

Computer Networks

LECTURE 17

Transport Layer

Tamal Biswas
Department of Computer Science
University of Rochester

Computer Networks (Transport Layer) 3-1

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

Computer Networks (Transport Layer) 3-2

Send and Receive Window

- Send and receive buffer sizes can be set on individual sockets by using the **SO_SNDBUF** and **SO_RCVBUF** socket options with the **setsockopt(2)** call.
- Ref: <http://man7.org/linux/man-pages/man7/tcp.7.html>
-

Computer Networks (Transport Layer) 3-3

SO_RCVBUF

- Sets or gets the maximum socket receive buffer in bytes.
- The kernel doubles this value (to allow space for bookkeeping overhead) when it is set using **setsockopt(2)**, and this doubled value is returned by **getsockopt(2)**.
- The default value is set by the `/proc/sys/net/core/rmem_default` file, and the maximum allowed value is set by the `/proc/sys/net/core/rmem_max` file.
- The minimum (doubled) value for this option is **256**.

Computer Networks (Transport Layer) 3-4

SO_SNDBUF

- Sets or gets the maximum socket send buffer in bytes.
- The kernel doubles this value (to allow space for bookkeeping overhead) when it is set using `setsockopt(2)`, and this doubled value is returned by `getsockopt(2)`.
- The default value is set by the `/proc/sys/net/core/wmem_default` file and the maximum allowed value is set by the `/proc/sys/net/core/wmem_max` file.
- The minimum (doubled) value for this option is **2048**.

Computer Networks (Transport Layer) 3-5

Example:

```
int getsockopt(int sockfd, int level, int optname, void *optval, socklen_t *optlen);
```

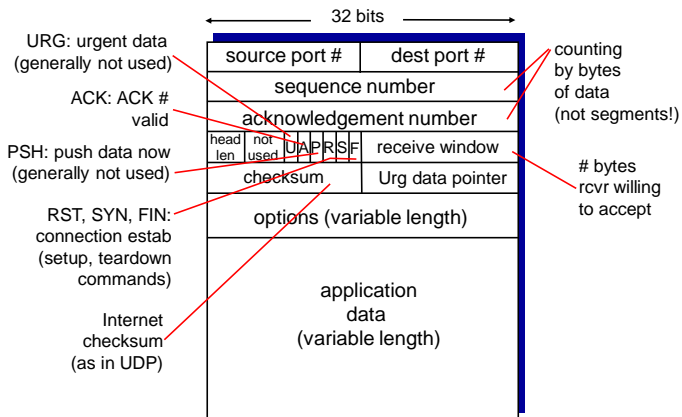
```
int setsockopt(int sockfd, int level, int optname, const void *optval, socklen_t optlen);
```

Example:

```
int n = 1024 * 3;
setsockopt(socket, SOL_SOCKET, SO_RCVBUF, &n, sizeof(n));
```

Computer Networks (Transport Layer) 3-6

TCP segment structure



Computer Networks (Transport Layer) 3-7

TCP seq. numbers, ACKs

sequence numbers:

- byte stream “number” of first byte in segment’s data

acknowledgements:

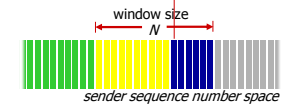
- seq # of next byte expected from other side
- cumulative ACK

Q: how receiver handles out-of-order segments

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

outgoing segment from sender

source port #	dest port #
sequence number	acknowledgement number
checksum	urg pointer



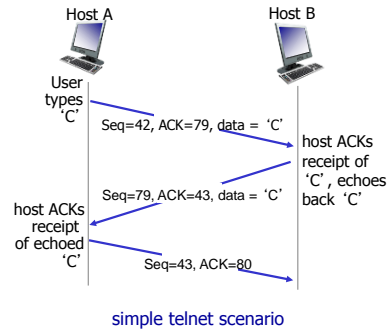
sent, not yet ACKed (“in-flight”)

incoming segment to sender

source port #	dest port #
sequence number	acknowledgement number
checksum	urg pointer

Computer Networks (Transport Layer) 3-8

TCP seq. numbers, ACKs



Computer Networks (Transport Layer) 3-9

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?

