

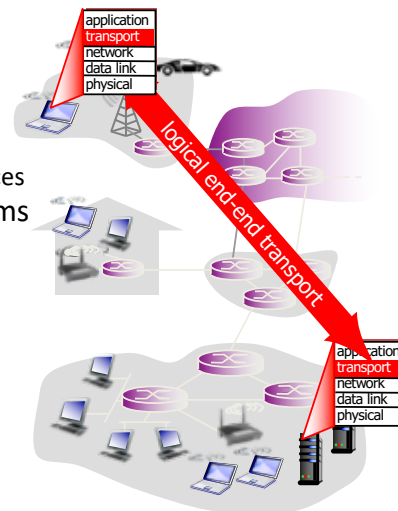
Transport Layer

- Reference
 - Chapter 03, *Computer Networking: A Top Down Approach*, 6/E, Jim Kurose, Keith Ross, Addison-Wesley
 - Adapted from part of the slides provided by the authors

CS3700 - Instructor: Dr. Zhu

Transport services and protocols

- Provide logical communication between app processes running on different hosts
 - Network layer: logical communication between hosts
 - Relies on, enhances, network layer services
- Transport protocols run in end systems
 - send side: breaks app messages into segments, passes to network layer
 - rcv side: reassembles segments into messages, passes to app layer
- more than one transport protocol available to apps
 - Internet: TCP and UDP



CS3700 - Instructor: Dr. Zhu

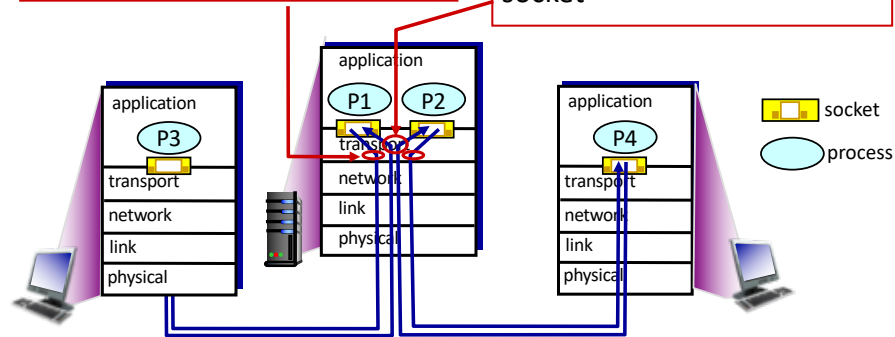
Multiplexing/demultiplexing

multiplexing at sender:

handle data from multiple sockets, add transport header (later used for demultiplexing)

demultiplexing at receiver:

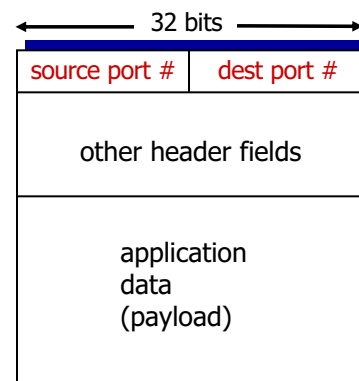
use header info to deliver received segments to correct socket



CS3700 - Instructor: Dr. Zhu

How demultiplexing works

- host receives IP datagrams
 - each datagram has source IP address, destination IP address
 - each datagram carries one transport-layer segment
 - each segment has source, destination port number
- host uses *IP addresses & port numbers* to direct segment to appropriate socket



TCP/UDP segment format

CS3700 - Instructor: Dr. Zhu

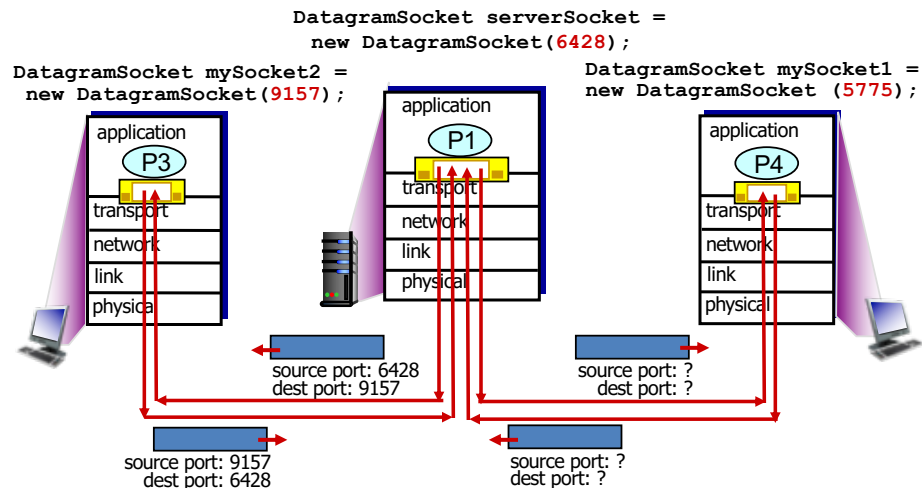
Connectionless demultiplexing

- *recall*: created socket has host-local port #:


```
DatagramSocket mySocket1 =
  new DatagramSocket(12534);
```
 - *recall*: when creating datagram to send into UDP socket, must specify
 - destination IP address
 - destination port #
 - when host receives UDP segment:
 - checks destination port # in segment
 - directs UDP segment to socket with that port #
- IP datagrams with *same dest. port #*, but different source IP addresses and/or source port numbers will be directed to *same socket* at dest

