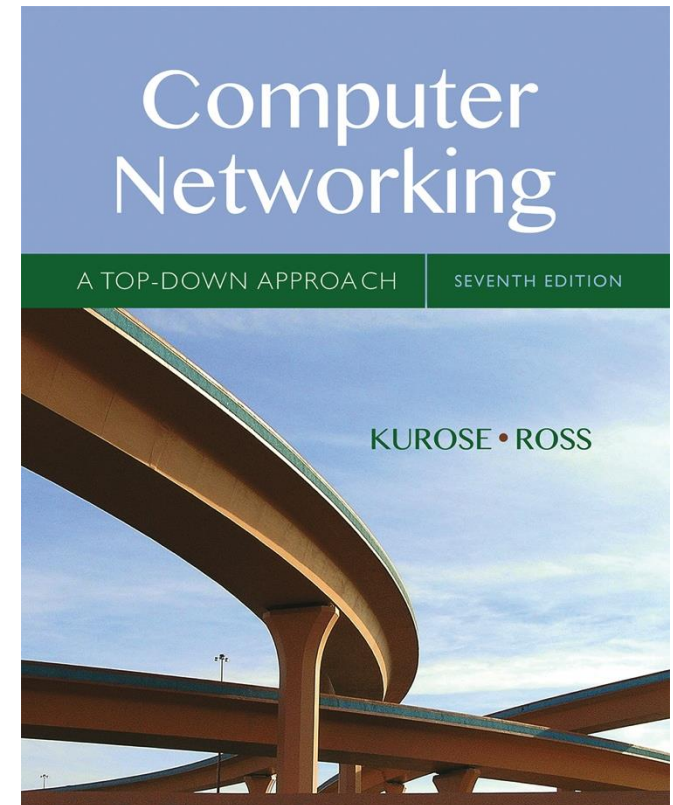# Chapter 6
# Transport Layer

A note on the use of these Powerpoint slides:

The notes used in this course are substantially based on Powerpoint slides developed and copyrighted by J.F. Kurose and K.W. Ross, 1996-2016

*Computer Networking: A Top Down Approach*

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

# Chapter 6: Transport Layer

## our goals:

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

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

# Chapter 6 outline

# 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 vs. network layer

- *network layer:* logical communication between hosts

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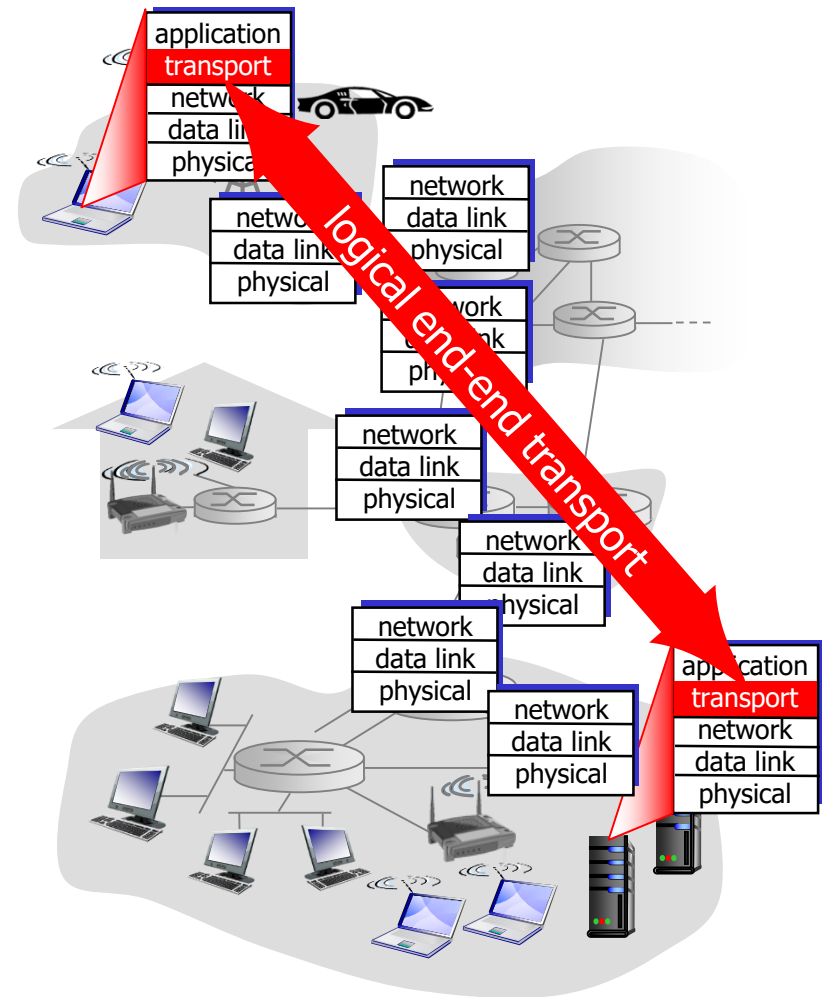  - relies on, enhances, network layer services

*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 6 outline

6.1 transport-layer services

6.2 multiplexing and demultiplexing

6.3 connectionless transport: UDP

6.4 principles of reliable data transfer

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

6.6 principles of congestion control

6.7 TCP congestion control
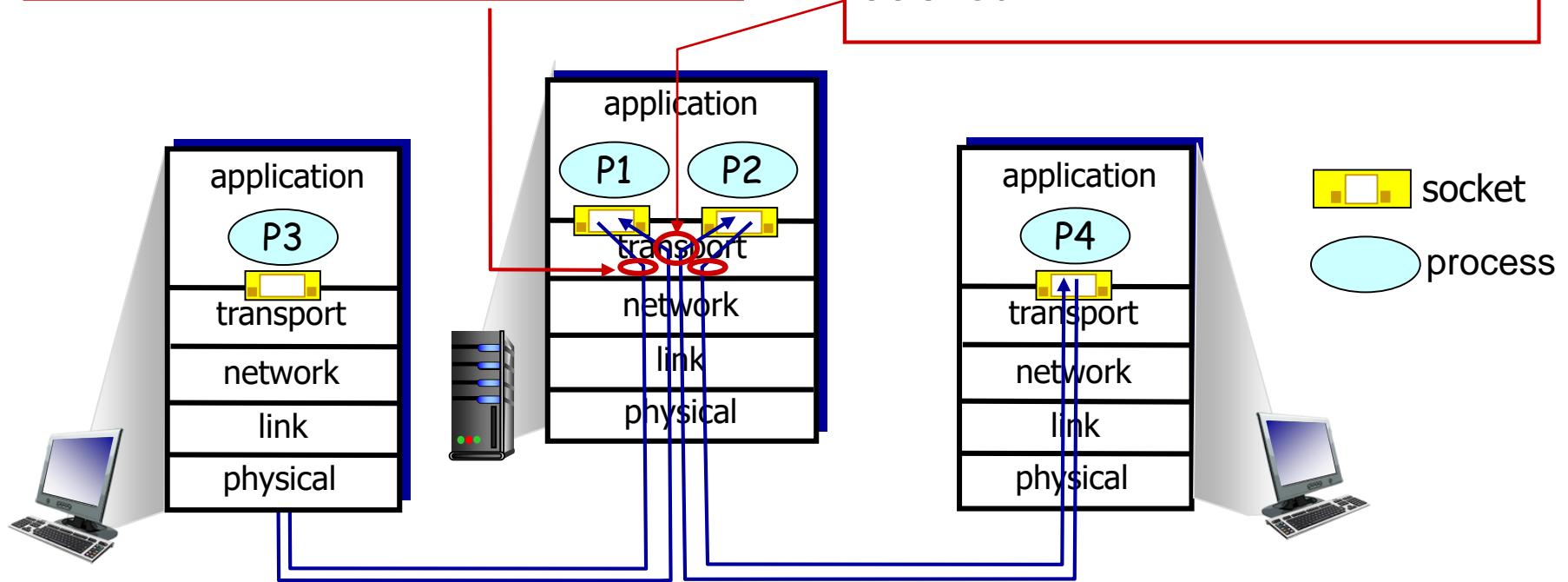
# Multiplexing/demultiplexing

*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

# Connectionless demultiplexing

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

```
clientSocket =
socket(AF_INET, SOCK_DGRAM)
clientSocket.bind(('',
10000))
```

- *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

`serverSocket.bind(('', 6428))`

`clientSocket.bind(('', 9157))`

`clientSocket.bind(('', 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

host: IP
address C

server: IP
address B

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

Transport Layer 6-13

# Connection-oriented demux: example

threaded server

application

P3

transport

network

link

physical

host: IP
address A

application

P4

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

# Chapter 6 outline

# 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

length, in bytes of UDP segment, including header

| 32 bits |
| source port # | dest port # |
| length | checksum |
| application data (payload) |

UDP segment format

# UDP: User Datagram Protocol [RFC 768]

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

# 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: One's (1s) complement of the 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

# Exercise

a. Suppose you have the following 2 bytes: 01011100 and 01100101. What is the 1s complement of the sum of these 2 bytes?

b. Suppose you have the following 2 bytes: 11011010 and 01100101. What is the 1s complement of the sum of these 2 bytes?

c. For the bytes in part (a), give an example where one bit is flipped in each of the 2 bytes and yet the 1s complement doesn't change.

# Chapter 6 outline

6.1 transport-layer services

6.2 multiplexing and demultiplexing

6.3 connectionless transport: UDP

6.4 principles of reliable data transfer

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

6.6 principles of congestion control

6.7 TCP congestion control

# Chapter 6 outline

6.1 transport-layer services

6.2 multiplexing and demultiplexing

6.3 connectionless transport: UDP

6.4 principles of reliable data transfer

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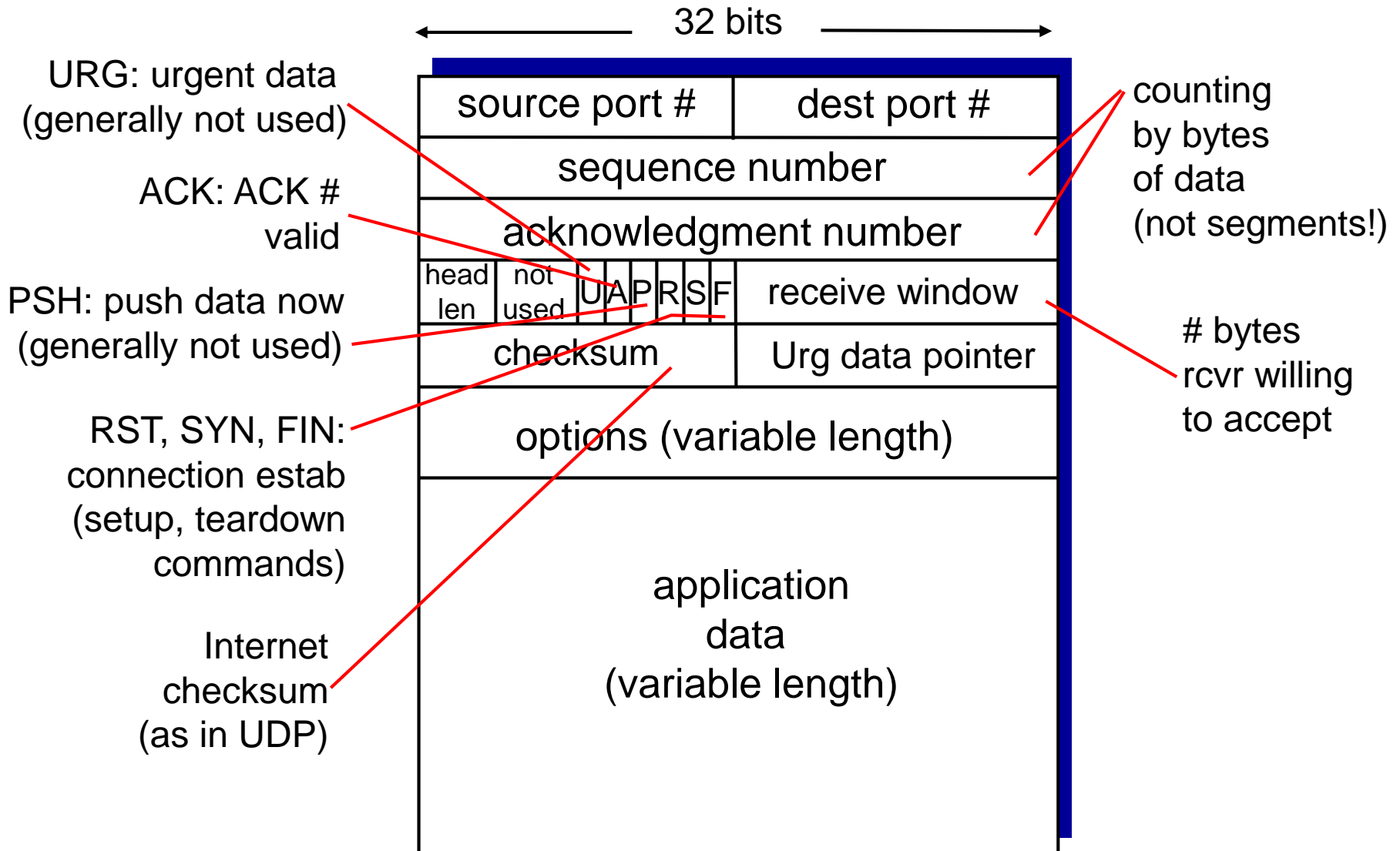
6.6 principles of congestion control

