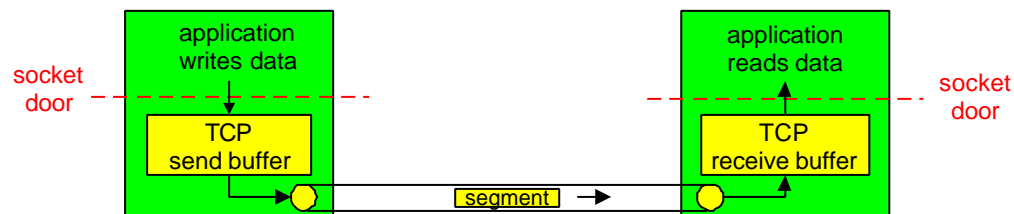


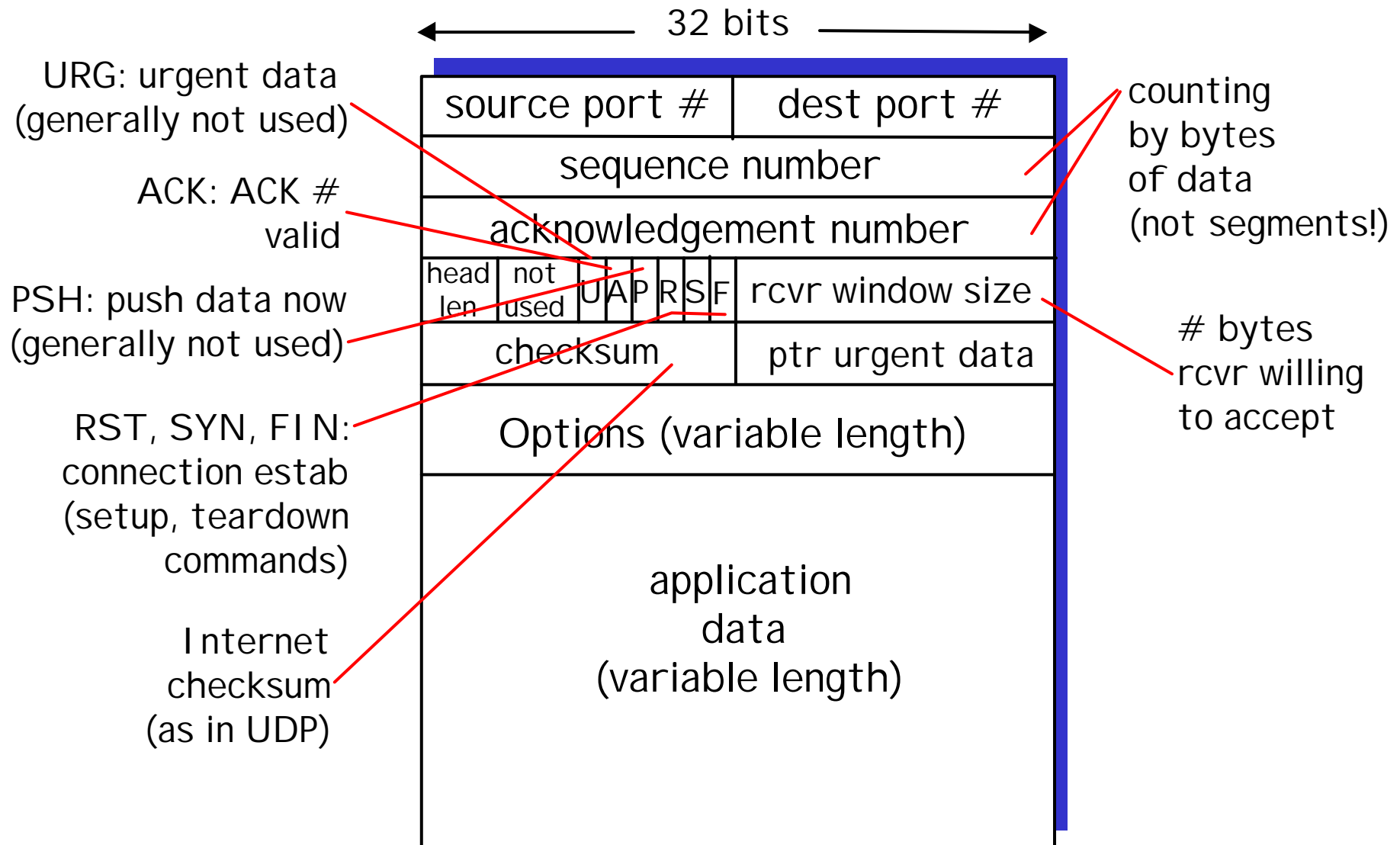
TCP: Overview

RFCs: 793, 1122, 1323, 2018, 2581

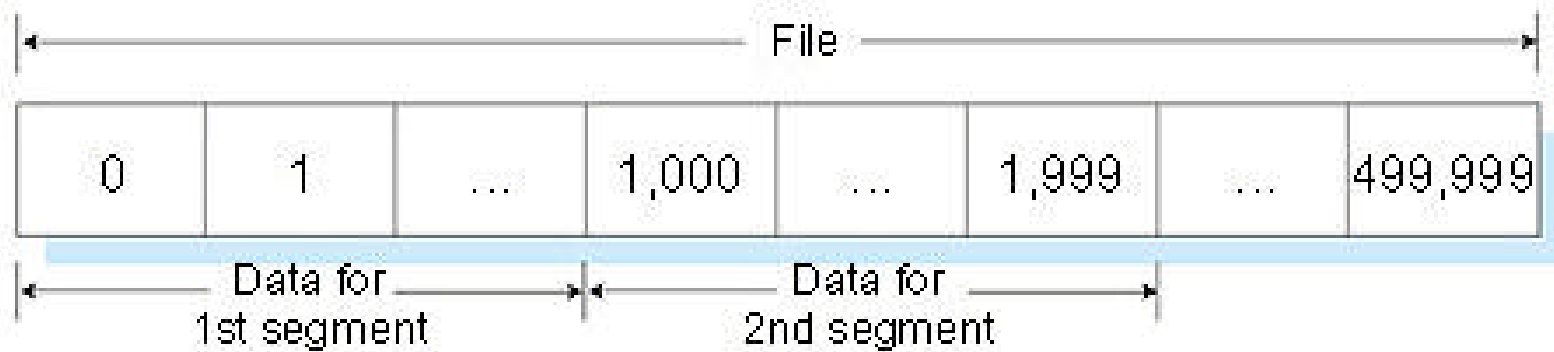
- ❑ **point-to-point:**
 - one sender, one receiver
 - ❑ **reliable, in-order *byte stream*:**
 - no “message boundaries”
 - ❑ **pipelined:**
 - TCP congestion and flow control set window size
 - ❑ ***send & receive buffers***
- ❑ **full duplex data:**
 - bi-directional data flow in same connection
 - MSS: maximum segment size
 - ❑ **connection-oriented:**
 - handshaking (exchange of control msgs) init's sender, receiver state before data exchange
 - ❑ **flow controlled:**
 - sender will not overwhelm receiver



TCP segment structure



Dividing file data into TCP segments



TCP seq. #'s and ACKs

Seq. #'s:

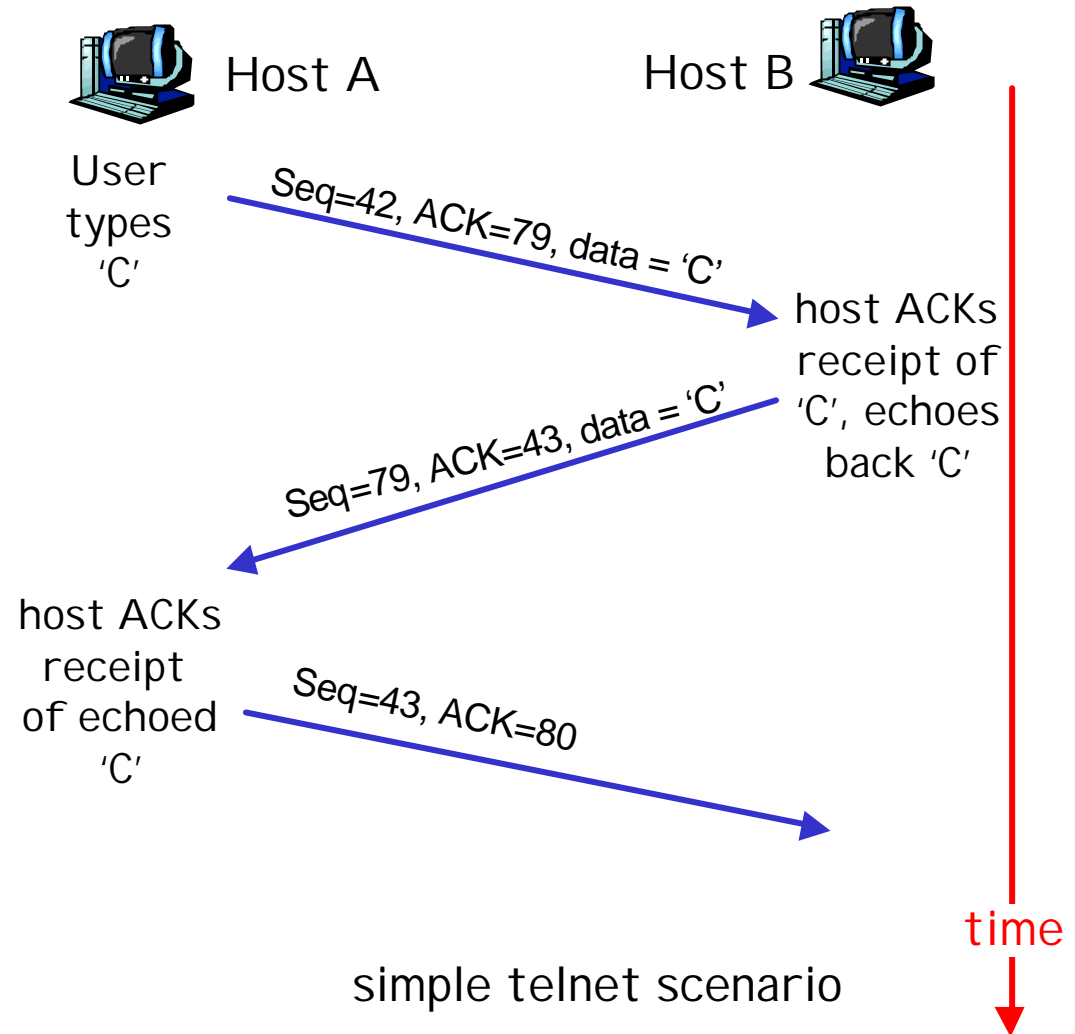
- byte stream
"number" of first byte in segment's data

ACKs:

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

Q: how receiver handles out-of-order segments

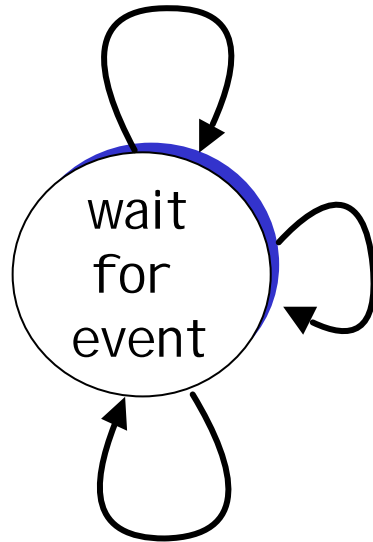
- A: TCP spec doesn't say, - up to implementor



TCP: reliable data transfer

event: data received
from application above

create, send segment



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

TCP: reliable data transfer

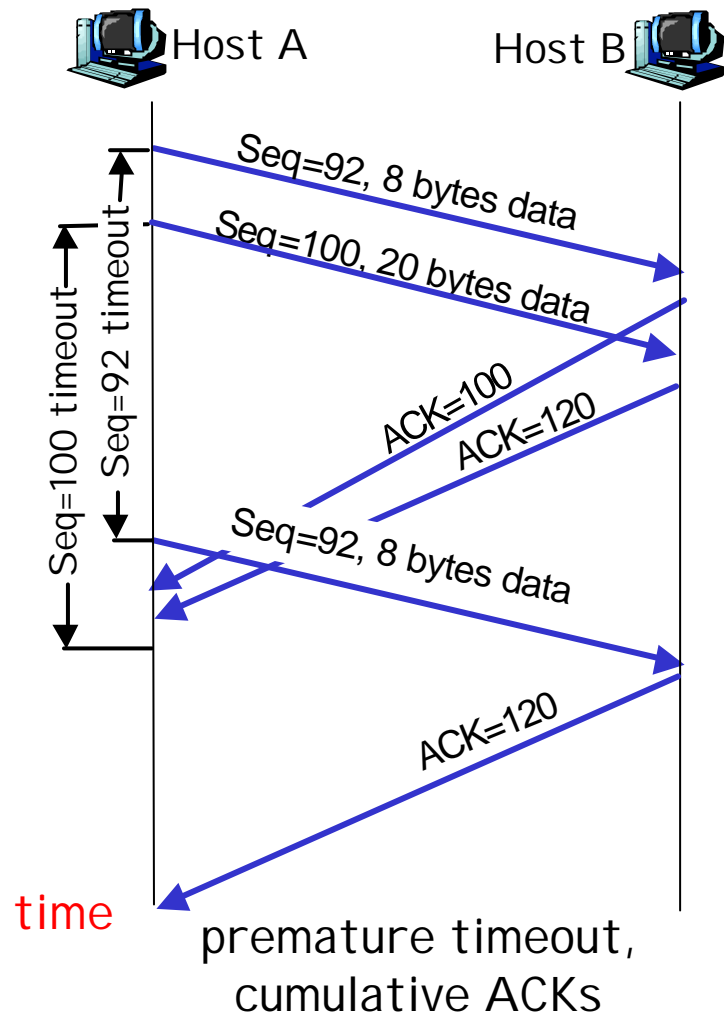
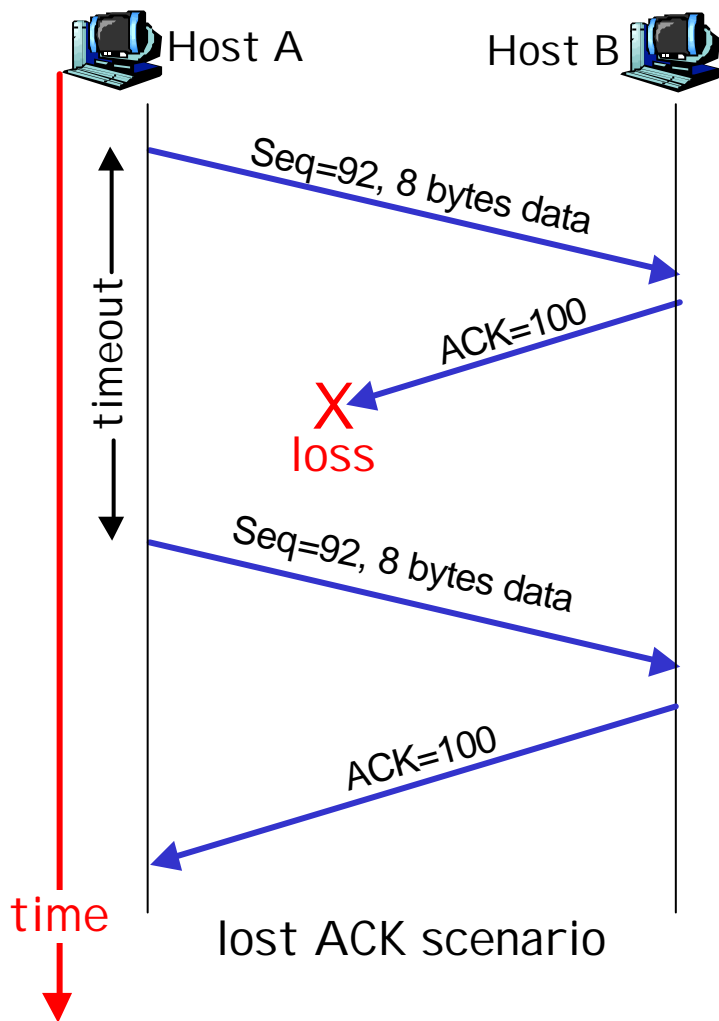
Simplified
TCP
sender

```
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             compute 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 */
```