CS3700 - Instructor: Dr. Zhu

Connectionless demux: example



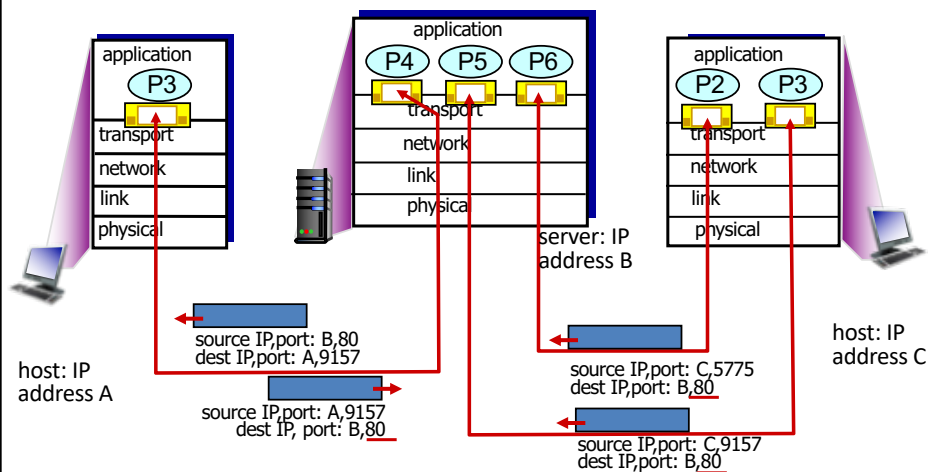
CS3700 - Instructor: Dr. Zhu

Connection-oriented demux

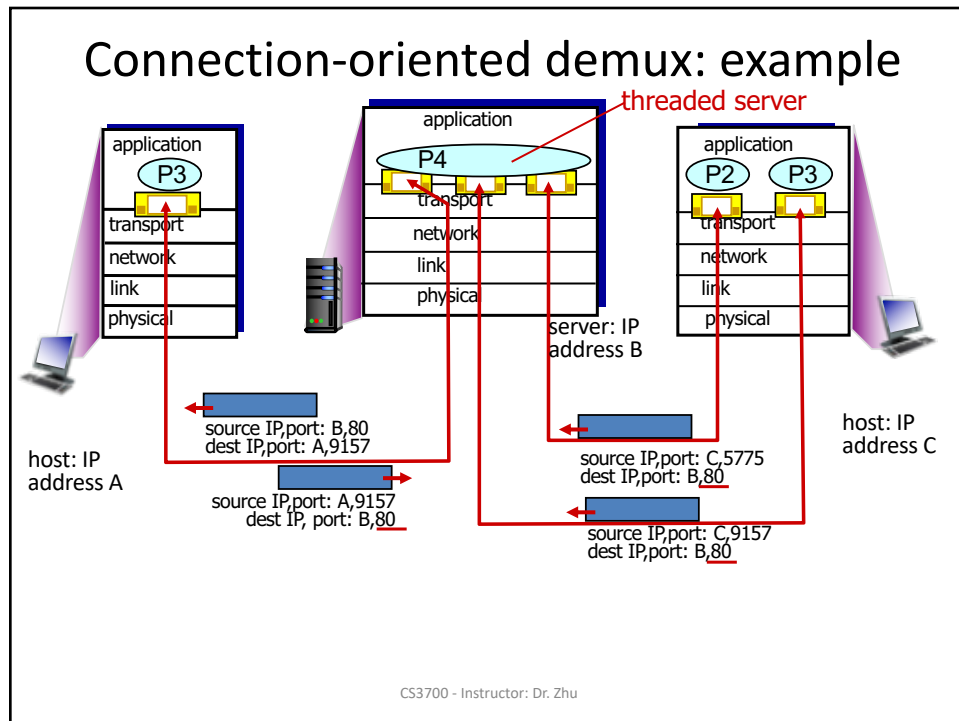
- TCP socket identified by 4-tuple:
 - source IP address
 - source port number
 - dest IP address
 - dest port number
- demux: receiver uses all four values to direct segment to appropriate socket
- server host may support many simultaneous TCP sockets:
 - each socket identified by its own 4-tuple
- web servers have different sockets for each connecting client
 - non-persistent HTTP will have different socket for each request

CS3700 - Instructor: Dr. Zhu

Connection-oriented demux: example



CS3700 - Instructor: Dr. Zhu

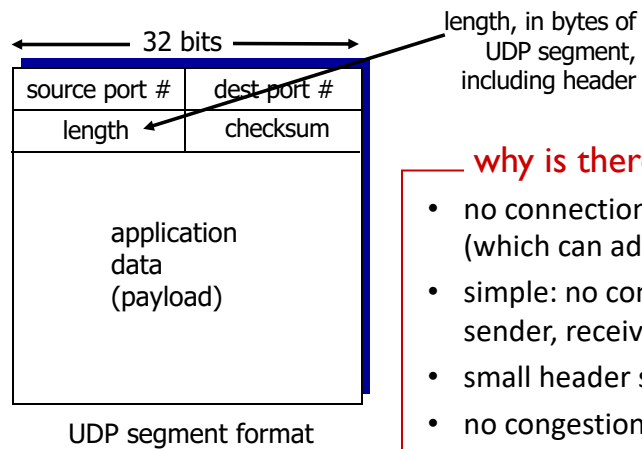


UDP: User Datagram Protocol [RFC 768]

- “no frills,” “bare bones” Internet transport protocol
- “best effort” service, UDP segments may be:
 - lost
 - delivered out-of-order to app
- connectionless:
 - no handshaking between UDP sender, receiver
 - each UDP segment handled independently of others
- UDP use:
 - streaming multimedia apps (loss tolerant, rate sensitive)
 - DNS
 - SNMP
- reliable transfer over UDP:
 - add reliability at application layer
 - application-specific error recovery!

CS3700 - Instructor: Dr. Zhu

UDP: segment header



why is there a UDP?

- no connection establishment (which can add delay)
- simple: no connection state at sender, receiver
- small header size
- no congestion control: UDP can blast away as fast as desired

CS3700 - Instructor: Dr. Zhu

UDP checksum

Goal: detect “errors” (e.g., flipped bits) in transmitted segment

sender:

- treat segment contents, including header fields, as sequence of 16-bit integers
- checksum: addition (one's complement sum) of segment contents
- sender puts checksum value into UDP checksum field

receiver:

- compute checksum of received segment
- check if computed checksum equals checksum field value:
 - NO - error detected
 - YES - no error detected. *But maybe errors nonetheless? More later*