6.7 TCP congestion control

# 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 | |
| acknowledgment 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
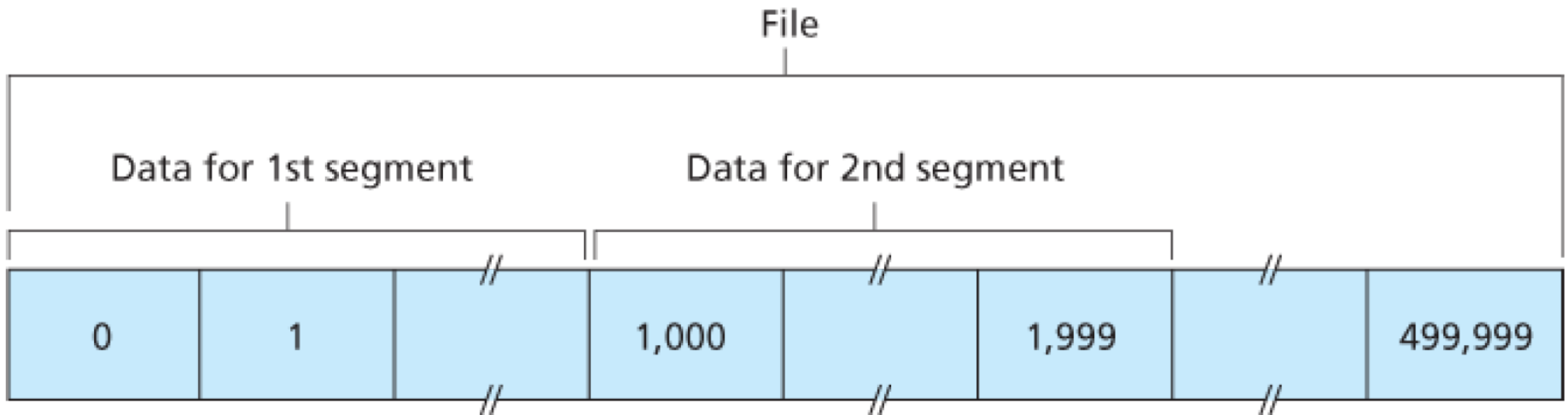
# TCP seq. numbers

<span style="color:red">sequence numbers:</span>

- byte stream "number" of first byte in segment's data
  - E.g., Suppose that a process in Host A wants to send a stream of data (a file of 500,000 bytes) to a process in Host B over a TCP connection. MSS is 1,000 bytes, and the 1st byte is numbered 0.

File

| 0 | 1 | // | 1,000 | // | 1,999 | // | 499,999 |

Data for 1st segment        Data for 2nd segment

- Seq # of the 1st segment: 0; Seq # of the 2nd segment: 1000;
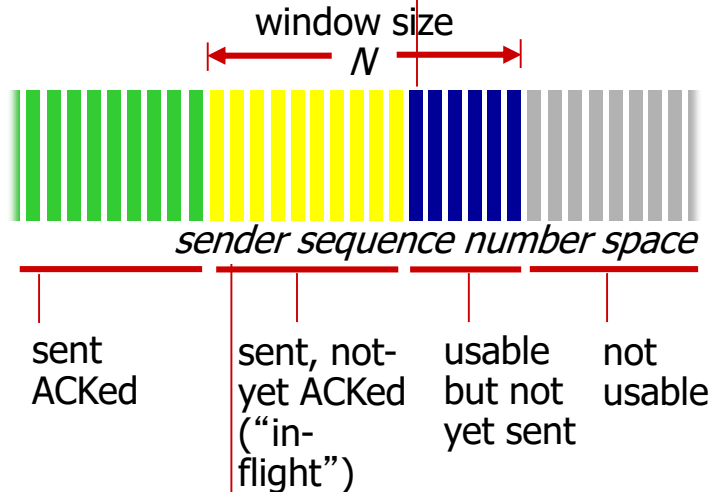- Seq # of the 3rd segment: ?

# TCP ACKs

acknowledgments:
- seq # of next byte expected from other side
  - Suppose that Host A has received all bytes numbered 0 through 535 from B and suppose that it is about to send a segment to Host B
  - Host A is waiting for byte 536 and all the subsequent bytes
  - So Host A puts 536 in the acknowledgment number field
- cumulative ACK
  - Suppose that Host A has received one segment from Host B containing bytes 0 ~ 535, and another segment containing bytes 900 ~ 1,000
  - For some reason Host A has not yet received bytes 536 ~ 899
  - Host A is still waiting for byte 536 (and beyond)
  - A's next segment to B will contain 536 in the acknowledgment number field

# TCP seq. numbers, ACKs

outgoing segment from sender

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

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

# 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
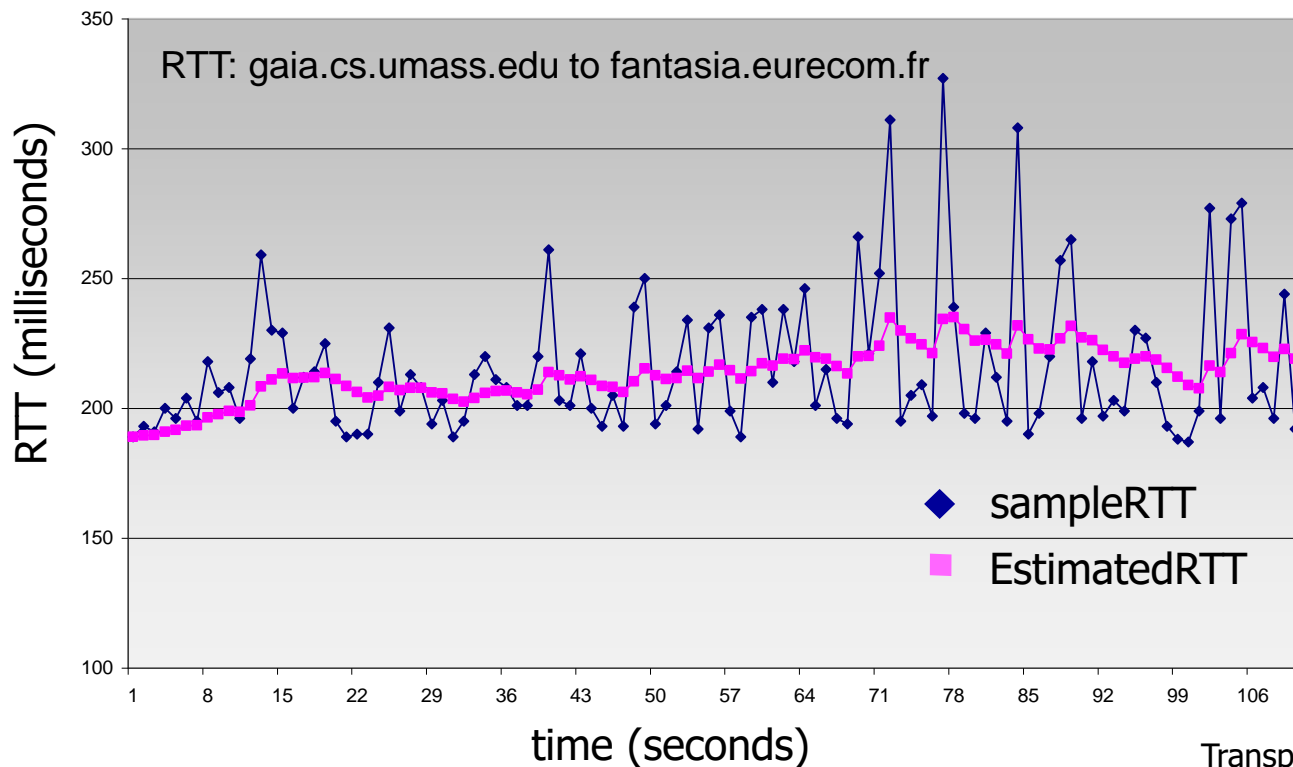- *too long:* slow reaction to segment loss

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

- ◆ sampleRTT
- ■ EstimatedRTT

RTT (milliseconds)

