestion last occurred

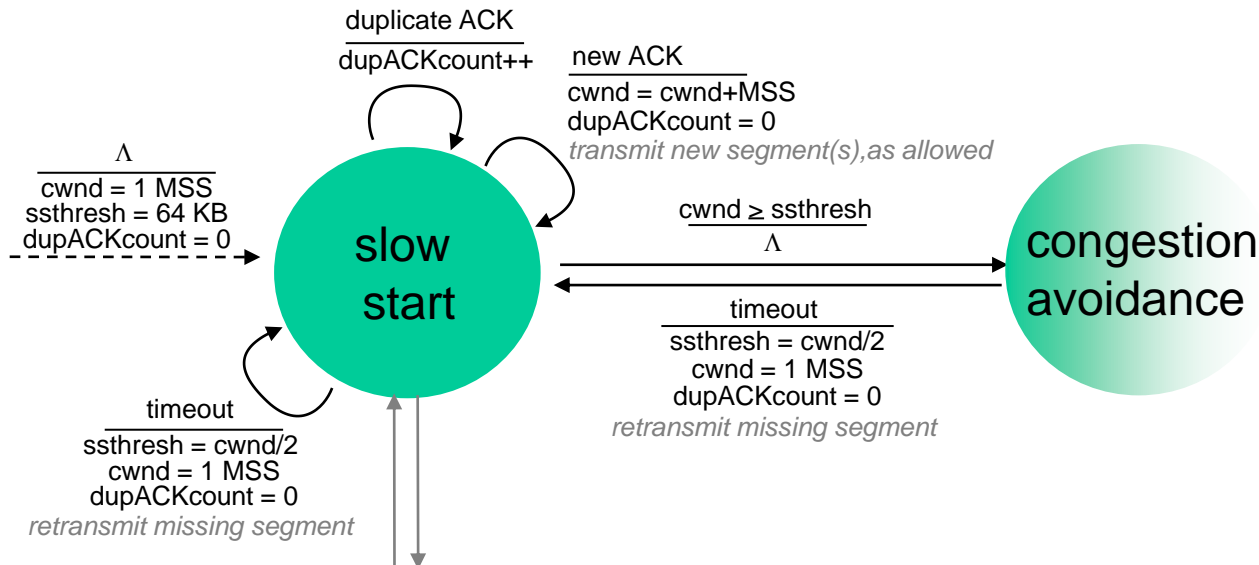❑ when **cwnd >= ssthresh**: transition from slowstart to congestion avoidance phase

duplicate ACK
————————
dupACKcount++

new ACK
————————
cwnd = cwnd+MSS
dupACKcount = 0
*transmit new segment(s),as allowed*

Λ
————————
cwnd = 1 MSS
ssthresh = 64 KB
dupACKcount = 0

**slow start**

cwnd ≥ ssthresh
————————
Λ

**congestion avoidance**

timeout
————————
ssthresh = cwnd/2
cwnd = 1 MSS
dupACKcount = 0
*retransmit missing segment*

timeout
————————
ssthresh = cwnd/2
cwnd = 1 MSS
dupACKcount = 0
*retransmit missing segment*

# TCP: congestion avoidance

❑ when `cwnd >= ssthresh`, `cwnd` grows linearly
  - o increase `cwnd` by 1 MSS per RTT
  - o approach possible congestion slower than in slowstart

A    B                                              A    B

CWND = 1

CWND = 2

CWND = 3
CWND = 4

CWND = 5
CWND = 6
CWND = 7
CWND = 8

CWND = 9
CWND = 10
CWND = 11
CWND = 12
CWND = 13
CWND = 14
CWND = 15
CWND = 16

CWND = 1

CWND = 2

CWND = 3
CWND = 4

CWND = 5
CWND = 6
CWND = 7
CWND = 8

CWND=9

Slow Start

Congestion avoidance

(a) Slow start, ending with a timeout    (b) Slow start followed by congestion avoidance

