# CS 305: Computer Networks
## Fall 2023

### Lecture 8: Transport Layer

**Ming Tang**

Department of Computer Science and Engineering
Southern University of Science and Technology (SUSTech)

# TCP Reliable Data Transfer

- ❖ Segment structure
- ❖ Round-trip time estimation
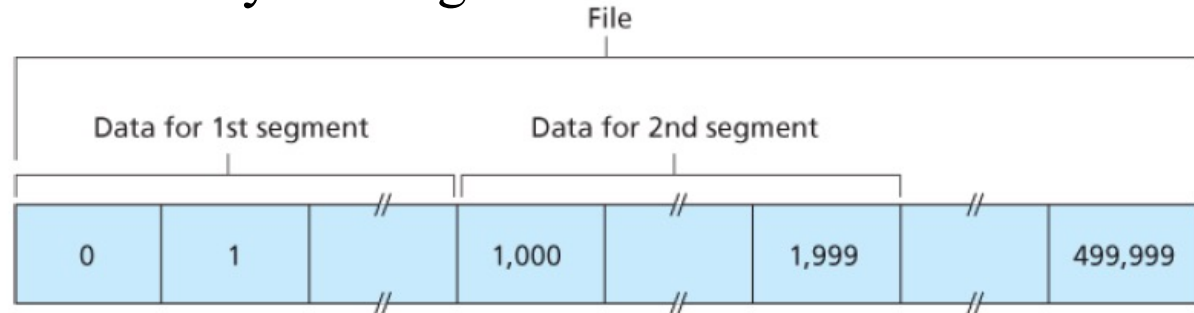- ❖ Reliable data transfer
- ❖ Flow control
- ❖ Control management

# TCP seq. numbers, ACKs

TCP views data as an unstructured, but ordered, stream of bytes.

- Sequence numbers are over the stream of transmitted bytes and *not* over the series of transmitted segments

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

File

| Data for 1st segment | | Data for 2nd segment | | |
|---|---|---|---|---|
| 0 | 1 | 1,000 | 1,999 | 499,999 |

acknowledgements:
- seq # of next byte expected from other side
  - E.g., receiver has received bytes numbered 0 through 535 and 900 through 1000; then, acknowledgement number is 536.
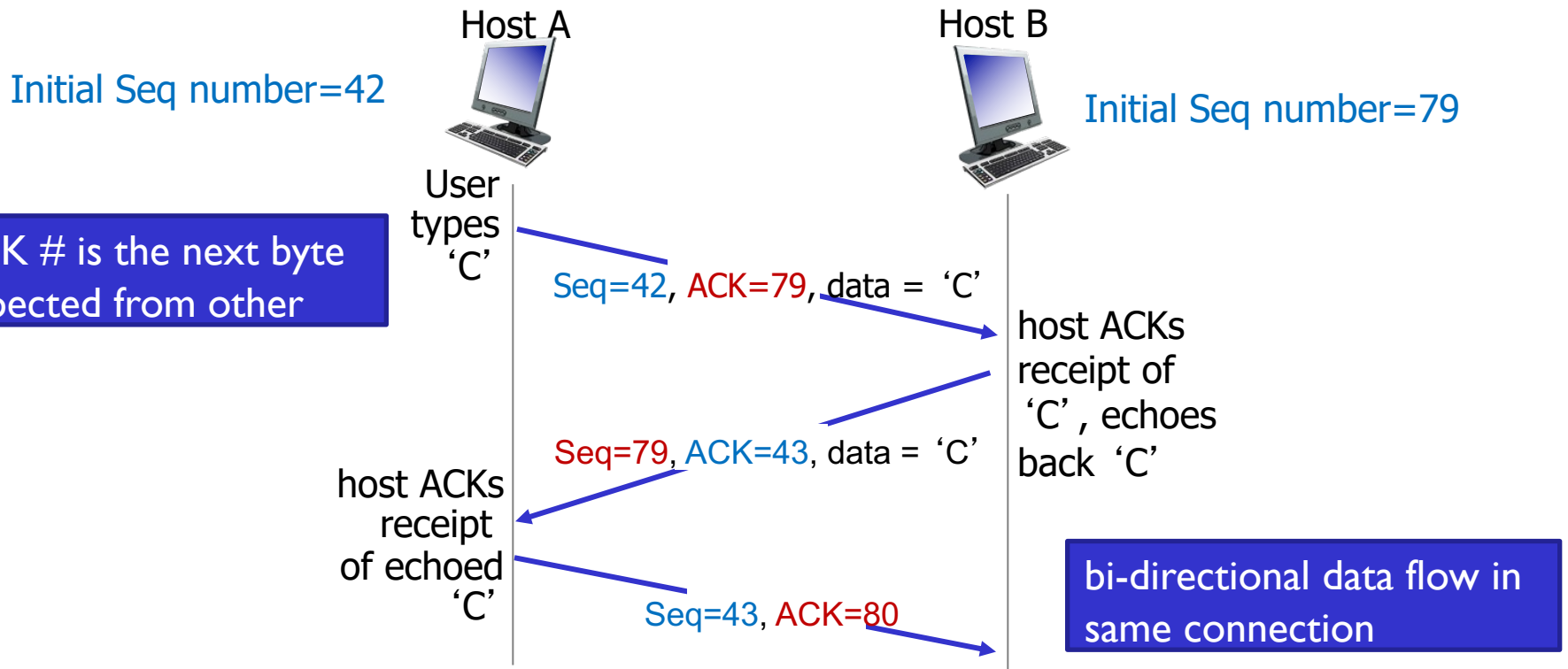- cumulative ACK

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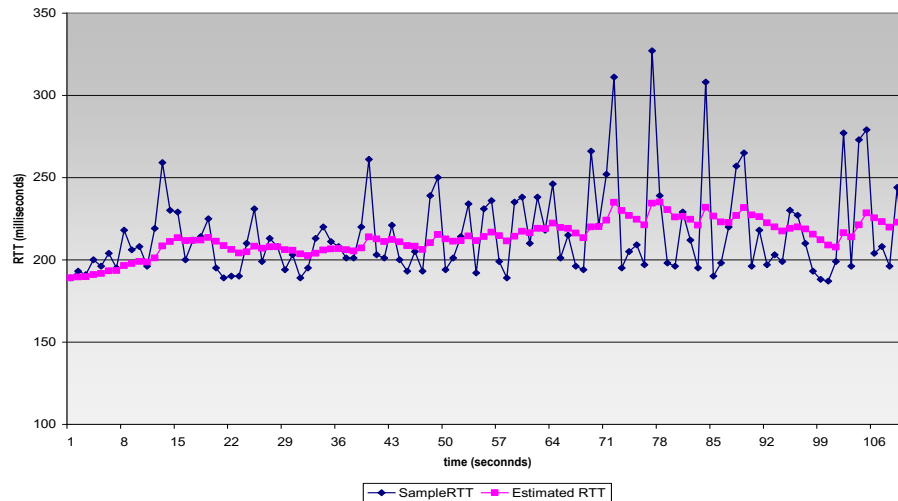- A: TCP spec doesn't say, - up to implementor

Initial sequence number is randomly chosen

# Telnet Case Study

❖ **full duplex data:**
  ▪ bi-directional data flow in same connection

Host A

Initial Seq number=42

Host B

Initial Seq number=79

User types 'C'

ACK # is the next byte expected from other

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

bi-directional data flow in same connection

# TCP round trip time, timeout



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

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

estimated RTT          "safety margin"

# TCP Reliable Data Transfer

- ❖ Segment structure
- ❖ Round-trip time estimation
- ❖ <span style="color:red">Reliable data transfer</span>
- ❖ Flow control
- ❖ Control management

# TCP reliable data transfer

❖ TCP creates rdt service on top of IP's unreliable service
  ▪ pipelined segments: window size, SendBase
  ▪ 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
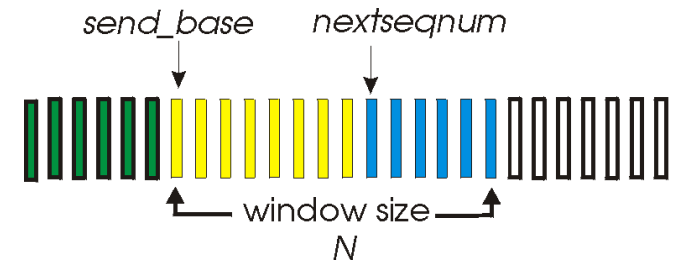
# TCP sender events:



*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 events:

```
NextSeqNum=InitialSeqNumber
SendBase=InitialSeqNumber

loop (forever) {
    switch(event)

        event: data received from application above
            create TCP segment with sequence number NextSeqNum
            if (timer currently not running)
                start timer
            pass segment to IP
            NextSeqNum=NextSeqNum+length(data)
            break;

        event: timer timeout
            retransmit not-yet-acknowledged segment with
                smallest sequence number
            start timer
            break;

        event: ACK received, with ACK field value of y
            if (y > SendBase) {
                SendBase=y
                if (there are currently any not-yet-acknowledged segments)
                    start timer
                }
            break;

    } /* end of loop forever */
```
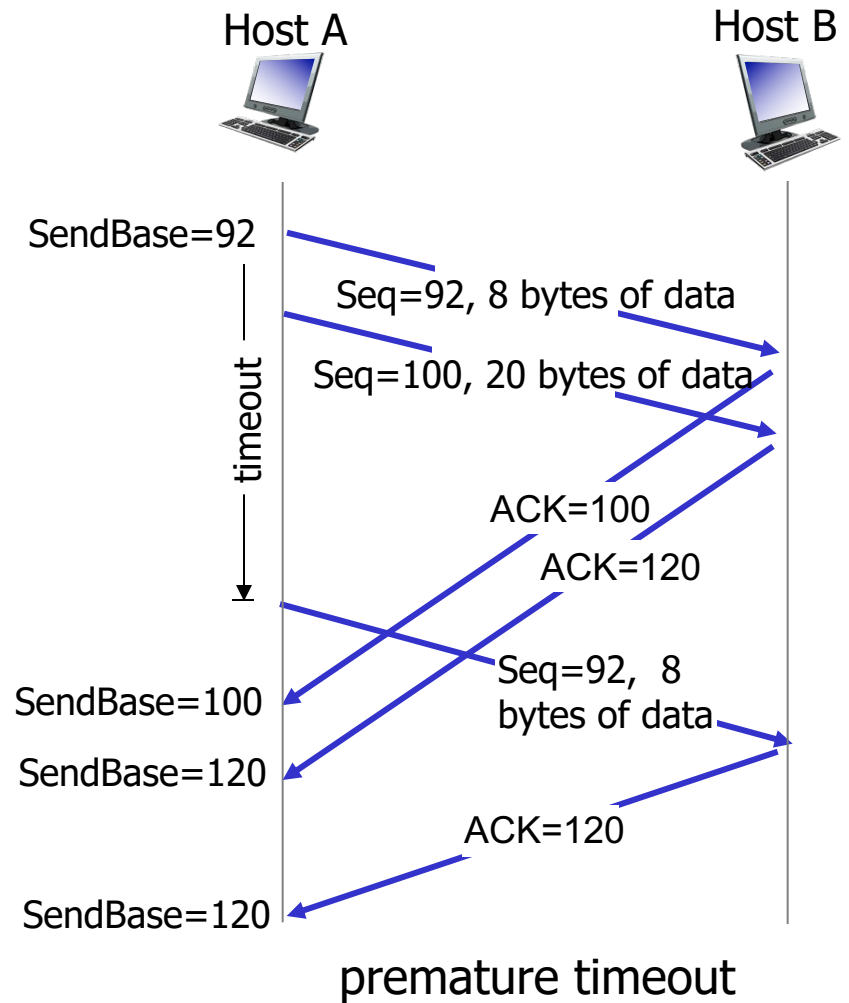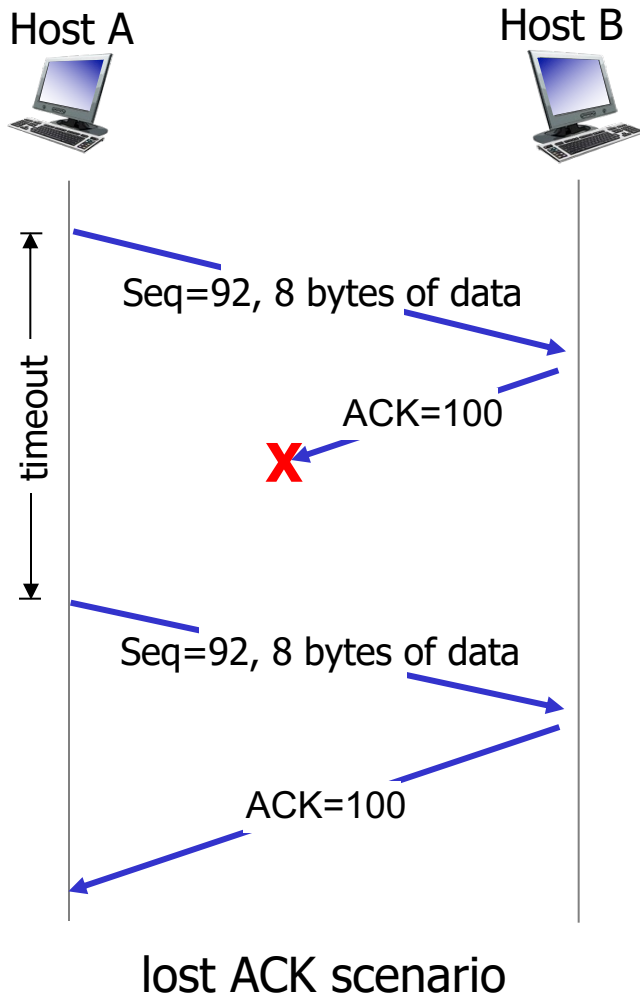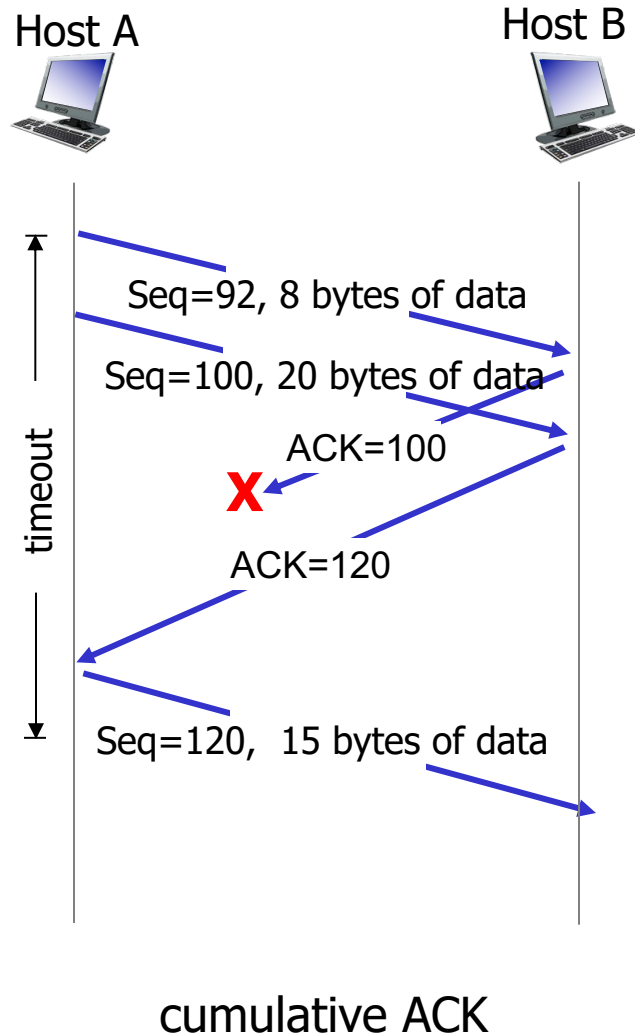
# TCP: retransmission scenarios



