

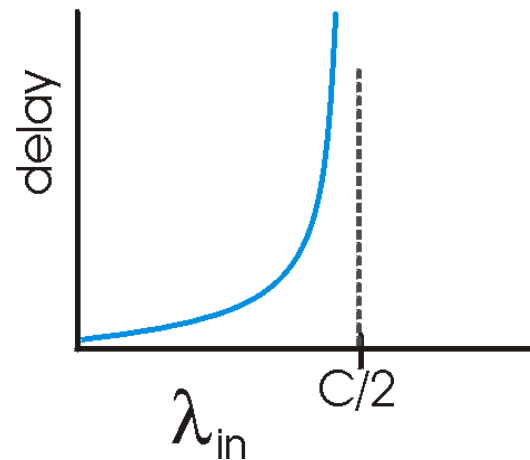
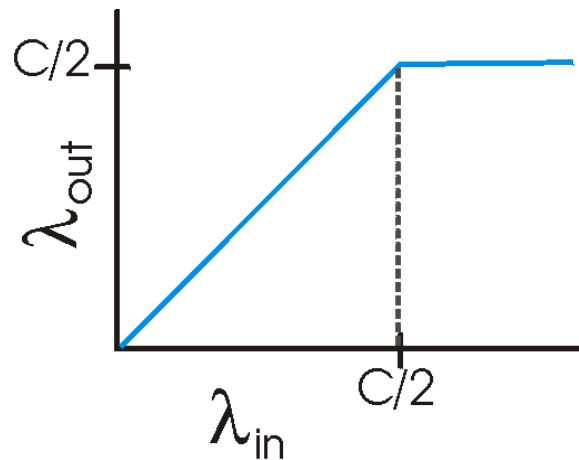
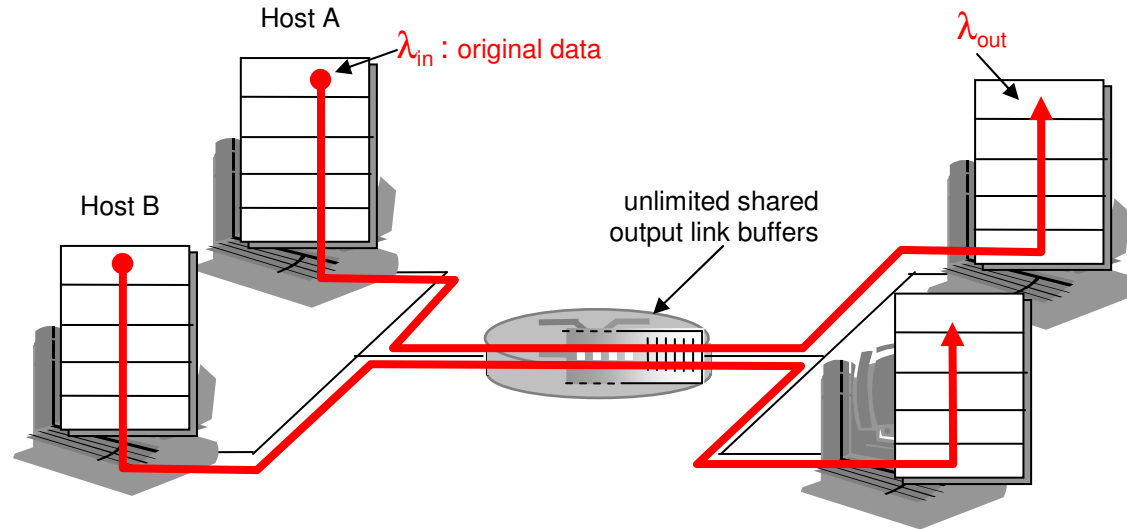
Principles of Congestion Control

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

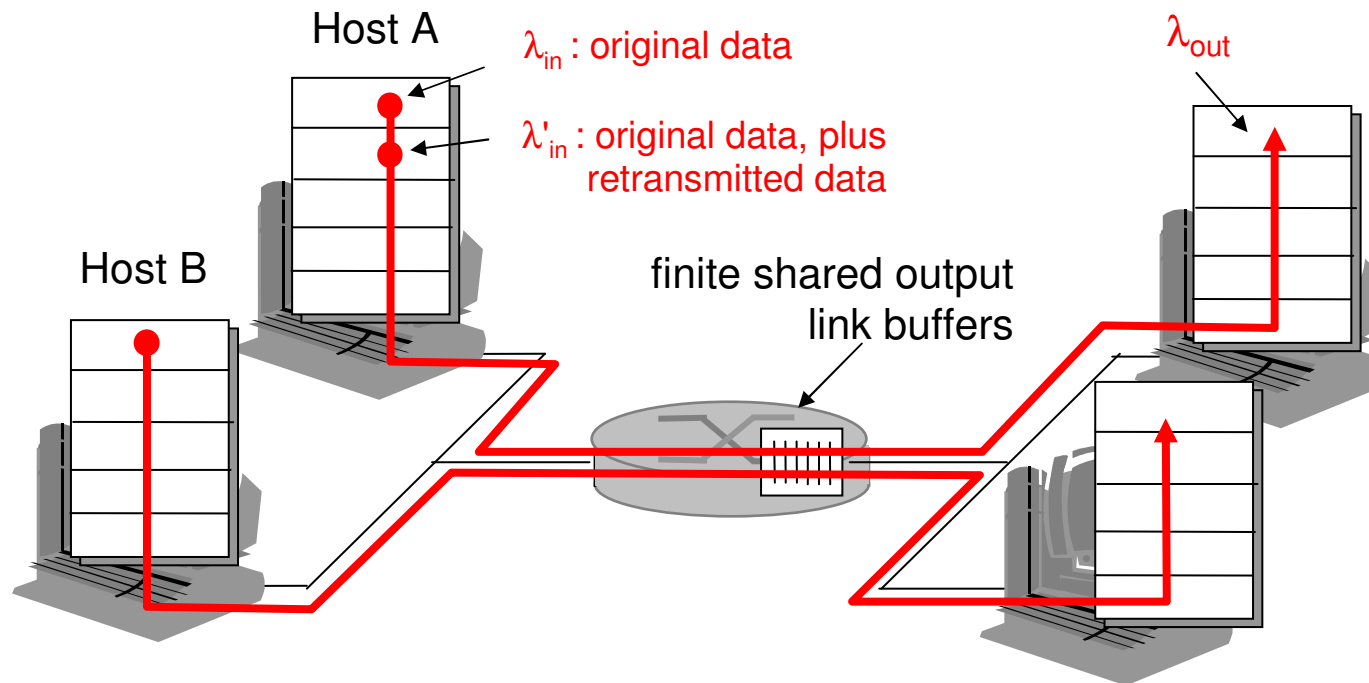
- ❑ two senders, two receivers
- ❑ one router, infinite buffers
- ❑ no retransmission