# TCP Congestion Control: segment loss event

❑ **Loss indicated by timeout:**
- o cut `cwnd` to 1 MSS
- o `cwnd` then grows exponentially (as in slow start) to threshold, then grows linearly
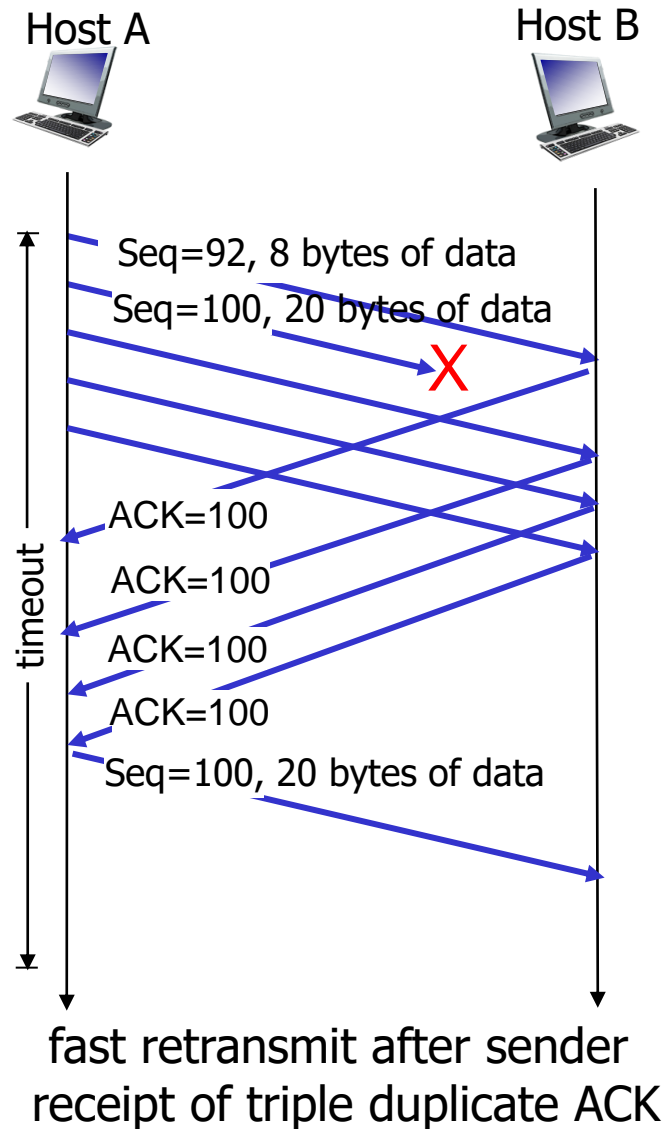
❑ **Loss indicated by 3 duplicate ACKs: TCP RENO**
- o `cwnd` is cut in half and `cwnd` then grows linearly: less aggressively than on timeout (Fast Recovery)
- o Note that TCP Tahoe always set `cwnd` to 1 (timeout or 3 duplicate acks)

---
**Philosophy:**

❑ 3 dup ACKs indicates network capable of delivering some segments (recall fast retransmit)

❑ timeout indicates a "more alarming/serious" congestion scenario

# TCP fast retransmit

Host A                                          Host B

Seq=92, 8 bytes of data
Seq=100, 20 bytes of data

X

ACK=100

ACK=100

ACK=100

ACK=100

Seq=100, 20 bytes of data

timeout

fast retransmit after sender
receipt of triple duplicate ACK

# Why Fast Recovery is used?

❑ When TCP retransmits a segment using Fast Retransmit, a segment was assumed lost

❑ Some congestion avoidance measures are appropriate at this point

❑ Slow Start may be unnecessarily conservative since multiple acks indicate segments are getting through (meaning congestion not so serious)

❑ Fast Recovery: retransmit lost segment, cut `CongWin` in half, and proceed with linear increase of `CongWin` (avoiding "slow" start-up)

# TCP Congestion Control: ACK received

## ACK received: increase CWND

❑ slowstart phase:
   o increase exponentially fast (despite name) at connection start or following timeout

❑ congestion avoidance:
   o increase linearly

# Variants of TCP Congestion Control Schemes

❑ TCP Tahoe: Slow Start + Congestion Avoidance.