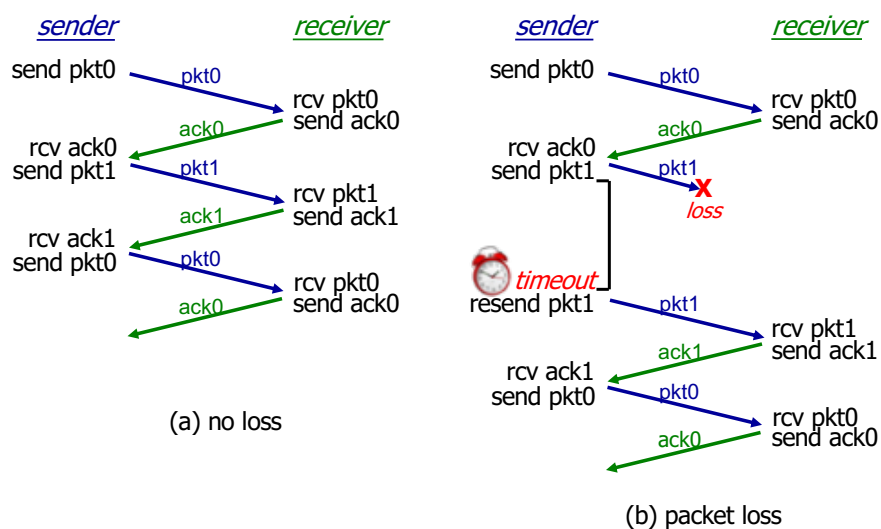
CS3700 - Instructor: Dr. Zhu

Reliable Data Transfer over Unreliable Channel

- underlying channel may flip bits in packet
 - checksum to detect bit errors
- the question: how to recover from errors:
 - acknowledgements (**ACKs**): receiver explicitly tells sender that pkt received OK
 - negative acknowledgements (**NAKs**): receiver explicitly tells sender that pkt had errors
 - sender retransmits pkt on receipt of NAK
- Underlying channel can also lose packets (data-pkt, ACKs)
 - sender waits “reasonable” amount of time for ACK
 - retransmits if no ACK received in this time
 - if pkt (or ACK) just delayed (not lost):
 - retransmission will be duplicate, but seq. #'s already handles this
 - receiver must specify seq # of pkt being ACKed
 - requires countdown timer

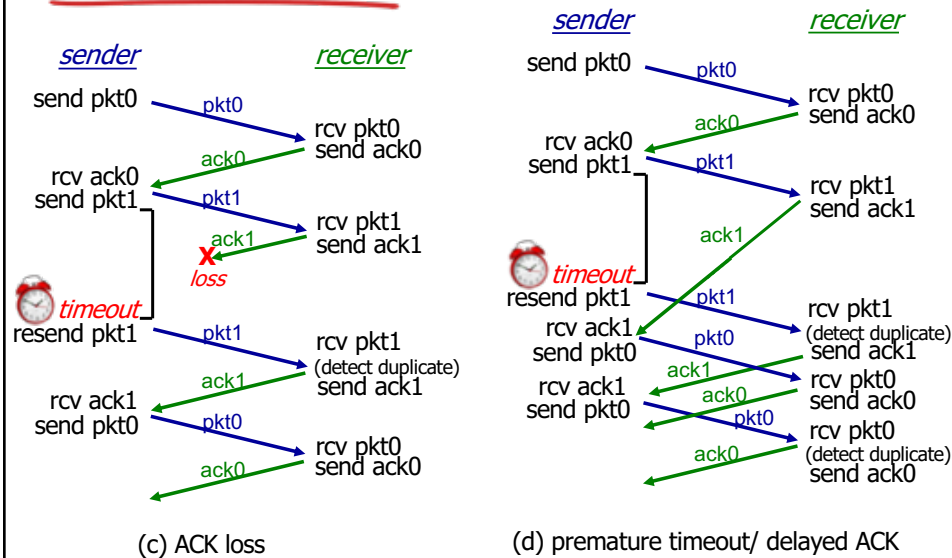
CS3700 - Instructor: Dr. Zhu

rdt3.0 in action



CS3700 - Instructor: Dr. Zhu

rdt3.0 in action



CS3700 - Instructor: Dr. Zhu

Performance of rdt3.0

- rdt3.0 is correct, but performance stinks
- e.g.: 1 Gbps link, 15 ms prop. delay, 8000 bit packet:

$$D_{trans} = \frac{L}{R} = \frac{8000 \text{ bits}}{10^9 \text{ bits/sec}} = 8 \text{ microsecs}$$

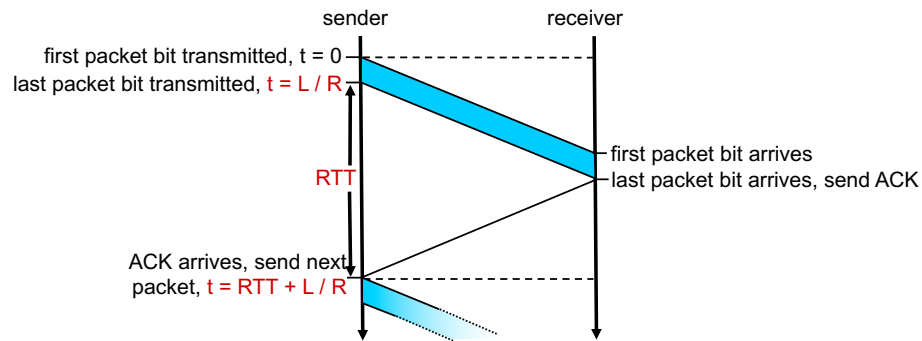
- U_{sender} : **utilization** – fraction of time sender busy sending

$$U_{\text{sender}} = \frac{L/R}{RTT + L/R} = \frac{.008}{30.008} = 0.00027$$

- if $RTT=30$ msec, 1KB pkt every 30 msec: 33kB/sec thrupt over 1 Gbps link
- network protocol limits use of physical resources!

