# CS450  Computer Networks

The slides used in class are derived from the slides available on
our text book companion website:
http://wps.pearsoned.com/ecs_kurose_compnetw_6/
copyright 1996-2012     J.F Kurose and K.W. Ross

# CS450  Computer Networks
# Lesson 10
# Transport Layer – TCP

# The organizing power of pure consciousness

# Lesson 10 -TCP

Our Goal:  Understand TCP.  A Connection-oriented
transport protocol

- segment structure
- reliable data transfer
- flow control
- connection management

# TCP: Overview   RFCs: 793, 1122, 1323, 2018, 2581

❖ **point-to-point:**
  - one sender, one receiver

❖ **reliable, in-order *byte steam:***
  - no "message boundaries"

❖ **pipelined:**
  - TCP congestion and flow control set window size

❖ **full duplex data:**
  - bi-directional data flow in same connection
  - MSS: maximum segment size
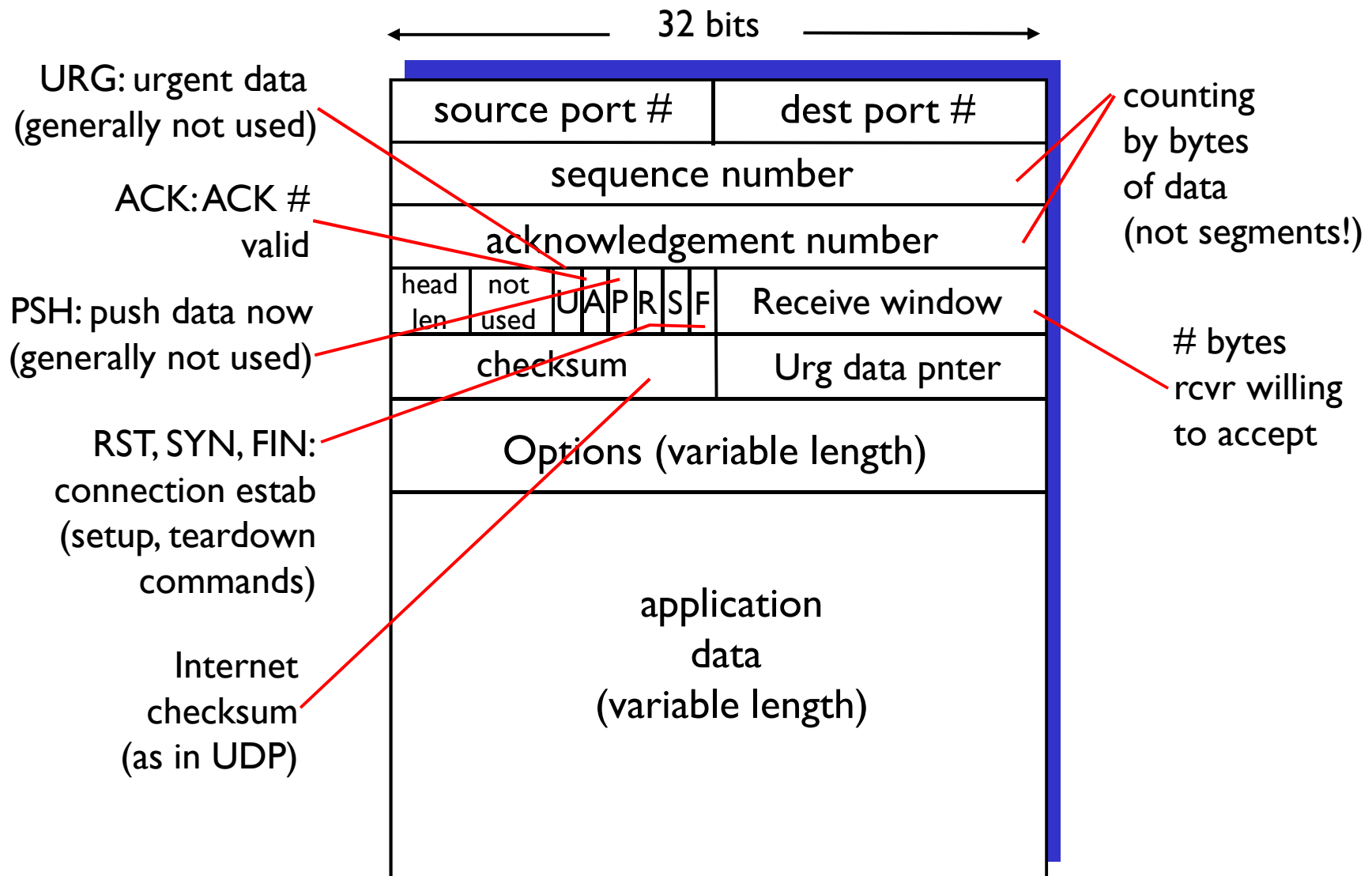
❖ **connection-oriented:**
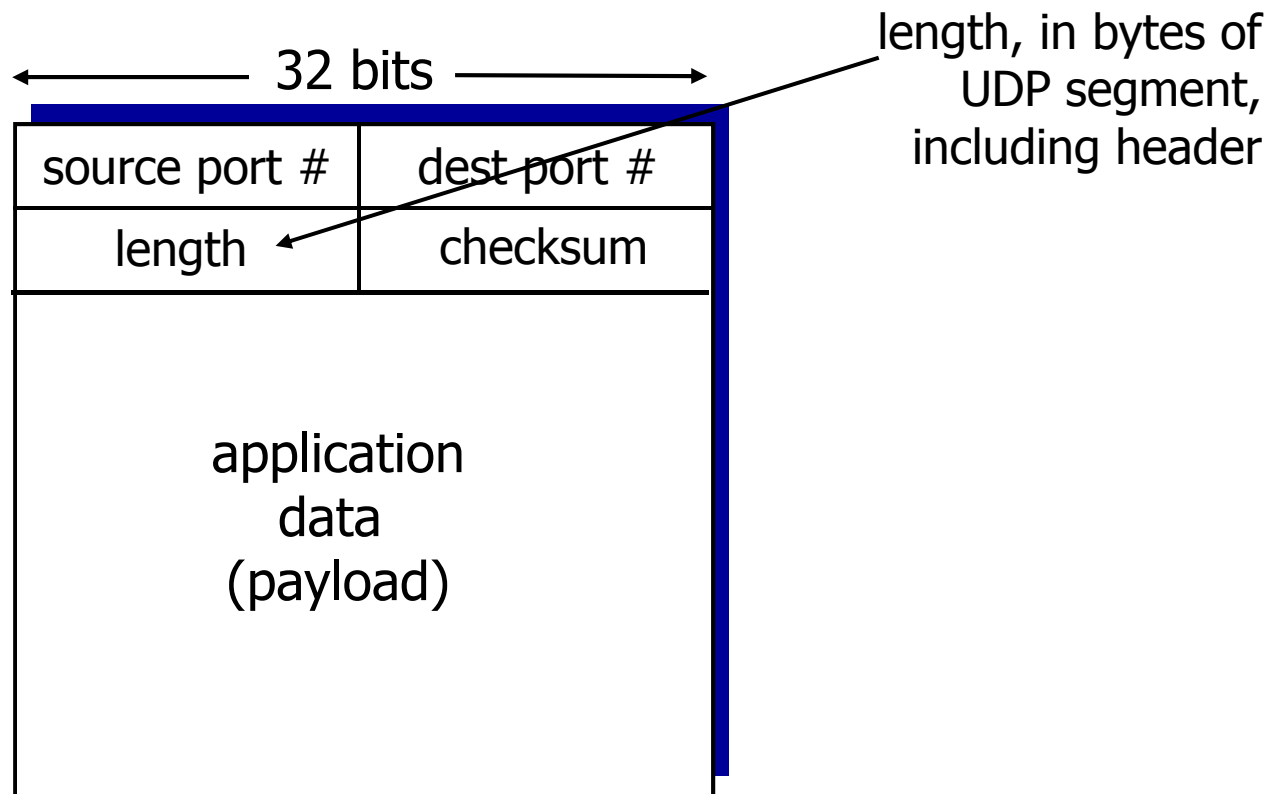  - handshaking (exchange of control msgs) inits sender, receiver state before data exchange

