

CSC358 Intro. to Computer Networks

Lecture 7: TCP, flow and congestion control

Amir H. Chinaei, Winter 2016

ahchinaei@cs.toronto.edu
http://www.cs.toronto.edu/~ahchinaei/

Many slides are (inspired/adapted) from the above source
© all material copyright; all rights reserved for the authors



Office Hours: T 17:00-18:00 R 9:00-10:00 BA4222

TA Office Hours: W 16:00-17:00 BA3201 R 10:00-11:00 BA7172

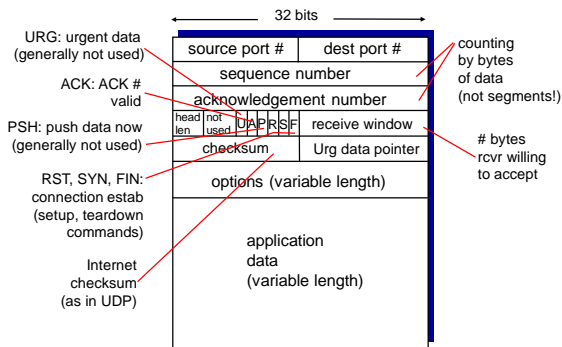
csc358ta@cdf.toronto.edu
http://www.cs.toronto.edu/~ahchinaei/teaching/2016jan/csc358/

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

Transport Layer 3-2

TCP segment structure



TCP seq. numbers, ACKs

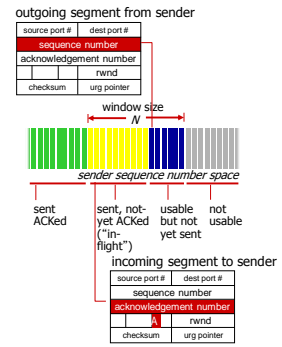
sequence numbers:

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

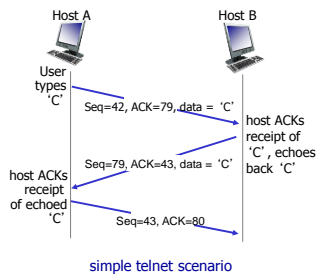
acknowledgements:

- 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 seq. numbers, ACKs



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

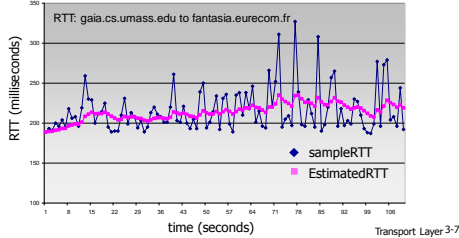
Q: how to estimate RTT?

- ❖ **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**
- Transport Layer 3-6

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: $\alpha = 0.125$



TCP round trip time, timeout

- ❖ **timeout interval:** EstimatedRTT plus “safety margin”
 - large variation in EstimatedRTT -> larger safety margin
- ❖ estimate SampleRTT deviation from EstimatedRTT:

$$\text{DevRTT} = (1 - \beta) * \text{DevRTT} + \beta * |\text{SampleRTT} - \text{EstimatedRTT}|$$

(typically, $\beta = 0.25$)

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



↑
estimated RTT

↑
“safety margin”

Transport Layer 3-8

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

Transport Layer 3-9

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:

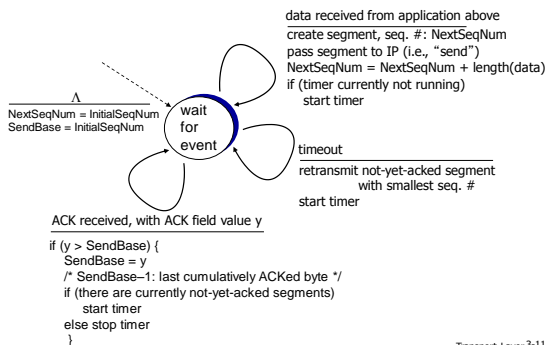
- ❖ retransmit segment that caused timeout
- ❖ restart timer

ack rcvd:

- ❖ if ack acknowledges previously unacked segments
 - update what is known to be ACKed
 - start timer if there are still unacked segments

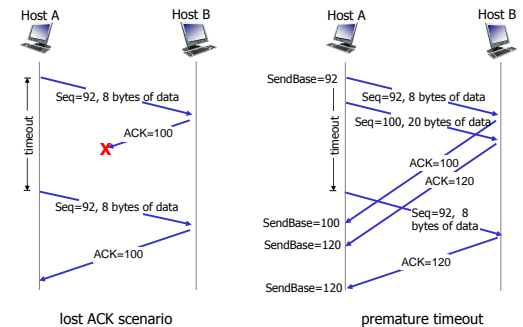
Transport Layer 3-10

TCP sender (simplified)



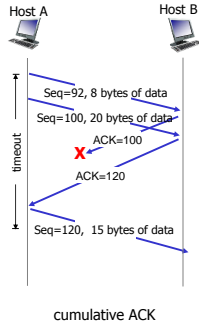
Transport Layer 3-11

TCP: retransmission scenarios



Transport Layer 3-12

TCP: retransmission scenarios



Transport Layer 3-13

TCP ACK generation [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-expected seq. #. Gap detected	immediately send <i>duplicate ACK</i> , 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

Transport Layer 3-14

TCP fast retransmit

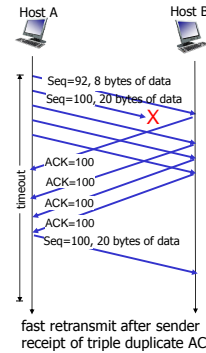
- time-out period often relatively long:
 - long delay before resending lost packet
- detect lost segments via duplicate ACKs.
 - sender often sends many segments back-to-back
 - if segment is lost, there will likely be many duplicate ACKs.

TCP fast retransmit
 if sender receives 3 ACKs for same data ("triple duplicate ACKs"), resend unacked segment with smallest seq #

- likely that unacked segment lost, so don't wait for timeout

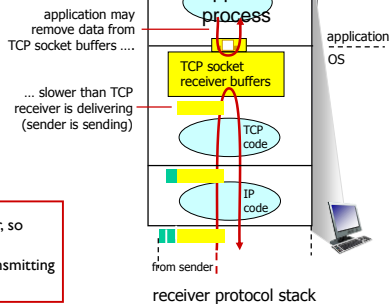
Transport Layer 3-15

TCP fast retransmit



Transport Layer 3-16

TCP flow control

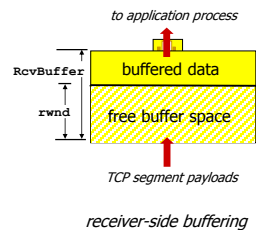


flow control
 receiver controls sender, so sender won't overflow receiver's buffer by transmitting too much, too fast

Transport Layer 3-17

TCP flow control

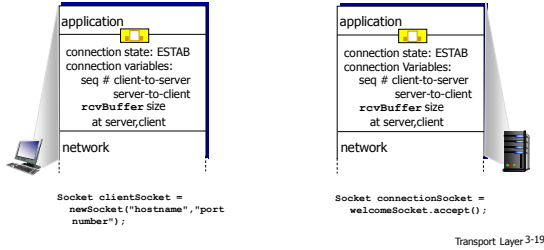
- receiver "advertises" free buffer space by including **rwnd** value in TCP header of receiver-to-sender segments
 - RcvBuffer** size set via socket options (typical default is 4096 bytes)
 - many operating systems autoadjust **RcvBuffer**
- sender limits amount of unacked ("in-flight") data to receiver's **rwnd** value
- guarantees receive buffer will not overflow