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-β)*DevRTT +
              β*|SampleRTT-EstimatedRTT|
         (typically, β = 0.25)
```

$$\mathtt{TimeoutInterval = EstimatedRTT + 4*DevRTT}$$

estimated RTT          "safety margin"

\* Check out the online interactive exercises for more
examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

# Chapter 6 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
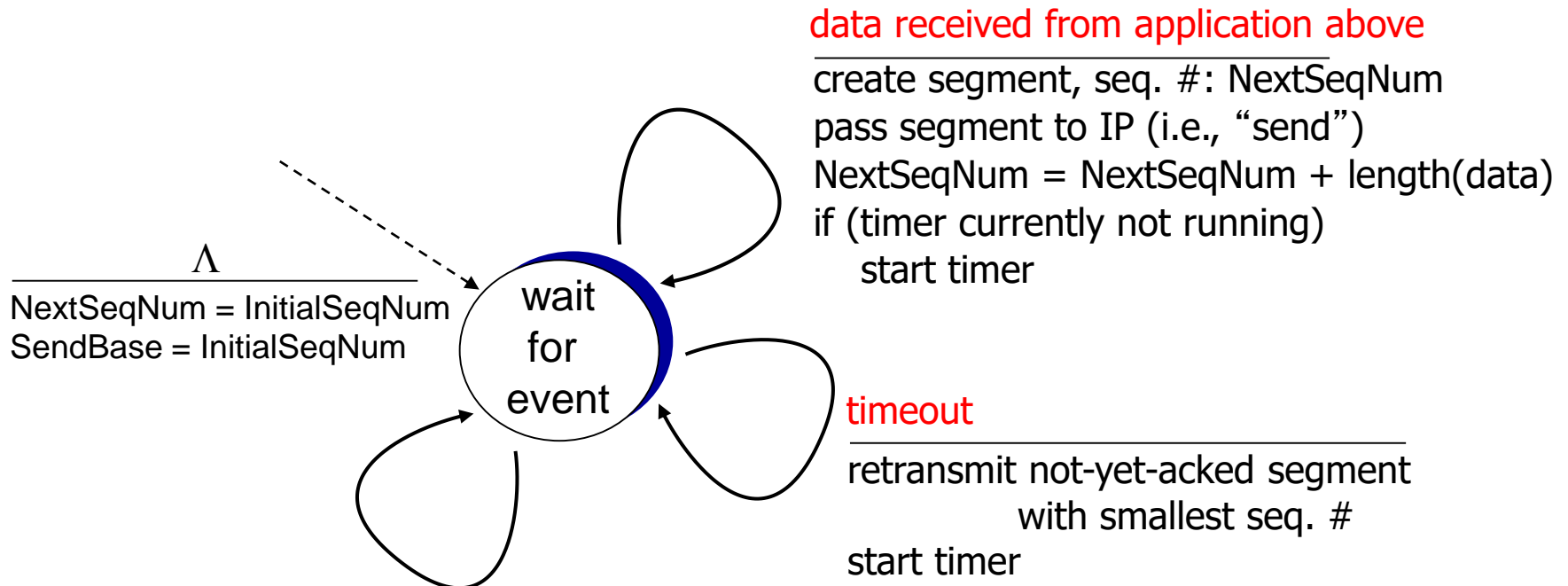  - expiration interval: `TimeOutInterval`

*timeout:*

- retransmit segment that caused timeout
- restart timer

*ack rcvd:*

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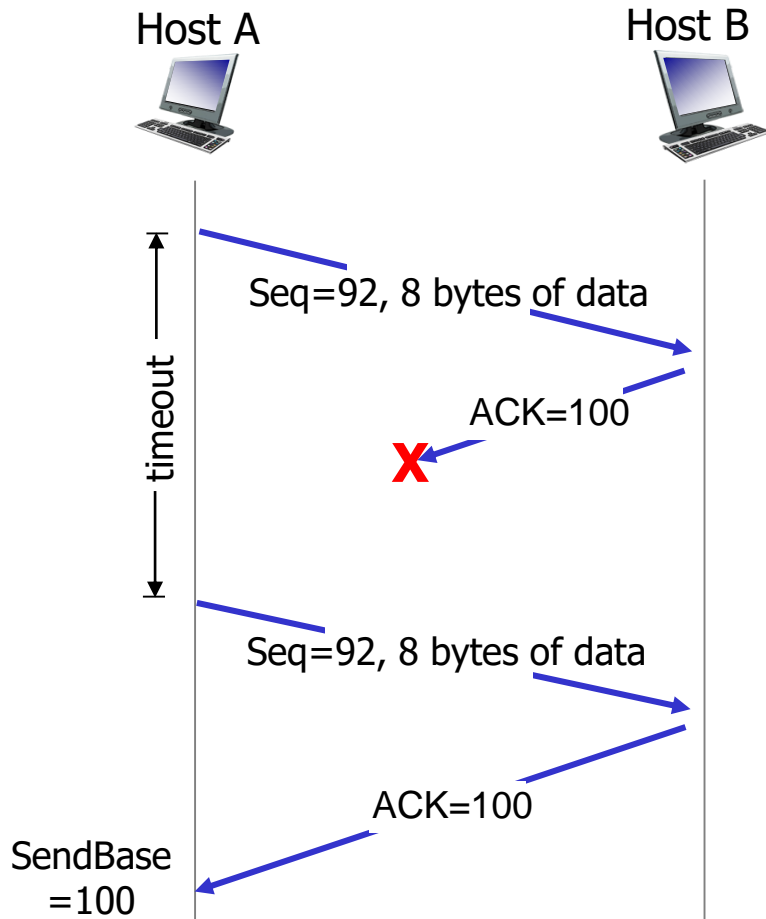  - start timer if there are still unacked segments
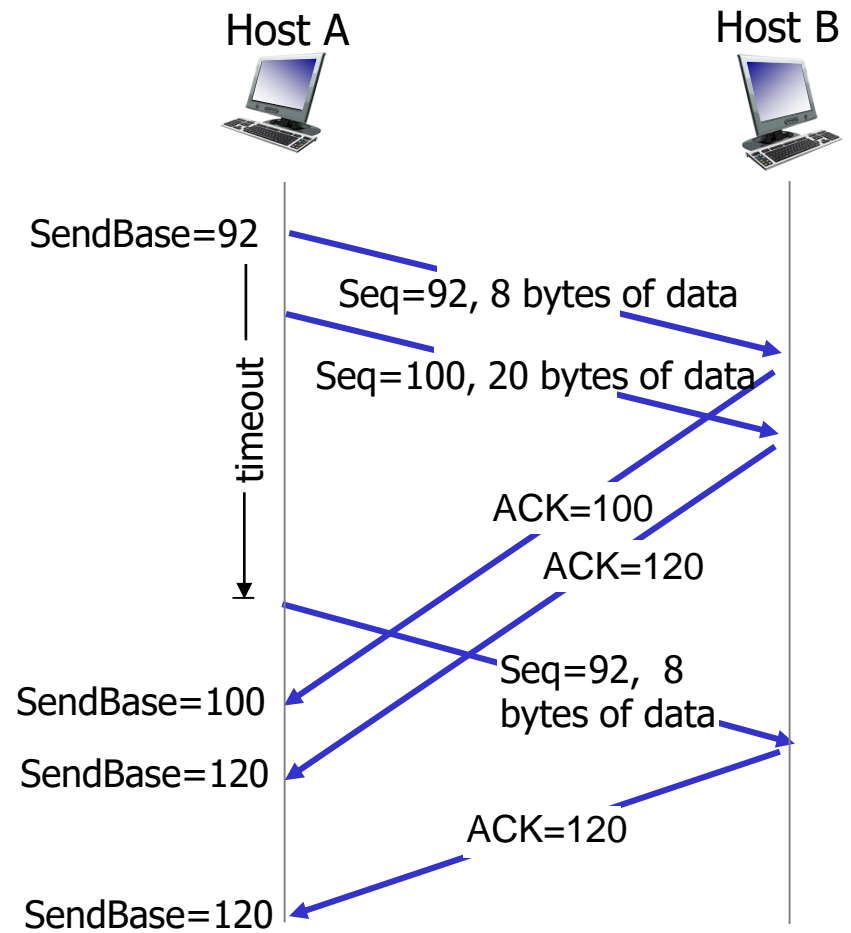
# TCP sender (simplified)

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  %cumulative acks
   /* 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

SendBase=92

Seq=92, 8 bytes of data

ACK=100

X

timeout

Seq=92, 8 bytes of data

ACK=100

SendBase
=100

lost ACK scenario

Host A                                    Host B

SendBase=92

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

timeout

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

Seq=120,  15 bytes of data

timeout

cumulative ACK

# TCP ACK generation [RFC 1122, RFC 5681]

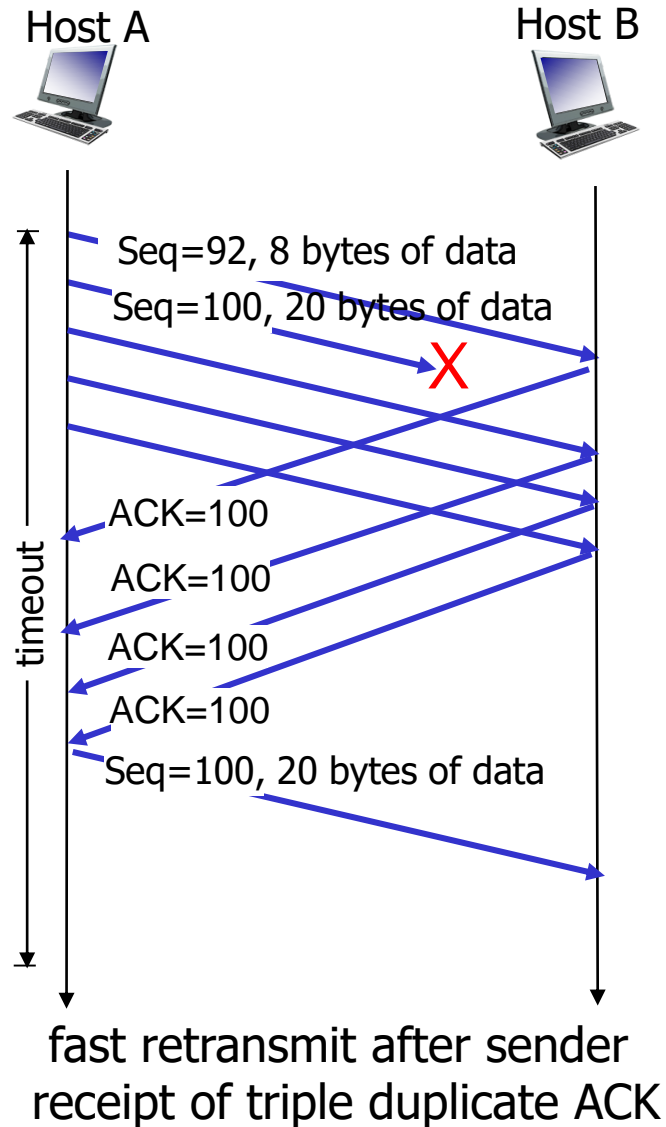| *event at receiver* | *TCP receiver action* |
|---|---|
| arrival of in-order segment with expected seq #. All data up to expected seq # already ACKed | delayed ACK. Wait up to 500ms for next segment. If no next segment, send ACK |
| arrival of in-order segment with expected seq #. One other segment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

ACK=100

ACK=100

Seq=100, 20 bytes of data

timeout

fast retransmit after sender
receipt of triple duplicate ACK

# Chapter 6 outline

6.1 transport-layer services

6.2 multiplexing and demultiplexing

6.3 connectionless transport: UDP

6.4 principles of reliable data transfer

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

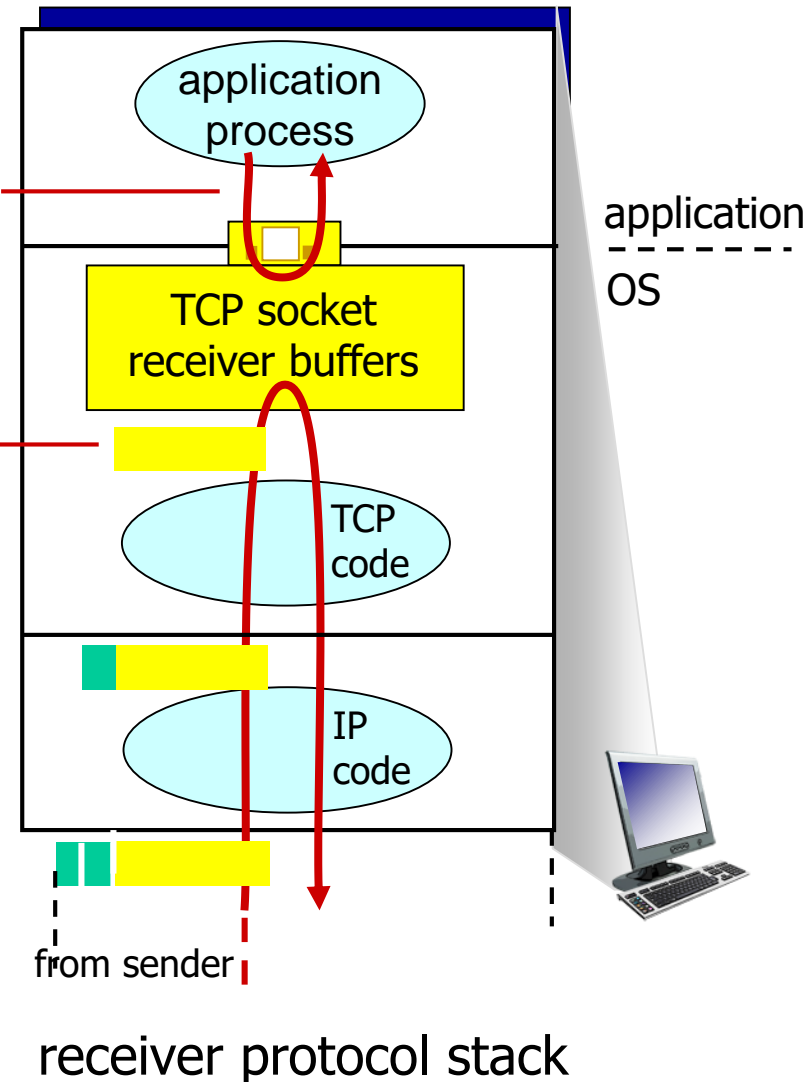6.6 principles of congestion control

6.7 TCP congestion control

# TCP flow control



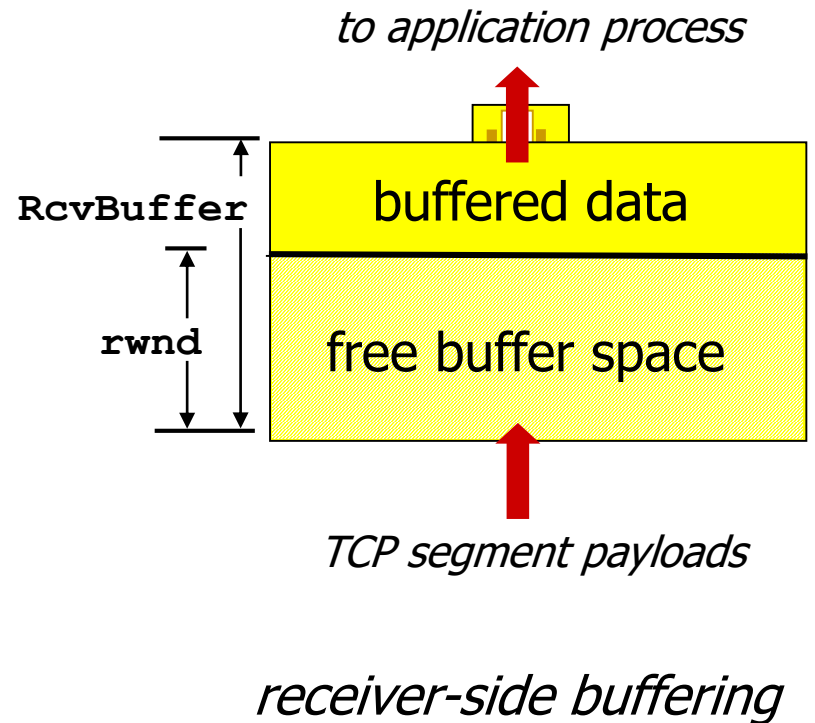application may remove data from TCP socket buffers ....

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

application process

TCP socket receiver buffers

TCP code

IP code

application
- - - - - - -
OS

from sender

receiver protocol stack

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

# 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*
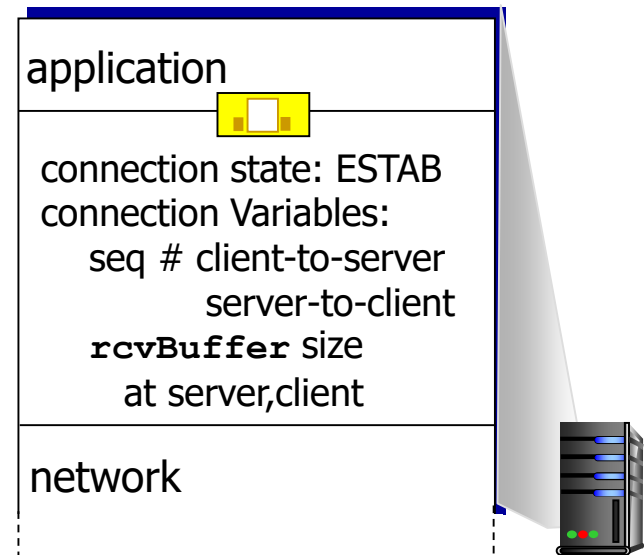
# Chapter 3 outline

# Connection Management

before exchanging data, sender/receiver "handshake":

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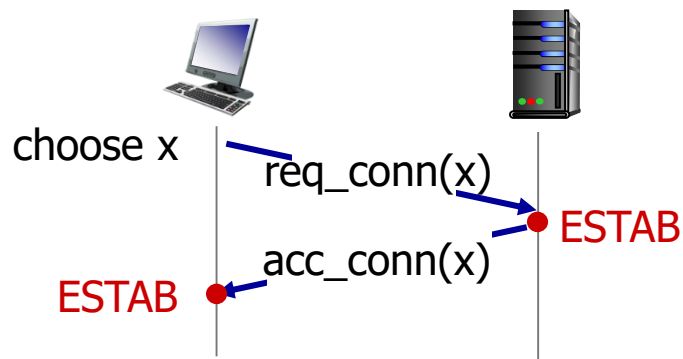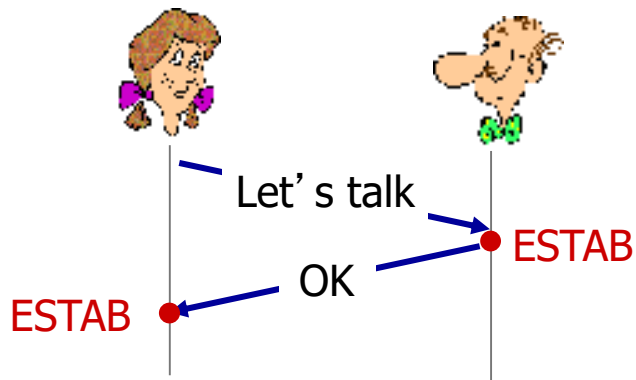- agree on connection parameters

application

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

network

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");
```

