# Chapter 3
# Transport Layer

*Computer Networking: A Top Down Approach*
6th edition
Jim Kurose, Keith Ross
Addison-Wesley
March 2012

# Chapter 3: Transport Layer

## our goals:

❖ understand principles behind transport layer services:
- multiplexing, demultiplexing
- reliable data transfer
- flow control
- congestion control

❖ learn about Internet transport layer protocols:
- UDP: connectionless transport
- TCP: connection-oriented reliable transport
- TCP congestion control

# Chapter 3 outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer

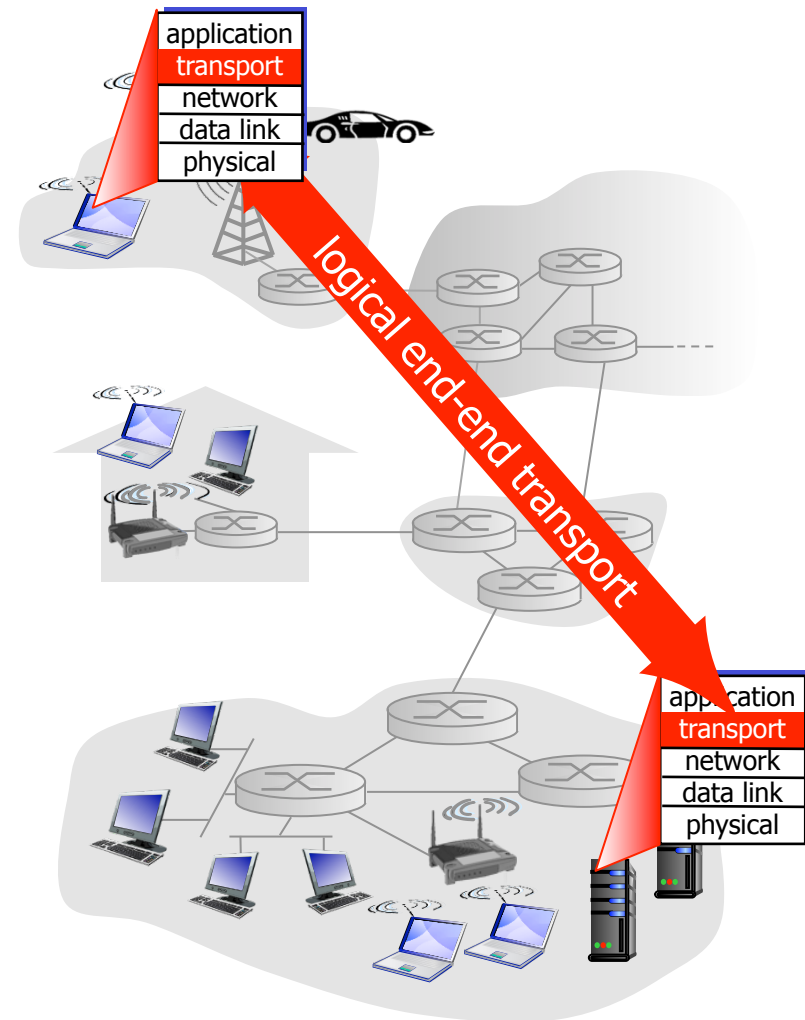3.5 connection-oriented transport: TCP

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

3.6 principles of congestion control

3.7 TCP congestion control

# Transport services and protocols

❖ provide *logical communication* between app processes running on different hosts

❖ transport protocols run in end systems

  ▪ send side: breaks app messages into *segments*, passes to network layer

  ▪ rcv side: reassembles segments into messages, passes to app layer

❖ more than one transport protocol available to apps

  ▪ Internet: TCP and UDP



application
transport
network
data link
physical

logical end-end transport

application
transport
network
data link
physical

# Transport vs. network layer

❖ *network layer:* logical communication between hosts

❖ *transport layer:* logical communication between processes
  - relies on, enhances, network layer services

---

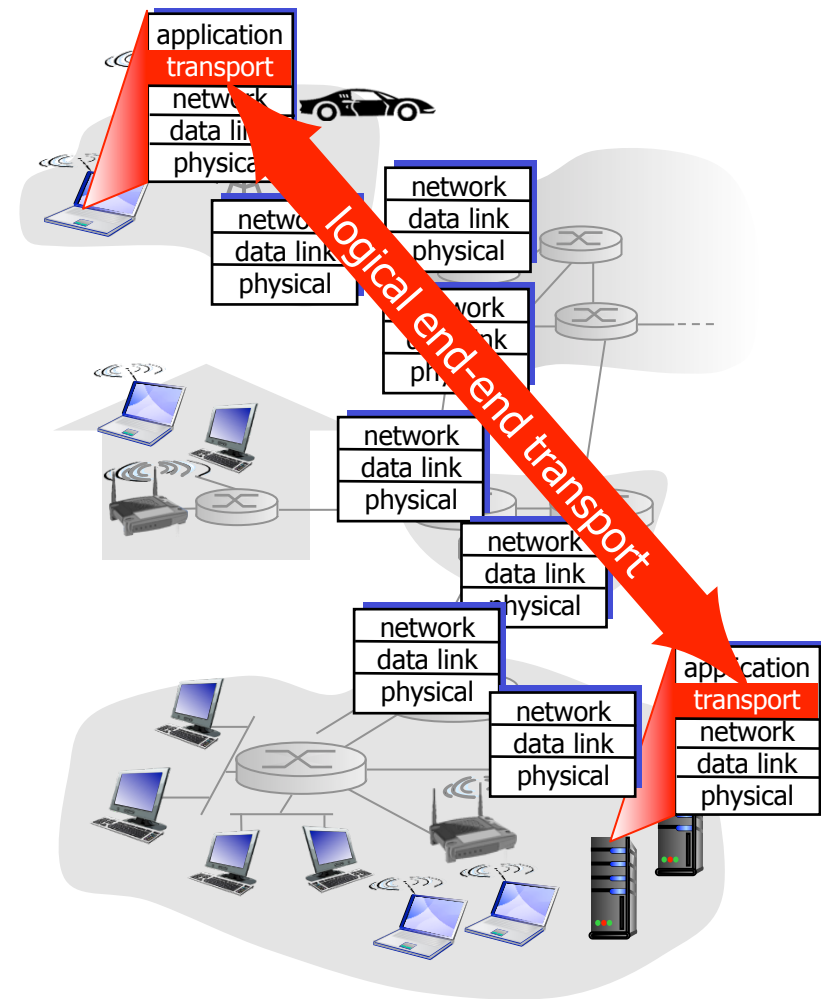*household analogy:*

12 kids in Ann's house sending letters to 12 kids in Bill's house:

❖ hosts = houses

❖ processes = kids

❖ app messages = letters in envelopes

❖ transport protocol = Ann and Bill who demux to in-house siblings

❖ network-layer protocol = postal service

# Internet transport-layer protocols

- ❖ **reliable, in-order delivery (TCP)**
  - congestion control
  - flow control
  - connection setup
- ❖ **unreliable, unordered delivery: UDP**
  - no-frills extension of "best-effort" IP
- ❖ **services not available:**
  - delay guarantees
  - bandwidth guarantees

