# Transmission Control Protocol (TCP)

Antonio Carzaniga

Faculty of Informatics
Università della Svizzera italiana

November 10, 2017

- Introduction to TCP

- Sequence numbers and acknowledgment numbers

- Timeouts and RTT estimation

- Reliable data transfer in TCP

- Connection management

# Transmission Control Protocol

- The Internet's primary transport protocol
  - defined in RFC 793, RFC 1122, RFC 1323, RFC 2018, and RFC 2581

# Transmission Control Protocol

- The Internet's primary transport protocol
  - defined in RFC 793, RFC 1122, RFC 1323, RFC 2018, and RFC 2581

- Connection-oriented service
  - endpoints "shake hands" to establish a connection
  - not a circuit-switched connection, nor a virtual circuit

# Transmission Control Protocol

- The Internet's primary transport protocol
  - defined in RFC 793, RFC 1122, RFC 1323, RFC 2018, and RFC 2581

- Connection-oriented service
  - endpoints "shake hands" to establish a connection
  - not a circuit-switched connection, nor a virtual circuit

- Full-duplex service
  - both endpoints can both send and receive, at the same time

- **_TCP segment:_** envelope for TCP data
  - ▸ TCP data are sent within TCP segments
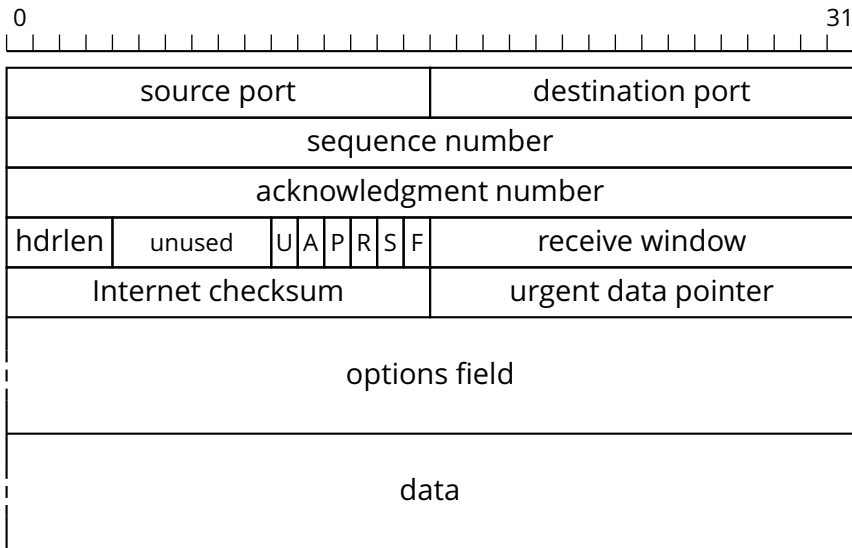  - ▸ TCP segments are usually sent within an IP packet

# Preliminary Definitions

- ***TCP segment:*** envelope for TCP data
  - ▸ TCP data are sent within TCP segments
  - ▸ TCP segments are usually sent within an IP packet

- ***Maximum segment size (MSS):*** maximum amount of application data transmitted in a single segment
  - ▸ typically related to the MTU of the connection, to avoid network-level fragmentation (we'll talk about all of this later)

- **■ *TCP segment:*** envelope for TCP data
  - ► TCP data are sent within TCP segments
  - ► TCP segments are usually sent within an IP packet

- **■ *Maximum segment size (MSS):*** maximum amount of application data transmitted in a single segment
  - ► typically related to the MTU of the connection, to avoid network-level fragmentation (we'll talk about all of this later)

- **■ *Maximum transmission unit (MTU):*** largest link-layer frame available to the sender host
  - ► *path MTU:* largest link-layer frame that can be sent on all links from the sender host to the receiver host

| 0 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | 31 |

| source port | destination port |
|---|---|
| sequence number | |
| acknowledgment number | |

| hdrlen | unused | U | A | P | R | S | F | receive window |
|---|---|---|---|---|---|---|---|---|

| Internet checksum | urgent data pointer |
|---|---|
| options field | |
| data | |

