# Chapter 3
# Transport Layer

**Computer Networking: A Top Down Approach**
6th edition
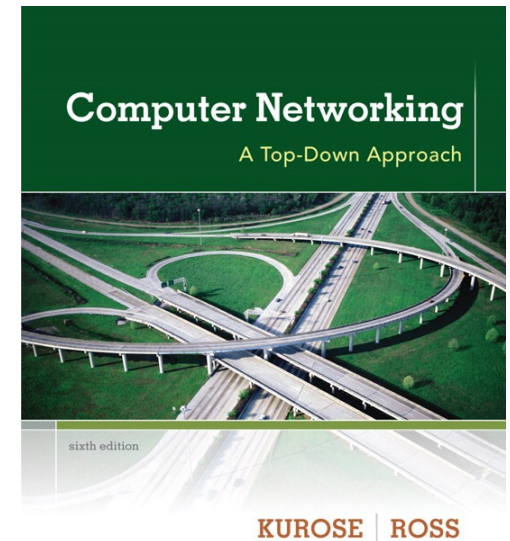Jim Kurose, Keith Ross
Addison-Wesley
March 2012

## A note on the use of these ppt 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

# Chapter 3 outline
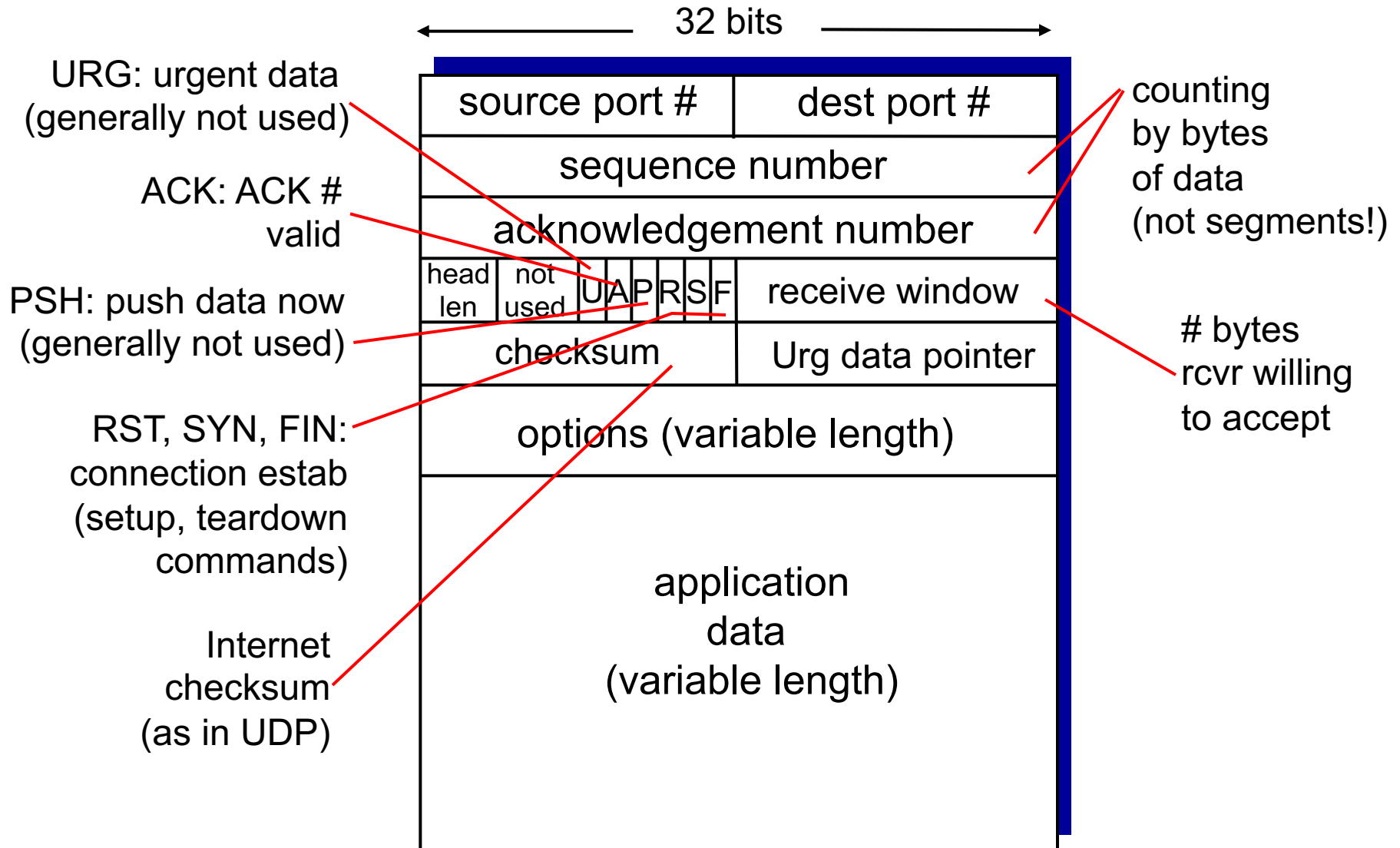
# TCP: Overview  RFCs: 793,1122,1323, 2018, 2581

❖ **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



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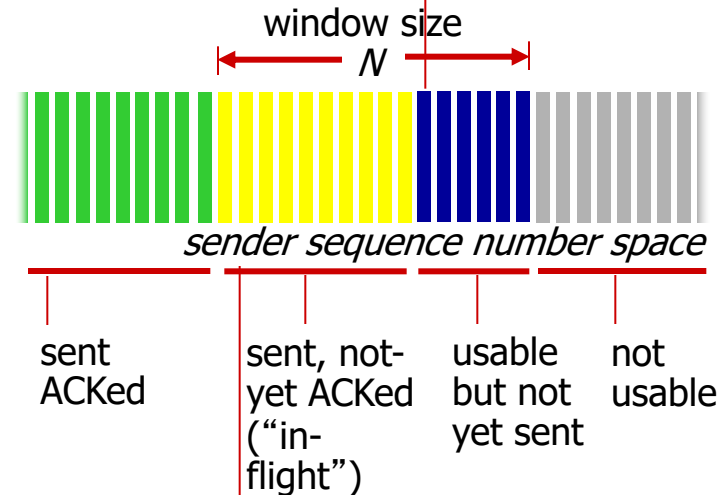
**acknowledgements:**
- seq # of next byte expected from other side
- cumulative ACK

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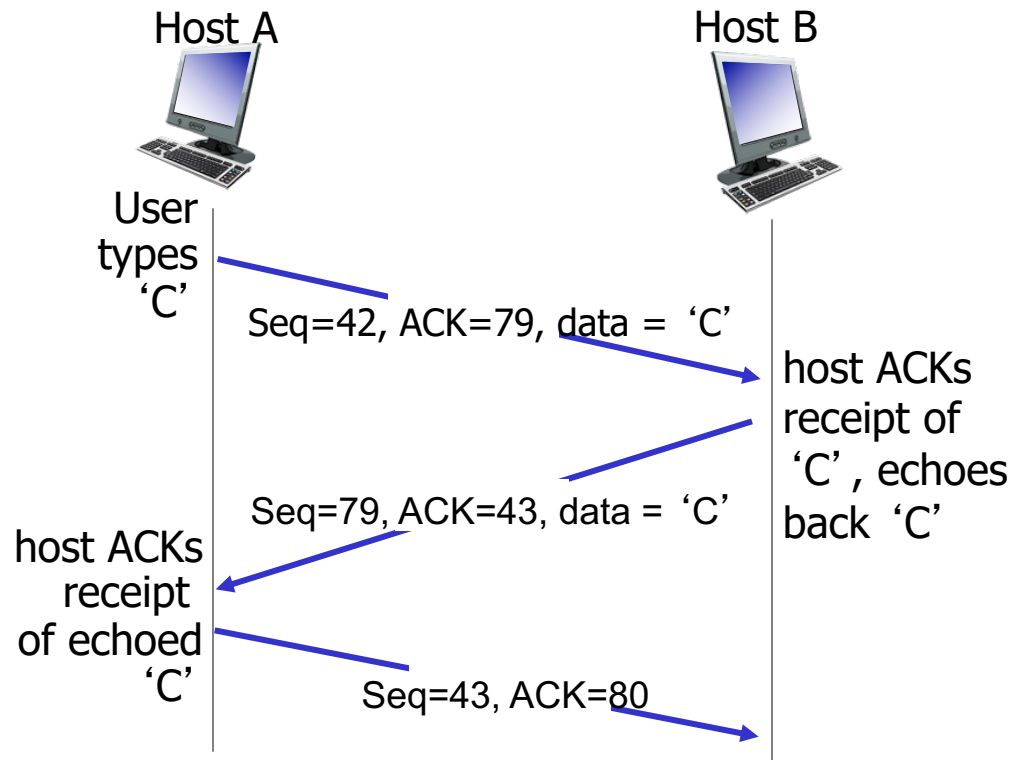- A: TCP spec doesn't say, - up to implementor

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
back 'C'

Seq=79, ACK=43, data = '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
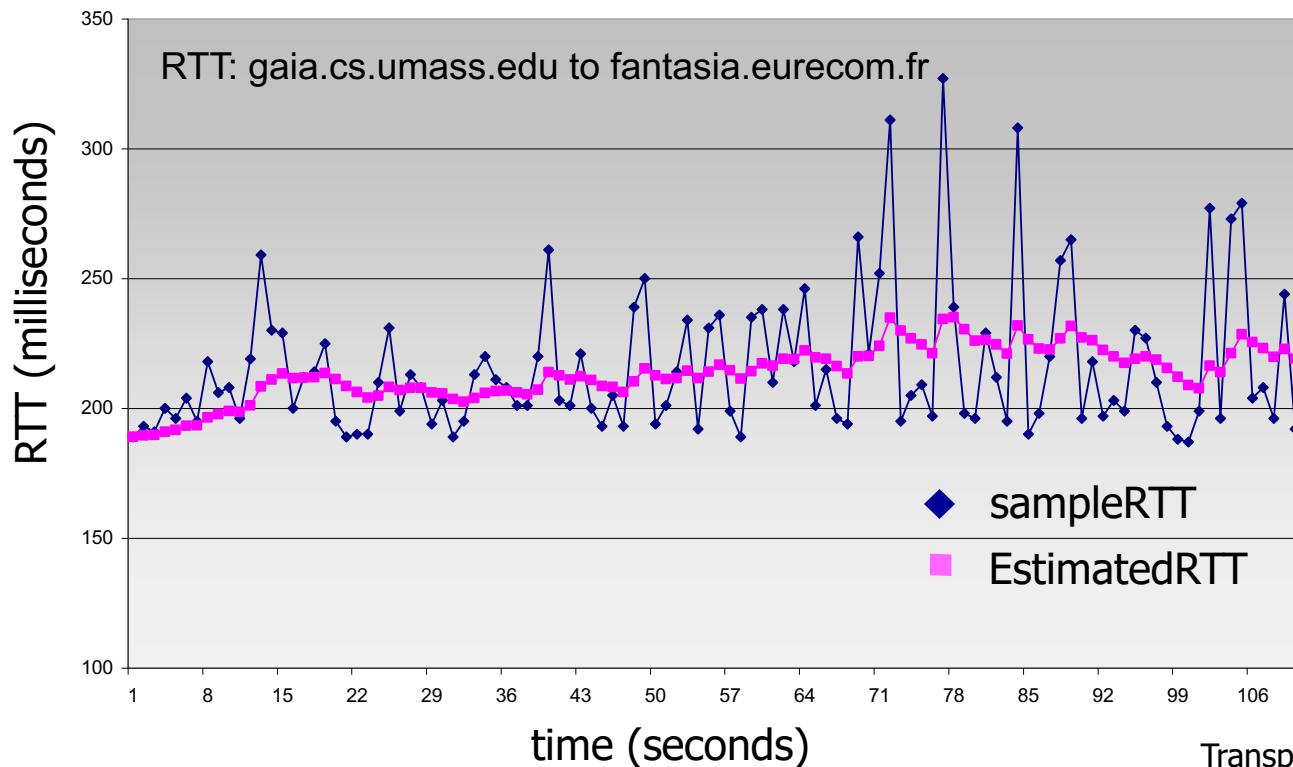❖ *too long:* slow reaction to segment loss