# Chapter 3 outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer

3.5 connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management

3.6 principles of congestion control

3.7 TCP congestion control
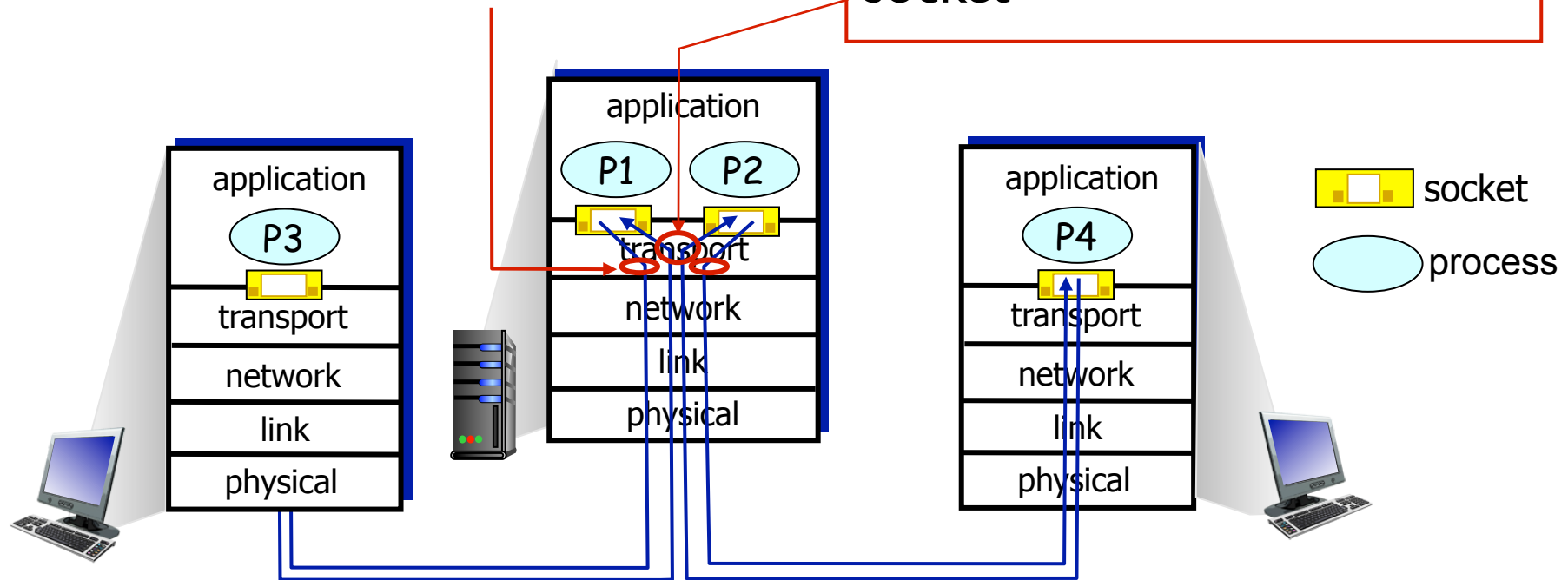
# Multiplexing/demultiplexing

*multiplexing at sender:*
handle data from multiple sockets, add transport header (later used for demultiplexing)

*demultiplexing at receiver:*
use header info to deliver received segments to correct socket

# How demultiplexing works

❖ host receives IP datagrams
  - each datagram has source IP address, destination IP address
  - each datagram carries one transport-layer segment
  - each segment has source, destination port number

❖ host uses *IP addresses & port numbers* to direct segment to appropriate socket

← 32 bits →

| source port # | dest port # |
|---|---|
| other header fields | |
| application data (payload) | |

TCP/UDP segment format

# Connectionless demultiplexing

❖ *recall:* created socket has host-local port #:
```
DatagramSocket mySocket1
= new DatagramSocket(12534);
```

❖ *recall:* when creating datagram to send into UDP socket, must specify
- destination IP address
- destination port #

---

❖ when host receives UDP segment:
- checks destination port # in segment
- directs UDP segment to socket with that port #

➡ IP datagrams with *same dest. port #,* but different source IP addresses and/or source port numbers will be directed to *same socket* at dest

# Connectionless demux: example

```
DatagramSocket
 serverSocket = new
 DatagramSocket
  (6428);
```

```
DatagramSocket
 mySocket2 = new
 DatagramSocket
  (9157);
```

```
DatagramSocket
 mySocket1 = new
 DatagramSocket
  (5775);
```

application

P3

transport

network

link

physical

application

P1

transport

network

link

physical

application

P4

transport

network

link

physical

source port: 6428
dest port: 9157

source port: ?
dest port: ?

source port: 9157
dest port: 6428

source port: ?
dest port: ?

# Connection-oriented demux

- ❖ TCP socket identified by 4-tuple:
  - ▪ source IP address
  - ▪ source port number
  - ▪ dest IP address
  - ▪ dest port number
- ❖ demux: receiver uses all four values to direct segment to appropriate socket

- ❖ server host may support many simultaneous TCP sockets:
  - ▪ each socket identified by its own 4-tuple
- ❖ web servers have different sockets for each connecting client
  - ▪ non-persistent HTTP will have different socket for each request

# Connection-oriented demux: example



application

P3

transport

network

link

physical

host: IP
address A

application

P4    P5    P6

transport

network

link

physical

server: IP
address B

application

P2    P3

transport

network

link

physical

host: IP
address C

source IP,port: B,80
dest IP,port: A,9157

source IP,port: A,9157
dest IP, port: B,80

source IP,port: C,5775
dest IP,port: B,80

source IP,port: C,9157
dest IP,port: B,80

three segments, all destined to IP address: B,
dest port: 80 are demultiplexed to *different* sockets

# Connection-oriented demux: example

threaded server



application

P4

transport

network

link

physical

server: IP
address B

application

P3

transport

network

link

physical

host: IP
address A

application

P2      P3

transport

network

link

physical

host: IP
address C

source IP,port: B,80
dest IP,port: A,9157

source IP,port: A,9157
dest IP, port: B,80

source IP,port: C,5775
dest IP,port: B,80

source IP,port: C,9157
dest IP,port: B,80

# Chapter 3 outline

# UDP: User Datagram Protocol [RFC 768]

* ❖ "no frills," "bare bones" Internet transport protocol
* ❖ "best effort" service, UDP segments may be:
  * ▪ lost
  * ▪ delivered out-of-order to app
* ❖ *connectionless:*
  * ▪ no handshaking between UDP sender, receiver
  * ▪ each UDP segment handled independently of others

* ❖ UDP use:
  * ▪ streaming multimedia apps (loss tolerant, rate sensitive)
  * ▪ DNS
  * ▪ SNMP
* ❖ reliable transfer over UDP:
  * ▪ add reliability at application layer
  * ▪ application-specific error recovery!

# UDP: segment header

32 bits

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

application
data
(payload)

length, in bytes of UDP segment, including header

UDP segment format

## why is there a UDP?

❖ no connection establishment (which can add delay)

❖ simple: no connection state at sender, receiver

❖ small header size

❖ no congestion control: UDP can blast away as fast as desired

# UDP checksum

*Goal:* detect "errors" (e.g., flipped bits) in transmitted segment

## sender:

- ❖ treat segment contents, including header fields, as sequence of 16-bit integers
- ❖ checksum: addition (one's complement sum) of segment contents
- ❖ sender puts checksum value into UDP checksum field

## receiver:

- ❖ compute checksum of received segment
- ❖ check if computed checksum equals checksum field value:
  - NO - error detected
  - YES - no error detected. *But maybe errors nonetheless?* More later ….

# Checksum Calculation

At the sender:

| 1087 | 13 |
|------|-----|
| 15 | |
| Application Data (Payload) | |

```
00000100  00111111  ──────►  1087
00000000  00001101  ──────►  13
00000000  00001111  ──────►  15

0000 0100 0101 1011        → SUM
```

1st compliment:
**1111 1011 1010 0100** → CHECKSUM
= FBA$_H$

At the receiver:

| 1087 | 13 |
|------|-----|
| 15 | FBA4$_H$ |
| Application Data (Payload) | |

```
    0000 0100 0011 1111 → Source Port
+  0000 0000 0000 1101 →    Dest  Port
+  0000 0000 0000 1111 →         Length
+  1111 1011 1010 0100 →      Checksum
   1111 1111 1111 1111     All 1's
                          No Error
```

# Internet checksum:
# When CARRYOUT occurs

example: add two 16-bit integers

```
                1 1 1 0 0 1 1 0 0 1 1 0 0 1 1 0
                1 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1
            _____
wraparound ⓵  1 0 1 1 1 0 1 1 1 0 1 1 1 0 1 1
            _____
      sum     1 0 1 1 1 0 1 1 1 0 1 1 1 1 0 0
 checksum     0 1 0 0 0 1 0 0 0 1 0 0 0 0 1 1
```