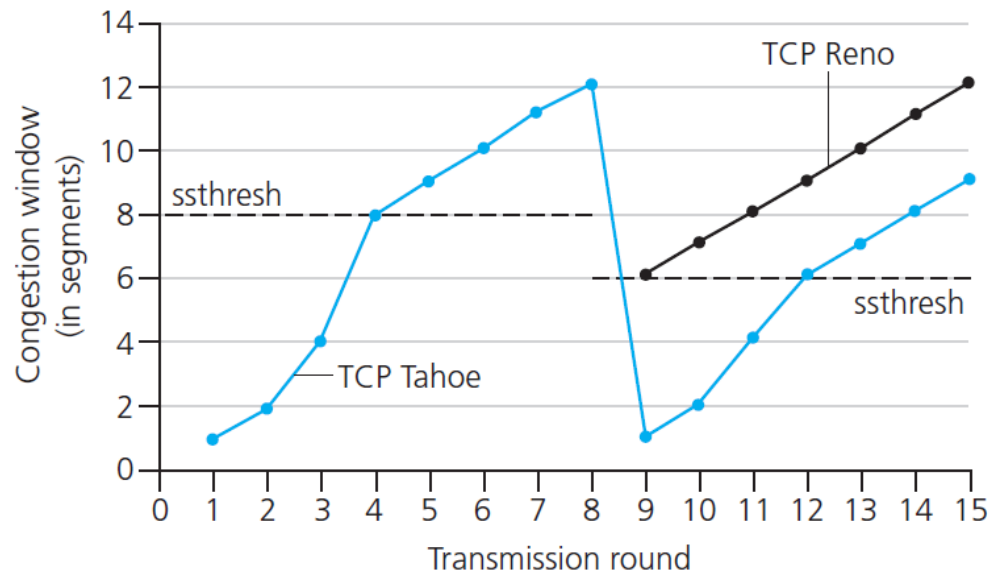❑ TCP Reno: TCP Tahoe + Fast Retransmit + Fast Recovery.