- *Source and destination ports:* (16-bit each) application identifiers

- *Source and destination ports:* (16-bit each) application identifiers

- *Sequence number:* (32-bit) used to implement reliable data transfer

- *Acknowledgment number:* (32-bit) used to implement reliable data transfer

- *Source and destination ports:* (16-bit each) application identifiers

- *Sequence number:* (32-bit) used to implement reliable data transfer

- *Acknowledgment number:* (32-bit) used to implement reliable data transfer

- *Receive window:* (16-bit) size of the "window" on the receiver end

# TCP Header Fields

- *Source and destination ports:* (16-bit each) application identifiers

- *Sequence number:* (32-bit) used to implement reliable data transfer

- *Acknowledgment number:* (32-bit) used to implement reliable data transfer

- *Receive window:* (16-bit) size of the "window" on the receiver end

- *Header length:* (4-bit) size of the TCP header in 32-bit words

## TCP Header Fields

- *Source and destination ports:* (16-bit each) application identifiers

- *Sequence number:* (32-bit) used to implement reliable data transfer

- *Acknowledgment number:* (32-bit) used to implement reliable data transfer

- *Receive window:* (16-bit) size of the "window" on the receiver end

- *Header length:* (4-bit) size of the TCP header in 32-bit words

- *Optional and variable-length options field:* may be used to negotiate protocol parameters

- *ACK flag:* (1-bit) signals that the value contained in the *acknowledgment number* represents a valid acknowledgment

- *ACK flag:* (1-bit) signals that the value contained in the *acknowledgment number* represents a valid acknowledgment

- *SYN flag:* (1-bit) used during connection setup and shutdown

- *ACK flag:* (1-bit) signals that the value contained in the *acknowledgment number* represents a valid acknowledgment

- *SYN flag:* (1-bit) used during connection setup and shutdown

- *RST flag:* (1-bit) used during connection setup and shutdown

- *ACK flag:* (1-bit) signals that the value contained in the *acknowledgment number* represents a valid acknowledgment

- *SYN flag:* (1-bit) used during connection setup and shutdown

- *RST flag:* (1-bit) used during connection setup and shutdown

- *FIN flag:* (1-bit) used during connection shutdown

- *ACK flag:* (1-bit) signals that the value contained in the *acknowledgment number* represents a valid acknowledgment

- *SYN flag:* (1-bit) used during connection setup and shutdown

- *RST flag:* (1-bit) used during connection setup and shutdown

- *FIN flag:* (1-bit) used during connection shutdown

- *PSH flag:* (1-bit) "push" flag, used to solicit the receiver to pass the data to the application immediately

- *ACK flag:* (1-bit) signals that the value contained in the *acknowledgment number* represents a valid acknowledgment

- *SYN flag:* (1-bit) used during connection setup and shutdown

- *RST flag:* (1-bit) used during connection setup and shutdown

- *FIN flag:* (1-bit) used during connection shutdown

- *PSH flag:* (1-bit) "push" flag, used to solicit the receiver to pass the data to the application immediately

- *URG flag:* (1-bit) "urgent" flag, used to inform the receiver that the sender has marked some data as "urgent". The location of this urgent data is marked by the *urgent data pointer* field

- *ACK flag:* (1-bit) signals that the value contained in the *acknowledgment number* represents a valid acknowledgment

- *SYN flag:* (1-bit) used during connection setup and shutdown

- *RST flag:* (1-bit) used during connection setup and shutdown

- *FIN flag:* (1-bit) used during connection shutdown

- *PSH flag:* (1-bit) "push" flag, used to solicit the receiver to pass the data to the application immediately

- *URG flag:* (1-bit) "urgent" flag, used to inform the receiver that the sender has marked some data as "urgent". The location of this urgent data is marked by the *urgent data pointer* field

- *Checksum:* (16-bit) used to detect transmission errors

- Sequence numbers are associated with *bytes* in the data stream
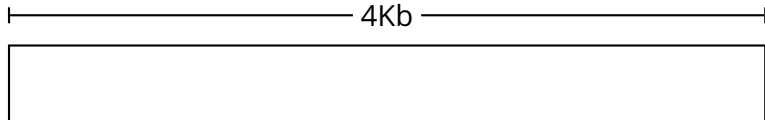  - not with segments, as we have used them before

- Sequence numbers are associated with *bytes* in the data stream
  - ▸ not with segments, as we have used them before

- The ***sequence number*** in a TCP segment indicates ***the sequence number of the first byte carried by that segment***

- Sequence numbers are associated with *bytes* in the data stream
  - not with segments, as we have used them before

- The *sequence number* in a TCP segment indicates *the sequence number of the first byte carried by that segment*
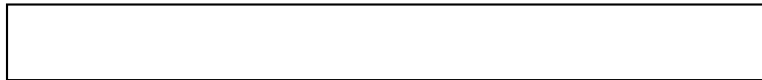
*application data stream*