CS3700 - Instructor: Dr. Zhu

rdt3.0: stop-and-wait operation



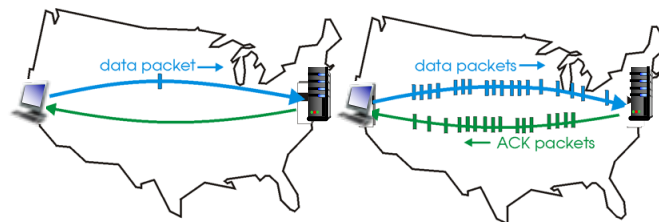
$$U_{\text{sender}} = \frac{L/R}{RTT + L/R} = \frac{.008}{30.008} = 0.00027$$

CS3700 - Instructor: Dr. Zhu

Pipelined protocols

pipelining: sender allows multiple, “in-flight”, yet-to-be-acknowledged pkts

- range of sequence numbers must be increased
- buffering at sender and/or receiver



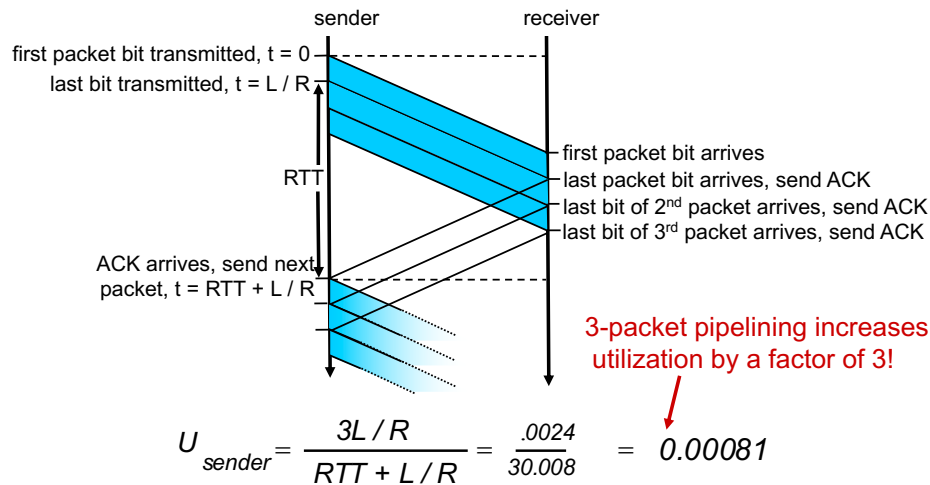
(a) a stop-and-wait protocol in operation

(b) a pipelined protocol in operation

- two generic forms of pipelined protocols: *go-Back-N*, *selective repeat*

CS3700 - Instructor: Dr. Zhu

Pipelining: increased utilization



CS3700 - Instructor: Dr. Zhu

Pipelined protocols: overview

Go-back-N:

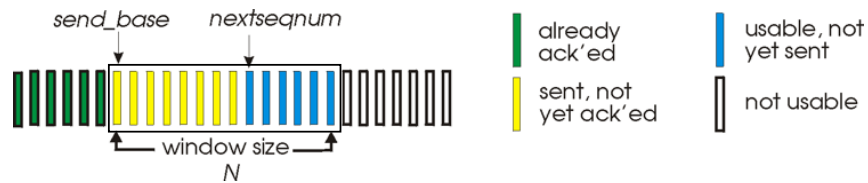
- sender can have up to N unacked packets in pipeline
- receiver only sends *cumulative ack*
 - doesn't ack packet if there's a gap
- sender has timer for oldest unacked packet
 - when timer expires, retransmit *all* unacked packets

Selective Repeat:

- sender can have up to N unack'ed packets in pipeline
- rcvr sends *individual ack* for each packet
- sender maintains timer for each unacked packet
 - when timer expires, retransmit only that unacked packet

Transport Layer 3-20

Go-Back-N: sender



- k-bit seq # in pkt header
- “window” of up to N, consecutive unack'd pkts allowed

Sender:

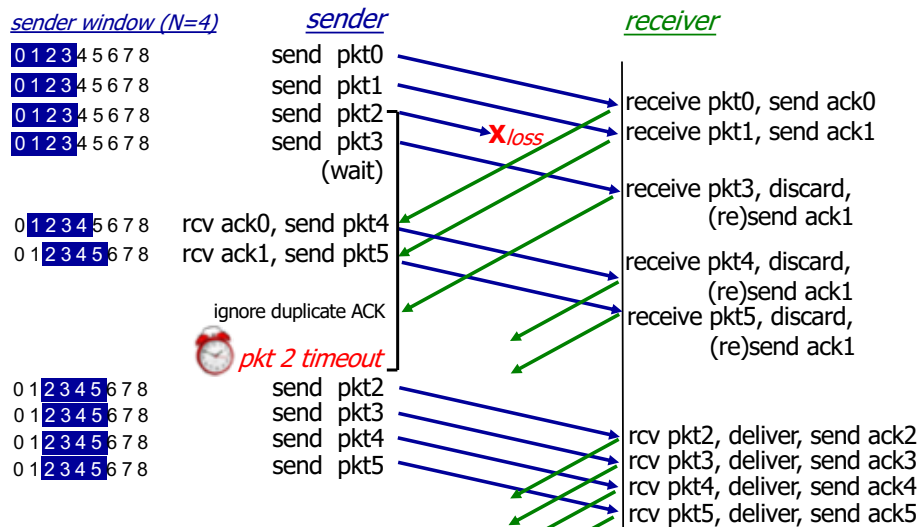
- ACK(n): ACKs all pkts up to, including seq # n - “cumulative ACK”
 - may receive duplicate ACKs (see receiver)
- Only one timer for oldest in-flight pkt
- timeout(n): retransmit packet n and all higher seq # pkts in window

Receiver:

- ACK-only: always send ACK for correctly-received pkt with highest in-order seq #
- out-of-order pkt:
 - discard (don't buffer): *no receiver buffering!*
 - re-ACK pkt with highest in-order seq #

