

## COMP 431

### Internet Services & Protocols

## The Transport Layer

Reliable data delivery & flow control in TCP

*Jasleen Kaur*

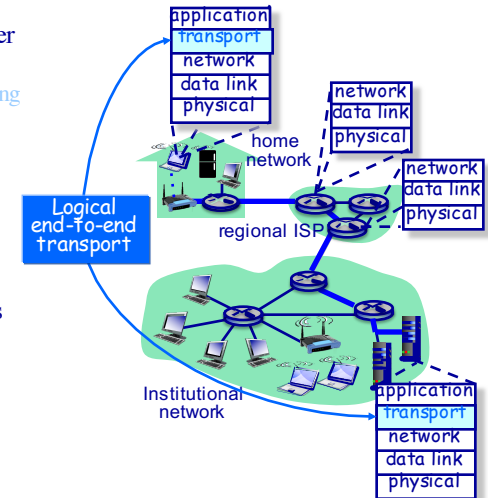
March 23, 2020

1

## Transport Layer Protocols & Services

### Outline

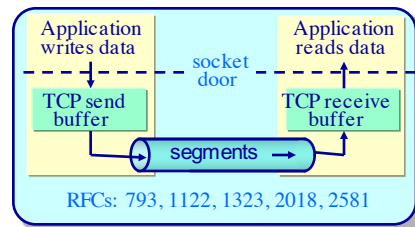
- ◆ Fundamental transport layer services
  - » Multiplexing/Demultiplexing
  - » Error detection
  - » Reliable data delivery
  - » Pipelining
  - » Flow control
  - » Congestion control
- ◆ Internet transport protocols
  - » UDP
  - » TCP



## TCP Overview

### TCP is...

- ◆ Point-to-point, full-duplex
  - » Bi-directional data flow within a connection
- ◆ Reliable, in-order *byte stream*
  - » No “message boundaries”
- ◆ Connection-oriented
  - » Handshaking initializes sender and receiver state before data exchange
- ◆ **Pipelined**
  - » Congestion and flow control determine window size
  - » Each endpoint has *two* buffers: a send and receive buffer
- ◆ Congestion controlled
  - » Internet would cease to function without this!
- ◆ Flow controlled
  - » Sender and receiver have synchronized windows to ensure receiver is not overwhelmed