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

TCP: retransmission scenarios



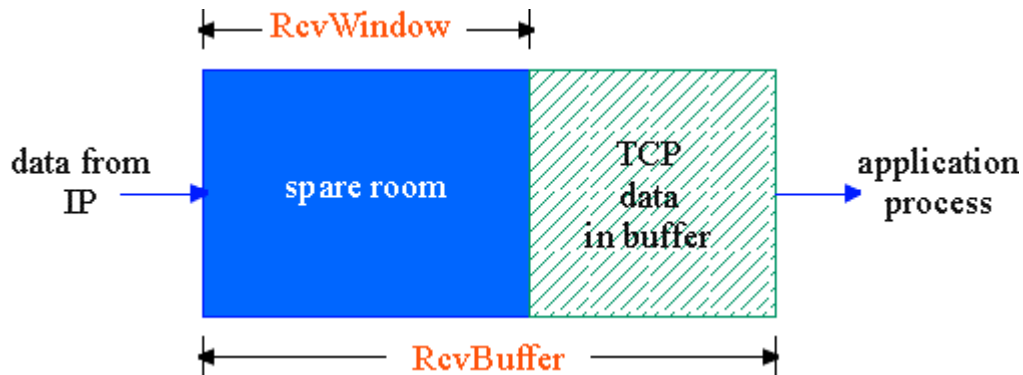
TCP Flow Control

flow control

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

RcvBuffer = size of TCP Receive Buffer

RcvWindow = amount of spare room in Buffer



receiver buffering

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

- **RcvWindow field** in TCP segment

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

TCP Round Trip Time and Timeout

Q: how to set TCP timeout value?

- ❑ longer than RTT
 - note: RTT will vary
- ❑ 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, cumulatively ACKed segments
- ❑ **SampleRTT** will vary, want estimated RTT "smoother"
 - use several recent measurements, not just current **SampleRTT**

TCP Round Trip Time and Timeout

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

- ❑ Exponential weighted moving average
- ❑ influence of given sample decreases exponentially fast
- ❑ typical value of x: 0.1

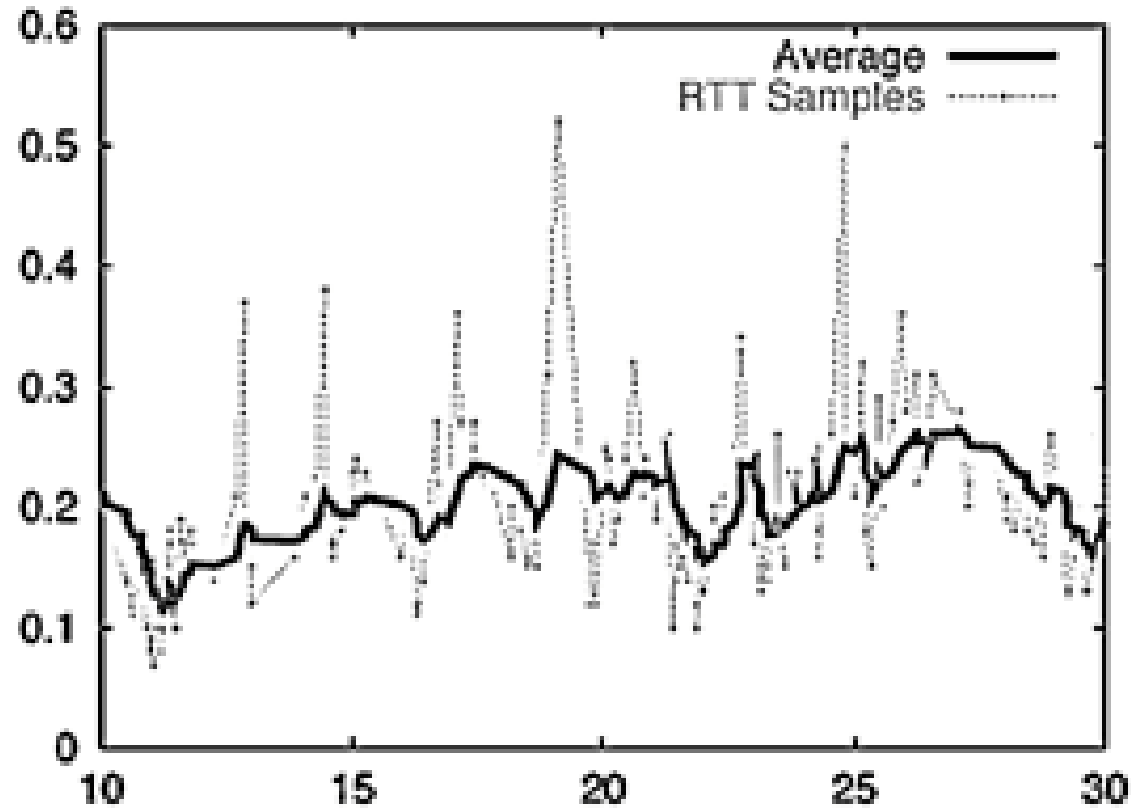
Setting the timeout

- ❑ **EstimtedRTT** plus "safety margin"
- ❑ large variation in **EstimatedRTT** -> larger safety margin

$$\text{Timeout} = \text{EstimatedRTT} + 4 * \text{Deviation}$$

$$\text{Deviation} = (1-x) * \text{Deviation} + x * |\text{SampleRTT} - \text{EstimatedRTT}|$$