lost ACK scenario

premature timeout

# TCP: retransmission scenarios

Host A                    Host B

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

ACK=100

X

ACK=120

Seq=120,  15 bytes of data

timeout

cumulative ACK

# TCP receiver [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 duplicate 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

```
NextSeqNum=InitialSeqNumber
SendBase=InitialSeqNumber

loop (forever) {
    switch(event)

        event: data received from application
            create TCP segment with sequence
            if (timer currently not running)
                start timer
            pass segment to IP
            NextSeqNum=NextSeqNum+length(data)
            break;

        event: timer timeout
            retransmit not-yet-acknowledged segment
                smallest sequence number
            start timer
            break;
```

```
event: ACK received, with ACK field value of y
        if (y > SendBase) {
        SendBase=y
        if (there are currently any not yet
                acknowledged segments)
            start timer

        }
        else {/* a duplicate ACK for already ACKed
                segment */
            increment number of duplicate ACKs
                received for y
            if (number of duplicate ACKS received
                for y==3)
                /* TCP fast retransmit */
                resend segment with sequence number y
            }
        break;
```

```
        event: ACK received, with ACK field value of y
            if (y > SendBase) {
                SendBase=y
                if (there are currently any not-yet-acknowledged segments)
                    start timer
                }
            break;
```
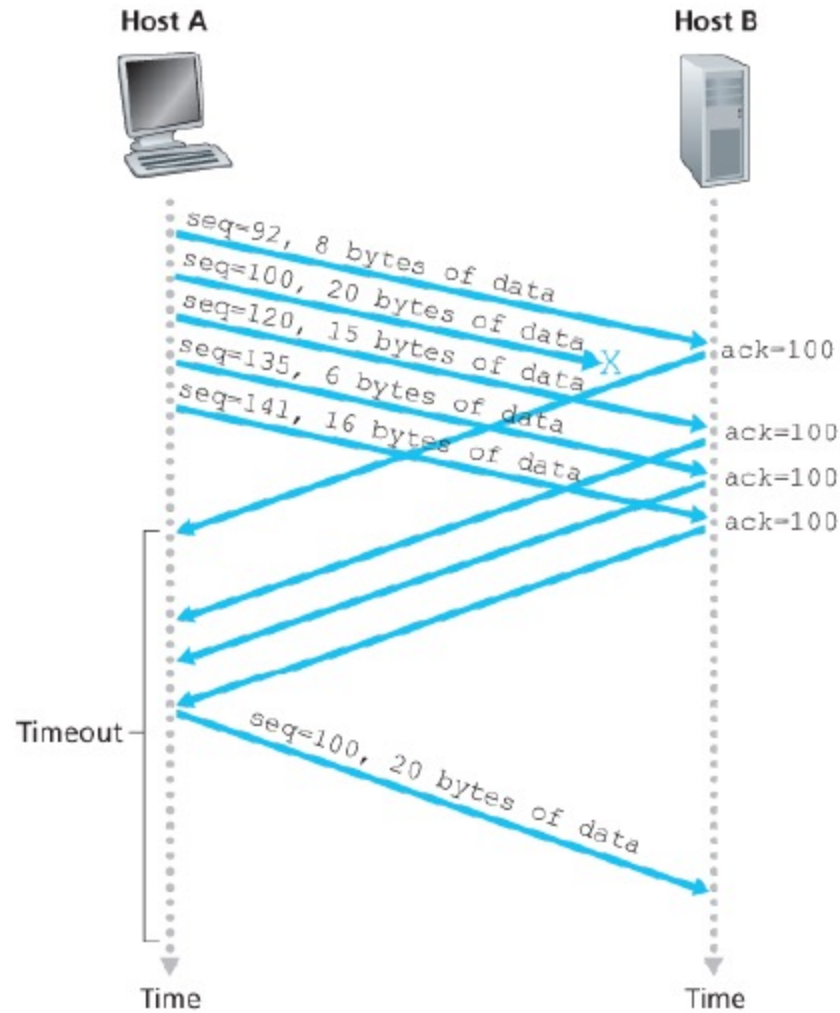
```
} /* end of loop forever */
```

# TCP fast retransmit
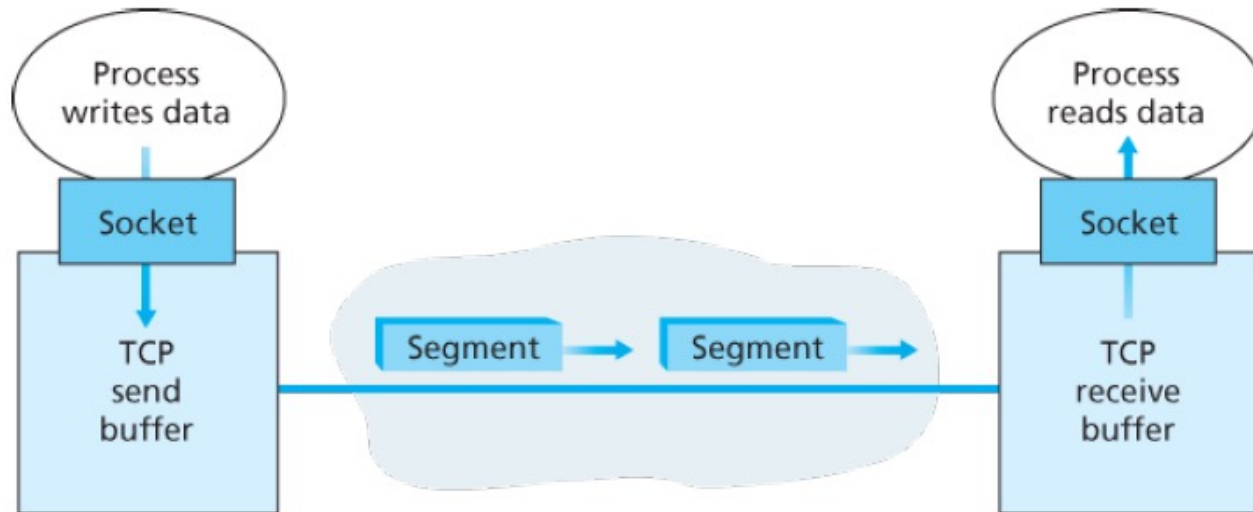


fast retransmit after sender
receipt of triple duplicate ACK

# TCP Reliable Data Transfer

❖ Segment structure
❖ Round-trip time estimation
❖ Reliable data transfer
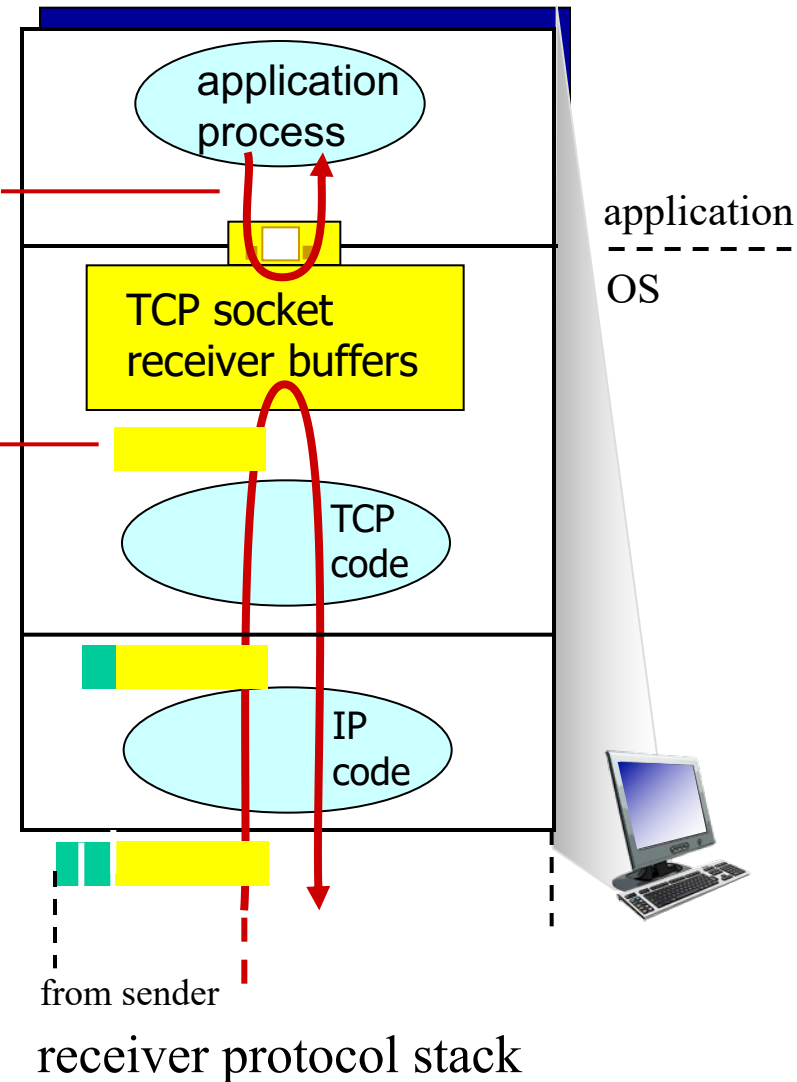❖ Flow control
❖ Control management

# TCP: Overview



- TCP connection
- TCP grab chunks of data from the sender buffer
- TCP receives a segment at the other end, place it in receiver buffer
- application reads the stream from the receive buffer

# TCP flow control

application may
remove data from
TCP socket buffers ….