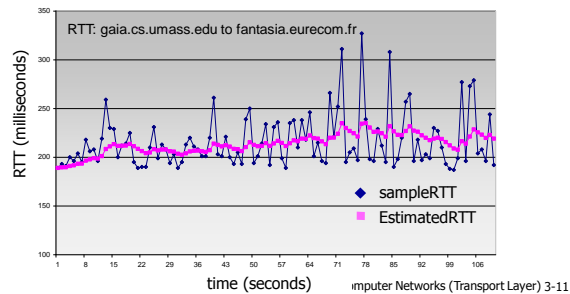
- SampleRTT:** measured time from segment transmission until ACK receipt
 - ignore retransmissions
- SampleRTT** will vary, want estimated RTT “smoother”
 - average several recent measurements, not just current **SampleRTT**

Computer Networks (Transport Layer) 3-10

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”

* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

Computer Networks (Transport Layer) 3-12

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

Computer Networks (Transport Layer) 3-13

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

Computer Networks (Transport Layer) 3-14

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

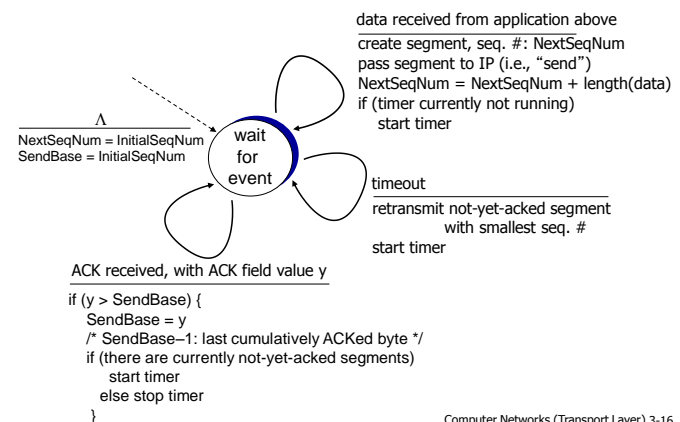
- restart timer

ack rcvd:

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

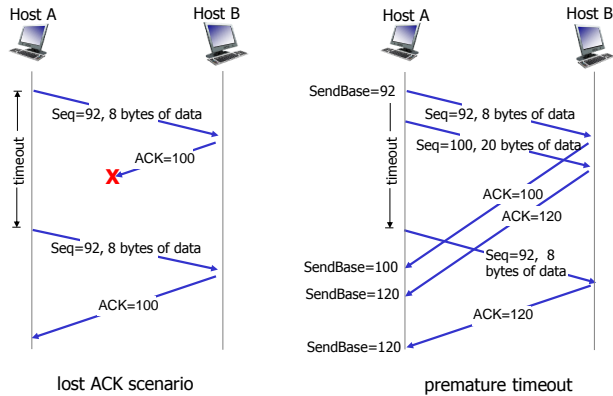
Computer Networks (Transport Layer) 3-15

TCP sender (simplified)



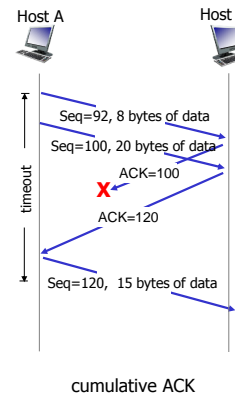
Computer Networks (Transport Layer) 3-16

TCP: retransmission scenarios



Computer Networks (Transport Layer) 3-17

TCP: retransmission scenarios



Computer Networks (Transport Layer) 3-18

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 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

Computer Networks (Transport Layer) 3-19

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"), resends unacked segment with smallest seq #
 likely that unacked segment lost, so don't wait for timeout

