

The Transport Layer

Reliable data delivery & flow control in TCP

Dr. Michele Weigle

Department of Computer Science

Old Dominion University

mweigle@cs.odu.edu

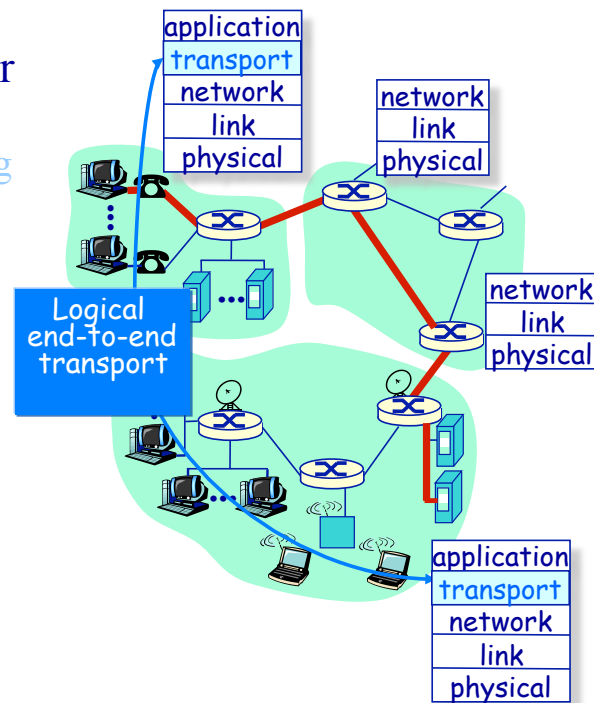
<http://www.cs.odu.edu/~mweigle/CS455-S13>

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

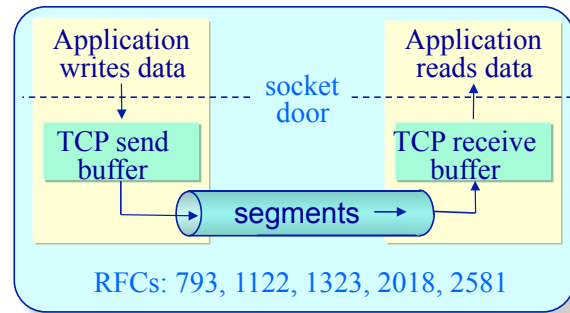


2

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

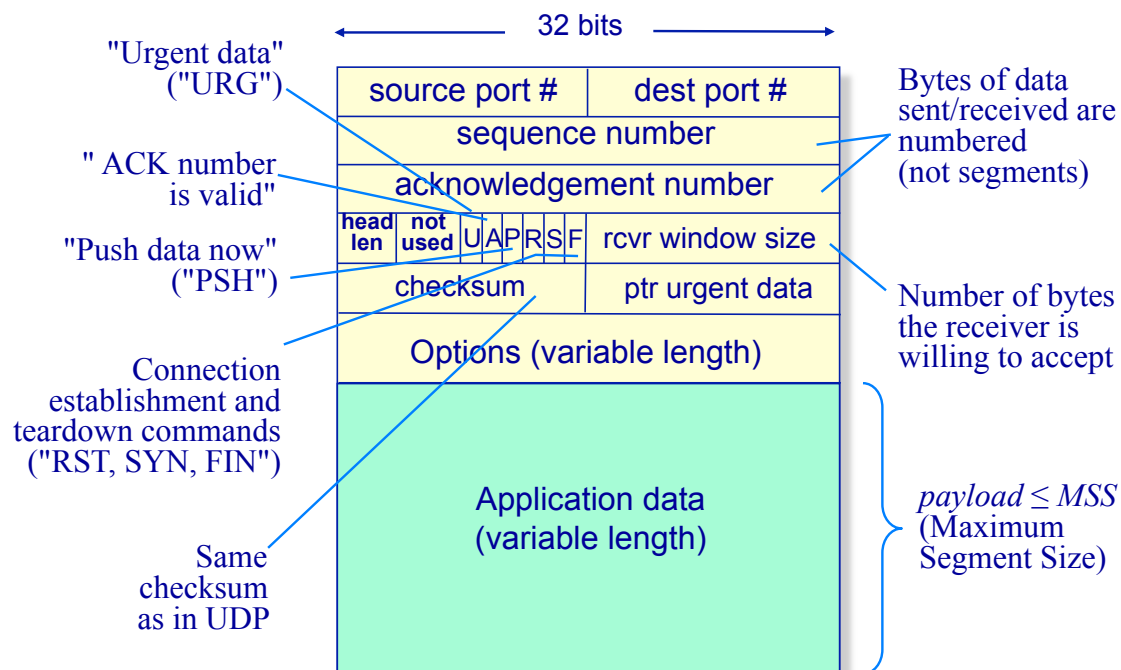


- ◆ Congestion controlled
 - » Internet would cease to function without this!
- ◆ Flow controlled
 - » Sender and receiver have synchronized windows to ensure receiver is not overwhelmed