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

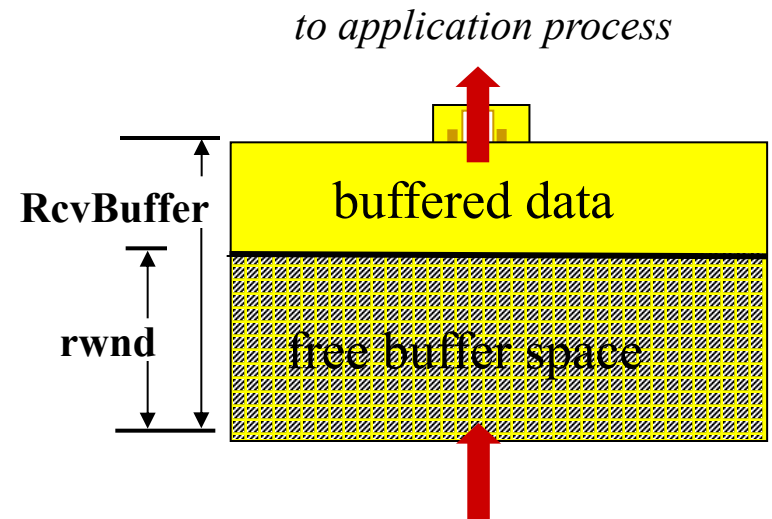*flow control*
Receiver controls sender, so
sender won't overflow
receiver's buffer by transmitting
too much, too fast

application
process

application
------
OS

TCP socket
receiver buffers

TCP
code

IP
code

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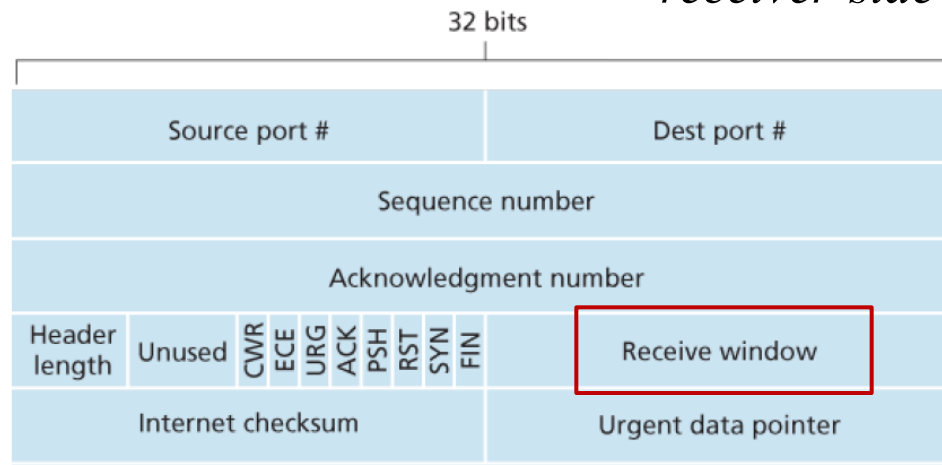
- ❖ **RcvBuffer** size set via socket options (typical default is 4096 bytes)
- ❖ many operating systems autoadjust **RcvBuffer**

**rwnd=RcvBuffer−[LastByteRcvd−LastByteRead]**

*to application process*

RcvBuffer

buffered data

rwnd

free buffer space

*TCP segment payloads*

*receiver-side buffering*

32 bits

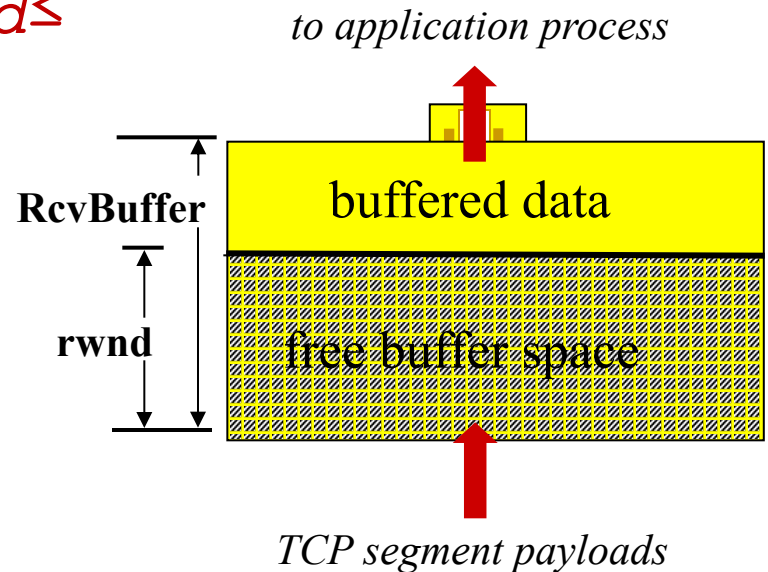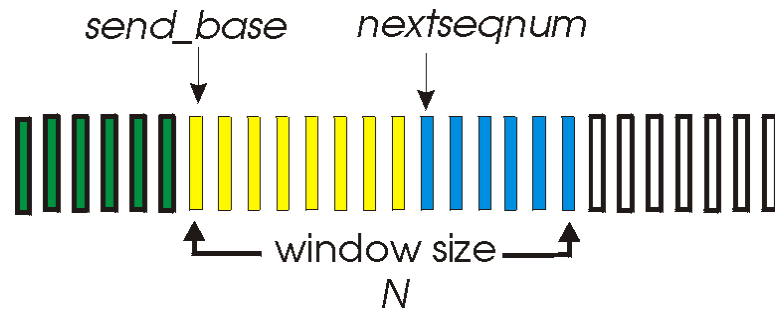| Source port # | Dest port # |
|---|---|
| Sequence number | |
| Acknowledgment number | |
| Header length / Unused / CWR ECE URG ACK PSH RST SYN FIN | Receive window |
| Internet checksum | Urgent data pointer |

# TCP flow control

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

$$LastByteSent-LastByteAcked\leq rwnd$$

Guarantees receive buffer will not overflow



*to application process*

**RcvBuffer**

buffered data

**rwnd**

free buffer space

*TCP segment payloads*

*receiver-side buffering*



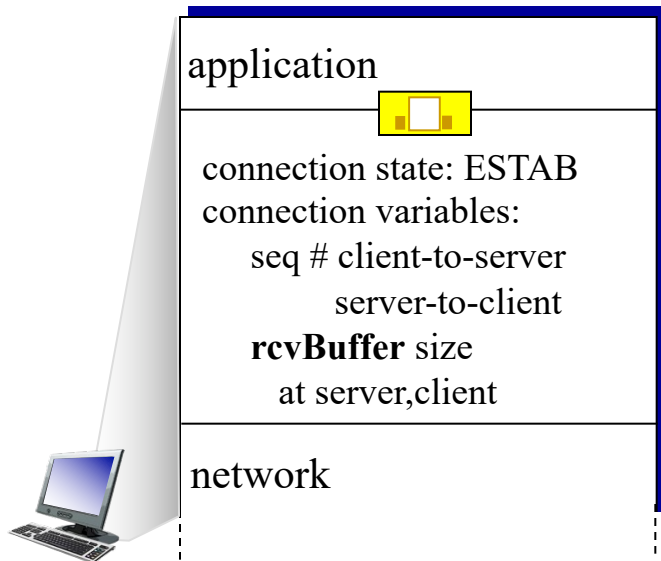send_base          nextseqnum

window size
N

# TCP Reliable Data Transfer

- ❖ Segment structure
- ❖ Round-trip time estimation
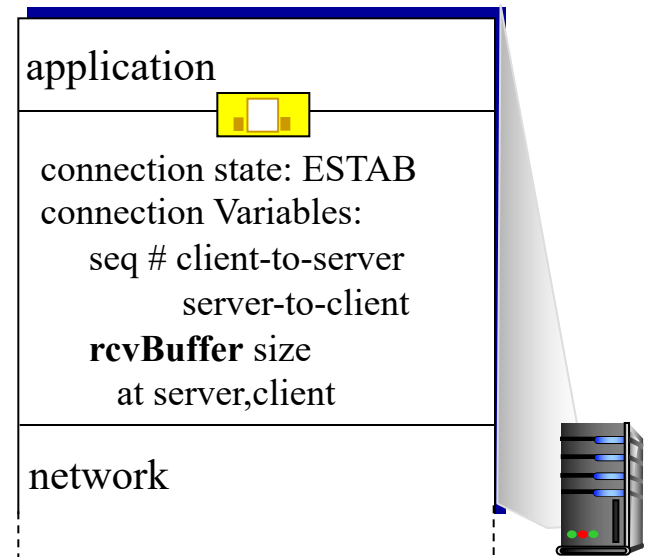- ❖ Reliable data transfer
- ❖ Flow control
- ❖ Control management

# Connection Management

before exchanging data, sender/receiver "handshake":

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

❖ agree on connection parameters

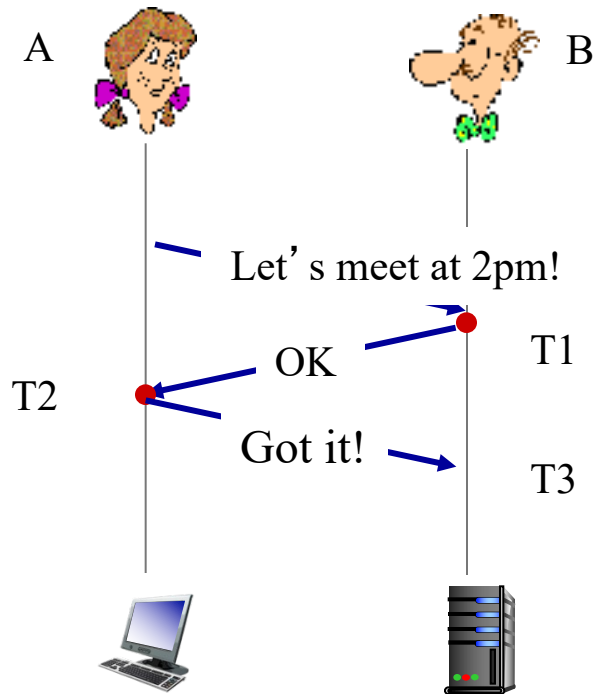| application | application |
|---|---|
| connection state: ESTAB<br>connection variables:<br>    seq # client-to-server<br>        server-to-client<br>**rcvBuffer** size<br>    at server,client | connection state: ESTAB<br>connection Variables:<br>    seq # client-to-server<br>        server-to-client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

```
clientSocket = socket(AF_INET, SOCK_STREAM);
clientSocket.connect((hostname,port number));
```
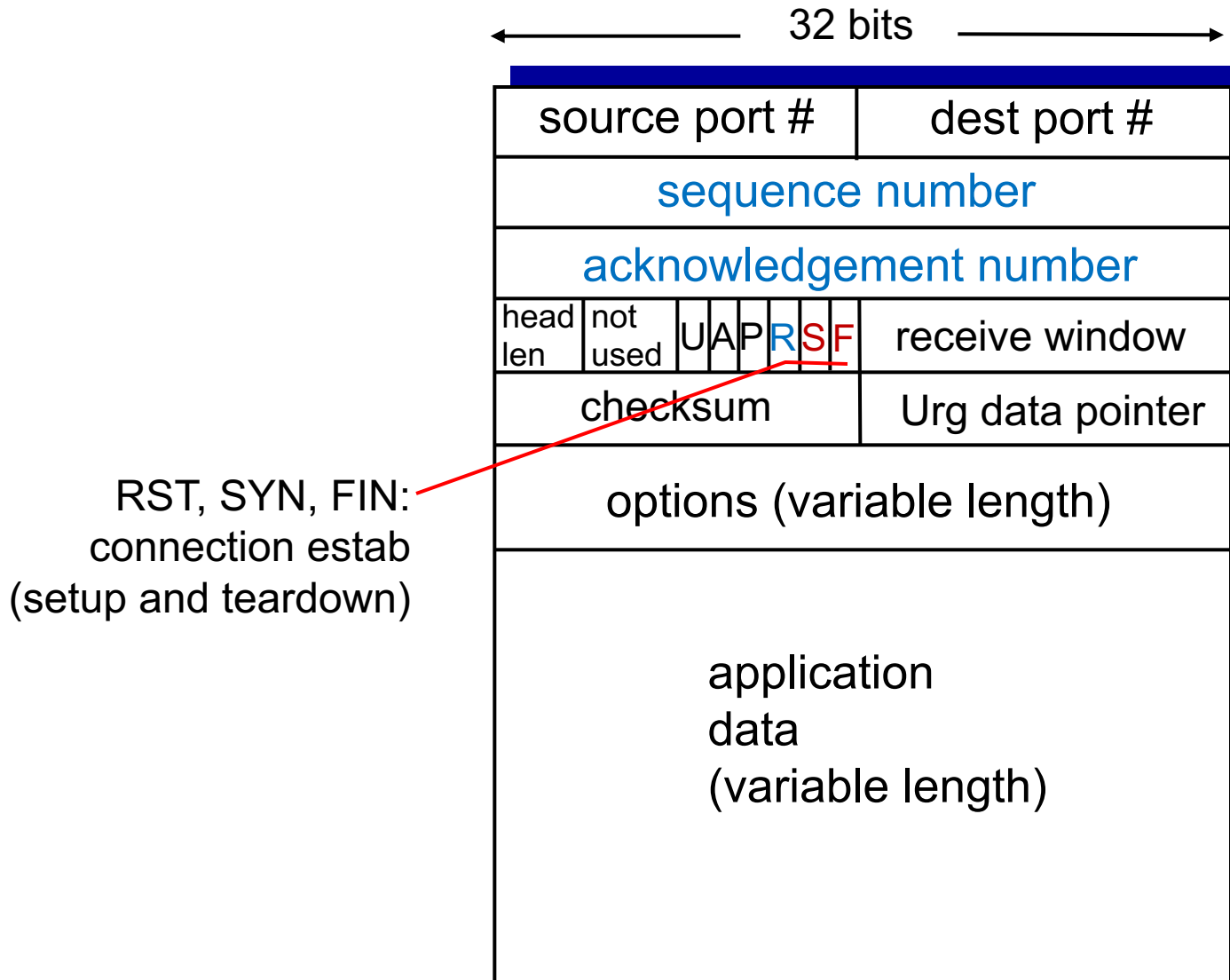
```
connectionSocket = welcomeSocket.accept();
```