- ❑ large delays when congested
- ❑ maximum achievable throughput

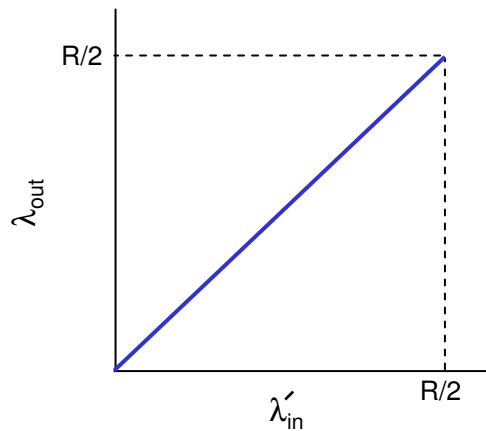
Causes/costs of congestion: scenario 2

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

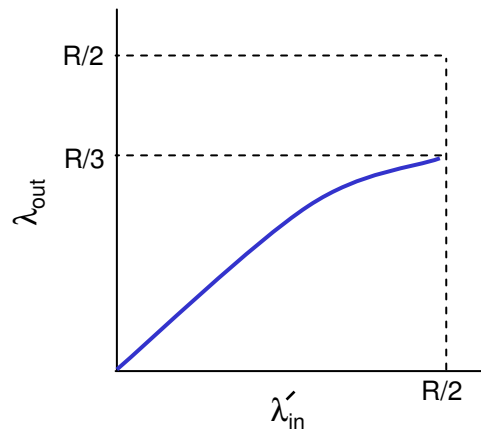


Causes/costs of congestion: scenario 2

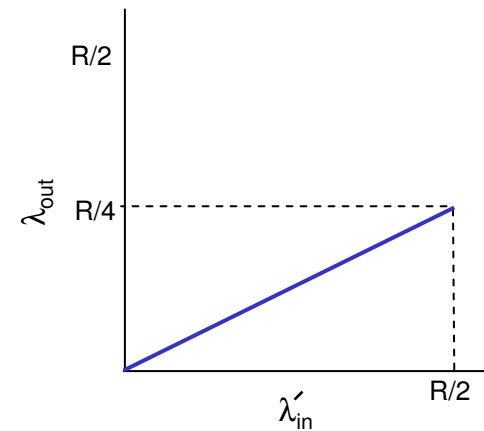
- always: $\lambda_{in} = \lambda_{out}$ (goodput)
- “perfect” retransmission only when loss: $\lambda'_{in} < \lambda_{out}$
- retransmission of delayed (not lost) packet makes λ'_{in} larger (than perfect case) for same λ_{out}



a.



b.



c.

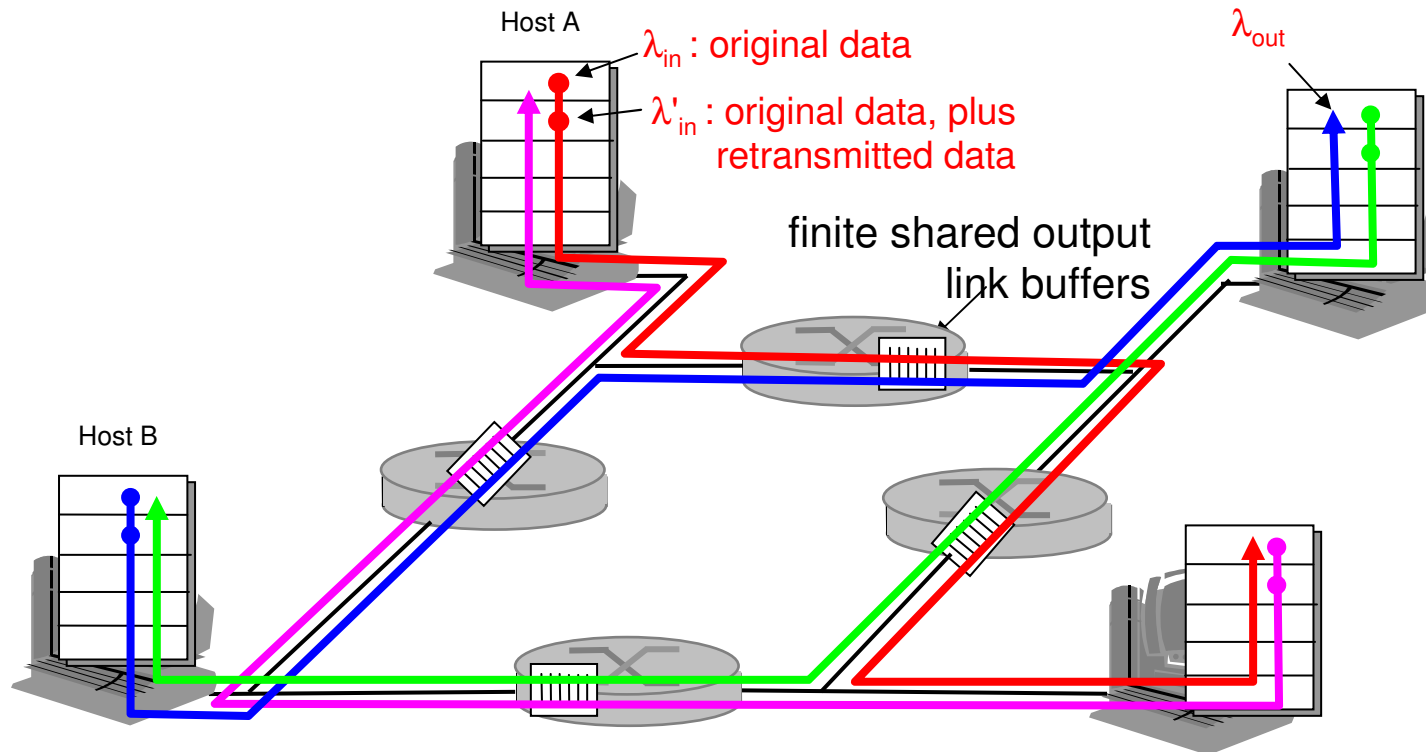
“costs” of congestion:

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

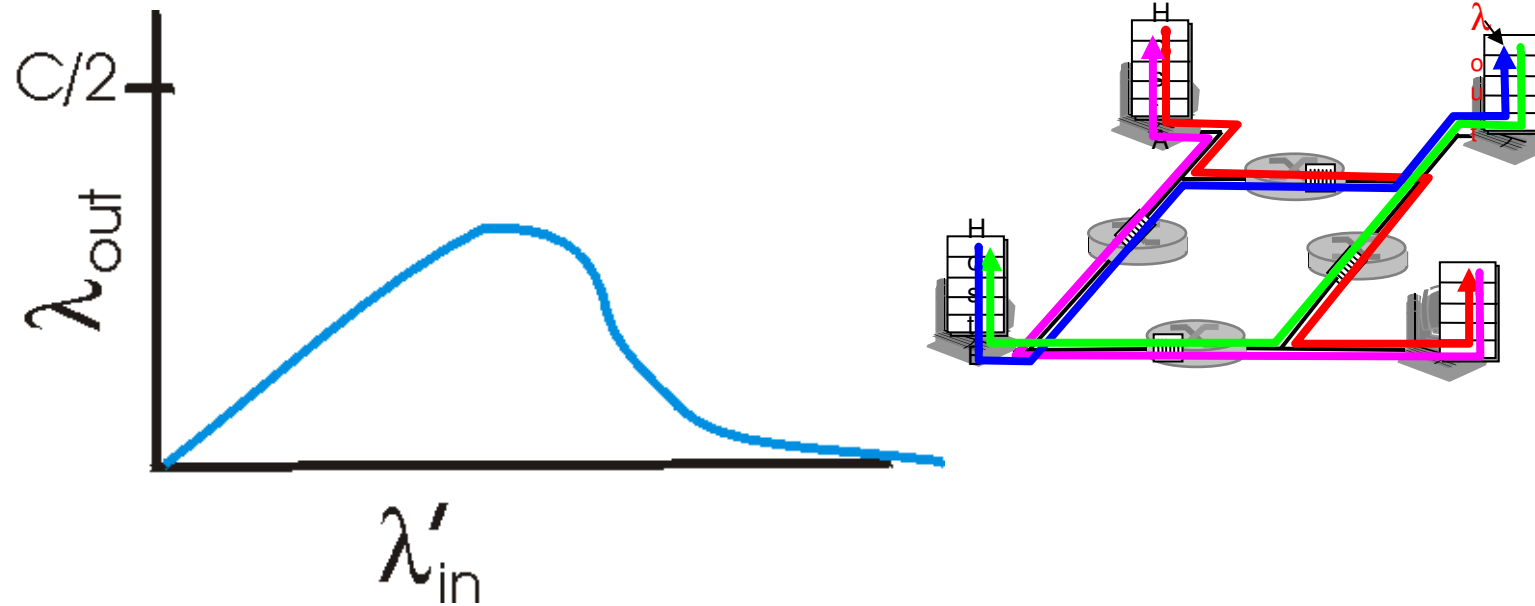
Causes/costs of congestion: scenario 3

- ❑ 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:

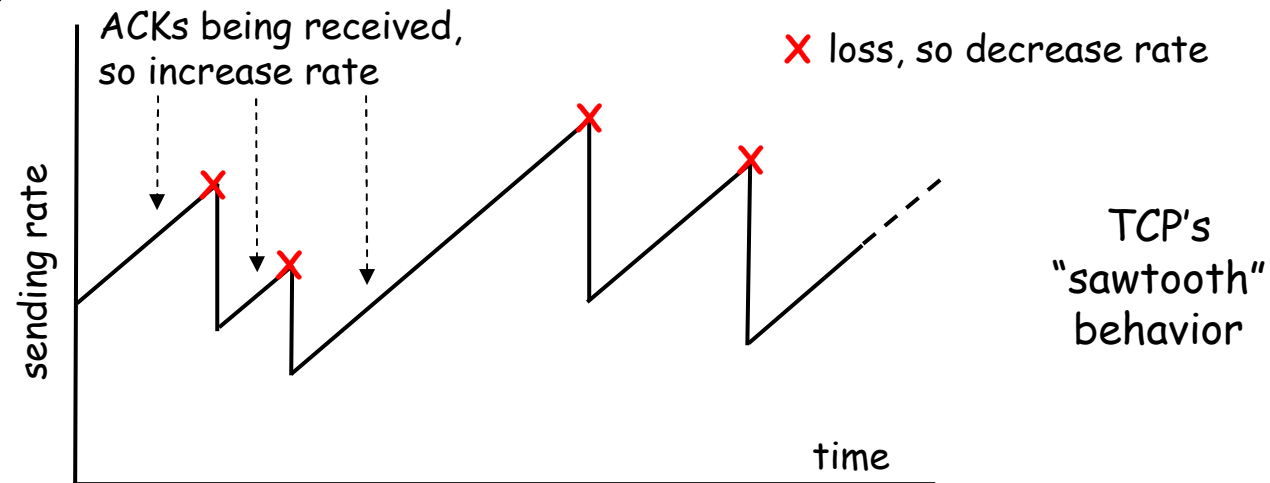
- ❑ routers provide feedback to end systems
 - single bit indicating congestion (SNA, DECbit, TCP/IP ECN, ATM)
 - explicit rate sender should send at