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

# Agreeing to establish a connection
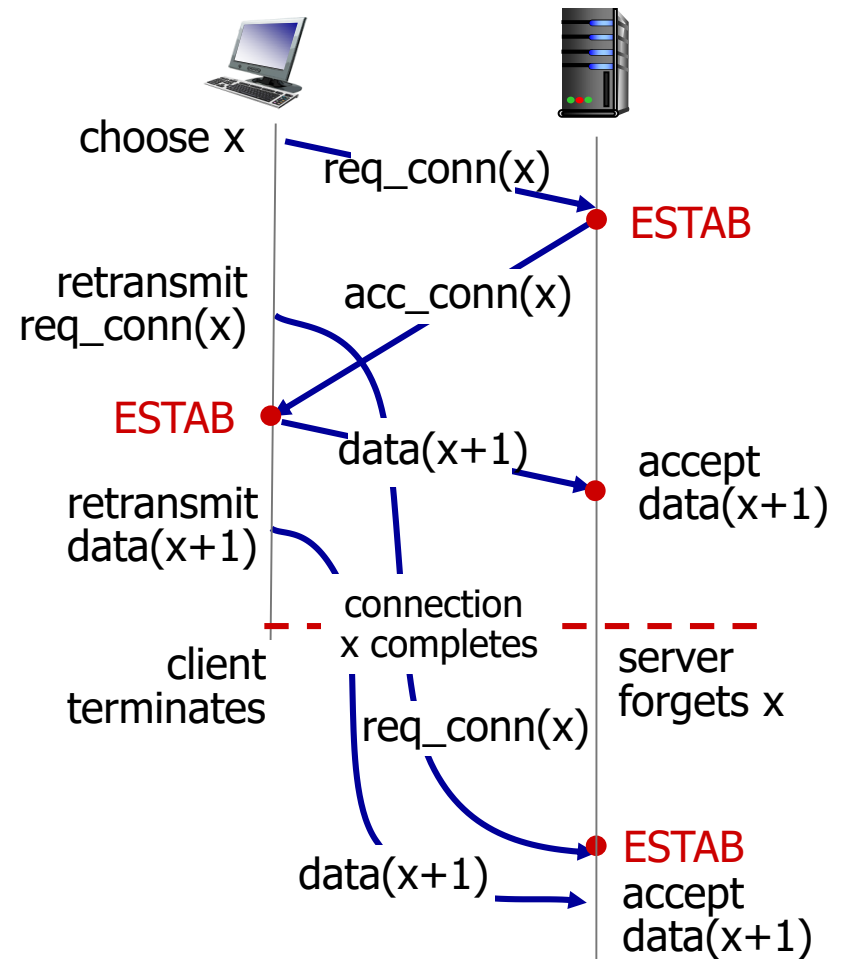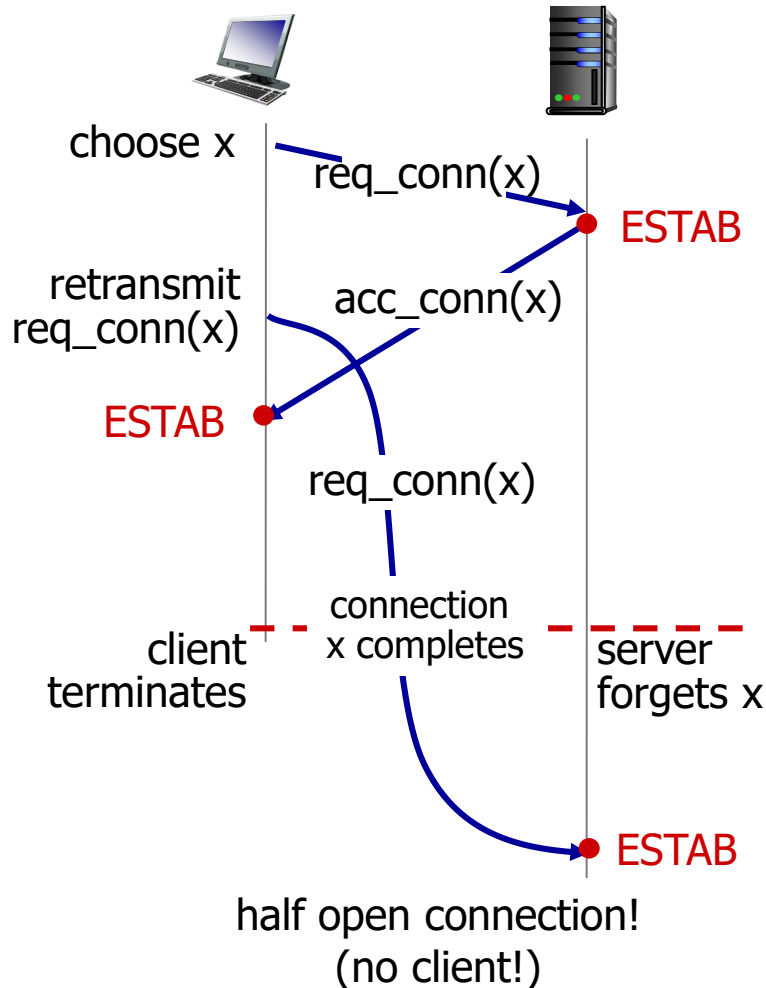
2-way handshake:



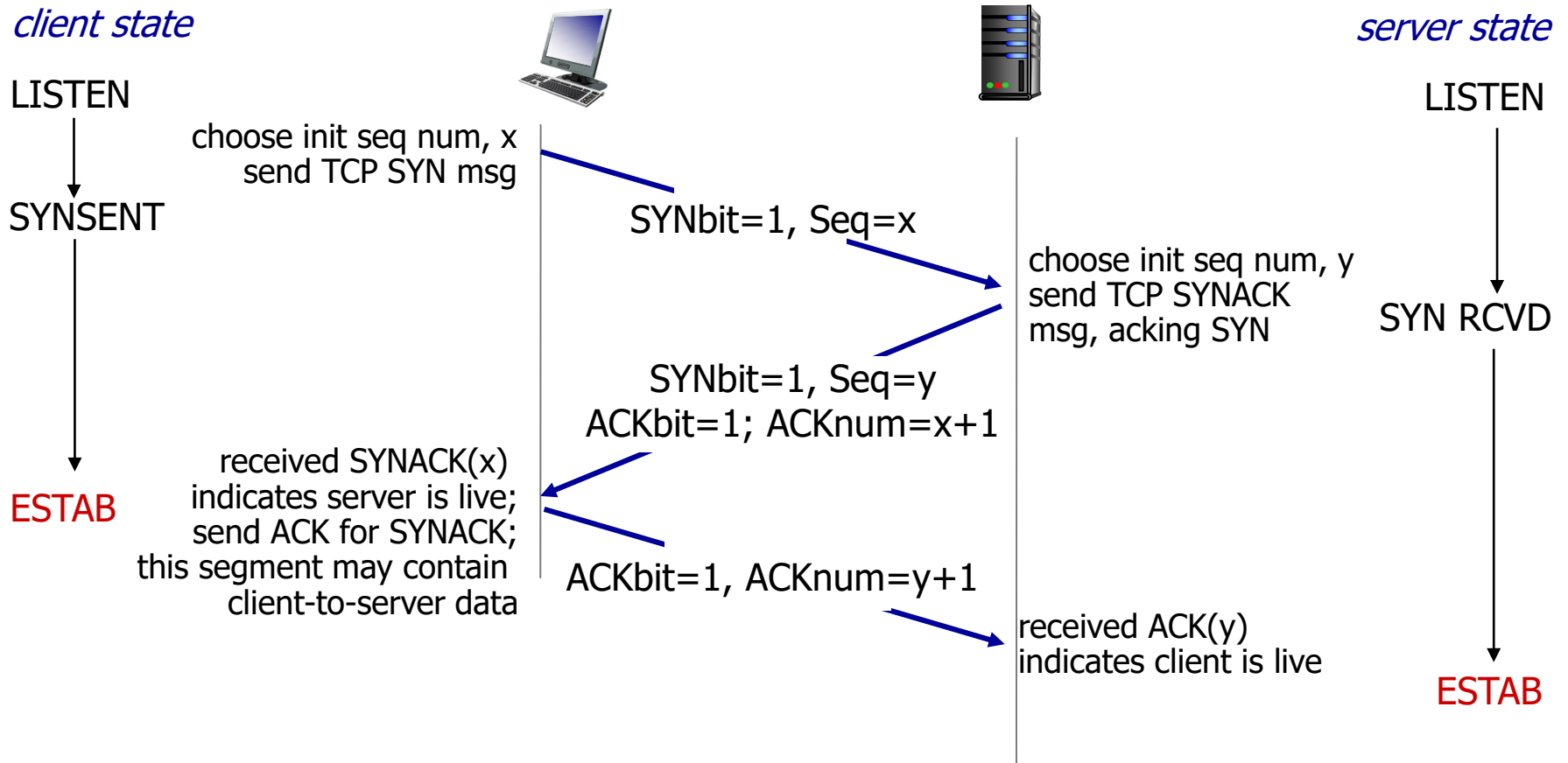*Q:* will 2-way handshake always work in network?