CS3700 - Instructor: Dr. Zhu

GBN in action



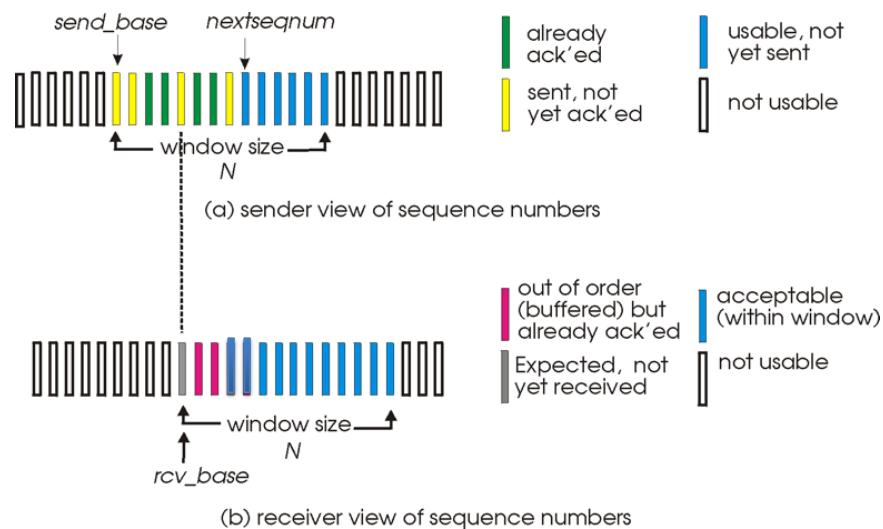
CS3700 - Instructor: Dr. Zhu

Selective repeat

- receiver *individually* acknowledges all correctly received pkts
 - buffers pkts, as needed, for eventual in-order delivery to upper layer
- sender only resends pkts for which ACK not received
 - sender timer for each unACKed pkt
- sender window
 - N consecutive seq #'s
 - limits seq #'s of sent, unACKed pkts

CS3700 - Instructor: Dr. Zhu

Selective repeat: sender, receiver windows



CS3700 - Instructor: Dr. Zhu

Selective repeat

sender

data from above:

- if next available seq # in window, send pkt

timeout(n):

- resend pkt n, restart timer

ACK(n) in [sendbase, sendbase+N-1]:

- mark pkt n as received
- if n smallest unACKed pkt, advance window base to next unACKed seq #

receiver

pkt n in [rcvbase, rcvbase+N-1]

- send ACK(n)
- out-of-order: buffer
- in-order: deliver (also deliver buffered, in-order pkts), advance window to next not-yet-received pkt

pkt n in [rcvbase-N, rcvbase-1]

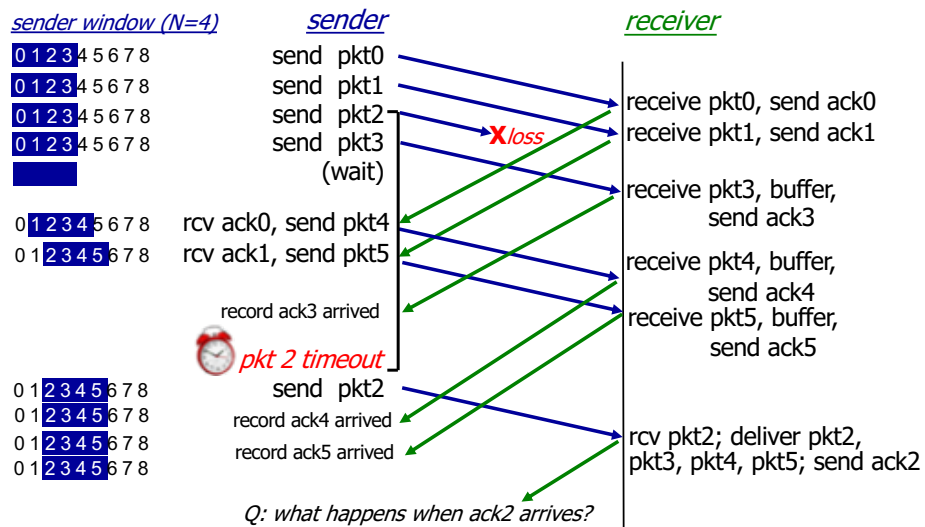
- ACK(n)

otherwise:

- ignore

CS3700 - Instructor: Dr. Zhu

Selective repeat in action



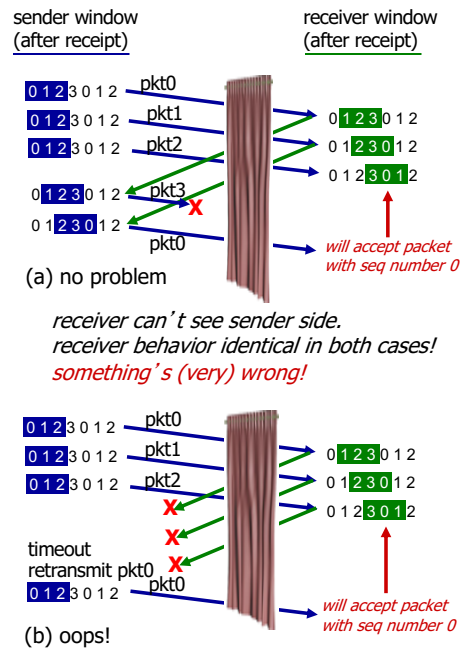
CS3700 - Instructor: Dr. Zhu

Selective repeat: dilemma

example:

- seq #'s: 0, 1, 2, 3
- window size=3
- receiver sees no difference in two scenarios!
- duplicate data accepted as new in (b)

Q: what relationship between seq # size and window size to avoid problem in (b)?



