# Transport Layer: Outline

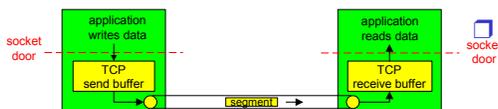- ❒ Transport-layer services
- ❒ Multiplexing and demultiplexing
- ❒ Connectionless transport: UDP
- ❒ Principles of reliable data transfer
- ❒ Connection-oriented transport: TCP
  - ❍ Segment structure
  - ❍ Reliable data transfer
  - ❍ Connection management
  - ❍ Flow control
- ❒ Principles of congestion control
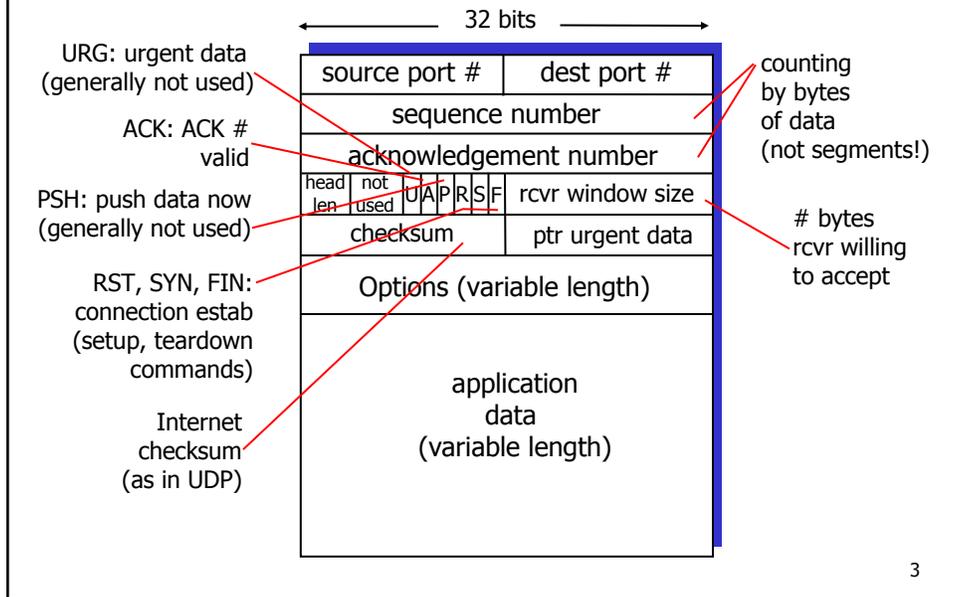- ❒ TCP congestion control

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
- ❒ Congestion controlled:
  - ❍ Sender will not overwhelm network



2

1

# TCP segment structure

32 bits

URG: urgent data
(generally not used)

ACK: ACK #
valid

PSH: push data now
(generally not used)

RST, SYN, FIN:
connection estab
(setup, teardown
commands)

Internet
checksum
(as in UDP)

| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |
| head len / not used / U A P R S F | rcvr window size |
| checksum | ptr urgent data |
| Options (variable length) | |
| application data (variable length) | |

counting
by bytes
of data
(not segments!)

# bytes
rcvr willing
to accept

3

---

# Transport layer: Outline

❐ Transport-layer services
❐ Multiplexing and demultiplexing
❐ Connectionless transport: UDP
❐ Principles of reliable data transfer

❐ Connection-oriented transport: TCP
  ❍ Segment structure
  ❍ Reliable data transfer
  ❍ Connection management
  ❍ Flow control
❐ Principles of congestion control
❐ TCP congestion control

4

# 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
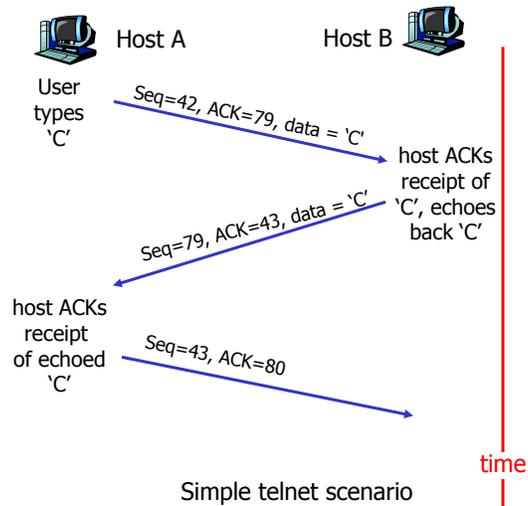  - One way dataflow

5

# 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

Q: How receiver handles out-of-order segments
- A: TCP spec doesn't say, - up to implementer