*Note:* when adding numbers, a carryout from the most significant bit needs to be added to the result

# Chapter 3 outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer

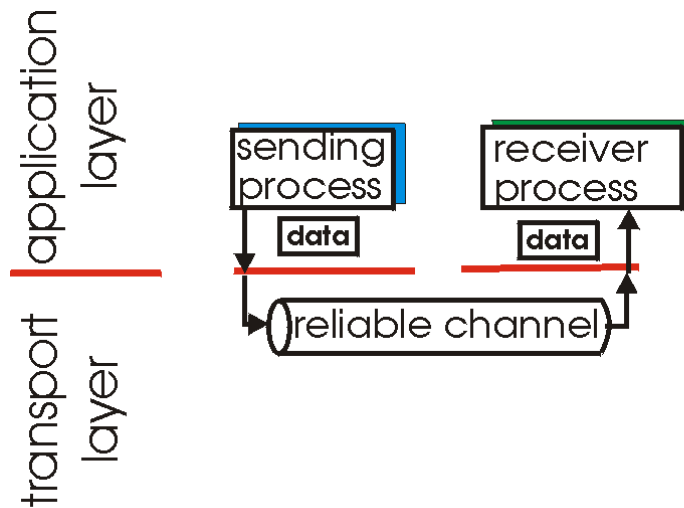3.5 connection-oriented transport: TCP
- segment structure
- reliable data transfer
- flow control
- connection management

3.6 principles of congestion control

3.7 TCP congestion control

# Principles of reliable data transfer

❖ **important in application, transport, link layers**
  ▪ top-10 list of important networking topics!



(a) provided service

❖ **characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)**

# Principles of reliable data transfer

- ❖ **important in application, transport, link layers**
  - ▪ top-10 list of important networking topics!



(a) provided service

(b) service implementation

- ❖ **characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)**

# Principles of reliable data transfer

❖ **important in application, transport, link layers**
  ▪ **top-10 list of important networking topics!**

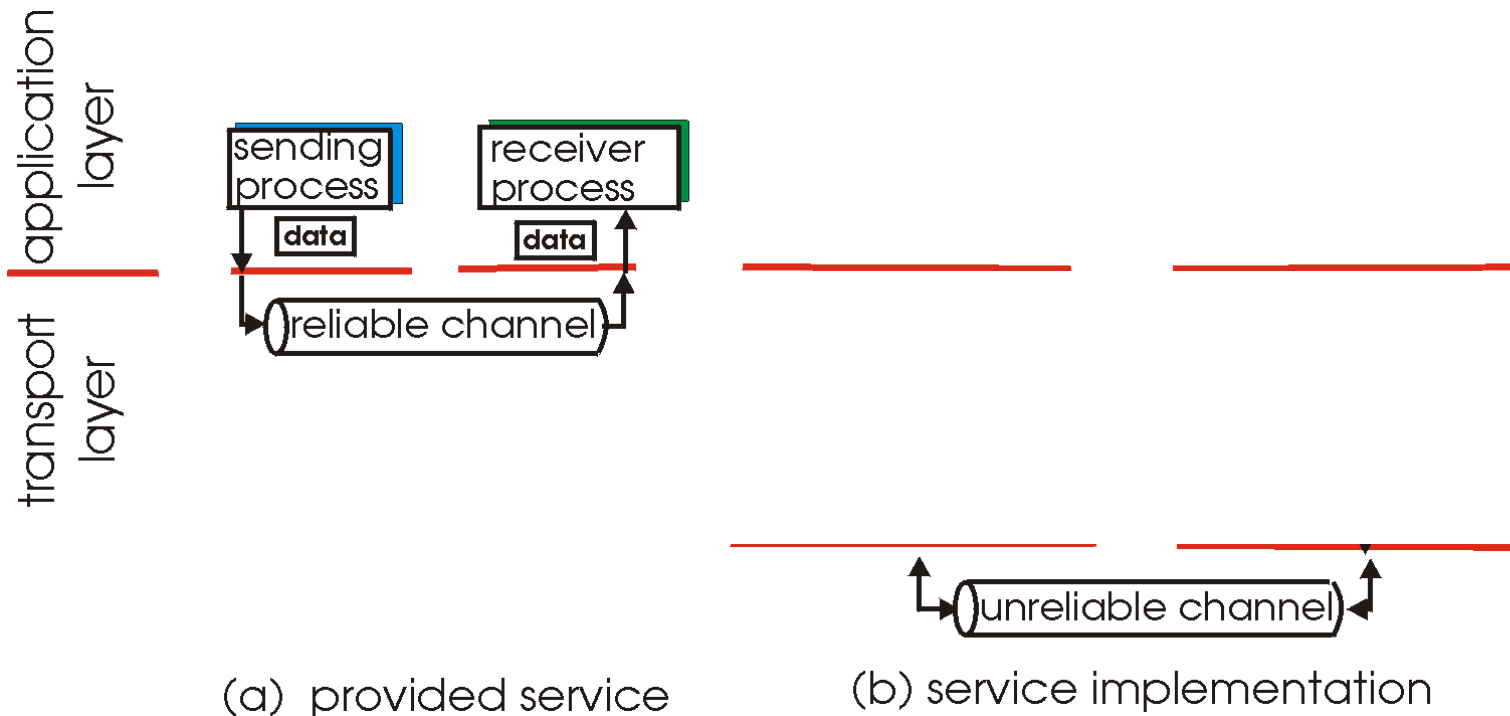

(a) provided service

(b) service implementation

❖ **characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)**
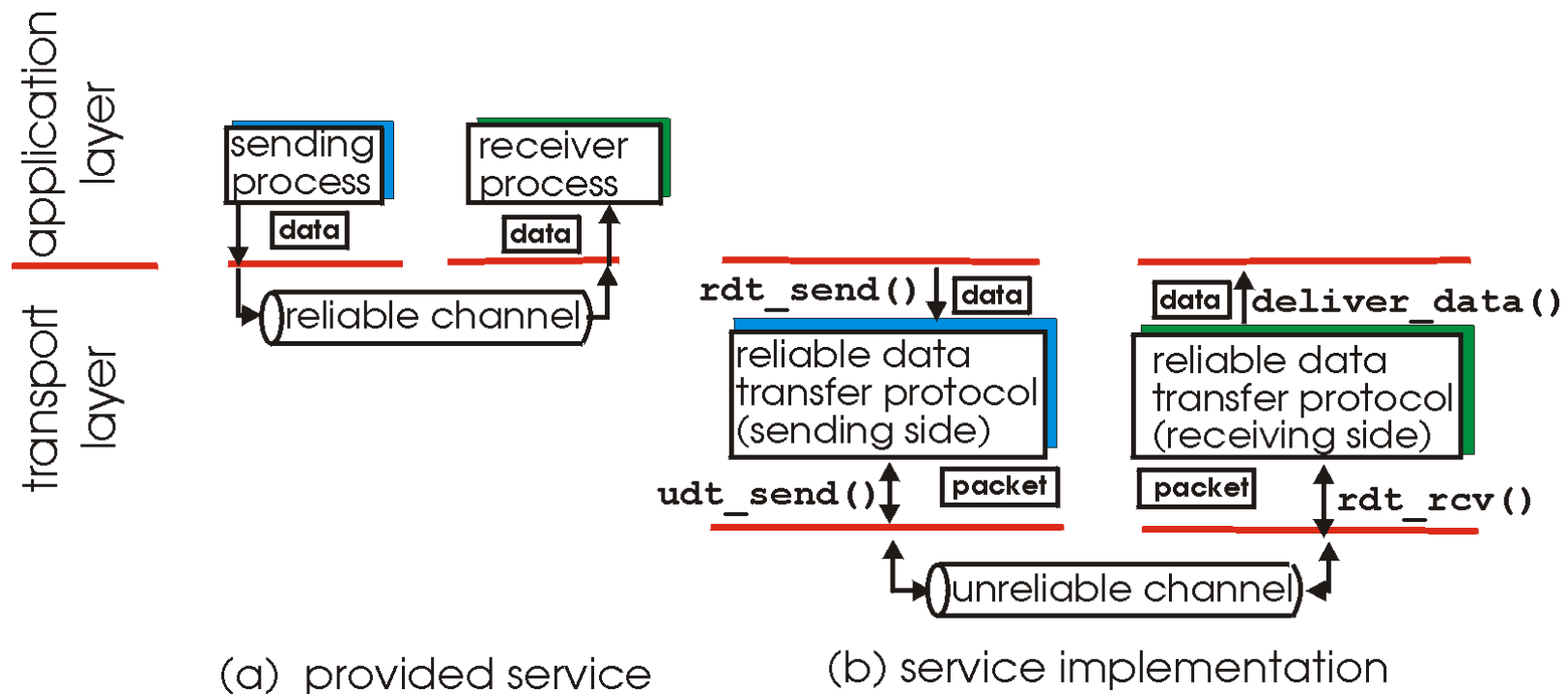
# Reliable data transfer (rdt):