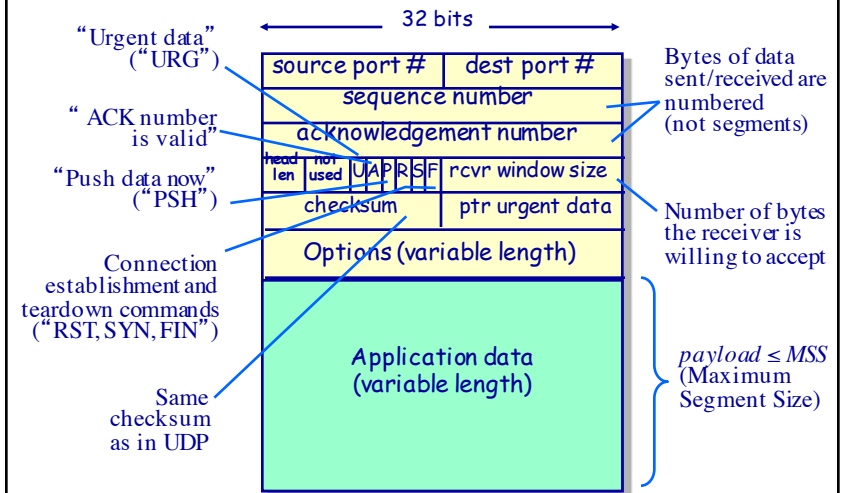
3

UDP: datagram oriented, connectionless, only between sender and receiver

Each endpoint in TCP is both a sender and a receiver

## TCP Segment Structure

### Header and payload format



4

Sequence number tells where in the byte stream that data belongs

U and P are not really used

R: reset, restart to terminate connection

S: synchronize, part of 3 way handshake

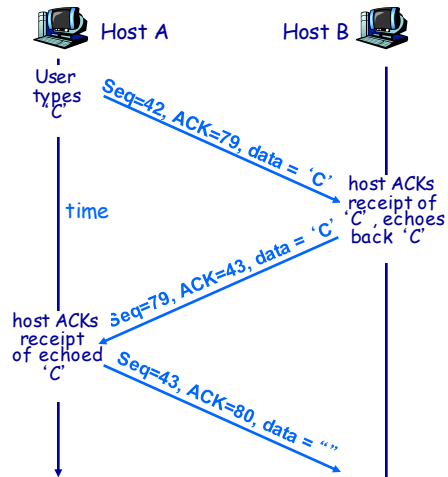
F: finish, use at the end of a connection

Window size plays a part in flow control

## TCP Sequence Numbers and ACKs

### Telnet example

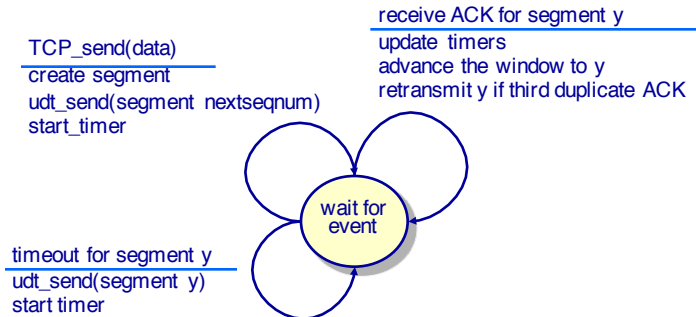
- ◆ Sequence numbers:
  - » Byte stream “index” of the first byte in the segment’s payload
- ◆ ACKs:
  - » Sequence number of next byte expected from the other side
  - » ACKs are cumulative
- ◆ How does receiver handle out-of-order segments?
  - » TCP spec doesn’t say, it’s up to the implementor



5

## Reliable Data Transfer in TCP

### Sender's state machine (simplified!)

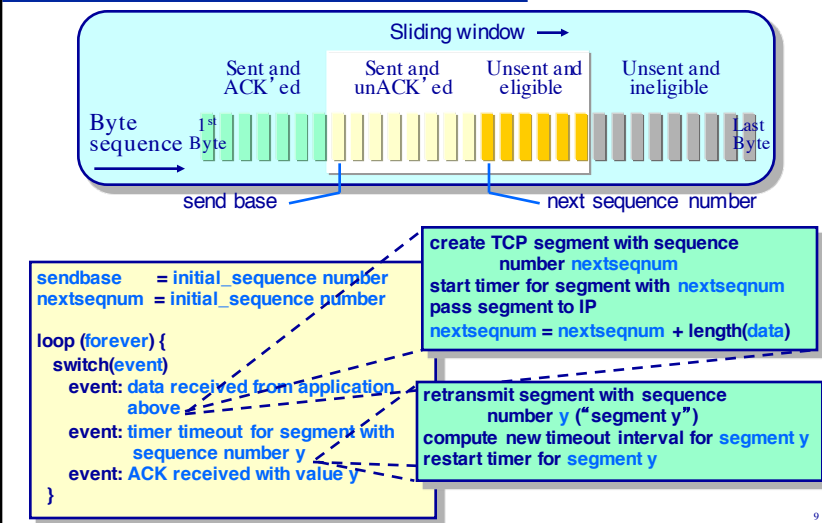


- ◆ TCP retransmits segments if:
  - » An expected ACK times out
  - » 3 duplicate ACKs for a segment are received