Host A                          Host B

User types 'C'

Seq=42, ACK=79, data = 'C'

host ACKs receipt of 'C', echoes back 'C'

Seq=79, ACK=43, data = 'C'

host ACKs receipt of echoed 'C'

Seq=43, ACK=80

time

Simple telnet scenario

6

3

# 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:

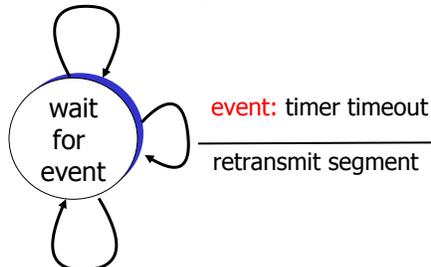❒ Retransmit the one segment that caused timeout

❒ Restart timer

### Ack rcvd:

❒ If acknowledges previously unacked segments

○ Update what is known to be acked

○ Restart timer if there are outstanding segments

7

---

# TCP: reliable data transfer

event: data received
from application above
create, send segment

Simplified sender, assuming
• One way data transfer
• No flow, congestion control

wait
for
event

event: timer timeout
retransmit segment

event: ACK received,
with ACK # y
ACK processing

8

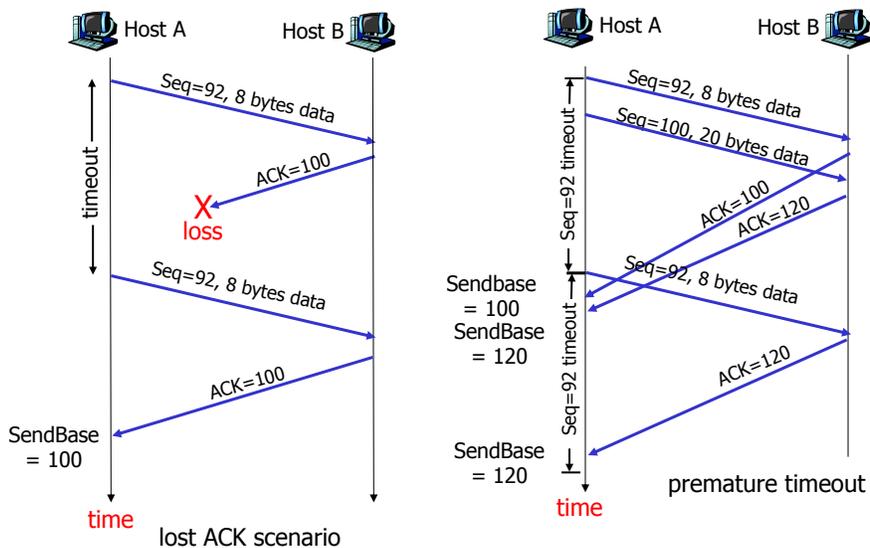## TCP sender (simplified)

```
00   sendbase = initial_sequence number
01   nextseqnum = initial_sequence number
02
03   loop (forever) {
04     switch(event)
05     event: data received from application above
06         create TCP segment with sequence number nextseqnum
07         if (timer currently not running) start timer
08         pass segment to IP
09         nextseqnum = nextseqnum + length(data)
10     event: timer timeout
11         retransmit not-yet-acknowledged segment with
12                 smallest sequence number
13         restart timer
14     event: ACK received, with ACK field value of y
15         if (y > sendbase) { /* cumulative ACK of all data up to y */
17             sendbase = y
18             if (currently not-yet-acknowledged segments) {
19                 restart timer
20             }
21         }
22   }  /* end of loop forever */
```
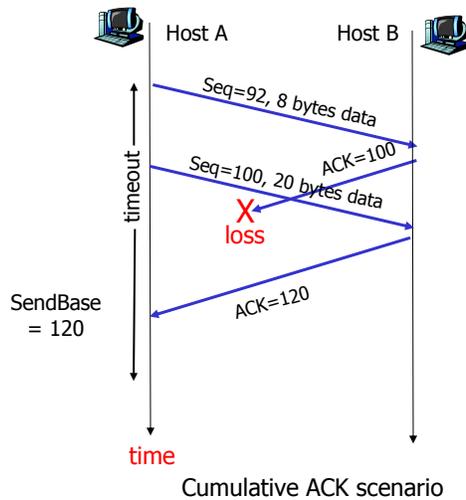
9