**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

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

- ❖ 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

time (seconds)

# TCP round trip time, timeout

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

❖ estimate SampleRTT deviation from EstimatedRTT:

$$DevRTT = (1-\beta)*DevRTT +$$
$$\beta*|SampleRTT-EstimatedRTT|$$

$$(typically, \beta = 0.25)$$

$$TimeoutInterval = EstimatedRTT + 4*DevRTT$$

estimated RTT          "safety margin"

# Chapter 3 outline

# 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
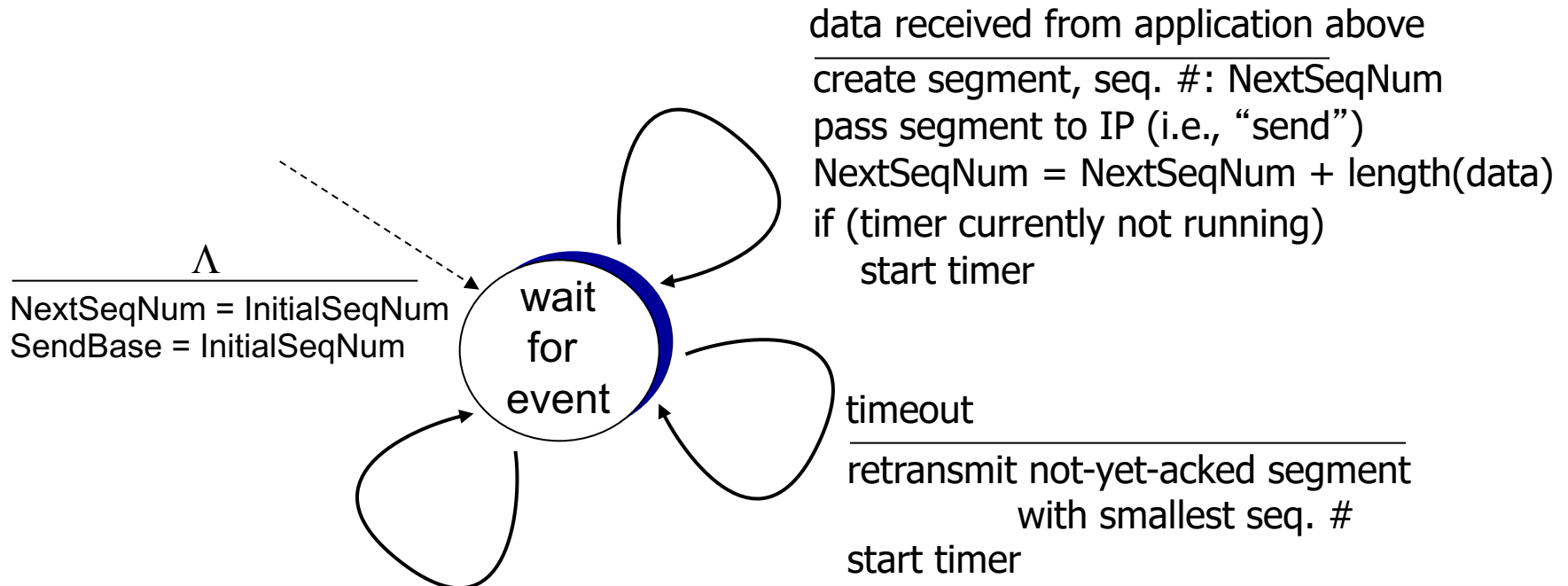  - expiration interval: `TimeOutInterval`

*timeout:*

- ❖ retransmit segment that caused timeout
- ❖ restart timer

*ack rcvd:*

- ❖ if ack acknowledges previously unacked segments
  - update what is known to be ACKed
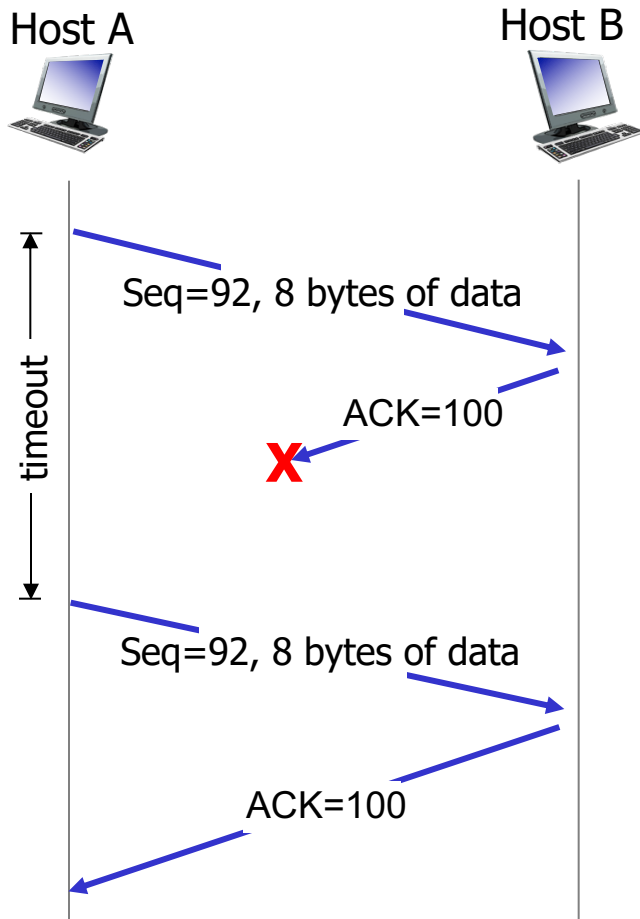  - start timer if there are still unacked segments
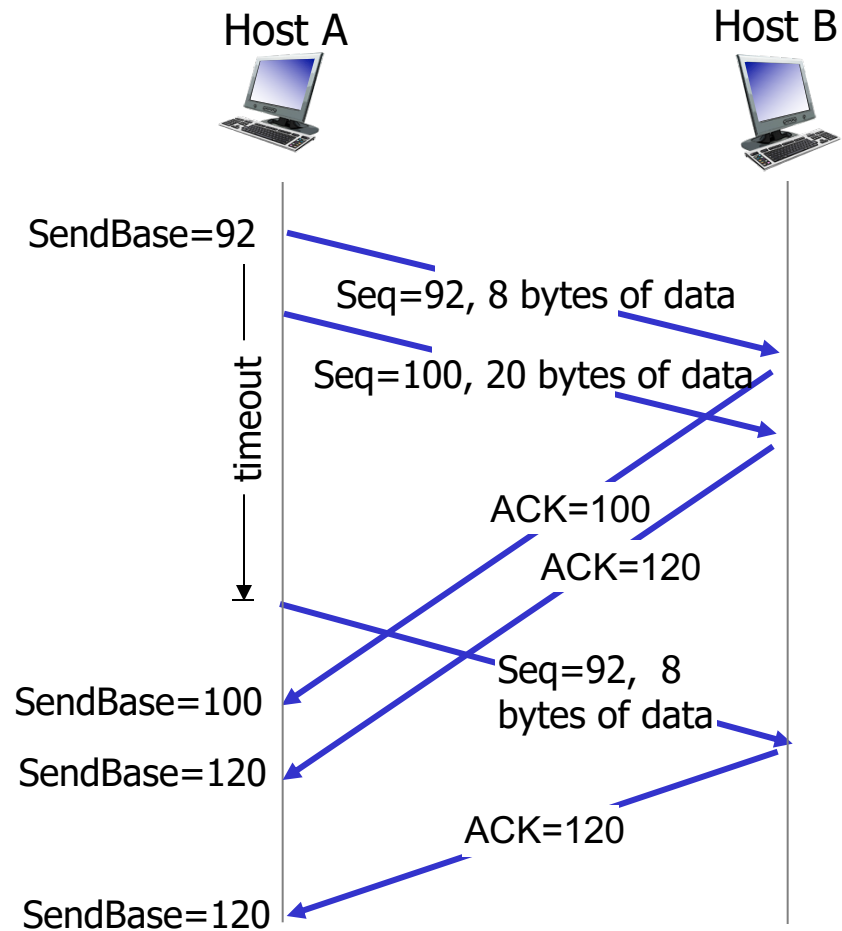
# TCP sender (simplified)

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
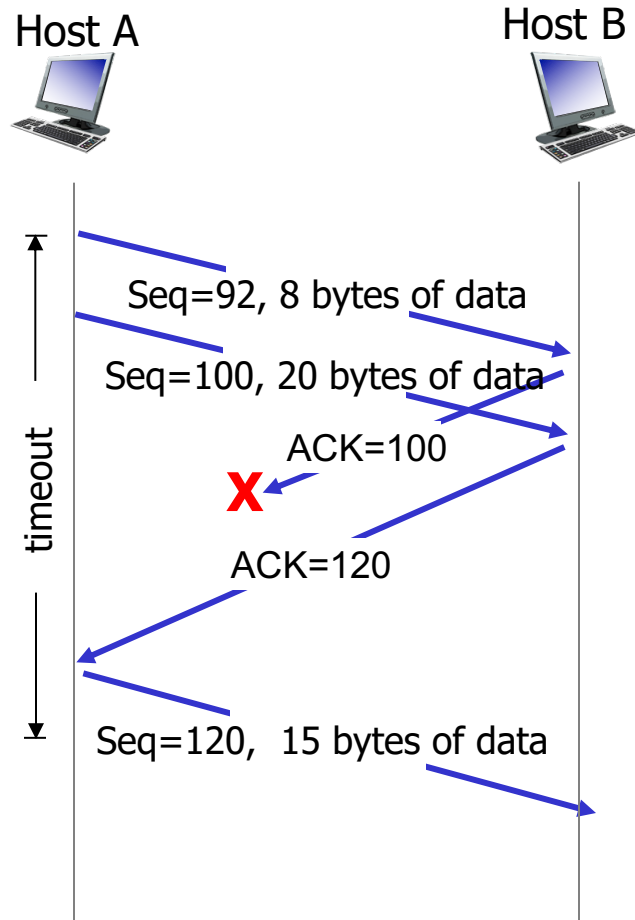      else stop timer
    }

# TCP: retransmission scenarios



Host A                    Host B

Seq=92, 8 bytes of data
                          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
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

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 3 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



fast retransmit after sender
receipt of triple duplicate ACK

# Chapter 3 outline
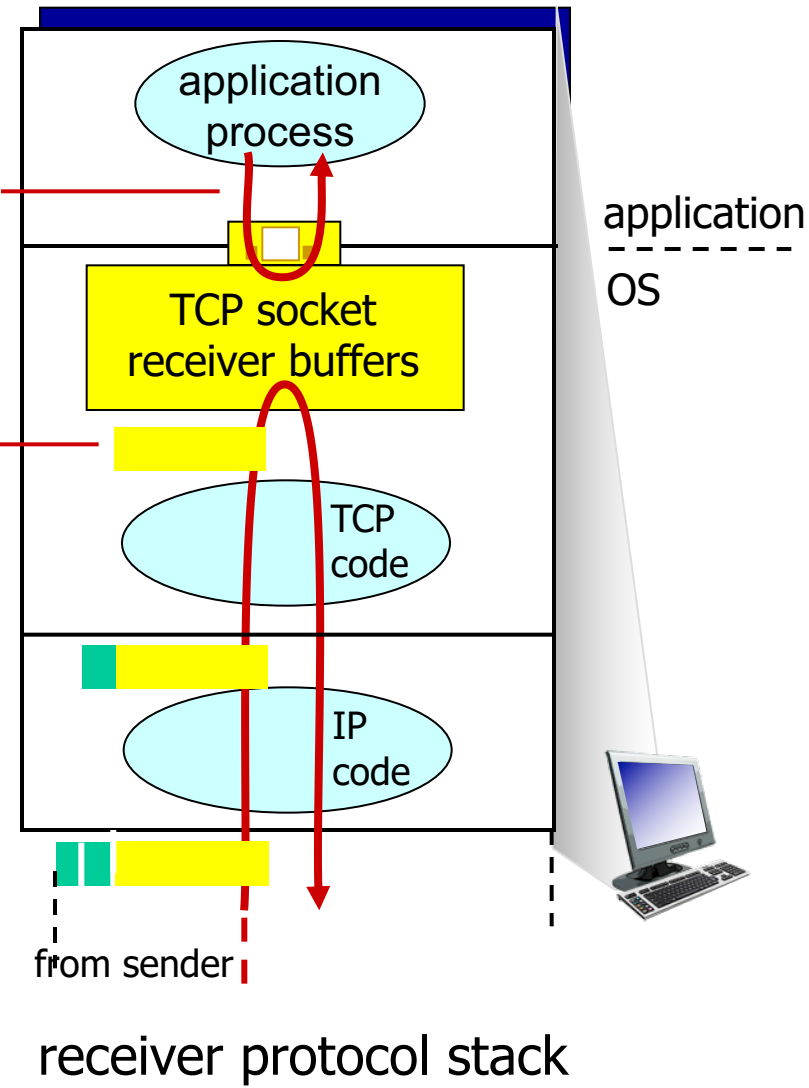
# TCP flow control

application may
remove data from
TCP socket buffers ….