7

## Reliable Data Transfer in TCP

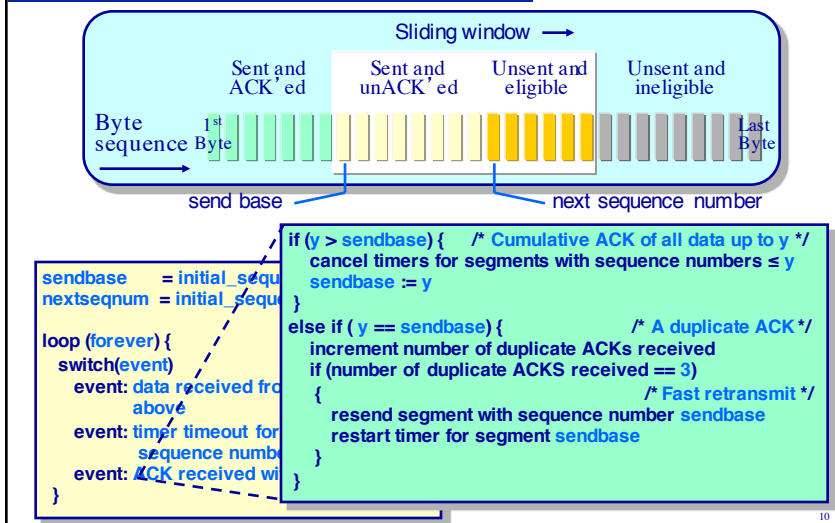
### Simplified sender's state machine



MSS: maximum segment size

## Reliable Data Transfer in TCP

### Simplified sender's state machine



## Reliable Data Transfer in TCP

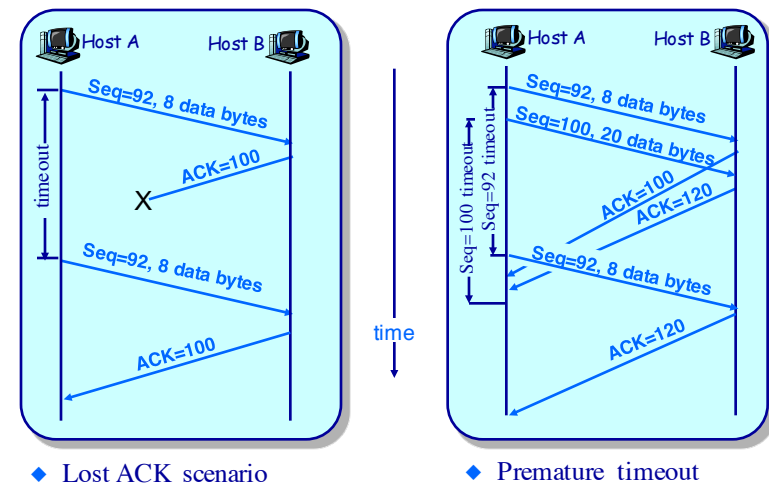
### Receiver ACK generation rules [RFC 1122, RFC 2581]

Event	TCP Receiver action
In-order segment arrival, no gaps, all previous data already ACKed	Delayed ACK. Wait 200 <i>ms</i> (up to 500 <i>ms</i> allowed) for next segment. If no segment received, send ACK
In-order segment arrival, no gaps, one delayed ACK pending	Immediately send single cumulative ACK
Out-of-order segment arrival (higher than expected sequence number) — Gap detected	Send duplicate ACK, indicating sequence number of next expected byte
Arrival of segment that partially or completely fills gap	Immediate ACK if segment starts at lower end of gap