❖ **flow controlled:**
  - sender will not overwhelm receiver

# TCP segment structure

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)

32 bits

| source port # | dest port # |
|---|---|
| sequence number ||
| acknowledgement number ||

| head len | not used | U | A | P | R | S | F | Receive window |
|---|---|---|---|---|---|---|---|---|
| checksum ||||||| Urg data pnter |

Options (variable length)

application
data
(variable length)

counting
by bytes
of data
(not segments!)

# bytes
rcvr willing
to accept

# UDP: segment header – much simpler!



32 bits

| source port # | dest port # |
|---|---|
| length | checksum |

application
data
(payload)

length, in bytes of
UDP segment,
including header

UDP segment format

# TCP seq. numbers, ACKs

sequence numbers:
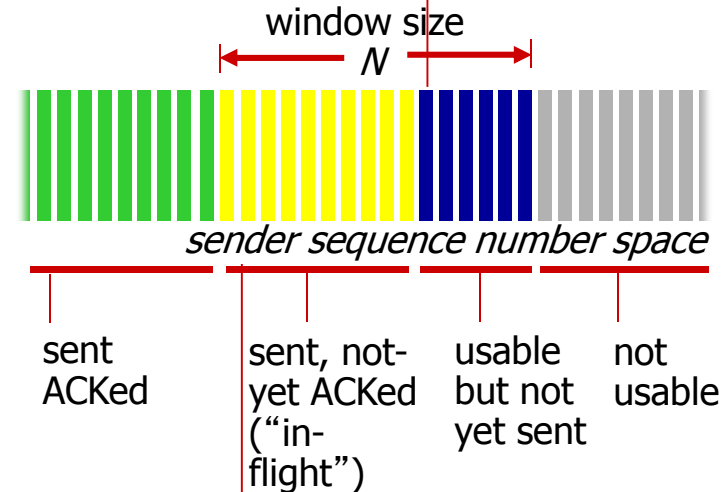- byte stream "number" of first byte in segment's data