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

application
process

application
--------
OS

TCP socket
receiver buffers

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*
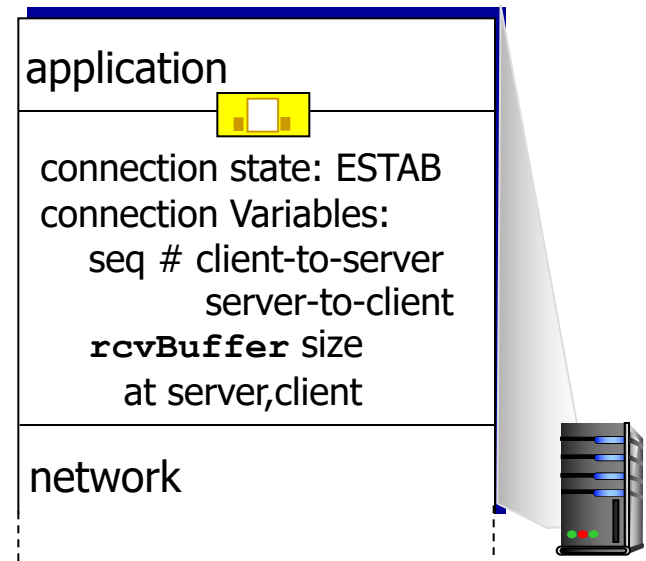
# Chapter 3 outline

# Connection Management

before exchanging data, sender/receiver "handshake":

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

❖ agree on connection parameters

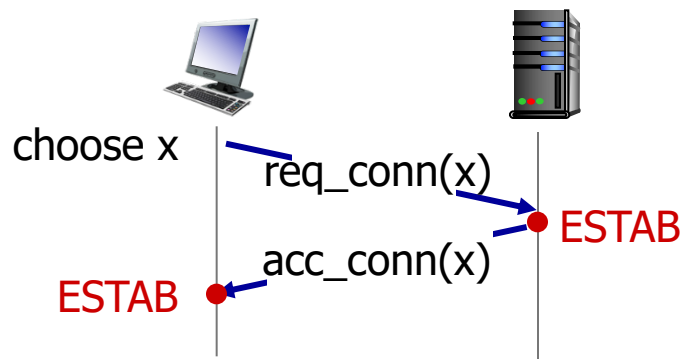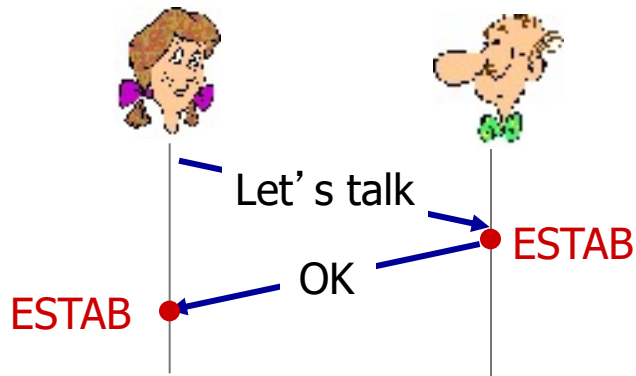| application | | application |
|---|---|---|
| connection state: ESTAB<br>connection variables:<br>   seq # client-to-server<br>       server-to-client<br>   **rcvBuffer** size<br>    at server,client | | connection state: ESTAB<br>connection Variables:<br>   seq # client-to-server<br>       server-to-client<br>   **rcvBuffer** size<br>    at server,client |
| network | | network |

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

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

# Agreeing to establish a connection

## 2-way handshake:



Let's talk → ESTAB

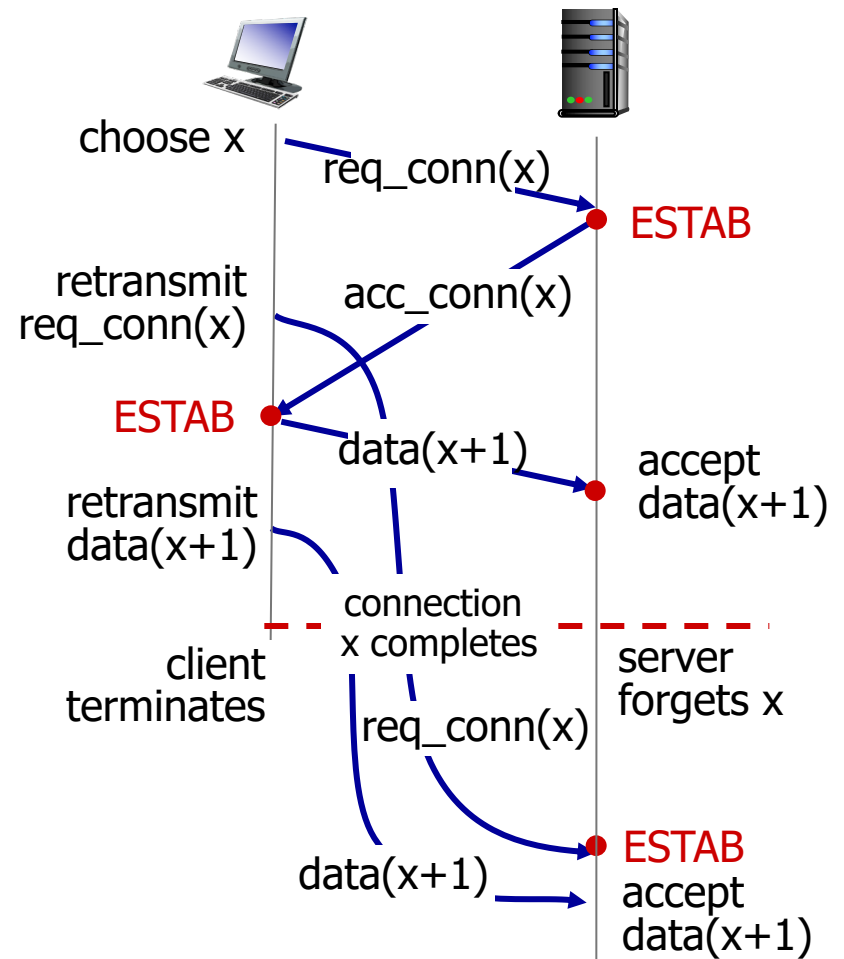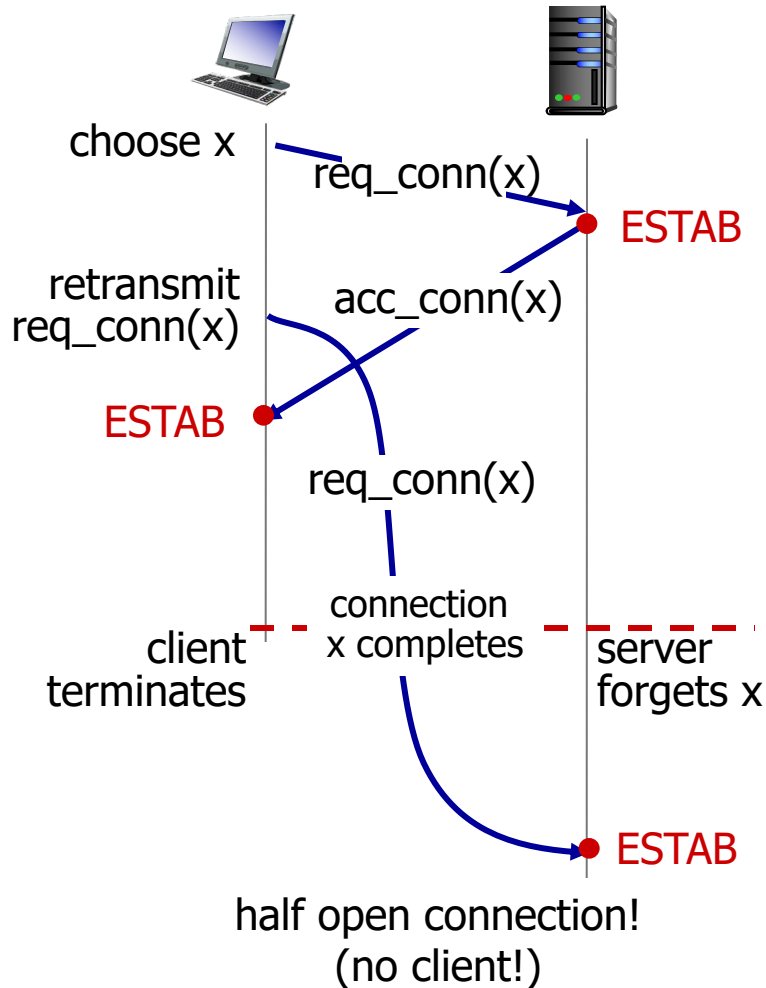OK ← ESTAB