---

# TCP retransmission scenarios



lost ACK scenario

premature timeout

10

# TCP retransmission scenarios (2.)

Host A                    Host B

Seq=92, 8 bytes data

ACK=100

Seq=100, 20 bytes data

X
loss

SendBase
= 120

ACK=120

timeout

time

Cumulative ACK scenario

# TCP round trip time and timeout

Q: How to set TCP timeout value?

❐ Longer than RTT
  ❍ Note: RTT will vary
❐ 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, cumulatively ACKed segments
❐ **SampleRTT** will vary, want estimated RTT "smoother"
  ❍ Use several recent measurements, not just current **SampleRTT**
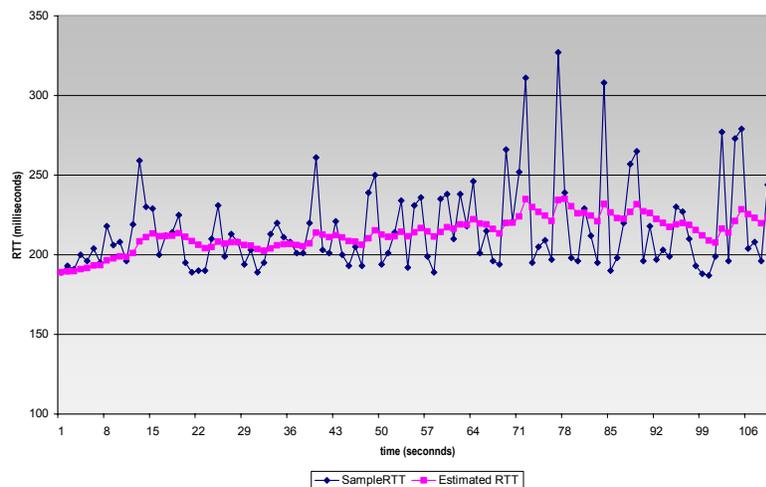
# TCP round trip time and timeout

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

- ❐ Exponential weighted moving average
- ❐ Influence of given sample decreases exponentially fast
- ❐ Typical value of $\alpha$: 0.125

- ❐ Key observation:
  - ❍ At high loads round trip variance is high

13

# Example RTT estimation

**RTT: gaia.cs.umass.edu to fantasia.eurecom.fr**



14

# TCP round trip time and timeout

## Setting the timeout

❐ **EstimtedRTT** plus "safety margin"

  ❍ Large variation in **EstimatedRTT** -> larger safety margin

❐ First estimate of how much SampleRTT deviates from EstimatedRTT:

```
DevRTT = (1-β)*DevRTT +
              β*|SampleRTT-EstimatedRTT|

(typically, β = 0.25)
```

Then set timeout interval:

```
TimeoutInterval = EstimatedRTT + 4*DevRTT
```

15

# Retransmission ambiguity



16

8

# Karn's RTT estimator

❐ Accounts for retransmission ambiguity
- ❍ If a segment has been retransmitted: Don't count RTT sample on ACKs for this segment

❐ If retransmission timer expires
- ❍ Double retransmission TimeoutInterval
- ❍ Do not use RTT estimate to calculate TimeoutInterval until successful retransmission

❐ Timer restarted (not due to timeout)
- ❍ Reuse RTT estimate

17

# Timestamp extension

❐ Used to improve timeout mechanism by more accurate measurement of RTT

❐ When sending a packet, insert current timestamp into option
- ❍ 4 bytes for seconds, 4 bytes for microseconds

❐ Receiver echoes timestamp in ACK
- ❍ Actually will echo whatever is in timestamp

❐ Removes retransmission ambiguity
- ❍ Can get RTT sample on any packet

18

# Timer granularity

❒ Many TCP implementations set RTO in multiples of 200, 500, 1000ms
❒ Why?
- ❍ Avoid spurious timeouts – RTTs can vary quickly due to cross traffic
- ❍ Make timers interrupts efficient