├──────────────────── 4Kb ────────────────────┤

- Sequence numbers are associated with *bytes* in the data stream
  - ▸ not with segments, as we have used them before

- The *sequence number* in a TCP segment indicates *the sequence number of the first byte carried by that segment*
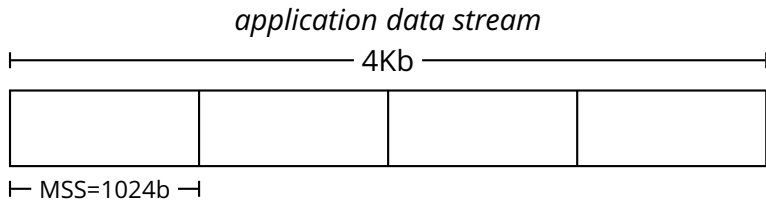
*application data stream*
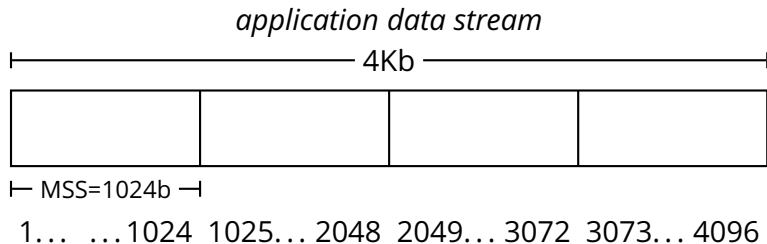
├───────────────── 4Kb ─────────────────┤

├─ MSS=1024b ─┤

- Sequence numbers are associated with *bytes* in the data stream
  - ▸ not with segments, as we have used them before

- The *sequence number* in a TCP segment indicates *the sequence number of the first byte carried by that segment*

*application data stream*

├────────────────────────── 4Kb ──────────────────────────┤

├─ MSS=1024b ─┤

- Sequence numbers are associated with *bytes* in the data stream
  - ▸ not with segments, as we have used them before

- The *sequence number* in a TCP segment indicates *the sequence number of the first byte carried by that segment*
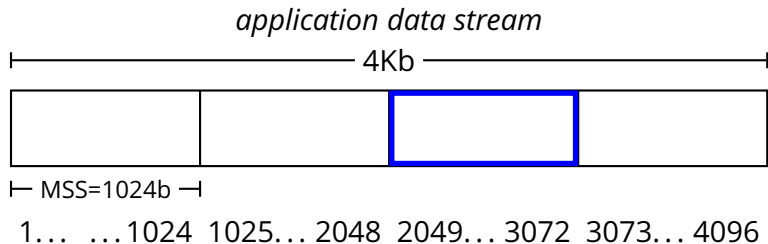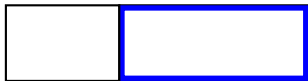
*application data stream*

⊢────────────────────── 4Kb ──────────────────────⊣

⊢ MSS=1024b ⊣

1... ...1024  1025... 2048  2049... 3072  3073... 4096

- Sequence numbers are associated with *bytes* in the data stream
  - ▶ not with segments, as we have used them before

- The *sequence number* in a TCP segment indicates *the sequence number of the first byte carried by that segment*

*application data stream*



├─────────────── 4Kb ───────────────┤

├─ MSS=1024b ─┤

1... ...1024  1025... 2048  2049... 3072  3073... 4096

*a TCP segment*

- Sequence numbers are associated with *bytes* in the data stream
  - ▸ not with segments, as we have used them before

- The **sequence number** in a TCP segment indicates **the sequence number of the first byte carried by that segment**
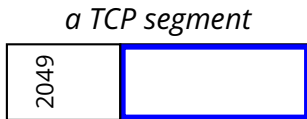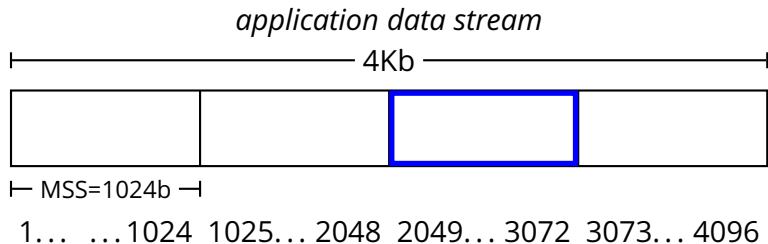
*application data stream*