choose x → req_conn(x) → ESTAB

acc_conn(x) → ESTAB

*Q:* will 2-way handshake always work in network?
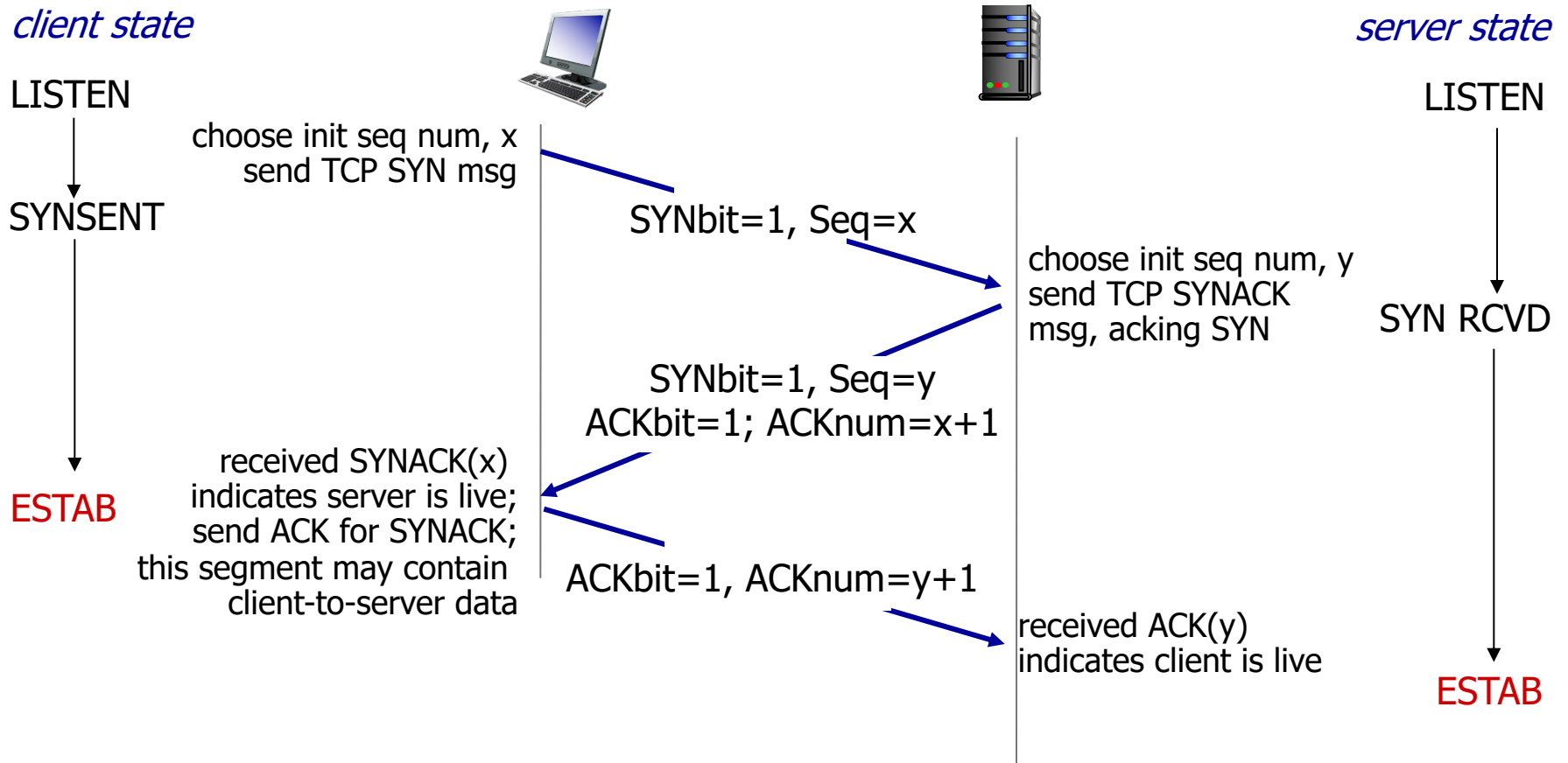
❖ variable delays
❖ retransmitted messages (e.g. req_conn(x)) due to message loss
❖ message reordering
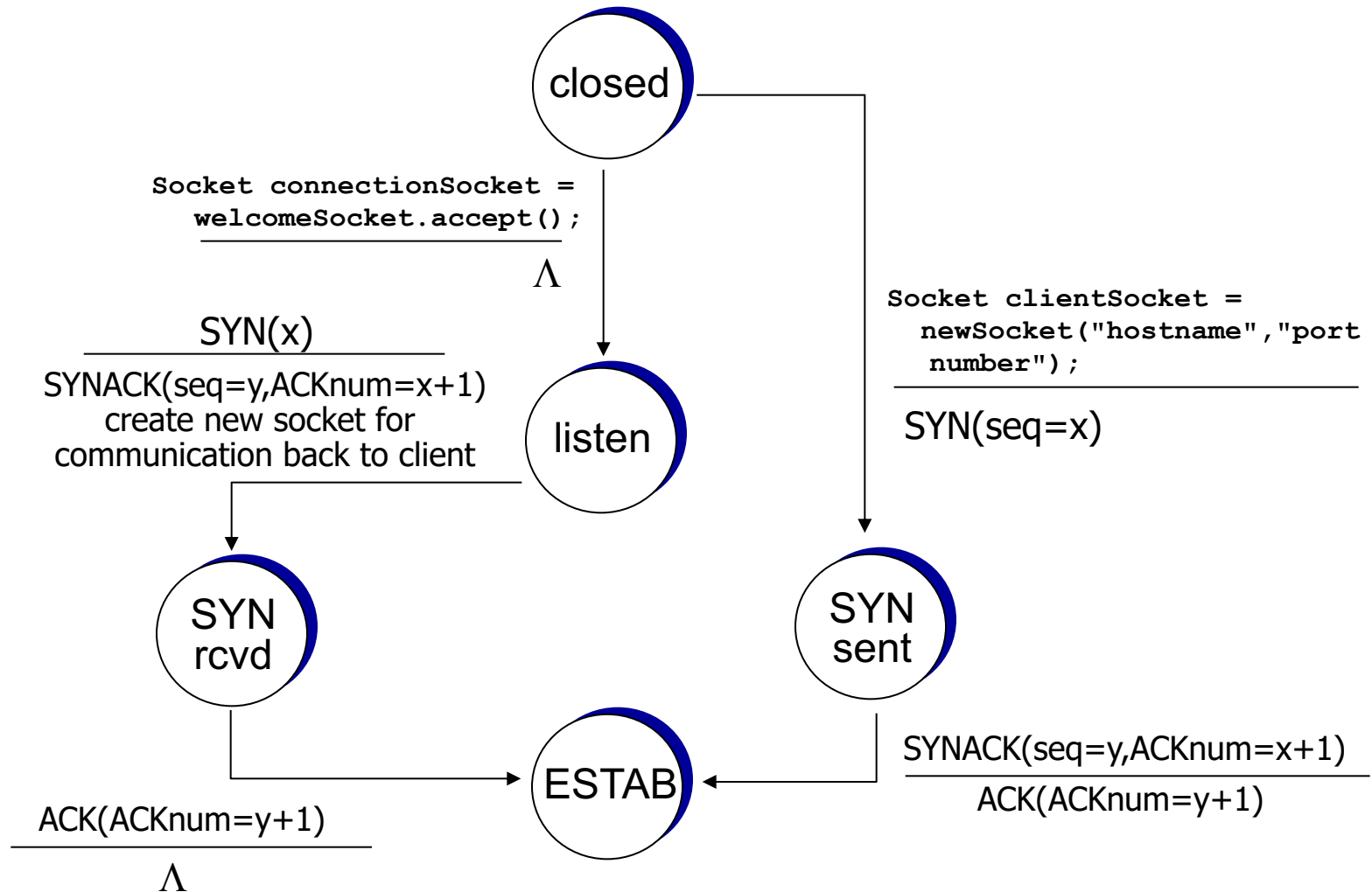❖ can't "see" other side

# Agreeing to establish a connection

2-way handshake failure scenarios:



choose x
req_conn(x)
ESTAB
retransmit req_conn(x)
acc_conn(x)
ESTAB
req_conn(x)
client terminates
connection x completes
server forgets x
ESTAB
half open connection! (no client!)

choose x
req_conn(x)
ESTAB
retransmit req_conn(x)
acc_conn(x)
ESTAB
data(x+1)
accept data(x+1)
retransmit data(x+1)
client terminates
connection x completes
server forgets x
req_conn(x)
data(x+1)
ESTAB
accept data(x+1)

# TCP 3-way handshake

*client state*

LISTEN

SYNSENT

ESTAB

choose init seq num, x
send TCP SYN msg

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

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 3-way handshake: FSM

closed

Socket connectionSocket =
    welcomeSocket.accept();
——————————————————————
Λ

Socket clientSocket =
    newSocket("hostname","port
    number");
——————————————————————
SYN(seq=x)

SYN(x)
——————————————————————
SYNACK(seq=y,ACKnum=x+1)
create new socket for
communication back to client

listen

SYN
rcvd

SYN
sent

ESTAB

ACK(ACKnum=y+1)
——————————————————————
Λ

SYNACK(seq=y,ACKnum=x+1)
——————————————————————
ACK(ACKnum=y+1)

# 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

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

can still
send data

FIN_WAIT_2   wait for server
             close

LAST_ACK

FINbit=1, seq=y

TIMED_WAIT

can no longer
send data

ACKbit=1; ACKnum=y+1

timed wait
for 2*max
segment lifetime

CLOSED

CLOSED

# Chapter 3 outline

# Principles of congestion control

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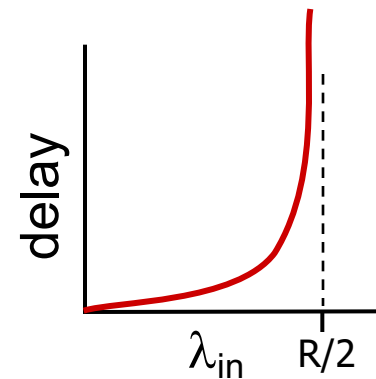
❖ informally: "too many sources sending too much data too fast for *network* to handle"

❖ different from flow control!

❖ manifestations:

  ▪ lost packets (buffer overflow at routers)

  ▪ long delays (queueing in router buffers)

❖ a top-10 problem!

# Causes/costs of congestion: scenario 1

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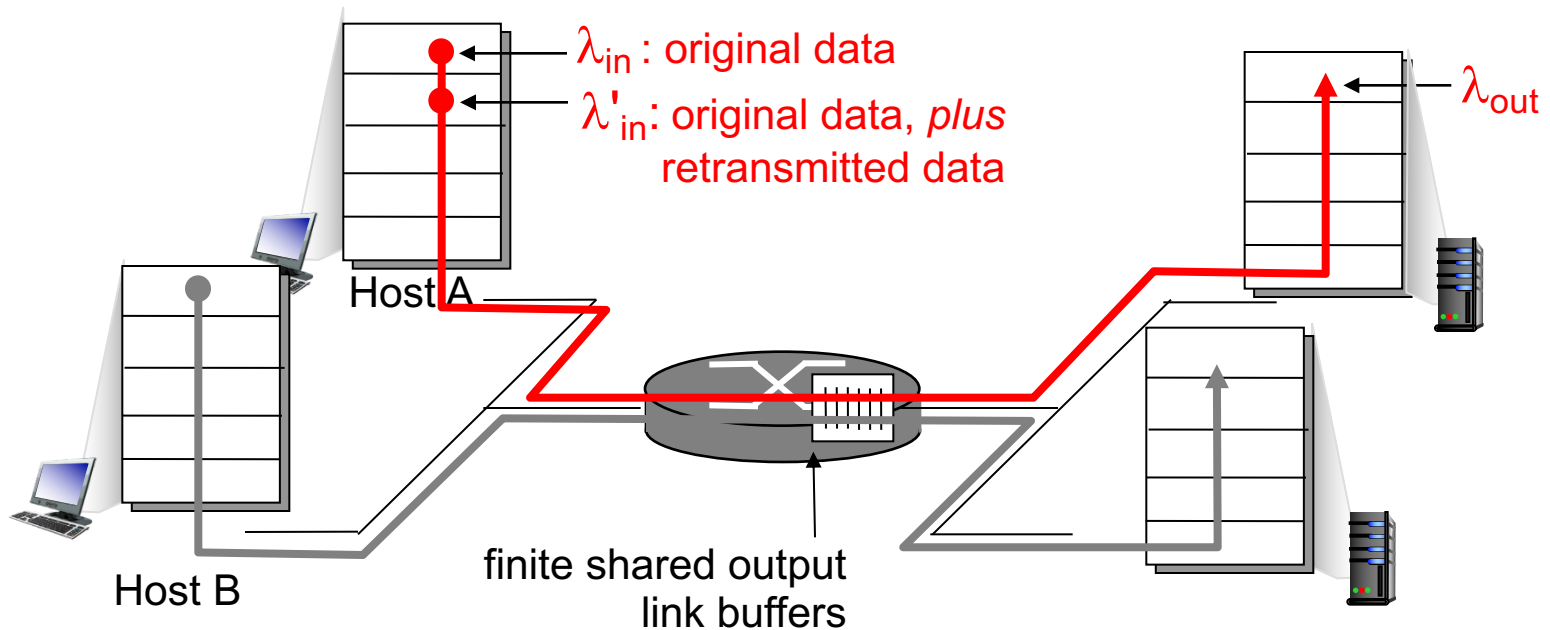- ❖ output link capacity: R
- ❖ no retransmission

