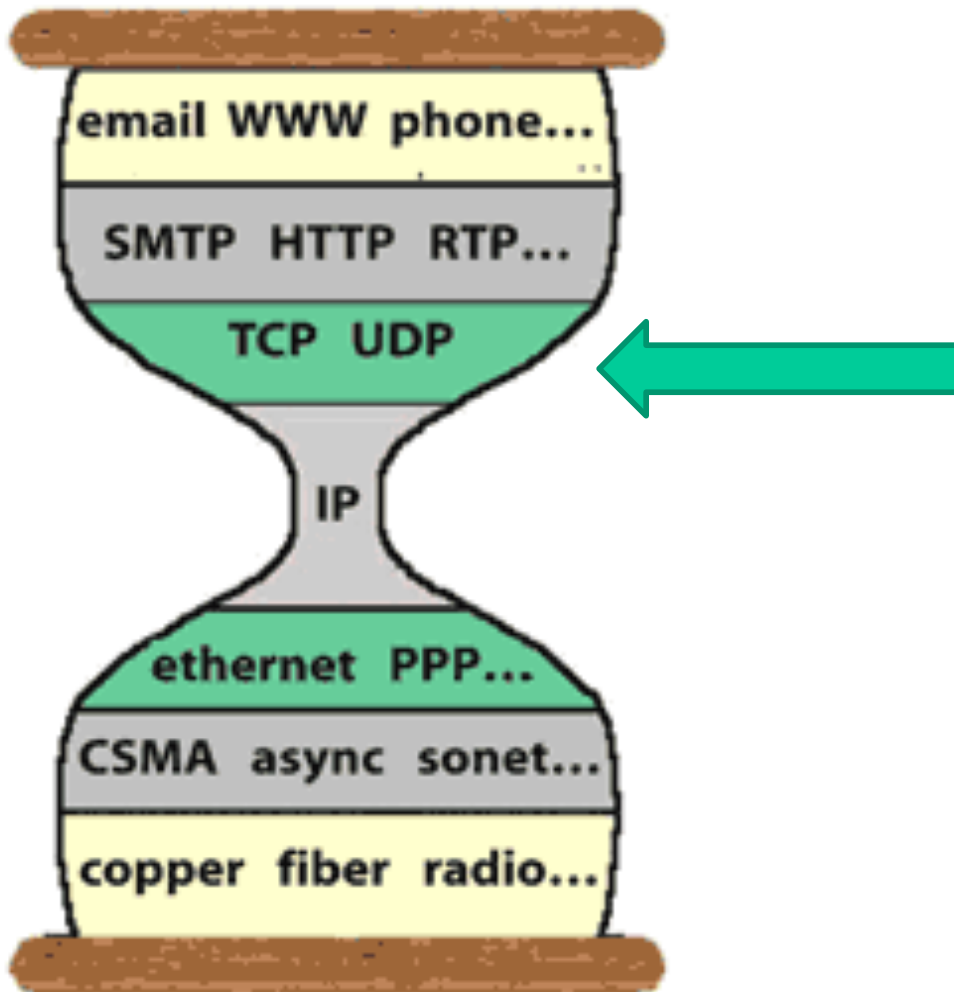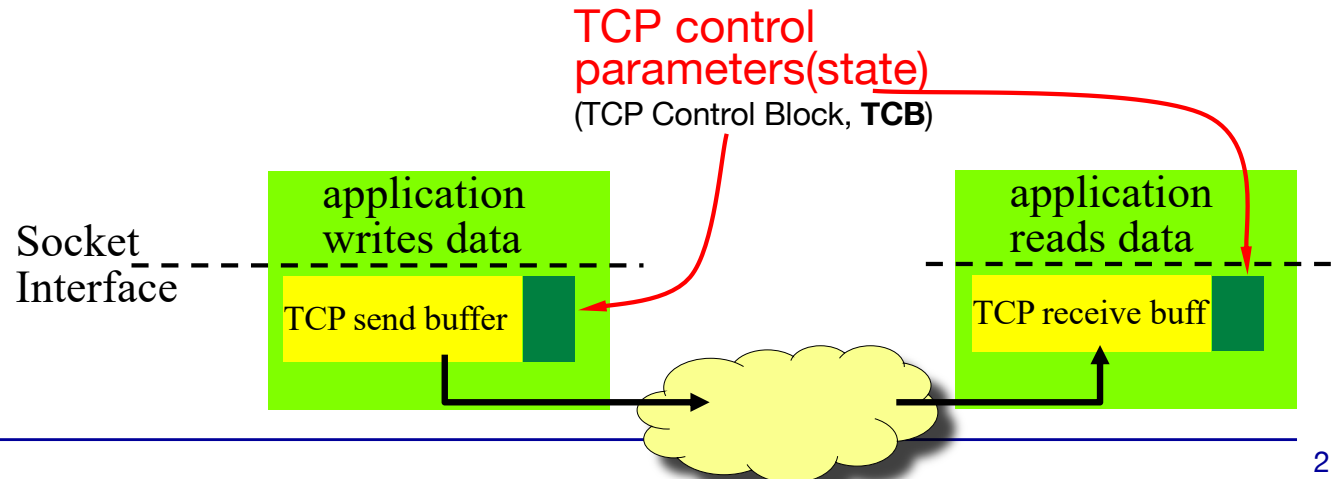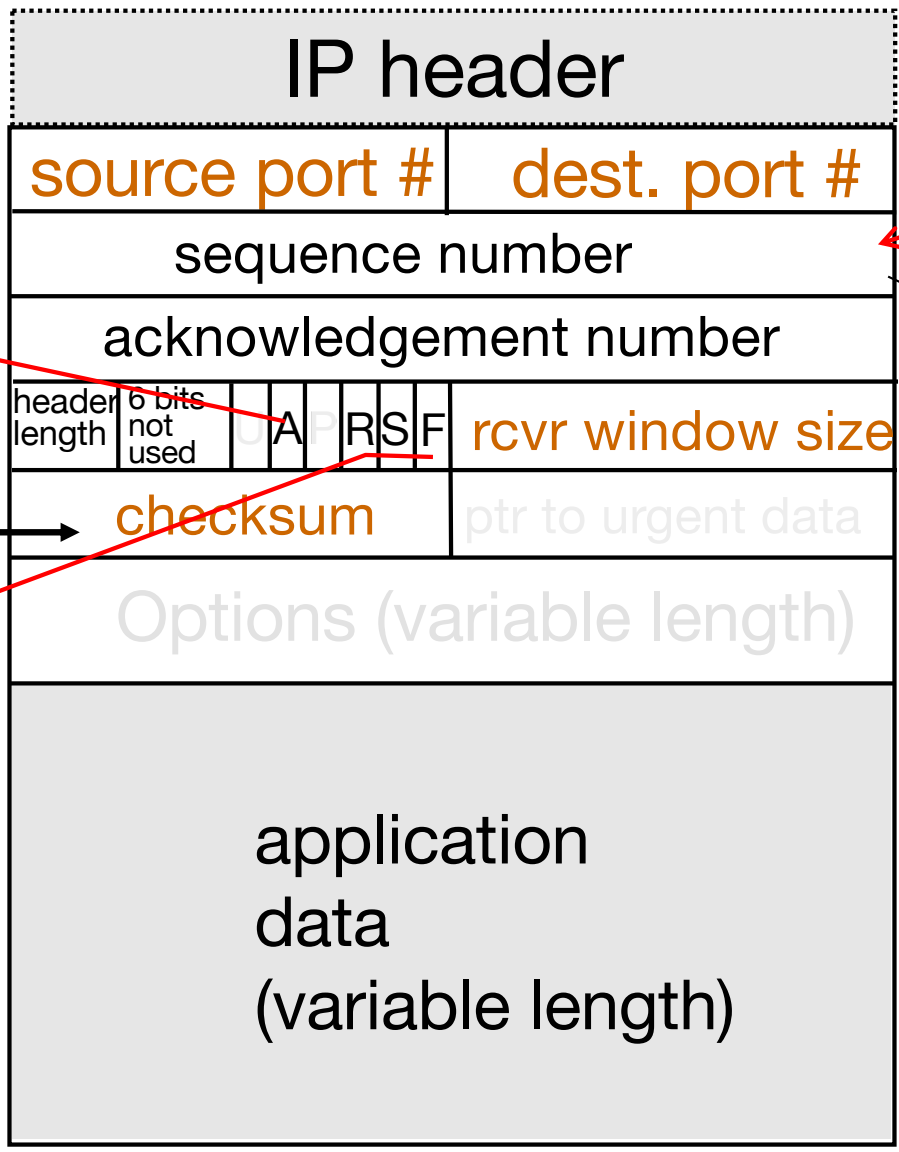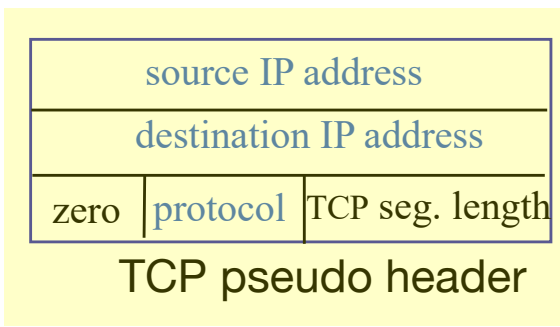# Lecture-6: TCP