❖ **Incrementally develop the sender and receiver sides with reliable data transfer protocol (rdt)**

❖ **rtd protocol version :**
  ➢ rdt1.0: reliable transfer over a reliable channel
    ❖ underlying channel perfectly reliable
      ▪ no bit errors
      ▪ no loss of packets
    ❖ no need to provide feedback to sender
    ❖ no need for the rcv to ask sender to slow down sending rate
  ➢ rdt2.0: channel with bit errors
  ➢ rdt3.0: channels with errors *and* loss

# rdt2.0: channel with bit errors

❖ **underlying channel may flip bits in packet**
  ▪ checksum to detect bit errors
❖ *the* question: how to recover from errors:

  ▪ *acknowledgements (ACKs):* receiver explicitly tells sender that pkt received OK
  ▪ *negative acknowledgements (NAKs):* receiver explicitly tells sender that pkt had errors
  ▪ sender retransmits pkt on receipt of NAK

❖ **new mechanisms in `rdt2.0` (beyond `rdt1.0`):**
  ▪ error detection
  ▪ feedback: control msgs (ACK,NAK) from receiver to sender

# rdt2.0 has a fatal flaw!

## what happens if ACK/NAK corrupted?

- ❖ sender doesn't know what happened at receiver!
- ❖ can't just retransmit: possible duplicate

## handling duplicates:

- ❖ sender retransmits current pkt if ACK/NAK corrupted
- ❖ sender adds *sequence number* to each pkt
- ❖ receiver discards (doesn't deliver up) duplicate pkt

**stop and wait**
sender sends one packet, then waits for receiver response

# rdt3.0: channels with errors *and* loss

**new assumption:**
underlying channel can also lose packets (data, ACKs)

- checksum, seq. #, ACKs, retransmissions will be of help … but not enough

**approach:** sender waits "reasonable" amount of time for ACK

❖ retransmits if no ACK received in this time
❖ if pkt (or ACK) just delayed (not lost):
  - retransmission will be duplicate, but seq. #'s already handles this
  - receiver must specify seq # of pkt being ACKed
❖ requires countdown timer

# rdt3.0 in action

sender                                                receiver

send pkt0
     pkt0
          rcv pkt0
          send ack0
     ack0
rcv ack0
send pkt1
     pkt1
          rcv pkt1
          send ack1
     ack1
rcv ack1
send pkt0
     pkt0
          rcv pkt0
          send ack0
     ack0

(a) no loss

sender                                                receiver

send pkt0
     pkt0
          rcv pkt0
          send ack0
     ack0
rcv ack0
send pkt1
     pkt1
        X
        *loss*

*timeout*
resend pkt1
     pkt1
          rcv pkt1
          send ack1
     ack1
rcv ack1
send pkt0
     pkt0
          rcv pkt0
          send ack0
     ack0

(b) packet loss

# rdt3.0 in action

**sender**                    **receiver**

send pkt0 → pkt0 → rcv pkt0
                               send ack0
rcv ack0 ← ack0
send pkt1 → pkt1 → rcv pkt1
                               send ack1
      ack1 ✗ loss
⏰ timeout
resend pkt1 → pkt1 → rcv pkt1
                               (detect duplicate)
                               send ack1
rcv ack1 ← ack1
send pkt0 → pkt0 → rcv pkt0
                               send ack0
      ← ack0

(c) ACK loss

**sender**                    **receiver**

send pkt0 → pkt0 → rcv pkt0
                               send ack0
rcv ack0 ← ack0
send pkt1 → pkt1 → rcv pkt1
                               send ack1
⏰ timeout ← ack1
resend pkt1 → pkt1 → rcv pkt1
                               (detect duplicate)
rcv ack1 ← pkt0                send ack1
send pkt0 → ← ack1 → rcv pkt0
rcv ack1 ← ack0                send ack0
send pkt0 → pkt0 → rcv pkt0
                               (detect duplicate)
      ← ack0                   send ack0

(d) premature timeout/ delayed ACK

# Performance of rdt3.0

❖ rdt3.0 is correct, but performance stinks
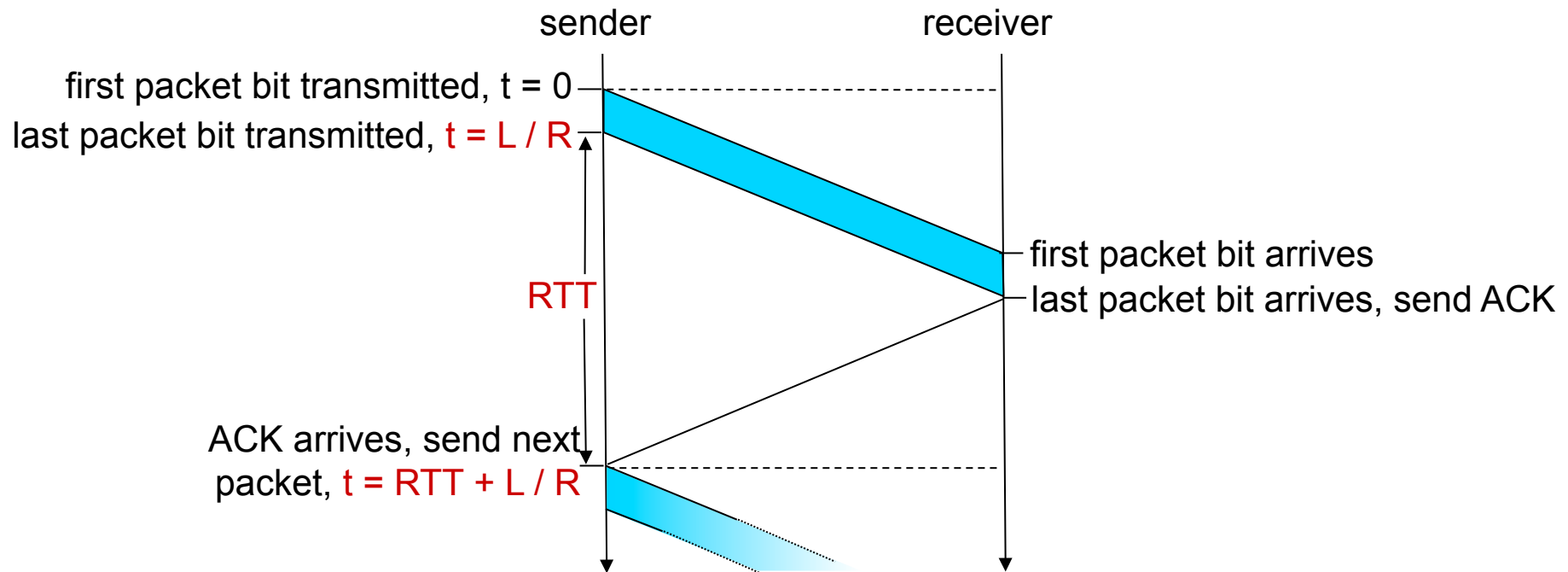
❖ e.g.: 1 Gbps link, 15 ms prop. delay, 8000 bit packet:

$$D_{trans} = \frac{L}{R} = \frac{8000\ bits}{10^9\ bits/sec} = 8\ microsecs$$

▪ U $_{sender}$: *utilization* – fraction of time sender busy sending

$$U_{sender} = \frac{L/R}{RTT + L/R} = \frac{.008}{30.008} = 0.00027$$

▪ Therefore, the throughput is 8Kb/30.008ms=267Kb/sec. If 1Kb pkt is transferred in every 30.008 msec, the throughput is 33Kb/sec over 1 Gbps link.