- 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:



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)

connection
x completes

client
terminates

req_conn(x)

server
forgets x

data(x+1)

ESTAB
accept
data(x+1)

# TCP 3-way handshake

client state

LISTEN

SYNSENT

ESTAB

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

choose init seq num, x
send TCP SYN msg

SYNbit=1, Seq=x

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

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

$$\frac{\texttt{Socket connectionSocket =}}{\texttt{welcomeSocket.accept();}}$$
$$\Lambda$$

$$\frac{\texttt{Socket clientSocket =}}{\texttt{newSocket("hostname","port}}{\texttt{ number");}}$$
SYN(seq=x)

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

listen

SYN rcvd

SYN sent

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

ESTAB

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

# 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

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

ESTAB

CLOSE_WAIT    can still
send data

LAST_ACK    can no longer
send data

CLOSED

# Sequences of TCP states

## TCP Client Lifecycle



Client application initiates a TCP connection

CLOSED

Send SYN

SYN_SENT

Receive SYN & ACK, send ACK

ESTABLISHED

Send FIN

Client application initiates close connection

FIN_WAIT_1

Receive ACK, send nothing

FIN_WAIT_2

Receive FIN, send ACK

TIME_WAIT

Wait 30 seconds

## TCP Server Lifecycle



Server application creates a listen socket

CLOSED

Receive ACK, send nothing

LAST_ACK

LISTEN

Receive SYN send SYN & ACK

SYN_RCVD

Receive ACK, send nothing

ESTABLISHED

Receive FIN, send ACK

CLOSE_WAIT

Send FIN

# Chapter 6 outline

6.1 transport-layer services

6.2 multiplexing and demultiplexing

6.3 connectionless transport: UDP

6.4 principles of reliable data transfer

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

6.6 principles of congestion control
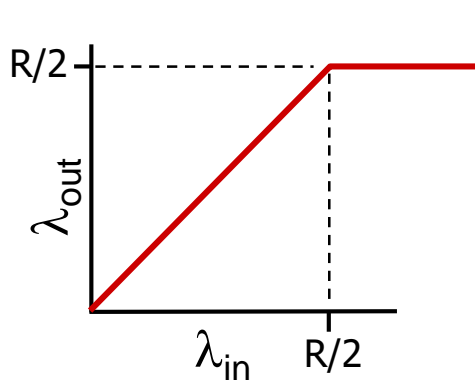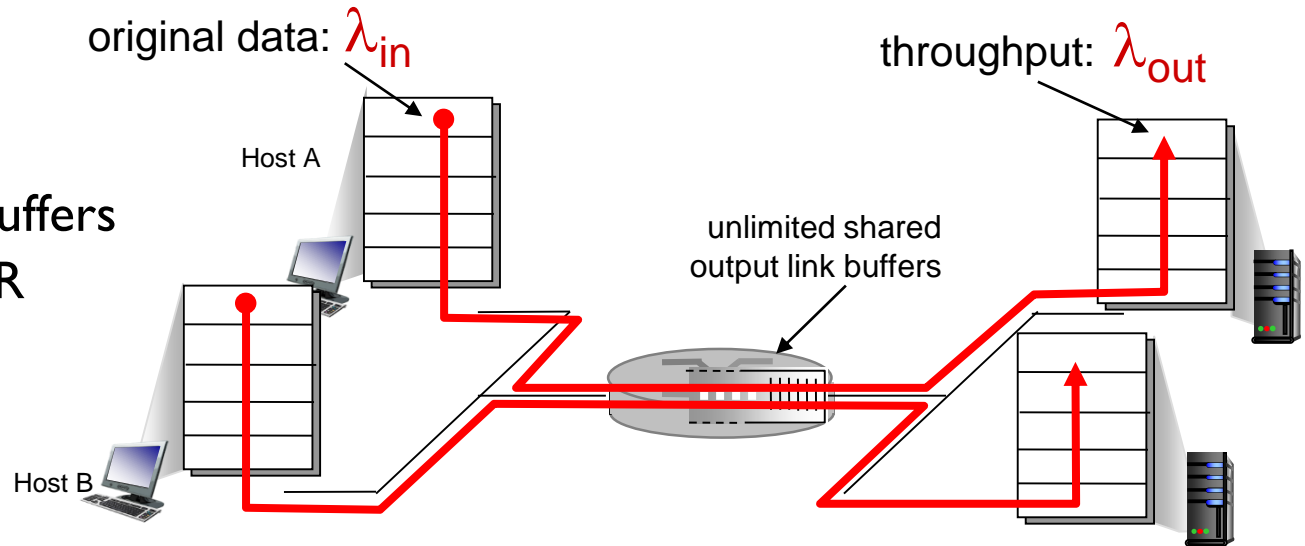
6.7 TCP congestion control
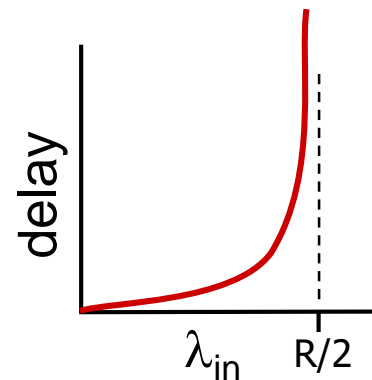
# Principles of congestion control

*congestion:*

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

- different from flow control!

- 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
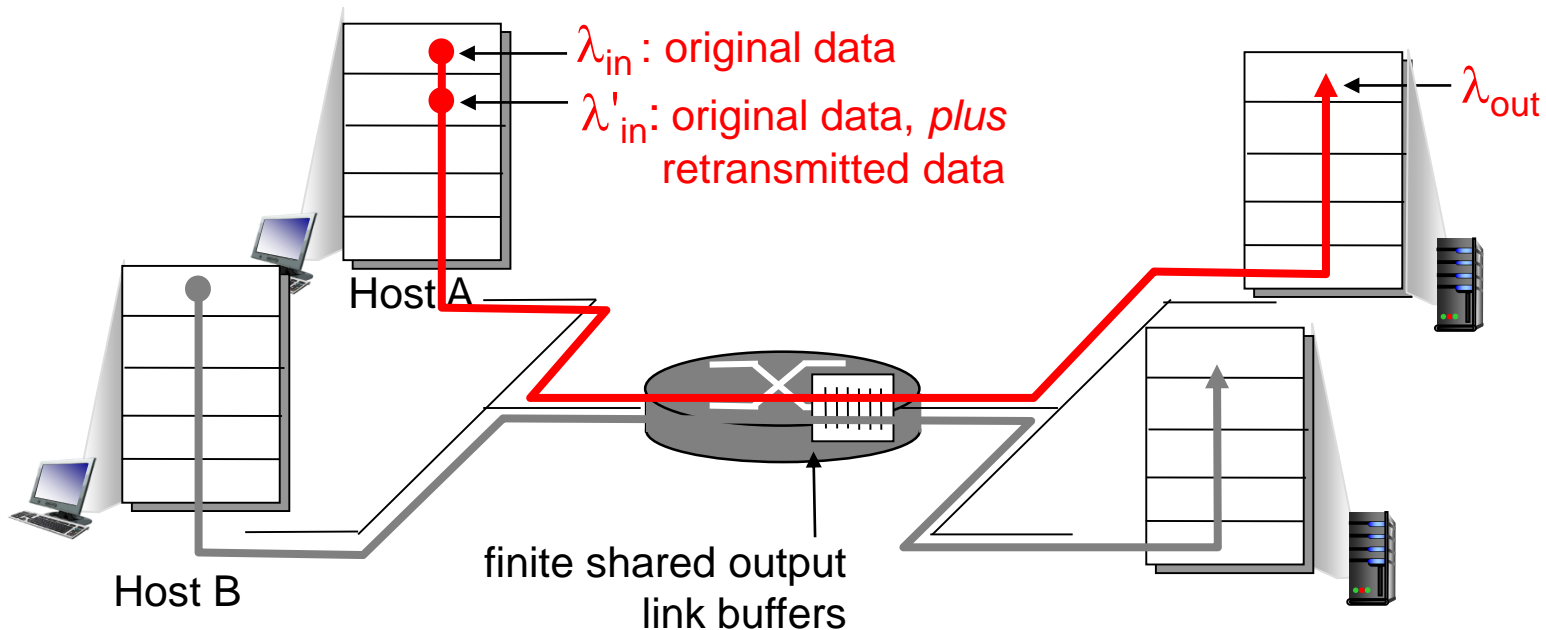- output link capacity: R
- no retransmission

original data: $\lambda_{in}$

throughput: $\lambda_{out}$

Host A

unlimited shared output link buffers

Host B



- maximum per-connection throughput: R/2

- large delays as arrival rate, $\lambda_{in}$, approaches capacity

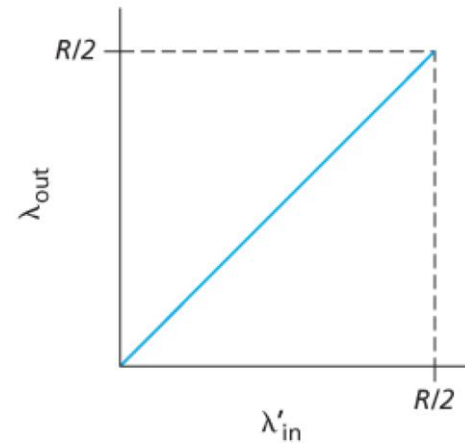# Causes/costs of congestion: scenario 2

- **one router, *finite* buffers**
- **sender retransmission of timed-out packet**
  - transport-layer input includes *retransmissions* : $\lambda'_{in} \geq \lambda_{in}$
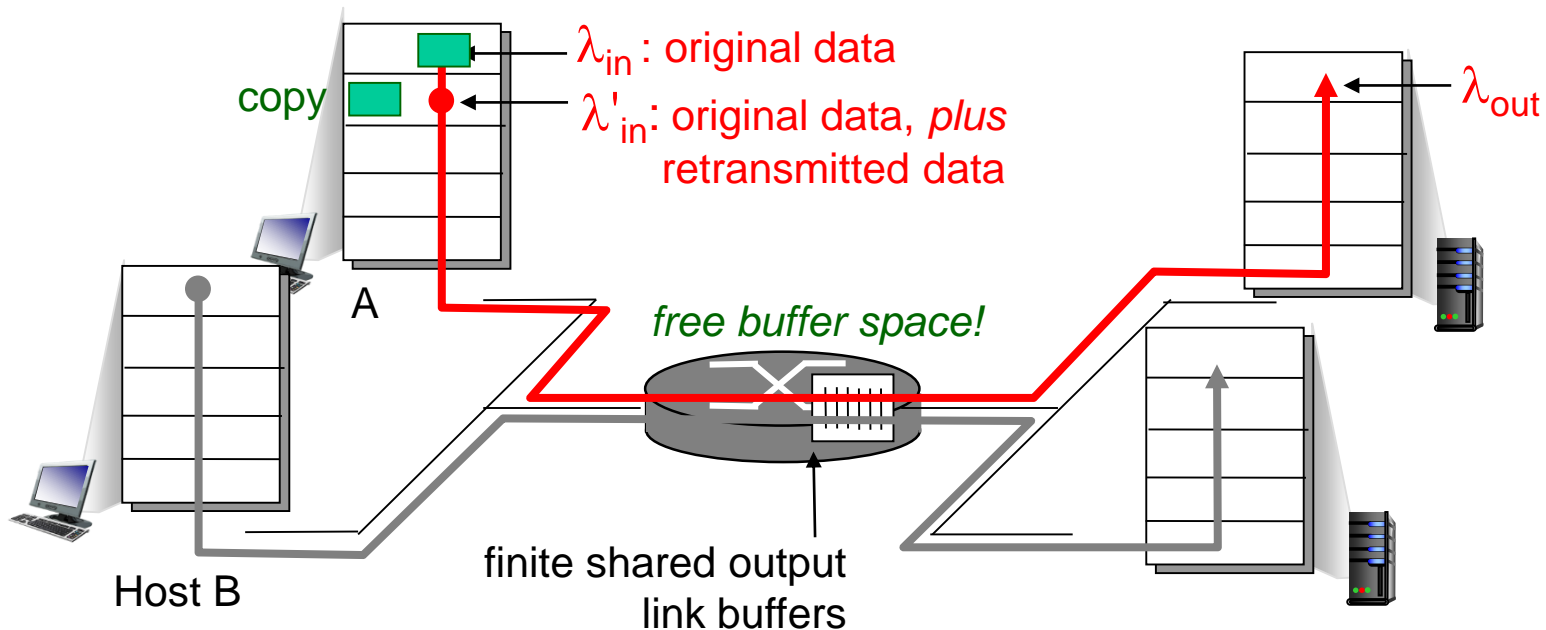