3

TCP Segment Structure

Header and payload format

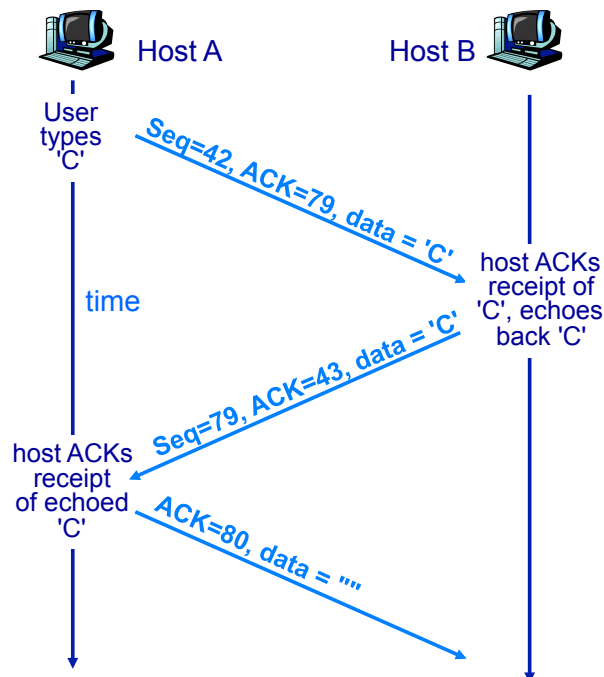


4

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

TCP Intro

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?
 - » Measure the time from segment transmission until receipt of ACK ("SampleRTT")
 - ❖ Ignore retransmissions - Karn's algorithm
 - ❖ Measure only one segment's RTT at a time
 - » SampleRTT will vary, so we compute an average RTT based on several recent RTT samples

6

TCP Intro

Estimating round-trip-time

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

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

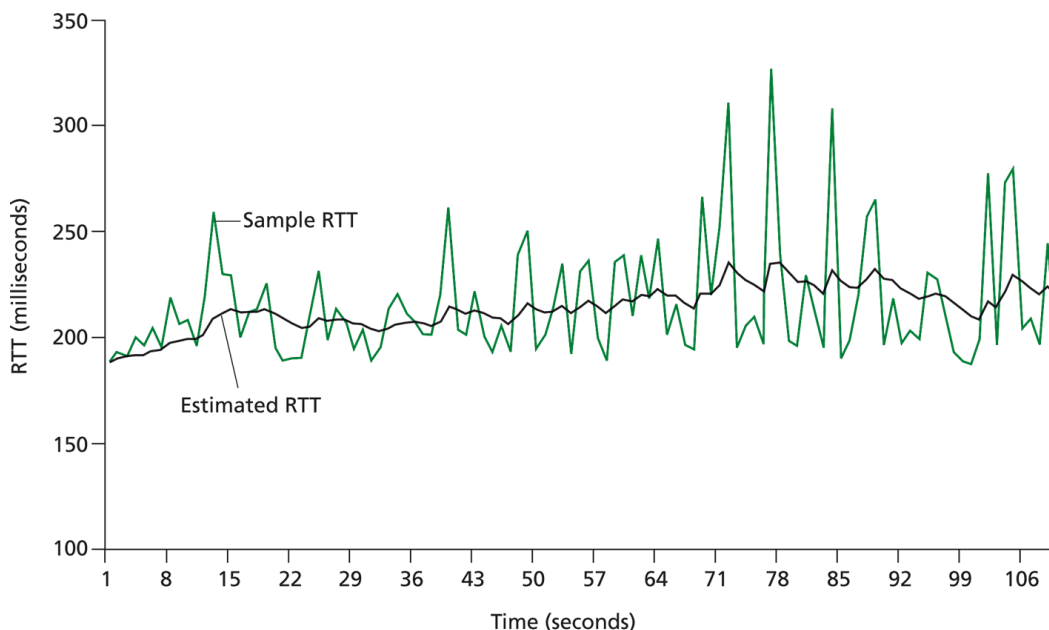
$$\text{Deviation} = (1-\beta) * \text{Deviation} + \beta * |\text{SampleRTT} - \text{EstimatedRTT}|$$

- ◆ The estimated RTT is an exponential weighted moving average (EWMA)
 - » Computes a "smooth" average
 - » Influence of a given sample decreases exponentially fast
 - » Typical value of α is 0.125
 - » Typical value of β is 0.25
- ◆ Timeout is EstimatedRTT plus "safety margin"
- ◆ Large variation in EstimatedRTT results in a larger safety margin