TCP congestion control:

- *goal:* TCP sender should transmit as fast as possible, but without congesting network
 - Q: how to find rate *just* below congestion level
- decentralized: each TCP sender sets its own rate, based on *implicit* feedback:
 - *ACK*: segment received (a good thing!), network not congested, so increase sending rate
 - *lost segment*: assume loss due to congested network, so decrease sending rate

TCP congestion control: bandwidth probing

- “probing for bandwidth”: increase transmission rate on receipt of ACK, until eventually loss occurs, then decrease transmission rate
 - continue to increase on ACK, decrease on loss (since available bandwidth is changing, depending on other connections in network)



- Q: how fast to increase/decrease?
 - details to follow

TCP Congestion Control: details

- sender limits rate by limiting number of unACKed bytes "in pipeline":

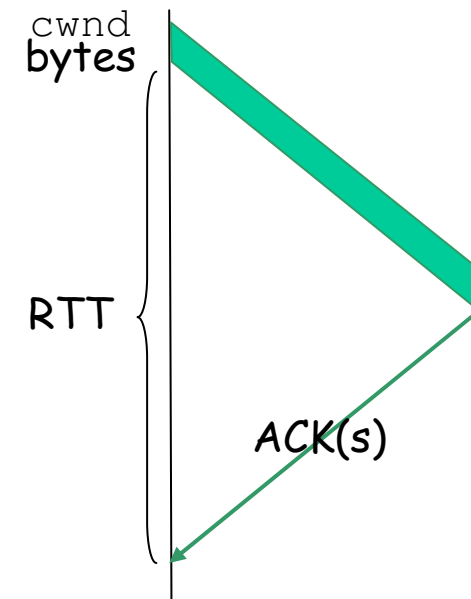
$\text{LastByteSent} - \text{LastByteAcked} \leq \text{cwnd}$

- cwnd: differs from rwnd (how, why?)
- sender limited by $\min(\text{cwnd}, \text{rwnd})$

- roughly,

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

- cwnd is dynamic, function of perceived network congestion



TCP Congestion Control: more details

segment loss event: reducing cwnd

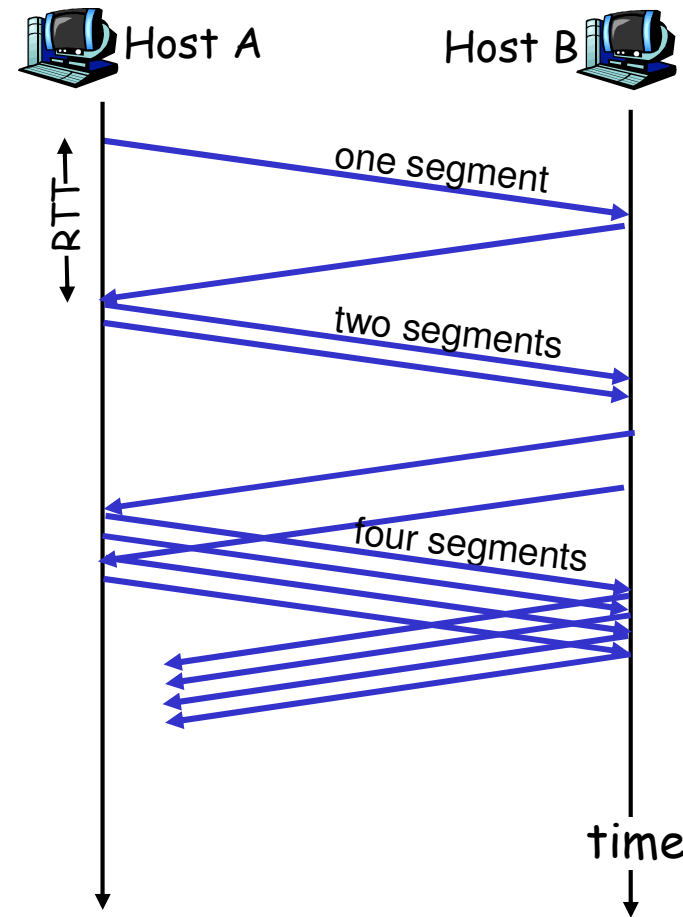
- ❑ timeout: no response from receiver
 - cut cwnd to 1
- ❑ 3 duplicate ACKs: at least some segments getting through (recall fast retransmit)
 - cut cwnd in half, less aggressively than on timeout

ACK received: increase cwnd

- ❑ slowstart phase:
 - increase exponentially fast (despite name) at connection start, or following timeout
- ❑ congestion avoidance:
 - increase linearly

