

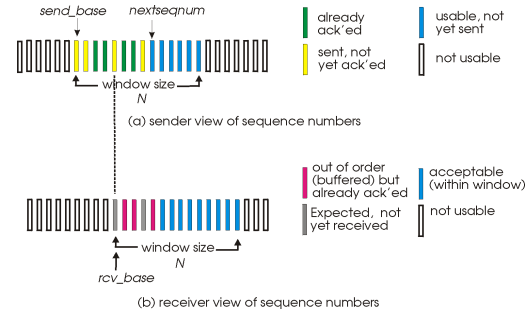
Selective Repeat

- 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
 - again limits seq #'s of sent, unACKed pkts

Transport Layer

3-1

Selective repeat: sender, receiver windows



Transport Layer

3-2

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]:

- 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

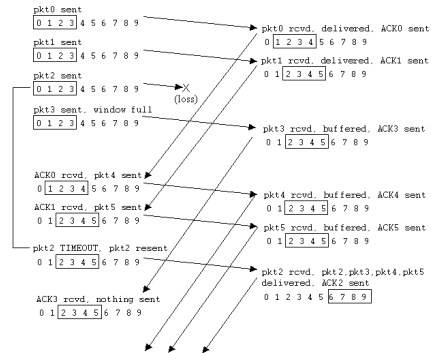
pkt n in [rcvbase-N, rcvbase-1]

- ACK(n)
- otherwise: ignore

Transport Layer

3-3

Selective repeat in action



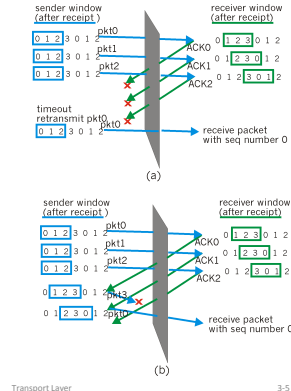
3-4

Selective repeat: dilemma

- Example:**
- seq #'s: 0, 1, 2, 3
 - window size=3

- receiver sees no difference in two scenarios!
- incorrectly passes duplicate data as new in (a)

Q: what relationship between seq # size and window size?