# TCP: from slow start to congestion avoidance

**Q:** when should the exponential increase switch to linear?
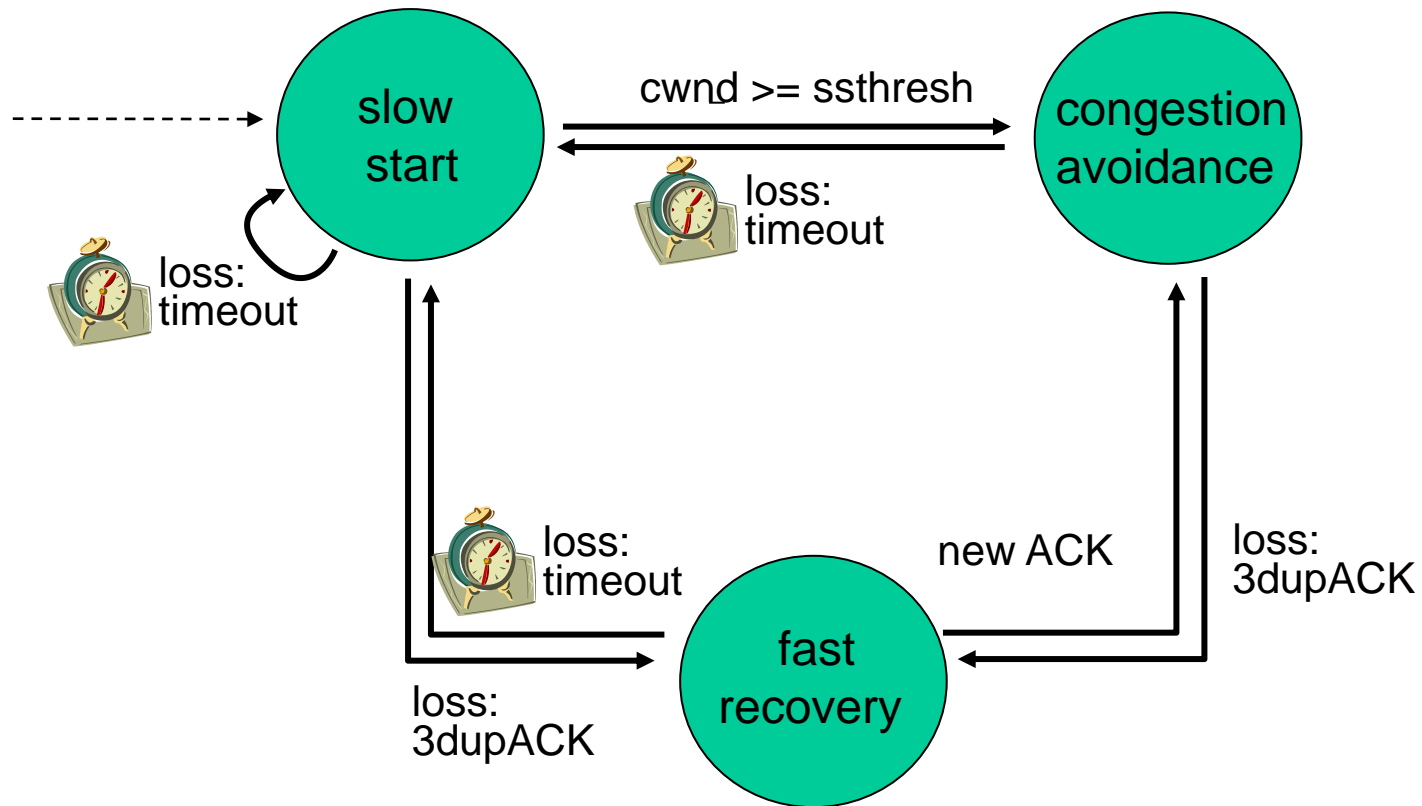
**A:** when **cwnd** gets to 1/2 of its value before last timeout.
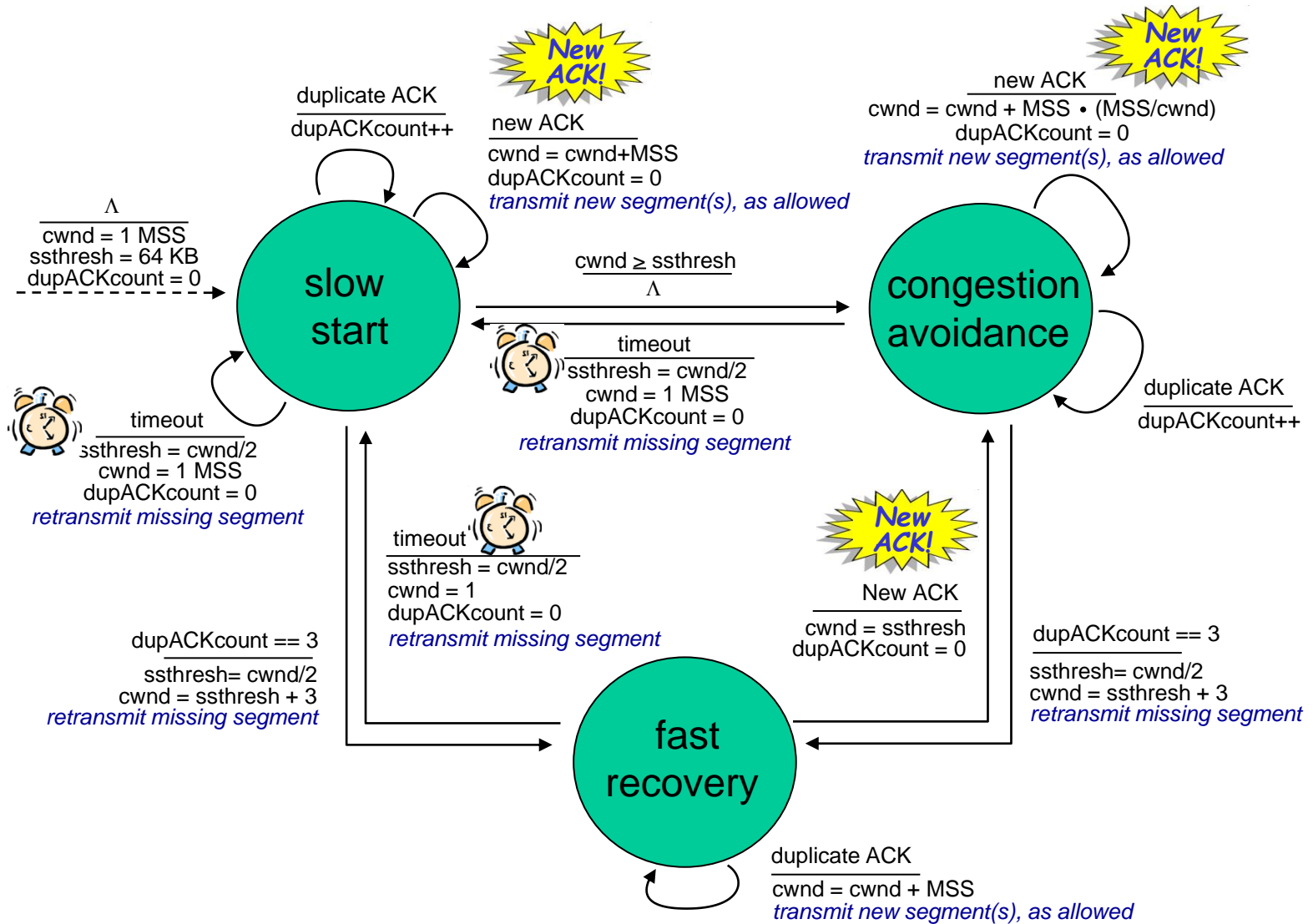
## Implementation:

❖ variable **ssthresh**

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

# TCP Reno congestion control FSM: overview

# TCP Reno congestion control FSM: details



duplicate ACK
$\overline{\text{dupACKcount++}}$

New ACK!

new ACK
$\overline{\text{cwnd = cwnd+MSS}}$
dupACKcount = 0
*transmit new segment(s), as allowed*

$\Lambda$
$\overline{\text{cwnd = 1 MSS}}$
ssthresh = 64 KB
dupACKcount = 0

**slow start**

New ACK!

new ACK
$\overline{\text{cwnd = cwnd + MSS} \cdot \text{(MSS/cwnd)}}$
dupACKcount = 0
*transmit new segment(s), as allowed*

cwnd ≥ ssthresh
$\overline{\Lambda}$

timeout
$\overline{\text{ssthresh = cwnd/2}}$
cwnd = 1 MSS
dupACKcount = 0
*retransmit missing segment*

**congestion avoidance**

duplicate ACK
$\overline{\text{dupACKcount++}}$

timeout
$\overline{\text{ssthresh = cwnd/2}}$
cwnd = 1 MSS
dupACKcount = 0
*retransmit missing segment*

timeout
$\overline{\text{ssthresh = cwnd/2}}$
cwnd = 1
dupACKcount = 0
*retransmit missing segment*

New ACK!