acknowledgements:
- seq # of next byte expected from other side
- cumulative ACK

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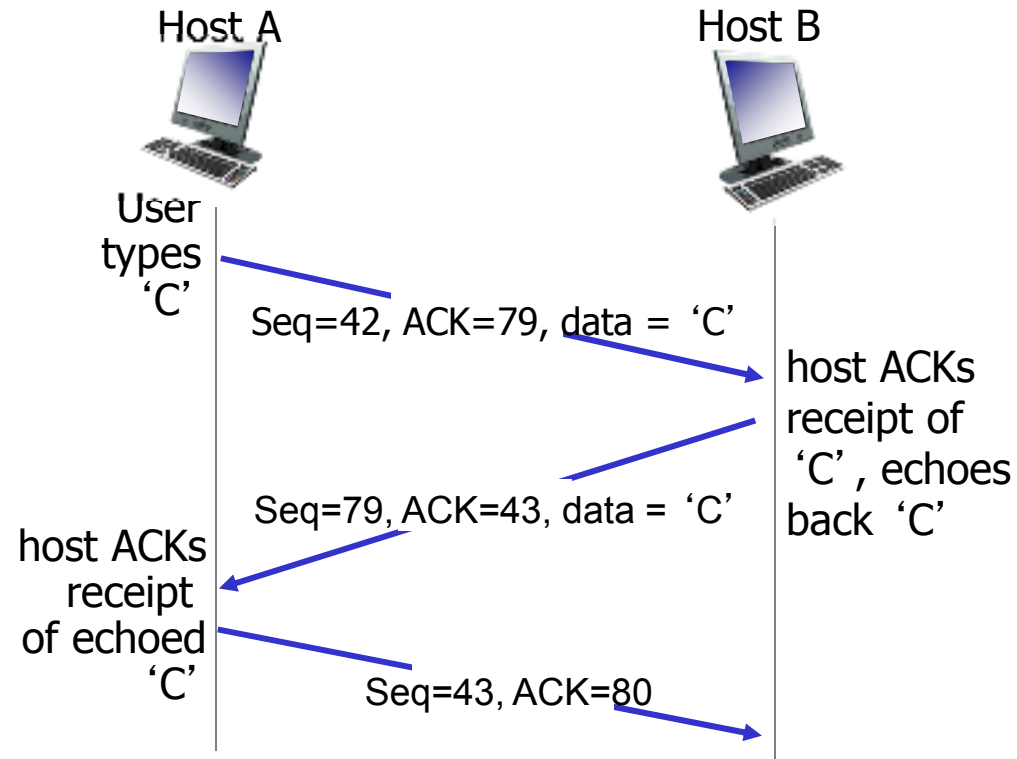- A: TCP spec doesn't say, - up to implementor

outgoing segment from sender

| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |
| rwnd | |
| checksum | urg pointer |

window size
N

sender sequence number space

| sent ACKed | sent, not-yet ACKed ("in-flight") | usable but not yet sent | not usable |

incoming segment to sender

| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |
| A | rwnd |
| checksum | urg pointer |

# TCP seq. numbers, ACKs

Host A                                        Host B

User
types
'C'

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

                     host ACKs
                     receipt of
    Seq=79, ACK=43, data = 'C'  'C', echoes
                     back 'C'

host ACKs
receipt
of echoed
'C'

    Seq=43, ACK=80

simple telnet scenario

# TCP Round Trip Time and Timeout

Q: how to set TCP
  timeout value?

# TCP Round Trip Time and Timeout

**Q:** how to set TCP timeout value?

❖ longer than RTT

  ▪ but RTT varies

❖ too short: premature timeout

  ▪ unnecessary retransmissions

❖ too long: slow reaction to segment loss

**Q:** how to estimate RTT?

# TCP Round Trip Time and Timeout

**Q:** how to set TCP timeout value?

* longer than RTT
  * but RTT varies
* too short: premature timeout
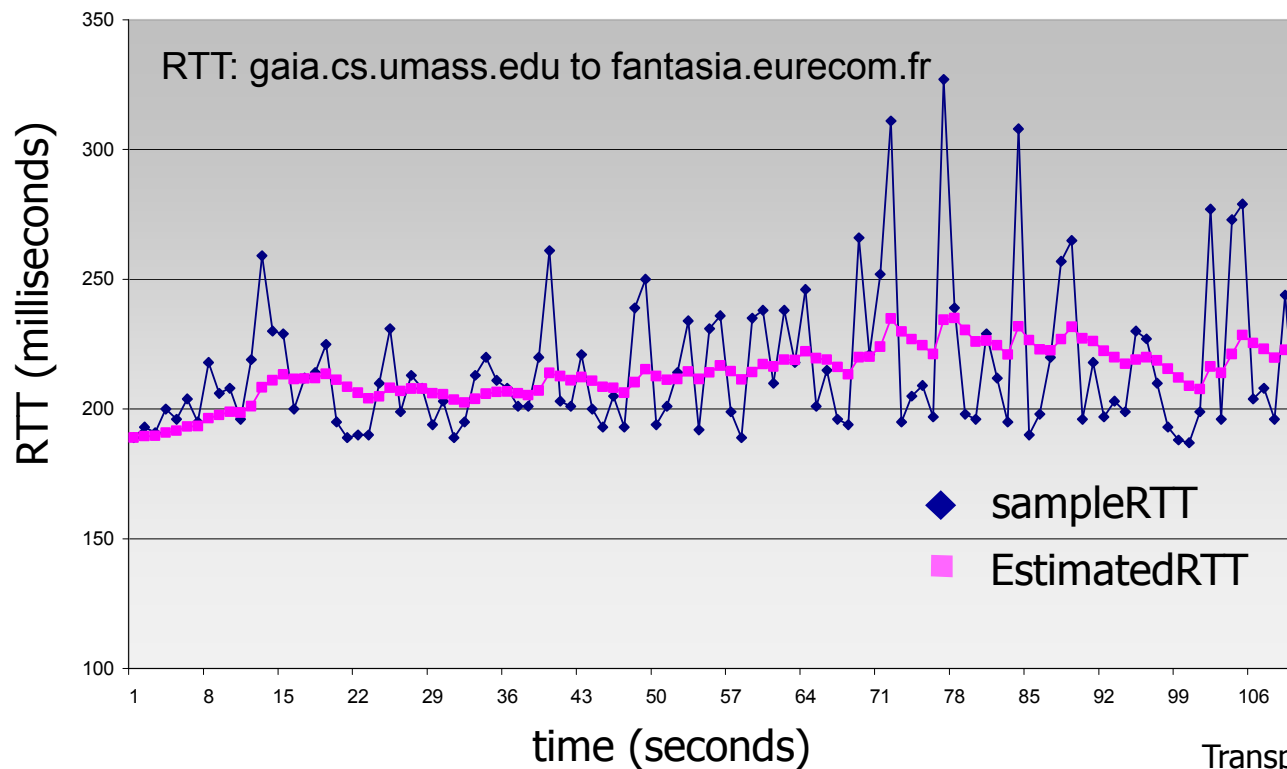  * unnecessary retransmissions
* too long: slow reaction to segment loss

