

# Communication Networks (0368-3030) / Spring 2011

The Blavatnik School of Computer Science,  
Tel-Aviv University

Allon Wagner

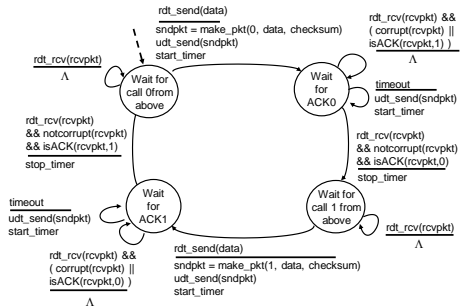


## Reliable Data Transfer

Kurose & Ross, Chapter 3.4 (5<sup>th</sup> ed.)

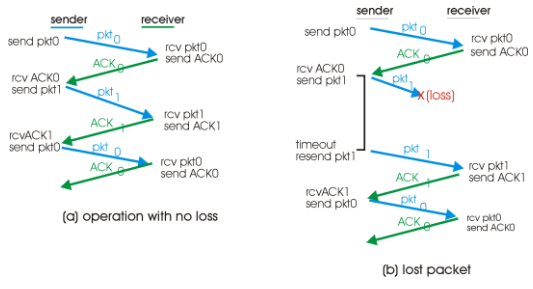
Slides adapted from:  
J. Kurose & K. Ross \ Computer Networking: A Top Down Approach (5<sup>th</sup> ed.) Addison-Wesley, April 2009.  
Copyright 1996-2010, J.F Kurose and K.W. Ross, All Rights Reserved.

### rdt3.0 sender



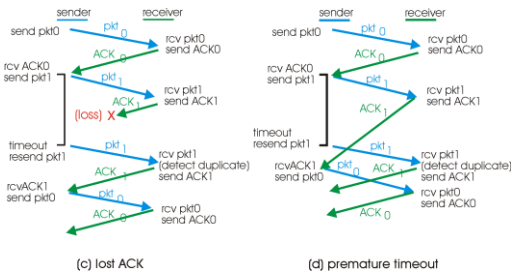
Transport Layer 3-3

### rdt3.0 in action



Transport Layer 3-4

### rdt3.0 in action



Transport Layer 3-5

### Performance of rdt3.0

- rdt3.0 works, but performance stinks
- ex: 1 Gbps link, 15 ms prop. delay, 8000 bit packet:

$$d_{trans} = \frac{L}{R} = \frac{8000\text{bits}}{10^9\text{bps}} = 8\text{microseconds}$$

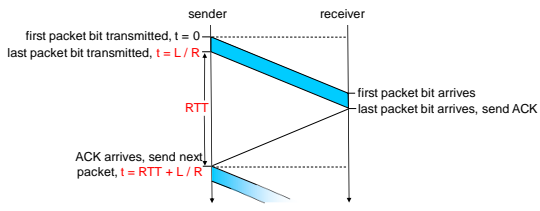
- U<sub>sender</sub>: utilization - fraction of time sender busy sending

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

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

Transport Layer 3-6

### rdt3.0: stop-and-wait operation



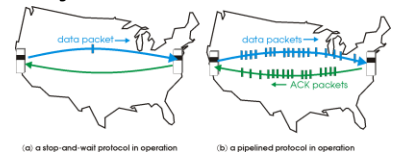
$$U_{\text{sender}} = \frac{L/R}{RTT + L/R} = \frac{.008}{30.008} = 0.00027$$

Transport Layer 3-7

### Pipelined protocols

**pipelining:** sender allows multiple, "in-flight", yet-to-be-acknowledged pkts

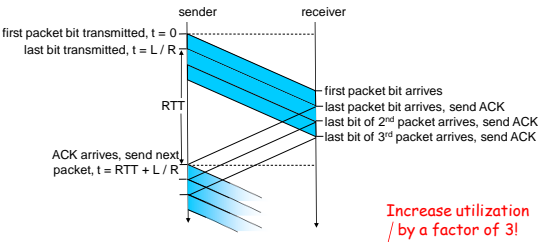
- range of sequence numbers must be increased
- buffering at sender and/or receiver



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

Transport Layer 3-8

### Pipelining: increased utilization



$$U_{\text{sender}} = \frac{3 * L/R}{RTT + L/R} = \frac{.024}{30.008} = 0.0008$$

Increase utilization by a factor of 3!