original data: $\lambda_{in}$

Host A

throughput: $\lambda_{out}$

unlimited shared output link buffers

Host B

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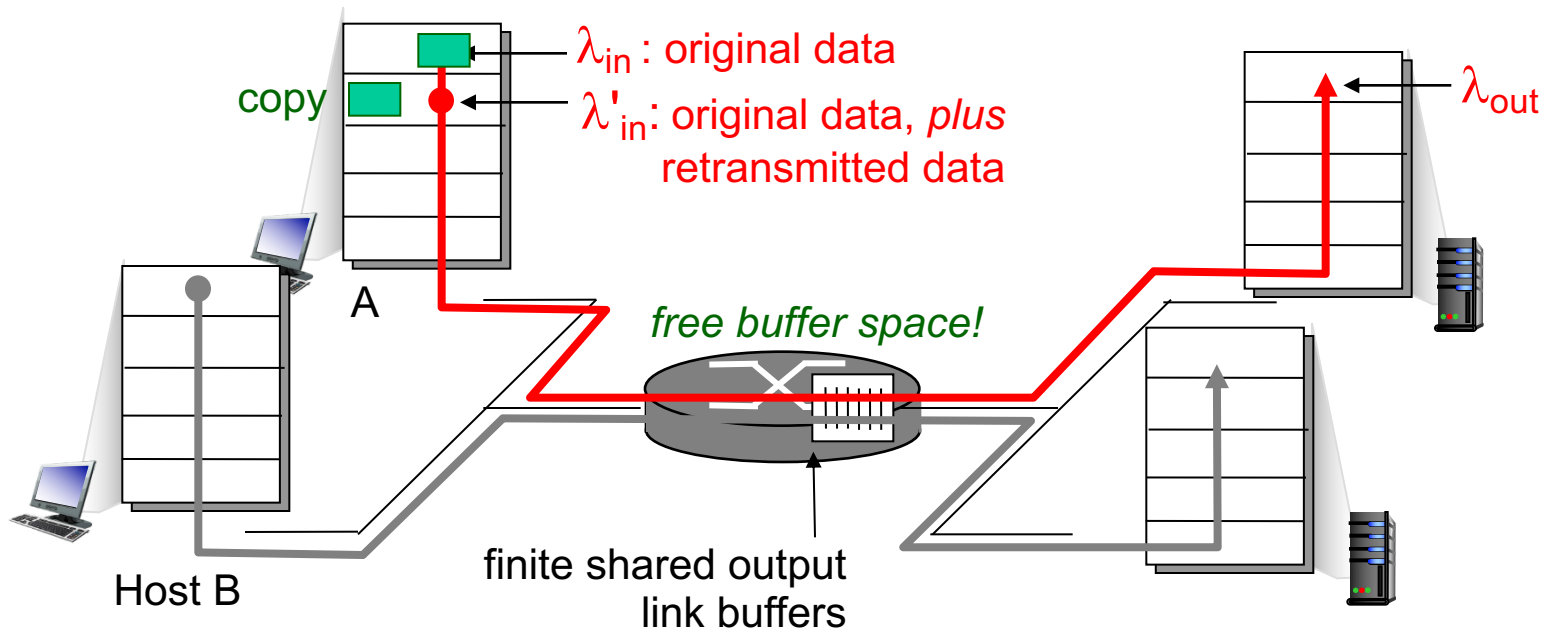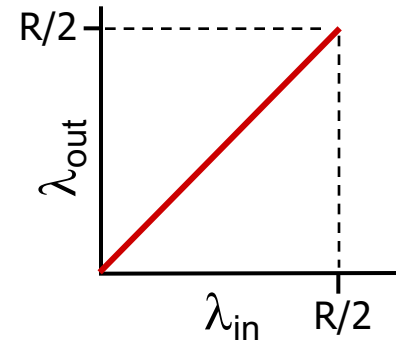
# Causes/costs of congestion: scenario 2

❖ one router, *finite* buffers

❖ sender retransmission of timed-out packet
  ▪ application-layer input = application-layer output: $\lambda_{in}$ = $\lambda_{out}$
  ▪ transport-layer input includes *retransmissions* : $\lambda'_{in} \geq \lambda_{in}$



$\lambda_{in}$ : original data

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

$\lambda_{out}$

Host A

Host B

finite shared output link buffers

# Causes/costs of congestion: scenario 2

idealization: perfect knowledge

❖ sender sends only when router buffers available



λ_in : original data

λ'_in: original data, *plus* retransmitted data

copy

A

Host B

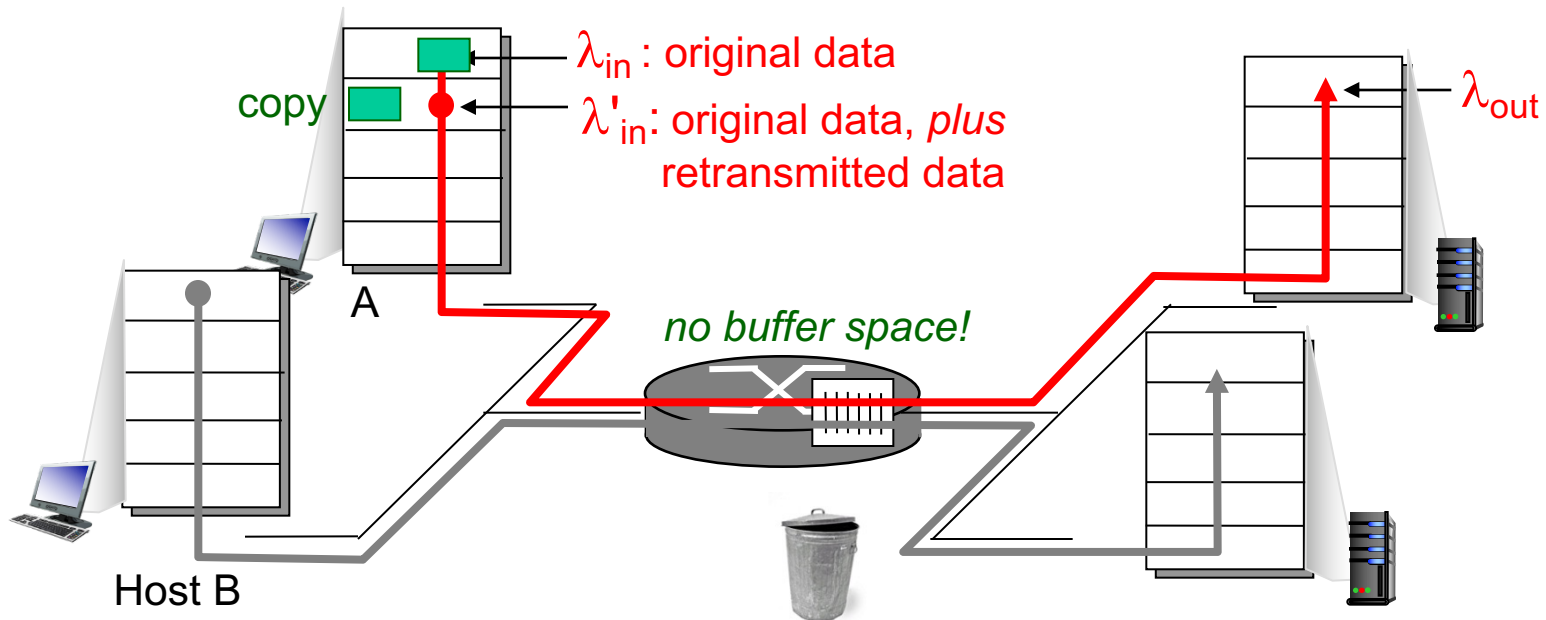free buffer space!

λ_out

finite shared output link buffers

# Causes/costs of congestion: scenario 2

*Idealization: known loss*
packets can be lost, dropped at router due to full buffers

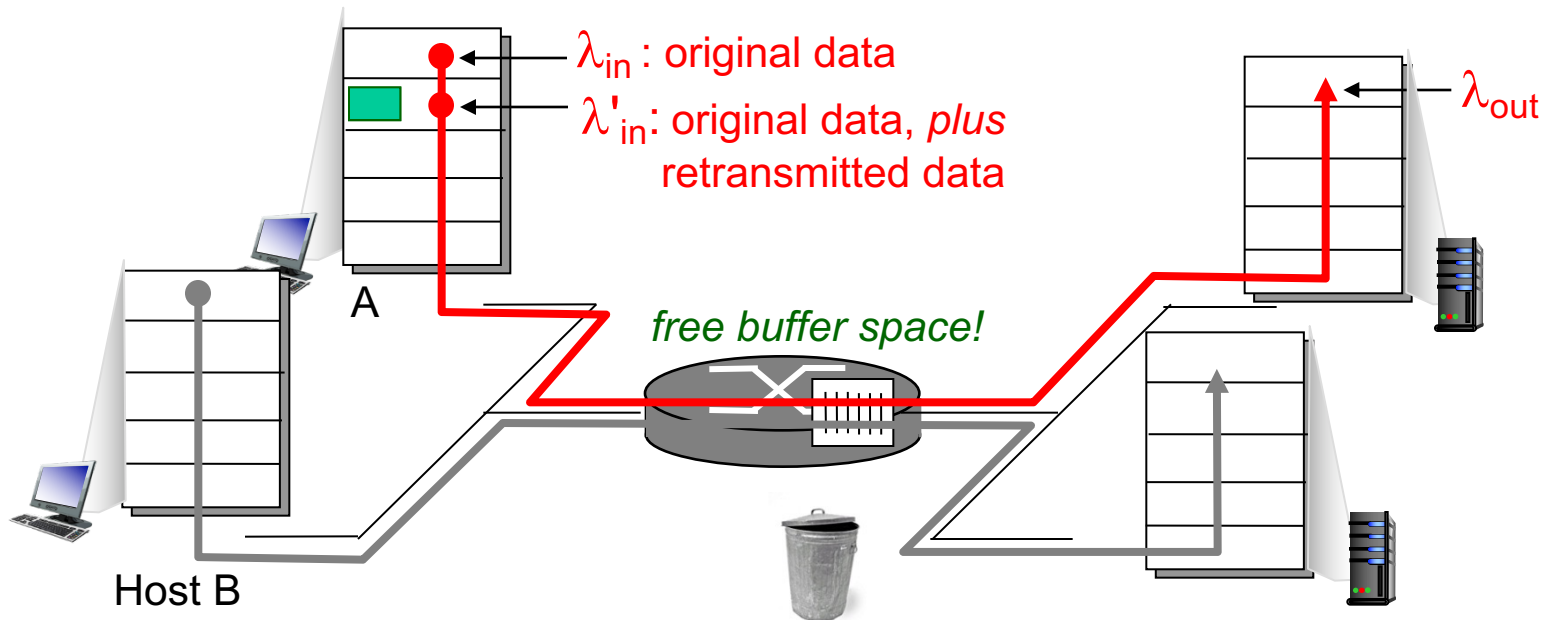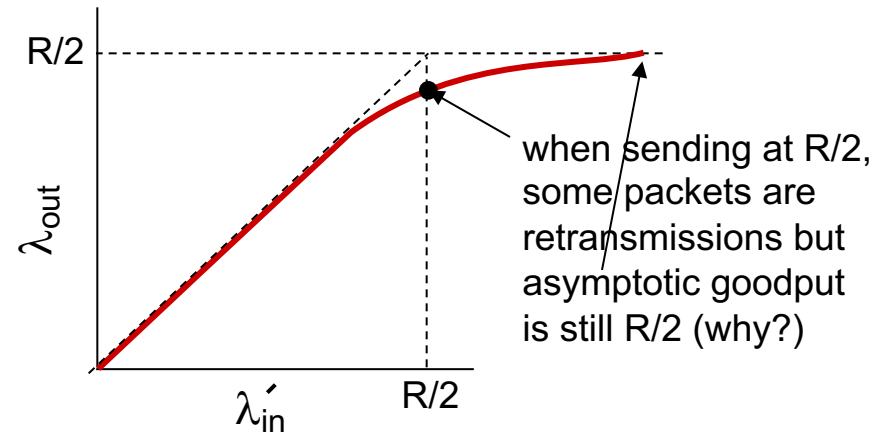❖ sender only resends if packet *known* to be lost



copy

$\lambda_{in}$ : original data

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

$\lambda_{out}$

A

*no buffer space!*

Host B

# Causes/costs of congestion: scenario 2

*Idealization: known loss*
packets can be lost, dropped at router due to full buffers