7

TCP Intro

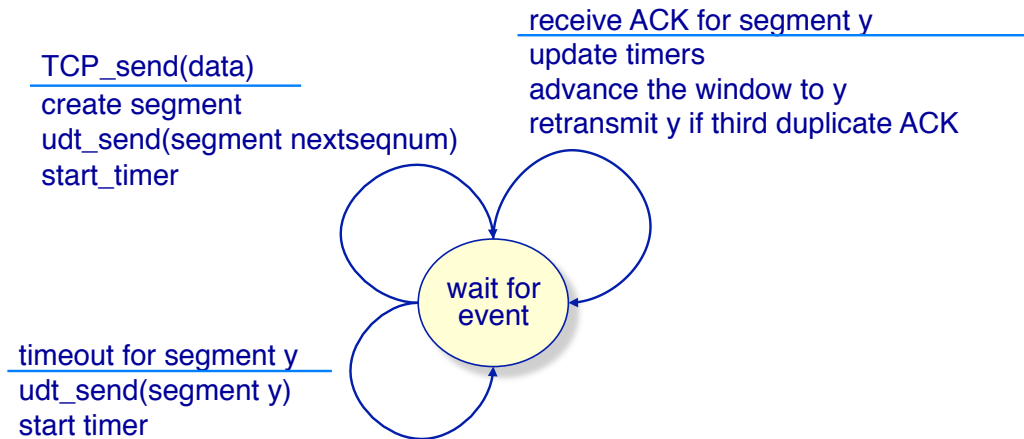
Estimating round-trip-time



8

Reliable Data Transfer in TCP

Sender's state machine

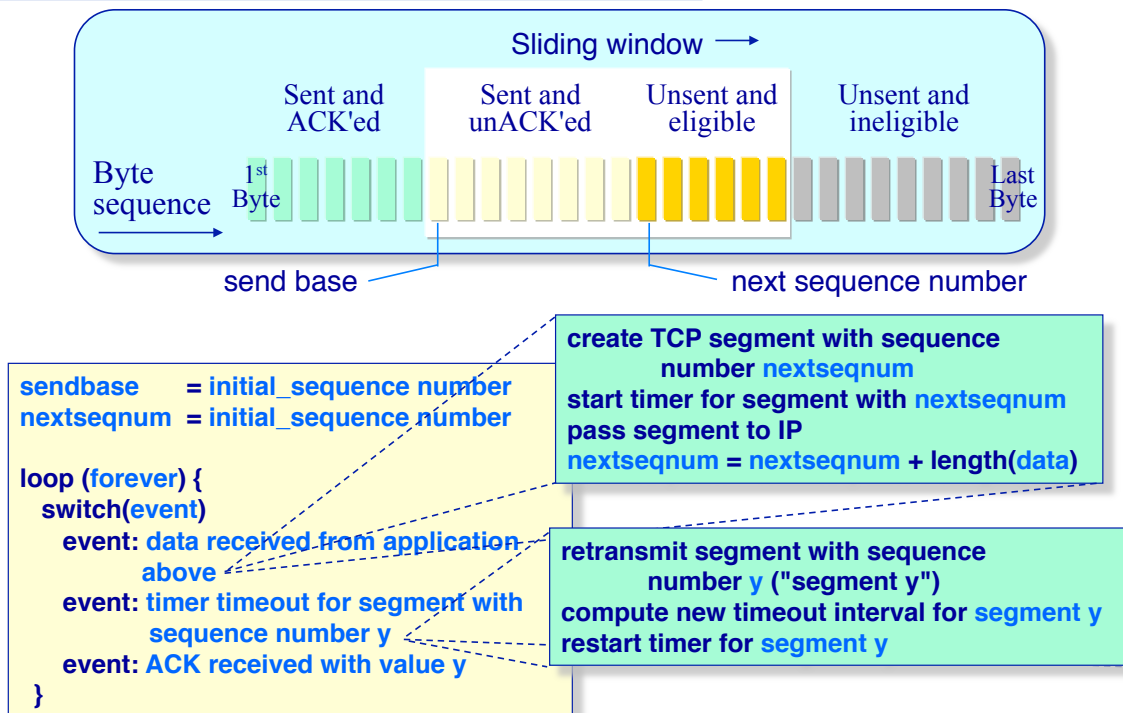


- ◆ TCP retransmits segments if:
 - » An expected ACK times out
 - » 3 duplicate ACKs for a segment are received