Computer Networks (Transport Layer) 3-20

Chapter 3 outline

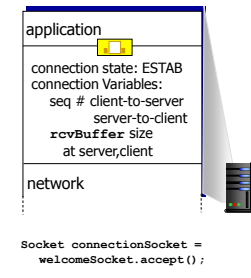
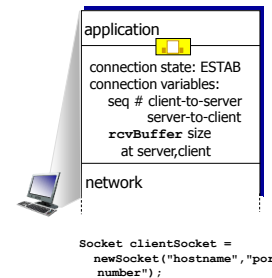
- 3.1 transport-layer services
- 3.2 multiplexing and demultiplexing
- 3.3 connectionless transport: UDP
- 3.4 principles of reliable data transfer
- 3.5 connection-oriented transport: TCP
 - segment structure
 - reliable data transfer
 - flow control
 - connection management
- 3.6 principles of congestion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

Computer Networks (Transport Layer) 3-25

Connection Management

before exchanging data, sender/receiver “handshake”:

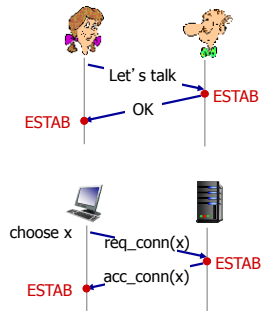
- agree to establish connection (each knowing the other willing to establish connection)
- agree on connection parameters



Computer Networks (Transport Layer) 3-26

Agreeing to establish a connection

2-way handshake:



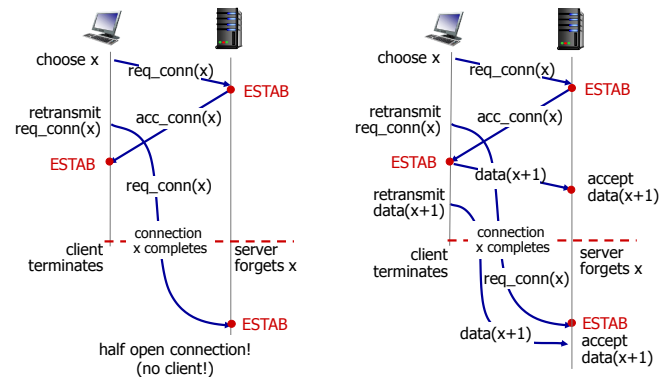
Q: will 2-way handshake always work in network?

- variable delays
- retransmitted messages (e.g. req_conn(x)) due to message loss
- message reordering
- can't "see" other side

Computer Networks (Transport Layer) 3-27

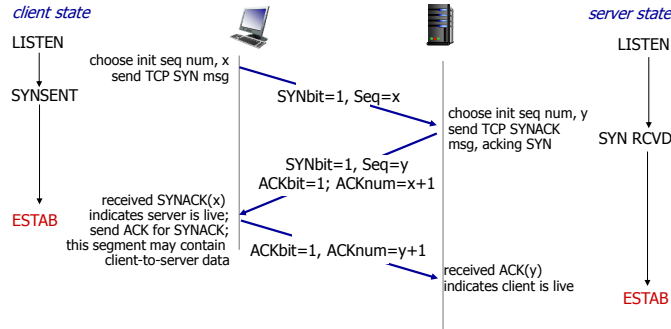
Agreeing to establish a connection

2-way handshake failure scenarios:



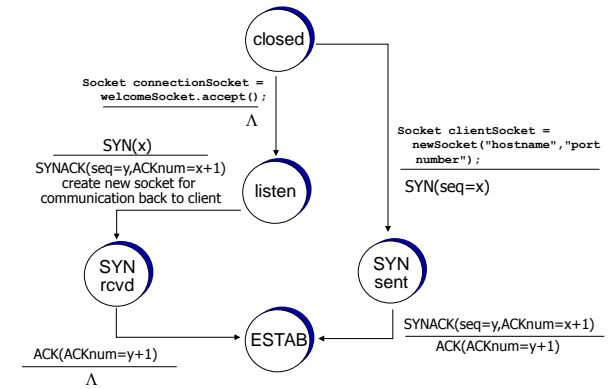
Computer Networks (Transport Layer) 3-28