├─────────────────── 4Kb ───────────────────┤

├─ MSS=1024b ─┤

1... ...1024  1025... 2048  2049... 3072  3073... 4096

*a TCP segment*

2049

- Sequence numbers are associated with *bytes* in the data stream
  - ▸ not with segments, as we have used them before

- The *sequence number* in a TCP segment indicates *the sequence number of the first byte carried by that segment*
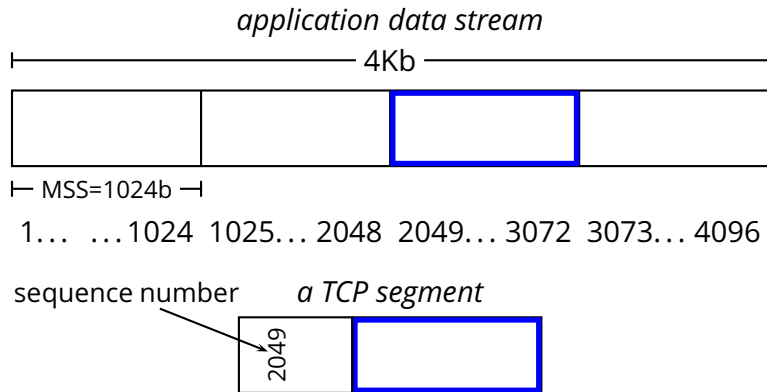
*application data stream*

├───────────────── 4Kb ─────────────────┤

├─ MSS=1024b ─┤

1... ...1024  1025... 2048  2049... 3072  3073... 4096

sequence number     *a TCP segment*

2049

- An *acknowledgment number* represents the *first sequence number not yet seen by the receiver*
  - TCP acknowledgments are *cumulative*

- An *acknowledgment number* represents the *first sequence number not yet seen by the receiver*
  - ► TCP acknowledgments are *cumulative*

A                                                                                          B
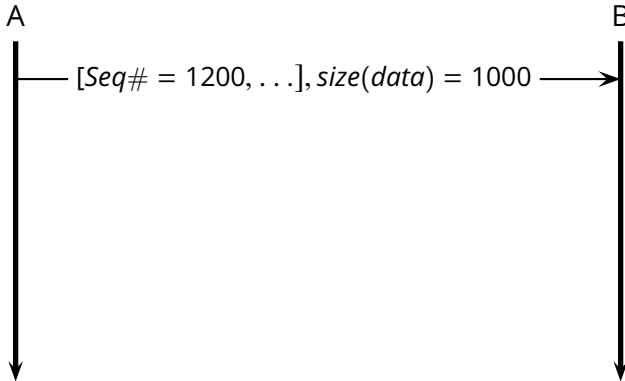
# Acknowledgment Numbers

- An *acknowledgment number* represents the *first sequence number not yet seen by the receiver*
  - TCP acknowledgments are *cumulative*

A                                                                        B

$[Seq\# = 1200, \ldots], size(data) = 1000$ →

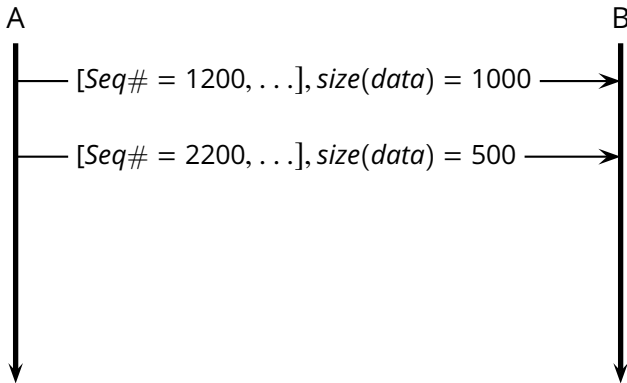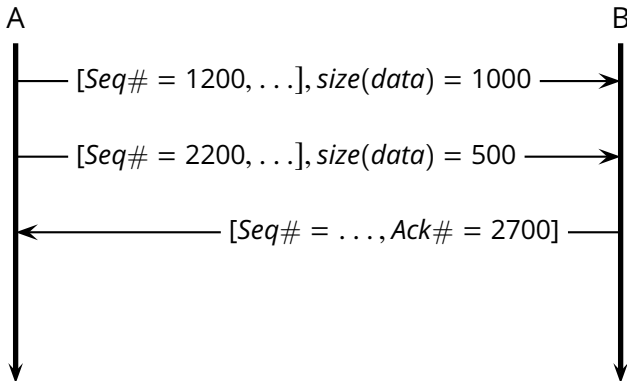# Acknowledgment Numbers