❖ network protocol limits use of physical resources!
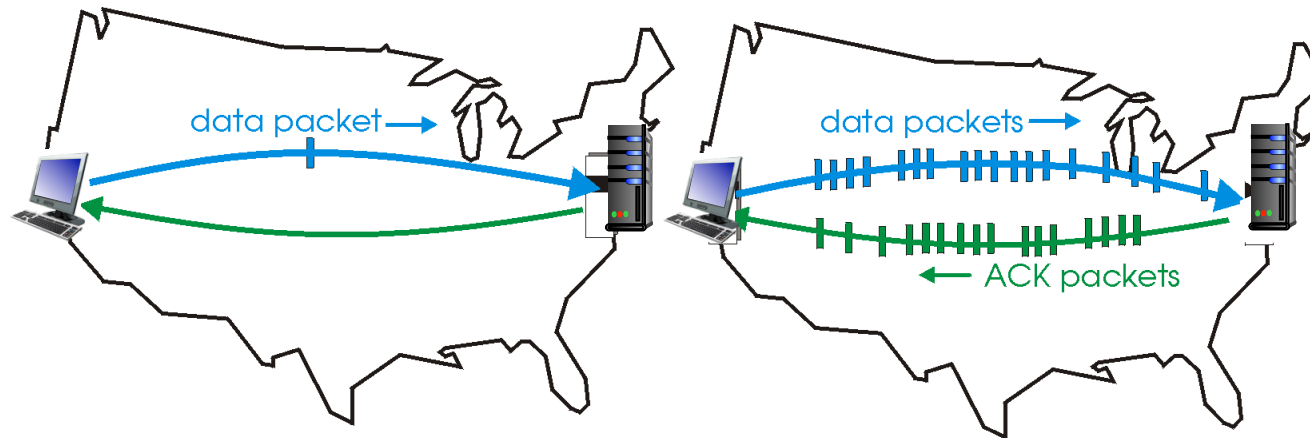
# rdt3.0: stop-and-wait operation

sender       receiver

first packet bit transmitted, t = 0

last packet bit transmitted, t = L / R

RTT

first packet bit arrives

last packet bit arrives, send ACK

ACK arrives, send next packet, t = RTT + L / R

$$U_{sender} = \frac{L/R}{RTT + L/R} = \frac{.008}{30.008} = 0.00027$$

# Pipelined protocols

pipelining: sender allows multiple, "in-flight", yet-to-be-acknowledged pkts
- range of sequence numbers must be increased
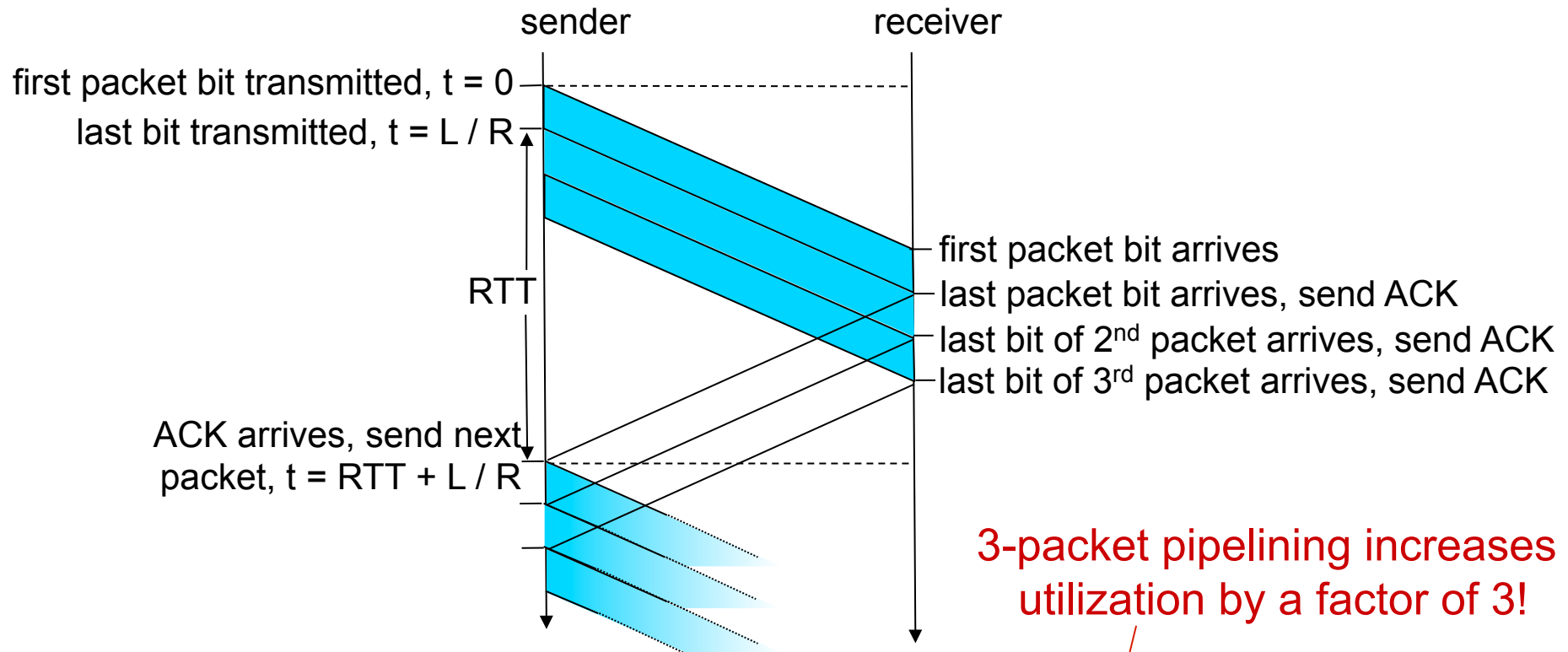- buffering at sender and/or receiver



(a) a stop-and-wait protocol in operation          (b) a pipelined protocol in operation

❖ two generic forms of pipelined protocols: *go-Back-N, selective repeat*

# Pipelining: increased utilization



sender       receiver

first packet bit transmitted, t = 0
last bit transmitted, t = L / R

RTT

first packet bit arrives
last packet bit arrives, send ACK
last bit of 2$^{nd}$ packet arrives, send ACK
last bit of 3$^{rd}$ packet arrives, send ACK

ACK arrives, send next packet, t = RTT + L / R

3-packet pipelining increases utilization by a factor of 3!