## Chapter 3

## 3.5 TCP

- Protocol format
- Connection management
- Flow control
- Retransmission timer

# TCP function Overview

◆ point-to-point: creating a virtual pipe between 2 processes

◆ connection-oriented: set up connection first before data transmission, tear down the connection after finish

◆ bi-directional, reliable byte steam delivery (figure illustrates one way only)
  ■ no "message" boundaries

◆ flow controlled: prevent sender from overwhelming <u>receiver</u>

◆ congestion controlled: mitigating traffic overload <u>*inside the network*</u>

TCP control parameters(state)
(TCP Control Block, **TCB**)

application writes data

Socket Interface

TCP send buffer

application reads data

TCP receive buff

# TCP segment format

**TCP pseudo header**

| source IP address | | |
|---|---|---|
| destination IP address | | |
| zero | protocol | TCP seg. length |

counting the number of **bytes**

| IP header | |
|---|---|
| source port # | dest. port # |
| sequence number | |
| acknowledgement number | |

| header length | 6 bits not used | U | A | P | R | S | F | rcvr window size |

| checksum | ptr to urgent data |

Options (variable length)

application data (variable length)

ACK flag: ACK# field valid

Checksum is computed over TCP segment plus pseudo header

SYN, FIN, Reset: connection management flags (Setup, Finish, Reset)

TCP header has no info for congestion control

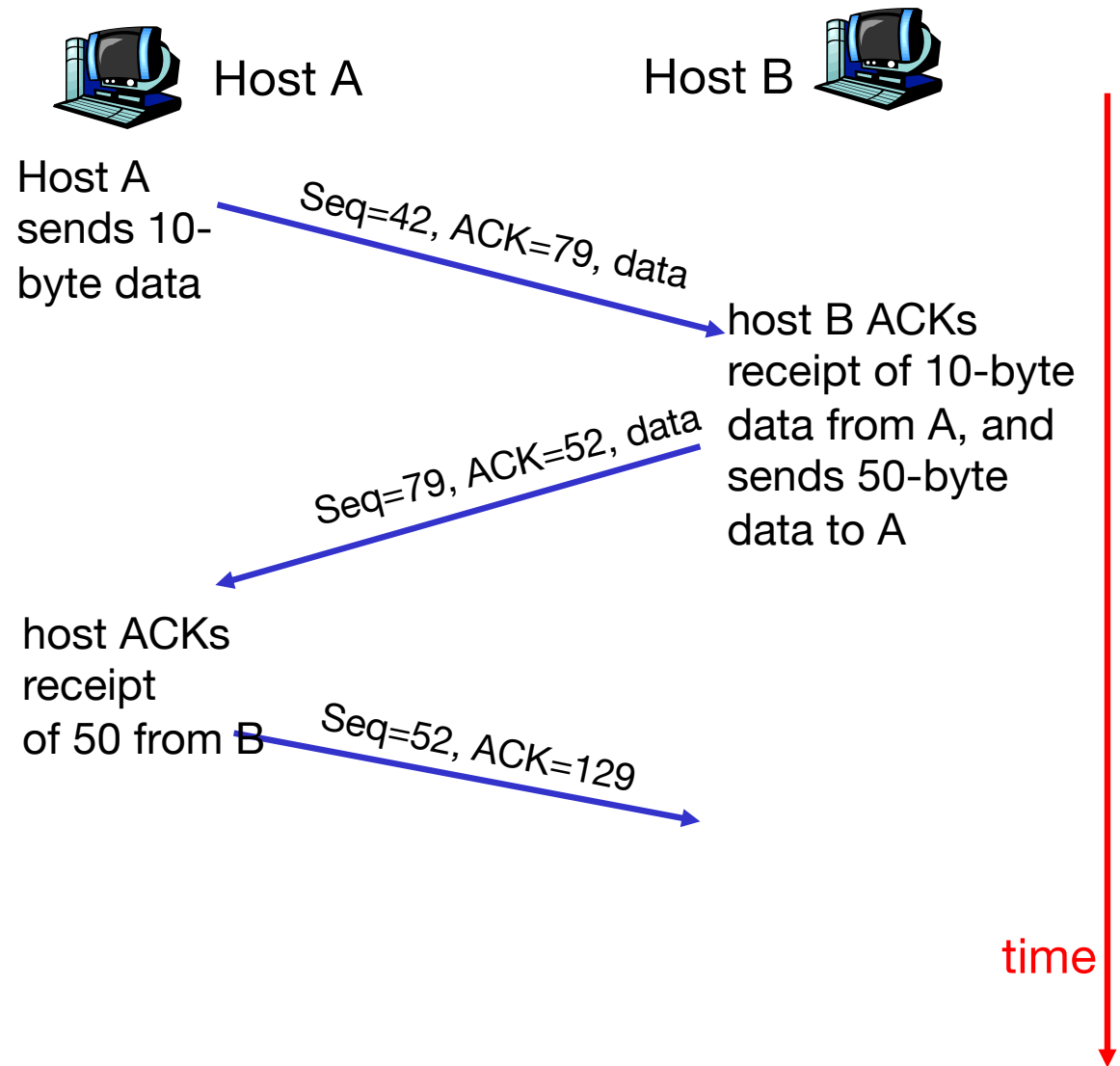Seq# of the first byte in the payload

32 bits

# TCP's seq. #s and ACK #s

Lets first assume that a TCP connection between A and B is already setup:

**Seq. #:** the seq# of the first byte in this segment's data

**ACK #:** the seq# of next byte expected from the other end

◆ TCP uses cumulative ACK

Host A         Host B

Host A sends 10-byte data

Seq=42, ACK=79, data

host B ACKs receipt of 10-byte data from A, and sends 50-byte data to A

Seq=79, ACK=52, data

host ACKs receipt of 50 from B

Seq=52, ACK=129

time

# TCP Connection Management

- ◆ Set up connection before starting data transmission
    - Each of the 2 ends reliably informs the other its initial data byte sequence number value

- ◆ Close connection after finishing data transmission
    - Each of the 2 ends reliably informs the other its final data byte sequence number value

- ◆ Abort connection
    - When receiving a RST segment
    - When a node may send a RST segment
        - receives a TCP segment of unknown connection
        - TCP retransmission count hits the upper-bound
        - need to reject a new connection request or close an existing TCP connection, due to resource limitation

# TCP Connection Setup

Initialize TCP connection variables to get ready before sending data
- Initial seq. # used in each direction
- Buffer size (rcvWindow)