- An *acknowledgment number* represents the *first sequence number not yet seen by the receiver*
  - TCP acknowledgments are *cumulative*

A                                                                    B

$[Seq\# = 1200, \ldots], size(data) = 1000$

$[Seq\# = 2200, \ldots], size(data) = 500$

- An *acknowledgment number* represents the *first sequence number not yet seen by the receiver*
  - TCP acknowledgments are *cumulative*



A ............................................................ B

$[Seq\# = 1200, \ldots], size(data) = 1000$ →

$[Seq\# = 2200, \ldots], size(data) = 500$ →

← $[Seq\# = \ldots, Ack\# = 2700]$

# Sequence Numbers and ACK Numbers

## Sequence Numbers and ACK Numbers

- Notice that a TCP connection is a *full-duplex* link
  - therefore, there are ***two streams***
  - two different sequence numbers

- Notice that a TCP connection is a *full-duplex* link
  - ► therefore, there are ***two streams***
  - ► two different sequence numbers

E.g., consider a simple "Echo" application:

A                                                    B

- Notice that a TCP connection is a *full-duplex* link
  - ▸ therefore, there are ***two streams***
  - ▸ two different sequence numbers

E.g., consider a simple "Echo" application:

A                                                                                           B
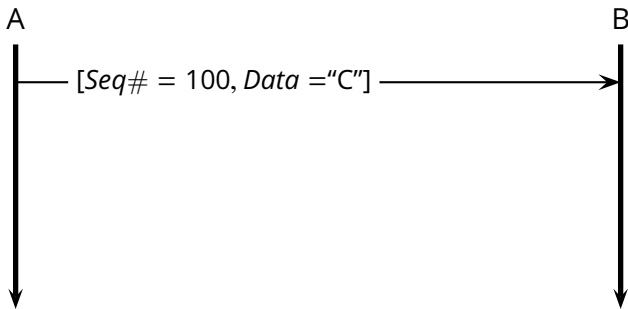
————— [*Seq#* = 100, *Data* ="C"] —————————→

# Sequence Numbers and ACK Numbers

- Notice that a TCP connection is a *full-duplex* link
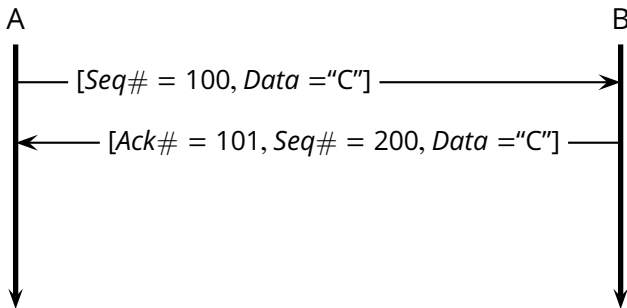  - therefore, there are ***two streams***
  - two different sequence numbers

E.g., consider a simple "Echo" application:

A                               B

$$\longrightarrow [Seq\# = 100, Data = \text{"C"}] \longrightarrow$$

$$\longleftarrow [Ack\# = 101, Seq\# = 200, Data = \text{"C"}] \longleftarrow$$

# **Sequence Numbers and ACK Numbers**

- Notice that a TCP connection is a *full-duplex* link
  - ▶ therefore, there are **two streams**
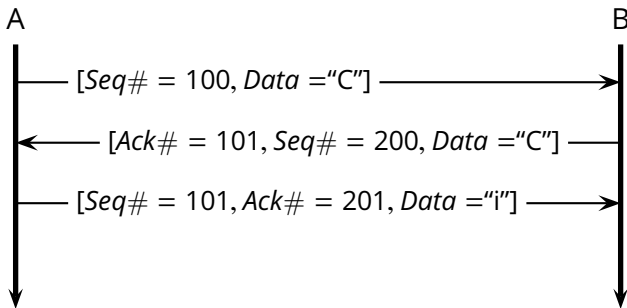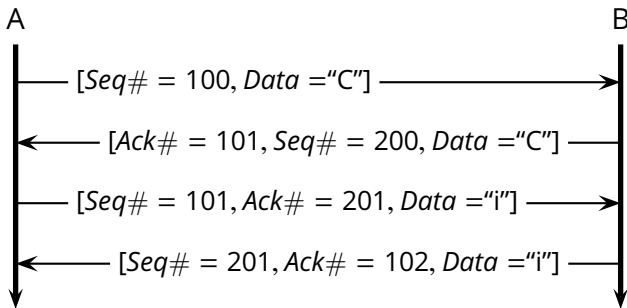  - ▶ two different sequence numbers