TCP 3-way handshake



Computer Networks (Transport Layer) 3-29

TCP 3-way handshake: FSM



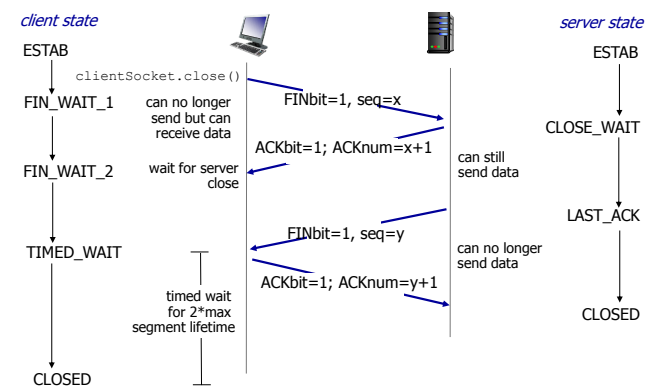
Computer Networks (Transport Layer) 3-30

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

Computer Networks (Transport Layer) 3-31

TCP: closing a connection



Computer Networks (Transport Layer) 3-32

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

Computer Networks (Transport Layer) 3-33

Principles of congestion control

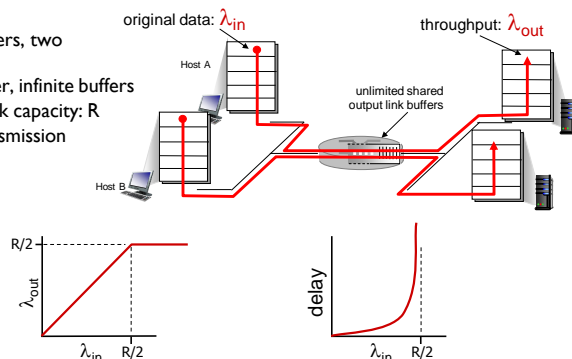
congestion:

- 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!

Computer Networks (Transport Layer) 3-34

Causes/costs of congestion: scenario 1

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

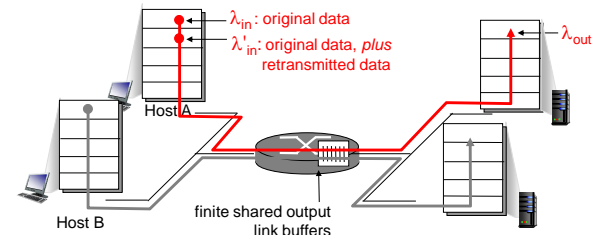


- maximum per-connection throughput: $R/2$
- ❖ large delays as arrival rate, λ_{in} , approaches capacity

Computer Networks (Transport Layer) 3-35

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

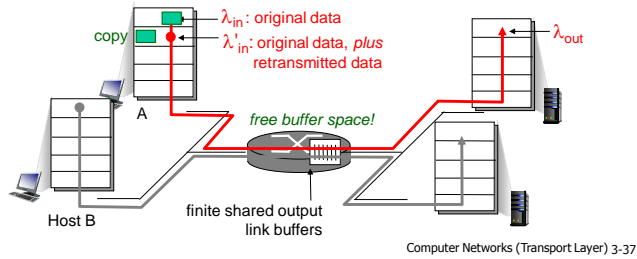
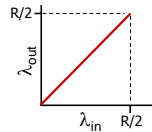