RTT samples/estimate



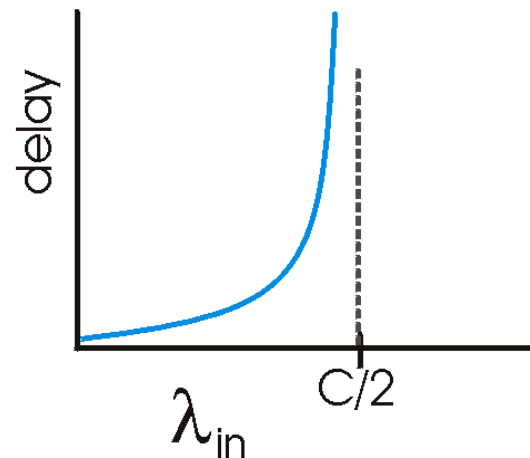
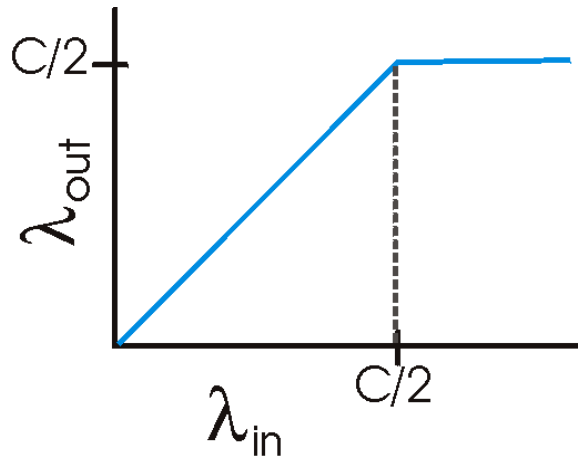
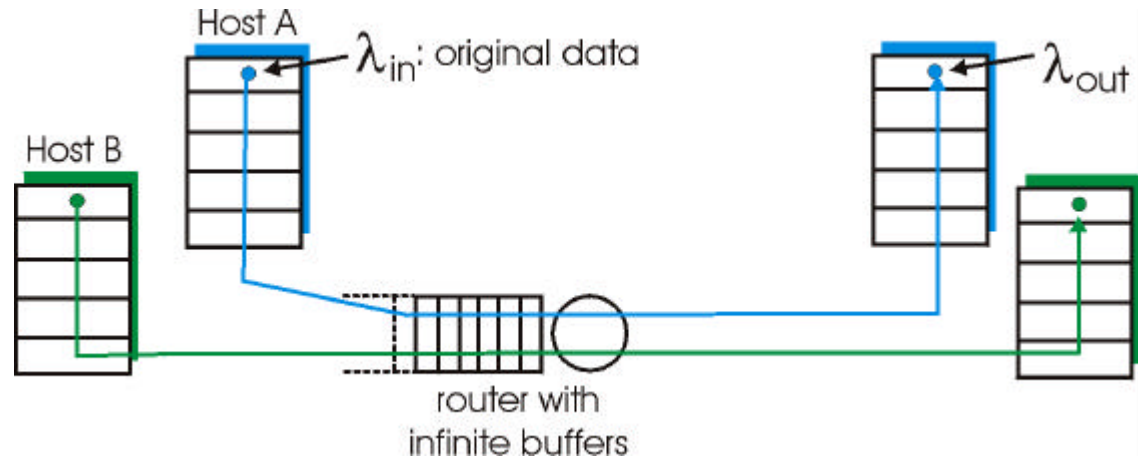
Principles of Congestion Control

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

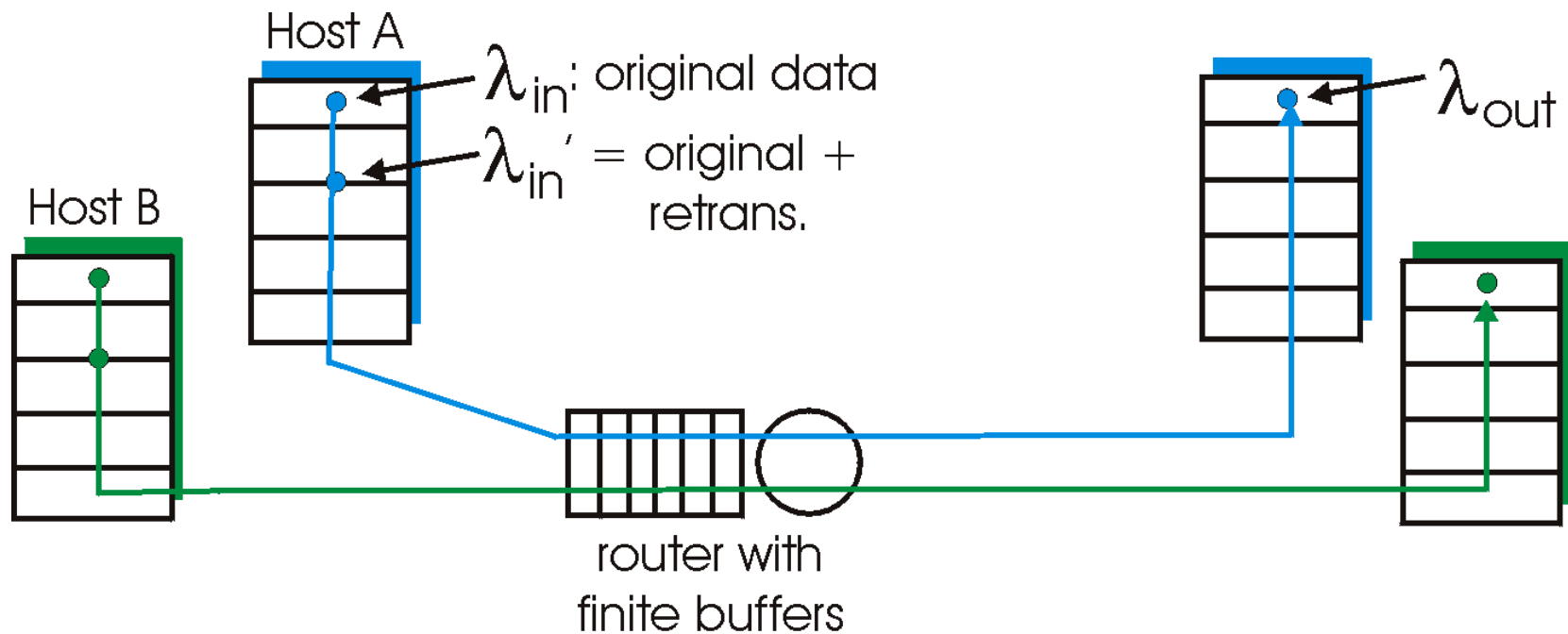
- two senders, two receivers
- one router, infinite buffers
- no retransmission