**Q:** how to estimate RTT?

* `SampleRTT`: measured time from segment transmission until ACK receipt
  * ignore retransmissions
* `SampleRTT` will vary, want estimated RTT "smoother"
  * average several recent measurements, not just current `SampleRTT`

# TCP round trip time, timeout

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

- ❖ exponential weighted moving average
- ❖ influence of past sample decreases exponentially fast
- ❖ typical value: $\alpha = 0.125$



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

◆ sampleRTT
■ EstimatedRTT

RTT (milliseconds) — time (seconds)

# TCP round trip time, timeout

❖ timeout interval: **EstimatedRTT** plus "safety margin"
  - large variation in **EstimatedRTT** -> larger safety margin

❖ estimate SampleRTT deviation from EstimatedRTT:

$$\texttt{DevRTT = (1-}\beta\texttt{)*DevRTT +}$$
$$\beta\texttt{*|SampleRTT-EstimatedRTT|}$$
$$\texttt{(typically, }\beta\texttt{ = 0.25)}$$

**TimeoutInterval = EstimatedRTT + 4*DevRTT**

estimated RTT          "safety margin"

# 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

❖ Let's consider simplified TCP sender:

- ignore duplicate acks
- ignore flow control, congestion control

# 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)
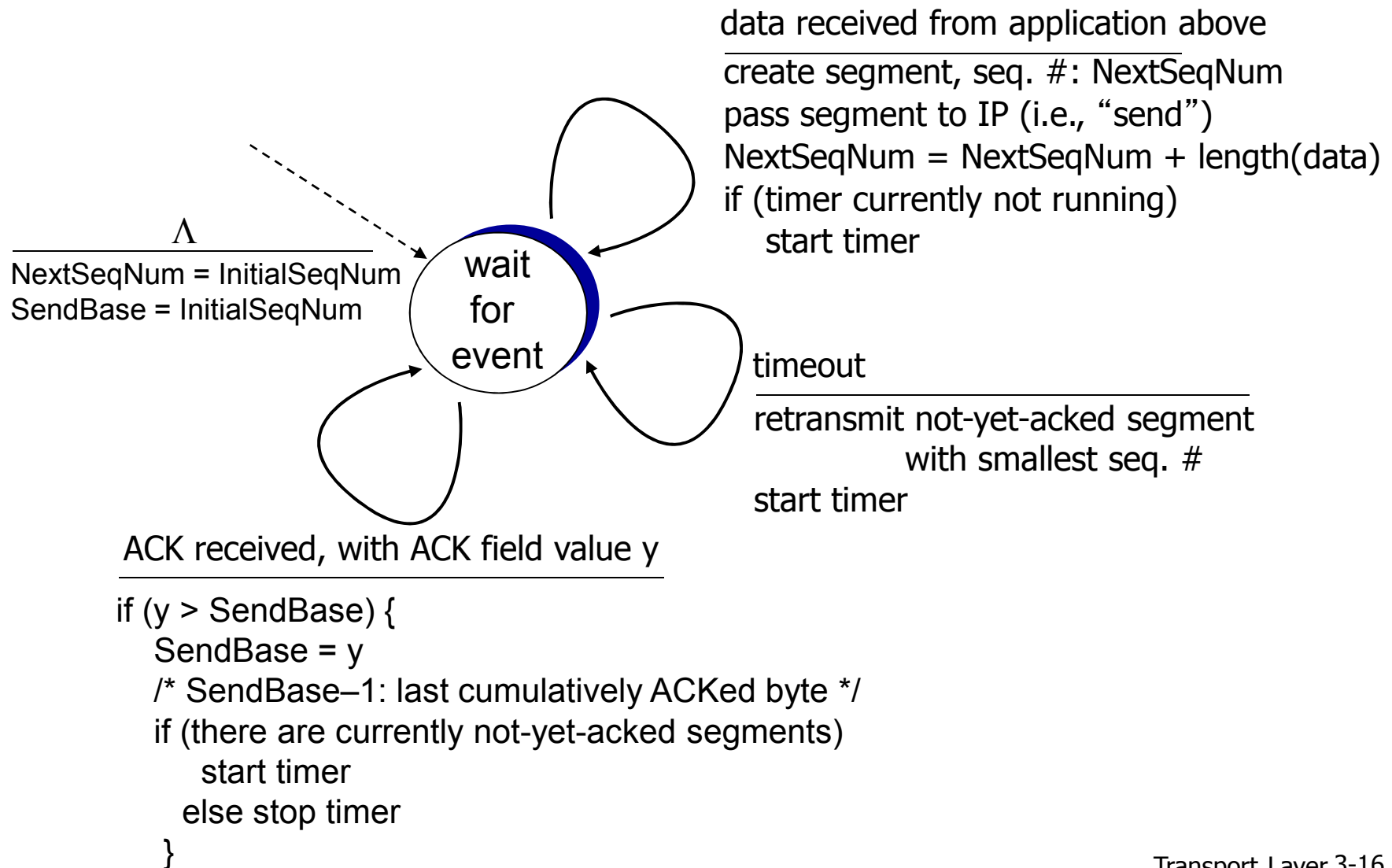
❖ expiration interval: `TimeOutInterval`