19

# 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

20

# 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)
                restart 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
            }
```

Duplicate ACK for
already ACKed segment

Fast retransmit

# Delayed ACK

❒ It is inefficient to send too many ACK only packets
❒ Why?
  ❍ No data => >40 Bytes for 1 byte of information
❒ Goal:
  ❍ Wait for additional data to piggy bag ACK on data pkt.
❒ Implementation
  ❍ Try to not ACK every packet but only ever second
  ❍ Wait for at most 200ms
  ❍ ACK any out of order data

# TCP ACK generation [RFC 1122, RFC 2581]

| Event at Receiver | TCP Receiver action |
|---|---|
| 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 |
| In-order segment with expected seq #. One other segment has ACK pending | Immediately send single cumulative ACK, ACKing both in-order segments |
| Out-of-order segment higher-than-expect seq. # . Gap detected | Immediately send *duplicate ACK*, indicating seq. # of next expected byte |
| Segment that partially or completely fills gap | Immediate send ACK, provided that segment starts at lower end of gap |

23

# Transport layer: Outline

❐ Transport-layer services
❐ Multiplexing and demultiplexing
❐ Connectionless transport: UDP
❐ Principles of reliable data transfer

❐ Connection-oriented transport: TCP
  ❍ Segment structure
  ❍ Reliable data transfer
  ❍ Connection management
  ❍ Flow control
❐ Principles of congestion control
❐ TCP congestion control

24

# 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");`
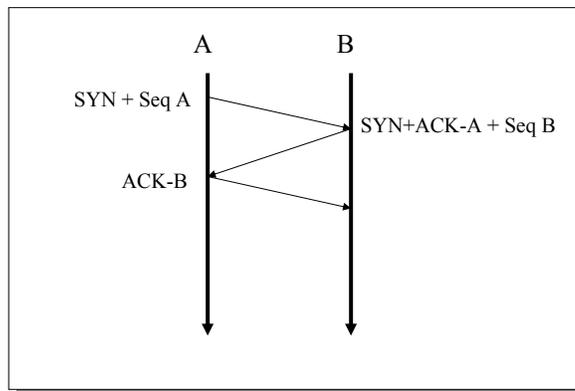- *server:* contacted by client

    `Socket connectionSocket = welcomeSocket.accept();`

25

# Connection establishment

- Use 3-way handshake



26

13

# Sequence number selection

❒ Why not simply chose 0?

❒ Must avoid overlap with earlier incarnation