TCP Slow Start

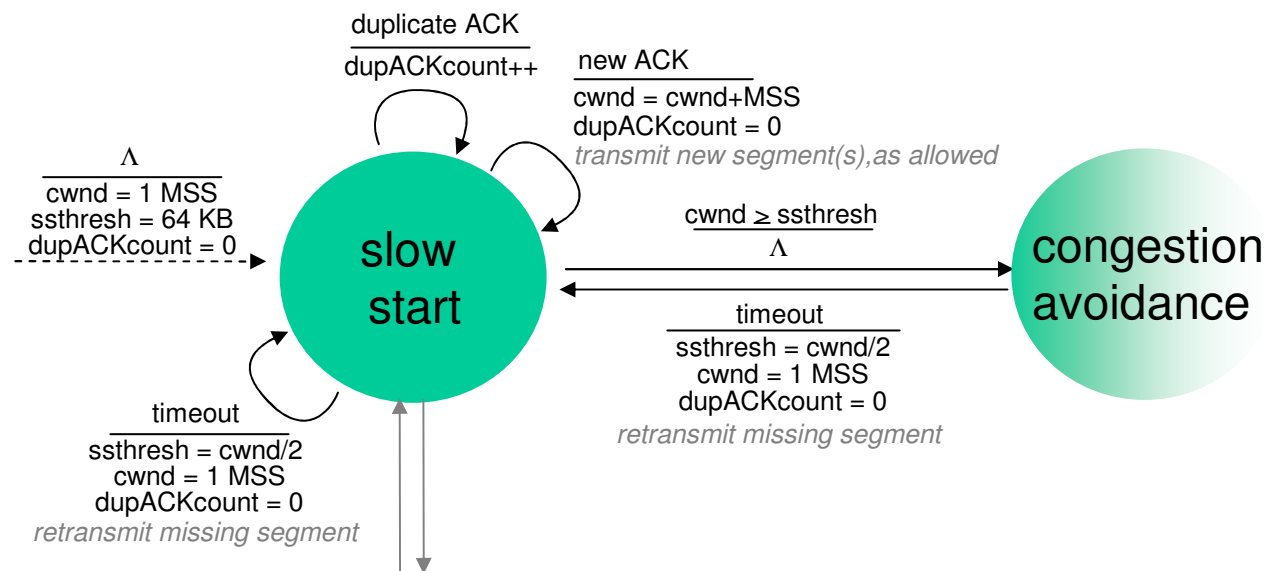
- when connection begins, `cwnd` = 1 MSS
 - example: `MSS` = 500 bytes & `RTT` = 200 msec
 - initial rate = 20 kbps
- available bandwidth may be \gg `MSS/RTT`
 - desirable to quickly ramp up to respectable rate
- increase rate exponentially until first loss event or when threshold reached
 - double `cwnd` every `RTT`
 - done by incrementing `cwnd` by 1 for every `ACK` received



Transitioning into/out of slowstart

ssthresh: cwnd threshold maintained by TCP

- ❑ on loss event: set ssthresh to cwnd/2
 - remember (half of) TCP rate when congestion last occurred
- ❑ when $\text{cwnd} \geq \text{ssthresh}$: transition from slowstart to congestion avoidance phase



TCP: congestion avoidance

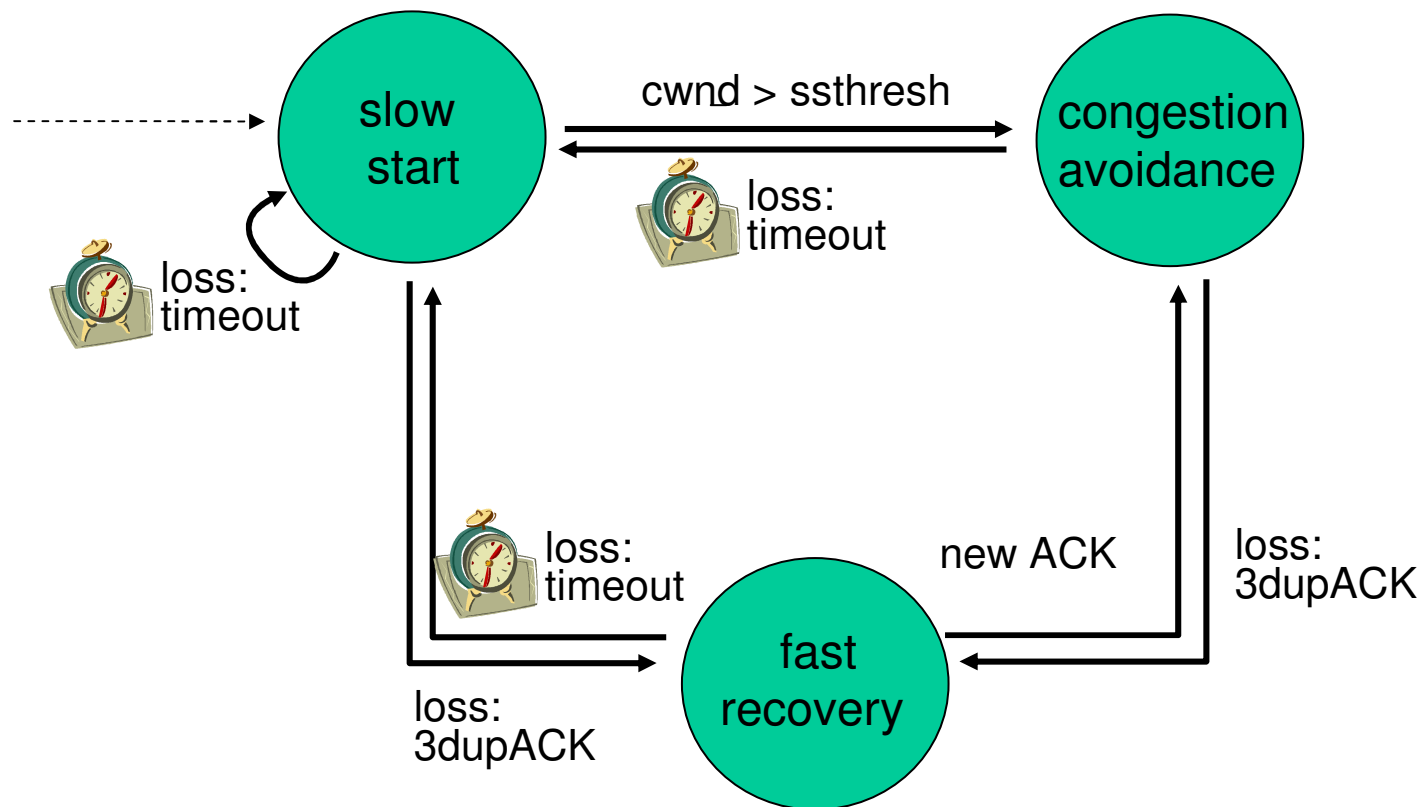
- when $cwnd > ssthresh$
grow $cwnd$ linearly
 - increase $cwnd$ by 1 MSS per RTT
 - approach possible congestion slower than in slowstart
 - implementation: $cwnd = cwnd + MSS/cwnd$ for each ACK received