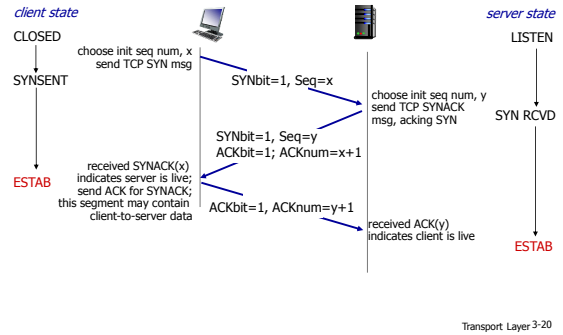
Transport Layer 3-18

Connection Management

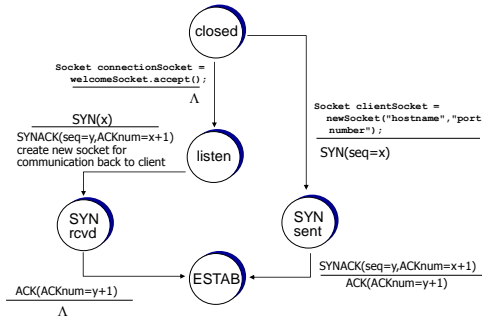
- before exchanging data, sender/receiver "handshake":
- ❖ agree to establish connection (each knowing the other willing to establish connection)
 - ❖ agree on connection parameters



TCP 3-way handshake



TCP 3-way handshake: FSM

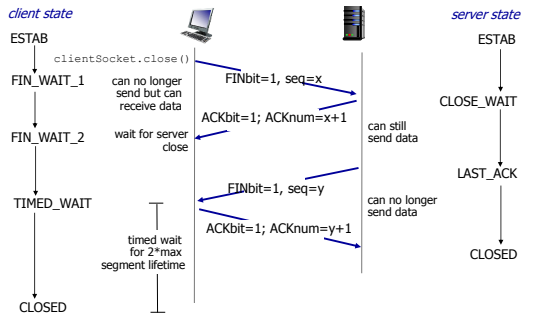


TCP: closing a connection

- ❖ 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

Transport Layer 3-22

TCP: closing a connection



Chapter 3 outline

- 3.1 transport-layer services
- 3.2 multiplexing and demultiplexing
- 3.3 connectionless transport: UDP
- 3.4 principles of reliable data transfer
- 3.5 connection-oriented transport: TCP
 - segment structure
 - reliable data transfer
 - flow control
 - connection management
- 3.6 principles of congestion control
- 3.7 TCP congestion control

Transport Layer 3-24

Principles of congestion control

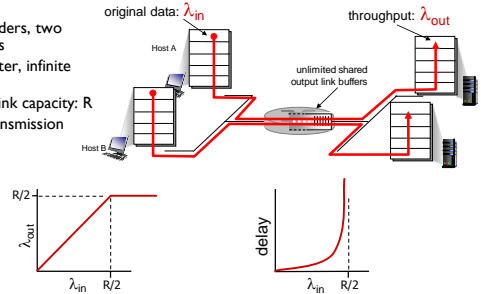
congestion:

- 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!

Transport Layer 3-25

Causes/costs of congestion: scenario 1

- two senders, two receivers
- one router, infinite buffers
- output link capacity: R
- no retransmission

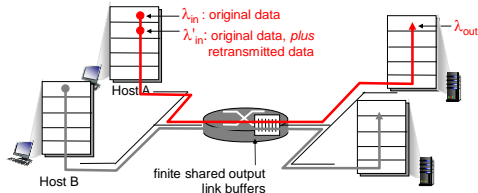


- maximum per-connection throughput: $R/2$
- large delays as arrival rate, λ_{in} , approaches capacity

Transport Layer 3-26

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: $\lambda'_{in} \geq \lambda_{in}$

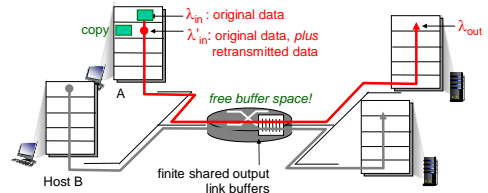
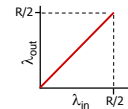