9

Reliable Data Transfer in TCP

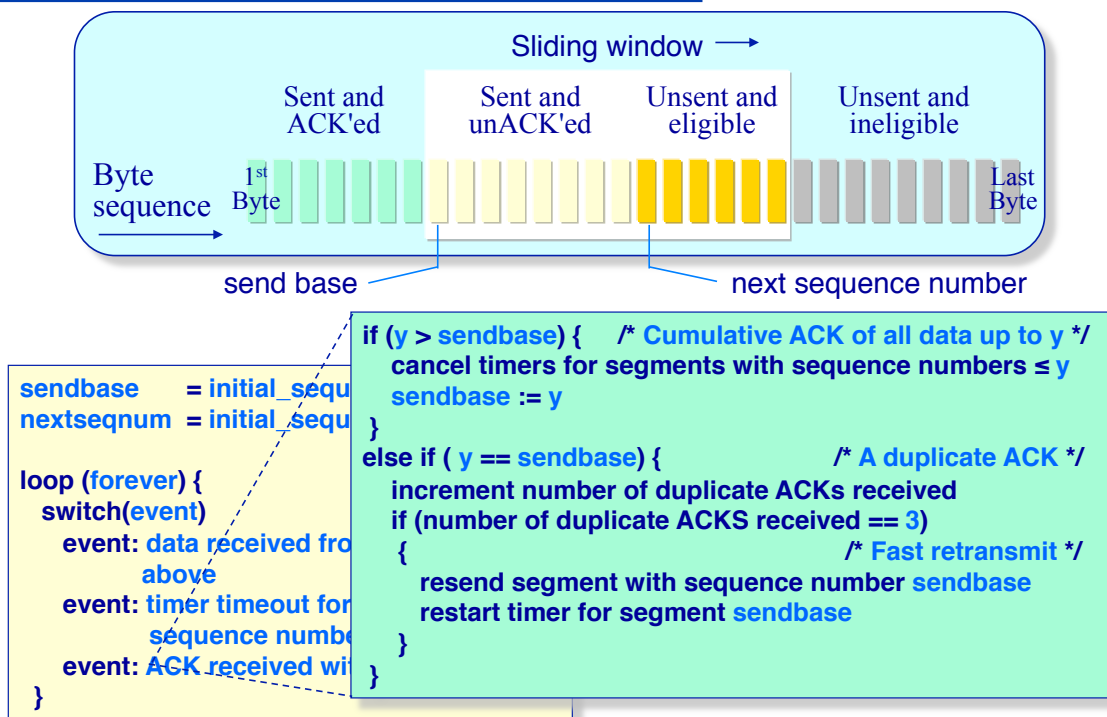
Simplified sender's state machine



10

Reliable Data Transfer in TCP

Simplified sender's state machine



11

Reliable Data Transfer in TCP

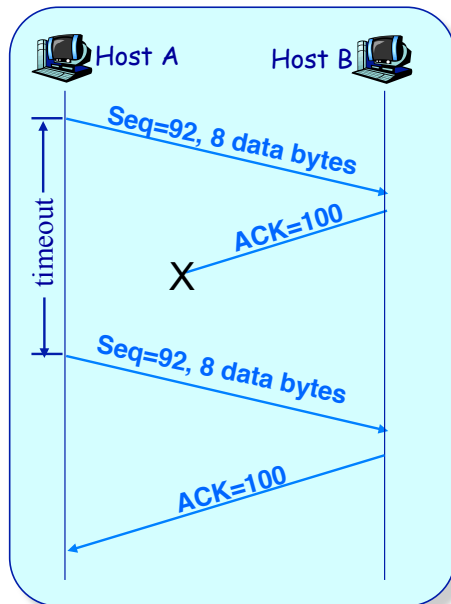
ACK generation rules [RFC 1122, RFC 2581]

Event	TCP Receiver action
In-order segment arrival, no gaps, all previous data already ACKed	Delayed ACK. Wait 200 ms (up to 500 ms 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