- large delays when congested
- maximum achievable throughput

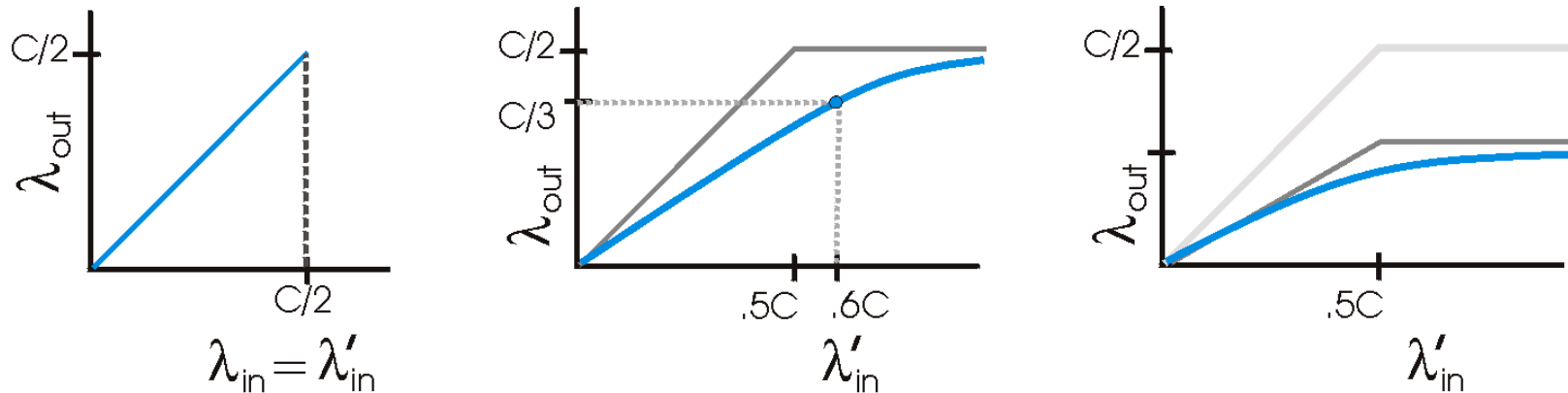
Causes/costs of congestion: scenario 2

- ❑ one router, *finite* buffers
- ❑ sender retransmission of lost packet



Causes/costs of congestion: scenario 2

- always: $\lambda_{in} = \lambda_{out}$ (goodput)
- “perfect” retransmission only when loss: $\lambda'_{in} > \lambda_{out}$
- retransmission of delayed (not lost) packet makes λ'_{in} larger (than perfect case) for same λ_{out}



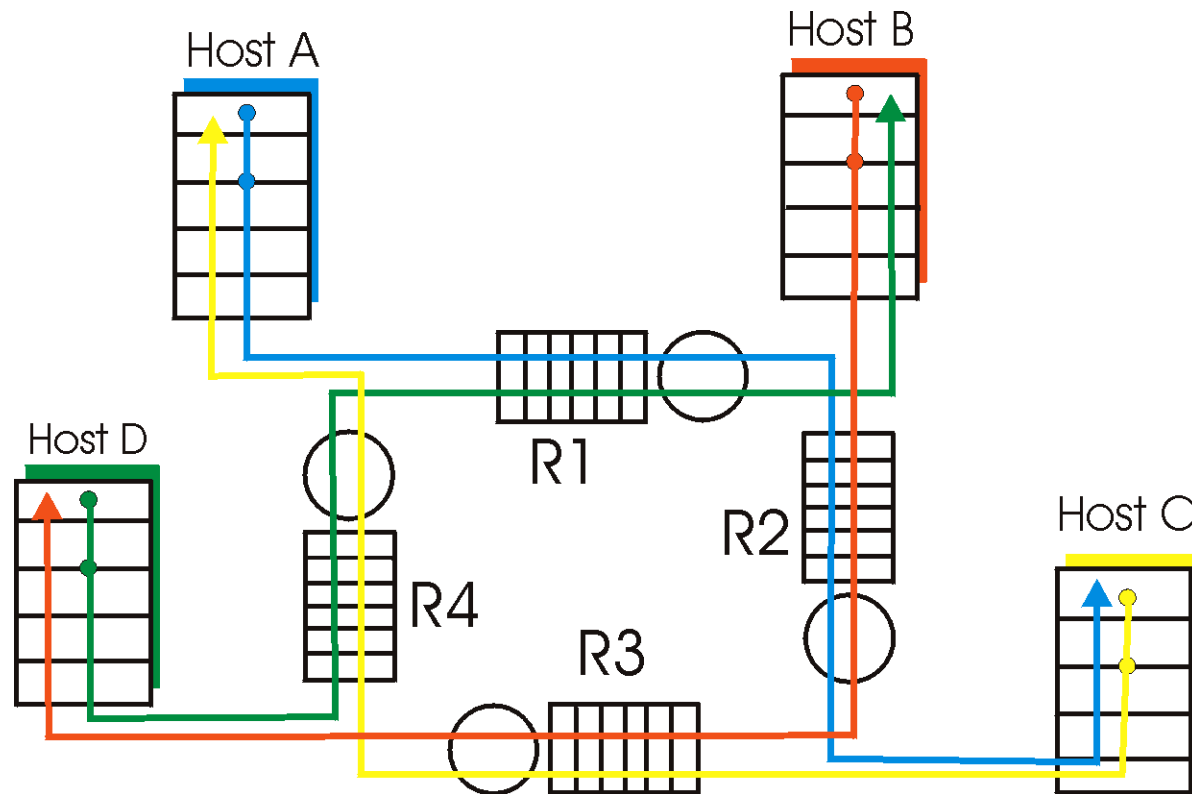
“costs” of congestion:

- more work (retrans) for given “goodput”
- unneeded retransmissions: link carries multiple copies of pkt

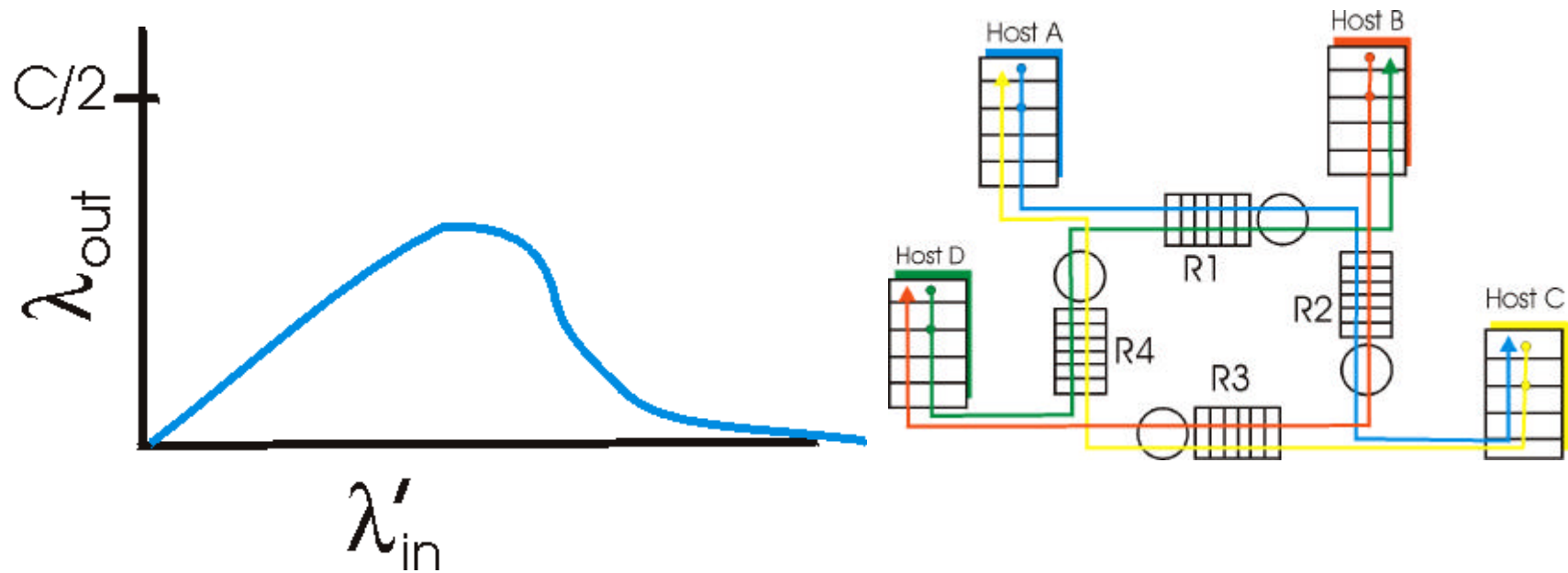
Causes/costs of congestion: scenario 3

- ❑ four senders
- ❑ multihop paths
- ❑ timeout/retransmit

Q: what happens as λ_{in} and λ'_{in} increase ?