# Agreeing to establish a connection

3-way handshake:



❖ T1 : B knows A's transmitter and B's receiver is OK

❖ T2: A knows A's transceiver and B's transceiver is OK, B has no more information than T1

❖ T3: Both A and B know their transceiver are OK, they can start the communication!

# TCP segment structure

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)

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

# TCP 3-way handshake

**client state**

LISTEN

SYNSENT

ESTAB

choose init seq num, x
send TCP SYN msg

SYNbit=1, Seq=x

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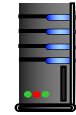send TCP SYNACK
msg, acking SYN

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

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

SYNbit=0, Seq=x+1
ACKbit=1, ACKnum=y+1

received ACK(y)
indicates client is live
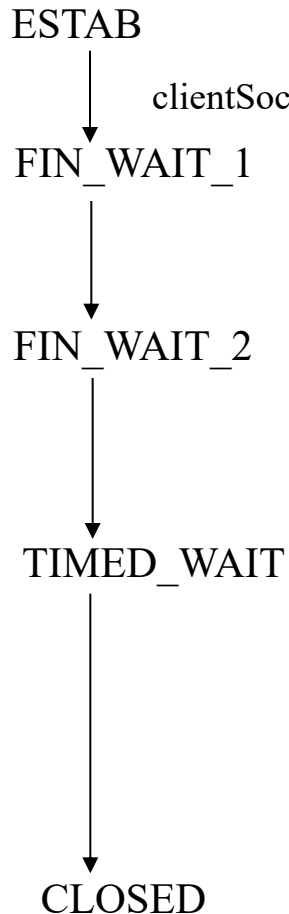
**server state**

LISTEN

SYN RCVD

ESTAB

Once these three steps have been completed, the client and server hosts can
send segments containing data to each other.

- In each of these future segments, SYNbit=0

# TCP: closing a connection

**client state**                              **server state**

ESTAB                                                ESTAB

clientSocket.close()

FIN_WAIT_1          can no longer
                   send but can        FINbit=1, seq=x
                   receive data                                CLOSE_WAIT

                                      ACKbit=1; ACKnum=x+1
FIN_WAIT_2         wait for server                             can still
                   close                                       send data

                                       FINbit=1, seq=y         LAST_ACK

TIMED_WAIT                                                     can no longer
                                                               send data

                   timed wait         ACKbit=1; ACKnum=y+1
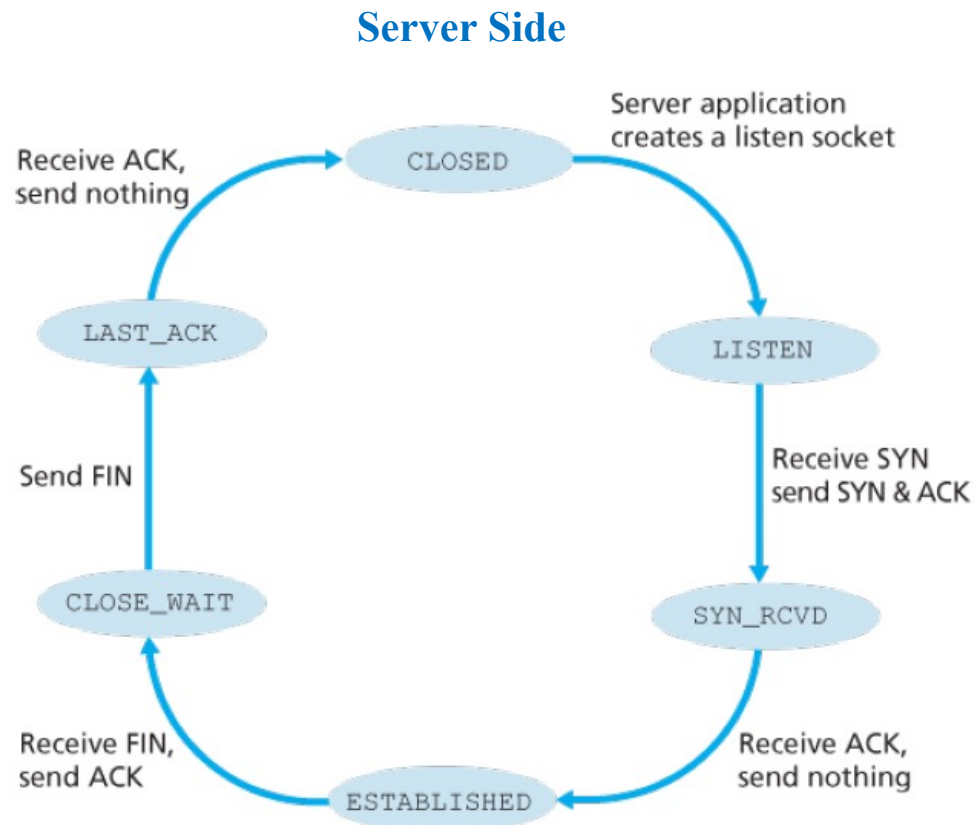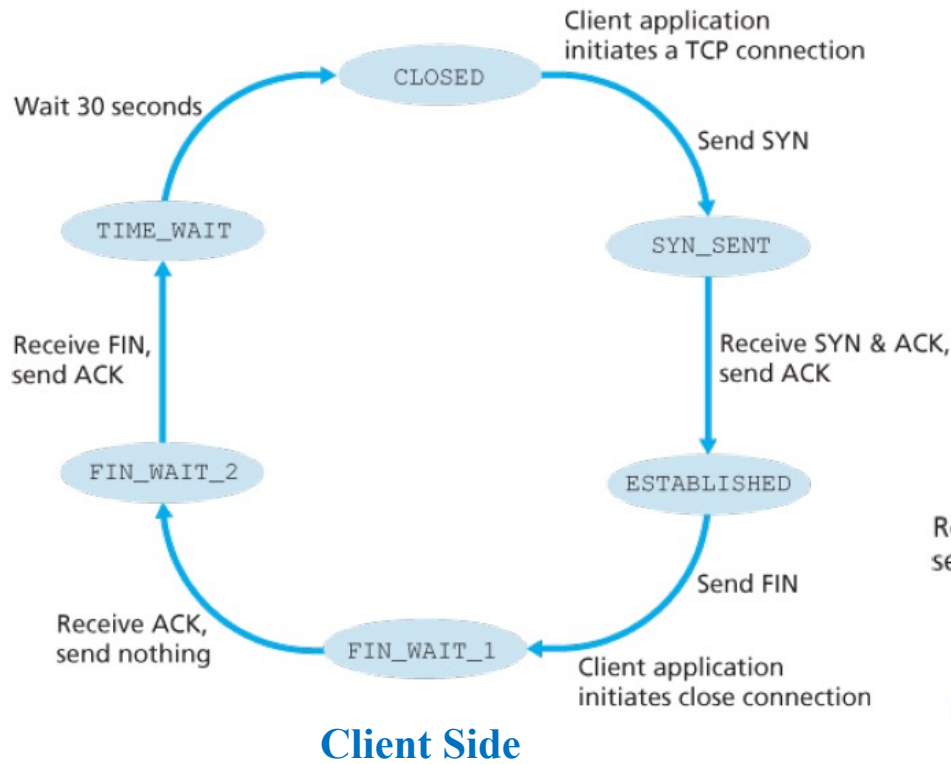                                                               CLOSED

CLOSED

The TIME_WAIT state lets the TCP client resend the final
acknowledgment in case the ACK is lost.

# TCP: closing a connection

❖ Four-way handshaking
- Either of the two processes participating in a TCP connection can end the connection.

❖ client, server each close their side of connection
- send TCP segment with FIN bit = 1

❖ respond to received FIN with ACK

❖ Why FIN and ACK can not be sent in one msg as SYNACK in connection establishment?
- The other side may still have packets need to be sent. It can not send FIN until the transmission is finished.

# TCP States

Client application
initiates a TCP connection

CLOSED

Wait 30 seconds

Send SYN

TIME_WAIT

SYN_SENT

Receive FIN,
send ACK

Receive SYN & ACK,
send ACK

**Server Side**

FIN_WAIT_2

ESTABLISHED

Server application
creates a listen socket

Receive ACK,
send nothing

CLOSED

Receive ACK,
send nothing

Send FIN

Receive ACK,
send nothing

FIN_WAIT_1

Client application
initiates close connection

LAST_ACK

LISTEN

**Client Side**

Receive SYN
send SYN & ACK

Send FIN

CLOSE_WAIT

SYN_RCVD

Receive FIN,
send ACK

Receive ACK,
send nothing

ESTABLISHED

# Reset Segment

When a host receives a TCP segment whose port numbers or source IP address do not match with any of the ongoing sockets.

❖ Then the host will send a special reset segment to the source. RST flag bit is set to 1.

❖ "I don't have a socket for that segment. Please do not resend the segment."
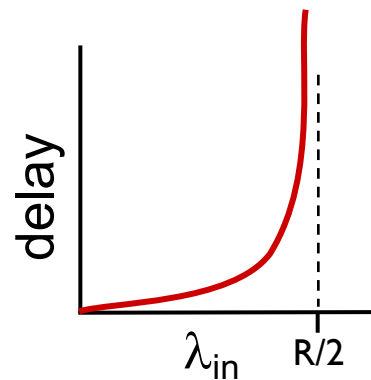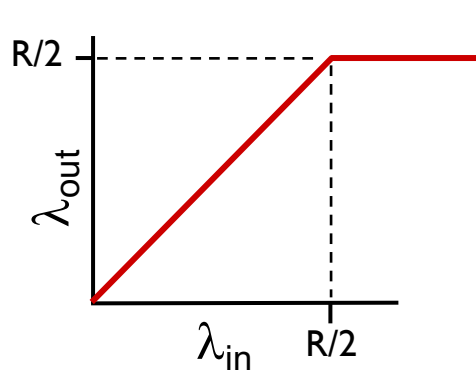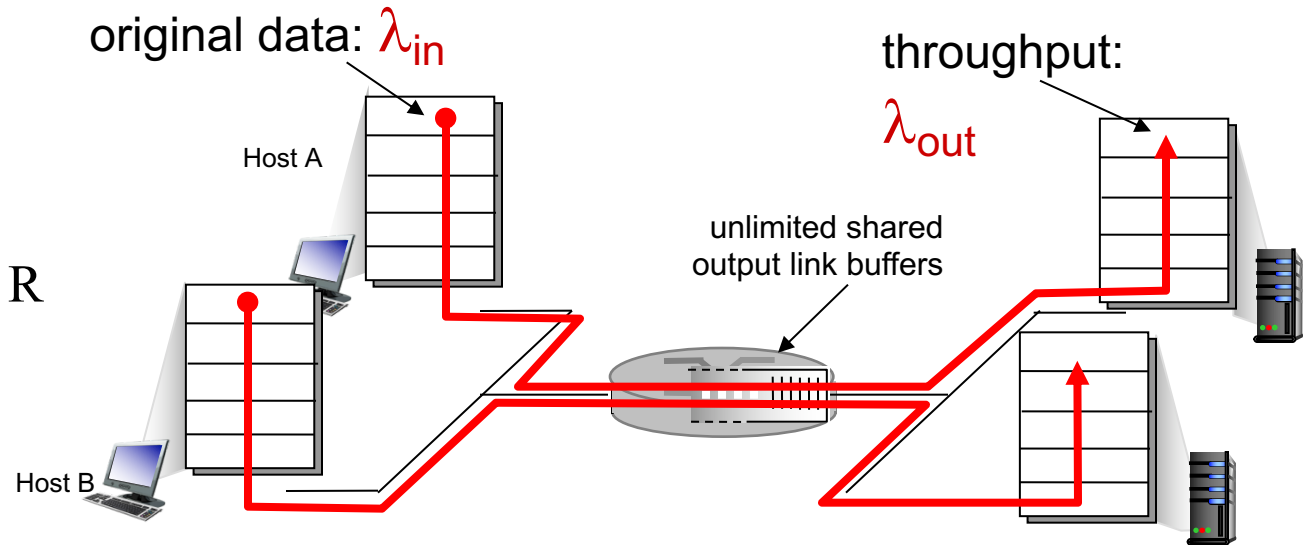
# Chapter 3 outline