$\lambda_{in}$ : original data

$\lambda'_{in}$: 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



$\lambda_{in}$ : original data

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

$\lambda_{out}$

copy

free buffer space!
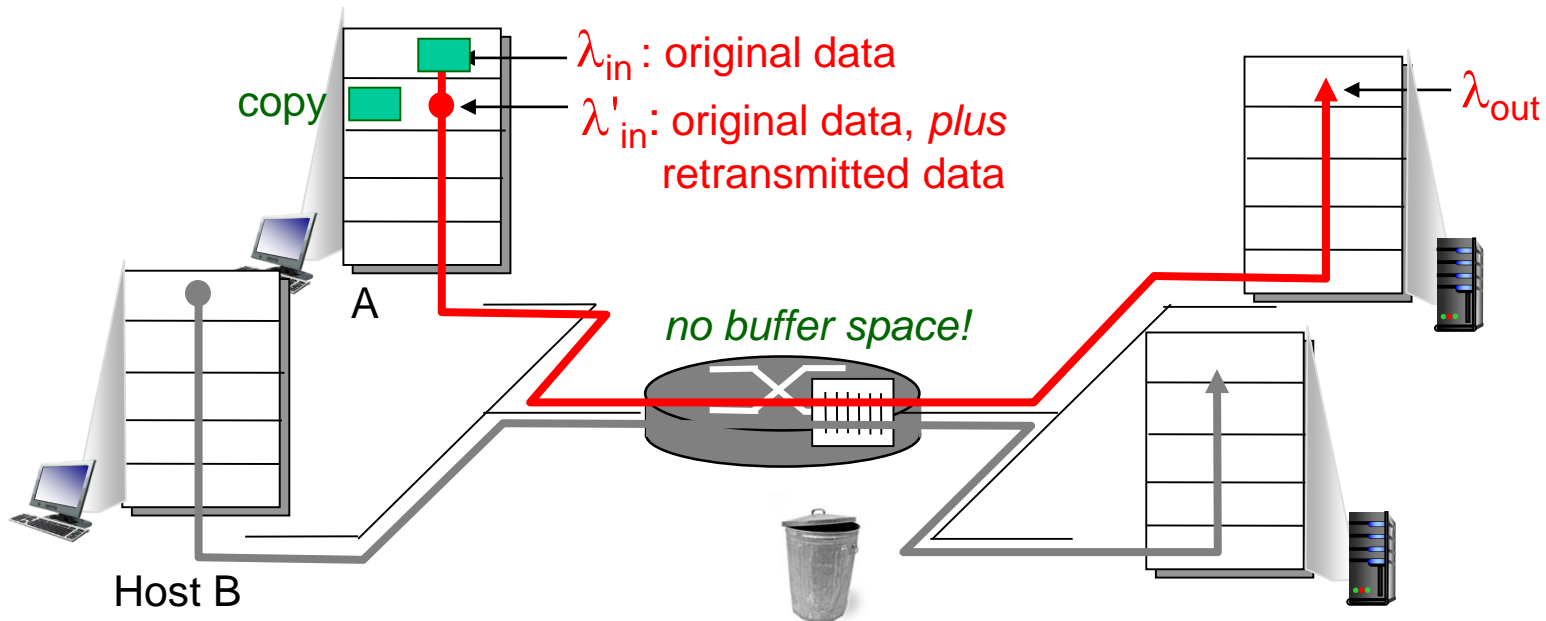
A

Host B

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

$\lambda_{in}$ : original data

copy

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

$\lambda_{out}$
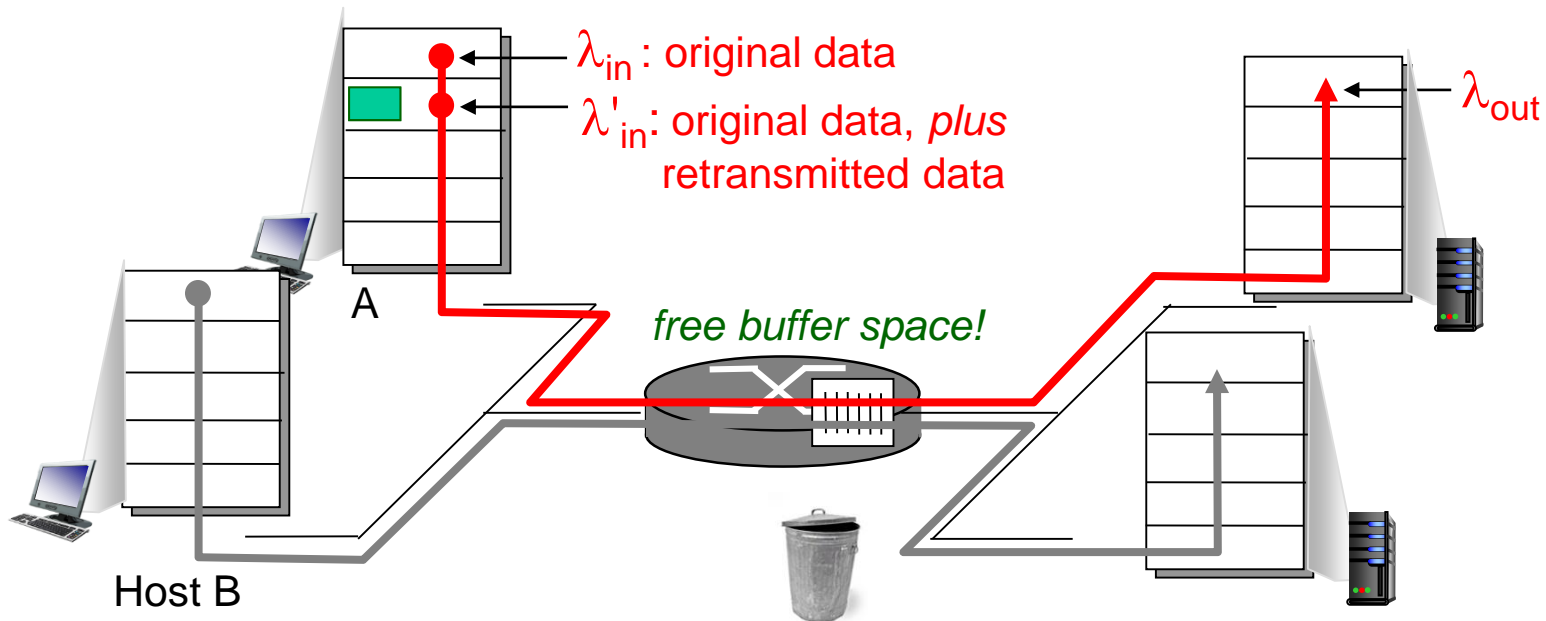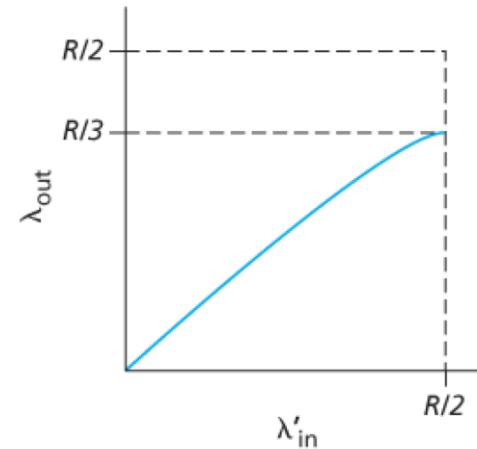
A

*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

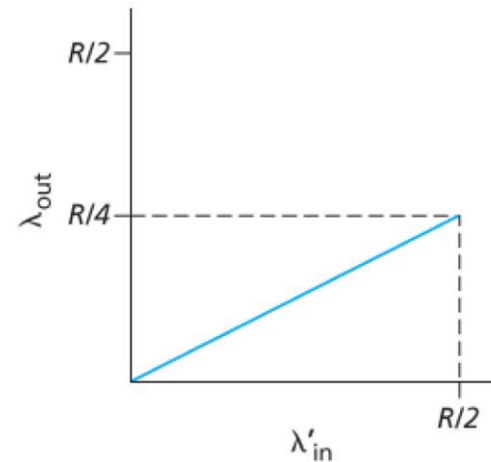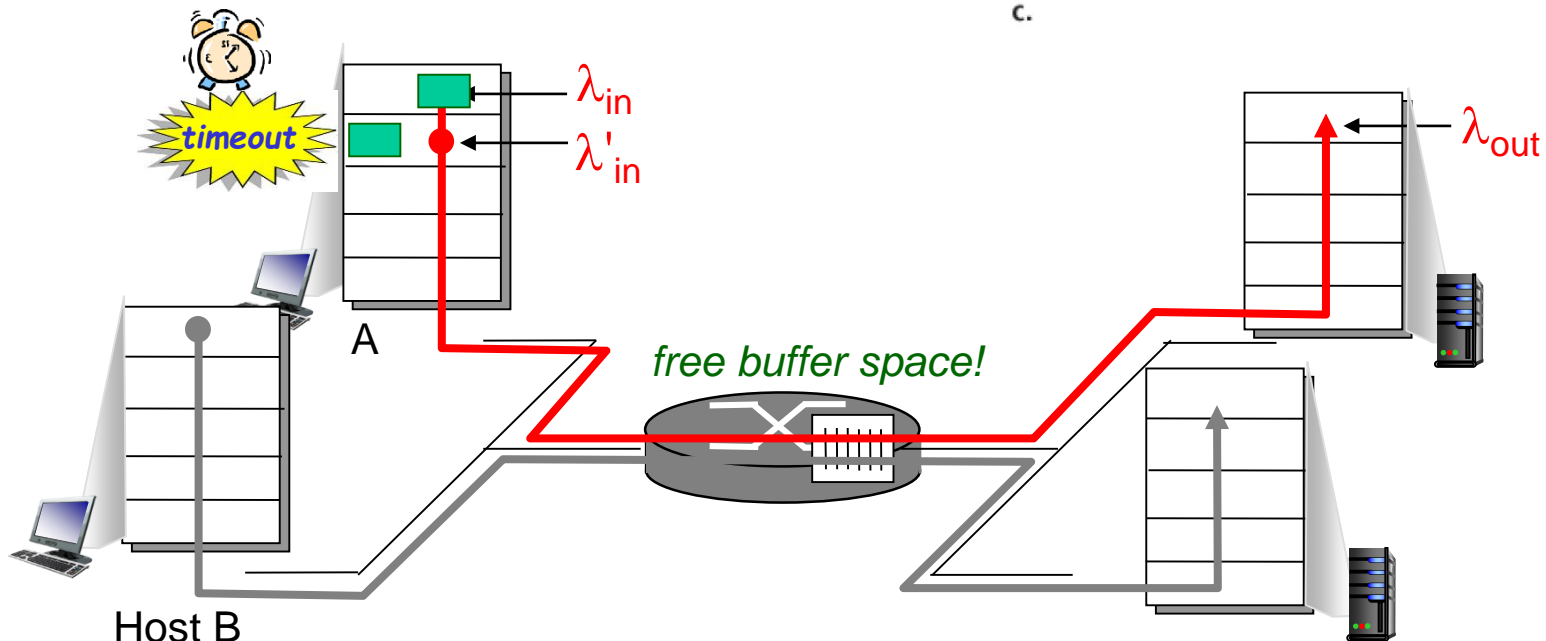- sender only resends if packet *known* to be lost



b.