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:

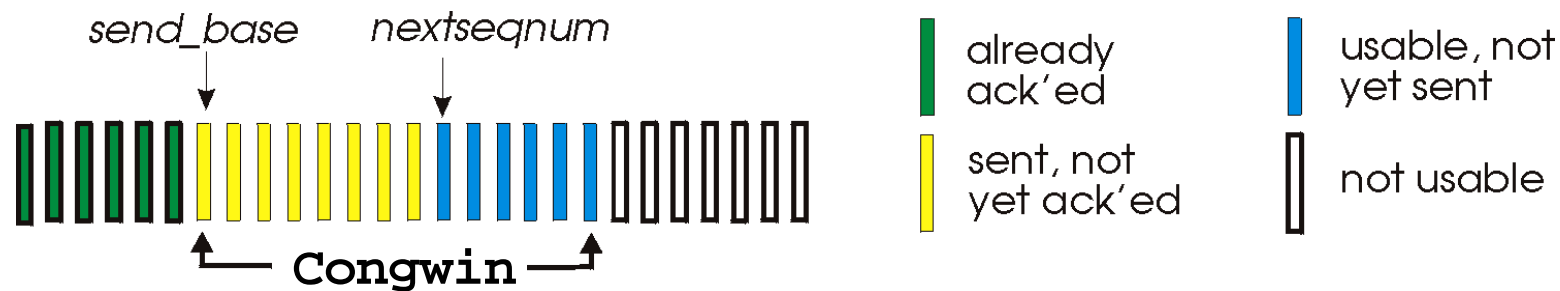
- ❑ no explicit feedback from network
- ❑ congestion inferred from end-system observed loss, delay
- ❑ approach taken by TCP

Network-assisted congestion control:

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

TCP Congestion Control

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



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

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

TCP congestion control:

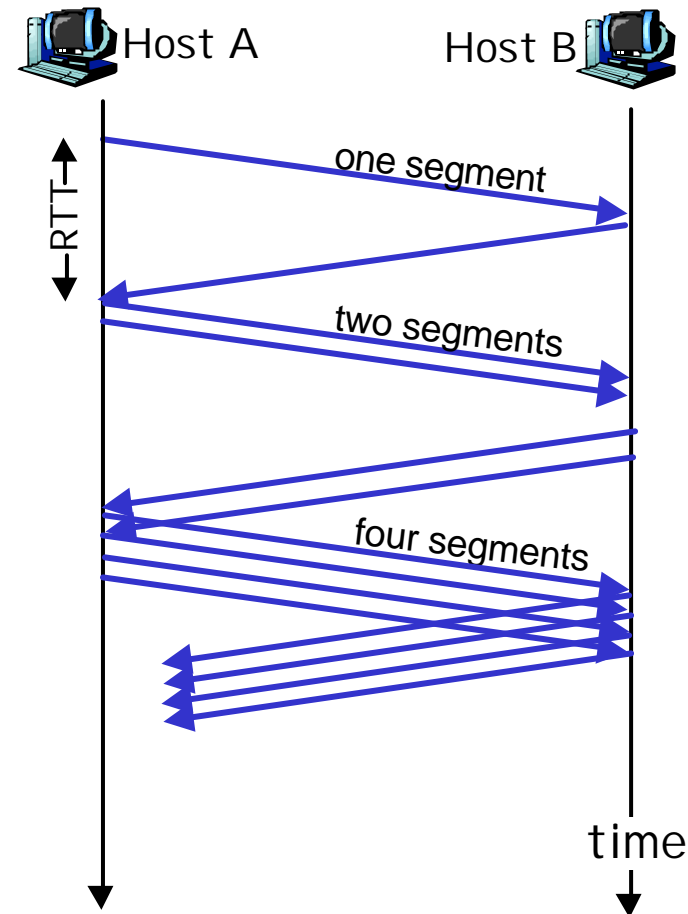
- “**probing**” for usable bandwidth:
 - **ideally**: transmit as fast as possible (**Congwin** as large as possible) without loss
 - *increase* **Congwin** until loss (congestion)
 - *loss*: *decrease* **Congwin**, then begin probing (increasing) again
- two “phases”
 - **slow start**
 - **congestion avoidance**
- important variables:
 - **Congwin**
 - **threshold**: defines threshold between two slow start phase, congestion control phase

TCP Slowstart

Slowstart algorithm

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

- exponential increase (per RTT) in window size (not so slow!)
- loss event: timeout (Tahoe TCP) and/or or three duplicate ACKs (Reno TCP)

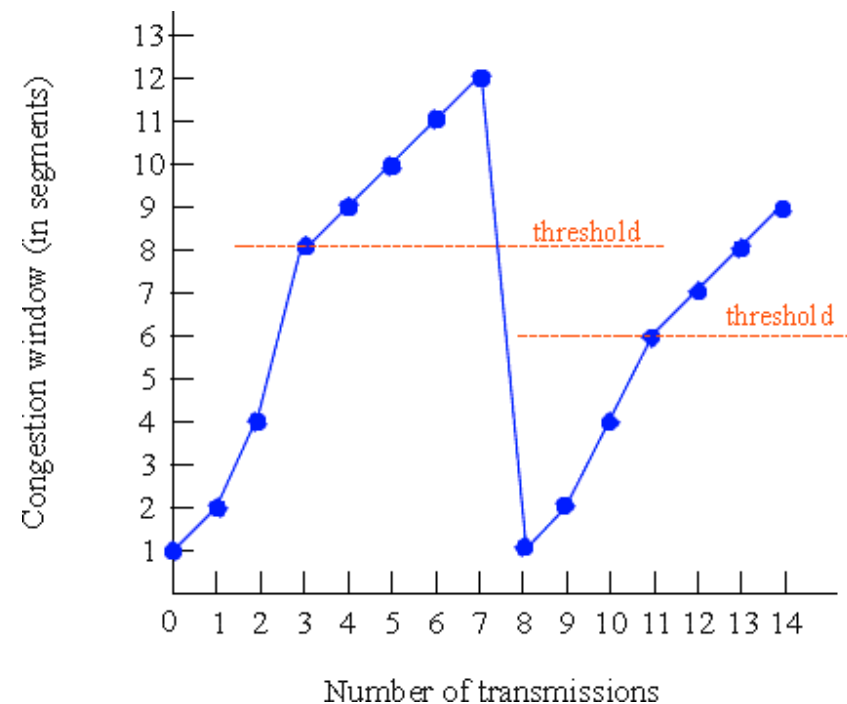


TCP Congestion Avoidance

Congestion avoidance

```
/* slowstart is over */
/* Congwin > threshold */
Until (loss event) {
    every w segments ACKed:
        Congwin++
}
threshold = Congwin/2
Congwin = 1
perform slowstart1
```