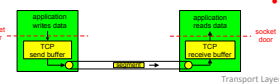
Transport Layer

3-5

TCP: Overview

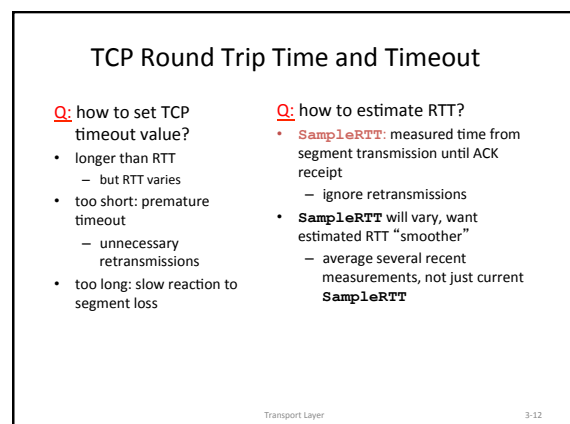
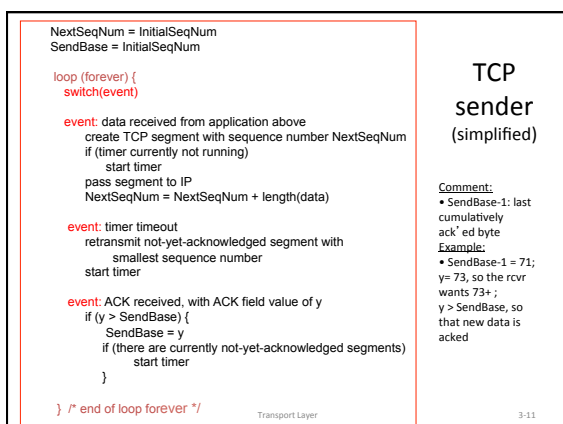
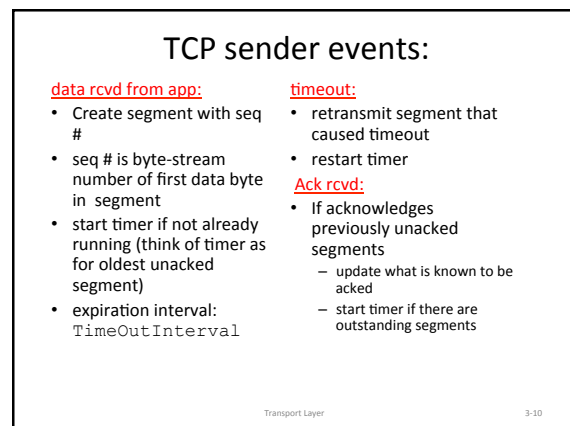
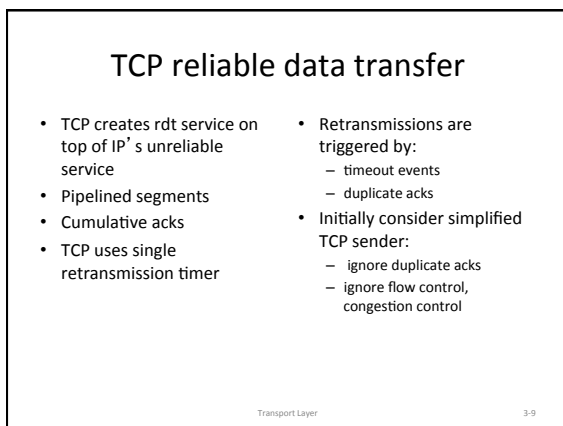
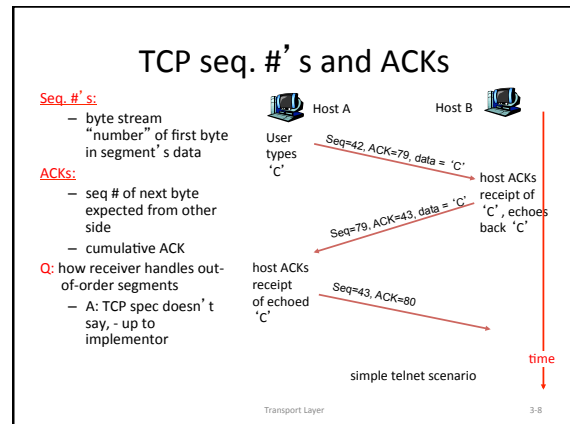
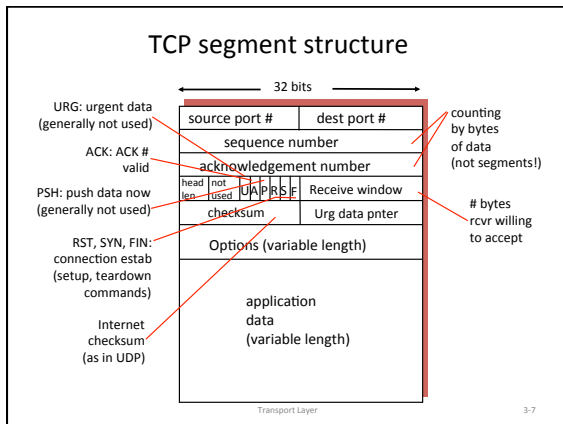
RFCs: 793, 1122, 1323, 2018, 2581

- 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



Transport Layer

3-6



TCP Round Trip Time and Timeout

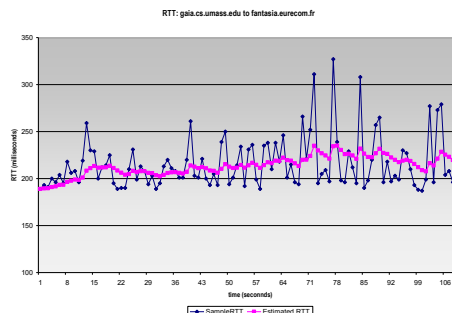
$$\text{EstimatedRTT} = (1 - \alpha) * \text{EstimatedRTT} + \alpha * \text{SampleRTT}$$

- Exponential weighted moving average
- influence of past sample decreases exponentially fast
- typical value: $\alpha = 0.125$

Transport Layer

3-13

Example RTT estimation:



Transport Layer

3-14

TCP Round Trip Time and Timeout

Setting the timeout

- **EstimatedRTT** plus "safety margin"
 - large variation in **EstimatedRTT** → larger safety margin
- first estimate of how much **SampleRTT** deviates from **EstimatedRTT**:

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

(typically, $\beta = 0.25$)

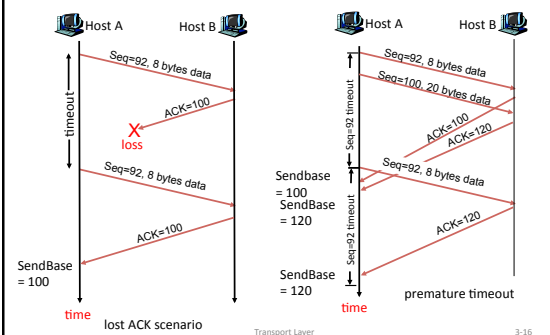
Then set timeout interval:

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

Transport Layer

3-15

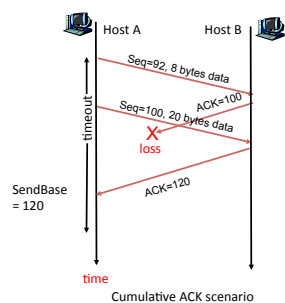
TCP: retransmission scenarios



Transport Layer

3-16

TCP retransmission scenarios (more)