New ACK
$\overline{\text{cwnd = ssthresh}}$
dupACKcount = 0

dupACKcount == 3
$\overline{\text{ssthresh= cwnd/2}}$
cwnd = ssthresh + 3
*retransmit missing segment*

dupACKcount == 3
$\overline{\text{ssthresh= cwnd/2}}$
cwnd = ssthresh + 3
*retransmit missing segment*

**fast recovery**

duplicate ACK
$\overline{\text{cwnd = cwnd + MSS}}$
*transmit new segment(s), as allowed*

# Summary: TCP Congestion Control
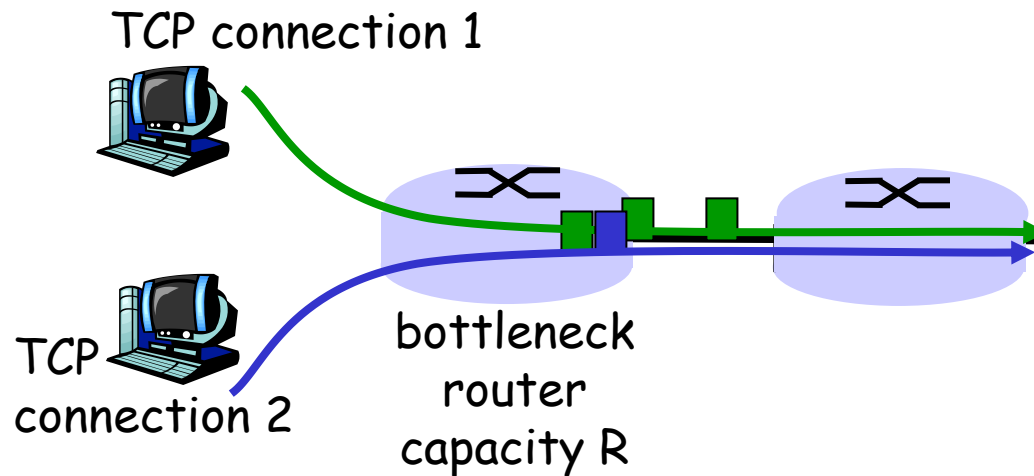
❑ when `cwnd < ssthresh`, sender is in slow-start phase, window grows exponentially.

❑ when `cwnd >= ssthresh`, sender is in congestion-avoidance phase, window grows linearly.

❑ when triple duplicate ACK occurs, `ssthresh set` to `cwnd/2,` `cwnd set` to `ssthresh` and go to congestion-avoidance phase.

❑ when timeout occurs, `ssthresh set` to `cwnd/2,` `cwnd set` to 1 MSS and go to slow-start phase.

# TCP sender congestion control

| State | Event | TCP Sender Action | Commentary |
|---|---|---|---|
| Slow Start (SS) | ACK receipt for previously unacked data | CongWin = CongWin + 1 MSS, If (CongWin >= Threshold)     set state to "Congestion Avoidance" | Resulting in a doubling of CongWin every RTT |
| Congestion Avoidance (CA) | ACK receipt for previously unacked data | CongWin = CongWin + 1 MSS * (MSS/CongWin) | Additive increase, resulting in increase of CongWin by 1 MSS every RTT |
| SS or CA | Loss event detected by triple duplicate ACK | Threshold = CongWin/2, CongWin = Threshold, Set state to "Congestion Avoidance" | Fast recovery, implementing multiplicative decrease. CongWin will not drop below 1 MSS. |
| SS or CA | Timeout | Threshold = CongWin/2, CongWin = 1 MSS, Set state to "Slow Start" | Enter slow start |
| SS or CA | Duplicate ACK | Increment duplicate ACK count for segment being acked | CongWin and Threshold not changed |

# TCP Fairness

fairness goal: if K TCP sessions share same bottleneck link of bandwidth R, each should have average rate of R/K