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



AIMD

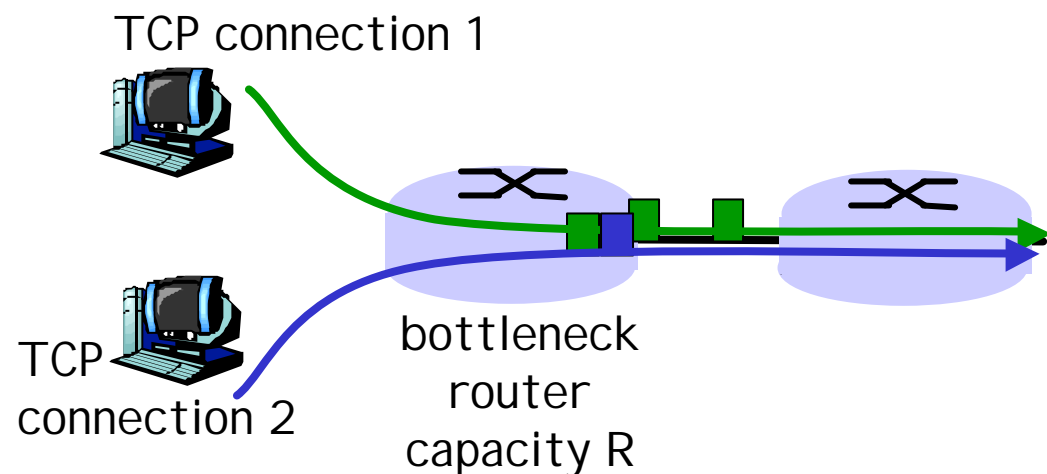
TCP congestion avoidance:

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

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

TCP Fairness

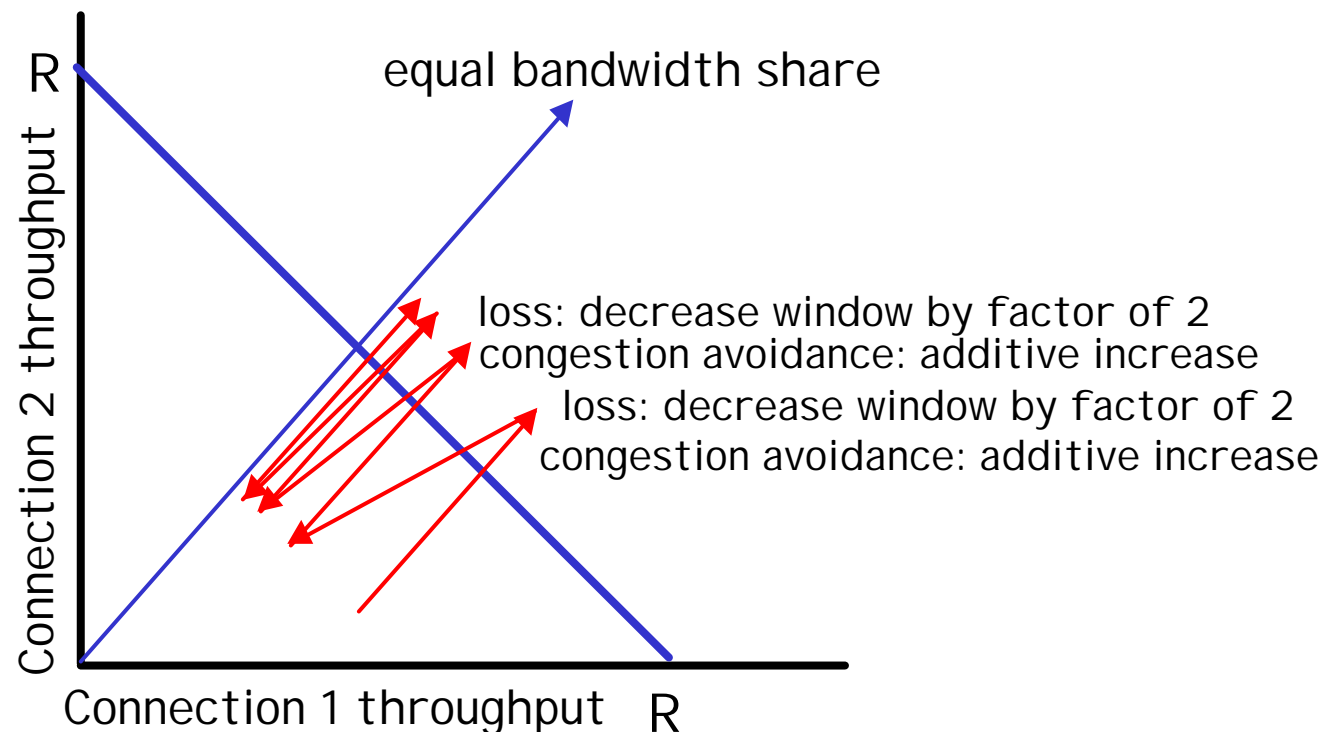
Fairness goal: if N TCP sessions share same bottleneck link, each should get $1/N$ of link capacity



Why is TCP fair?

Two competing sessions:

- Additive increase gives slope of 1, as throughput increases
- multiplicative decrease decreases throughput proportionally



TCP latency modeling

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

- ❑ TCP connection establishment
- ❑ data transfer delay

Notation, assumptions:

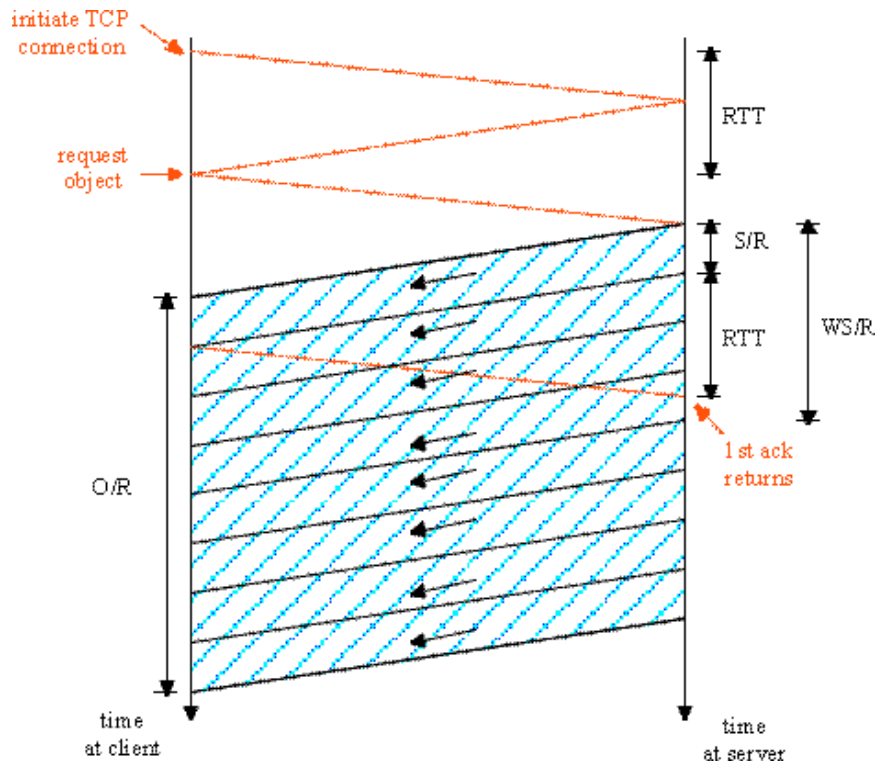
- ❑ Assume one link between client and server of rate R
- ❑ Assume: fixed congestion window, W segments
- ❑ S : MSS (bits)
- ❑ O : object size (bits)
- ❑ no retransmissions (no loss, no corruption)