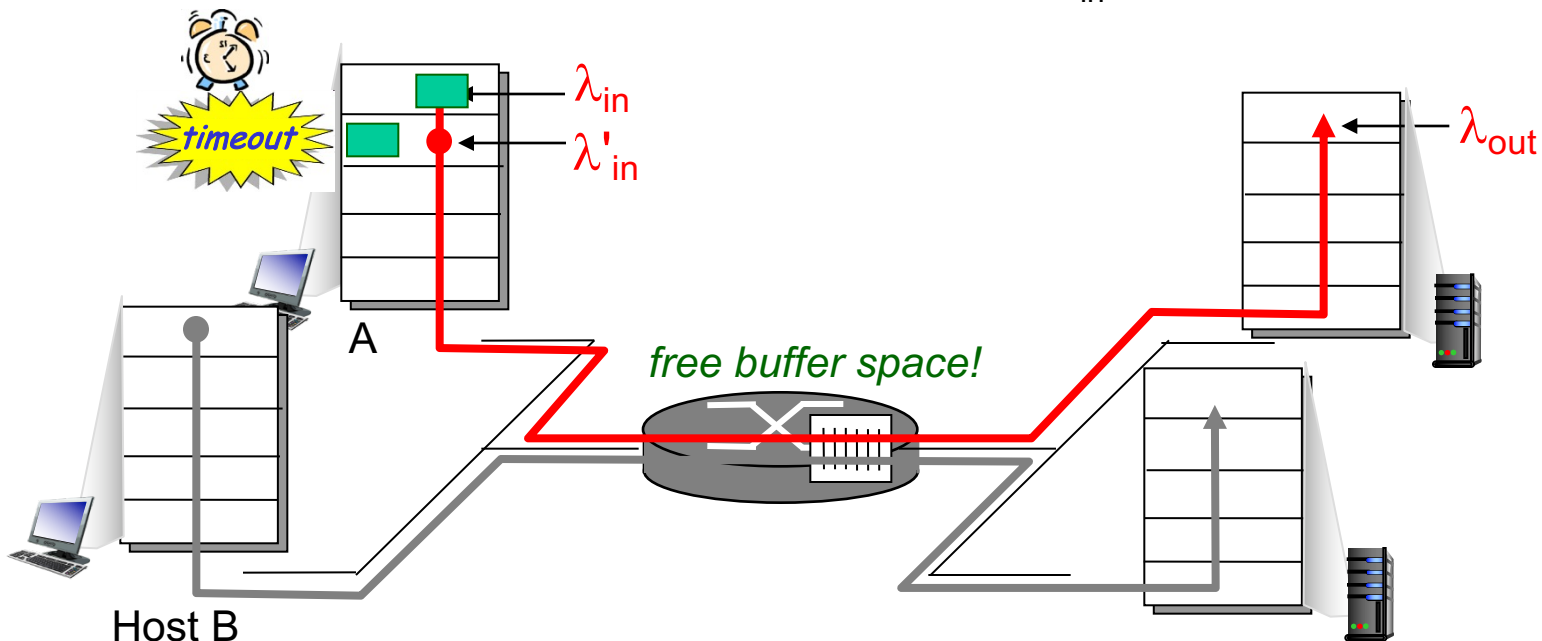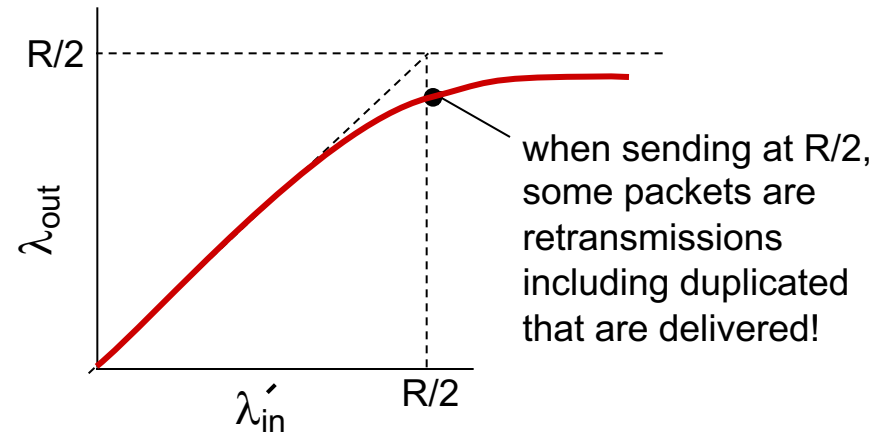❖ sender only resends if packet *known* to be lost



when sending at R/2, some packets are retransmissions but asymptotic goodput is still R/2 (why?)

$\lambda_{in}$ : original data

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

$\lambda_{out}$

*free buffer space!*

A

Host B

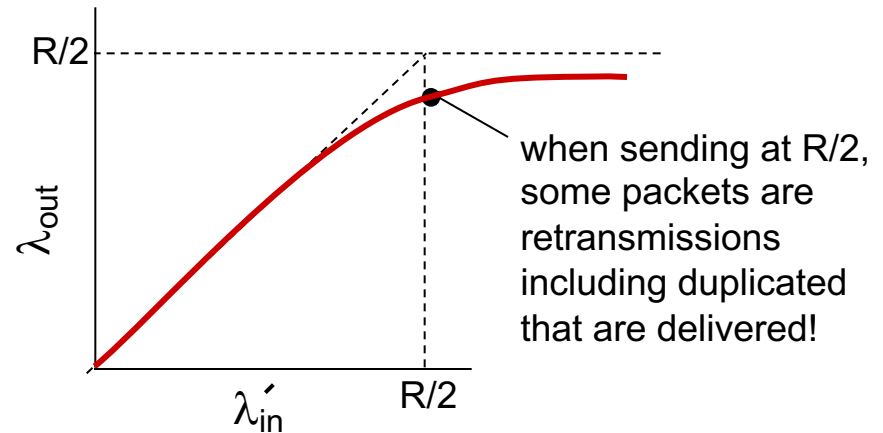# Causes/costs of congestion: scenario 2

## *Realistic: duplicates*

❖ packets can be lost, dropped at router due to full buffers

❖ sender times out prematurely, sending *two* copies, both of which are delivered

when sending at R/2, some packets are retransmissions including duplicated that are delivered!

# Causes/costs of congestion: scenario 2

*Realistic: duplicates*

❖ packets can be lost, dropped at router due to full buffers

❖ sender times out prematurely, sending *two* copies, both of which are delivered



when sending at R/2, some packets are retransmissions including duplicated that are delivered!
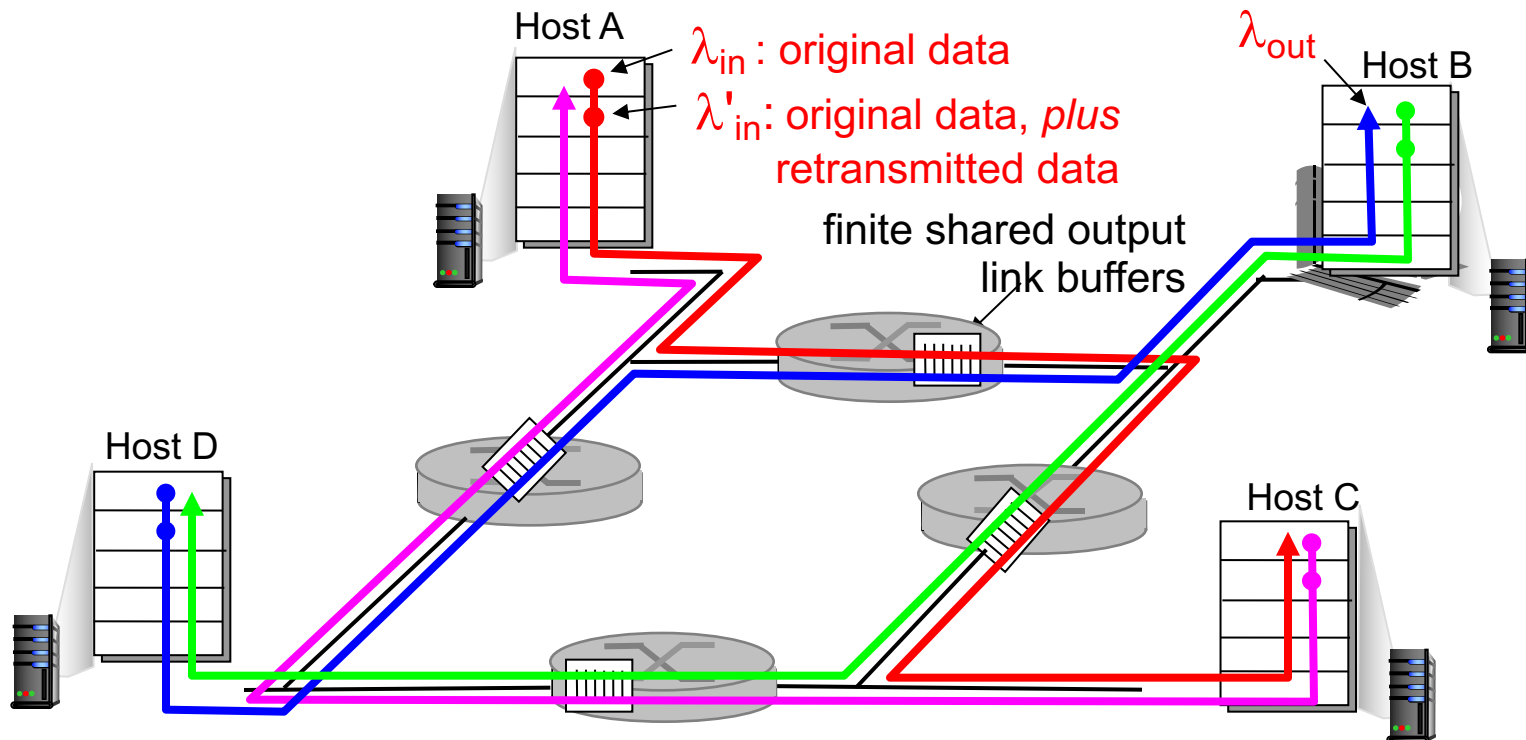
"costs" of congestion:

❖ more work (retrans) for given "goodput"

❖ unneeded retransmissions: link carries multiple copies of pkt
  ▪ decreasing goodput

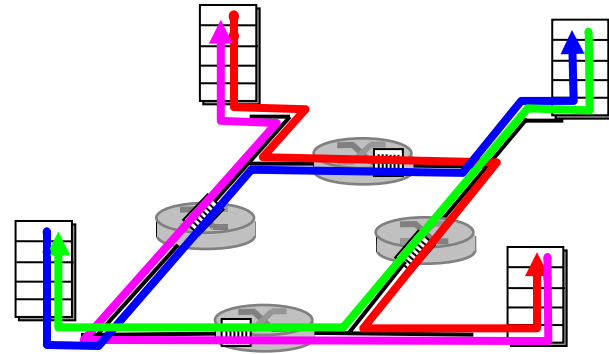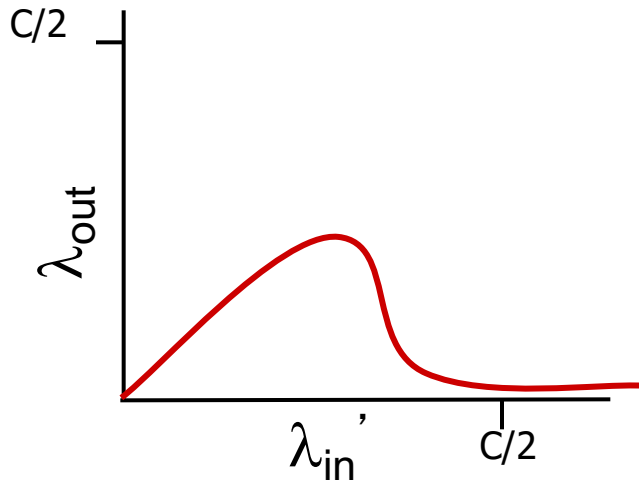# Causes/costs of congestion: scenario 3