# 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: $\lambda_{out}$
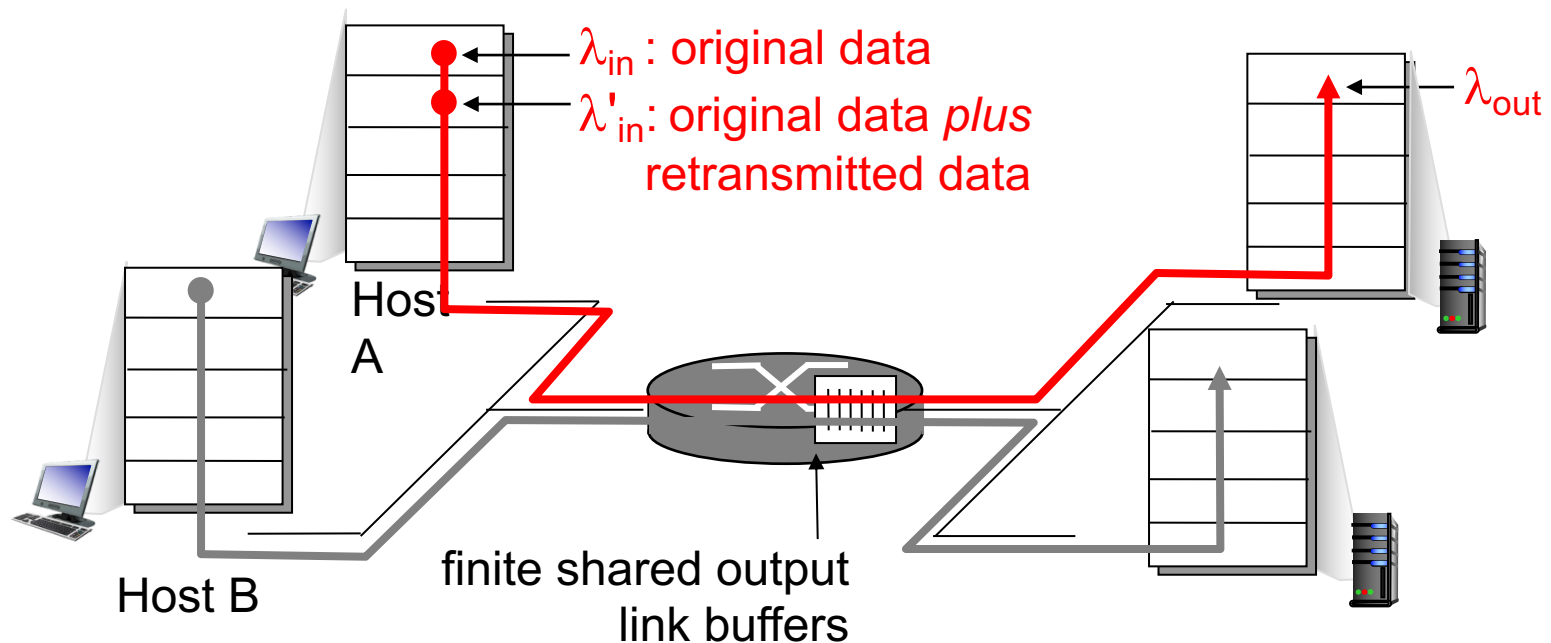
unlimited shared output link buffers

Host B



large queuing delays are experienced as the packet-arrival rate nears the link capacity.

- maximum per-connection throughput: R/2

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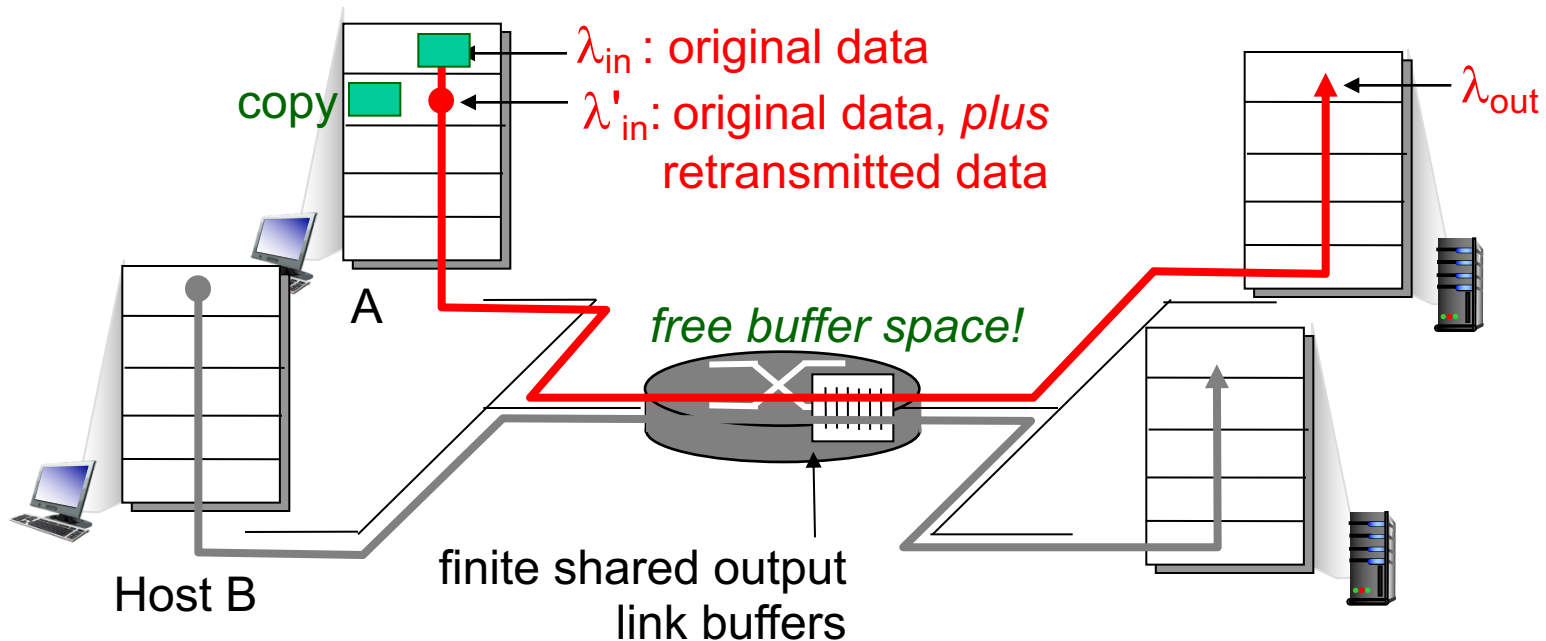
# Causes/costs of congestion: scenario 2

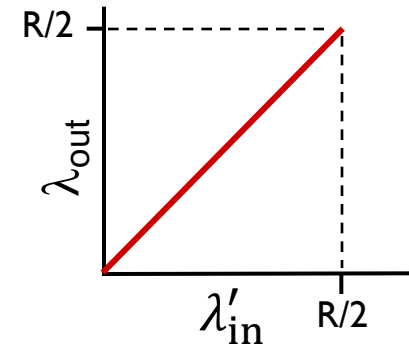❖ one router, *finite* buffers
❖ sender retransmission of timed-out packet
  ▪ application-layer input = application-layer output: $\lambda_{in} = \lambda_{out}$
  ▪ transport-layer input includes *retransmissions* : $\lambda'_{in} \geq \lambda_{in}$



$\lambda_{in}$ : 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

❖ No loss occurs: $\lambda'_{in} = \lambda_{in}$



$\lambda_{in}$ : original data

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

copy

A

free buffer space!

Host B

finite shared output link buffers

$\lambda_{out}$

# 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

❖ $\lambda'_{in} \geq \lambda_{in}$



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

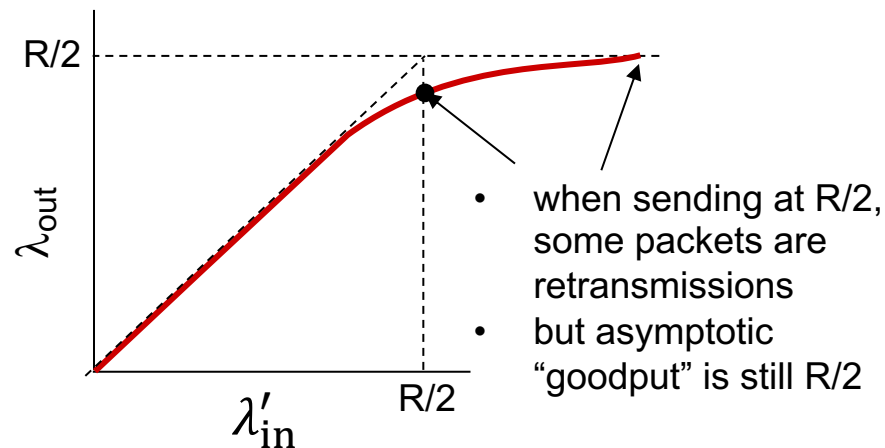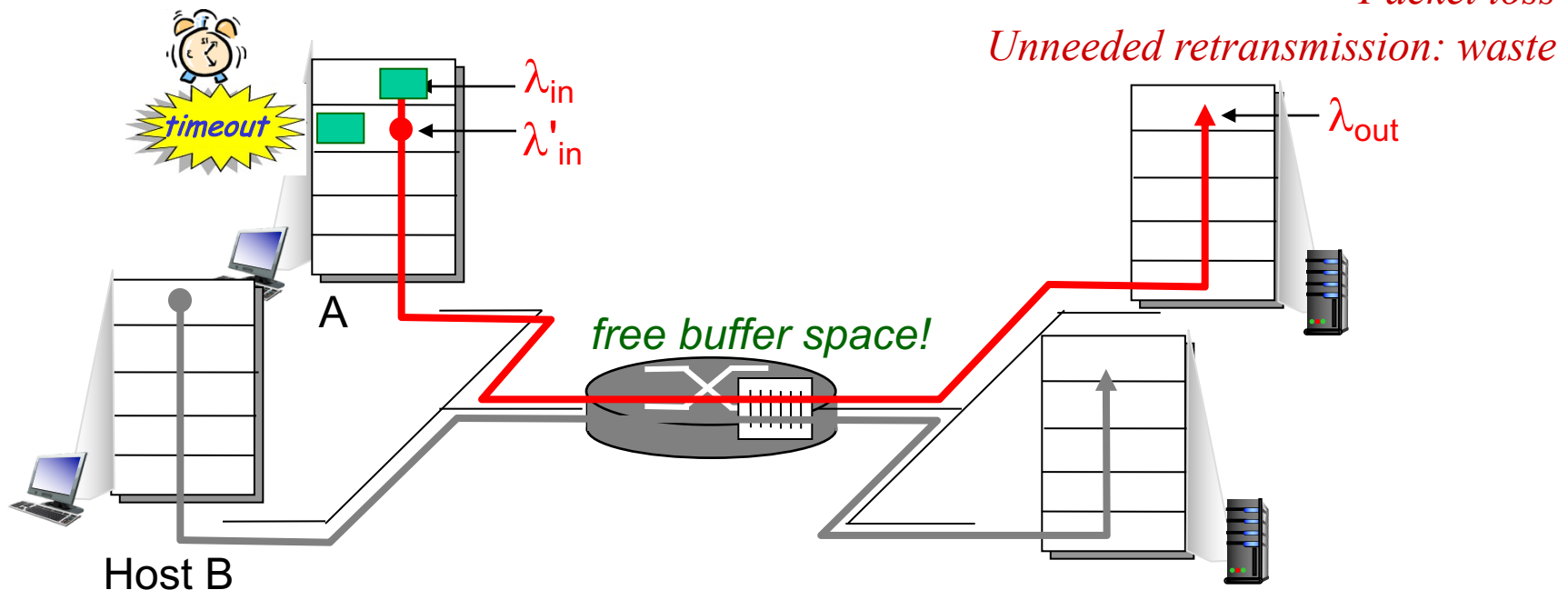❖ sender only resends if packet *known* to be lost