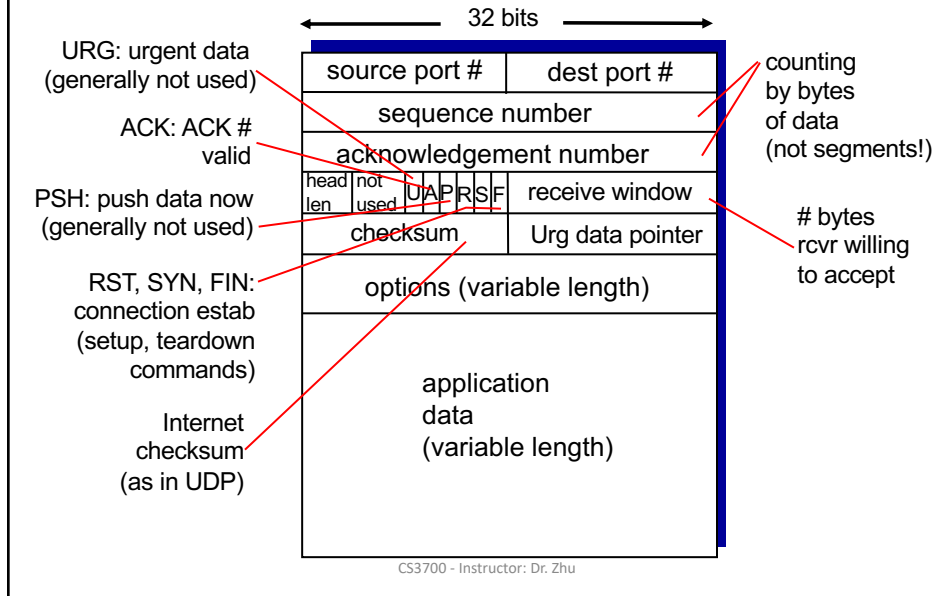
CS3700 - Instructor: Dr. Zhu

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, which is the max of application-layer data (payload) in the segment and does not include the header length.
- **connection-oriented:**
 - handshaking (exchange of control msgs) initiates sender, receiver state before data exchange
- **flow controlled:**
 - sender will not overwhelm receiver

CS3700 - Instructor: Dr. Zhu

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

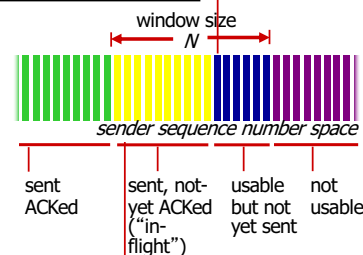
Both sequence # and Ack # may be included in one segment

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

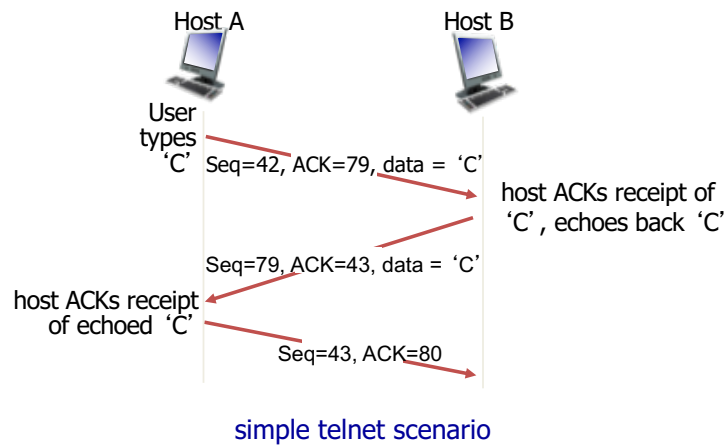


incoming segment to sender

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

CS3700 - Instructor: Dr. Zhu

TCP seq. numbers, ACKs



CS3700 - Instructor: Dr. Zhu

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

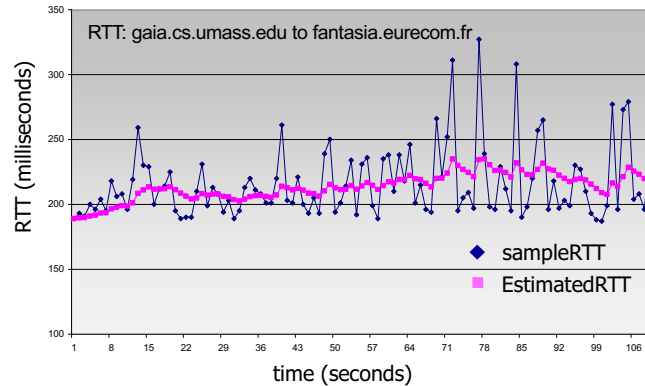
Q: how to estimate RTT?

- **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**

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:

$$\begin{aligned} \text{DevRTT} &= (1 - \beta) * \text{DevRTT} + \\ &\quad \beta * |\text{SampleRTT} - \text{EstimatedRTT}| \\ &\quad \text{(typically, } \beta = 0.25) \end{aligned}$$

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



↑
estimated RTT

↑
“safety margin”

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

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:

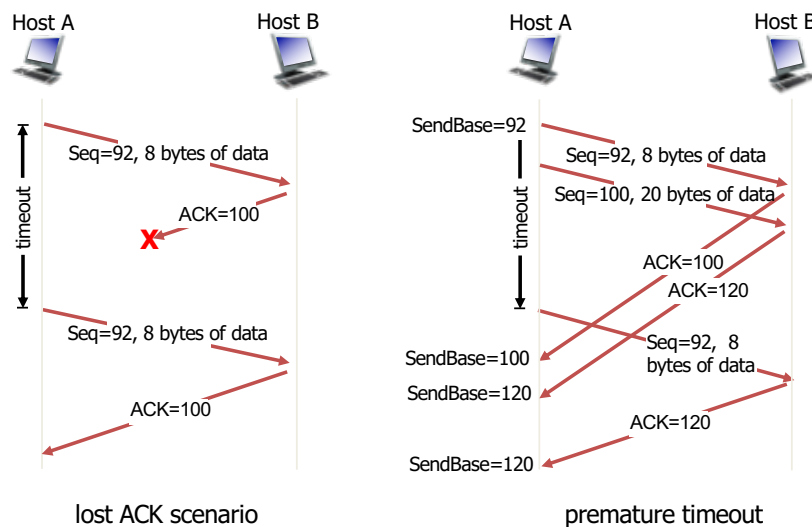
- retransmit ONLY THE segment that caused timeout
- restart timer

ack rcvd:

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

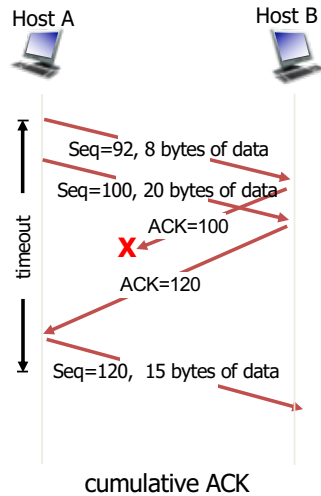
CS3700 - Instructor: Dr. Zhu

TCP: retransmission scenarios



CS3700 - Instructor: Dr. Zhu

TCP: retransmission scenarios



CS3700 - Instructor: Dr. Zhu

TCP ACK generation [RFC 1122, RFC 2581]