$$U_{sender} = \frac{3L / R}{RTT + L / R} = \frac{.0024}{30.008} = 0.00081$$
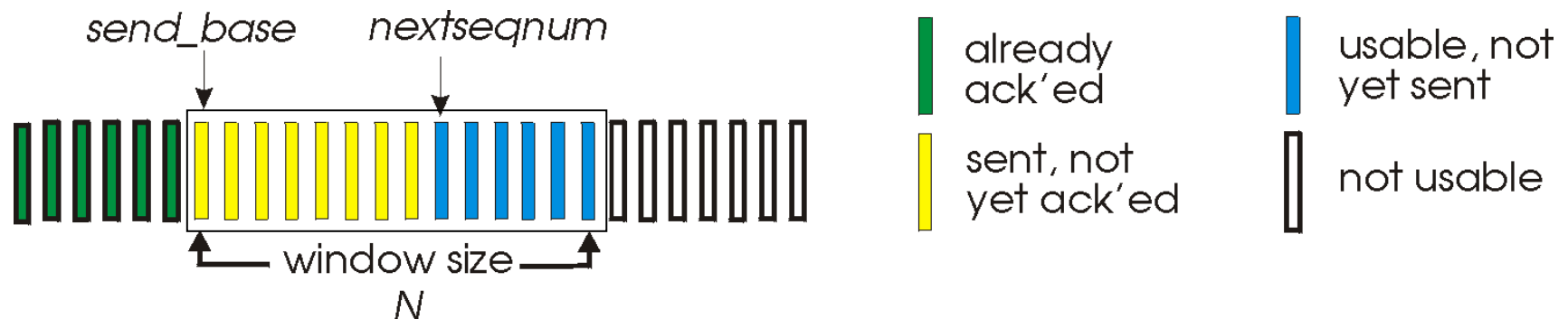
# Pipelined protocols: overview

## Go-back-N:

❖ sender can have up to N unack'ed packets in pipeline

❖ receiver only sends *cumulative ack*

  ▪ doesn't ack packet if there's a gap

❖ sender has timer for oldest unacked packet

  ▪ when timer expires, retransmit *all* unacked packets

## Selective Repeat:

❖ sender can have up to N unack'ed packets in pipeline

❖ rcvr sends *individual ack* for each packet

❖ sender maintains timer for each unacked packet

  ▪ when timer expires, retransmit only that unacked packet

# Go-Back-N: sender

❖ "window" size N and each k-bit has seq # in pkt header
❖ "window" of up to N, consecutive unack'ed pkts allowed



❖ ACK(n): ACKs all pkts up to, including seq # n - *"cumulative ACK"*

   ▪ may receive duplicate ACKs (see receiver)

❖ timer for oldest in-flight pkt

❖ *timeout(n):* retransmit packet n and all higher seq # pkts in window

# GBN in action

sender window (N=4)                    sender                           receiver

`0 1 2 3` 4 5 6 7 8          send  pkt0
`0 1 2 3` 4 5 6 7 8          send  pkt1
`0 1 2 3` 4 5 6 7 8          send  pkt2                          receive pkt0, send ack0
`0 1 2 3` 4 5 6 7 8          send  pkt3          **X**loss       receive pkt1, send ack1
                             (wait)

                                                                 receive pkt3, discard,
0 `1 2 3 4` 5 6 7 8          rcv ack0, send pkt4                        (re)send ack1
0 1 `2 3 4 5` 6 7 8          rcv ack1, send pkt5
                                                                 receive pkt4, discard,
                                                                        (re)send ack1
                      ignore duplicate ACK                       receive pkt5, discard,
                                                                        (re)send ack1
                   pkt 2 timeout

0 1 `2 3 4 5` 6 7 8          send  pkt2
0 1 `2 3 4 5` 6 7 8          send  pkt3
0 1 `2 3 4 5` 6 7 8          send  pkt4          rcv pkt2, deliver, send ack2
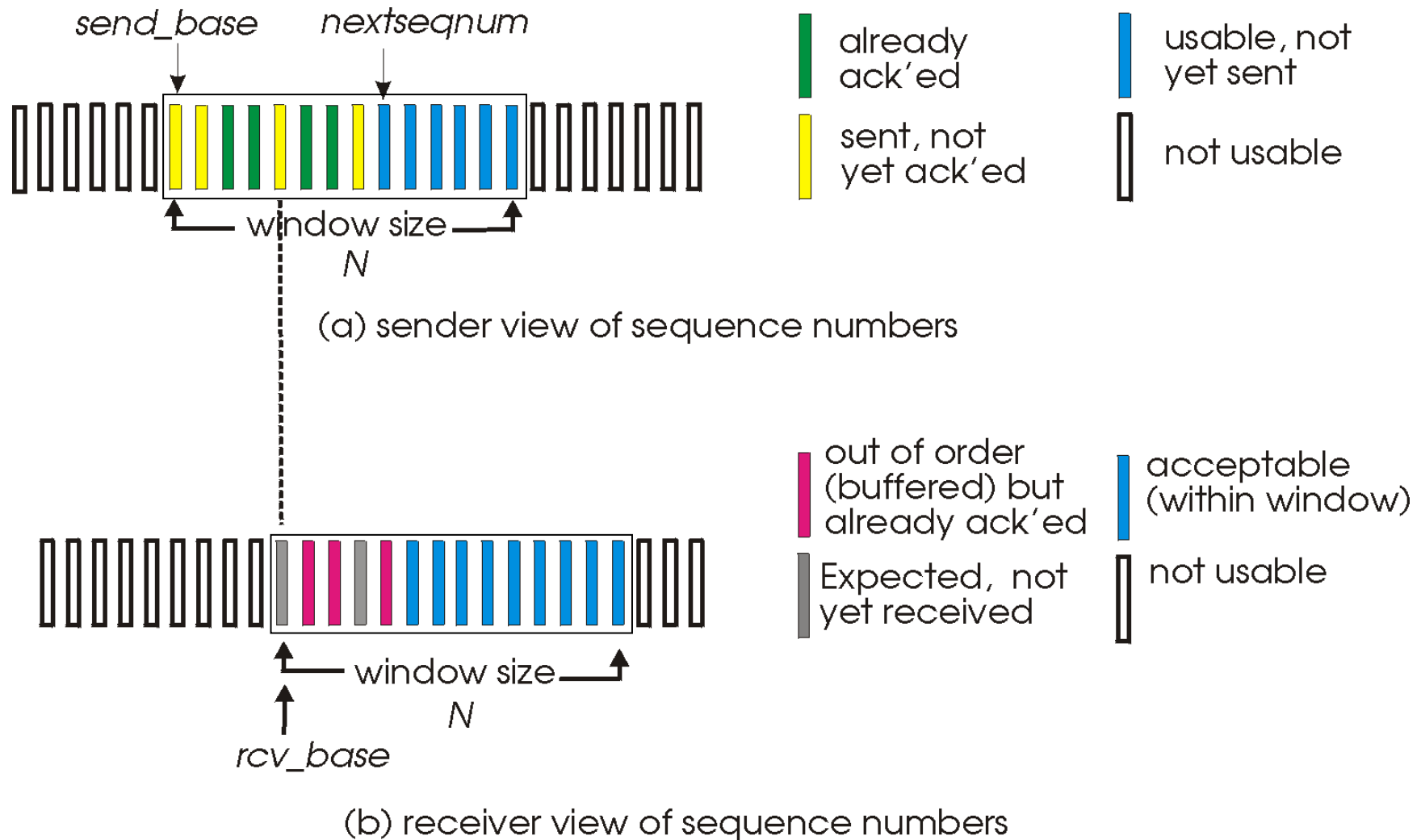0 1 `2 3 4 5` 6 7 8          send  pkt5          rcv pkt3, deliver, send ack3
                                                 rcv pkt4, deliver, send ack4
                                                 rcv pkt5, deliver, send ack5

# Selective repeat

❖ receiver *individually* acknowledges all correctly received pkts

- buffers pkts, as needed, for eventual in-order delivery to upper layer

❖ sender only resends pkts for which ACK not received

- sender timer for each unACKed pkt

❖ sender window

- has *N* consecutive seq #'s
- limits seq #s of sent, unACKed pkts (up-to "window size *N*"

# Selective repeat: sender, receiver windows

send_base      nextseqnum

| | already ack'ed | | usable, not yet sent |
| | sent, not yet ack'ed | | not usable |

window size N

(a) sender view of sequence numbers

| | out of order (buffered) but already ack'ed | | acceptable (within window) |
| | Expected, not yet received | | not usable |

window size N

rcv_base

(b) receiver view of sequence numbers

# Selective repeat (how it works?)

## sender

**data from above:**

❖ if next available seq # in window, send pkt

**timeout(n):**

❖ resend pkt n, restart timer

**ACK(n)** in [sendbase,sendbase+N]:

❖ mark pkt n as received

❖ if n smallest unACKed pkt, advance window base to next unACKed seq #

## receiver

**pkt n in** [rcvbase, rcvbase+N-1]

❖ send ACK(n)