$\lambda_{in}$ : original data

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

$\lambda_{out}$

A

*free buffer space!*

Host B

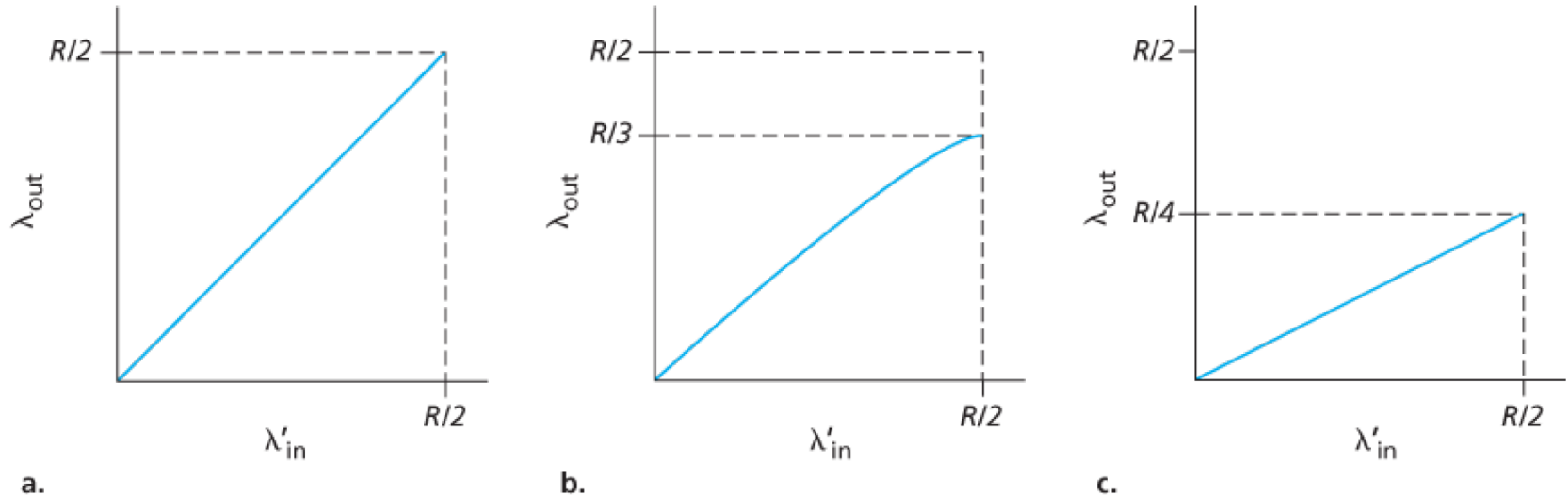# Causes/costs of congestion: scenario 2

## *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



$\lambda_{in}$

$\lambda'_{in}$

$\lambda_{out}$

timeout

A

Host B

free buffer space!

# Causes/costs of congestion: scenario 2
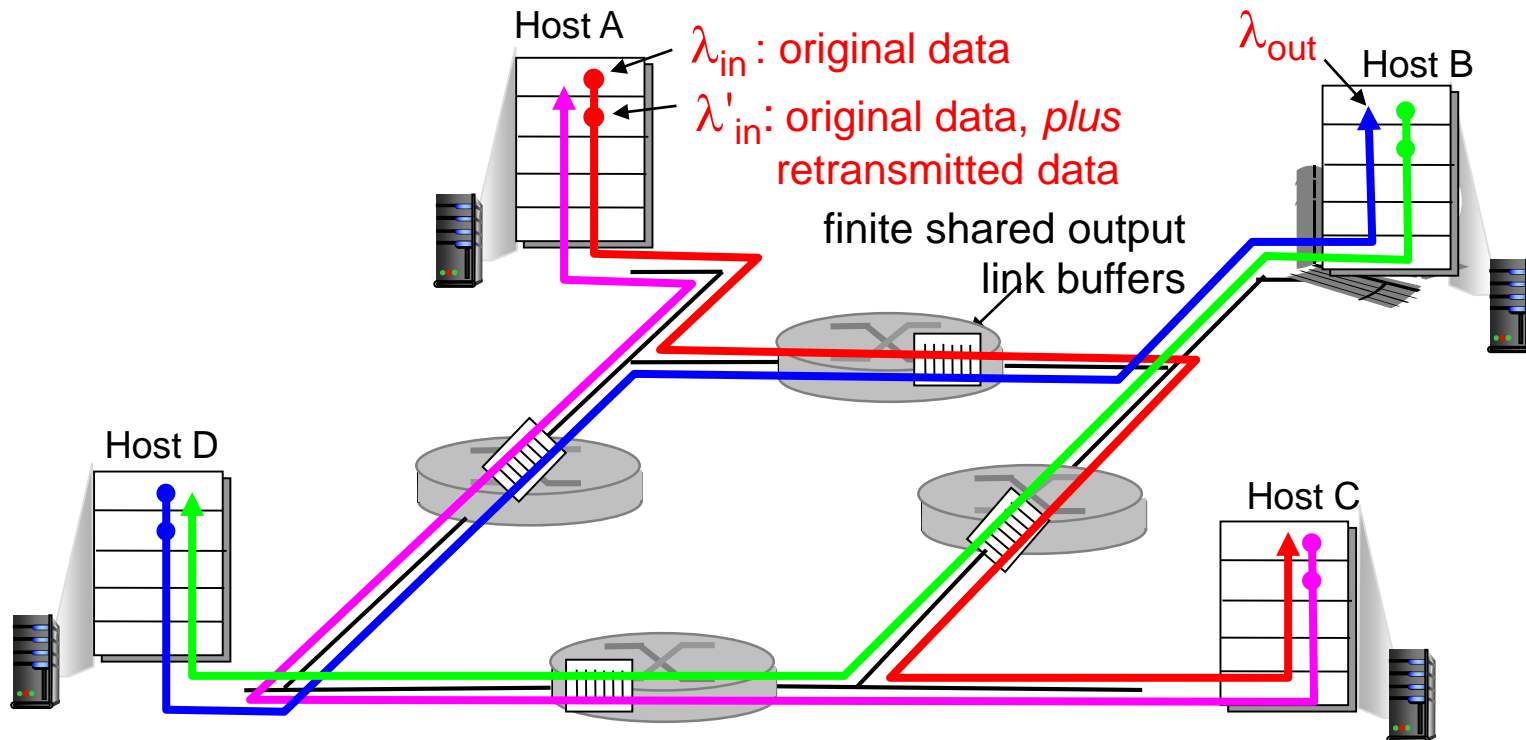


## "costs" of congestion:

- queuing delays
- more work (retrans) to compensate for lost packets
- 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 $\rightarrow$ 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 3



## another "cost" of congestion:

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

# Chapter 6 outline

6.1 transport-layer services

6.2 multiplexing and demultiplexing

6.3 connectionless transport: UDP

6.4 principles of reliable data transfer

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

6.6 principles of congestion control

6.7 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

$\longleftarrow$ cwnd $\longrightarrow$

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 $\approx \dfrac{\text{cwnd}}{\text{RTT}}$ 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

Host A                                            Host B

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
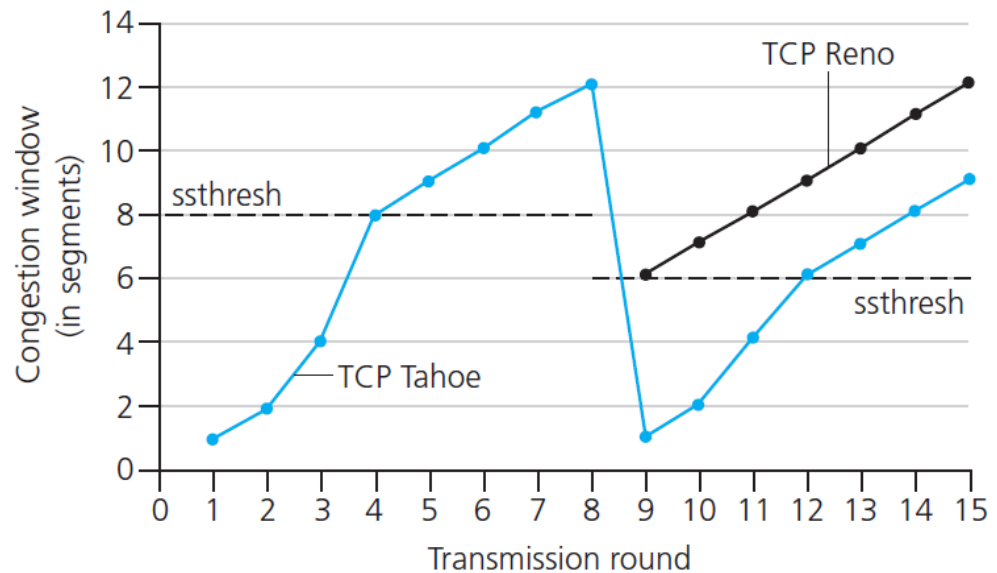- TCP Tahoe always sets **cwnd** to 1 (timeout or 3 duplicate acks)

# TCP: switching from slow start to CA
## (congestion avoidance)

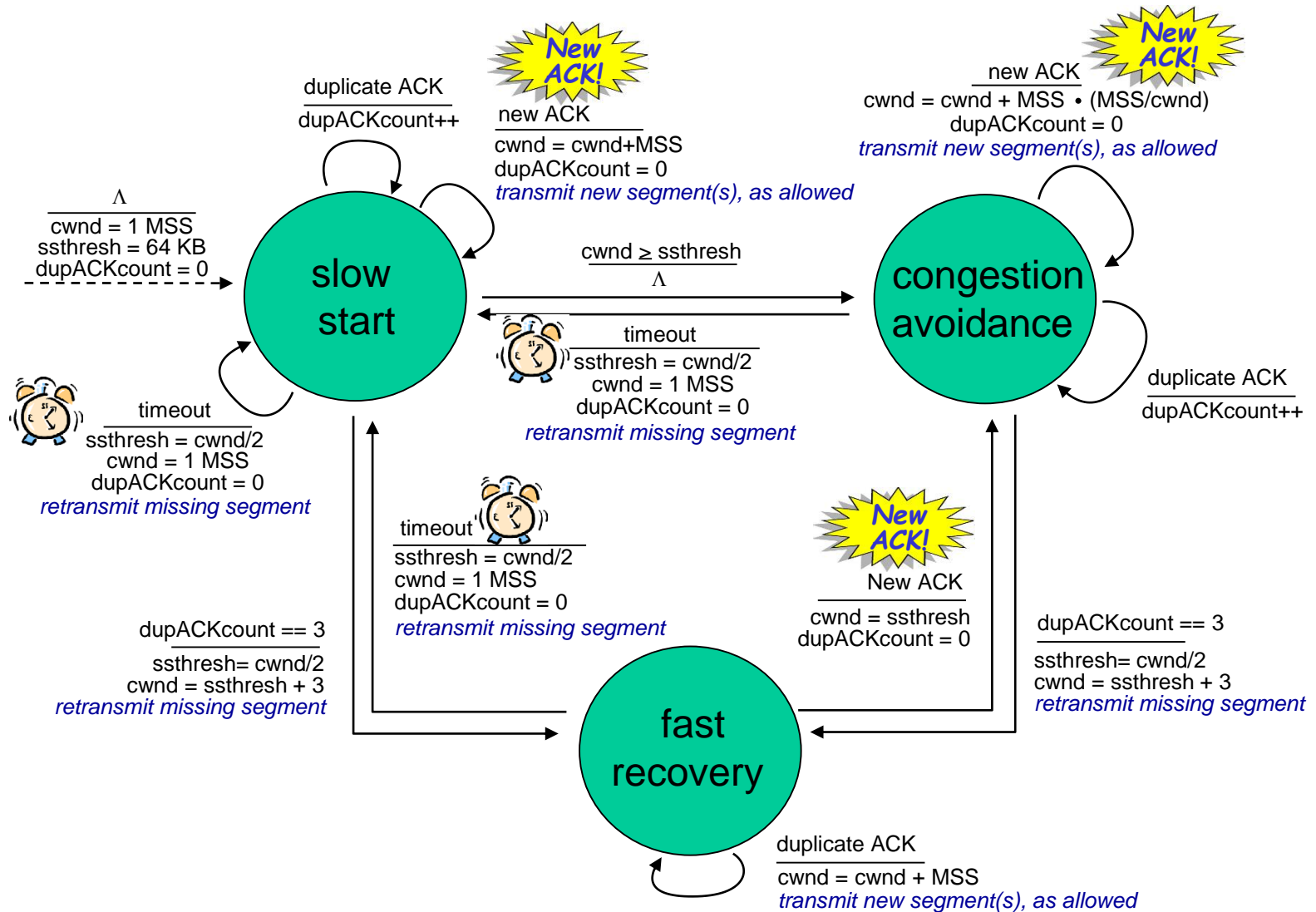Q: when should the exponential increase switch to linear?

