# TCP: Overview     RFCs: 793, 1122, 1323, 2018, 2581
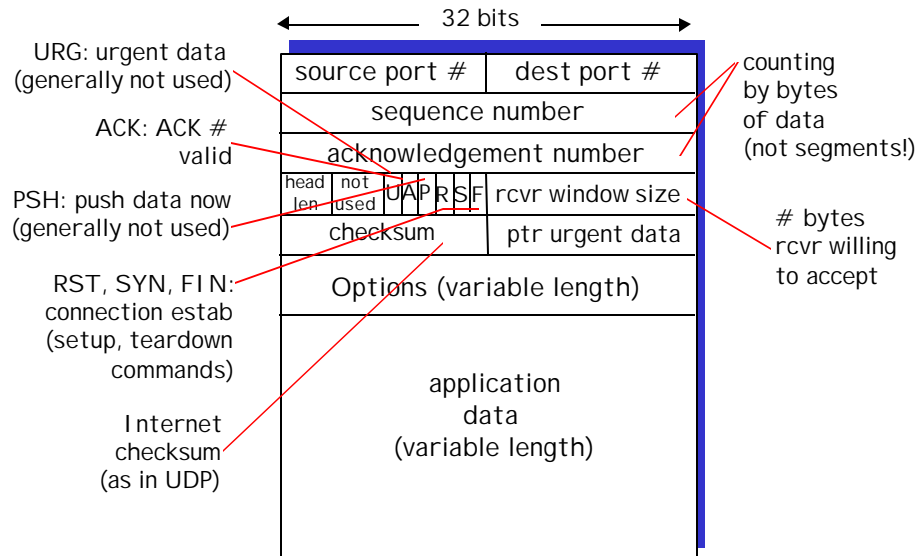
r **point-to-point:**
  - m one sender, one receiver

r **reliable, in-order *byte steam:***
  - m no "message boundaries"

r **pipelined:**
  - m TCP congestion and flow control set window size

r ***send & receive buffers***

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

r **connection-oriented:**
  - m handshaking (exchange of control msgs) init's sender, receiver state before data exchange

r **flow controlled:**
  - m sender will not overwhelm receiver

3: Transport Layer     3b-1

---

# TCP segment structure

32 bits

URG: urgent data (generally not used)

ACK: ACK # valid

PSH: push data now (generally not used)

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

Internet checksum (as in UDP)

| source port # | dest port # |
| sequence number | |
| acknowledgement number | |
| head len | not used | U A P R S F | rcvr window size |
| checksum | ptr urgent data |
| Options (variable length) | |
| application data (variable length) | |

counting by bytes of data (not segments!)
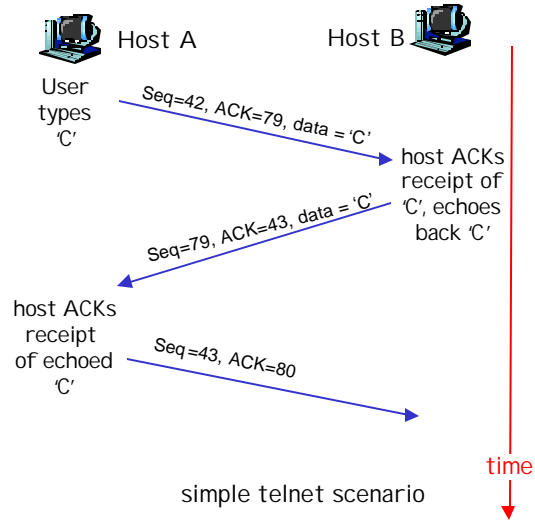
# bytes rcvr willing to accept

3: Transport Layer     3b-2

# TCP seq. #'s and ACKs

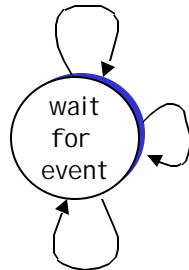Seq. #'s:
- m byte stream "number" of first byte in segment's data

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

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

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

time

simple telnet scenario

3: Transport Layer    3b-3

# TCP: reliable data transfer

event: data received from application above

create, send segment

wait for event

event: timer timeout for segment with seq # y

retransmit segment

event: ACK received, with ACK # y

ACK processing

simplified sender, assuming
- •one way data transfer
- •no flow, congestion control