**timeout:**

❖ retransmit segment that caused timeout

❖ restart timer

**Ack rcvd:**

❖ If acknowledges previously unacked segments
  ▪ update what is known to be acked
  ▪ start timer if there are outstanding segments

# TCP sender (simplified)

data received from application above
_____
create segment, seq. #: NextSeqNum
pass segment to IP (i.e., "send")
NextSeqNum = NextSeqNum + length(data)
if (timer currently not running)
   start timer

$\Lambda$
_____
NextSeqNum = InitialSeqNum
SendBase = InitialSeqNum

**wait for event**

timeout
_____
retransmit not-yet-acked segment
               with smallest seq. #
start timer

ACK received, with ACK field value y
_____
if (y > SendBase) {
   SendBase = y
   /* SendBase–1: last cumulatively ACKed byte */
   if (there are currently not-yet-acked segments)
      start timer
    else stop timer
  }

```
NextSeqNum = InitialSeqNum
SendBase = InitialSeqNum

loop (forever) {
  switch(event)

  event: data received from application above
      create TCP segment with sequence number NextSeqNum
      if (timer currently not running)
          start timer
      pass segment to IP
      NextSeqNum = NextSeqNum + length(data)

  event: timer timeout
      retransmit not-yet-acknowledged segment with
          smallest sequence number
      start timer

  event: ACK received, with ACK field value of y
      if (y > SendBase) {
          SendBase = y
          if (there are currently not-yet-acknowledged segments)
              start timer
      }

} /* end of loop forever */
```