Transport Layer 3-9

### Pipelined Protocols

**Go-back-N: big picture:**

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

**Selective Repeat: big pic**

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

Transport Layer 3-10

### Go-Back-N

**Sender:**

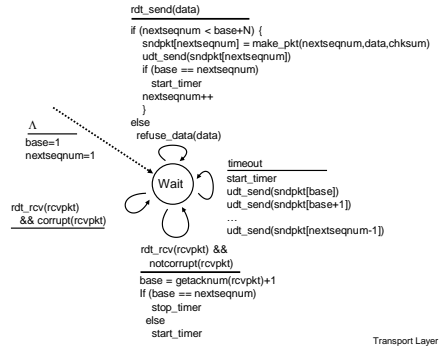
- k-bit seq # in pkt header
- "window" of up to N, consecutive unack'd pkts allowed



- ACK(n): ACKs all pkts up to, including seq # n - "cumulative ACK"
  - may receive duplicate ACKs (see receiver)
- timer for each in-flight pkt
- timeout(n): retransmit pkt n and all higher seq # pkts in window

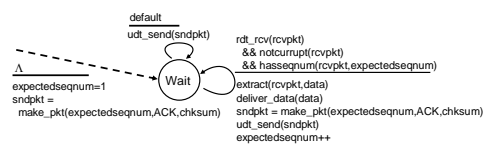
Transport Layer 3-11

### GBN: sender extended FSM



Transport Layer 3-12

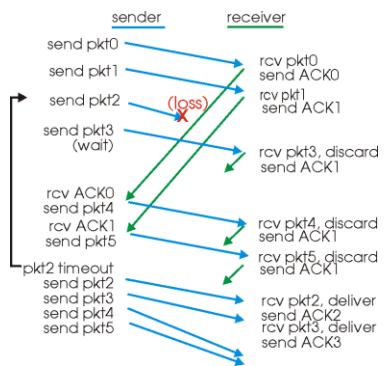
**GBN: receiver extended FSM**



- ACK-only:** always send ACK for correctly-received pkt with highest *in-order* seq #
- may generate duplicate ACKs
  - need only remember **expectedseqnum**
- ❖ **out-of-order pkt:**
- discard (don't buffer) -> **no receiver buffering!**
  - Re-ACK pkt with highest in-order seq #

Transport Layer 3-13

**GBN in action**



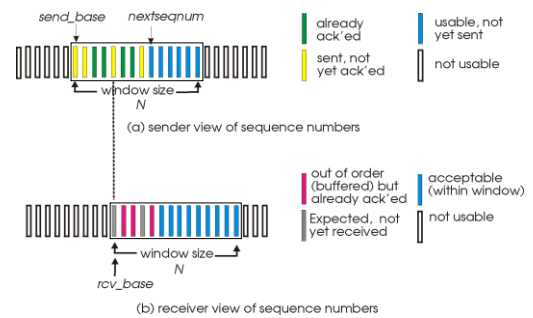
Transport Layer 3-14

**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, unACK'ed pkts

Transport Layer 3-15

**Selective repeat: sender, receiver windows**



Transport Layer 3-16

**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]

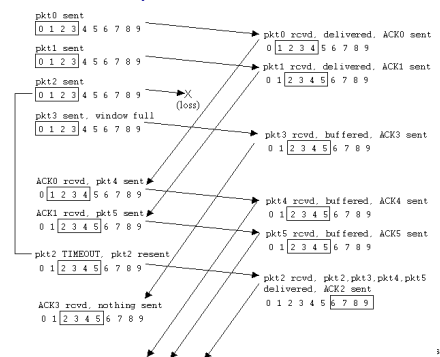
- ❖ ACK(n)