AIMD

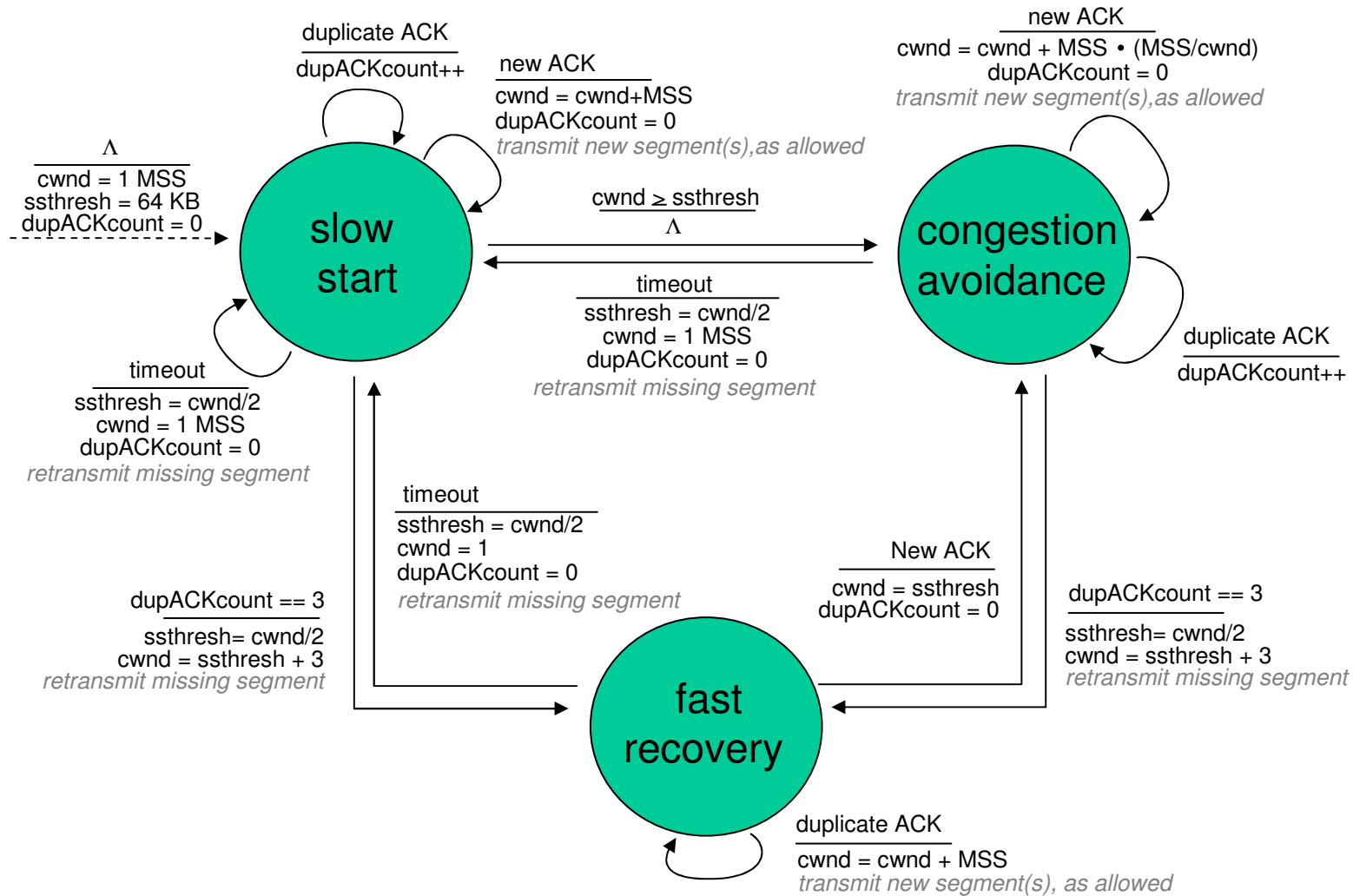
- **ACKs**: increase $cwnd$ by 1 MSS per RTT: additive increase
- **loss**: cut $cwnd$ in half (non-timeout-detected loss): multiplicative decrease