E.g., consider a simple "Echo" application:

```
A                                                              B
│                                                              │
│───── [Seq# = 100, Data ="C"] ──────────────────────────────▶│
│                                                              │
│◀───── [Ack# = 101, Seq# = 200, Data ="C"] ──────────────────│
│                                                              │
│───── [Seq# = 101, Ack# = 201, Data ="i"] ──────────────────▶│
│                                                              │
│                                                              │
▼                                                              ▼
```

# Sequence Numbers and ACK Numbers

- Notice that a TCP connection is a *full-duplex* link
  - ▶ therefore, there are **two streams**
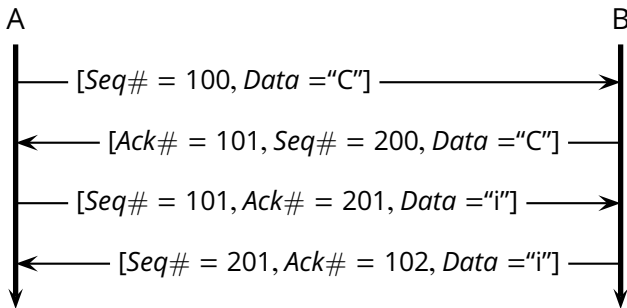  - ▶ two different sequence numbers

E.g., consider a simple "Echo" application:

A                                                                    B

── [*Seq#* = 100, *Data* ="C"] ────────────→

←──── [*Ack#* = 101, *Seq#* = 200, *Data* ="C"] ────

── [*Seq#* = 101, *Ack#* = 201, *Data* ="i"] ──────→

←──── [*Seq#* = 201, *Ack#* = 102, *Data* ="i"] ────

# Sequence Numbers and ACK Numbers

- Notice that a TCP connection is a *full-duplex* link
  - therefore, there are ***two streams***
  - two different sequence numbers

E.g., consider a simple "Echo" application:

```
A                                                          B
│─── [Seq# = 100, Data ="C"] ──────────────────────────────│
│                                                          │
│◄── [Ack# = 101, Seq# = 200, Data ="C"] ──────────────────│
│                                                          │
│─── [Seq# = 101, Ack# = 201, Data ="i"] ──────────────────►│
│                                                          │
│◄── [Seq# = 201, Ack# = 102, Data ="i"] ──────────────────│
▼                                                          ▼
```

- Acknowledgments are "piggybacked" on data segments

- TCP provides reliable data transfer using a *timer* to detect lost segments
  - timeout without an ACK → lost packet → retransmission

- TCP provides reliable data transfer using a *timer* to detect lost segments
  - timeout without an ACK → lost packet → retransmission

- How long to wait for acknowledgments?

- TCP provides reliable data transfer using a *timer* to detect lost segments
  - timeout without an ACK → lost packet → retransmission

- How long to wait for acknowledgments?

- Retransmission timeouts should be larger than the round-trip time $RTT = 2L$
  - as close as possible to the $RTT$

- TCP provides reliable data transfer using a *timer* to detect lost segments
  - timeout without an ACK → lost packet → retransmission

- How long to wait for acknowledgments?

- Retransmission timeouts should be larger than the round-trip time $RTT = 2L$
  - as close as possible to the *RTT*

- TCP controls its timeout by continuously *estimating the current RTT*

- RTT is measured using ACKs
  - only for packets transmitted once

- Given a single sample *S* at any given time

- *Exponential weighted moving average (EWMA)*

$$\overline{RTT} = (1 - \alpha)\overline{RTT}' + \alpha S$$

- RTT is measured using ACKs
  - only for packets transmitted once

- Given a single sample *S* at any given time

- *Exponential weighted moving average (EWMA)*

$$\overline{RTT} = (1 - \alpha)\overline{RTT}' + \alpha S$$

  - RFC 2988 recommends $\alpha = 0.125$

- RTT is measured using ACKs
  - only for packets transmitted once

- Given a single sample *S* at any given time

- *Exponential weighted moving average (EWMA)*

$$\overline{RTT} = (1 - \alpha)\overline{RTT}' + \alpha S$$

  - RFC 2988 recommends $\alpha = 0.125$

- TCP also measures the *variability of RTT*

$$\overline{DevRTT} = (1 - \beta)\overline{DevRTT}' + \beta|\overline{RTT}' - S|$$