Computer Networks (Transport Layer) 3-36

Causes/costs of congestion: scenario 2

idealization: perfect knowledge

- sender sends only when router buffers available

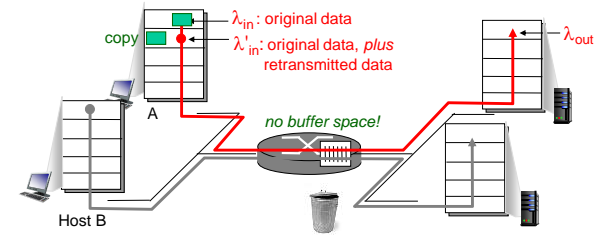


Causes/costs of congestion: scenario 2

Idealization: *known loss*

packets can be lost, dropped at router due to full buffers

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

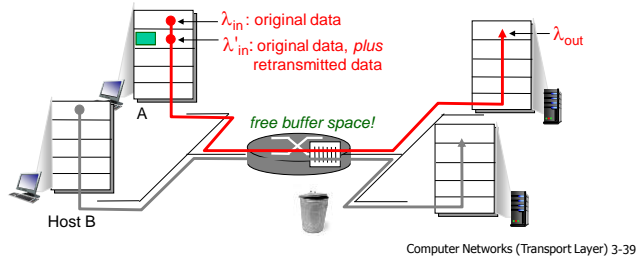
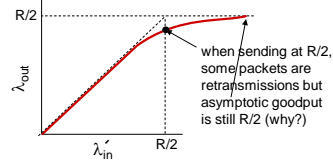


Causes/costs of congestion: scenario 2

Idealization: *known loss*

packets can be lost, dropped at router due to full buffers

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

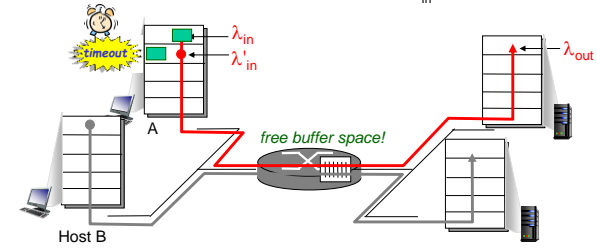
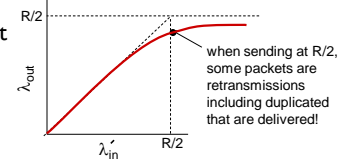


Causes/costs of congestion: scenario 2

Realistic: *duplicates*

- packets can be lost, dropped at router due to full buffers

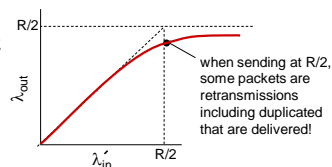
- sender times out prematurely, sending *two* copies, both of which are delivered



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”
- unnecessary retransmissions: link carries multiple copies of pkt
 - decreasing goodput

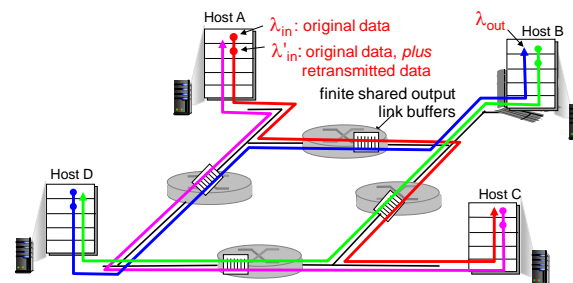
Computer Networks (Transport Layer) 3-41

Causes/costs of congestion: scenario 3

- four senders
- multihop paths
- timeout/retransmit

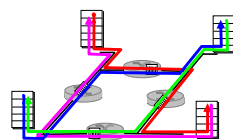
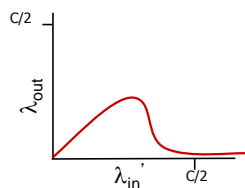
Q: what happens as λ_{in} and λ_{in}' increase?

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



Computer Networks (Transport Layer) 3-42

Causes/costs of congestion: scenario 3



another “cost” of congestion:

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

Computer Networks (Transport Layer) 3-43

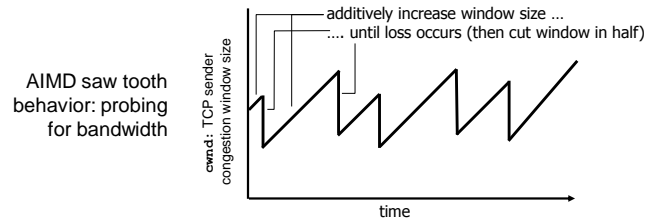
Chapter 3 outline

- 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

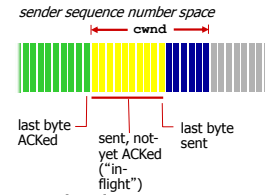
Computer Networks (Transport Layer) 3-44

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



TCP Congestion Control: details



TCP sending rate:

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

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

- sender limits transmission:

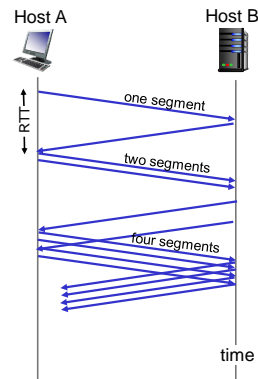
$$\text{LastByteSent} - \text{LastByteAked} \leq \text{cwnd}$$

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

Computer Networks (Transport Layer) 3-46

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



Computer Networks (Transport Layer) 3-47

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)

Computer Networks (Transport Layer) 3-48

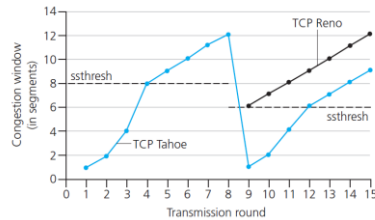
TCP: switching from slow start to CA

Q: when should the exponential increase switch to linear?

A: when **cwnd** gets to 1/2 of its value before timeout.

Implementation:

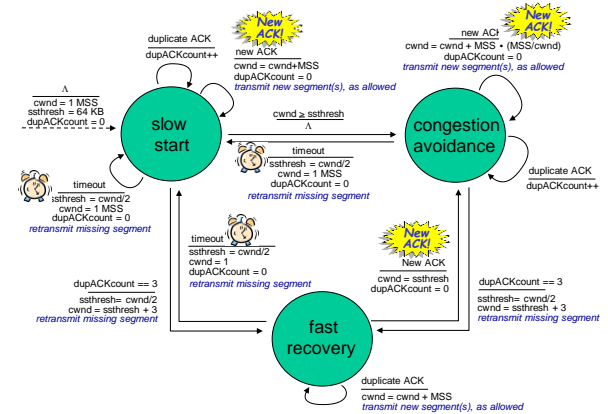
- variable **ssthresh**
- on loss event, **ssthresh** is set to 1/2 of **cwnd** just before loss event



* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

Computer Networks (Transport Layer) 3-49

Summary: TCP Congestion Control



Computer Networks (Transport Layer) 3-50

TCP throughput

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

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



Computer Networks (Transport Layer) 3-51

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]:

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

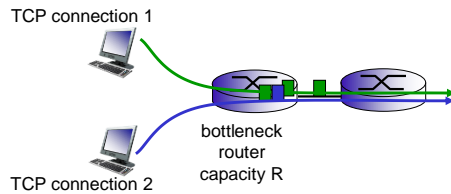
→ to achieve 10 Gbps throughput, need a loss rate of $L = 2 \cdot 10^{-10}$ – a very small loss rate!

- new versions of TCP for high-speed
- Ref: <http://ccr.sigcomm.org/archive/1997/jul97/ccr-9707-mathis.pdf>

Computer Networks (Transport Layer) 3-52

TCP Fairness

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

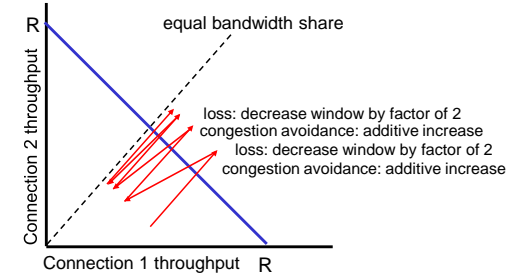


Computer Networks (Transport Layer) 3-53

Why is TCP fair?

two competing sessions:

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



Computer Networks (Transport Layer) 3-54

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 $R/2$

Computer Networks (Transport Layer) 3-55

Summary

- principles behind transport layer services:
 - multiplexing, demultiplexing
 - reliable data transfer
 - flow control
 - congestion control
- instantiation, implementation in the Internet
 - UDP
 - TCP

Computer Networks (Transport Layer) 3-56

Disclaimer

- Parts of the lecture slides contain original work of James Kurose and Keith Ross. The slides are intended for the sole purpose of instruction of computer networks at the University of Rochester. All copyrighted materials belong to their original owner(s).