AIMD: Additive Increase
Multiplicative Decrease

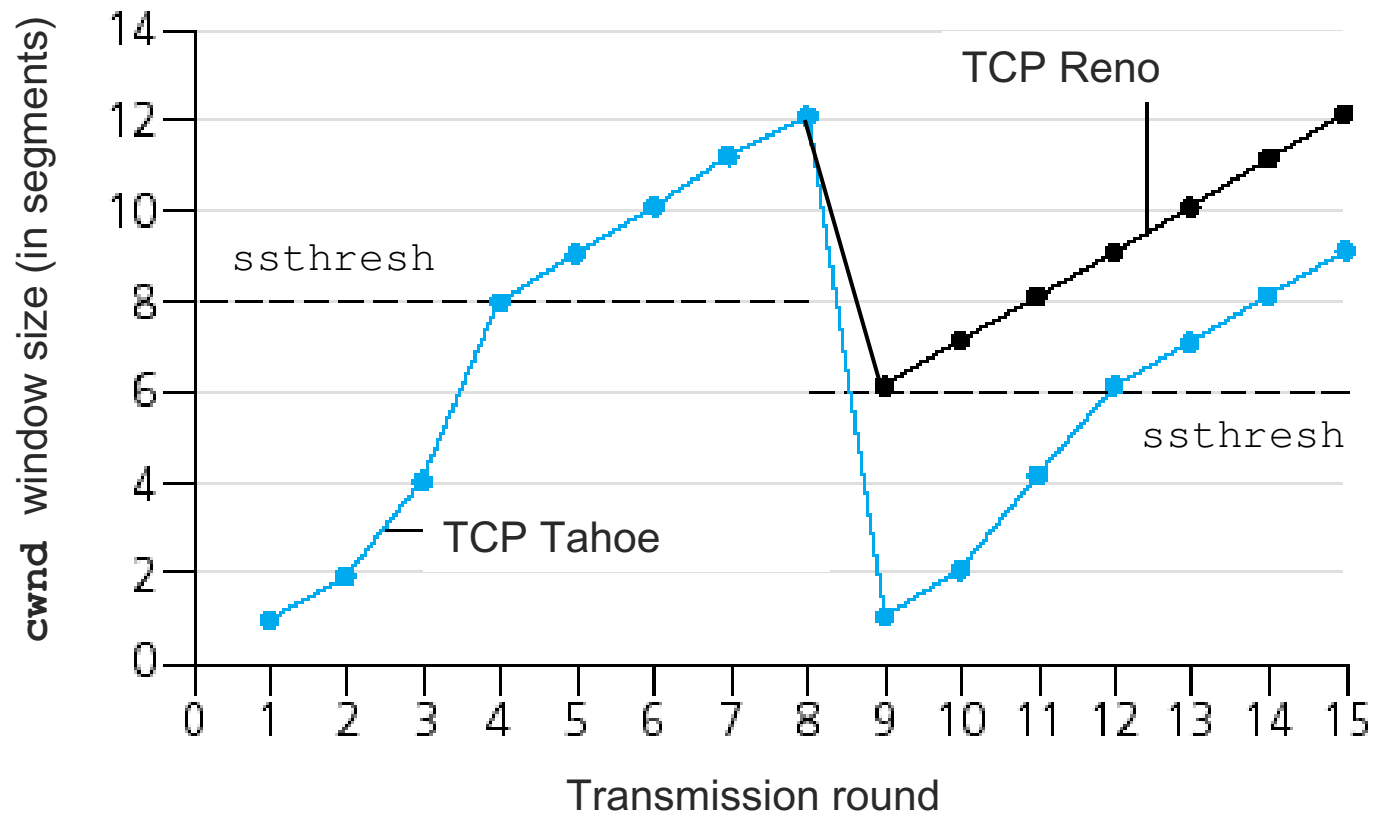
TCP congestion control FSM: overview



TCP congestion control FSM: details



Popular "flavors" of TCP



Summary: TCP Congestion Control

- when $cwnd < ssthresh$, sender in **slow-start** phase, window grows exponentially.
- when $cwnd \geq ssthresh$, sender is in **congestion-avoidance** phase, window grows linearly.
- when **triple duplicate ACK** occurs, $ssthresh$ set to $cwnd/2$, $cwnd$ set to $\sim ssthresh$
- when **timeout** occurs, $ssthresh$ set to $cwnd/2$, $cwnd$ set to 1 MSS.

TCP throughput

- Q: what's average throughput of TCP as function of window size, RTT?
 - ignoring slow start
- let W be window size when loss occurs.
 - when window is W , throughput is W/RTT
 - just after loss, window drops to $W/2$, throughput to $W/2RTT$.
 - average throughput: $.75 W/RTT$

TCP Futures: TCP over "long, fat pipes"

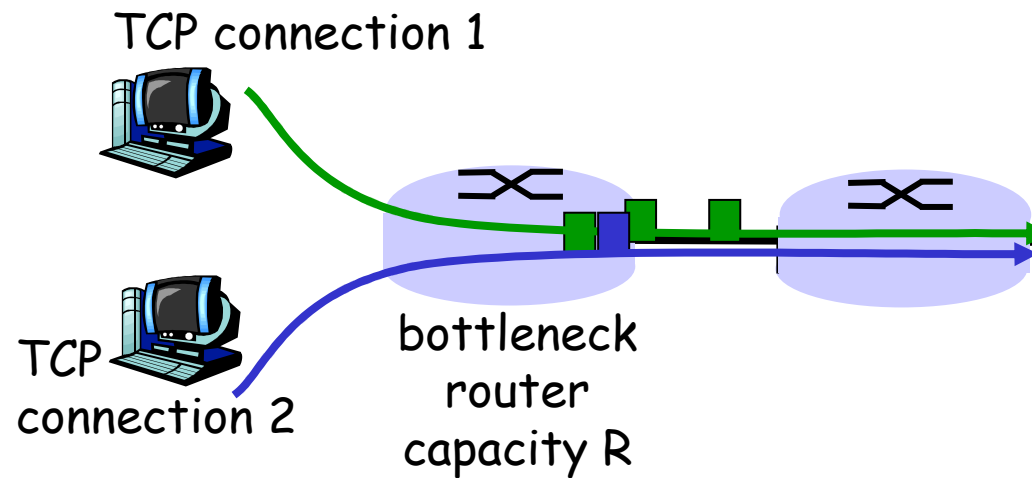
- ❑ example: 1500 byte segments, 100ms RTT, want 10 Gbps throughput
- ❑ requires window size $W = 83,333$ in-flight segments
- ❑ throughput in terms of loss rate:

$$\frac{1.22 \cdot MSS}{RTT \sqrt{L}}$$

- ❑ $\rightarrow L = 2 \cdot 10^{-10}$ *Wow*
- ❑ new versions of TCP for high-speed

