

TCP

Transport Layer 3-1

TCP: Overview

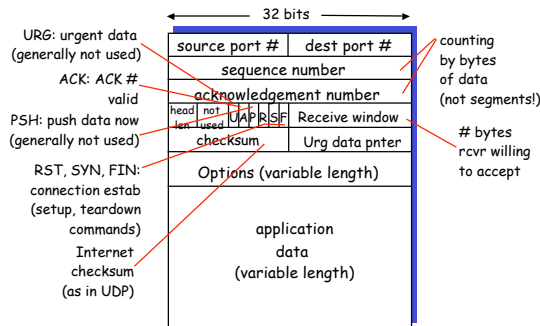
RFCs: 793, 1122, 1323, 2018, 2581

- **point-to-point:**
 - one sender, one receiver
- **reliable, in-order byte stream:**
 - no "message boundaries"
- **pipelined:**
 - TCP congestion and flow control set window size
- **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 3-2

TCP segment structure



Transport Layer 3-3

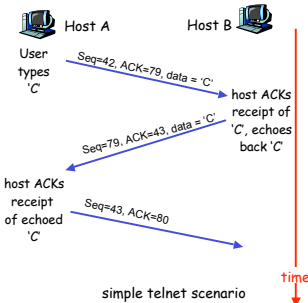
TCP seq. #'s and ACKs

Seq. #'s:

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

ACKs:

- seq # of next byte expected from other side
- cumulative ACK
- piggybacking
- how receiver handles out-of-order segments
- A: TCP spec doesn't say, - up to implementor



Transport Layer 3-4

TCP Round Trip Time and Timeout

Q: how to set TCP timeout value?

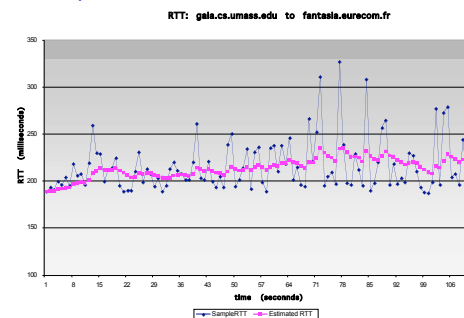
- longer than RTT
 - but RTT varies
- too short: premature timeout
 - unnecessary retransmissions
- too long: slow reaction to segment loss

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 3-5

Example RTT estimation:



Transport Layer 3-6

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 3-7

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 3-8

```
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
  }
```

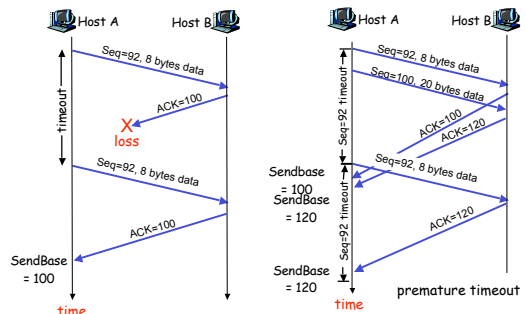
```
} /* end of loop forever */
```

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

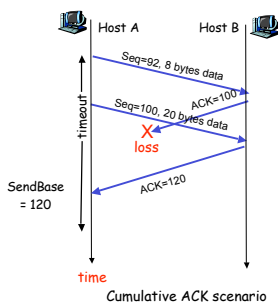
Transport Layer 3-9

TCP: retransmission scenarios



Transport Layer 3-10

TCP retransmission scenarios (more)



Transport Layer 3-11

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 3-12

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 3-13

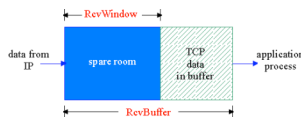
TCP Flow Control

- receive side of TCP connection has a receive 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
- app process may be slow at reading from buffer

Transport Layer 3-14

TCP Flow control: how it works



(Suppose TCP receiver discards 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 $RcvWindow$
 - guarantees receive buffer doesn't overflow

Transport Layer 3-15

TCP Connection Management

Recall: TCP sender, receiver establish "connection" before exchanging data segments

- initialize TCP variables:
 - seq. #s
 - buffers, flow control info (e.g. $RcvWindow$)
- **client**: connection initiator


```
Socket clientSocket = new Socket("hostname", "port number");
```
- **server**: contacted by client


```
Socket connectionSocket = welcomeSocket.accept();
```

Transport Layer 3-16

TCP Connection Management

Three way handshake:

Step 1: client host sends TCP SYN segment to server

- specifies initial seq #
- no data

Step 2: server host receives SYN, replies with SYNACK segment

- server allocates buffers
- specifies server initial seq. #

Step 3: client receives SYNACK, replies with ACK segment, which may contain data

Transport Layer 3-17

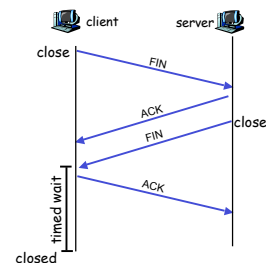
TCP Connection Management (cont.)

Closing a connection:

client closes socket:
`clientSocket.close();`

Step 1: client end system sends TCP FIN control segment to server

Step 2: server receives FIN, replies with ACK. Closes connection, sends FIN.



Transport Layer 3-18

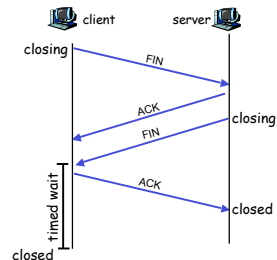
TCP Connection Management (cont.)

Step 3: client receives FIN, replies with ACK.

- Enters "timed wait" - will respond with ACK to received FINs

Step 4: server, receives ACK. Connection closed.

Note: with small modification, can handle simultaneous FINs.



Transport Layer 3-19