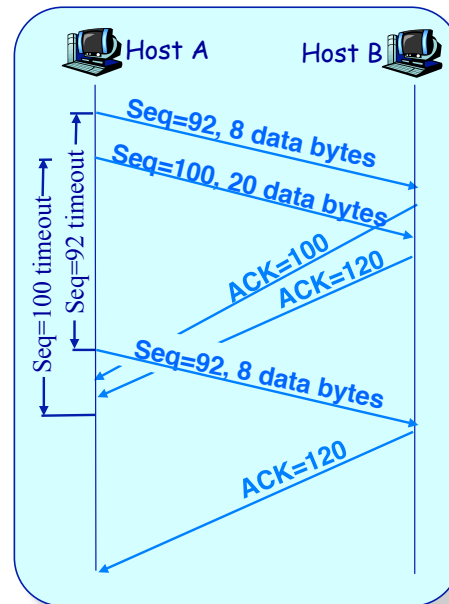
12

Reliable Data Transfer in TCP

Retransmission examples



◆ Lost ACK scenario

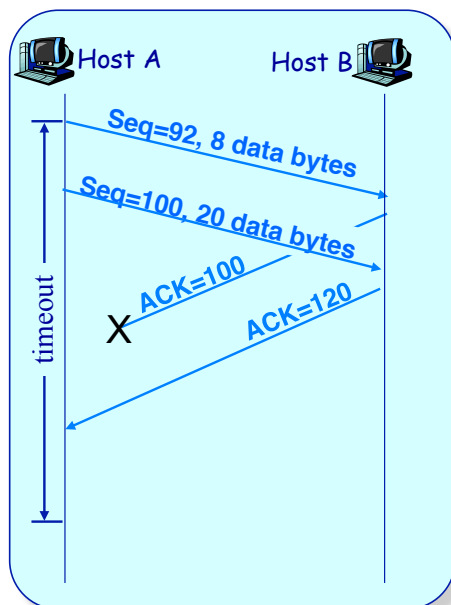


◆ Premature timeout

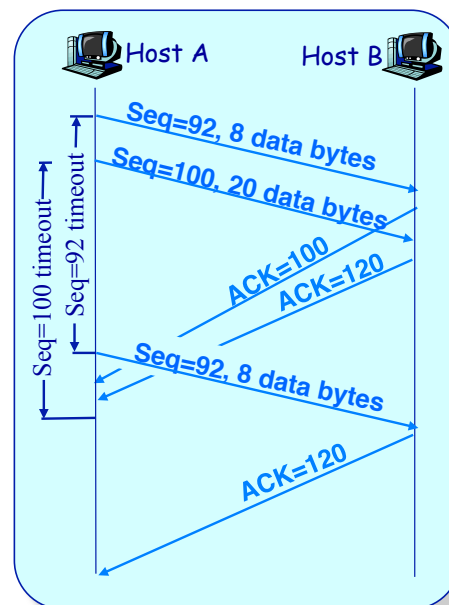
13

Reliable Data Transfer in TCP

Retransmission examples



◆ Cumulative ACKs potentially avoid retransmissions



◆ Premature timeout

14

Reliable Data Transfer in TCP

Timeout Interval

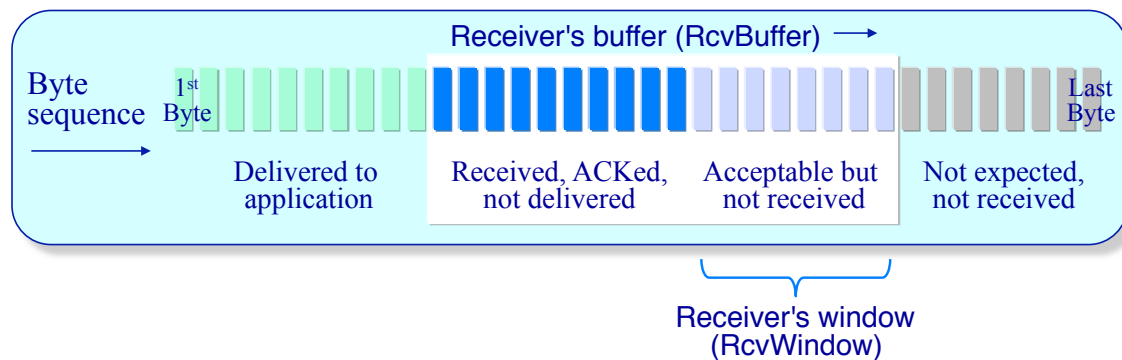
Whenever a timeout occurs

- » TCP retransmits the non-yet-acknowledged segment with the smallest sequence number
- » Sets the timeout interval to twice the previous value
 - ❖ timeout interval grows *exponentially* after every consecutive retransmission

15

TCP Flow Control

Window control

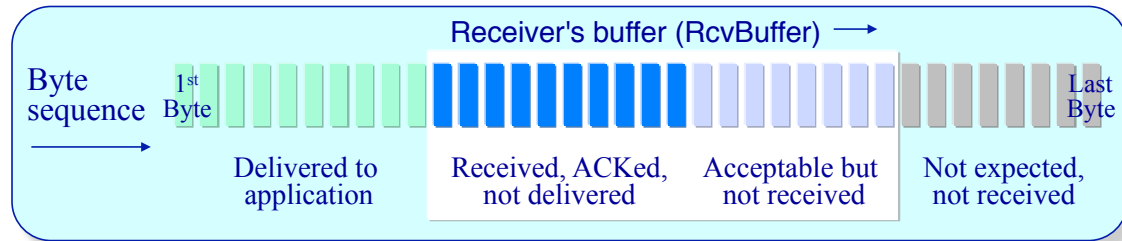


- ◆ Flow control is the problem of ensuring the receiver is not overwhelmed
 - » The receiver can become overwhelmed if the application reads too slow or the sender transmits too fast
- ◆ The receiver's window represents its remaining buffer capacity
- ◆ The window advances as the application reads received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