**otherwise:**

- ❖ ignore

Transport Layer 3-17

**Selective repeat in action**



Transport Layer 3-18

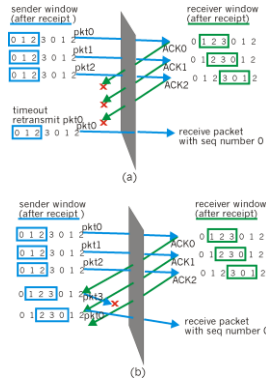
### Selective repeat: dilemma

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 3-19

### TCP Round Trip Time and 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 3-20

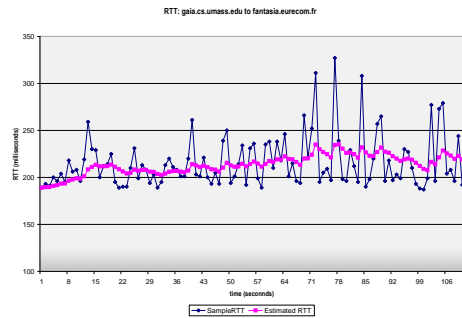
### TCP Round Trip Time and Timeout

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

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

Transport Layer 3-21

### Example RTT estimation:



Transport Layer 3-22

### 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 3-23

## Extra slides

if time permits

## 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 3-25

## TCP: Overview

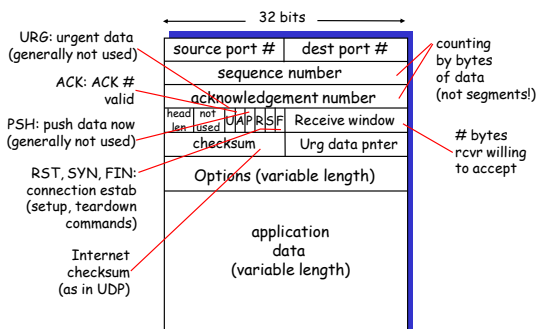
RFCs: 793, 1122, 1323, 2018, 2581

- ❖ **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) initiates sender, receiver state before data exchange
- ❖ **flow controlled:**
  - sender will not overwhelm receiver

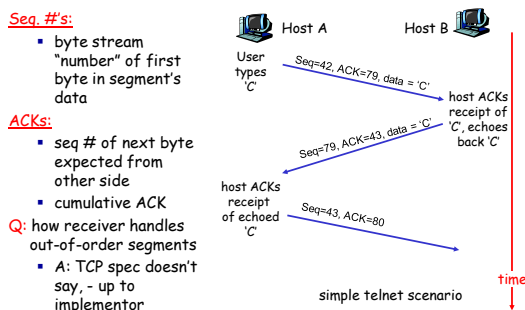


Transport Layer 3-26

## TCP segment structure



## TCP seq. #'s and ACKs



## TCP reliable data transfer

- ❖ TCP creates rdt service on top of IP's unreliable service
- ❖ pipelined segments
- ❖ cumulative acks
- ❖ TCP uses single retransmission timer
- ❖ retransmissions are triggered by:
  - timeout events
  - duplicate acks
- ❖ initially consider simplified TCP sender:
  - ignore duplicate acks
  - ignore flow control, congestion control

Transport Layer 3-29

## 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 acknowledges previously unacked segments
    - update what is known to be acked
    - start timer if there are outstanding segments

Transport Layer 3-30

```

NextSeqNum = InitialSeqNum
SendBase = InitialSeqNum

loop (forever) {
  switch(event)

  event: data received from application above
  create TCP segment with sequence number NextSeqNum
  if (timer currently not running)
    start timer
  pass segment to IP
  NextSeqNum = NextSeqNum + length(data)

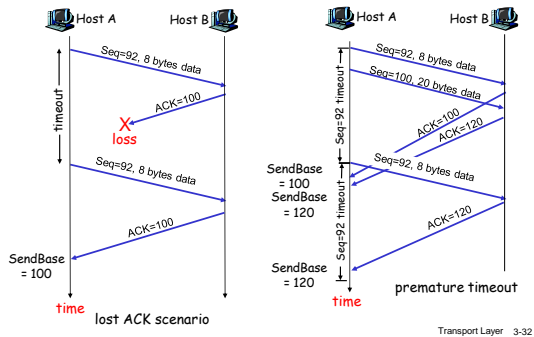
  event: timer timeout
  retransmit not-yet-acknowledged segment with
  smallest sequence number
  start timer

  event: ACK received, with ACK field value of y
  if (y > SendBase) {
    SendBase = y
    if (there are currently not-yet-acknowledged segments)
      start timer
  }
} /* end of loop forever */
    
```

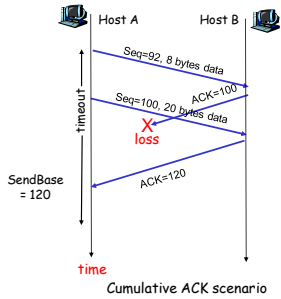
### TCP sender (simplified)

**Comment:**  
 • SendBase-1: last cumulatively acked byte  
**Example:**  
 • SendBase-1 = 71; y = 73, so the rcvr wants 73+ ; y > SendBase, so that new data is acked

### TCP: retransmission scenarios



### TCP retransmission scenarios (more)



### 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

### 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

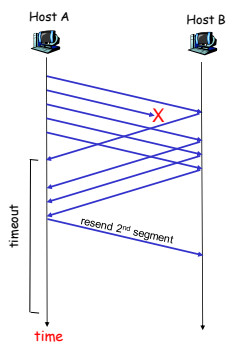


Figure 3.37 Resending a segment after triple duplicate ACK

## 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 3-37