<u>3-way handshake</u> in setting up a connection

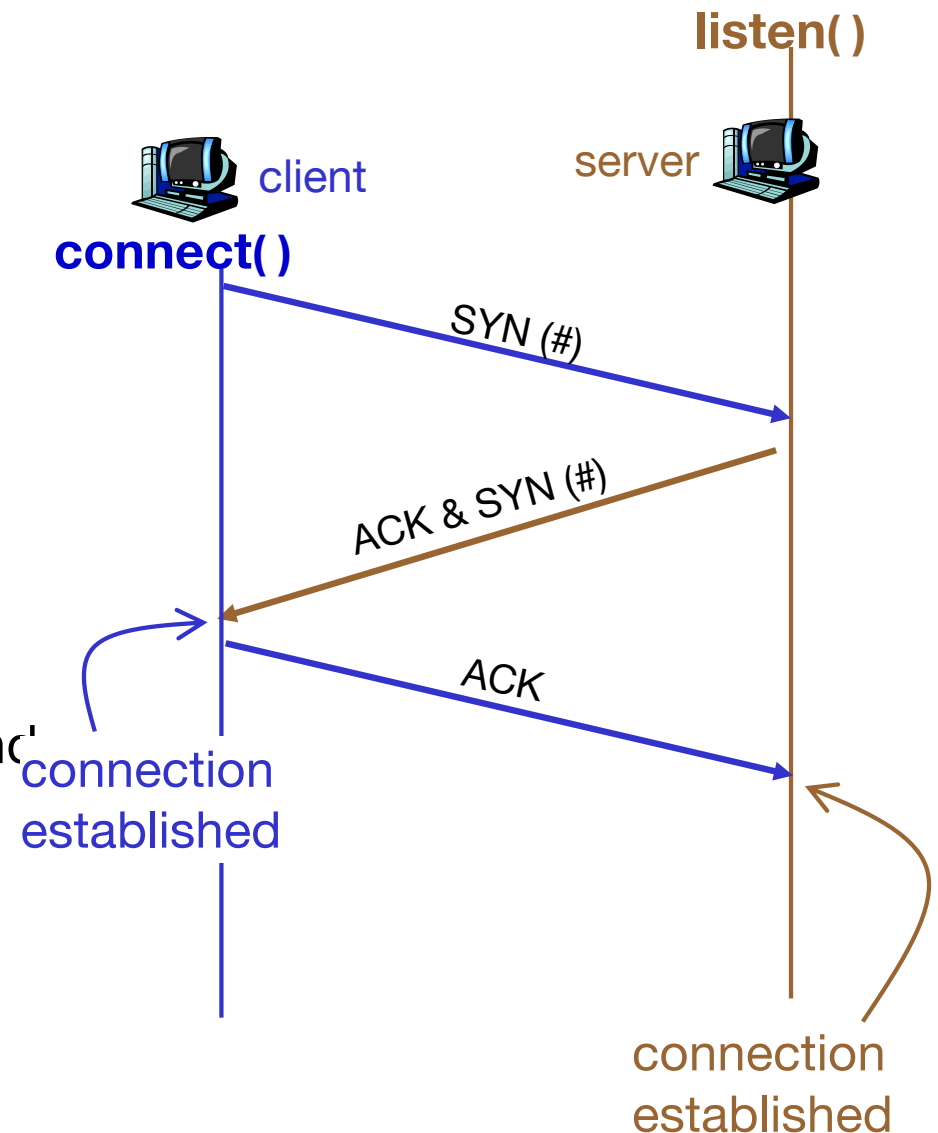<u>1:</u> client host sends TCP SYN segment to server
- SYN flag sets to 1
- specifies client's initial seq #
  - a random number
- does *not* carry data

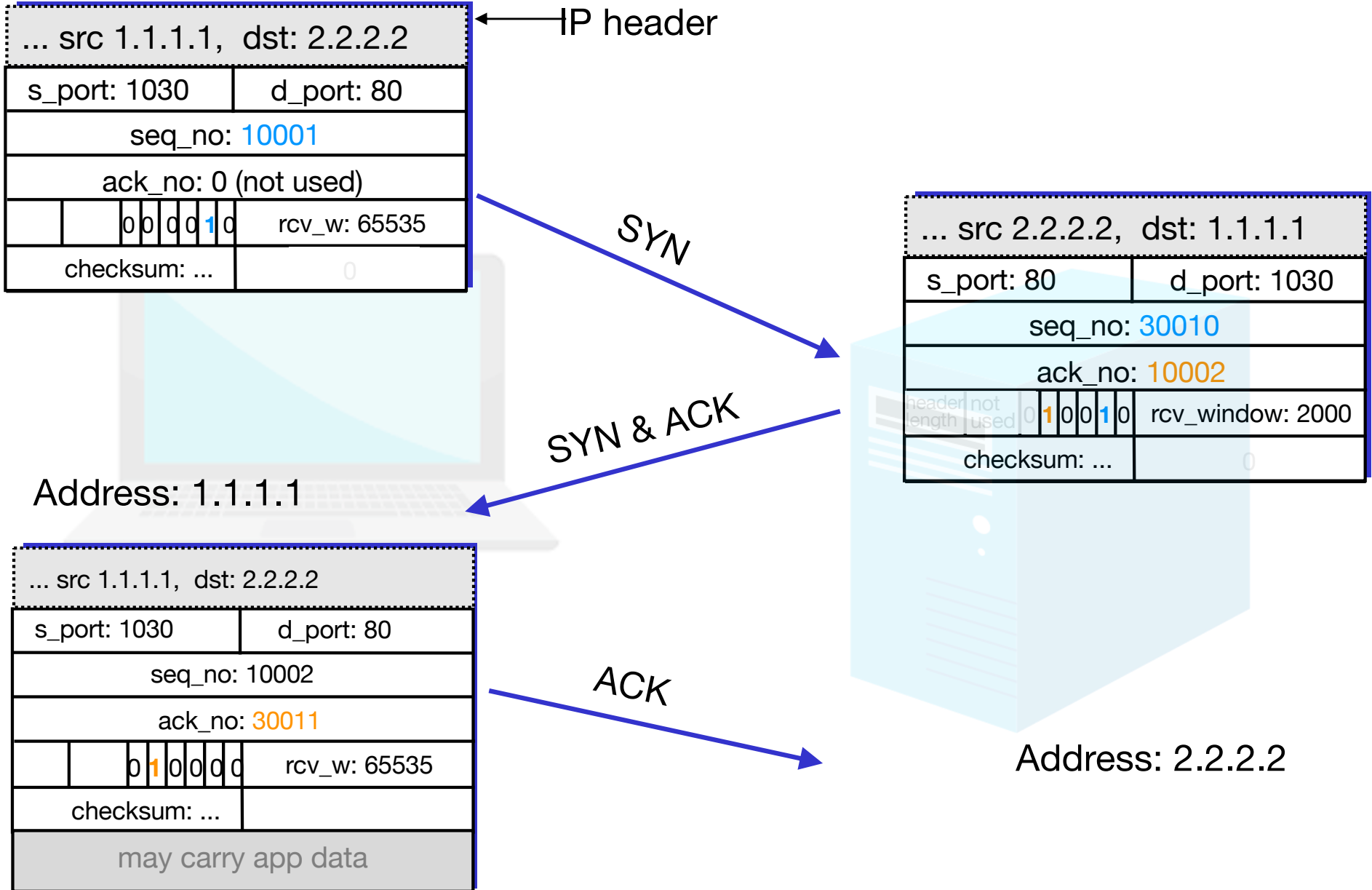<u>2:</u> server receives SYN, replies with **ACK** and **SYN** control segment
- SYN and ACK flags set to 1
  - ACK received seq#
- Specifies its own initial seq #
  - also selected randomly