11

## Reliable Data Transfer in TCP

### Retransmission examples



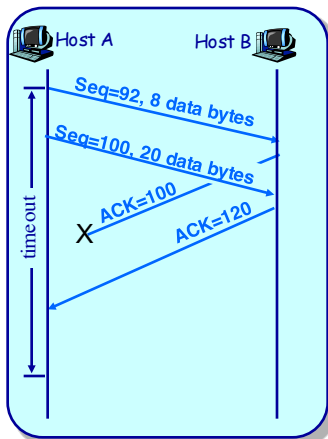
◆ Lost ACK scenario

◆ Premature timeout

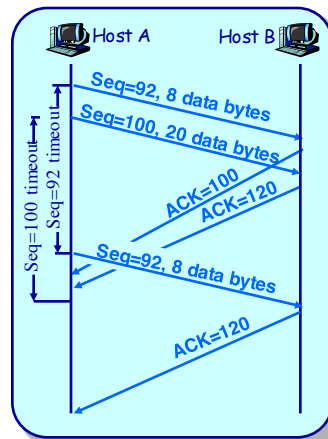
12

## Reliable Data Transfer in TCP

### Retransmission examples



- ◆ Cumulative ACKs potentially avoid retransmissions



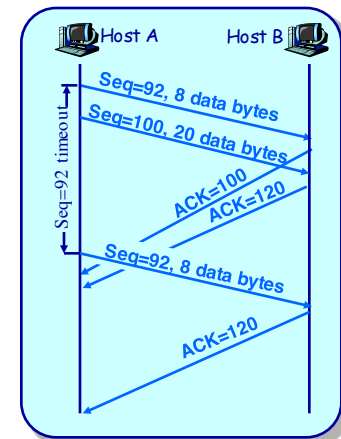
- ◆ Premature timeout

13

## Reliable Data Transfer in TCP

### Setting the ACK timer

- ◆ How large should the ACK timeout value be?
  - » Too short: Premature timeouts result in unnecessary retransmissions
  - » Too long: Slow reaction to loss results in poor performance because the sender's windows stops advancing
- ◆ Timer should be longer than the RTT, but how do we estimate RTT?

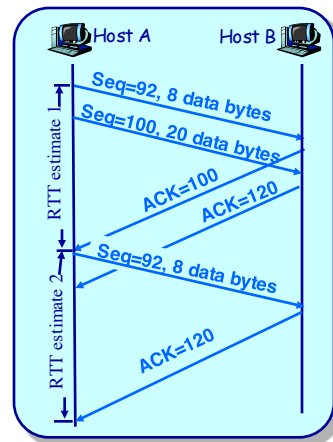


14

## Reliable Data Transfer in TCP

### Setting the ACK timer

- ◆ Measure the time from segment transmission until receipt of ACK (“**SampleRTT**”)
  - » Ignore retransmissions
  - » Measure only one segment’s RTT at a time
- ◆ **SampleRTT** will vary, so we compute an average RTT based on several recent RTT samples



15

## Reliable Data Transfer in TCP

### Estimating round-trip-time

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

$$\text{Timeout} = \text{EstimatedRTT} + 4 * \text{Deviation}$$

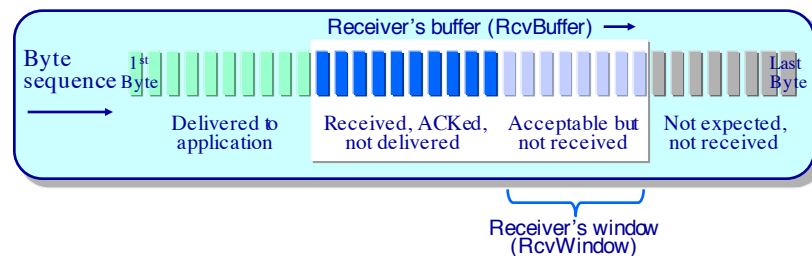
$$\text{Deviation} = (1-x) * \text{Deviation} + x * |\text{SampleRTT} - \text{EstimatedRTT}|$$

- ◆ The estimated RTT is an exponential weighted moving average (EWMA)
  - » Computes a “smooth” average
  - » Influence of a given sample decreases exponentially fast
 
$$E_n = x * S_n + x(1-x)S_{n-1} + x(1-x)^2S_{n-2} + \dots + x(1-x)^iS_{n-i} + \dots$$
  - » Typical value of  $x$  is 0.125
- ◆ Timeout is **EstimatedRTT** plus “safety margin”
- ◆ Large variation in **EstimatedRTT** results in a larger safety margin