# TCP sender (simplified)
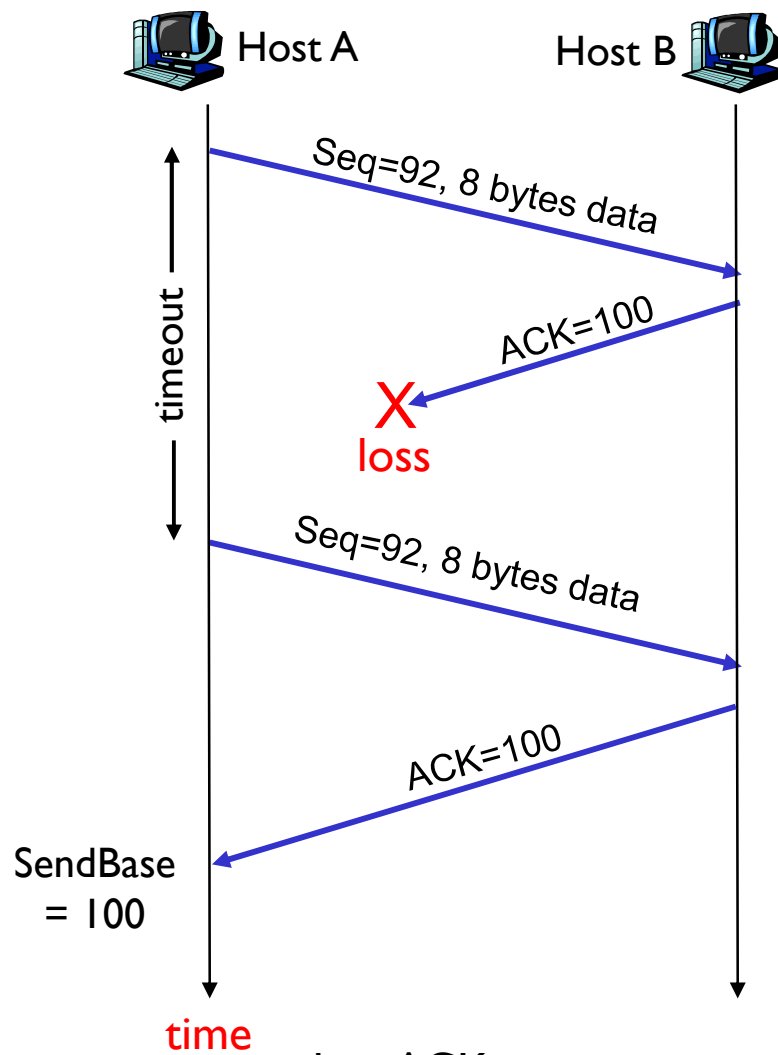
Comment:
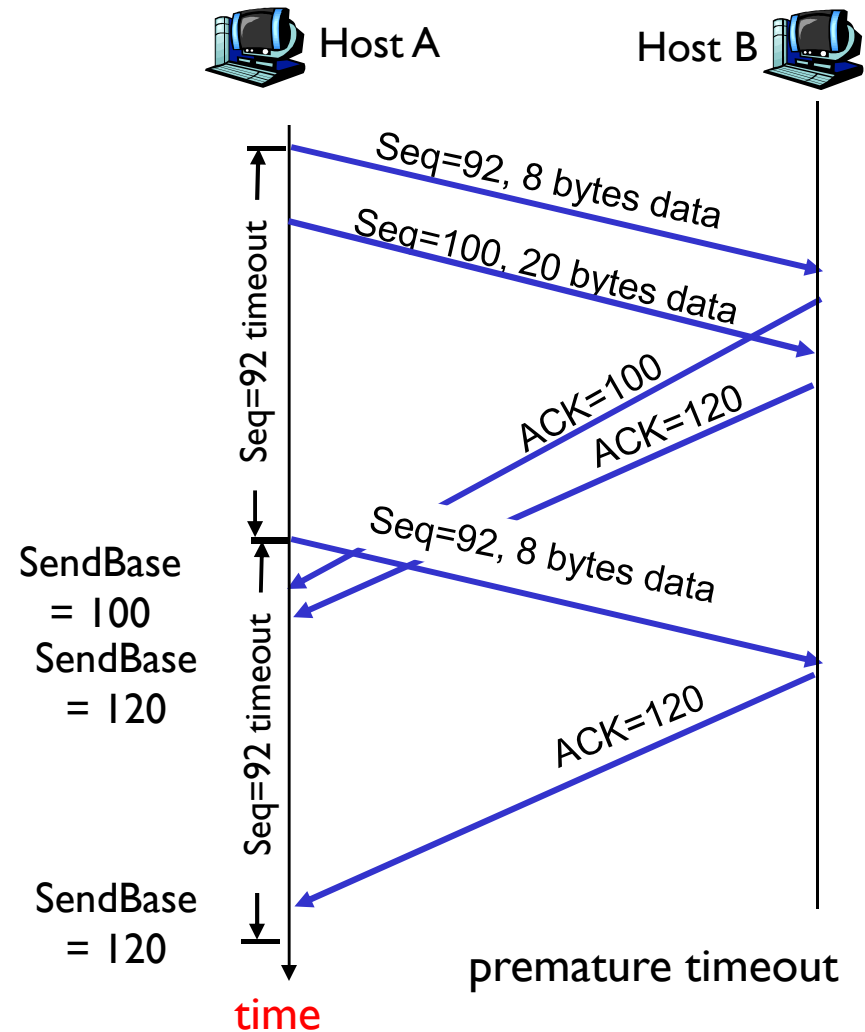• SendBase-1: last cumulatively acked byte
Example:
• SendBase-1 = 71; y= 73, so the rcvr wants 73+ ; y > SendBase, so that new data is acked
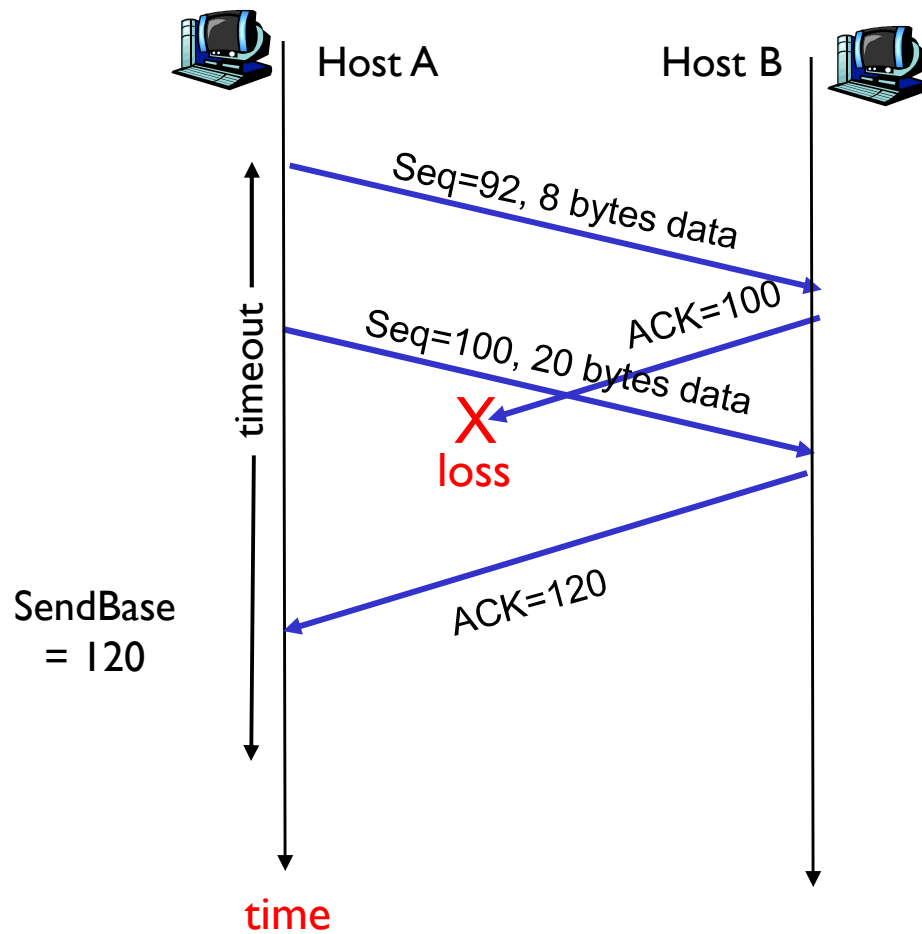
# TCP: retransmission scenarios