❖ $\lambda'_{in} \geq \lambda_{in}$



• when sending at R/2, some packets are retransmissions
• but asymptotic "goodput" is still R/2

# 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!

$\lambda_{out}$

$\lambda'_{in}$      R/2

R/2

*Packet loss*

*Unneeded retransmission: waste*

timeout

$\lambda_{in}$

$\lambda'_{in}$

$\lambda_{out}$
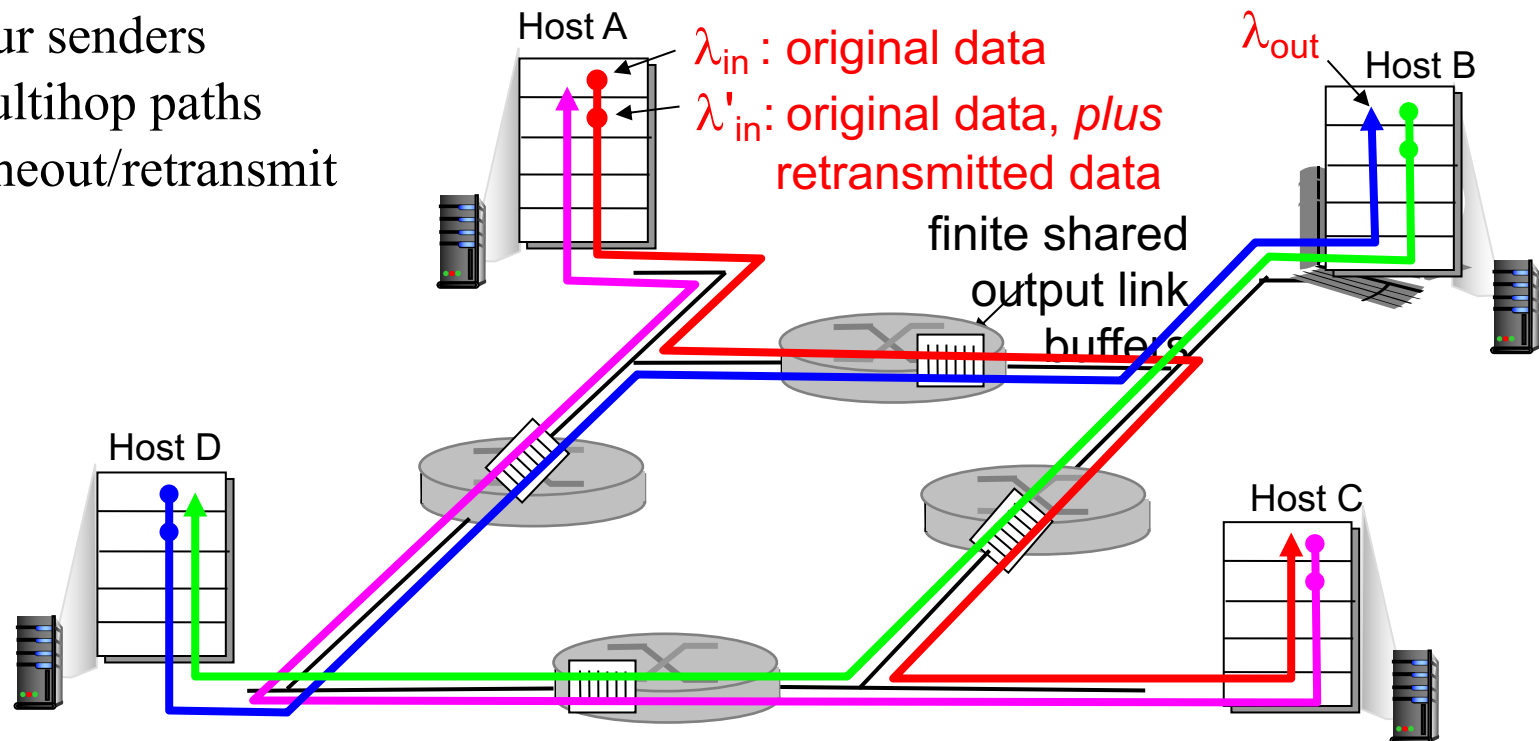
A

*free buffer space!*

Host B

# Causes/costs of congestion: scenario 3

For small values of $\lambda_{in}$:
- buffer overflows are rare
- the throughput $\lambda_{out}$ approximately equals the offered load $\lambda'_{in} = \lambda_{in}$.
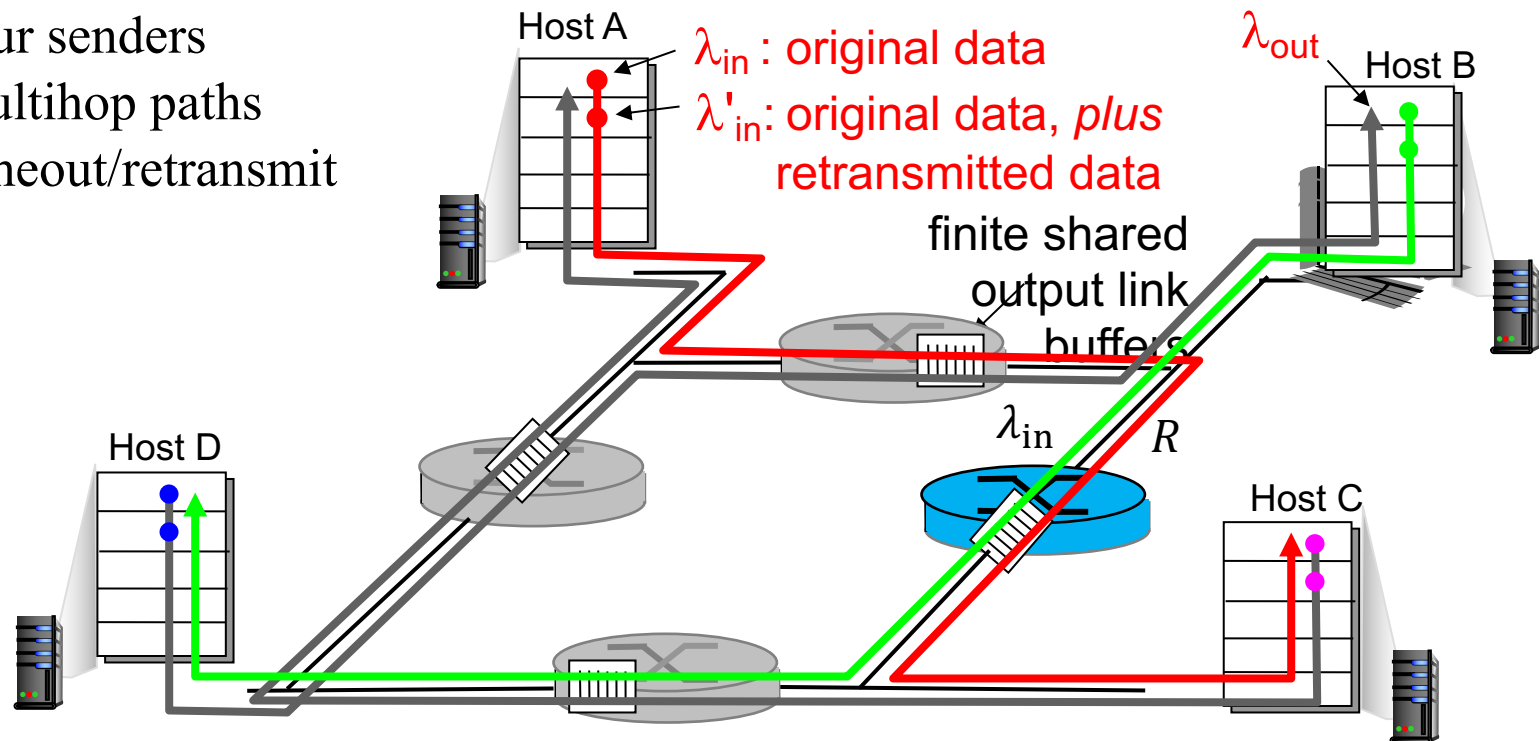
- ❖ four senders
- ❖ multihop paths
- ❖ timeout/retransmit



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

Q: what happens as $\lambda_{in}$ increases ?
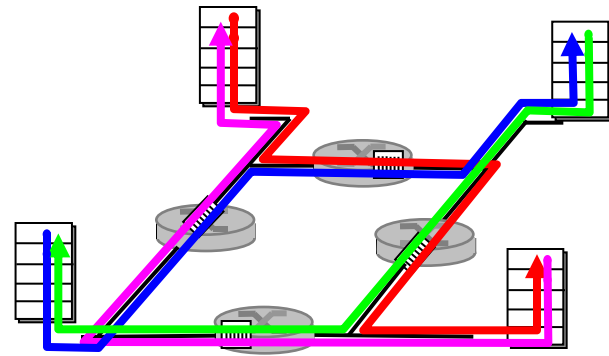
A: as green $\lambda_{in}$ increases, all arriving red pkts at upper queue are dropped, red throughput goes 0

- four senders
- multihop paths
- timeout/retransmit



Host A

$\lambda_{in}$ : original data

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

$\lambda_{out}$

Host B

finite shared output link buffers

$\lambda_{in}$     R

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!

# Cause and Cost of Congestion

**Cause**
- Shared link; limited link capacity
- Sending at a high rate

**Cost of Congestion**
- Delay
- Packet lost and retransmission
- Unneeded retransmission: waste
- "upstream" transmission capacity was wasted

# Approaches to Congestion Control

**End-to-end congestion control:**
- TCP segment loss or round-trip segment delay
- TCP decreases its window size accordingly

**Network-assisted congestion control:**
- routers provide feedback to the sender and/or receiver
- a single bit indicating congestion at a link; the maximum host sending rate the router can support

# Chapter 3 outline

# TCP congestion control: additive increase multiplicative decrease

- TCP use end-to-end congestion control
- have each sender limit the rate at which it sends traffic into its connection as a function of perceived network congestion

**Questions for achieving congestion control:**

Q1: How does a TCP sender limit the rate at which it sends traffic into its connection?

Q2: How does a TCP sender perceive that there is congestion on the path between itself and the destination?

Q3: What algorithm should the sender use to change its send rate as a function of perceived end-to-end congestion?

# TCP congestion control: additive increase multiplicative decrease

**Questions for achieving congestion control:**

Q1: How does a TCP sender limit the rate at which it sends traffic into its connection?

Congestion window: **cwnd**

LastByteSent – LastByteAcked ≤ min{cwnd, rwnd}

Q2: How does a TCP sender perceive that there is congestion on the path between itself and the destination?

Timeout; three duplicate ACKs

Q3: What algorithm should the sender use to change its send rate as a function of perceived end-to-end congestion?

# Congestion window

*sender sequence number space*



last byte ACKed

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

last byte sent



sender    receiver

RTT

❖ sender limits transmission:

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

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

*TCP sending rate:*

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

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

# TCP congestion control: additive increase multiplicative decrease

Q3: What algorithm should the sender use to change its send rate as a function of perceived end-to-end congestion?