❖ four senders
❖ multihop paths
❖ timeout/retransmit

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

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

$\lambda_{in}$ : original data

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

finite shared output link buffers

Host A

Host B

Host C

Host D

$\lambda_{out}$

# Causes/costs of congestion: scenario 3



another "cost" of congestion:

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

# Approaches towards congestion control

two broad approaches towards congestion control:

## end-end congestion control:

- ❖ 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
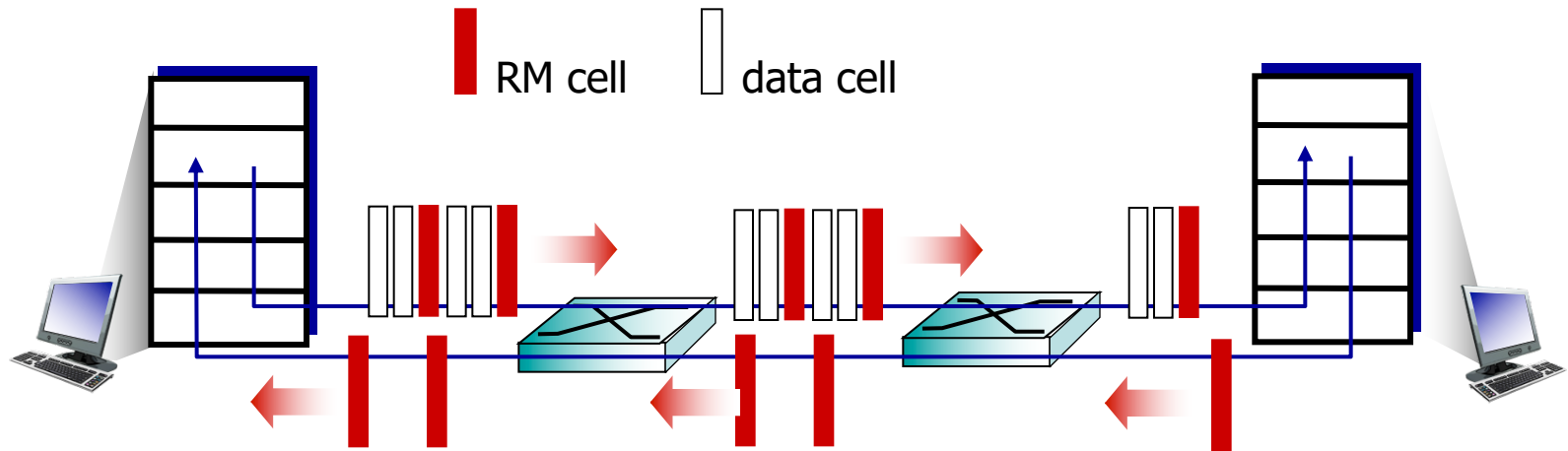
# Case study: ATM ABR congestion control

## ABR: available bit rate:

❖ "elastic service"
❖ if sender's path "underloaded":
  ▪ sender should use available bandwidth
❖ if sender's path congested:
  ▪ sender throttled to minimum guaranteed rate

## RM (resource management) cells:

❖ sent by sender, interspersed with data cells
❖ bits in RM cell set by switches ("*network-assisted*")
  ▪ *NI bit:* no increase in rate (mild congestion)
  ▪ *CI bit:* congestion indication
❖ RM cells returned to sender by receiver, with bits intact

# Case study: ATM ABR congestion control



RM cell    data cell

- ❖ **two-byte ER (explicit rate) field in RM cell**
  - congested switch may lower ER value in cell
  - senders' send rate thus max supportable rate on path
- ❖ **EFCI bit in data cells: set to 1 in congested switch**
  - if data cell preceding RM cell has EFCI set, receiver sets CI bit in returned RM cell

# 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
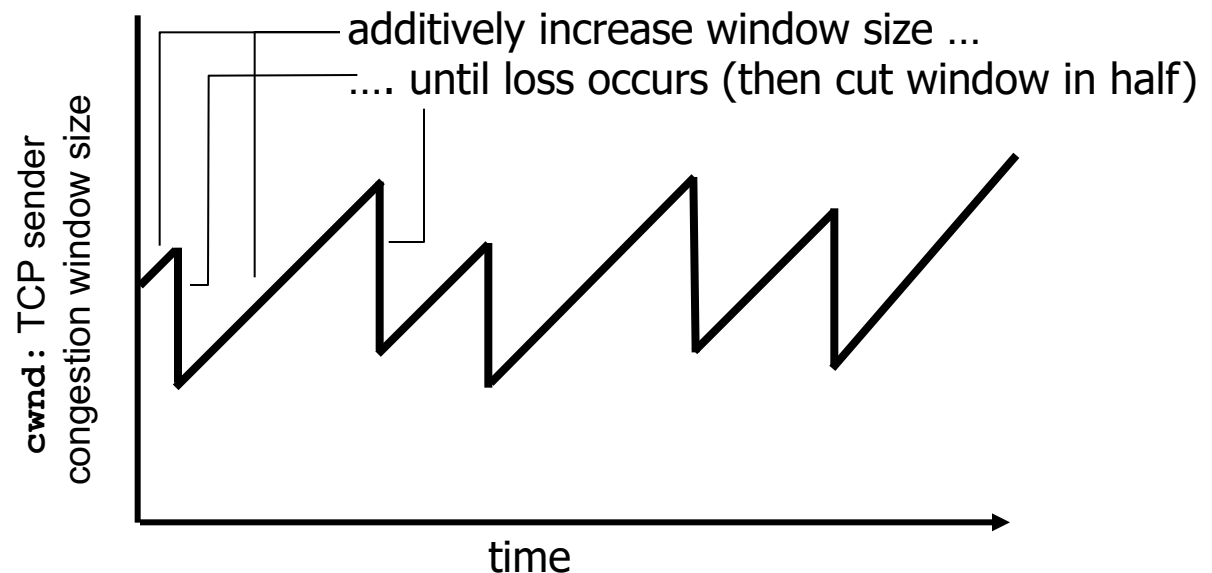- reliable data transfer
- flow control
- connection management

3.6 principles of congestion control
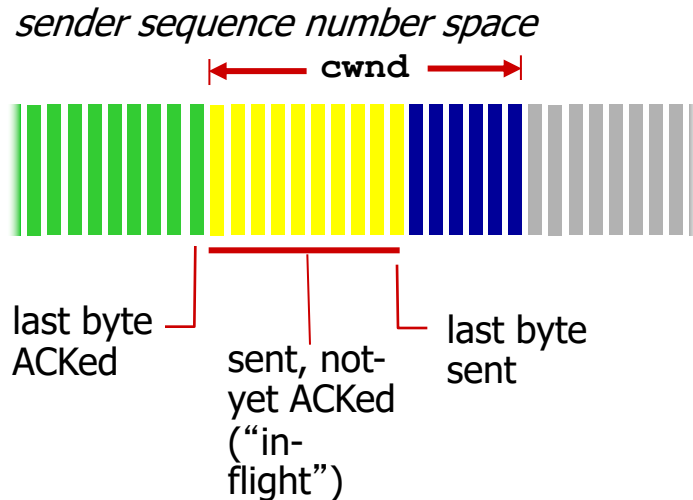
3.7 TCP congestion control

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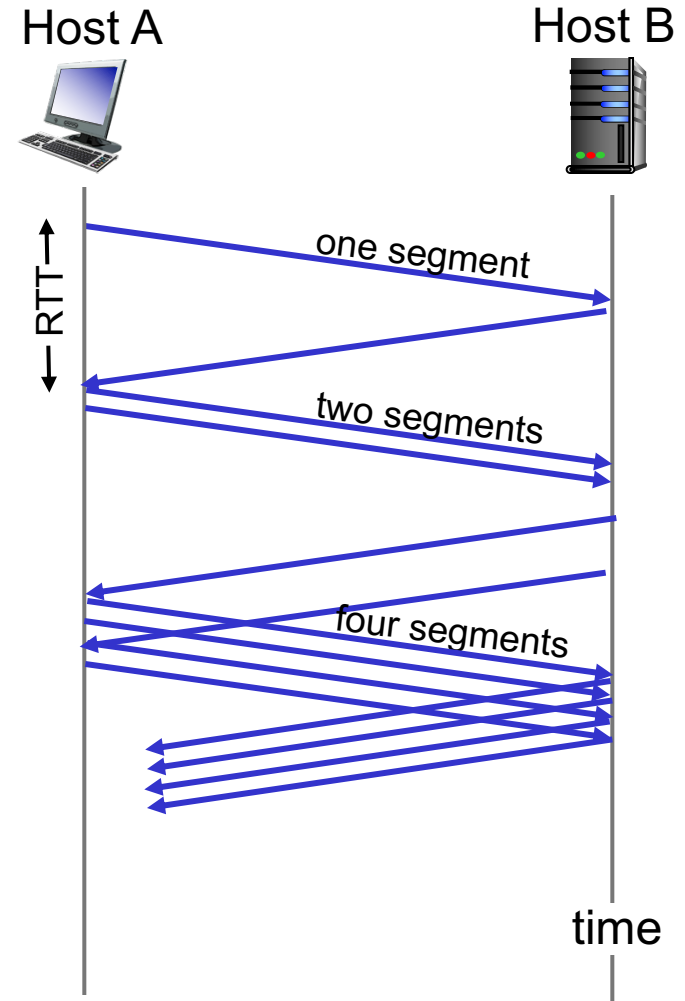
❖ *roughly:* send cwnd bytes, wait RTT for ACKS, then send more bytes

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

# TCP 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

# 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)
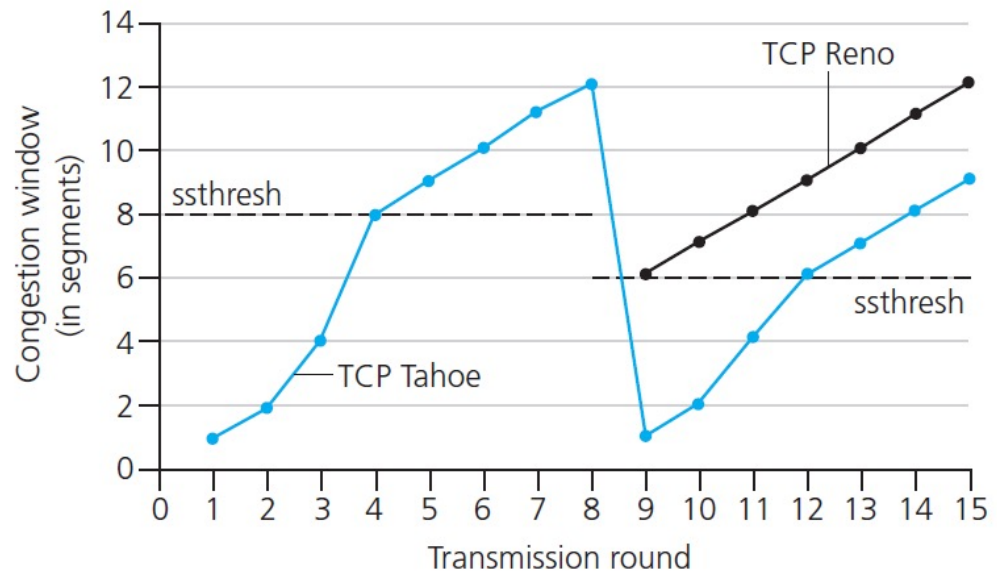
# TCP: switching from slow start to CA