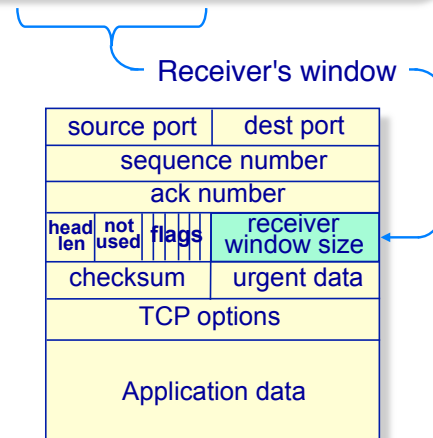
16

TCP Flow Control

Window control



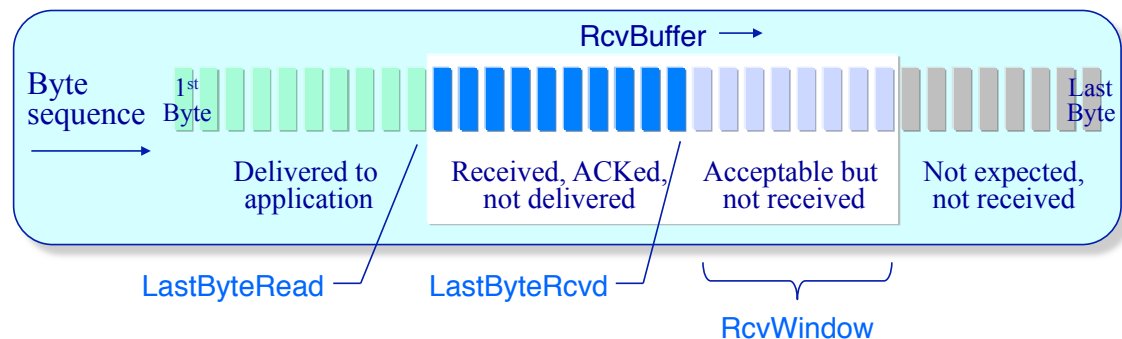
- ◆ The receiver explicitly informs the sender of the amount of free buffer space in RcvBuffer
 - » RcvWindow field in TCP segment
- ◆ The sender keeps the amount of transmitted, unACKed data less than most recently received RcvWindow



17

TCP Flow Control

Window control



- ◆ The goal is to ensure:

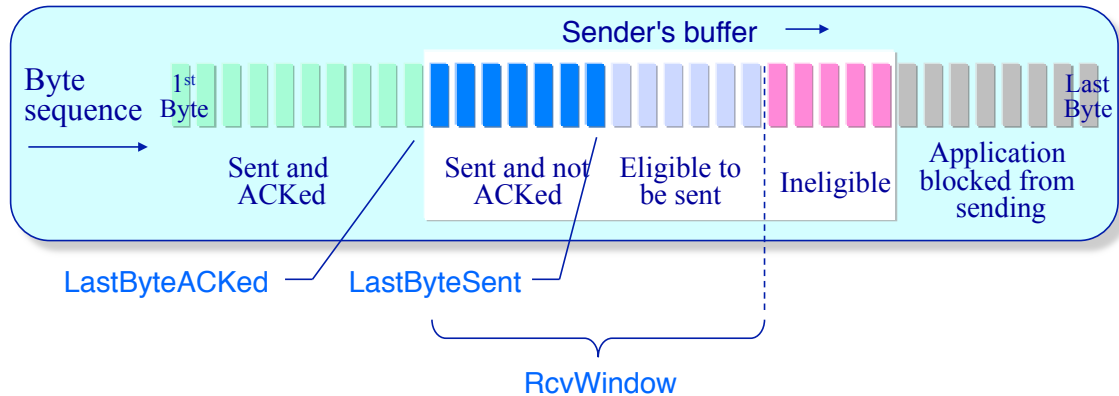
$$\text{LastByteRcvd} - \text{LastByteRead} \leq \text{RcvBuffer}$$
- ◆ Sender is sent:

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

18

TCP Flow Control

Window control



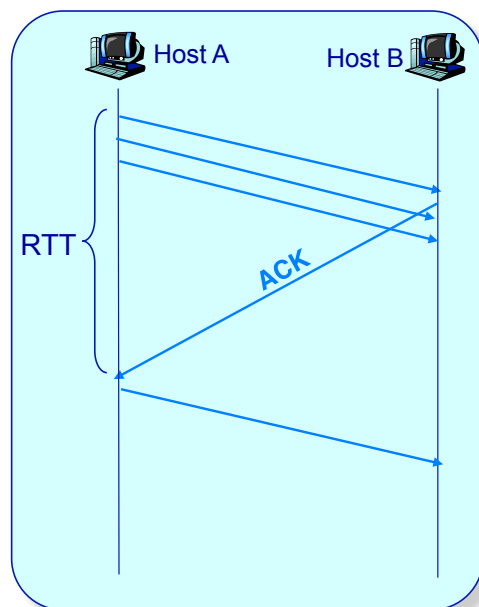
- ◆ The sender ensures:

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

19

TCP Flow Control

Example: How big should receiver window be?