16

## TCP Flow Control

### Receiver Window control



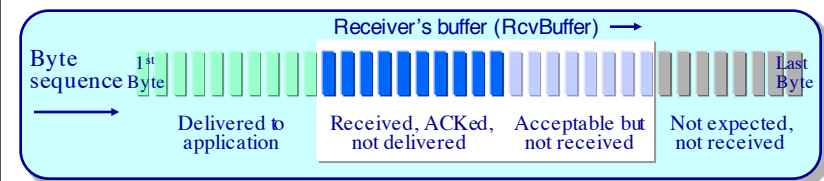
- ◆ Flow control is the problem of ensuring the receiver is not overwhelmed
  - » The receiver can become overwhelmed if the application reads too slow or the sender transmits too fast
- ◆ The receiver's window represents its remaining buffer capacity
- ◆ The window advances as the application reads received data



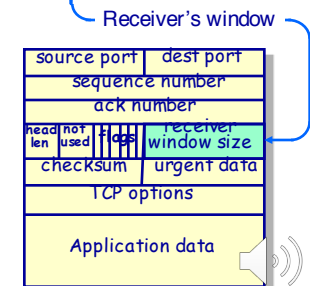
19

## TCP Flow Control

### Receiver window control



- ◆ The receiver explicitly informs the sender of the amount of free buffer space in RcvBuffer
  - » RcvWindow field in TCP segment
- ◆ The sender keeps the amount of transmitted, unACKed data less than most recently received RcvWindow

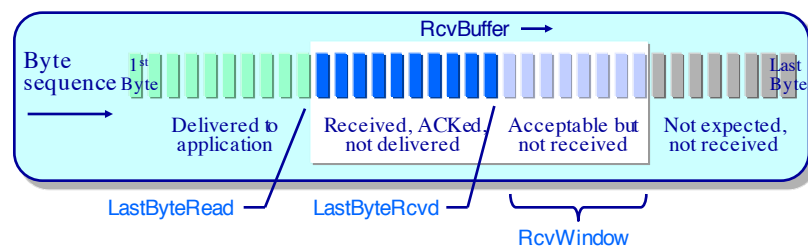


20



## TCP Flow Control

### Receiver window control



- ◆ The goal is to ensure:

$$\text{LastByteRcvd} - \text{LastByteRead} \leq \text{RcvBuffer}$$

- ◆ Sender is sent:

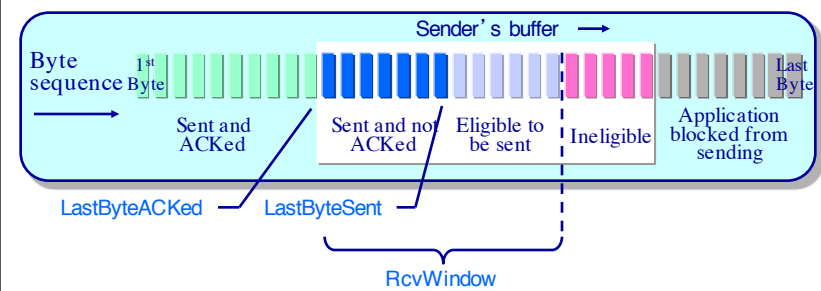
$$\text{RcvWindow} = \text{RcvBuffer} - (\text{LastByteRcvd} - \text{LastByteRead})$$



21

## TCP Flow Control

### Sender Window control



- ◆ The sender ensures:

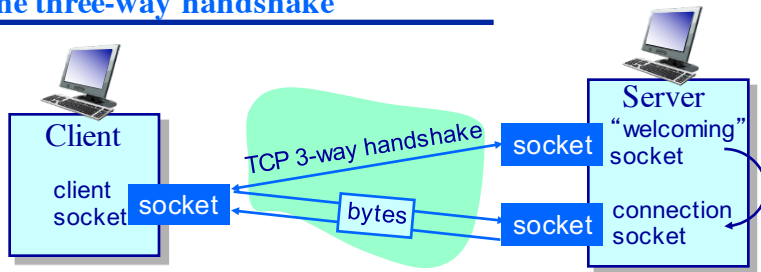
$$\text{LastByteSent} - \text{LastByteACKed} \leq \text{RcvWindow}$$