3: Transport Layer    3b-4

# TCP: reliable data transfer

Simplified
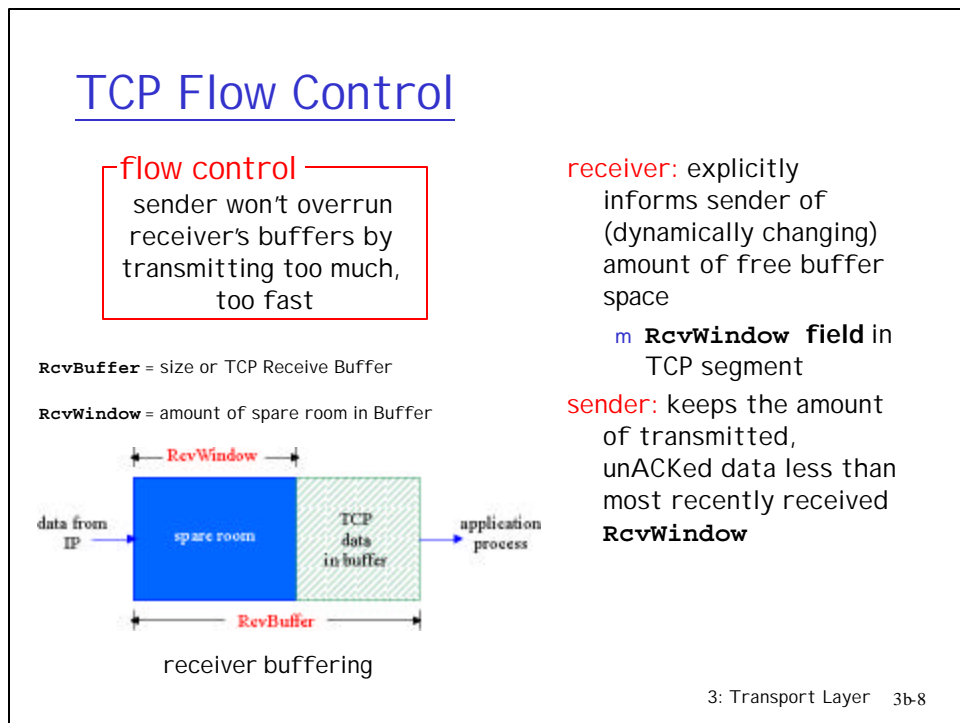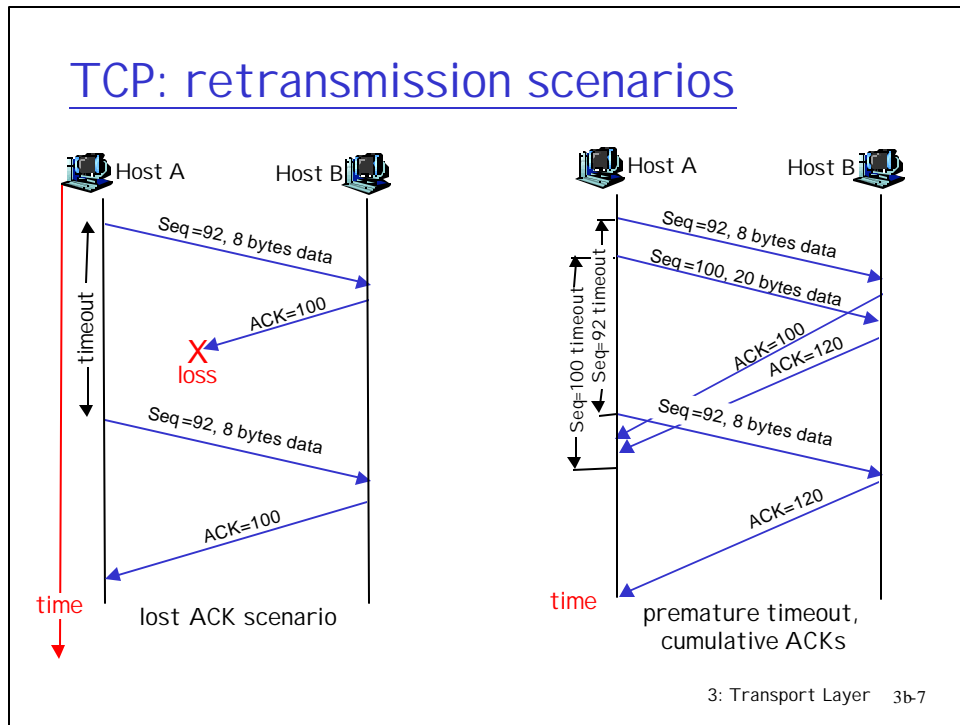TCP
sender

```
00   sendbase = initial_sequence number
01   nextseqnum = initial_sequence number
02
03   loop (forever) {
04     switch(event)
05     event: data received from application above
06         create TCP segment with sequence number nextseqnum
07         start timer for segment  nextseqnum
08         pass segment to IP
09         nextseqnum = nextseqnum + length(data)
10     event: timer timeout for segment with sequence number y
11         retransmit segment with sequence number y
12         compue new timeout interval for segment y
13         restart timer for sequence number y
14     event: ACK received, with ACK field value of y
15         if (y > sendbase) { /* cumulative ACK of all data up to y */
16             cancel all timers for segments with sequence numbers < y
17             sendbase = y
18             }
19         else { /* a duplicate ACK for already ACKed segment */
20             increment number of duplicate ACKs received for y
21             if (number of duplicate ACKS received for y == 3) {
22                 /* TCP fast retransmit */
23                 resend segment with sequence number y
24                 restart timer for segment y
25             }
26   } /* end of loop forever */
```