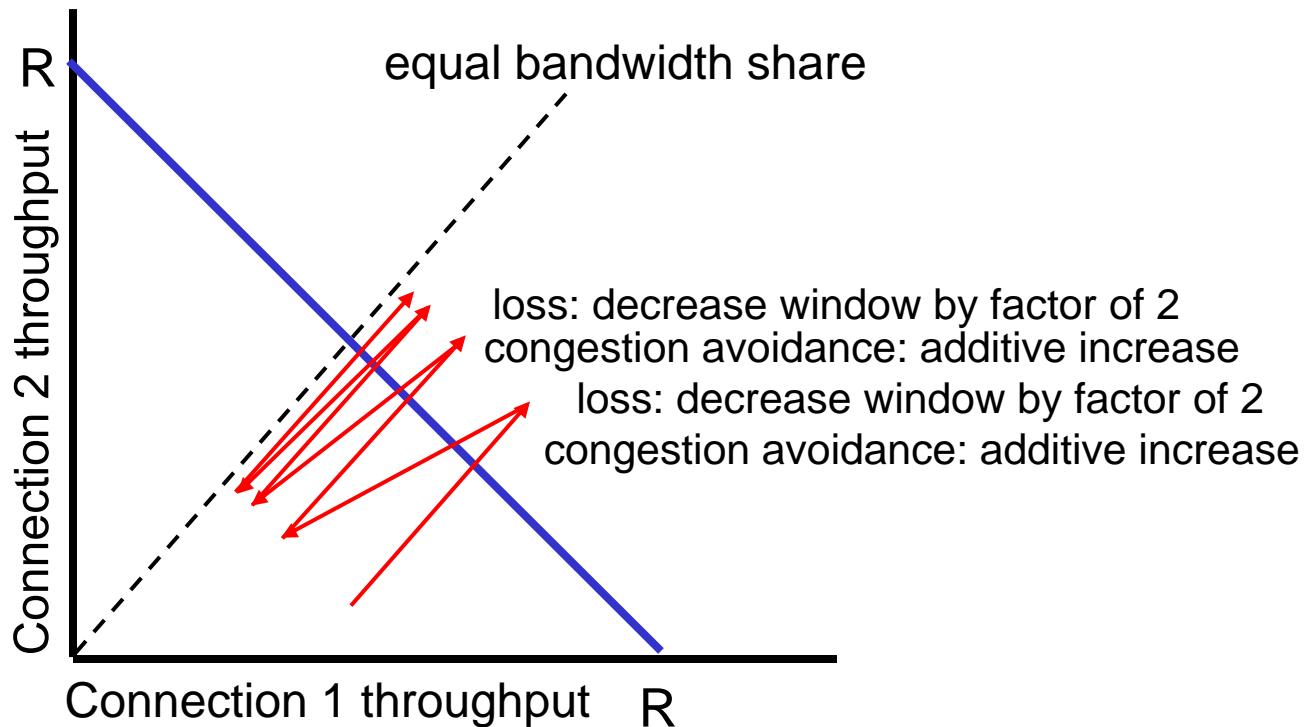TCP connection 1

TCP connection 2

bottleneck router capacity R

# Why is TCP fair?

two competing sessions:

❖ additive increase gives slope of 1, as throughput increases
❖ multiplicative decrease decreases throughput proportionally



equal bandwidth share

Connection 2 throughput

loss: decrease window by factor of 2
congestion avoidance: additive increase
loss: decrease window by factor of 2
congestion avoidance: additive increase

Connection 1 throughput    R

# Fairness (more)

## Fairness and UDP

- ❑ multimedia apps often do not use TCP
  - o do not want rate slowed down by congestion control
- ❑ instead use UDP:
  - o pump audio/video at constant rate, tolerate packet loss

## Fairness and parallel TCP connections

- ❑ nothing prevents app from opening parallel connections between 2 hosts.
- ❑ web browsers do this
- ❑ For example: link of rate R supporting 9 connections;
  - o new app asks for 1 TCP, gets rate R/10
  - o new app asks for 9 TCPs, gets rate R/2 !

# Summary

❑ Transport-layer services
❑ Connectionless transport: UDP
❑ Connection-oriented transport: TCP
   o segment structure
   o reliable data transfer
   o flow control
   o connection management
❑ Principles of congestion control
❑ TCP congestion control

# Q & A