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



Implementation:
- variable `ssthresh`
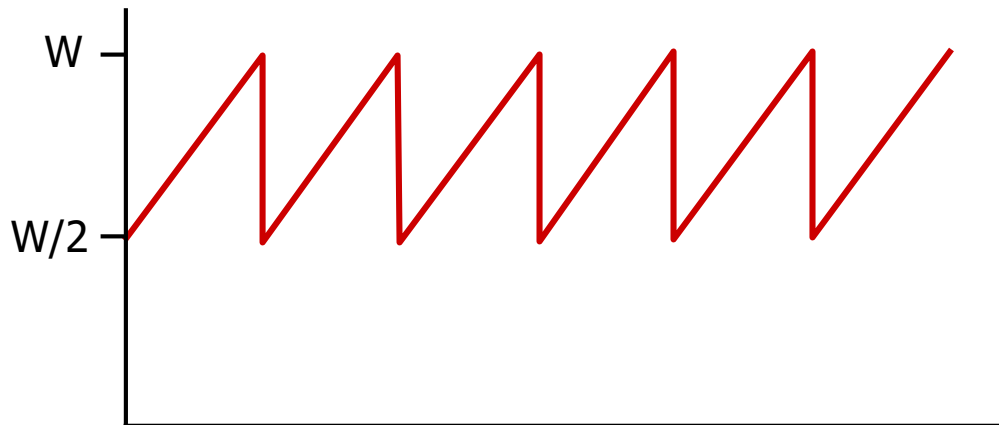- on loss event, `ssthresh` is set to 1/2 of `cwnd` just before loss event

# Summary: TCP Congestion Control

duplicate ACK
—————————
dupACKcount++

**New ACK!**

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

**New ACK!**

new ACK
—————————
cwnd = cwnd + MSS • (MSS/cwnd)
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*

duplicate ACK
—————————
dupACKcount++

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

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

**New ACK!**

New ACK
—————————
cwnd = ssthresh
dupACKcount = 0

dupACKcount == 3
—————————
ssthresh= cwnd/2
cwnd = ssthresh + 3
*retransmit missing segment*

dupACKcount == 3
—————————
ssthresh= cwnd/2
cwnd = ssthresh + 3
*retransmit missing segment*

**fast recovery**

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

# 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

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

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

- example: 1500 byte segments, 100ms RTT, want 10 Gbps throughput
- requires $\overline{W}$ = 83,333 in-flight segments
- throughput in terms of segment loss probability, L [Mathis 1997]:
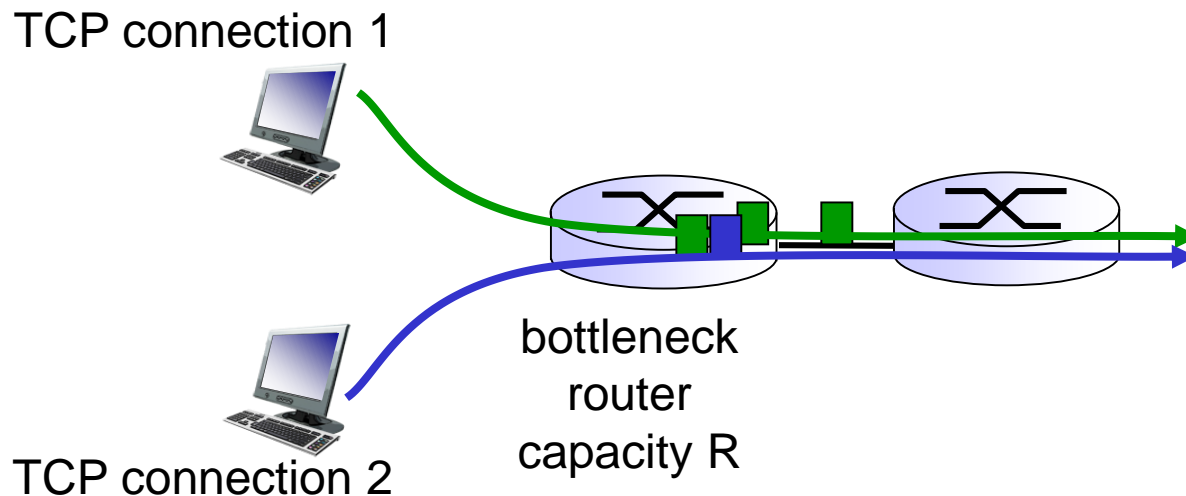
$$\text{TCP throughput} = \frac{1.22 \cdot \text{MSS}}{\text{RTT} \sqrt{L}}$$

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

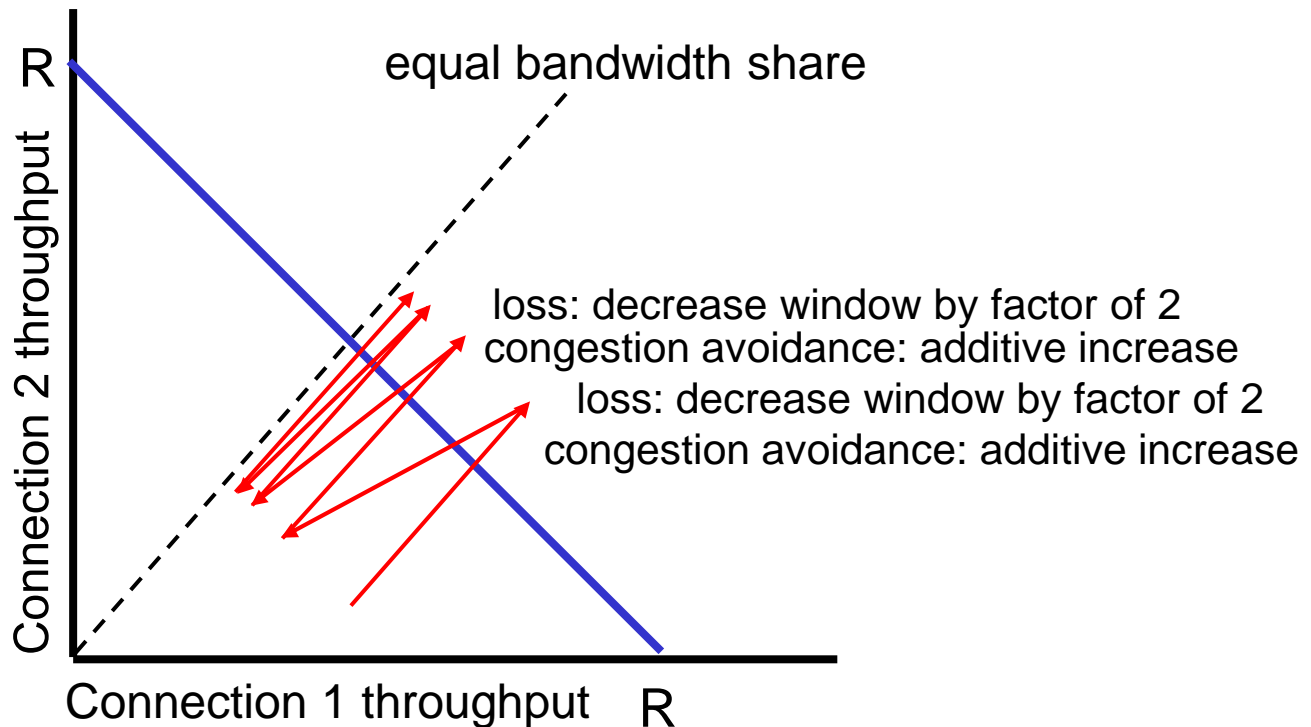- new versions of TCP for high-speed

# 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

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

Connection 2 throughput

Connection 1 throughput

R

R

# Fairness (more)

## Fairness and UDP

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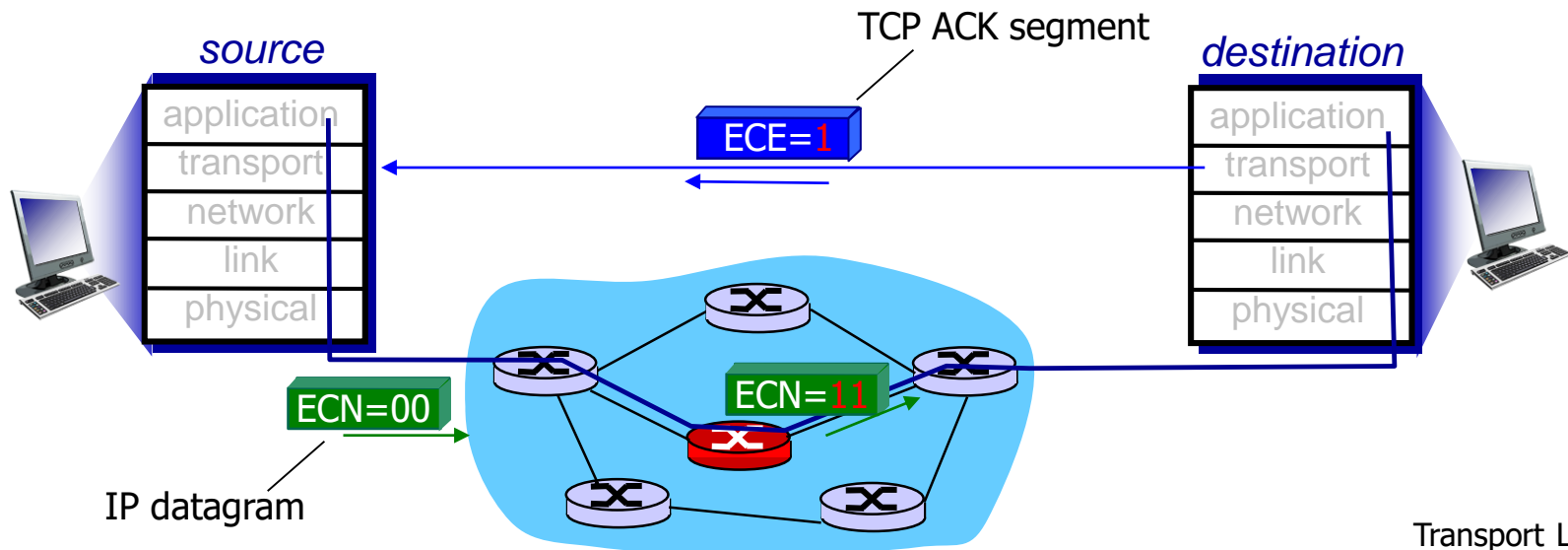  - send audio/video at constant rate, tolerate packet loss

## 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

# Explicit Congestion Notification (ECN)

*network-assisted congestion control:*

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



TCP ACK segment

*source*    application transport network link physical

ECE=1

*destination*    application transport network link physical

ECN=00

ECN=11

IP datagram

# Chapter 6: summary

- **principles behind transport layer services:**
  - multiplexing, demultiplexing
  - reliable data transfer
  - flow control
  - congestion control
- **instantiation, implementation in the Internet**
  - UDP
  - TCP