# Round-Trip Time Estimation

- RTT is measured using ACKs
  - only for packets transmitted once

- Given a single sample *S* at any given time

- *Exponential weighted moving average (EWMA)*

$$\overline{RTT} = (1 - \alpha)\overline{RTT}' + \alpha S$$

  - RFC 2988 recommends $\alpha = 0.125$

- TCP also measures the *variability of RTT*

$$\overline{DevRTT} = (1 - \beta)\overline{DevRTT}' + \beta|\overline{RTT}' - S|$$

  - RFC 2988 recommends $\beta = 0.25$

- The timeout interval *T* must be larger than the RTT
  - so as to avoid unnecessary retransmission

- However, *T* should not be too far from RTT
  - so as to detect (and retransmit) lost segments as quickly as possible

- The timeout interval *T* must be larger than the RTT
  - so as to avoid unnecessary retransmission

- However, *T* should not be too far from RTT
  - so as to detect (and retransmit) lost segments as quickly as possible

- TCP sets its timeouts using the estimated RTT ($\overline{RTT}$) and the variability estimate $\overline{DevRTT}$:

$$T = \overline{RTT} + 4\overline{DevRTT}$$

A simplified TCP sender

- r_send(*data*)

  **if** (timer not running)
    start_timer()
  u_send([*data*,*next_seq_num*])
  *next_seq_num* ← *next_seq_num* + *length*(*data*)

A simplified TCP sender

- r_send(*data*)

  **if** (timer not running)
    start_timer()
  u_send([*data*,*next_seq_num*])
  *next_seq_num* ← *next_seq_num* + *length*(*data*)

- timeout

  u_send(pending segment with smallest sequence number)
  start_timer()

A simplified TCP sender

- r_send(*data*)

  **if** (timer not running)
     start_timer()
  u_send([*data*,*next_seq_num*])
  *next_seq_num* ← *next_seq_num* + *length*(*data*)

- timeout

  u_send(pending segment with smallest sequence number)
  start_timer()

- u_recv([ACK,*y*])

  **if** (*y* > *base*)
     *base* ← *y*
     **if** (there are pending segments)
        start_timer()
  **else** . . .

# Acknowledgment Generation (Receiver)

- Arrival of in-order segment with expected sequence number; all data up to expected sequence number already acknowledged

- Arrival of in-order segment with expected sequence number; all data up to expected sequence number already acknowledged
  - *Delayed ACK:* wait 500ms for another in-order segment; If that does not arrive, send ACK

# Acknowledgment Generation (Receiver)

- Arrival of in-order segment with expected sequence number; all data up to expected sequence number already acknowledged
  - *Delayed ACK:* wait 500ms for another in-order segment; If that does not arrive, send ACK

- Arrival of in-order segment with expected sequence number. One other in-order segment waiting for ACK (see above)

# Acknowledgment Generation (Receiver)

- Arrival of in-order segment with expected sequence number; all data up to expected sequence number already acknowledged
  - *Delayed ACK:* wait 500ms for another in-order segment; If that does not arrive, send ACK

- Arrival of in-order segment with expected sequence number. One other in-order segment waiting for ACK (see above)
  - *Cumulative ACK:* immediately send cumulative ACK (for both segments)

# Acknowledgment Generation (Receiver)

- Arrival of in-order segment with expected sequence number; all data up to expected sequence number already acknowledged
  - *Delayed ACK:* wait 500ms for another in-order segment; If that does not arrive, send ACK

- Arrival of in-order segment with expected sequence number. One other in-order segment waiting for ACK (see above)
  - *Cumulative ACK:* immediately send cumulative ACK (for both segments)

- Arrival of out of order segment with higher-than-expected sequence number (gap detected)

# **Acknowledgment Generation (Receiver)**

- Arrival of in-order segment with expected sequence number; all data up to expected sequence number already acknowledged
  - *Delayed ACK:* wait 500ms for another in-order segment; If that does not arrive, send ACK

- Arrival of in-order segment with expected sequence number. One other in-order segment waiting for ACK (see above)
  - *Cumulative ACK:* immediately send cumulative ACK (for both segments)

- Arrival of out of order segment with higher-than-expected sequence number (gap detected)
  - *Duplicate ACK:* immediately send duplicate ACK

# Acknowledgment Generation (Receiver)

- Arrival of in-order segment with expected sequence number; all data up to expected sequence number already acknowledged
  - *Delayed ACK:* wait 500ms for another in-order segment; If that does not arrive, send ACK