<u>3:</u> client host sends ACK
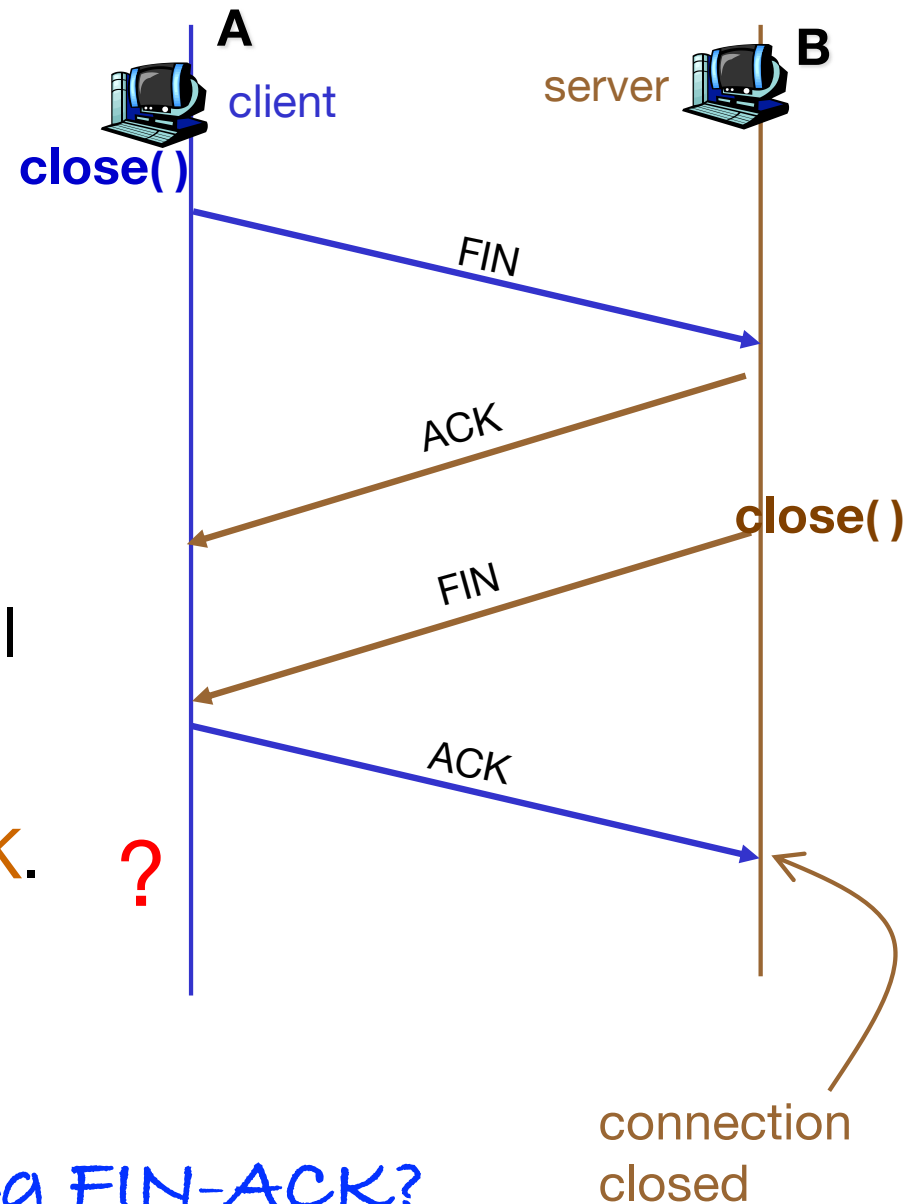- ACK flag sets to 1
  - ACK received seq#
- May carry data

**listen( )**

client

server

**connect( )**

SYN (#)

ACK & SYN (#)

ACK

connection established

connection established

# A TCP connection setup example

IP header

**Packet 1 (SYN):**

| ... src 1.1.1.1, dst: 2.2.2.2 | |
|---|---|
| s_port: 1030 | d_port: 80 |
| seq_no: 10001 | |
| ack_no: 0 (not used) | |
| 0 0 0 0 1 0 | rcv_w: 65535 |
| checksum: ... | 0 |

SYN →

Address: 1.1.1.1

**Packet 2 (SYN & ACK):**

| ... src 2.2.2.2, dst: 1.1.1.1 | |
|---|---|
| s_port: 80 | d_port: 1030 |
| seq_no: 30010 | |
| ack_no: 10002 | |
| header length / not used / 0 1 0 0 1 0 | rcv_window: 2000 |
| checksum: ... | 0 |

← SYN & ACK

Address: 2.2.2.2

**Packet 3 (ACK):**

| ... src 1.1.1.1, dst: 2.2.2.2 | |
|---|---|
| s_port: 1030 | d_port: 80 |
| seq_no: 10002 | |
| ack_no: 30011 | |
| 0 1 0 0 0 0 | rcv_w: 65535 |
| checksum: ... | |
| may carry app data | |

ACK →

# TCP Connection Close

Either end can initiate the close of **its end** of the connection at any time

**1:** A sends TCP FIN control segment to the other
- FIN flag sets to 1
- This segment must not carry data

**2:** the other end (B) receives FIN segment, replies with ACK
- regular ACK, ACK A's FIN

3: later when B finishes sending all its data and ready to close, it sends FIN segment

**4:** A receives FIN, replies with ACK.
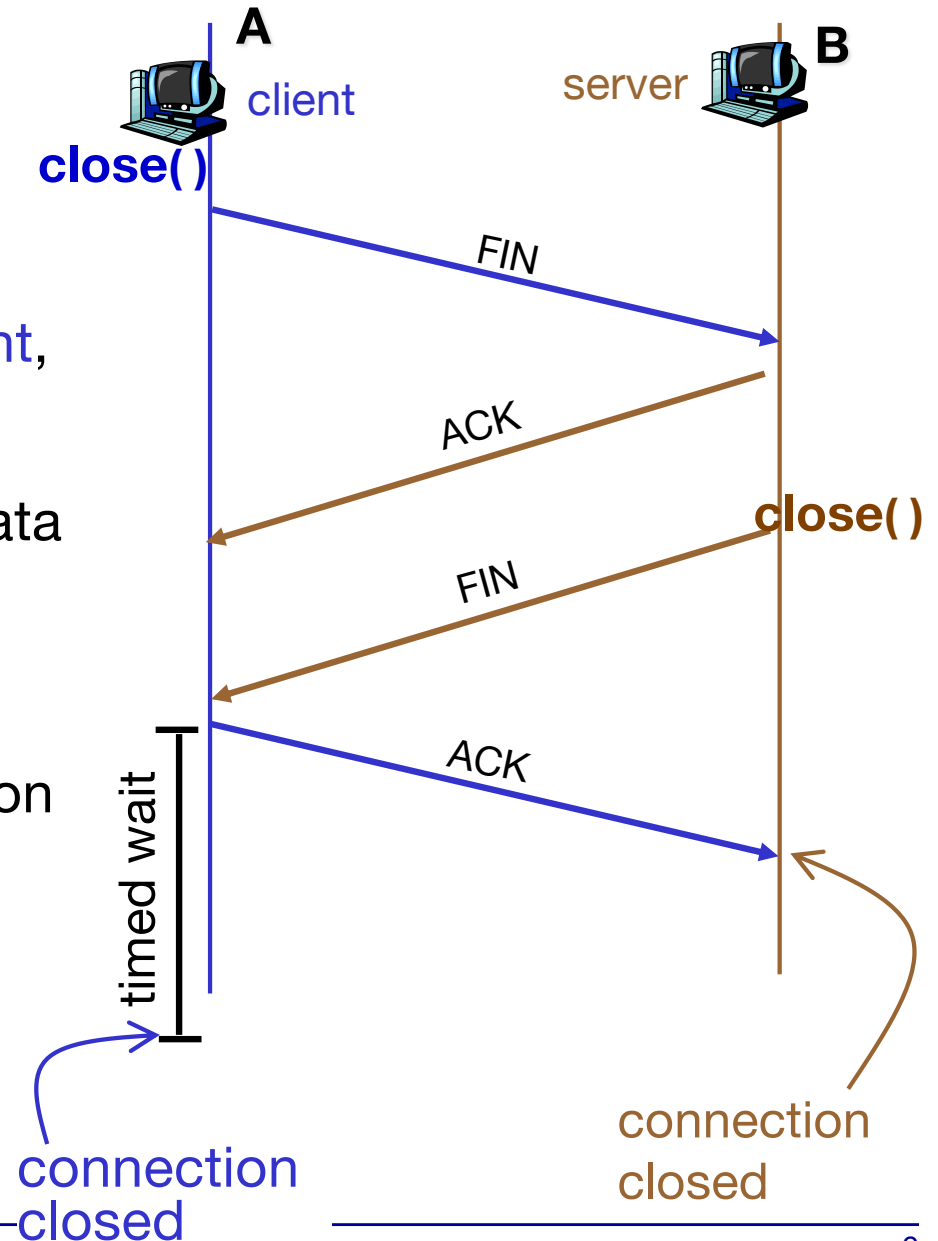
**5:** B receives FIN-ACK, closes connection

*what should A do after sending FIN-ACK?*

A client    B server

close( )

FIN

ACK

close( )

FIN

ACK

?

connection closed

# TCP Connection Close

Either end can initiate the close of **its end** of the connection at any time

**1:** A sends TCP FIN control segment to the other
  - FIN flag sets to 1
  - This segment must not carry data

**2:** the other end (B) receives FIN segment, replies with ACK
  - regular ACK, ACK A's FIN

**3:** later when B finishes sending all its data and ready to close, it sends FIN segment

**4:** A receives FIN, replies with ACK.

**5:** B receives FIN-ACK, closes connection

**6:** A closes the connection after waiting for "long enough" time w/o receiving retransmitted FIN
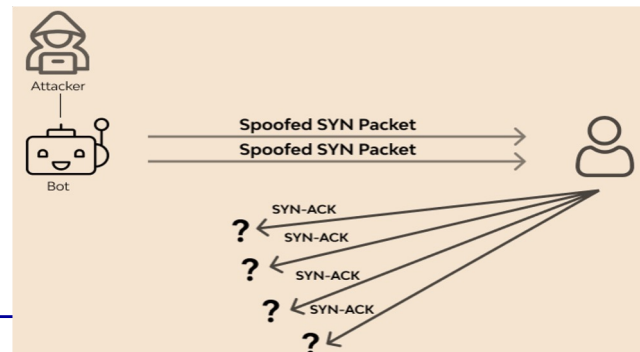  - Long enough = 2 x Max. Seg. Lifetime
    Max. Seg. Lifetime = 2 minutes