Transport Layer 3-27

Causes/costs of congestion: scenario 2

idealization: perfect knowledge

- sender sends only when router buffers available



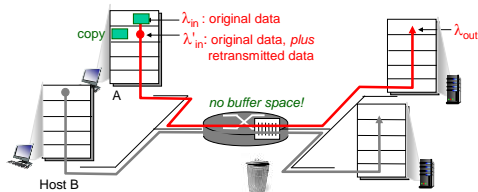
Transport Layer 3-28

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



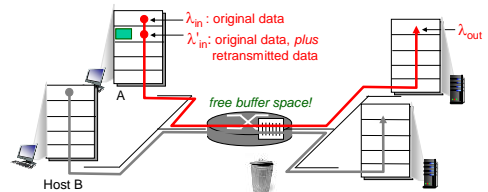
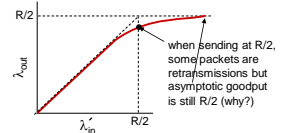
Transport Layer 3-29

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

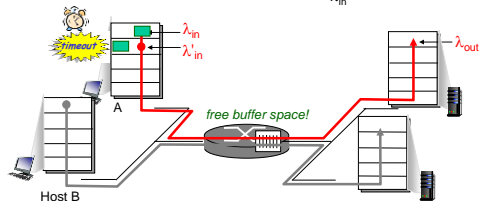


Transport Layer 3-30

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

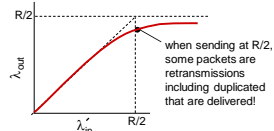


Transport Layer 3-31

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



“costs” of congestion:

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

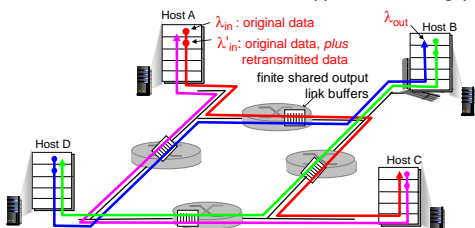
Transport Layer 3-32

Causes/costs of congestion: scenario 3

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

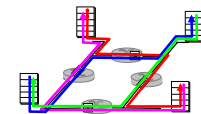
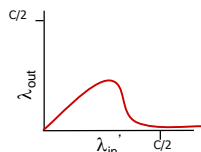
Q: what happens as λ_{in} and λ_{in}' increase?

A: as red λ_{in}' increases, all arriving blue pkts at upper queue are dropped, blue throughput $\rightarrow 0$



Transport Layer 3-33

Causes/costs of congestion: scenario 3



another “cost” of congestion:

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

Transport Layer 3-34

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 for sender to send at

Transport Layer 3-35

Case study: ATM ABR congestion control

ABR: available bit rate:

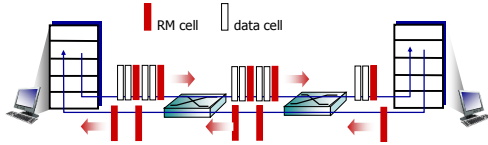
- ❖ “elastic service”
- ❖ if sender’s path “underloaded”:
 - sender should use available bandwidth
- ❖ if sender’s path congested:
 - sender throttled to minimum guaranteed rate

RM (resource management) cells:

- ❖ sent by sender, interspersed with data cells
- ❖ bits in RM cell set by switches (“network-assisted”)
 - NI bit: no increase in rate (mild congestion)
 - CI bit: congestion indication
- ❖ RM cells returned to sender by receiver, with bits intact