TCP's bandwidth-delay product (BDP):
bottleneck bandwidth x RTT

10 Mbps bottleneck bandwidth
RTT = 20 ms

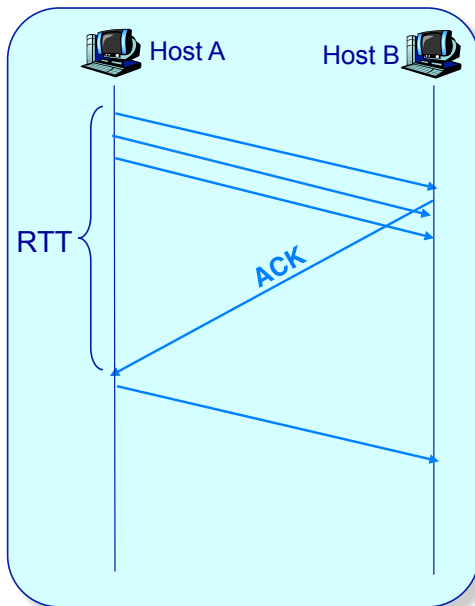
$$\text{BDP} = 10 \text{ Mbps} \times 20 \text{ ms} = \mathbf{25,000 \text{ B}}$$

So, receiver's window size (W)
should be at least 25,000 bytes.

20

TCP Flow Control

Example: What's max TCP throughput?



$$R \times RTT = W$$

2 KB (2048 B) receive window
20 ms RTT

$$R = W / RTT$$

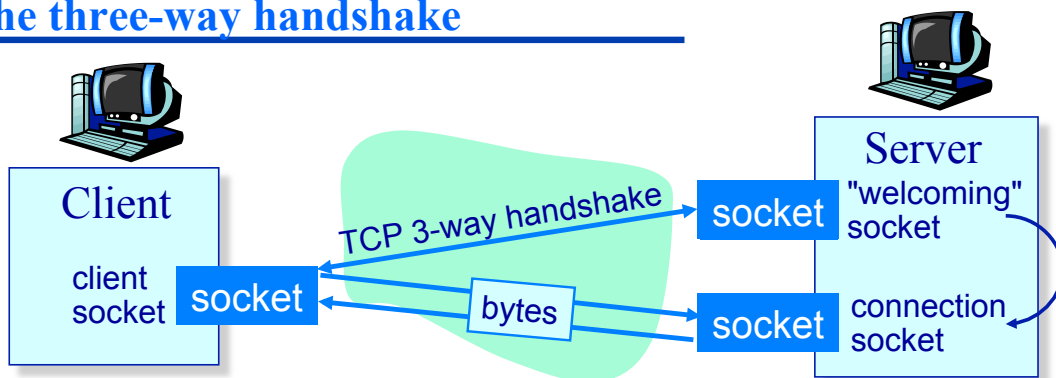
$$\begin{aligned} R \text{ bps} &= 2048 \text{ B} / 20 \text{ ms} \\ &= 16,384 \text{ b} / 0.02 \text{ seconds} \\ &= \mathbf{819,200 \text{ bps}} \end{aligned}$$

So, max TCP throughput is
819.2 kbps

21

TCP Connection Management

The three-way handshake



- ◆ TCP endpoints establish a "connection" before exchanging data segments

» *client*: connection initiator

```
mySocket = new Socket (hostname, port);
```

» *server*: contacted by client

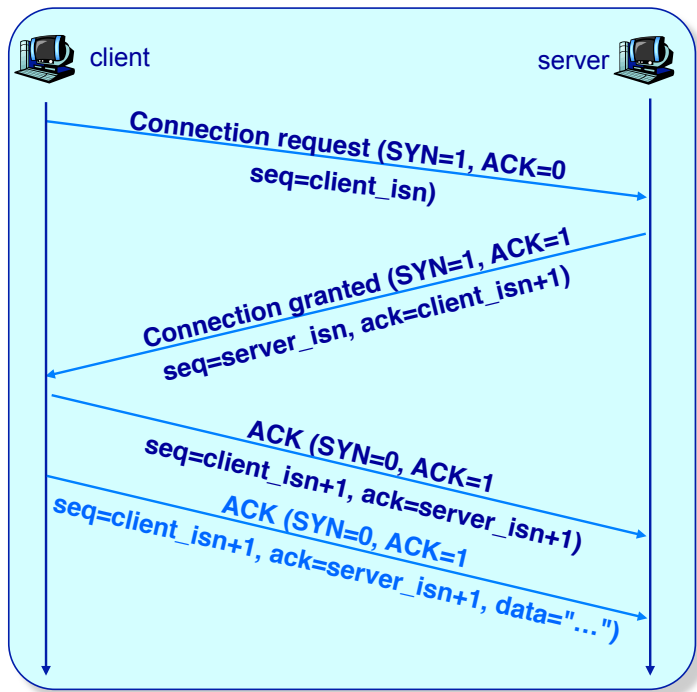
```
clientSocket = welcomeSocket.accept();
```