Two cases to consider:

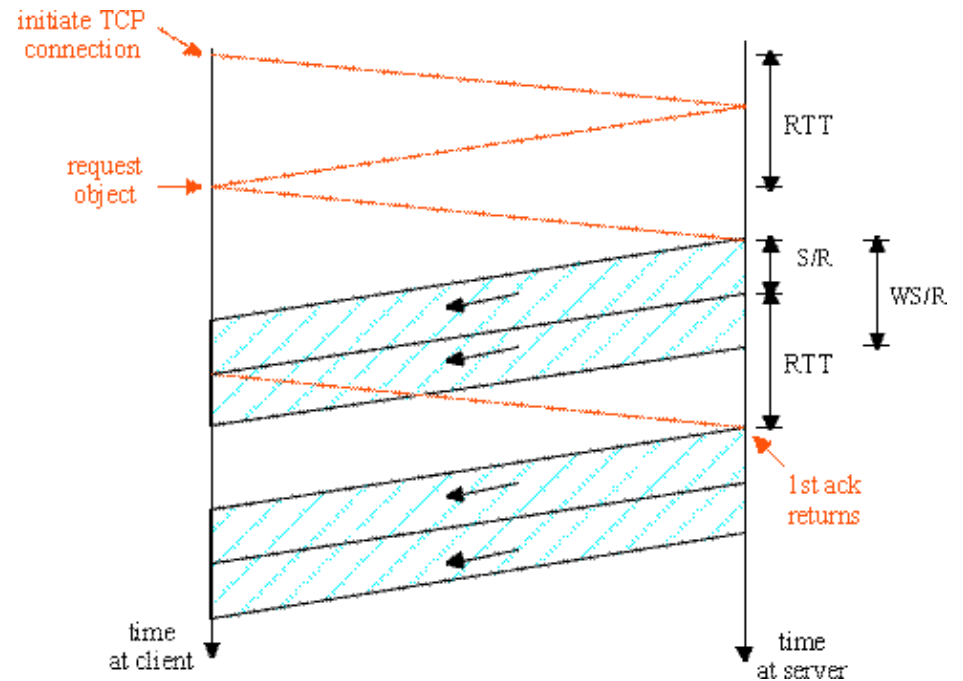
- ❑ $WS/R > RTT + S/R$: ACK for first segment in window returns before window's worth of data sent
- ❑ $WS/R < RTT + S/R$: wait for ACK after sending window's worth of data sent

TCP latency Modeling

$$K := O/WS$$



Case 1: latency = $2RTT + O/R$



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

TCP Latency Modeling: Slow Start

- ❑ Now suppose window grows according to slow start.
- ❑ 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.

TCP Latency Modeling: Slow Start (cont.)

Example:

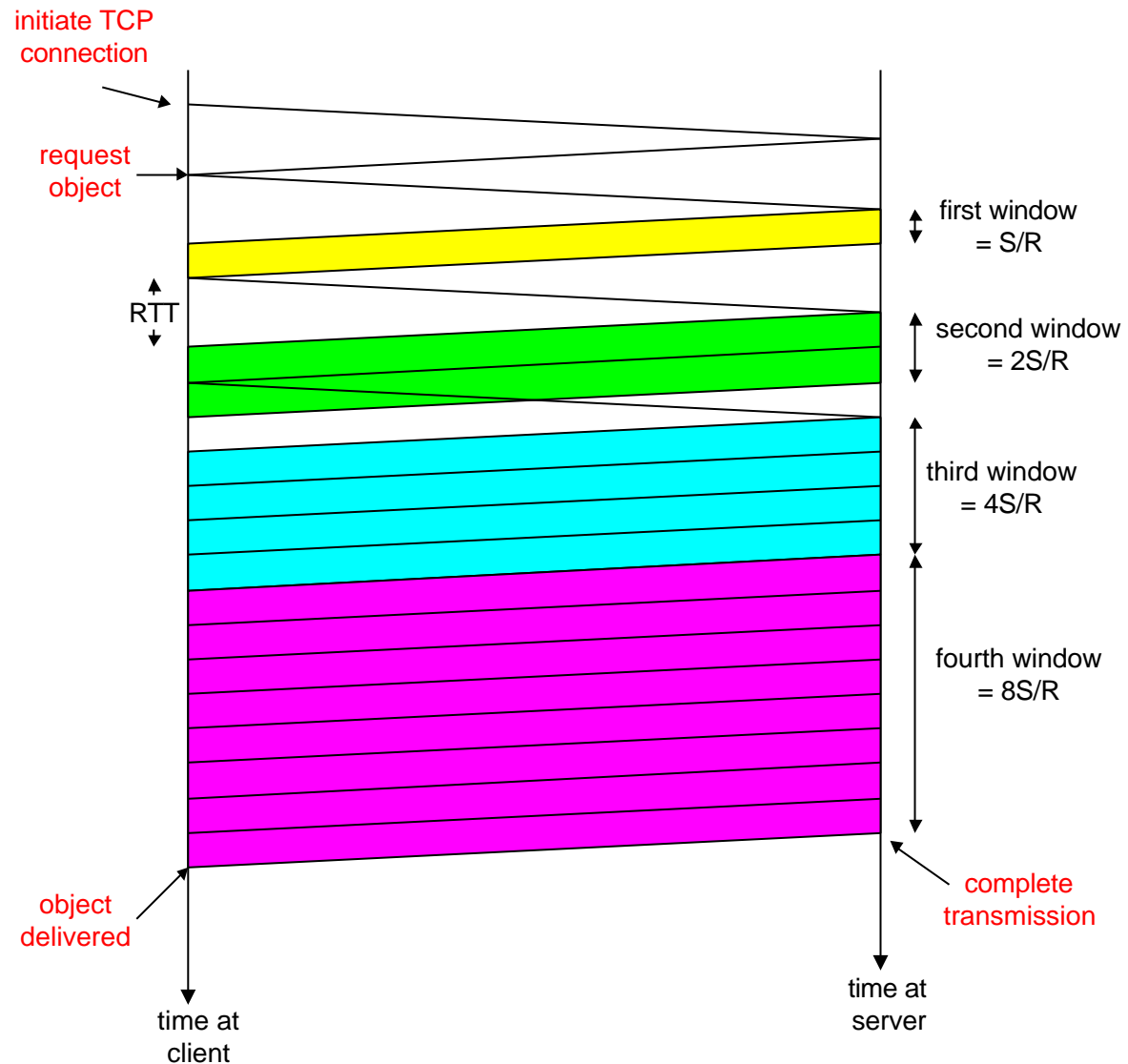
$O/S = 15$ segments

$K = 4$ windows

$Q = 2$

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

Server stalls $P=2$ times.



TCP Latency Modeling: Slow Start (cont.)

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

until server receives acknowledgement

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

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

$$\begin{aligned} \text{latency} &= \frac{O}{R} + 2RTT + \sum_{p=1}^P \text{stallTime}_p \\ &= \frac{O}{R} + 2RTT + \sum_{k=1}^P \left[\frac{S}{R} + RTT - 2^{k-1} \frac{S}{R} \right] \\ &= \frac{O}{R} + 2RTT + P \left[RTT + \frac{S}{R} \right] - (2^P - 1) \frac{S}{R} \end{aligned}$$