lost ACK scenario

premature timeout

# TCP retransmission scenarios (more)



Cumulative ACK scenario

# TCP ACK generation [RFC 1122, RFC 2581]

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

# TCP 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.

> *TCP fast retransmit*
>
> if sender receives 3 ACKs for same data ("triple duplicate ACKs"),
>
> resend unacked segment with smallest seq #
>   - likely that unacked segment lost, so don't wait for timeout

# TCP fast retransmit



Host A                                                    Host B

Seq=92, 8 bytes of data
Seq=100, 20 bytes of data
X

ACK=100

ACK=100

ACK=100

ACK=100

Seq=100, 20 bytes of data

timeout

fast retransmit after sender
receipt of triple duplicate ACK

# 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)
                start 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
            }
        }
```

a duplicate ACK for
already ACKed segment

fast retransmit

# TCP reliable data transfer

❖ TCP reliable data transfer – a GBN or an SR protocol?

❖ Discuss and present your analysis to class.

# TCP reliable data transfer - Sender

❖ Like GBN:

- The tcp sender keeps the smallest sequence number of a transmitted byte that is unack'ed == SendBase.

- The tcp sender keeps the sequence number of the next byte to be sent == NextSeqNum.

- Like SR:

- when timer fires the sender only transmits a single segment -- The one that is missing it's ack.

# TCP reliable data transfer - Receiver

❖ Like SR:

- Receiver can buffer out-of-order packets.

- TCP implementations **_may_** provide
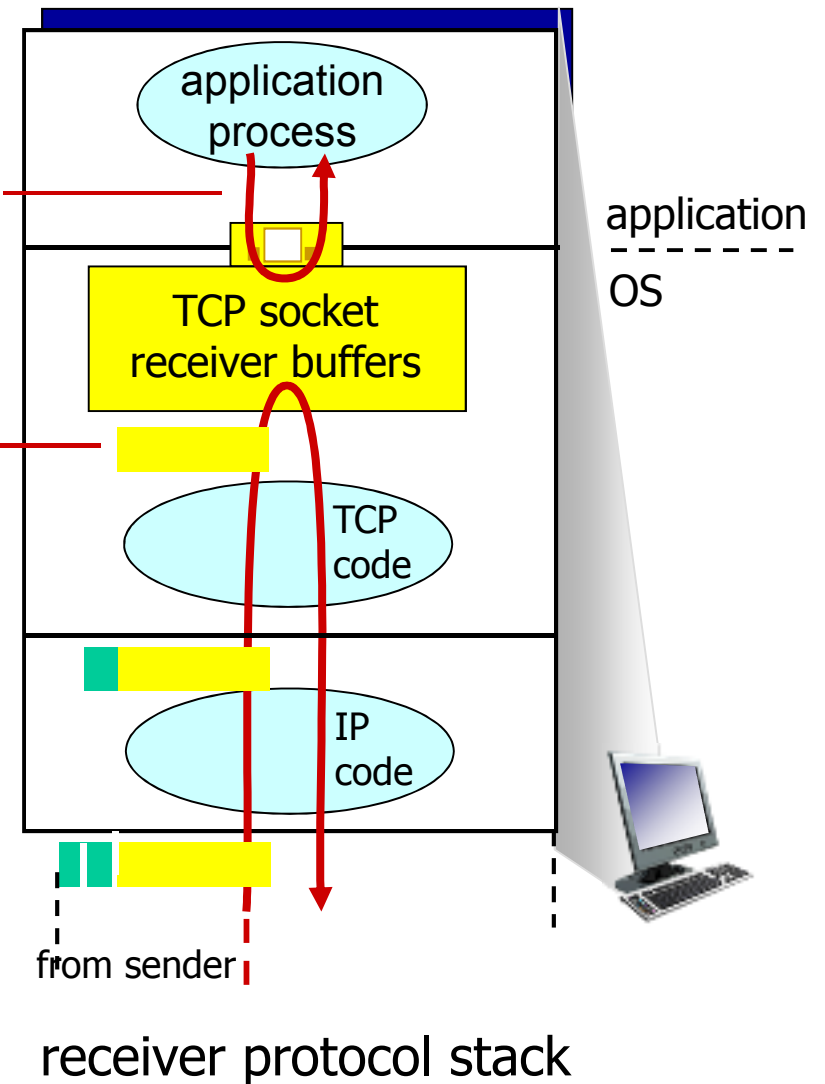
    **_selective acknowledgement_**

as opposed to cumulatively ack'ing the last correctly received, in-order segmen**t**

# TCP flow control

application may
remove data from
TCP socket buffers ....

... slower than TCP
receiver is delivering
(sender is sending)

application
process

TCP socket
receiver buffers

TCP
code

IP
code

application
OS

*flow control*

receiver controls sender, so
sender won't overflow
receiver's buffer by transmitting
too much, too fast

from sender

receiver protocol stack

# TCP flow control

❖ receiver "advertises" free buffer space by including **rwnd** value in TCP header of receiver-to-sender segments

  ▪ **RcvBuffer** size set via socket options (typical default is 4096 bytes)

  ▪ many operating systems autoadjust **RcvBuffer**

❖ sender limits amount of unacked ("in-flight") data to receiver's **rwnd** value

❖ guarantees receive buffer will not overflow

*to application process*

RcvBuffer

buffered data

rwnd

free buffer space

*TCP segment payloads*

*receiver-side buffering*

# Connection Management

before exchanging data, sender/receiver "handshake":

- ❖ agree to establish connection (each knowing the other willing to establish connection)
- ❖ agree on connection parameters

| application | application |
|---|---|
| connection state: ESTAB<br>connection variables:<br>    seq # client-to-server<br>        server-to-client<br>    **rcvBuffer** size<br>      at server,client | connection state: ESTAB<br>connection Variables:<br>    seq # client-to-server<br>        server-to-client<br>    **rcvBuffer** size<br>      at server,client |
| network | network |