Transport Layer 3-36

Case study: ATM ABR congestion control

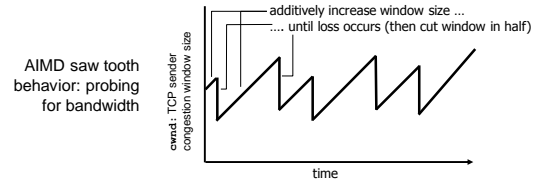


- ❖ two-byte ER (explicit rate) field in RM cell
 - congested switch may lower ER value in cell
 - senders' send rate thus max supportable rate on path
- ❖ EFCI bit in data cells: set to 1 in congested switch
 - if data cell preceding RM cell has EFCI set, receiver sets CI bit in returned RM cell

Transport Layer 3-37

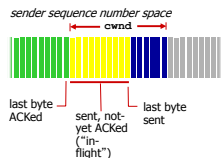
TCP congestion control: additive increase multiplicative decrease

- ❖ **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



Transport Layer 3-38

TCP Congestion Control: details



TCP sending rate:

- ❖ **roughly:** send $cwnd$ bytes, wait RTT for ACKs, then send more bytes

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

- ❖ sender limits transmission:

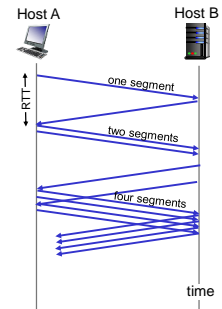
$$\text{LastByteSent} - \text{LastByteAcked} \leq cwnd$$

- ❖ $cwnd$ is dynamic, function of perceived network congestion

Transport Layer 3-39

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



Transport Layer 3-40

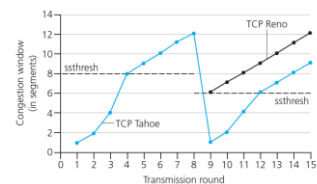
TCP: detecting, reacting to loss

- ❖ 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)

Transport Layer 3-41

TCP: switching from slow start to CA

- Q: when should the exponential increase switch to linear?
- A: when $cwnd$ gets to 1/2 of its value before timeout.

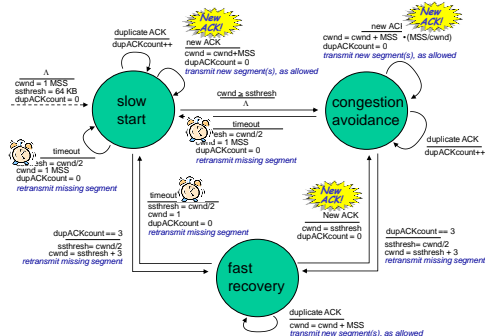


Implementation:

- ❖ variable $ssthresh$
- ❖ on loss event, $ssthresh$ is set to 1/2 of $cwnd$ just before loss event

Transport Layer 3-42

Summary: TCP Congestion Control

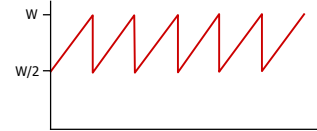


Transport Layer 3-43

TCP throughput

- ❖ avg. TCP thrupt 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 $\frac{3}{4} W$
 - avg. thrupt is $\frac{3}{4} W$ per RTT

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



Transport Layer 3-44

TCP Futures: TCP over "long, fat pipes"

- ❖ 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]:

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

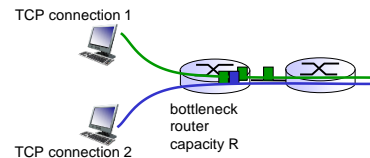
→ to achieve 10 Gbps throughput, need a loss rate of $L = 2 \cdot 10^{-10}$ — a very small loss rate!

- ❖ new versions of TCP for high-speed

Transport Layer 3-45

TCP Fairness

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

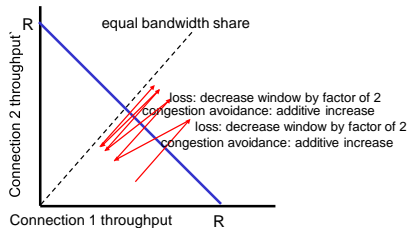


Transport Layer 3-46

Why is TCP fair?

two competing sessions:

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



Transport Layer 3-47

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 $R/2$

Transport Layer 3-48