22

## TCP Connection Management

### The three-way handshake



- ◆ TCP endpoints establish a “connection” before exchanging data segments
  - » *client*: connection initiator
 

```
clientSocket = socket(AFNET, SOCK_STREAM)
clientSocket.connect(serverName, serverPort)
```
  - » *server*: contacted by client
 

```
connectionSocket, addr = serverSocket.accept()
```



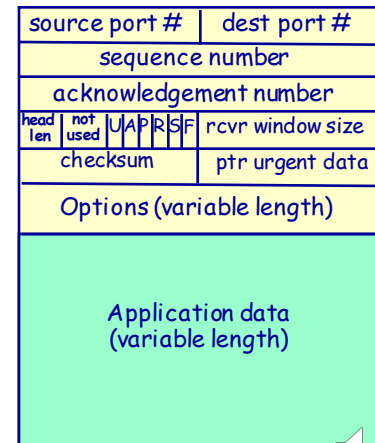
23

Client reaches out first

## TCP Connection Management

### The three-way handshake

- ◆ Client sends SYN segment to server
  - » The SYN specifies the client's initial sequence number
  - » The ACK number in the SYN will be 0

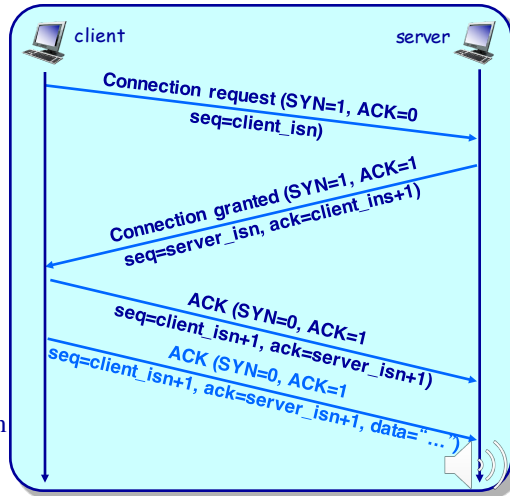


24

## TCP Connection Management

### The three-way handshake

- ◆ Client sends SYN segment to server
  - » The SYN specifies the client's initial sequence number
- ◆ Server receives SYN, replies with SYN+ACK segment
  - » ACKs received SYN
  - » Allocates buffers
  - » Specifies server's initial sequence number
- ◆ Third segment may be an ACK only or an ACK+data

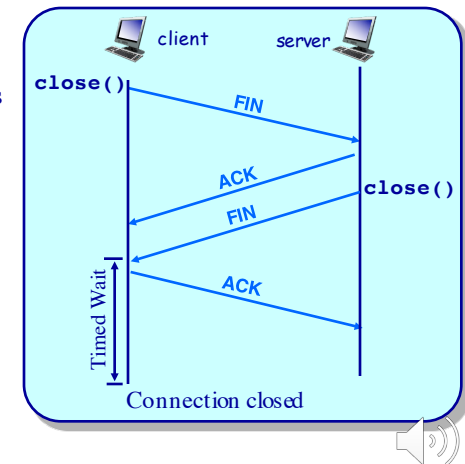


25

## TCP Connection Management

### Closing a connection

- ◆ Client sends FIN segment to server
- ◆ Server receives FIN, replies with ACK
  - » Server closes connection, sends FIN
- ◆ Client receives FIN, replies with ACK
- ◆ Client enters "timed wait" state
  - » Client will ACK any received FIN



26

# TCP Connection Management

## Client/Server lifecycles