3: Transport Layer    3b-5

# TCP ACK generation [RFC 1122, RFC 2581]

| Event | TCP Receiver action |
|---|---|
| in-order segment arrival, no gaps, everything else already ACKed | delayed ACK. Wait up to 500ms for next segment. If no next segment, send ACK |
| in-order segment arrival, no gaps, one delayed ACK pending | immediately send single cumulative ACK |
| out-of-order segment arrival higher-than-expect seq. # gap detected | send duplicate ACK, indicating seq. # of next expected byte |
| arrival of segment that partially or completely fills gap | immediate ACK if segment starts at lower end of gap |

3: Transport Layer    3b-6

# TCP: retransmission scenarios



lost ACK scenario

premature timeout, cumulative ACKs

3: Transport Layer    3b-7

# TCP Flow Control

flow control

sender won't overrun receiver's buffers by transmitting too much, too fast

**RcvBuffer** = size or TCP Receive Buffer

**RcvWindow** = amount of spare room in Buffer



receiver buffering

receiver: explicitly informs sender of (dynamically changing) amount of free buffer space

- m **RcvWindow field** in TCP segment

sender: keeps the amount of transmitted, unACKed data less than most recently received **RcvWindow**

3: Transport Layer    3b-8

# TCP Round Trip Time and Timeout

**Q:** how to set TCP timeout value?

r  longer than RTT

   m  note: RTT will vary

r  too short: premature timeout

   m  unnecessary retransmissions

r  too long: slow reaction to segment loss

**Q:** how to estimate RTT?

r  **SampleRTT**: measured time from segment transmission until ACK receipt

   m  ignore retransmissions, cumulatively ACKed segments

r  **SampleRTT** will vary, want estimated RTT "smoother"

   m  use several recent measurements, not just current **SampleRTT**

3: Transport Layer   3b-9

# TCP Round Trip Time and Timeout

**EstimatedRTT = (1-x)*EstimatedRTT + x*SampleRTT**

   r  Exponential weighted moving average

   r  influence of given sample decreases exponentially fast

   r  typical value of x: 0.1

## Setting the timeout

r  **EstimtedRTT** plus "safety margin"

r  large variation in **EstimatedRTT ->** larger safety margin

     **Timeout = EstimatedRTT + 4*Deviation**

    **Deviation = (1-x)*Deviation +**
               **x*|SampleRTT-EstimatedRTT|**

3: Transport Layer   3b-10

## TCP Connection Management

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

r   initialize TCP variables:

   m   seq. #s

   m   buffers, flow control info (e.g. **RcvWindow**)

r   *client:* connection initiator
   **Socket clientSocket = new Socket("hostname","port number");**

r   *server:* contacted by client
   **Socket connectionSocket = welcomeSocket.accept();**