<i>event at receiver</i>	<i>TCP receiver action</i>
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

CS3700 - Instructor: Dr. Zhu

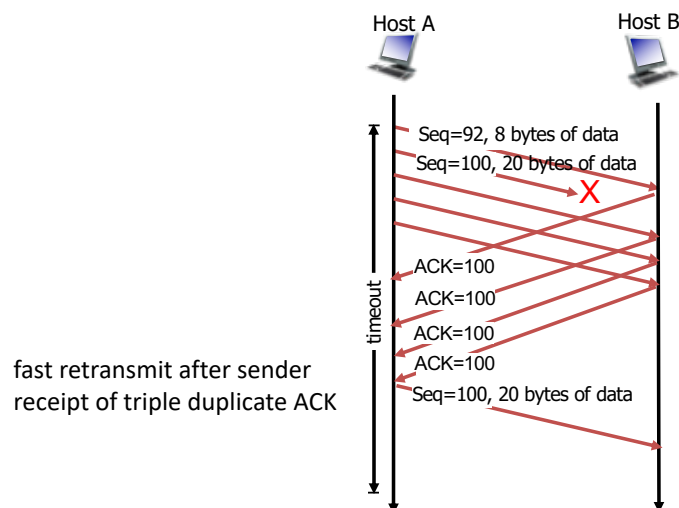
TCP fast retransmit

- 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 **more** ACKs for same data (“triple duplicate ACKs”), resend unacked segment with smallest seq #
 - likely that unacked segment lost, so don't wait for timeout

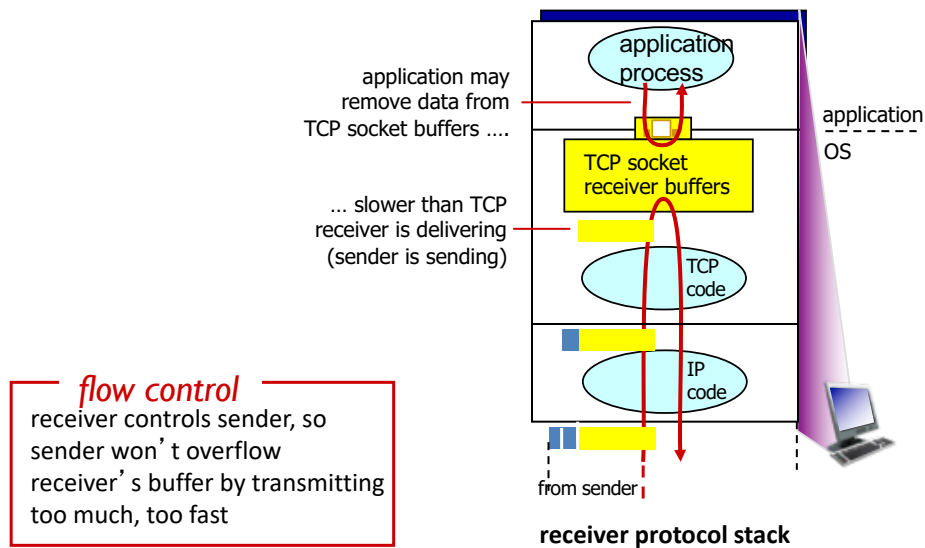
CS3700 - Instructor: Dr. Zhu

TCP fast retransmit



CS3700 - Instructor: Dr. Zhu

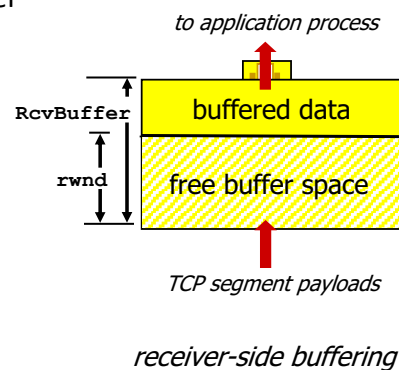
TCP flow control



CS3700 - Instructor: Dr. Zhu

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

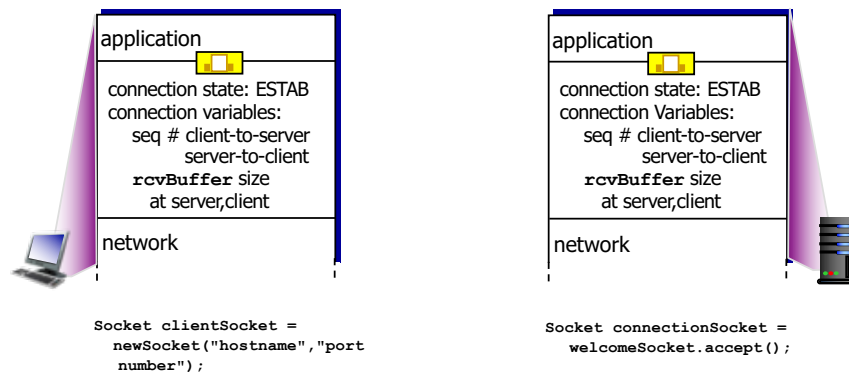


CS3700 - Instructor: Dr. Zhu

Connection Management

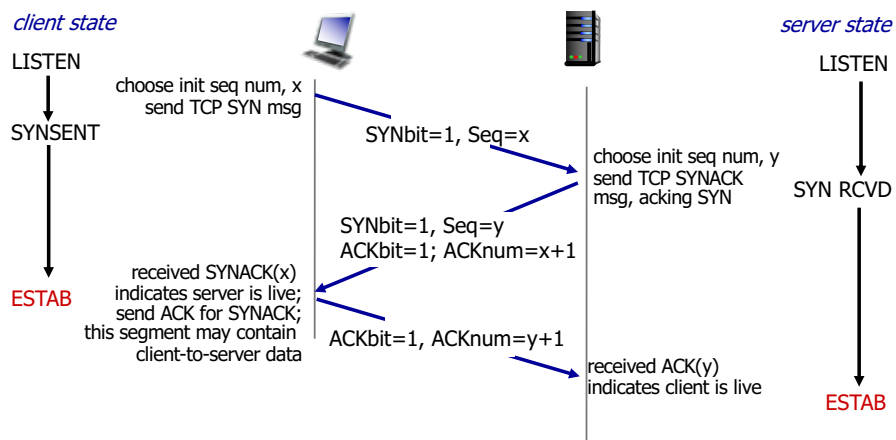
before exchanging data, sender/receiver “handshake”:

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



CS3700 - Instructor: Dr. Zhu

Establish a connection: TCP 3-way handshake



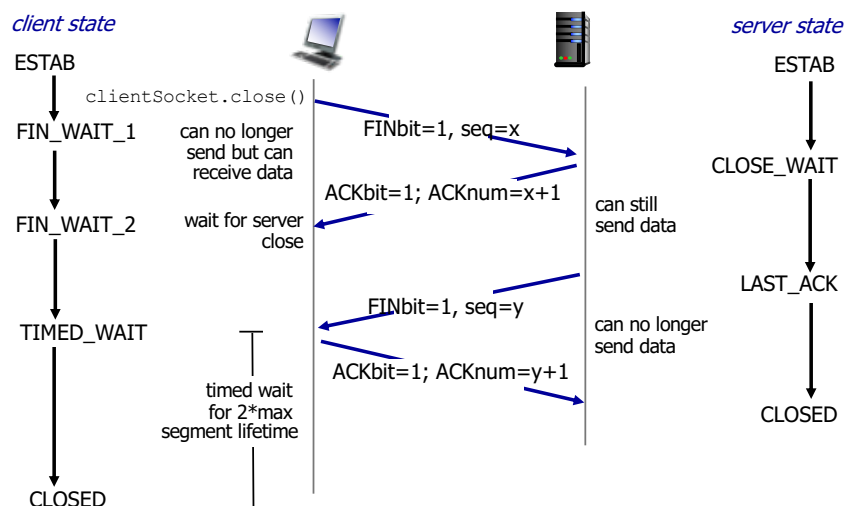
CS3700 - Instructor: Dr. Zhu

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

CS3700 - Instructor: Dr. Zhu

TCP: closing a connection



CS3700 - Instructor: Dr. Zhu

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)
- long delays (queueing in router buffers)
- a top-10 problem!

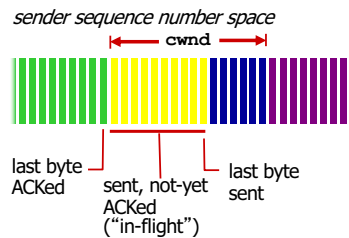
CS3700 - Instructor: Dr. Zhu

Two approaches toward 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

CS3700 - Instructor: Dr. Zhu

TCP Congestion Control: cwnd



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:

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

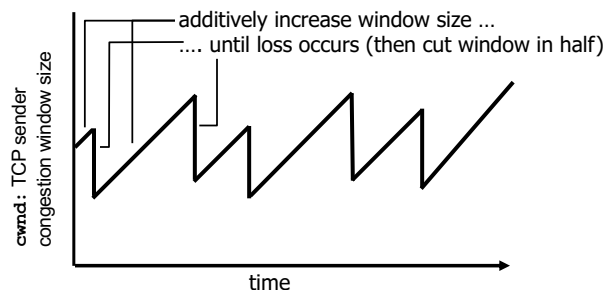
- **cwnd** is dynamic, function of perceived network congestion

CS3700 - Instructor: Dr. Zhu

AIMD in TCP congestion control

- AIMD: 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(Maximum Segment Size) every RTT until loss detected
 - *multiplicative decrease*: cut **cwnd** in half after loss

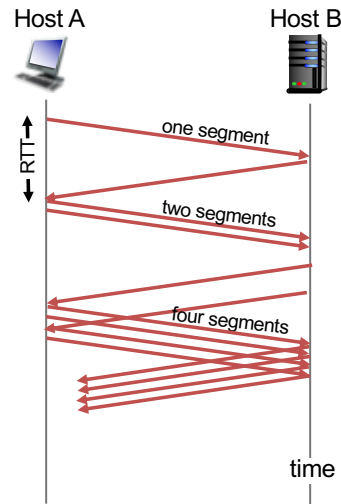
AIMD saw tooth behavior: probing for bandwidth



CS3700 - Instructor: Dr. Zhu

Slow Start in TCP Congestion Control

- when connection begins, increase rate exponentially until first loss event:
 - initially **cwnd** = 1 MSS
 - double **cwnd** every RTT
 - done by incrementing **cwnd** by 1 MSS for every ACK received
- summary: initial rate is slow but ramps up exponentially fast



CS3700 - Instructor: Dr. Zhu

TCP: detecting, reacting to loss

- **ssthresh** is always set as half of the **cwnd** just before the loss
- loss indicated by timeout (for both TCP RENO and TCP Tahoe):
 - **cwnd** set to 1 MSS;
 - window then grows exponentially (slow start) to threshold, then grows linearly (congestion avoidance).
- loss indicated by 3 duplicate ACKs: TCP RENO (newer version of TCP)
 - dup ACKs indicate network capable of delivering some segments
 - **cwnd** is cut in half (**plus 3 for good measure to account for the triple duplicate ACKs received**) then grows linearly
- loss indicated by 3 duplicate ACKs: TCP Tahoe (an early version of TCP)
 - Same as timeout: sets **cwnd** to 1 MSS and then grows exponentially to threshold and then grows linearly

CS3700 - Instructor: Dr. Zhu

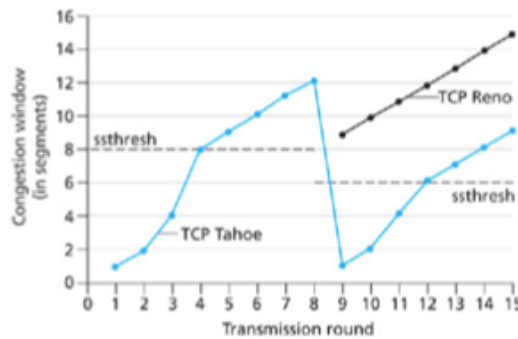
TCP: switching from slow start to CA

Q: When should the exponential increase switch to linear? Or how to choose the threshold value?

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



1 transmission round is about 1 RTT

CS3700 - Instructor: Dr. Zhu

Summary: TCP Congestion Control