# TCP connection: Three way handshake

Step 1: Client end system sends TCP SYN control segment to server
  - ❍ Specifies initial seq #
  - ❍ Specifies initial window #

Step 2: Server end system receives SYN, replies with SYNACK control segment

  - ❍ ACKs received SYN
  - ❍ Allocates buffers
  - ❍ Specifies server-> receiver initial seq. #
  - ❍ Specifies initial window #

Step 3: Client system receives SYNACK, replies with ACK segment which may contain data
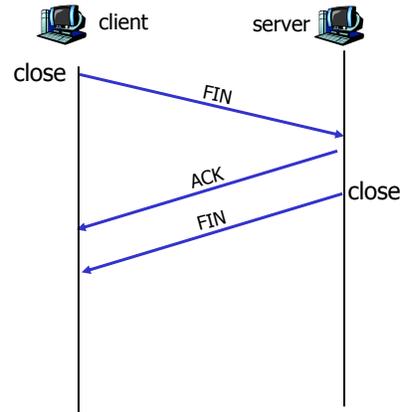
# TCP connection management (2.)

**Closing a connection:**

client closes socket:

`clientSocket.close();`

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

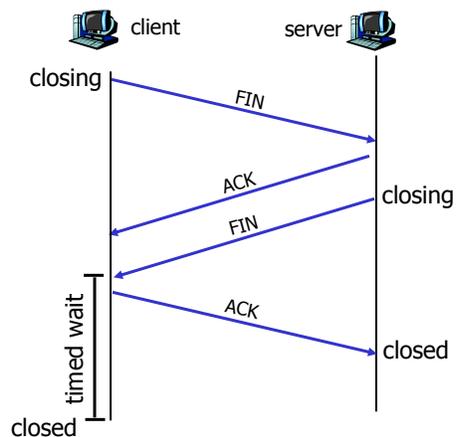**Step 2:** Server receives FIN, replies with ACK. Closes connection, sends FIN.

client     server

close    *FIN*

*ACK*   close

*FIN*

29

---

# TCP connection management (3.)

**Step 3:** Client receives FIN, replies with ACK.

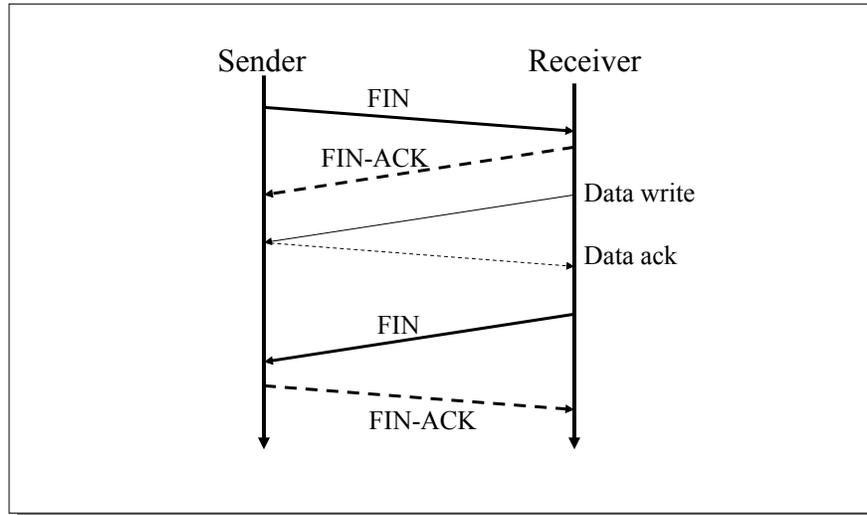- ❍ Enters "timed wait" - will respond with ACK to received FINs

**Step 4:** Server, receives ACK. Connection closed.

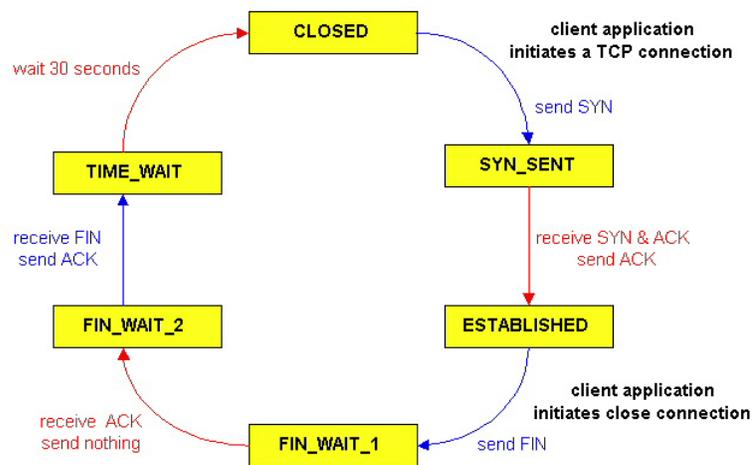**Note:** With small modification, can handle simultaneous FINs.

client     server

closing   *FIN*
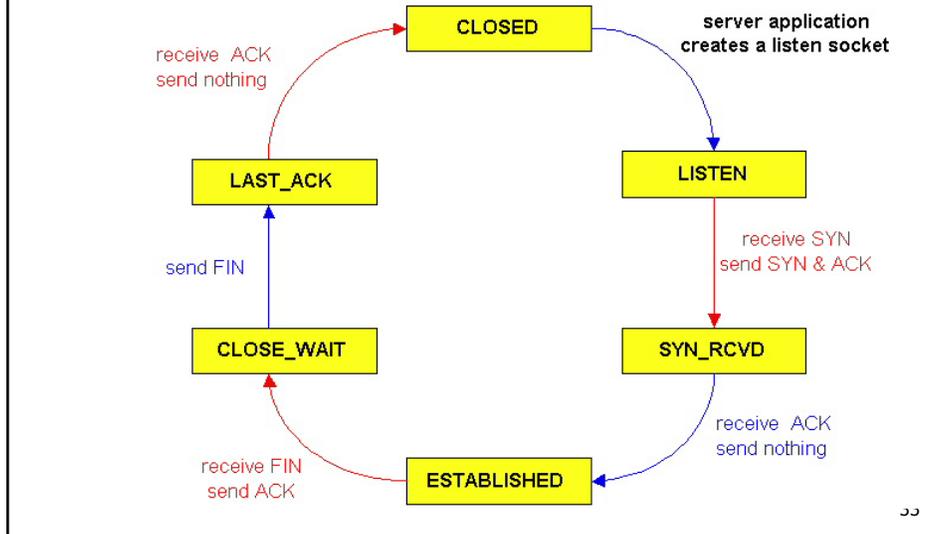
*ACK*   closing

*FIN*

timed wait   *ACK*   closed

closed

30

# Tear-down packet exchange

Sender                    Receiver

FIN

FIN-ACK

Data write

Data ack

FIN

FIN-ACK

# TCP connection management (cont.)

TCP client lifecycle

# TCP connection management (cont.)

## TCP server lifecycle



(CLOSED → LISTEN: server application creates a listen socket; LISTEN → SYN_RCVD: receive SYN / send SYN & ACK; SYN_RCVD → ESTABLISHED: receive ACK / send nothing; ESTABLISHED → CLOSE_WAIT: receive FIN / send ACK; CLOSE_WAIT → LAST_ACK: send FIN; LAST_ACK → CLOSED: receive ACK / send nothing)

33

---

# Detecting half-open connections