### Three way handshake:

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

   m   specifies initial seq #

Step 2: server end system receives SYN, replies with SYNACK control segment

   m   ACKs received SYN

   m   allocates buffers

   m   specifies server-> receiver initial seq. #

3: Transport Layer   3b-11

---

## TCP Connection Management (cont.)

### Closing a connection:

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

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

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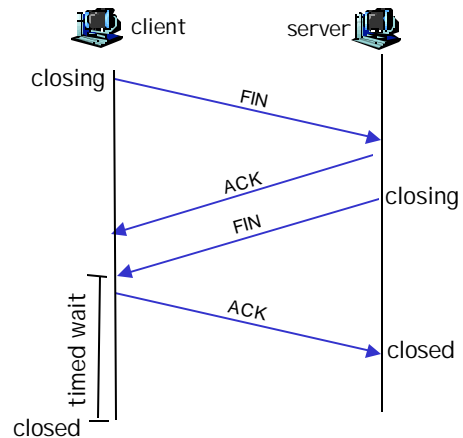


3: Transport Layer   3b-12

# TCP Connection Management (cont.)

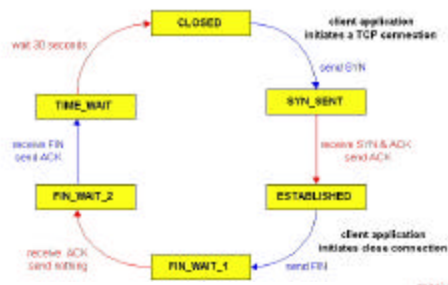**Step 3:** client receives FIN, replies with ACK.

- m Enters "timed wait" - will respond with ACK to received FINs
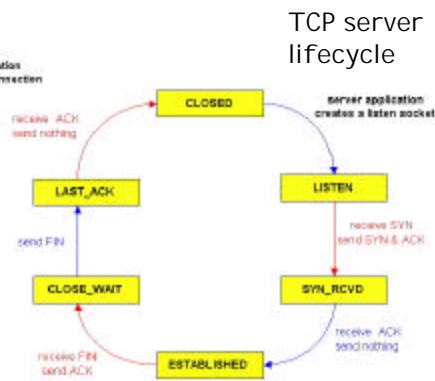
**Step 4:** server, receives ACK. Connection closed.

**Note:** with small modification, can handly simultaneous FINs.



client       server

closing

FIN

ACK

closing

FIN

timed wait

ACK

closed

closed

3: Transport Layer  3b-13

# TCP Connection Management (cont)



TCP client lifecycle

TCP server lifecycle

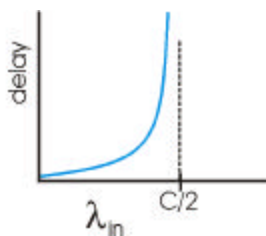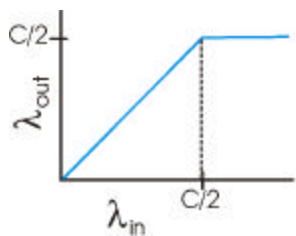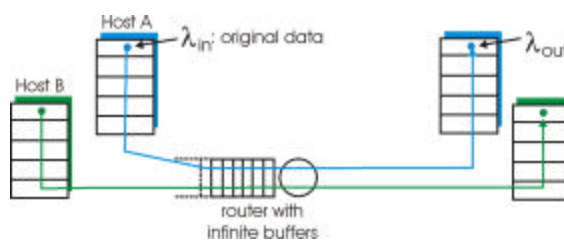3: Transport Layer  3b-14

# Principles of Congestion Control

Congestion:
- r informally: "too many sources sending too much data too fast for *network* to handle"
- r different from flow control!
- r manifestations:
    - m lost packets (buffer overflow at routers)
    - m long delays (queueing in router buffers)
- r a top-10 problem!