A  client

close( )

FIN

ACK

close( )

FIN

ACK

timed wait

connection closed

server  B

connection closed

# When to send "connection reset"

- ◆ Upon a TCP connection setup request: the system sets up a TCP Control Block (TCB)
  - ▪ Identified by: source+dest. addresses, source+dest. ports
  - ▪ Connection state includes info such as
    - ● Receiver flow control window size
    - ● Seq# of data
      - ▲ oldest sent but unacked
      - ▲ Latest sent, unacked
    - ● segments that arrived out of order
    - ● etc

- ◆ If TCP receives a segment (other than SYN) it *cannot* find corresponding TCB: reply with RST
  - ▪ Receiver of RST aborts the connection, all data on the connection considered lost
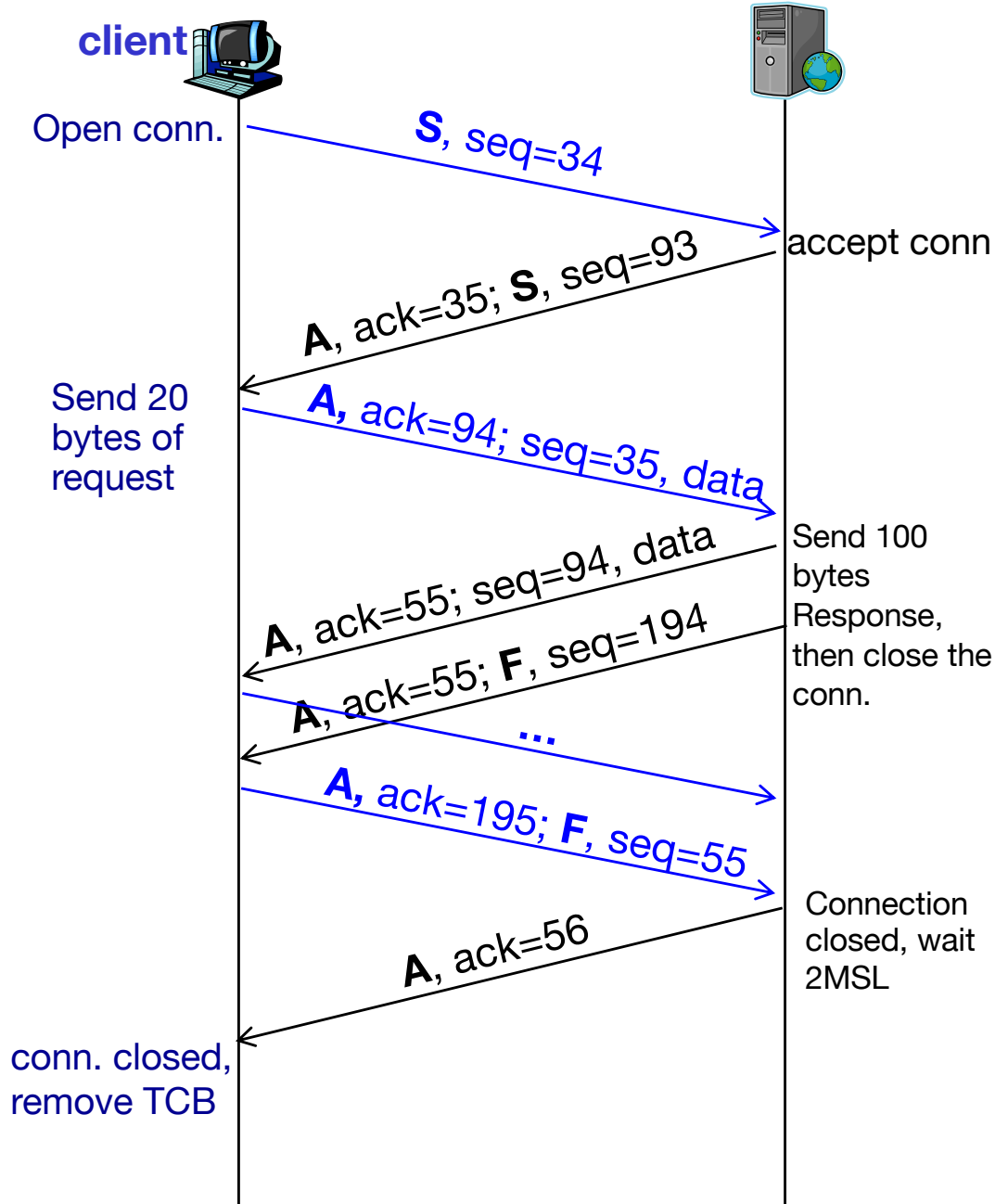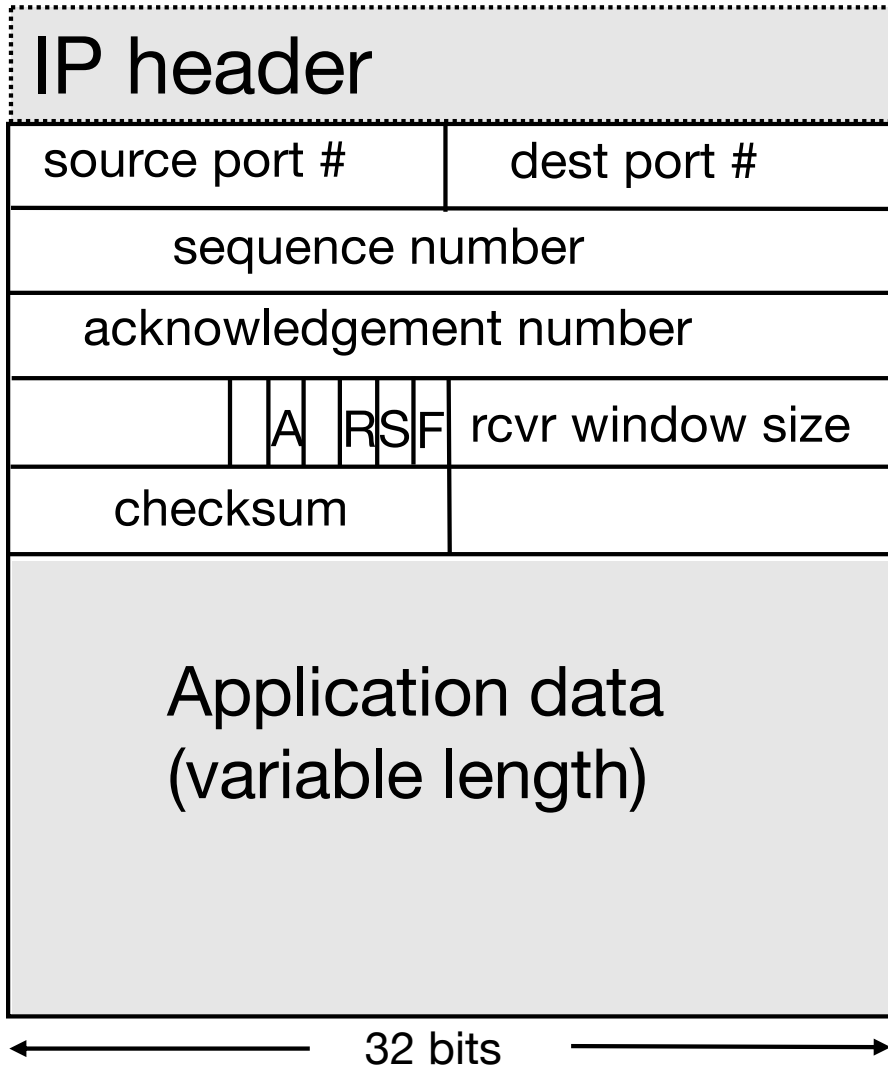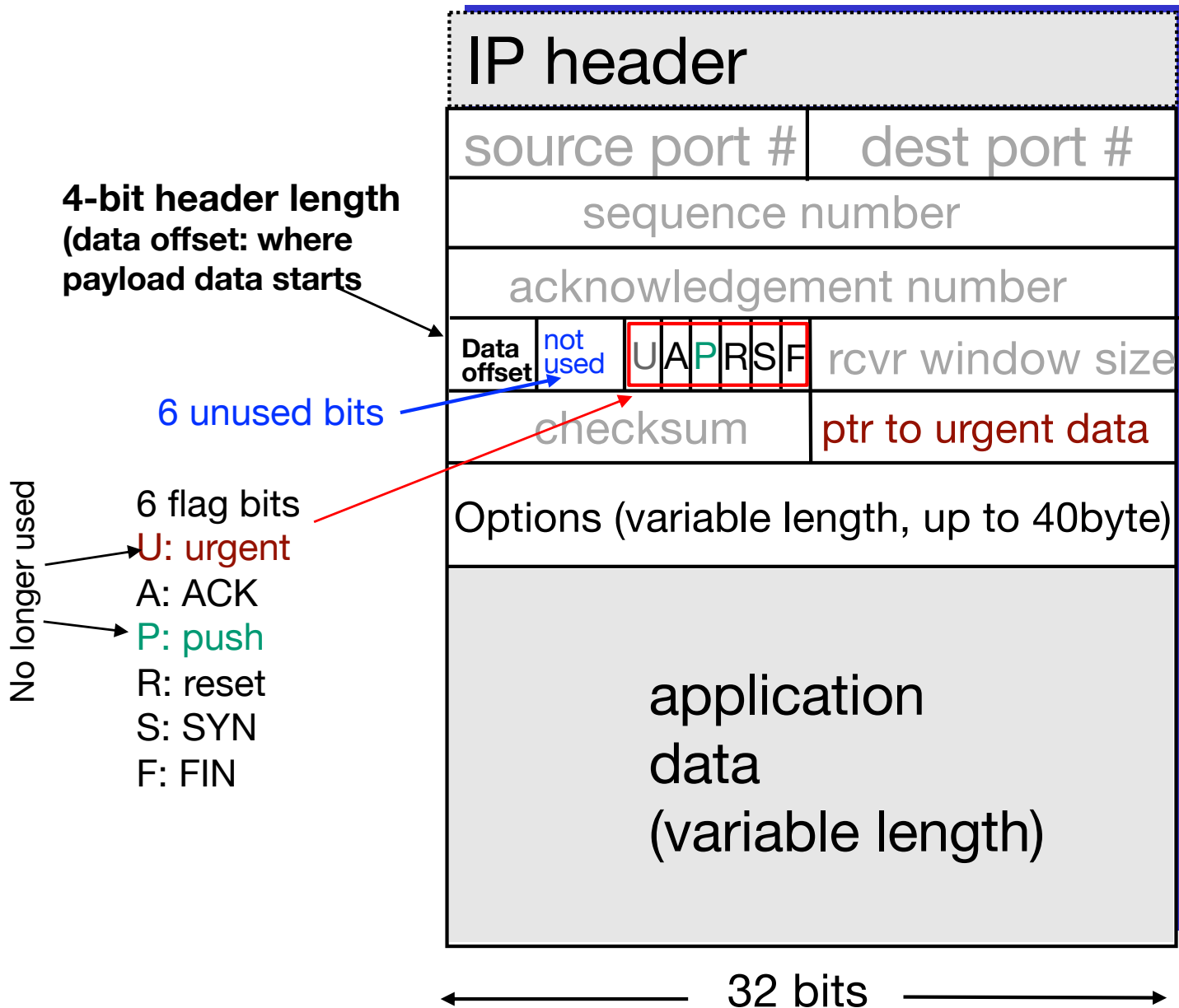
This can happen
  - ▪ Due to bit errors
  - ▪ By attacks:



Attacker

Bot

Spoofed SYN Packet
Spoofed SYN Packet

SYN-ACK
? SYN-ACK
? SYN-ACK
? SYN-ACK
?

# An HTTP 1.0 connection example

*important*

IP header

| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |
| A R S F | rcvr window size |
| checksum | |
| Application data (variable length) | |

← 32 bits →

**client**

Open conn. → **S**, seq=34 → accept conn

**A**, ack=35; **S**, seq=93 ←

Send 20 bytes of request → **A,** ack=94; seq=35, data →

Send 100 bytes Response, then close the conn.

**A**, ack=55; seq=94, data ←

**A**, ack=55; **F**, seq=194 ←

**...**

**A,** ack=195; **F**, seq=55 →

Connection closed, wait 2MSL

**A**, ack=56 ←

conn. closed, remove TCB

# TCP segment format: the remaining parts

**4-bit header length (data offset: where payload data starts**

6 unused bits

No longer used

6 flag bits
U: urgent
A: ACK
P: push
R: reset
S: SYN
F: FIN

| IP header | |
|---|---|
| source port # | dest port # |
| sequence number | |
| acknowledgement number | |
| Data offset / not used / U A P R S F | rcvr window size |
| checksum | ptr to urgent data |
| Options (variable length, up to 40byte) | |
| application data (variable length) | |

32 bits

# TCP Flow Control

**Flow control:** Prevent sender from overrunning receiver by transmitting too much data too fast

**receiver:** informs sender of amount of free buffer space
- Carried in `RcvWindow` field of TCP header of every arriving segment, **can change dynamically**

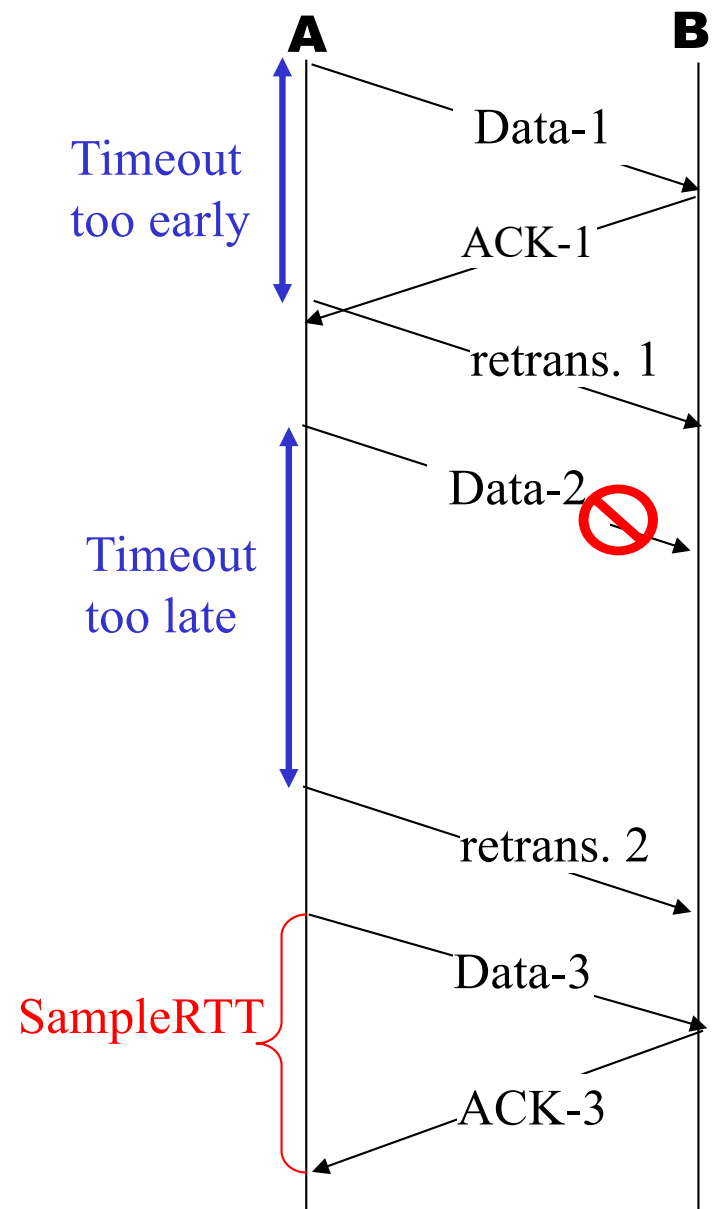**sender:** keeps the amount of transmitted, unACKed data no more than most recently received `RcvWindow` value