```
        TCP A                                              TCP B

1.   (CRASH)                                        (send 300, receive 100)
2.   CLOSED                                         ESTABLISHED
3.   SYN-SENT → <SEQ=400><CTL=SYN>           →    (??)
4.   (!!)        ← <SEQ=300><ACK=100><CTL=ACK>  ←   ESTABLISHED
5.   SYN-SENT → <SEQ=100><CTL=RST>           →    (Abort!!)
6.   SYN-SENT                                       CLOSED
7.   SYN-SENT → <SEQ=400><CTL=SYN>           →
```
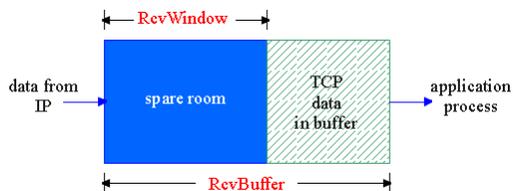
34

17

# Transport layer: outline

❒ Transport-layer services
❒ Multiplexing and demultiplexing
❒ Connectionless transport: UDP
❒ Principles of reliable data transfer

❒ Connection-oriented transport: TCP
   ❍ Segment structure
   ❍ Reliable data transfer
   ❍ Connection management
   ❍ Flow control
❒ Principles of congestion control
❒ TCP congestion control

35

# 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



❒ App process may be slow at reading from buffer

❒ Speed-matching service: match the send rate to the receiving app's drain rate

36

18

# 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**
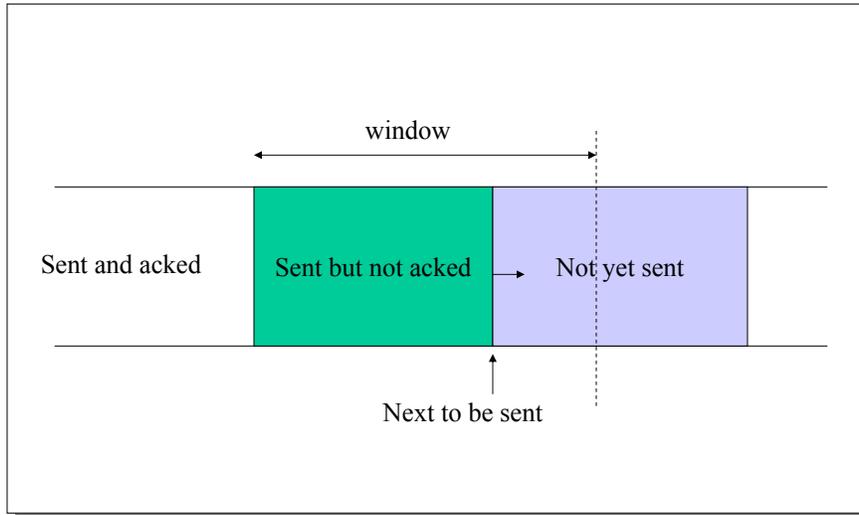  ❍ Guarantees receive buffer doesn't overflow

37

---

# TCP flow control: How it works (2.)

❐ TCP is a sliding window protocol
  ❍ For window size *n*, can send up to *n* bytes without receiving an acknowledgement
  ❍ When the data is acknowledged then the window slides forward

❐ Each packet advertises a window size
  ❍ Indicates number of bytes the receiver has space for

❐ Original TCP always sent entire window
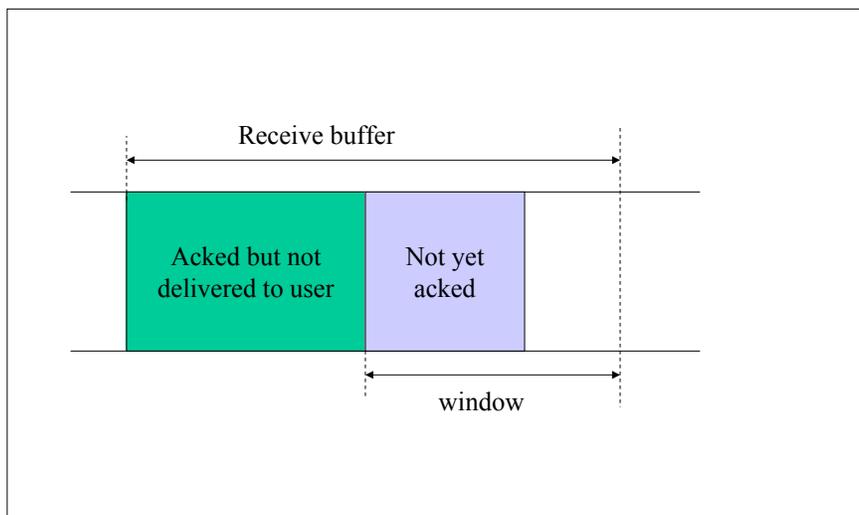  ❍ Congestion control now limits this

38

# Window flow control: Sender side

window

Sent and acked | Sent but not acked | Not yet sent

Next to be sent

39

# Window flow control: Receiver side

Receive buffer

Acked but not delivered to user | Not yet acked

window

40

# TCP persist

❒ What happens if window is 0?