- A lost segment → congestion → decrease rate
- An acknowledged segment → the network is fine → increase rate
- Bandwidth probing: network condition may change

# TCP Congestion Control: details

The congestion control algorithm has three major components:

❖ Slow start : exponentially increase

❖ Congestion avoidance: linearly increase

❖ Fast recovery:

# TCP Slow Start

❖ when connection begins or timeout occurs, increase rate exponentially until packet lost:
- initially **cwnd** = 1 MSS
- double **cwnd** every RTT
- done by incrementing **cwnd** for every ACK received

❖ Summary: initial rate is slow but ramps up exponentially fast

Host A

Host B

RTT

one segment

two segments

four segments

time

# TCP Congestion Control: FSM



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

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

duplicate ACK
dupACKcount++

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

cwnd≥ssthresh
Λ

**Slow start**

**Congestion avoidance**

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

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

new ACK
cwnd=ssthresh
dupACKcount=0

**when *cwnd ≥ ssthresh*: slow start → congestion avoidance**

dupACKcount==3
ssthresh=cwnd/2
cwnd=ssthresh+3·MSS
retransmit missing segment

**Fast recovery**

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

# TCP: detecting, reacting to loss

❖ loss indicated by timeout:
  ▪ **cwnd** set to 1 MSS; **ssthresh** = **cwnd**/2
  ▪ window then grows exponentially (as in slow start) to threshold, then grows linearly

❖ loss indicated by 3 duplicate ACKs:
  ▪ TCP Tahoe always sets **cwnd** to 1 (timeout or 3 duplicate acks) → Slow start
    • **cwnd** set to 1 MSS; **ssthresh** = **cwnd**/2
  ▪ TCP RENO
    • dup ACKs indicate network capable of delivering some segments → Fast Recovery
    • ssthresh = cwnd / 2; cwnd = ssthresh + 3MSS

# TCP: switching 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 timeout

**Thus, when timeout occurs, ssthresh = cwnd/2**

# TCP Congestion Avoidance

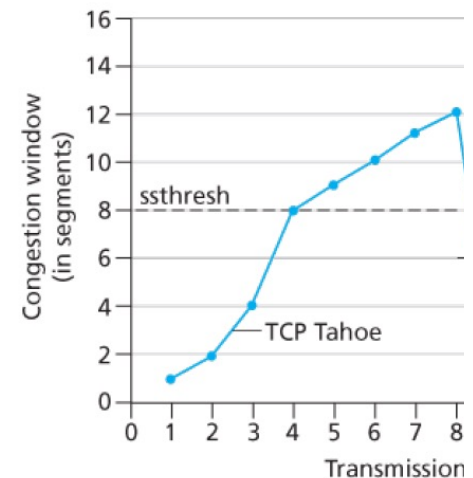Trigger: $cwnd \geq ssthresh$

Increases $cwnd$ linearly: by one MSS every RTT

❖ Increase $cwnd$ by (MSS/$cwnd$)MSS bytes whenever a new acknowledgment arrives.

❖ E.g., if MSS is 1,460 bytes and $cwnd$ is 14,600 bytes, then 10 segments are being sent within an RTT.

  ▪ Each arriving ACK (assuming one ACK per segment) increases the congestion window size by 1/10 MSS,

# TCP Congestion Control: FSM

# TCP Fast Recovery

- Trigger（RENO）: triple duplicate ACKs
- *ssthresh = cwnd / 2; cwnd = ssthresh + 3MSS*

- The value of *cwnd* is increased by 1 MSS for every duplicate ACK received for the missing segment that caused TCP to enter the fast-recovery state
- when an ACK arrives for the missing segment, enter congestion avoidance

# Summary: TCP Congestion Control

# Summary: TCP Congestion Control

duplicate ACK
---
dupACKcount++

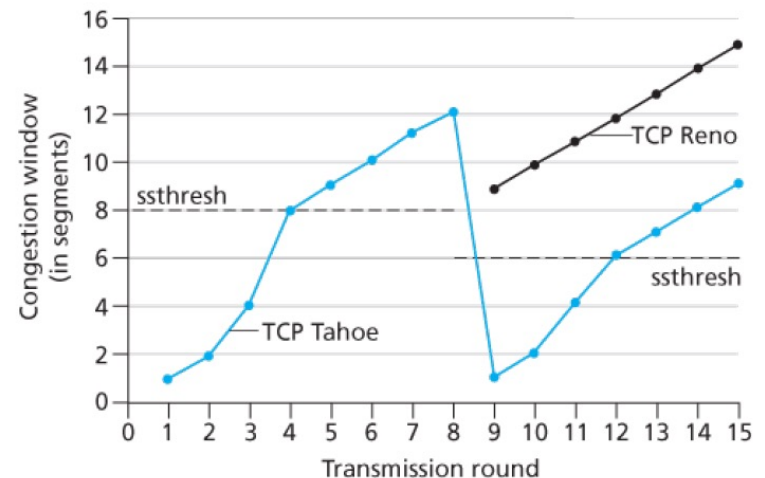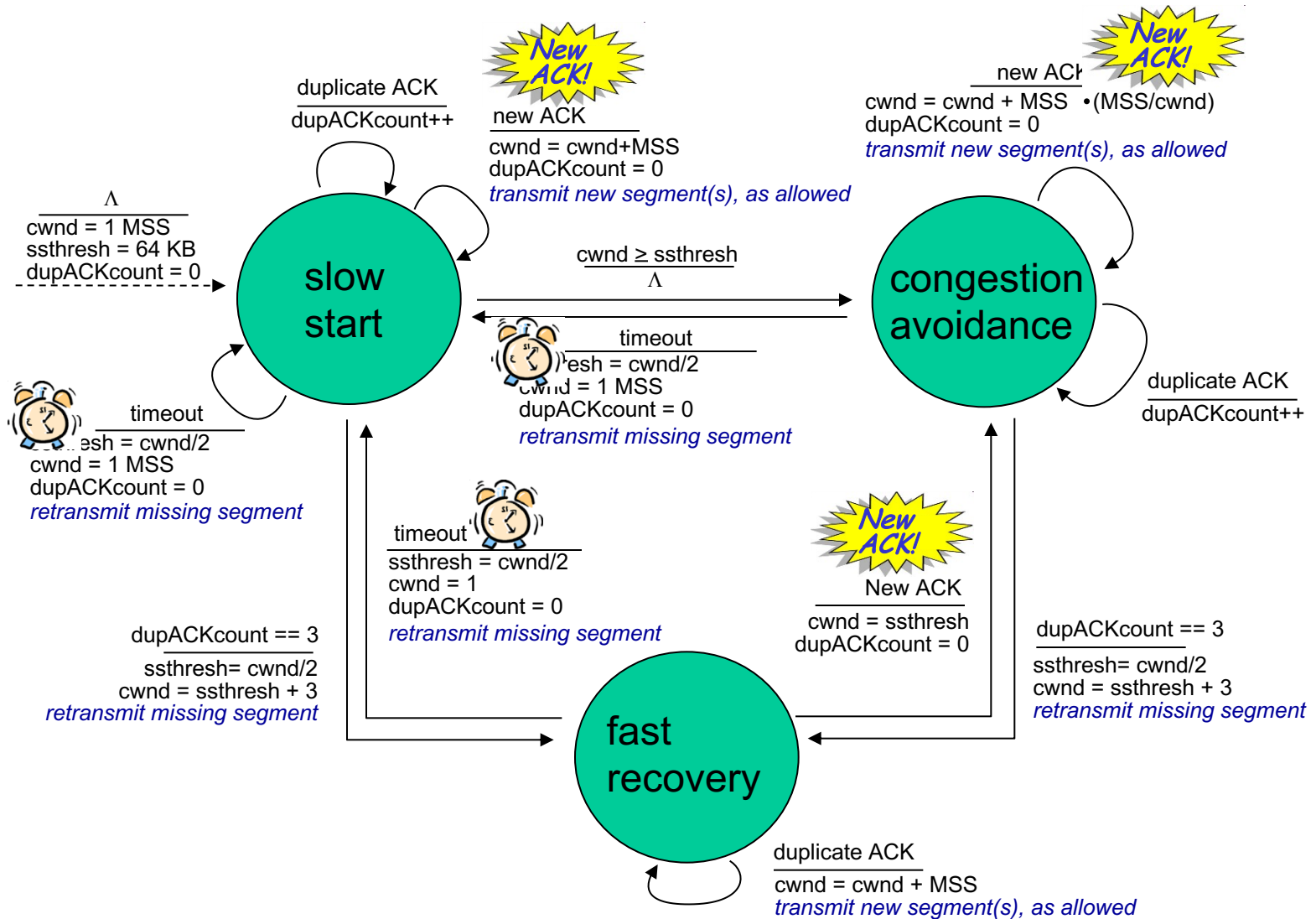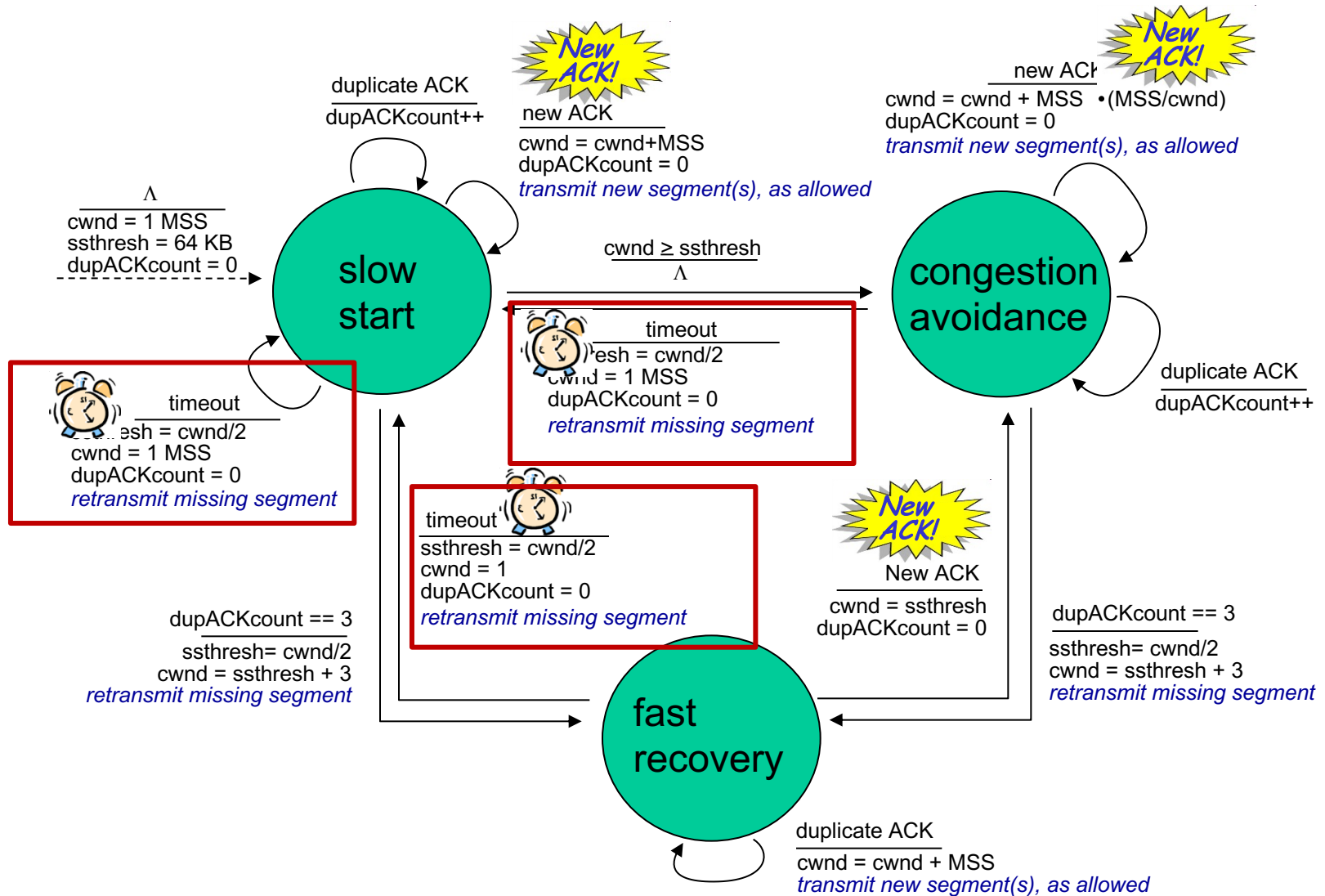**New ACK!**

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
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*

# Summary: TCP Congestion Control

# Summary: TCP Congestion Control



**New ACK!**

duplicate ACK
———————
dupACKcount++

**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

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

**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
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*

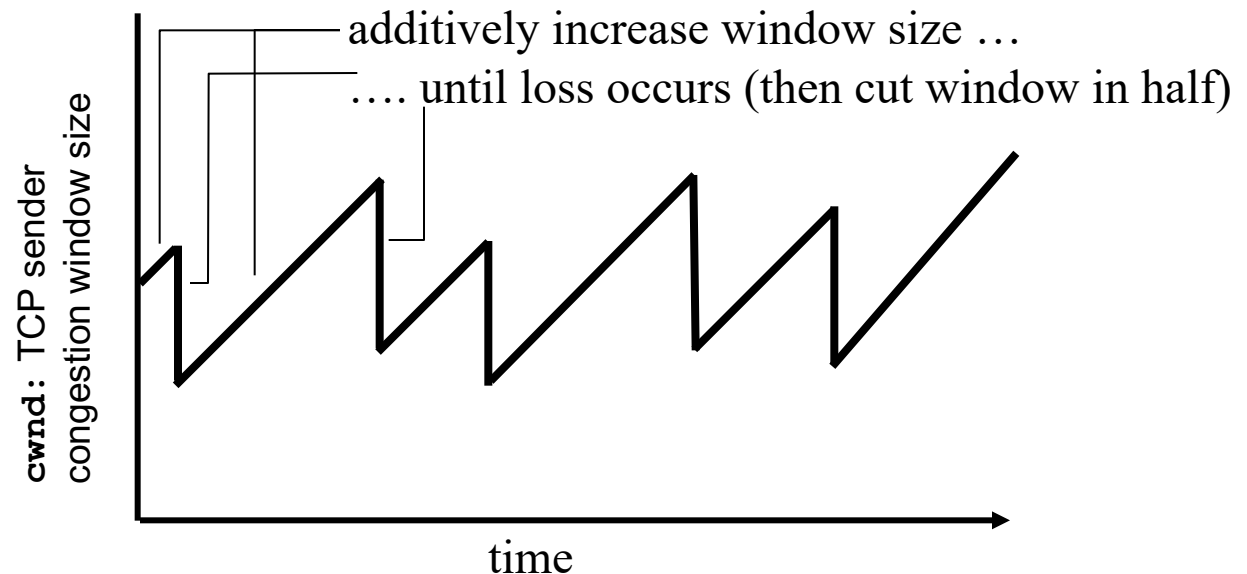# TCP congestion control: additive increase multiplicative decrease