❖ out-of-order: buffer

❖ in-order: deliver (also deliver buffered, in-order pkts), advance window to next not-yet-received pkt

**pkt n in** [rcvbase-N,rcvbase-1]

❖ ACK(n)

**otherwise:**

❖ ignore

# Selective repeat in action

sender window (N=4)

sender

receiver

0 1 2 3 4 5 6 7 8    send pkt0

0 1 2 3 4 5 6 7 8    send pkt1

0 1 2 3 4 5 6 7 8    send pkt2             receive pkt0, send ack0

0 1 2 3 4 5 6 7 8    send pkt3    **X**_loss_    receive pkt1, send ack1

(wait)

receive pkt3, buffer,
         send ack3

0 1 2 3 4 5 6 7 8    rcv ack0, send pkt4

0 1 2 3 4 5 6 7 8    rcv ack1, send pkt5

receive pkt4, buffer,
         send ack4

record ack3 arrived

receive pkt5, buffer,
         send ack5

_pkt 2 timeout_

0 1 2 3 4 5 6 7 8    send pkt2

0 1 2 3 4 5 6 7 8    record ack4 arrived

0 1 2 3 4 5 6 7 8    record ack5 arrived    rcv pkt2; deliver pkt2,
      pkt3, pkt4, pkt5; send ack2

0 1 2 3 4 5 6 7 8

_Q: what happens when ack2 arrives?_

# Chapter 3 outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer

3.5 connection-oriented transport: TCP
- segment structure
- reliable data transfer
- flow control
- connection management

3.6 principles of congestion control

3.7 TCP congestion control

# TCP: Overview RFCs: 793,1122,1323, 2018, 2581
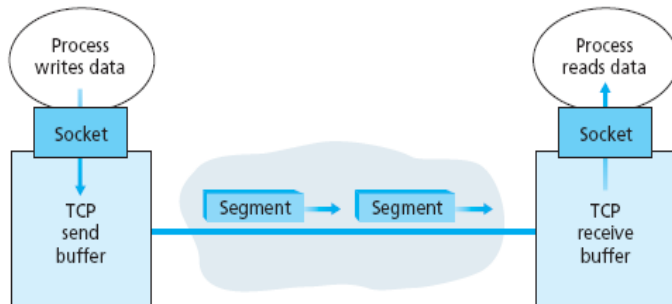
❖ **point-to-point:**
  ▪ one sender, one receiver



Figure 3.28 ♦ TCP send and receive buffers

❖ **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
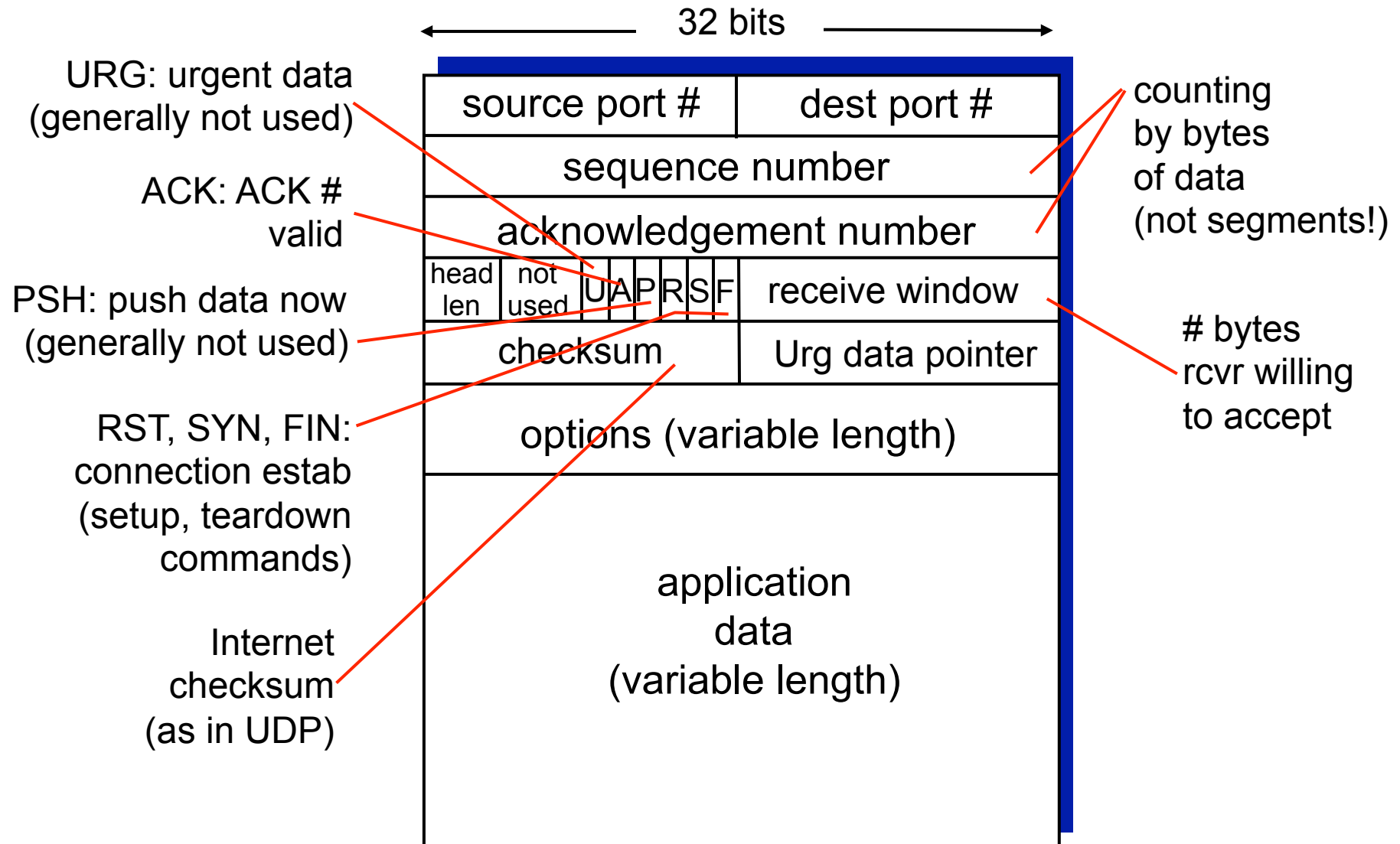  ▪ E.g. File size=500KB, MSS=1KB, so TCP construct 500 segments out of data stream.

❖ **connection-oriented:**
  ▪ handshaking (exchange of control msgs) inits sender, receiver state before data exchange

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

# 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)

counting
by bytes
of data
(not segments!)

# bytes
rcvr willing
to accept

| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |
| head len | not used | U A P R S F | receive window |
| checksum | Urg data pointer |
| options (variable length) | |
| application data (variable length) | |

# TCP seq. numbers, ACKs

sequence numbers (seq #):

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

acknowledgements (ACK):