```
Socket clientSocket =
  newSocket("hostname","port
  number");
```

```
Socket connectionSocket =
  welcomeSocket.accept();
```

# Agreeing to establish a connection

2-way handshake:



choose x

req_conn(x)

ESTAB

acc_conn(x)

ESTAB

*Q:* will 2-way handshake always work in network?

# Agreeing to establish a connection

## 2-way handshake:



*Q:* will 2-way handshake always work in network?

❖ variable delays
❖ retransmitted messages (e.g. req_conn(x)) due to message loss
❖ message reordering
❖ can't "see" other side

# Agreeing to establish a connection

2-way handshake failure scenarios:



choose x

req_conn(x)

ESTAB

retransmit
req_conn(x)

acc_conn(x)

ESTAB

req_conn(x)

connection
x completes

client
terminates

server
forgets x

ESTAB

half open connection!
(no client!)

choose x

req_conn(x)

ESTAB

retransmit
req_conn(x)

acc_conn(x)

ESTAB

data(x+1)

accept
data(x+1)

retransmit
data(x+1)

connection
x completes

client
terminates

server
forgets x

req_conn(x)

data(x+1)

ESTAB
accept
data(x+1)

# TCP 3-way handshake

*client state*

LISTEN

SYNSENT

ESTAB

*server state*

LISTEN

SYN RCVD

ESTAB

choose init seq num, x
send TCP SYN msg

SYNbit=1, Seq=x

choose init seq num, y
send TCP SYNACK
msg, acking SYN

SYNbit=1, Seq=y
ACKbit=1; ACKnum=x+1

received SYNACK(x)
indicates server is live;
send ACK for SYNACK;
this segment may contain
client-to-server data

ACKbit=1, ACKnum=y+1

received ACK(y)
indicates client is live

# TCP 3-way handshake: FSM

```
                              ┌─────────┐
                              │ closed  │
                              └─────────┘
      Socket connectionSocket =
         welcomeSocket.accept();
         ─────────────────────                    Socket clientSocket =
                  Λ                                  newSocket("hostname","port
                                                       number");
         SYN(x)                                     ────────────────────────
         ──────────────────
      SYNACK(seq=y,ACKnum=x+1)       ┌────────┐      SYN(seq=x)
         create new socket for       │ listen │
      communication back to client   └────────┘

   ┌────────┐                                    ┌────────┐
   │  SYN   │                                    │  SYN   │
   │  rcvd  │                                    │  sent  │
   └────────┘                                    └────────┘

                        ┌────────┐      SYNACK(seq=y,ACKnum=x+1)
                        │ ESTAB  │      ─────────────────────────
      ACK(ACKnum=y+1)   └────────┘         ACK(ACKnum=y+1)
      ──────────────
            Λ
```

# TCP: closing a connection

- ❖ client, server each close their side of connection
    - ▪ send TCP segment with FIN bit = 1
- ❖ respond to received FIN with ACK
    - ▪ on receiving FIN, ACK can be combined with own FIN
- ❖ simultaneous FIN exchanges can be handled

# TCP: closing a connection

client state

*client state*

ESTAB

```
clientSocket.close()
```

FIN_WAIT_1    can no longer
              send but can
              receive data

FIN_WAIT_2    wait for server
              close

TIMED_WAIT

              timed wait
              for 2*max
              segment lifetime

CLOSED

FINbit=1, seq=x

ACKbit=1; ACKnum=x+1

can still
send data

FINbit=1, seq=y

can no longer
send data

ACKbit=1; ACKnum=y+1

*server state*

ESTAB

CLOSE_WAIT

LAST_ACK

CLOSED

# Lesson 10: Summary

❖ TCP is an excellent example of the instantiation of the principles of transport layer services:

- multiplexing, demultiplexing

- reliable data transfer – a hybrid GBN and SR approach

- flow control

- congestion control  - to be covered in next lesson.