3: Transport Layer  3b-15

---

# Causes/costs of congestion: scenario 1

- r two senders, two receivers
- r one router, infinite buffers
- r no retransmission



Host A  $\lambda_{in}$: original data  $\lambda_{out}$
Host B
router with infinite buffers



$C/2$
$\lambda_{out}$
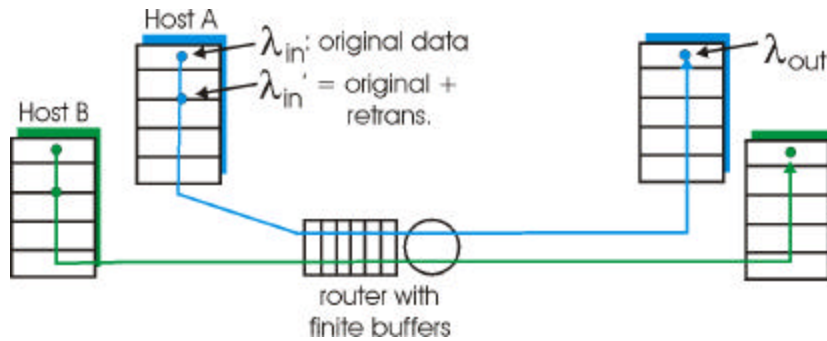$\lambda_{in}$
$C/2$

delay
$\lambda_{in}$
$C/2$

- r large delays when congested
- r maximum achievable throughput

3: Transport Layer  3b-16

## Causes/costs of congestion: scenario 2

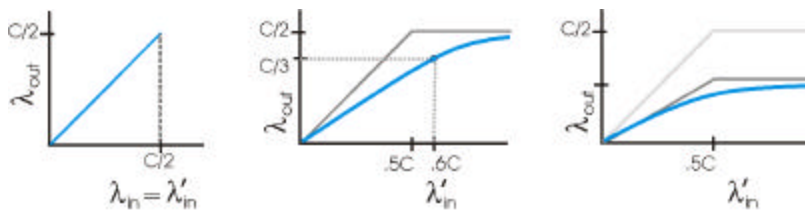r   one router, *finite* buffers
r   sender retransmission of lost packet



Host A — $\lambda_{in}$: original data
$\lambda_{in}' =$ original + retrans.

Host B

$\lambda_{out}$

router with finite buffers

## Causes/costs of congestion: scenario 2

r   always:  $\lambda_{in} = \lambda_{out}$  (goodput)
r   "perfect" retransmission only when loss:  $\lambda_{in}' > \lambda_{out}$
r   retransmission of delayed (not lost) packet makes  $\lambda_{in}'$  larger
    (than perfect case) for same  $\lambda_{out}$
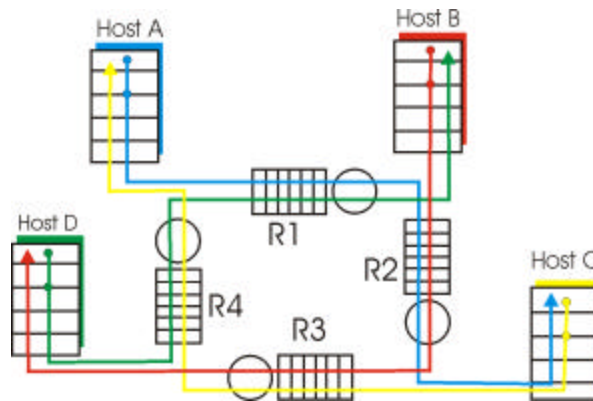


"costs" of congestion:
r   more work (retrans) for given "goodput"
r   unneeded retransmissions: link carries multiple copies of pkt
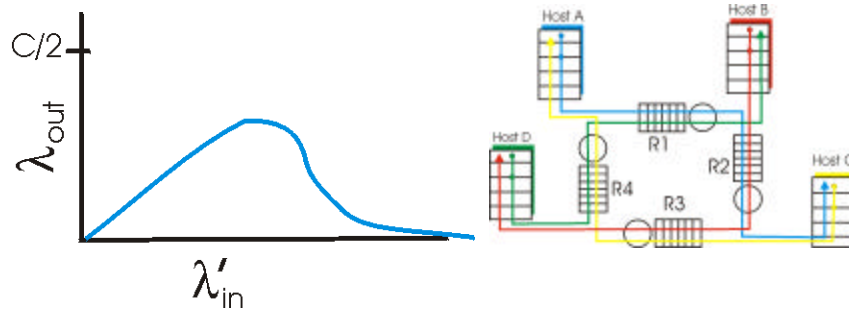
## Causes/costs of congestion: scenario 3

r   four senders
r   multihop paths
r   timeout/retransmit

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



3: Transport Layer 3b-19

## Causes/costs of congestion: scenario 3



Another "cost" of congestion:

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

3: Transport Layer 3b-20

## Approaches towards congestion control

Two broad approaches towards congestion control:

**End-end congestion control:**

r no explicit feedback from network

r congestion inferred from end-system observed loss, delay

r approach taken by TCP

**Network-assisted congestion control:**

r routers provide feedback to end systems

  m single bit indicating congestion (SNA, DECbit, TCP/IP ECN, ATM)

  m explicit rate sender should send at

3: Transport Layer  3b-21

## Case study: ATM ABR congestion control

**ABR: available bit rate:**

r "elastic service"

r if sender's path "underloaded":

  m sender should use available bandwidth

r if sender's path congested:

  m sender throttled to minimum guaranteed rate

**RM (resource management) cells:**

r sent by sender, interspersed with data cells

r bits in RM cell set by switches ("network-assisted")

  m NI bit: no increase in rate (mild congestion)

  m CI bit: congestion indication

r RM cells returned to sender by receiver, with bits intact

3: Transport Layer  3b-22

## Case study: ATM ABR congestion control



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

3: Transport Layer   3b-23

---

# TCP Congestion Control

r end-end control (no network assistance)
r transmission rate limited by congestion window size, **Congwin**, over segments:



r w segments, each with MSS bytes sent in one RTT:

$$\text{throughput} = \frac{w * MSS}{RTT} \text{ Bytes/sec}$$

3: Transport Layer   3b-24

# TCP congestion control:

r  *"probing"* for usable bandwidth:

  m  ideally: transmit as fast as possible (**Congwin** as large as possible) without loss

  m  *increase* **Congwin** until loss (congestion)

  m  loss: *decrease* **Congwin**, then begin probing (increasing) again

r  two *"phases"*

  m  slow start

  m  congestion avoidance

r  important variables:

  m  **Congwin**

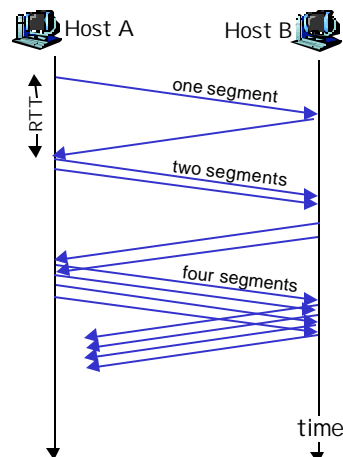  m  **threshold:** defines threshold between two slow start phase, congestion control phase

3: Transport Layer  3b-25

# TCP Slowstart

Slowstart algorithm

initialize: Congwin = 1
for (each segment ACKed)
     Congwin++
until (loss event OR
     CongWin > threshold)

r  exponential increase (per RTT) in window size (not so slow!)

r  loss event: timeout (Tahoe TCP) and/or or three duplicate ACKs (Reno TCP)



3: Transport Layer  3b-26

# TCP Congestion Avoidance

Congestion avoidance

```
/* slowstart is over     */
/* Congwin > threshold */
Until (loss event) {
   every w segments ACKed:
      Congwin++
   }
threshold = Congwin/2
Congwin = 1
perform slowstart[1]
```



1: TCP Reno skips slowstart (fast recovery) after three duplicate ACKs

3: Transport Layer  3b-27

---

# AIMD

TCP congestion avoidance:

r **AIMD**: *additive increase, multiplicative decrease*

   m increase window by 1 per RTT
   m decrease window by factor of 2 on loss event

# TCP Fairness

Fairness goal: if N TCP sessions share same bottleneck link, each should get 1/N of link capacity



TCP connection 1

TCP connection 2

bottleneck router capacity R

3: Transport Layer  3b-28

# Why is TCP fair?

Two competing sessions:

r   Additive increase gives slope of 1, as throughout increases

r   multiplicative decrease decreases throughput proportionally

R

equal bandwidth share

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

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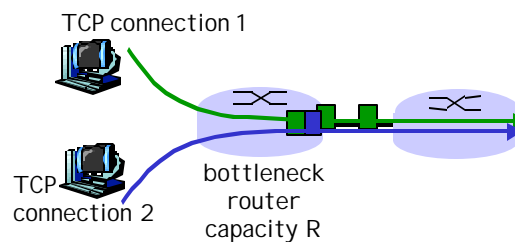congestion avoidance: additive increase