*Sender's output buffer*     *Receiver's input buffer*

*8 free spaces (rwnd)*

*5 new data packets (wnd)*

*rwnd=3*

# TCP loss detection and recovery

♦ TCP sets a retransmission timer (RTO) to detect packet losses

♦ A TCP connections sets one retransmission timer on the earliest sent, but unACKed segment $S$
  - If $S$ gets ACKed, restart the timer on next unACKed segment
  - (reset timer when receiving ACK for new data)

♦ When the timer expires, retransmit starting from $S$

♦ *How many segments to retransmit?*
  - Receiver flow control window, rwnd
  - Congestion control window, cwnd (next lecture)
  - the number of segments that can be retransmitted:
    ```
    min[cwnd, rwnd]
    ```
    • Dependent on how segment loss is detected, see next lecture

# Setting TCP Retransmission Timer

*important*

- ◆ TCP sets retrains. timer (RTO) based on estimated RTT
  - plus a "safety margin" (DevRTT)

- ◆ SRTT: estimated "smoothed" RTT
  - $SRTT = (1 - \alpha) \cdot SRTT + \alpha \cdot SampleRTT$
  - Exception: for the first measurement, $SRTT = SampleRTT$

- ◆ DevRTT: estimated RTT deviation
  - $DevRTT = (1 - \beta) \cdot DevRTT + \beta \cdot |SRTT - SampleRTT|$
  - Exception: for the first measurement: $DevRTT = \frac{SRTT}{2}$

- ◆ RTO: Retransmission timeout
  - $RTO = SRTT + 4 \cdot DevRTT$

- ◆ Typical parameters:
  - $\alpha = \frac{1}{8}, \quad \beta = \frac{1}{4},$

**A**      **B**

Timeout too early

Data-1

ACK-1

retrans. 1

Data-2

Timeout too late

retrans. 2

Data-3

SampleRTT

ACK-3

# No need to remember details
# Just understand the basic idea

- ◆ Network delay: random

- ◆ How to set retransmission timer:
  - Take measurements
  - Set the timer based on both average, and the variation

- ◆ Start the ball rolling: how to set the retransmission timer for the first packet of a connection?
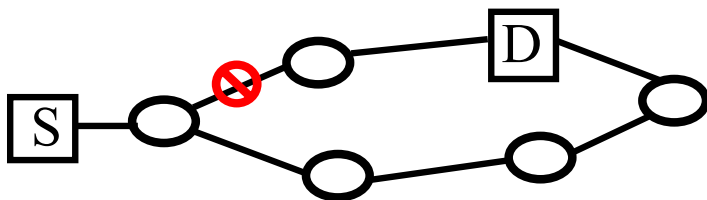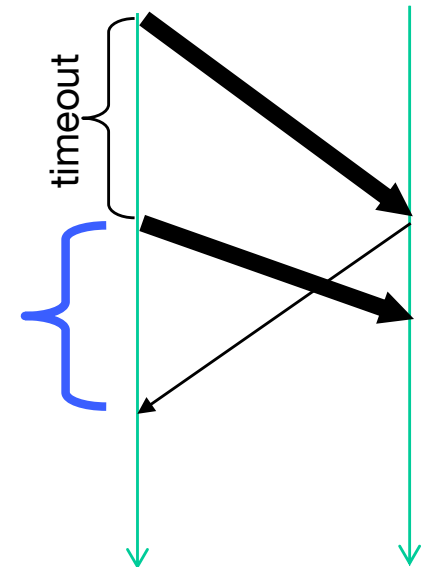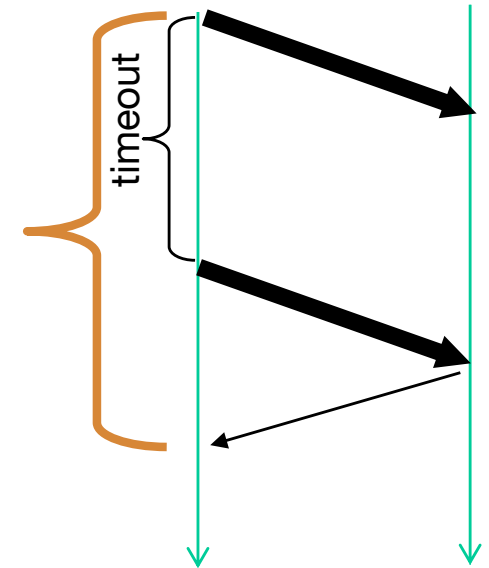
# One more question

How to set the RTO value for the first segment?

Set a default value by some engineered guessing

- what if the guessed value too small?
  - Unnecessary retransmissions

- what if the guessed value too large?
  - In case of first, or first few, packets being lost, wait longer than necessary before retransmission

- Current practice:
  - initial RTO = 1 sec (see RFC6298)
  - Once get the first sample RTT: SRTT←sample RTT, DevRTT =SRTT/2

# What to do in cases of *retransmissions*

◆ Taking measurement seems infeasible

- take the delay between first transmission and final ACK?
- take the delay between last retransmission of segment(n) and ACK(n)?

◆ Don't measure?

- Original path failed
- New path is much longer
- Without taking measurement, RTO got stuck with being too short

# Karn's algorithm

in case of segment retransmission:

- ◆ do not take the RTT sample (i.e. no update to SRTT or DevRTT)

- ◆ double the retransmission timeout value (RTO) after each timeout

- ◆ Take RTT measure again upon next successful data transmission (receiving ACK without retransmission)

# Computing RTO: an example

(from the earlier HTTP 1.0 connection example)

difference = SampleRTT − SRTT

SRTT = SRTT + 1/8 x difference

DevRTT = DevRTT +
        1/4 (|difference| - DevRTT)

RTO = SRTT + 4 x DevRTT

Initialize: RTO = 1 second

Upon receiving first packet:
  SRTT = sample RTT
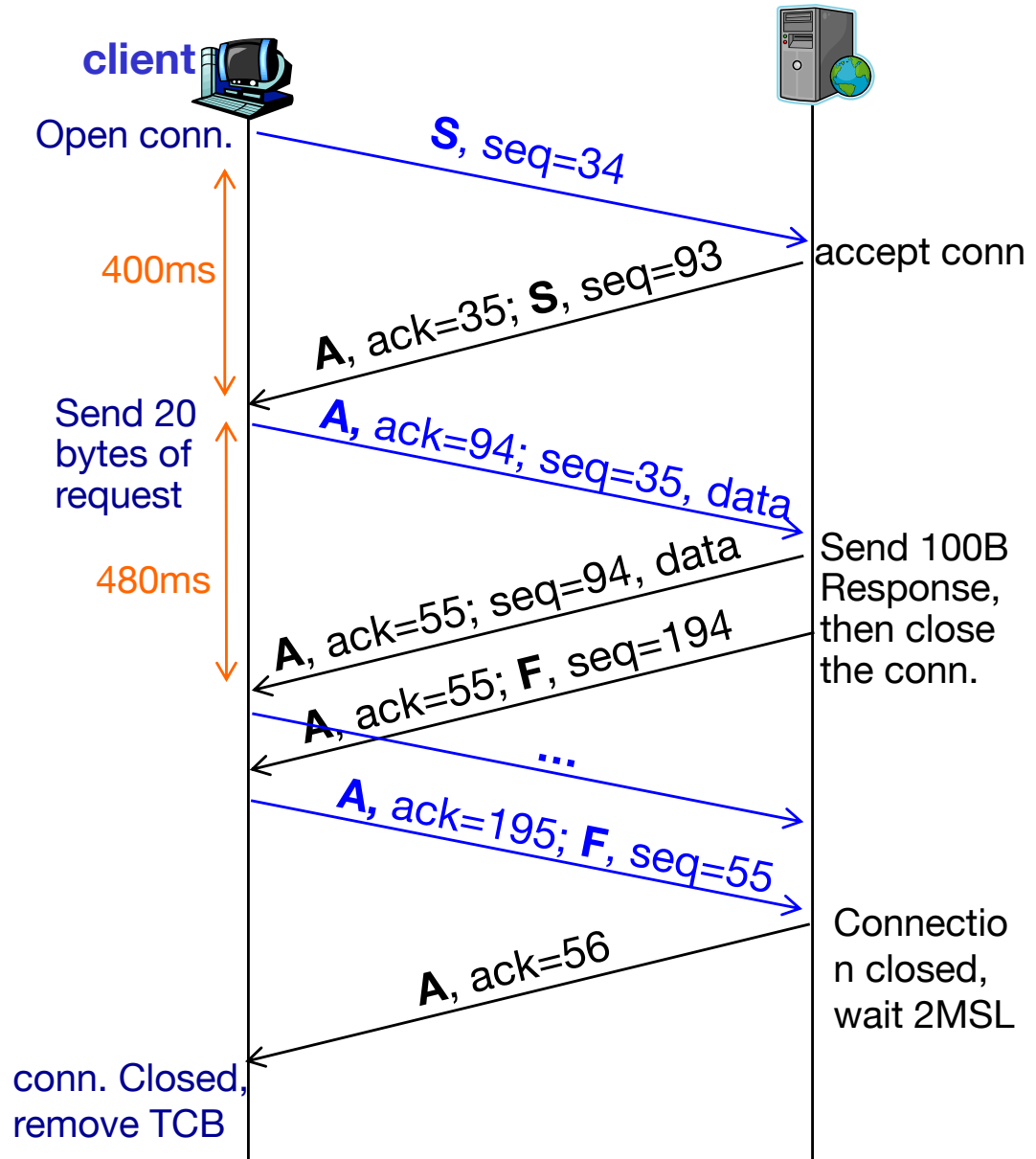  DevRTT = sample RTT / 2
  SRTT = 400, DevRTT = 200