22

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



23

Reliable Data Transfer in TCP

Timeout Interval

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

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

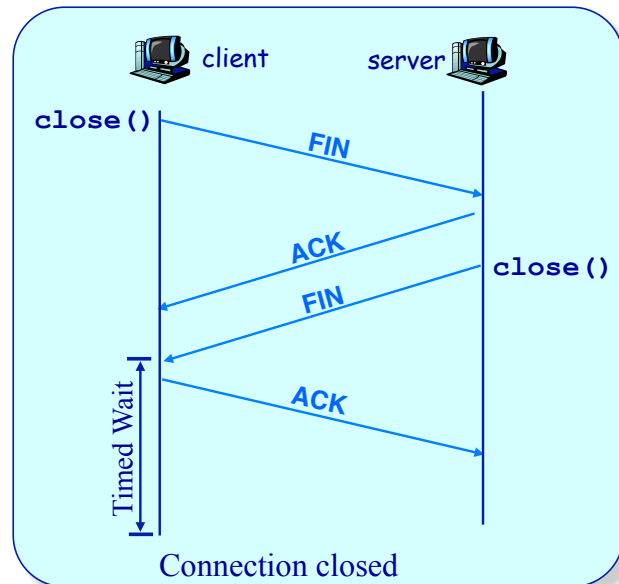
- ◆ What happens if the SYN or SYN/ACK is lost?
 - » Initial value of `Timeout` is 3 seconds
- ◆ What happens if the retransmission of the SYN is lost?
 - » Exponential increase

24

TCP Connection Management

Closing a connection

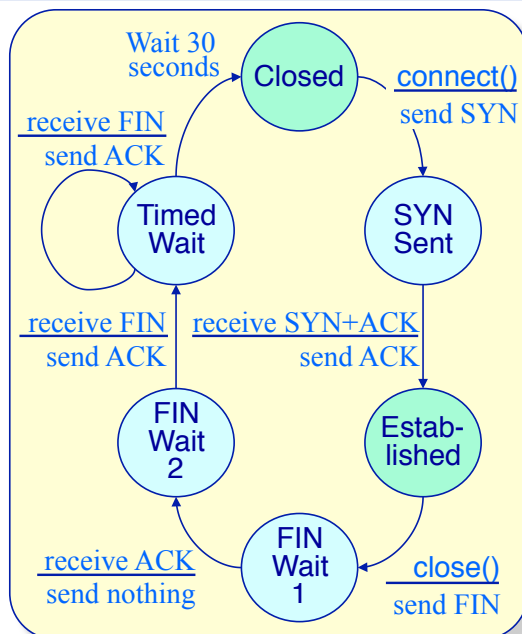
- ◆ Client sends FIN segment to server
- ◆ Server receives FIN, replies with ACK
 - » Server closes connection, sends FIN
- ◆ Client receives FIN, replies with ACK
- ◆ Client enters "timed wait" state
 - » Client will ACK any received FIN



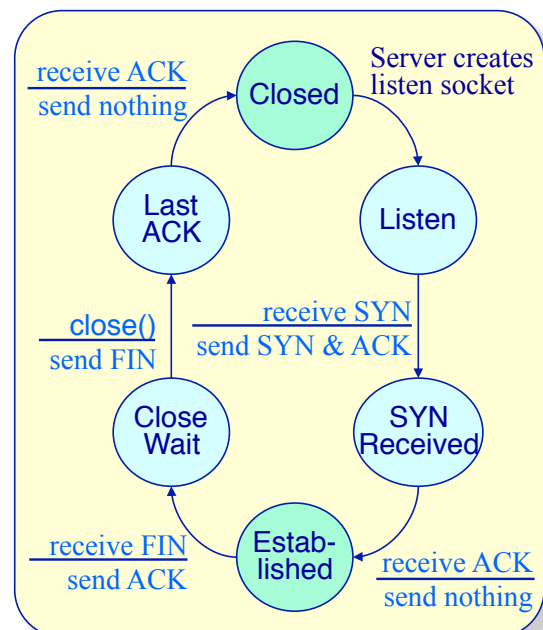
25

TCP Connection Management

Client/Server lifecycles



- ◆ TCP client lifecycle



- ◆ TCP server lifecycle

26

TCP Connection Management

SYN Flood Attack

- ◆ Setup: TCP Connection Setup
 - » TCP allocates and initializes connection variables and buffers in response to a received SYN
 - » If TCP never receives 3rd part of handshake, it will deallocate buffers (after a minute or more)
- ◆ Attack: SYN Flood
 - » Bad guy sends a large number of SYN segments, but never completes the handshake
 - » Once resources are exhausted, legitimate connections are refused
- ◆ Remedy: SYN Cookies
 - » Server does not allocate buffers on receipt of SYN
 - » Server sends SYN/ACK packet with special sequence number ("cookie")
 - » Client returns ACK with SYN/ACK seqno+1 in ackno field (like normal)
 - » Server verifies that this is valid and allocates buffers

27

TCP Connection Management

Port Scanning

nmap port-mapping tool sends SYN to a particular port on a host)

1. source receives a TCP SYN/ACK -- port is open
2. source receives a TCP RST segment -- SYN reached the host, but the port is closed (not blocked by firewall)
3. source receives nothing -- blocked by firewall

28

Wireshark Example

- ◆ See handout

- ◆ Wireshark
 - » network protocol analyzer
 - » <http://www.wireshark.org/>

- ◆ Capture examples
 - » <http://wiki.wireshark.org/SampleCaptures/>