TCP Fairness

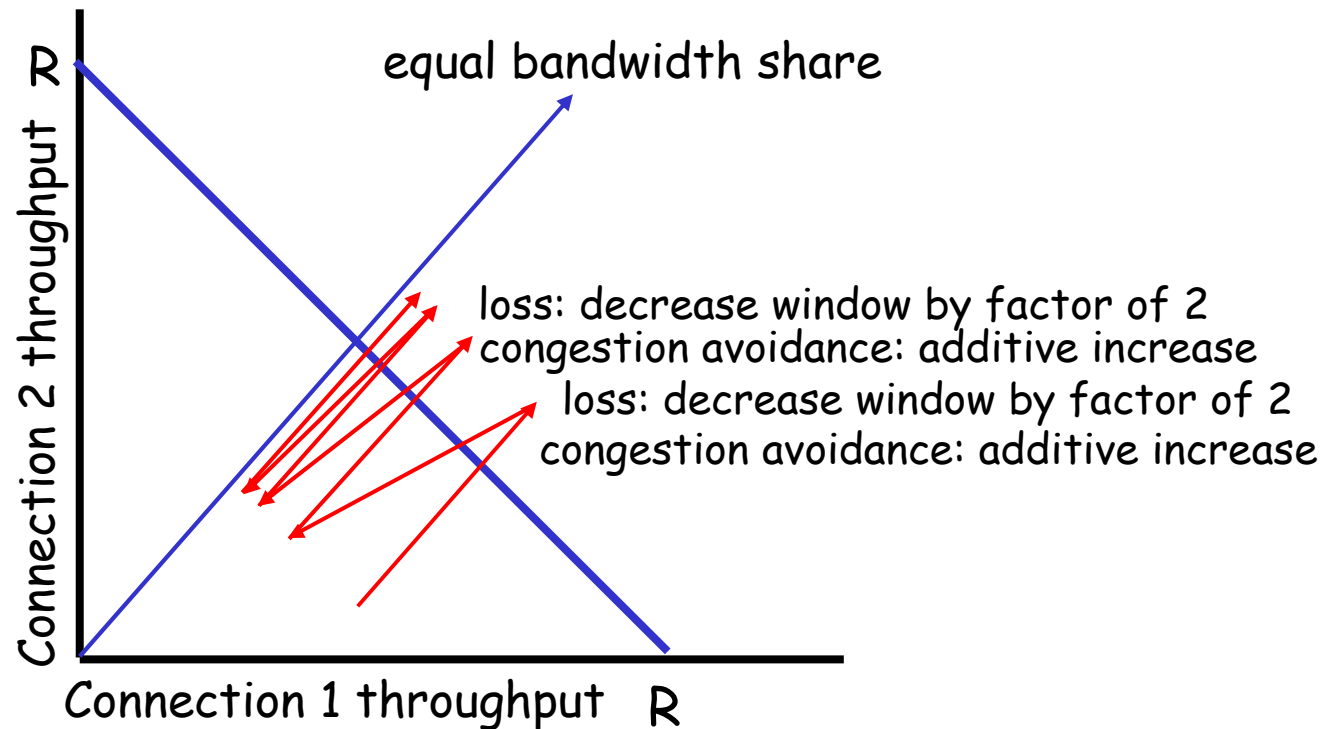
fairness goal: if K TCP sessions share same bottleneck link of bandwidth R , each should have average rate of R/K



Why is TCP fair?

Two competing sessions:

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



Fairness (more)

Fairness and UDP

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

Fairness and parallel TCP connections

- ❑ nothing prevents app from opening parallel connections between 2 hosts.
- ❑ web browsers do this
- ❑ example: link of rate R supporting 9 connections;
 - new app asks for 1 TCP, gets rate $R/10$
 - new app asks for 11 TCPs, gets $R/2$!

Chapter 3

Transport Layer

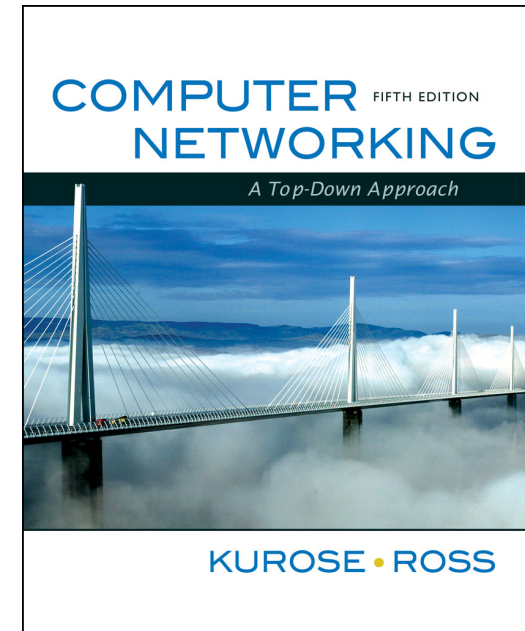
A note on the use of these ppt slides:

We're making these slides freely available to all (faculty, students, readers). They're in PowerPoint form so you can add, modify, and delete slides (including this one) and slide content to suit your needs. They obviously represent a *lot* of work on our part. In return for use, we only ask the following:

- If you use these slides (e.g., in a class) in substantially unaltered form, that you mention their source (after all, we'd like people to use our book!)
- If you post any slides in substantially unaltered form on a www site, that you note that they are adapted from (or perhaps identical to) our slides, and note our copyright of this material.

Thanks and enjoy! JFK/KWR

All material copyright 1996-2009
J.F Kurose and K.W. Ross, All Rights Reserved



*Computer Networking:
A Top Down Approach
5th edition.*

*Jim Kurose, Keith Ross
Addison-Wesley, April
2009.*