Transport Layer

3-17

TCP ACK generation [RFC 1122, RFC 2581]

Event at Receiver

Arrival of in-order segment with expected seq #. All data up to expected seq # already ACKed

Arrival of in-order segment with expected seq #. One other segment has ACK pending

Arrival of out-of-order segment higher-than-expected seq. #. Gap detected

Arrival of segment that partially or completely fills gap

TCP Receiver action

Delayed ACK. Wait up to 500ms for next segment. If no next segment, send ACK

Immediately send single cumulative ACK, ACKing both in-order segments

Immediately send duplicate ACK, indicating seq. # of next expected byte

Immediate send ACK, provided that segment starts at lower end of gap

Transport Layer

3-18

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.
- If sender receives 3 ACKs for the same data, it supposes that segment after ACKed data was lost:
 - fast retransmit**: resend segment before timer expires

Transport Layer

3-19

Fast retransmit algorithm:

```

event: ACK received, with ACK field value of y
  if (y > SendBase) {
    SendBase = y
    if (there are currently not-yet-acknowledged segments)
      start timer
  }
  else {
    increment count of dup ACKs received for y
    if (count of dup ACKs received for y = 3) {
      resend segment with sequence number y
    }
  }

```

a duplicate ACK for
already ACKed segment

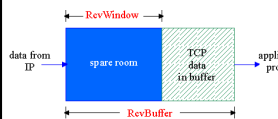
fast retransmit

Transport Layer

3-20

TCP Flow Control

- receive side of TCP connection has a receive buffer:



- app process may be slow at reading from buffer

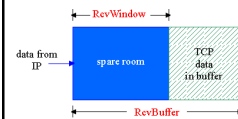
flow control
sender won't overflow receiver's buffer by transmitting too much, too fast

- speed-matching service: matching the send rate to the receiving app's drain rate

Transport Layer

3-21

TCP Flow control: how it works



(Suppose TCP receives out-of-order segments)

- spare room in buffer

= **RcvWindow**

= **RcvBuffer** - [**LastByteRcvd** - **LastByteRead**]

- Rcvr advertises spare room by including value of **RcvWindow** in segments
- Sender limits unACKed data to **RcvWindow**
- trap.rar guarantees receive buffer doesn't overflow

Transport Layer

3-22

Principles of Congestion Control

Congestion:

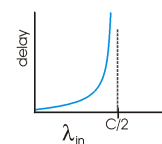
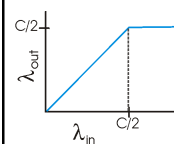
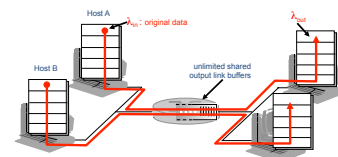
- informally: "too many sources sending too much data too fast for *network* to handle"
- different from flow control!
- manifestations:
 - lost packets (buffer overflow at routers)
 - long delays (queueing in router buffers)
- a top-10 problem!

Transport Layer

3-23

Causes/costs of congestion: scenario 1

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



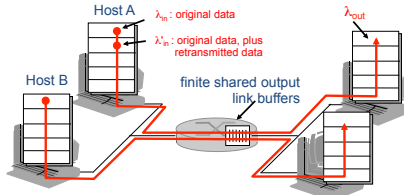
- large delays when congested
- maximum achievable throughput

Transport Layer

3-24

Causes/costs of congestion: scenario 2

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

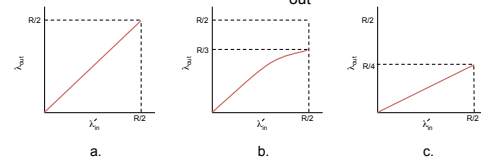


Transport Layer

3-25

Causes/costs of congestion: scenario 2

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



"costs" of congestion:

- more work (retrans) for given "goodput"
- unnecessary retransmissions: link carries multiple copies of pkt

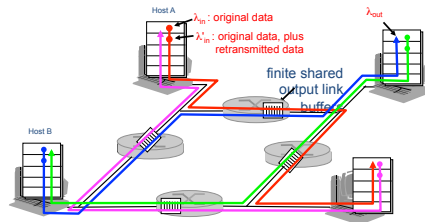
Transport Layer

3-26

Causes/costs of congestion: scenario 3

- four senders
- multihop paths
- timeout/retransmit

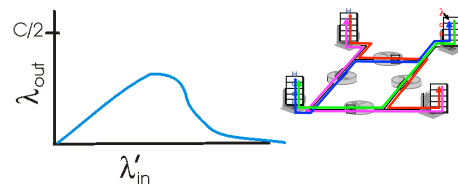
Q: what happens as λ_{in} increase?



Transport Layer

3-27

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

3-28