  ❍ Receiver updates window (i.e., sends ACK with new window size) when application reads data
  ❍ What if this update is lost?
❒ TCP persist state
  ❍ Sender periodically sends 1 byte packets
  ❍ Receiver responds with ACK even if it can't store the packet

# Observed TCP problems

❒ Too many small packets
  ❍ Silly window syndrome
  ❍ Nagel's algorithm
❒ Initial sequence number selection
❒ Amount of state maintained

# Silly window syndrome

❒ Problem: (Clark, 1982)
  ❍ If receiver advertises small increases in the receive window then the sender may waste time sending lots of small packets
❒ Solution
  ❍ Receiver must not advertise small window increases
  ❍ Increase window by min(MSS,RecvBuffer/2)

43

# Nagel's algorithm

❒ Small packet problem:
  ❍ Don't want to send a 41 byte packet for each keystroke
  ❍ How long to wait for more data?
❒ Solution:
  ❍ Allow only one outstanding small (not full sized) segment that has not yet been acknowledged

44

# Why is selecting ISN important?

❏ Suppose machine X selects ISN based on predictable sequence

❏ Fred has .rhosts to allow login to X from Y

❏ Evil Ed attacks
  ○ Disables host Y – denial of service attack
  ○ Make a bunch of connections to host X
  ○ Determine ISN pattern and guess next ISN
  ○ Fake pkt1: [<src Y><dst X>, guessed ISN]
  ○ Fake pkt2: desired command

45

# Time Wait issues

❏ Web servers not clients close connection first
  ○ Established → Fin-Waits → Time-Wait → Closed
  ○ Why would this be a problem?

❏ Time-Wait state lasts for 2 * MSL
  ○ MSL is should be 120 seconds (is often 60s)
  ○ Servers often have order of magnitude more connections in Time-Wait

46

# Transport layer: Outline

❒ Transport-layer services
❒ Multiplexing and demultiplexing
❒ Connectionless transport: UDP
❒ Principles of reliable data transfer

❒ Connection-oriented transport: TCP
  ❍ Segment structure
  ❍ Reliable data transfer
  ❍ Flow control
  ❍ Connection management
❒ Principles of congestion control
❒ TCP congestion control

47