C                Connection 1 throughput   R

o

3: Transport Layer  3b-29

n

e

# TCP latency modeling

Q: How long does it take to receive an object from a Web server after sending a request?

r   TCP connection establishment

r   data transfer delay

Notation, assumptions:

r   Assume one link between client and server of rate R

r   Assume: fixed congestion window, W segments

r   S: MSS (bits)

r   O: object size (bits)

r   no retransmissions (no loss, no corruption)

Two cases to consider:

r   WS/R > RTT + S/R: ACK for first segment in window returns before window's worth of data sent

r   WS/R < RTT + S/R: wait for ACK after sending window's worth of data sent

3: Transport Layer  3b-30

# TCP latency Modeling

K:= O/WS



Case 1: latency = 2RTT + O/R

Case 2: latency = 2RTT + O/R
+ (K-1)[S/R + RTT - WS/R]

3: Transport Layer  3b-31

---

# TCP Latency Modeling: Slow Start

r   Now suppose window grows according to slow start.

r   Will show that the latency of one object of size O is:

$$Latency = 2RTT + \frac{O}{R} + P\left[RTT + \frac{S}{R}\right] - (2^P - 1)\frac{S}{R}$$

where *P* is the number of times TCP stalls at server:

$$P = \min\{Q, K - 1\}$$

- where Q is the number of times the server would stall
  if the object were of infinite size.

- and  K is the number of windows that cover the object.

3: Transport Layer  3b-32

## TCP Latency Modeling: Slow Start (cont.)

Example:

O/S = 15 segments

K = 4 windows

Q = 2

P = min{K-1,Q} = 2

Server stalls P=2 times.

initiate TCP connection

request object

RTT

first window = S/R

second window = 2S/R

third window = 4S/R

fourth window = 8S/R

object delivered

complete transmission

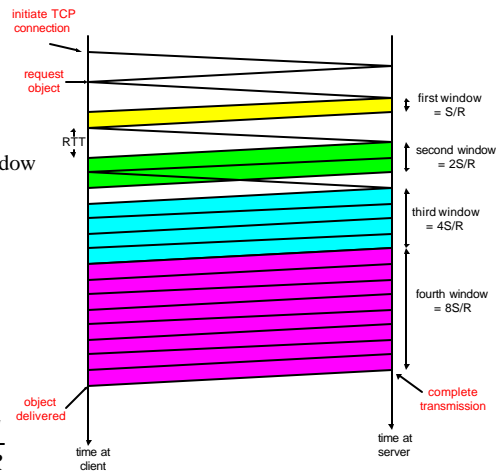time at client

time at server

3: Transport Layer   3b-33

## TCP Latency Modeling: Slow Start (cont.)

$$\frac{S}{R} + RTT = \text{time from when server starts to send segment}$$
$$\text{until server receives acknowledgement}$$

$$2^{k-1}\frac{S}{R} = \text{time to transmit the kth window}$$

$$\left[\frac{S}{R} + RTT - 2^{k-1}\frac{S}{R}\right]^{+} = \text{stall time after the kth window}$$

$$\text{latency} = \frac{O}{R} + 2RTT + \sum_{p=1}^{P} stallTime_p$$
$$= \frac{O}{R} + 2RTT + \sum_{k=1}^{P}[\frac{S}{R} + RTT - 2^{k-1}\frac{S}{R}]$$
$$= \frac{O}{R} + 2RTT + P[RTT + \frac{S}{R}] - (2^P - 1)\frac{S}{R}$$

initiate TCP connection

request object

RTT

first window = S/R

second window = 2S/R

third window = 4S/R

fourth window = 8S/R

object delivered

complete transmission

time at client

time at server

3: Transport Layer   3b-34

# Chapter 3: Summary

r  principles behind
   transport layer services:
   m  multiplexing/demultiplexing
   m  reliable data transfer
   m  flow control
   m  congestion control
r  instantiation and
   implementation in the Internet
   m  UDP
   m  TCP

Next:
r  leaving the network
   "edge" (application
   transport layer)
r  into the network "core"

3: Transport Layer  3b-35