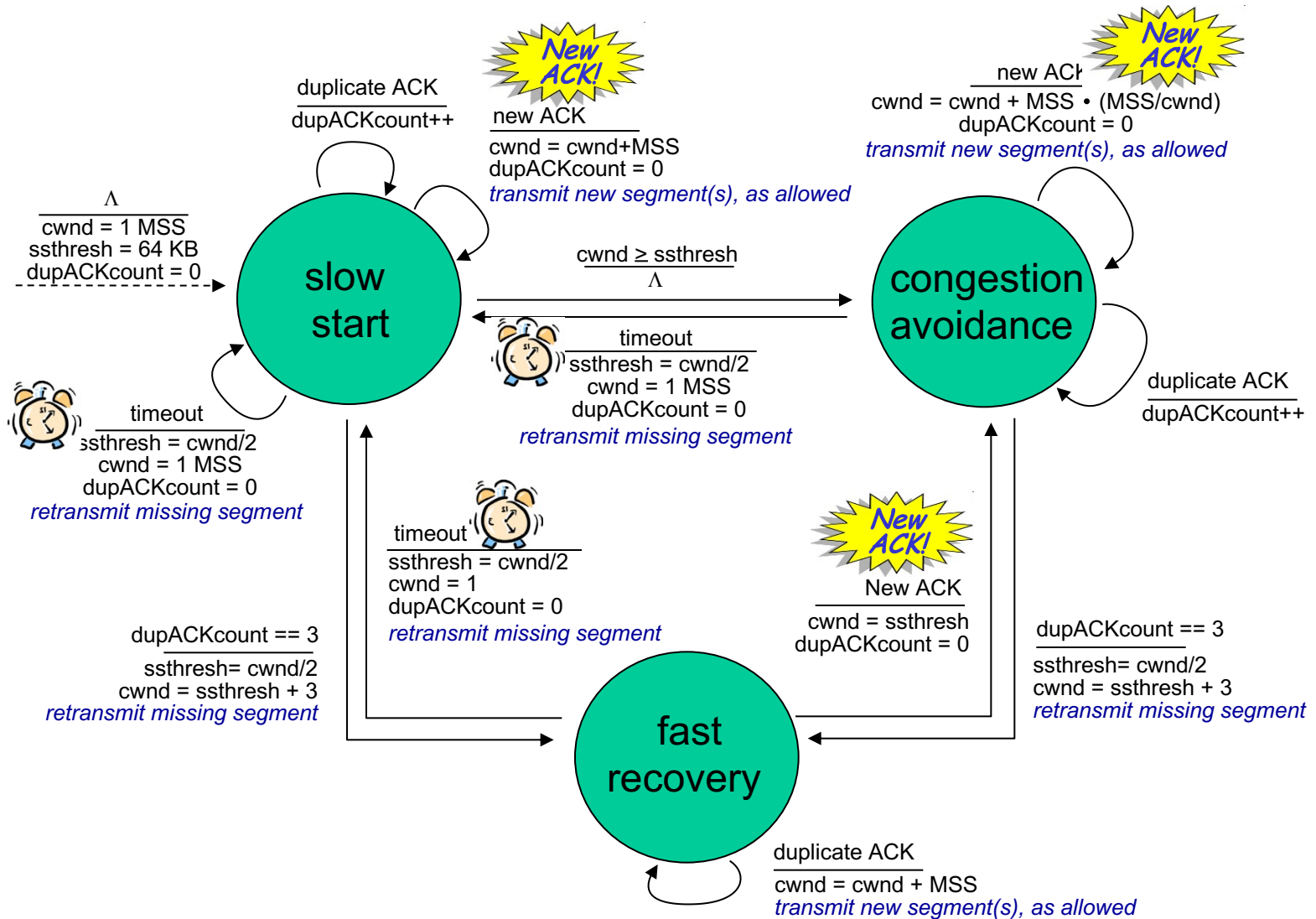Q: when should the exponential increase switch to linear?

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

## Implementation:

❖ variable **ssthresh**

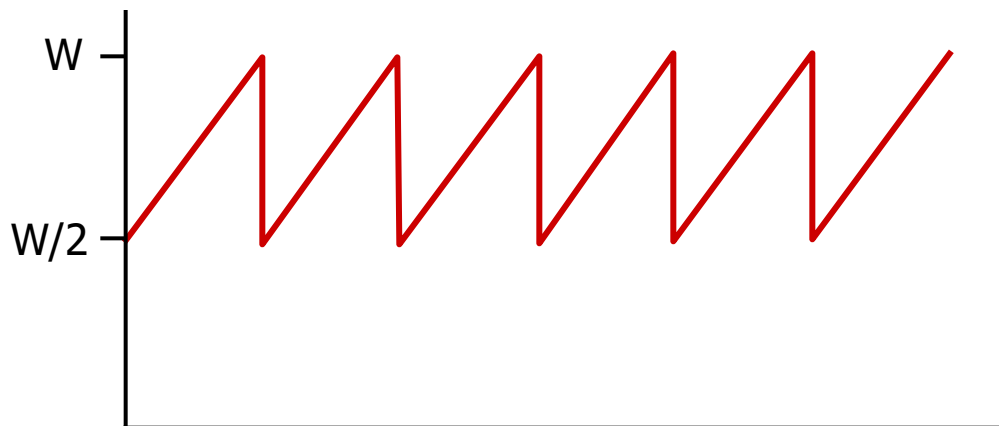❖ on loss event, **ssthresh** is set to 1/2 of **cwnd** just before loss event

# Summary: TCP Congestion Control

# TCP throughput

❖ avg. TCP thruput as function of window size, RTT?
  ▪ ignore slow start, assume always data to send
❖ W: window size (measured in bytes) where loss occurs
  ▪ avg. window size (# in-flight bytes) is ¾ W
  ▪ avg. thruput is 3/4W per RTT

$$\text{avg TCP thruput} = \frac{3}{4}\frac{W}{RTT} \text{ bytes/sec}$$

# TCP Futures: TCP over "long, fat pipes"

❖ example: 1500 byte segments, 100ms RTT, want 10 Gbps throughput

❖ requires W = 83,333 in-flight segments

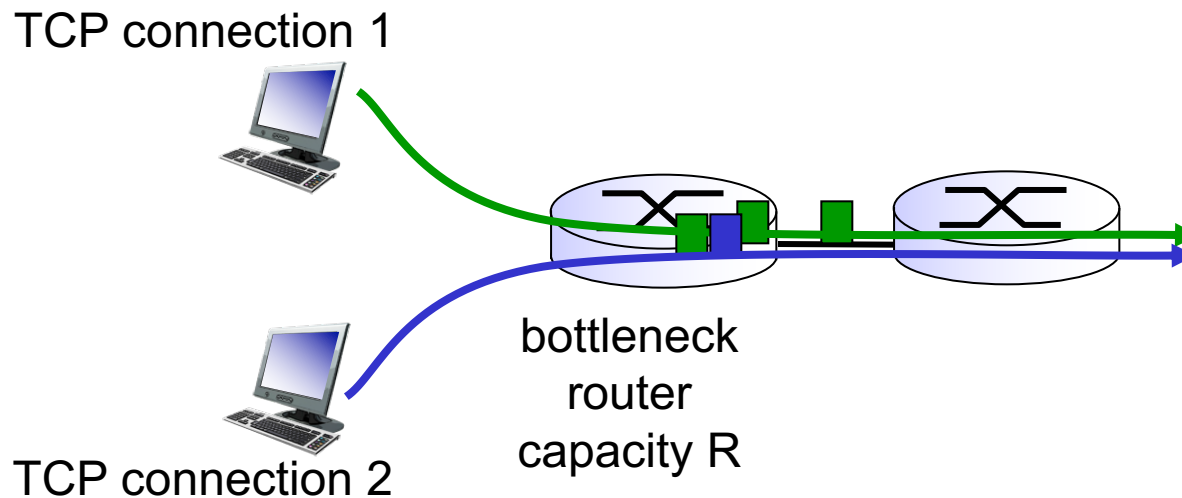❖ throughput in terms of segment loss probability, L [Mathis 1997]:

$$\text{TCP throughput} = \frac{1.22 \cdot MSS}{RTT \sqrt{L}}$$

➜ to achieve 10 Gbps throughput, need a loss rate of L = $2 \cdot 10^{-10}$   *– a very small loss rate!*

❖ new versions of TCP for high-speed
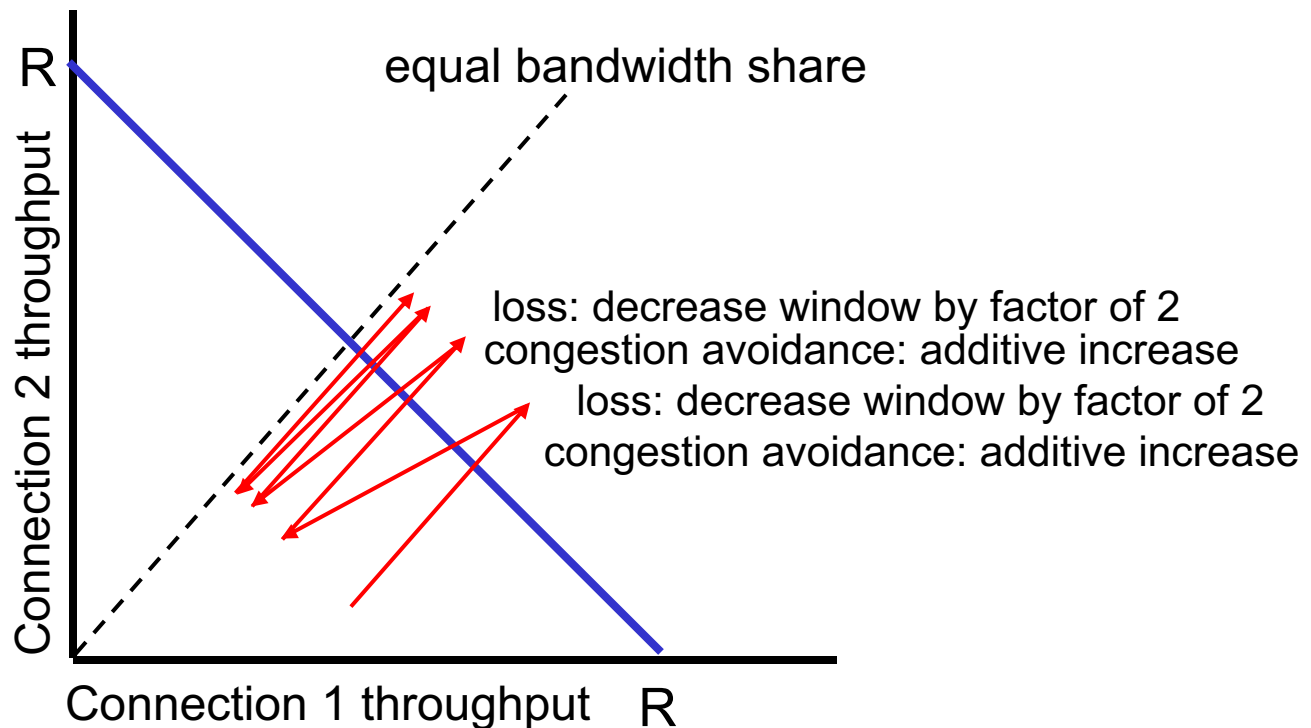
# TCP Fairness

*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 throughout increases
❖ multiplicative decrease decreases throughput proportionally



equal bandwidth share

Connection 2 throughput

Connection 1 throughput    R

R

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

# Fairness (more)

## Fairness and UDP

- multimedia apps often do not use TCP
  - do not want rate throttled by congestion control
- instead use UDP:
  - send audio/video at constant rate, tolerate packet loss

## Fairness, parallel TCP connections

- application can open multiple parallel connections between two hosts
- web browsers do this
- e.g., link of rate R with 9 existing connections:
  - new app asks for 1 TCP, gets rate R/10
  - new app asks for 11 TCPs, gets R/2

# Chapter 3: summary

❖ principles behind transport layer services:
- multiplexing, demultiplexing
- reliable data transfer
- flow control
- congestion control

❖ instantiation, implementation in the Internet
- UDP
- TCP

next:

❖ leaving the network "edge" (application, transport layers)

❖ into the network "core"