- Arrival of in-order segment with expected sequence number. One other in-order segment waiting for ACK (see above)
  - *Cumulative ACK:* immediately send cumulative ACK (for both segments)

- Arrival of out of order segment with higher-than-expected sequence number (gap detected)
  - *Duplicate ACK:* immediately send duplicate ACK

- Arrival of segment that (partially or completely) fills a gap in the received data

- Arrival of in-order segment with expected sequence number; all data up to expected sequence number already acknowledged
  - *Delayed ACK:* wait 500ms for another in-order segment; If that does not arrive, send ACK

- Arrival of in-order segment with expected sequence number. One other in-order segment waiting for ACK (see above)
  - *Cumulative ACK:* immediately send cumulative ACK (for both segments)

- Arrival of out of order segment with higher-than-expected sequence number (gap detected)
  - *Duplicate ACK:* immediately send duplicate ACK

- Arrival of segment that (partially or completely) fills a gap in the received data
  - *Immediate ACK:* immediately send ACK if the packet start at the lower end of the gap

- u_recv([ACK,*y*])

  **if** (*y* > *base*)
     *base* ← *y*
     **if** (there are pending segments)
        start_timer()

- u_recv([ACK,*y*])

  **if** (*y* > *base*)
    *base* ← *y*
    **if** (there are pending segments)
      start_timer()
  **else**
    *ack_counter*[*y*] ← *ack_counter*[*y*] + 1
    **if** (*ack_counter*[*y*] = 3)
      u_send(segment with sequence number *y*)

Three-way handshake

Three-way handshake

client                                                              server

Three-way handshake

client                                                              server

[*SYN*, *Seq#* = *cli_init_seq*] ─────────────────────────→

Three-way handshake



client                              server

$[SYN, Seq\# = cli\_init\_seq]$

$[SYN, ACK, Ack\# = cli\_init\_seq + 1, Seq\# = srv\_init\_seq]$

Three-way handshake

client                                                                          server

$[SYN, Seq\# = cli\_init\_seq]$ ⟶

⟵ $[SYN, ACK, Ack\# = cli\_init\_seq + 1, Seq\# = srv\_init\_seq]$

$[ACK, Seq\# = cli\_init\_seq + 1, Ack\# = srv\_init\_seq + 1]$ ⟶

"This is it."
"Okay, Bye now."
"Bye."

"This is it."
"Okay, Bye now."
"Bye."

client                                              server

"This is it."
"Okay, Bye now."
"Bye."

client                                                              server

|— [*FIN*] ————————————————————————→|

"This is it."
"Okay, Bye now."
"Bye."

"This is it."
"Okay, Bye now."
"Bye."

```
        client                                              server

          ├ [FIN] ─────────────────────────────────────▶
          
          ◀─────────────────────────────────── [ACK] ┤
          
          ◀─────────────────────────────────── [FIN] ┤
          
          ▼                                              ▼
```

"This is it."
"Okay, Bye now."
"Bye."

CLOSED

CLOSED

application
opens connection
send SYN

SYN_SENT

CLOSED

application
opens connection
send SYN

SYN_SENT

receive SYN,ACK
send ACK

ESTABLISHED

CLOSED

application
opens connection
send SYN

SYN_SENT

receive SYN,ACK
send ACK

ESTABLISHED

application
closes connection
send FIN

FIN_WAIT_1

application
opens connection
send SYN

CLOSED

SYN_SENT

receive SYN,ACK
send ACK

FIN_WAIT_2

ESTABLISHED

receive ACK

FIN_WAIT_1

application
closes connection
send FIN

CLOSED

application
opens connection
send SYN

SYN_SENT

receive SYN,ACK
send ACK

ESTABLISHED

application
closes connection
send FIN

FIN_WAIT_1

receive ACK

FIN_WAIT_2

receive FIN
send ACK

TIME_WAIT

CLOSED

CLOSED

application
opens server socket

LISTEN

# The TCP State Machine (Server)

CLOSED

application
opens server socket

LISTEN

receive SYN
send SYN,ACK

SYN_RCVD

```
   CLOSED                 application
      │              opens server socket
      │              ──────────────────
      ▼
   LISTEN
      │
      │           receive SYN
      │           ──────────
      ▼           send SYN,ACK
  SYN_RCVD
      │
      │           receive ACK
      │           ──────────
      ▼
 ESTABLISHED
```