Upon receiving second packet:
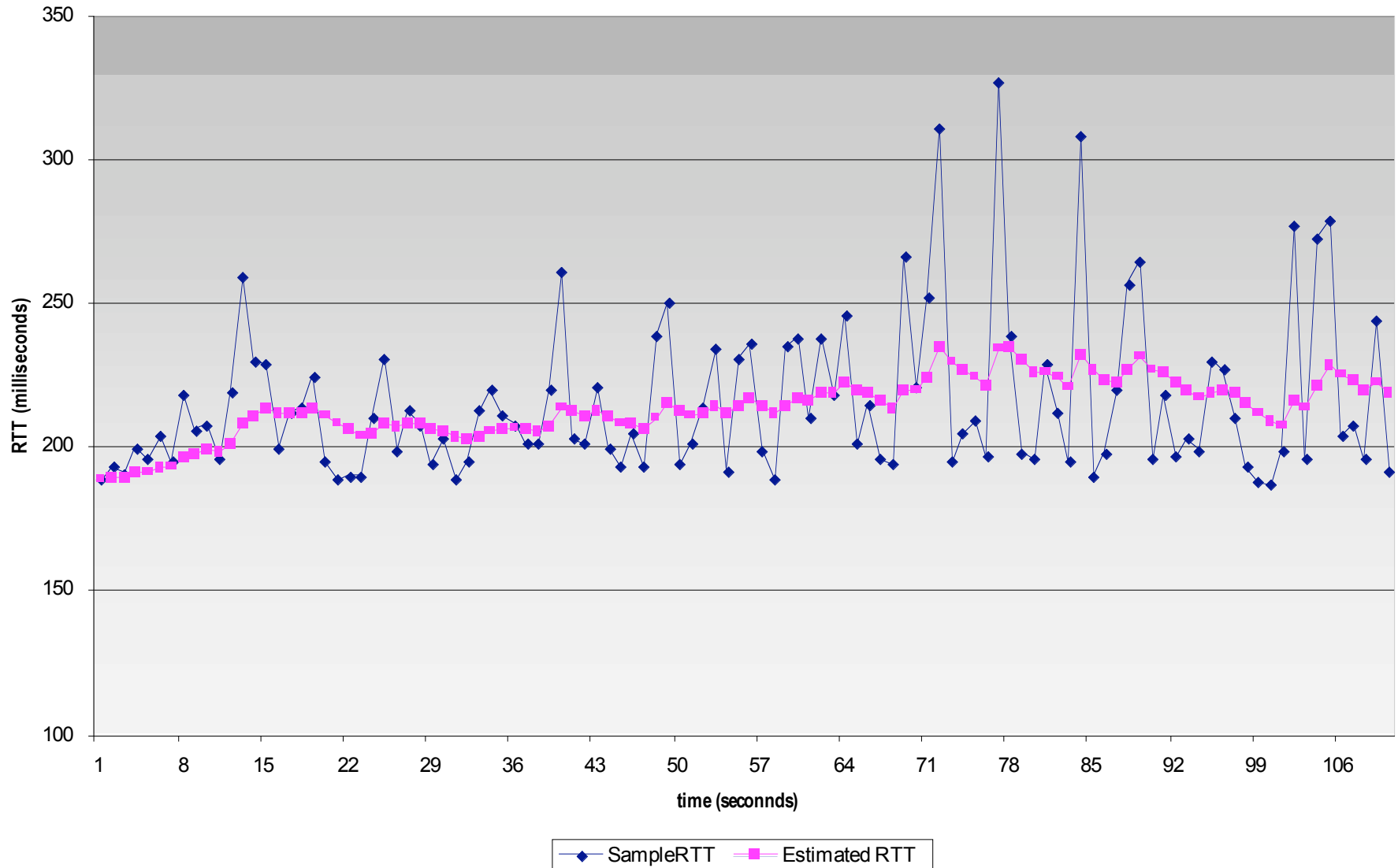  diff = 480 − 400 = 80
  SRTT = 400 + 10 = 410
  DevRTT = 200 + ¼ (80-200) = 170

**client**

Open conn.

S, seq=34

accept conn

400ms

A, ack=35; S, seq=93

Send 20 bytes of request

A, ack=94; seq=35, data

Send 100B Response, then close the conn.

480ms

A, ack=55; seq=94, data

A, ack=55; F, seq=194

...

A, ack=195; F, seq=55

Connection closed, wait 2MSL

A, ack=56

conn. Closed, remove TCB

# Example RTT estimation:

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

# TCP Fast  Retransmit

- RTO is set to a relatively long value
  - Aim at minimizing superfluous retransmission
  - long delay before resending lost packet

- Can detect lost segments via duplicate ACKs.
  - When a segment is lost, next arrival at receiver is out of order
  - When a segment arrives out of order, receiver can immediately sends an ACK indicating seq. # of next byte it is expecting

- When sender receives 3 duplicate ACKs for the same seq#(n), it assumes the segment with seq#(n) was lost

→fast retransmit: start retransmitting without waiting for the timer to expire
  - How many segments to retransmit? One only

# TCP fast retransmit example

# Yet another tweak of TCP: delayed ACK

*FYI*

- If a TCP connection carries traffic in both directions: ACKs are piggybacked on data segments

- For one-way data flow: If receiver sends an ACK after receiving everyone segment → double the packet count across the Internet

- **Delayed ACK**: after connection setup, upon receive one data segment $S_1$ :
  - wait a bit, see if next segment $S_2$ will arrive soon
  - If yes: sends an ACK for both
  - If no: send an ACK for $S_1$

*Does this delayed-ACK screw up RTT measurement? Maybe a little*

# TCP Receiver: when to send ACK?

| Event at TCP receiver | TCP Receiver action |
|---|---|
| in-order segment arrival, no gaps, everything earlier already ACKed | delayed ACK: wait up to 500ms, If nothing arrived, send ACK |
| in-order segment arrival, no gaps, one delayed ACK pending | immediately send one cumulative ACK |
| out-of-order arrival: higher-than-expect seq. #, gap detected | Immediately send ACK, indicating <u>seq. # of next expected byte</u> |
| arrival of segment that partially or completely fills a gap | immediate send ACK if segment starts at <u>the lower end of the gap</u> |

# Summary

- Connection management (SYN, FIN)

- Flow control for reliable delivery (sequence numbers, ACK)
  - ACK is a flag in the header; ACK flag == 0, ACK number in the header makes no sense (value ignored)

- Two-way communication
  - Separate sequence number management for both directions

- Error detection and recovery
  - Retransmission timer
  - Fast retransmit

- Receiver's flow control
  - Avoid overwhelming the receiver

- Congestion control
  - Avoid overwhelming the network

# After obtain a new RTT sample:

*important*

- difference = SampleRTT - SRTT

- SRTT' = (1-$\alpha$) x SRTT + $\alpha$ x SampleRTT

    = SRTT + $\alpha$ x difference

- DevRTT' = (1-$\beta$) x DevRTT + $\beta$ x |difference|

    = DevRTT + $\beta$ (|difference| - DevRTT)

- Retransmission Timer (RTO) = SRTT + 4 x DevRTT

Typically: $\alpha$ = 1/8, $\beta$ = 1/4