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

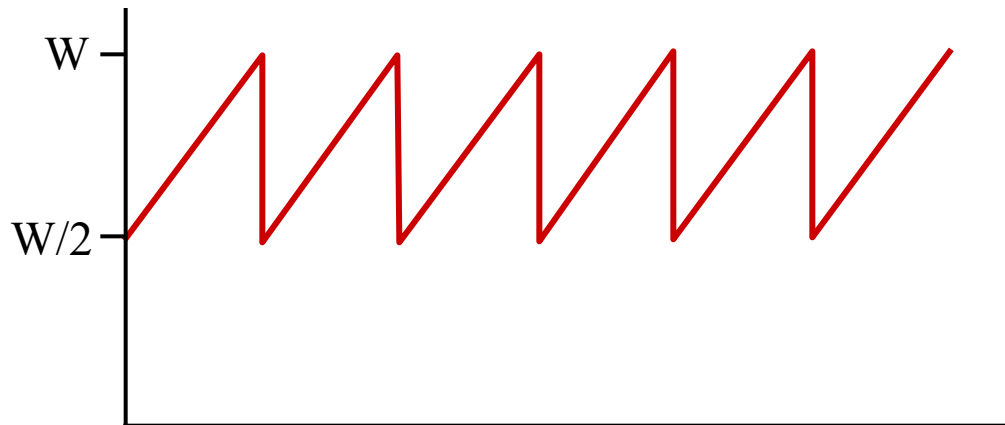- *multiplicative decrease*: cut **cwnd** in half after loss

AIMD: probing for bandwidth

additively increase window size …

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

cwnd: TCP sender congestion window size

time

# TCP throughput

❖ Avg. TCP throughput 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. throughput 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"

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

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

❖ requires W = 83,333 in-flight segments

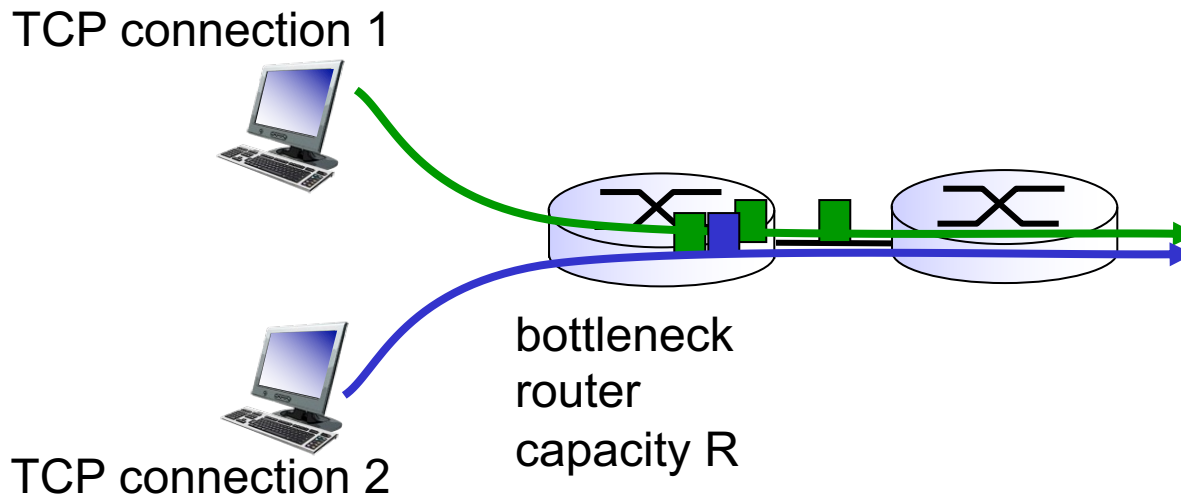❖ throughput in terms of segment loss probability, L [Mathis 1997]:

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

➜ to achieve 10 Gbps throughput, need a loss rate of L = $2 \cdot 10^{-10}$ *– a very small loss rate!*
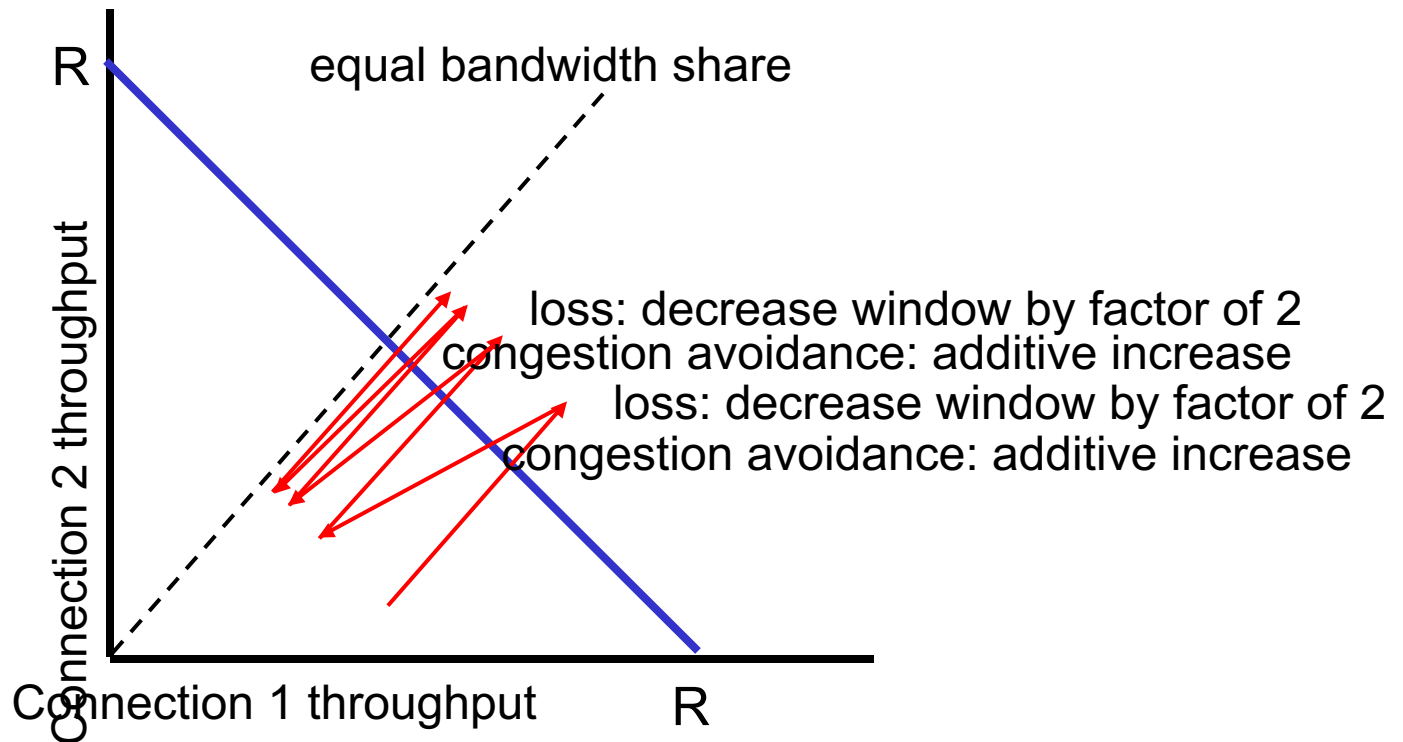
❖ new versions of TCP for high-speed

# 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

R 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

  ▪ do not want rate throttled by congestion control

❖ use UDP:

  ▪ send audio/video at constant rate, tolerate packet loss
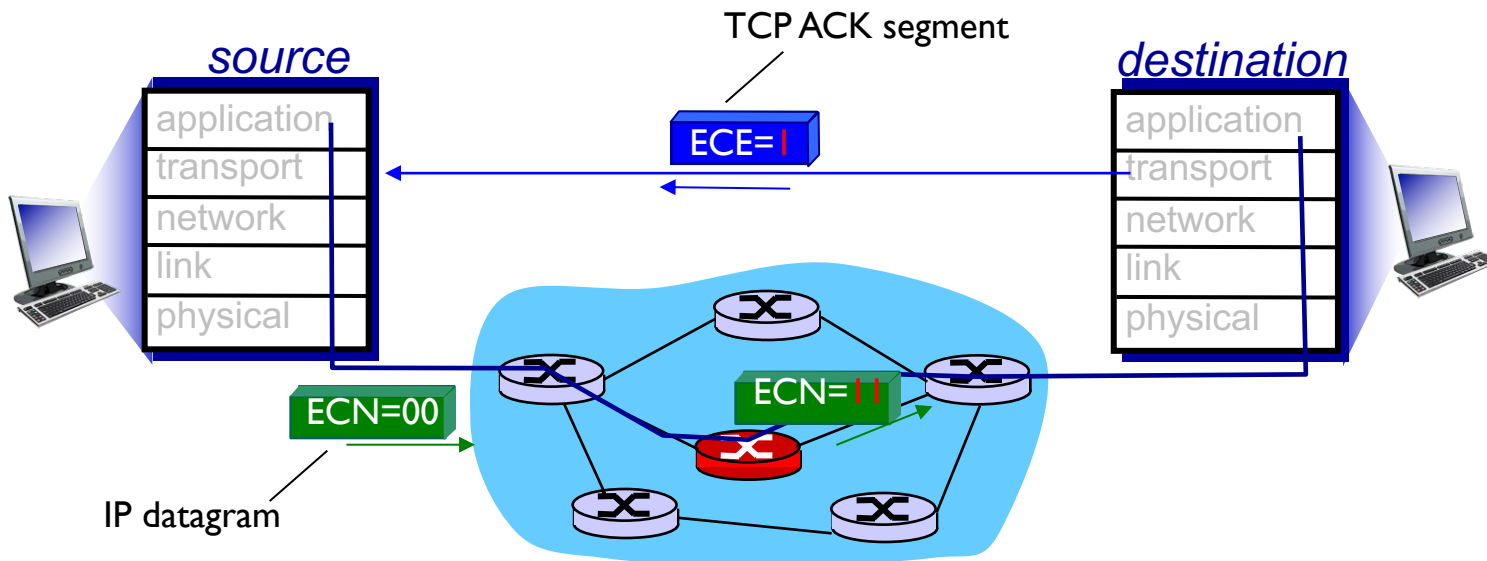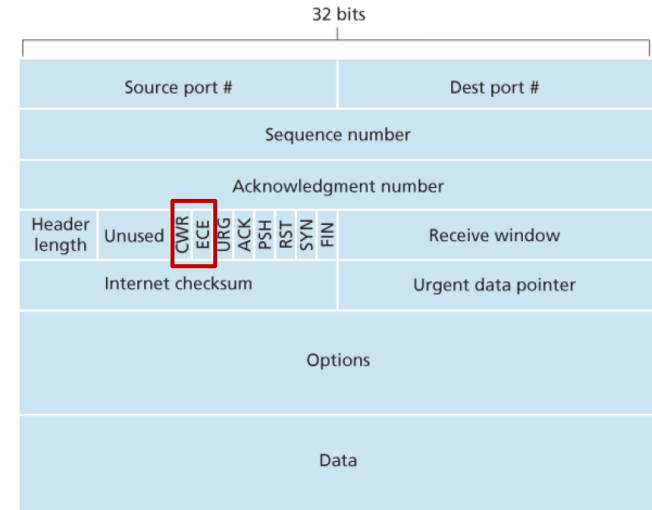
❖ UDP sources to crowd out TCP traffic

*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 more than R/2

# Explicit Congestion Notification (ECN)

*network-assisted congestion control:*

- two bits in IP header (ToS field) marked *by network router* to indicate congestion

- congestion indication carried to receiving host

- receiver sets ECE bit on receiver-to-sender ACK segment to notify sender of congestion

# 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"