- 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
  - E.g. use GBN or SR method
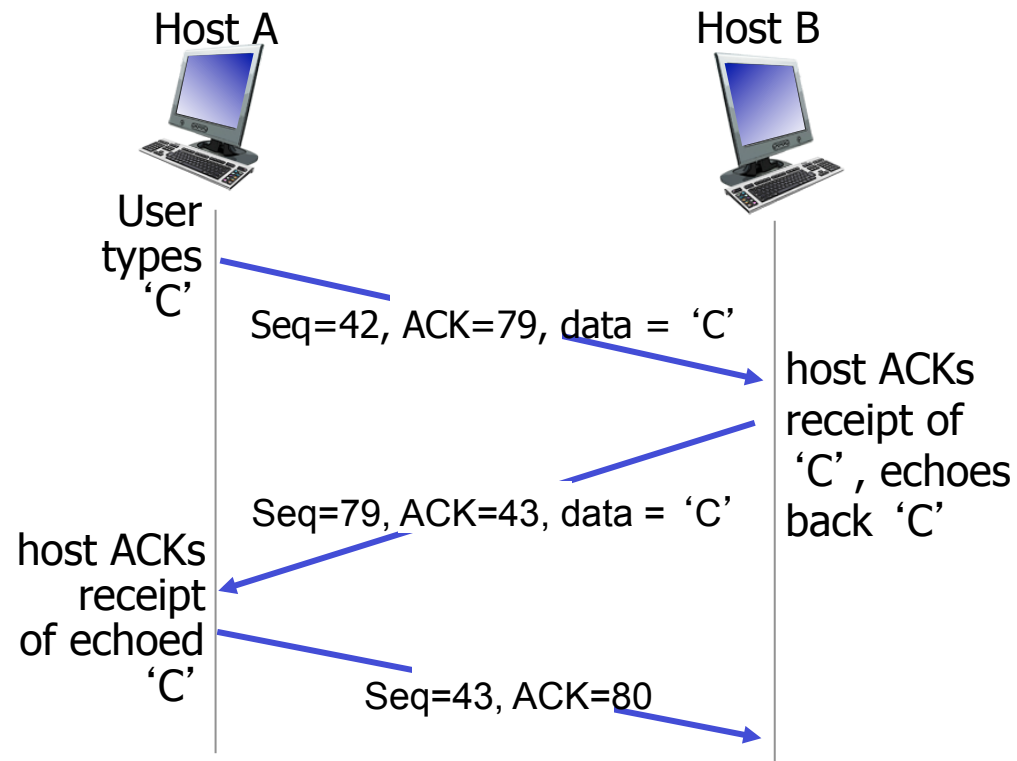
outgoing segment from sender

| 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 'C', echoes back 'C'

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

host ACKs receipt of echoed 'C'

Seq=43, ACK=80

simple telnet scenario

# TCP round trip time, 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**

# Chapter 3 outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer

3.5 connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management

3.6 principles of congestion control

3.7 TCP congestion control

# TCP reliable data transfer (rdt)

❖ **TCP creates rdt service on top of IP's unreliable service by implementing:**
  - pipelined segments
  - cumulative acks
  - single retransmission timer (refer to timer for oldest in-flight pkt)

let's initially consider simplified TCP sender:
  - ignore duplicate acks
  - ignore flow control, congestion control

❖ **retransmissions triggered by:**
  - timeout events
  - duplicate acks

*duplicate ACK,*
indicating seq. # of next expected byte

(Due to some reason expected seq. # is not received at receiver)

# 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
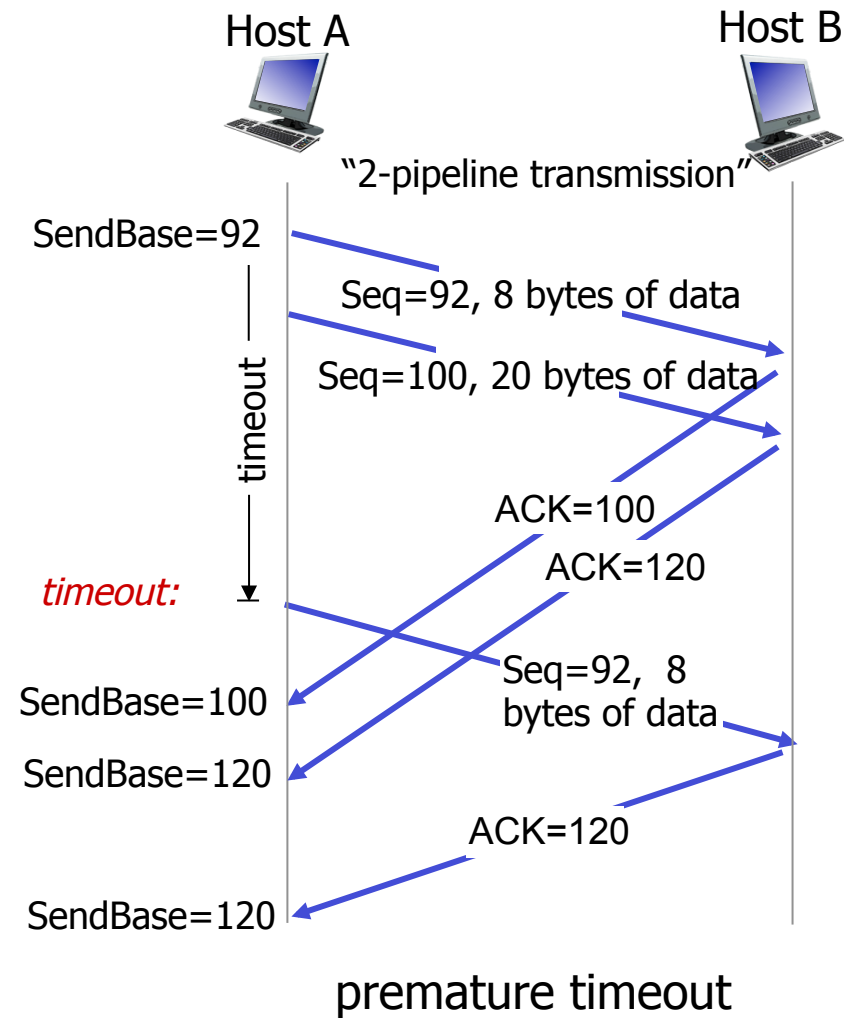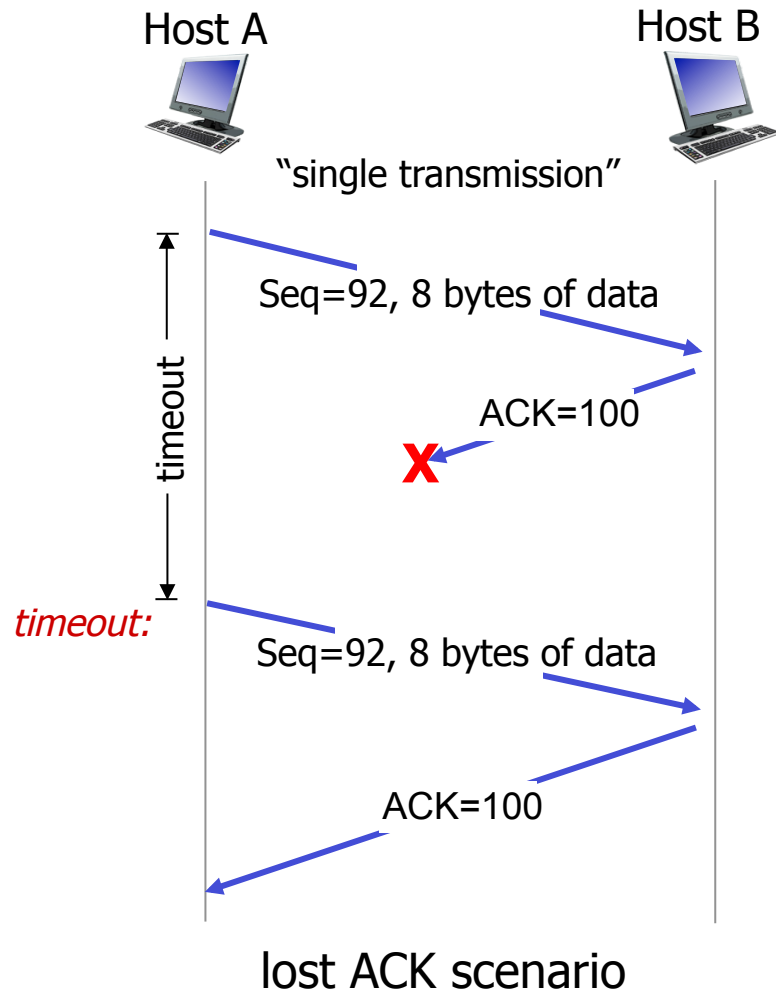  - ▪ expiration interval: `TimeOutInterval`

**timeout:**

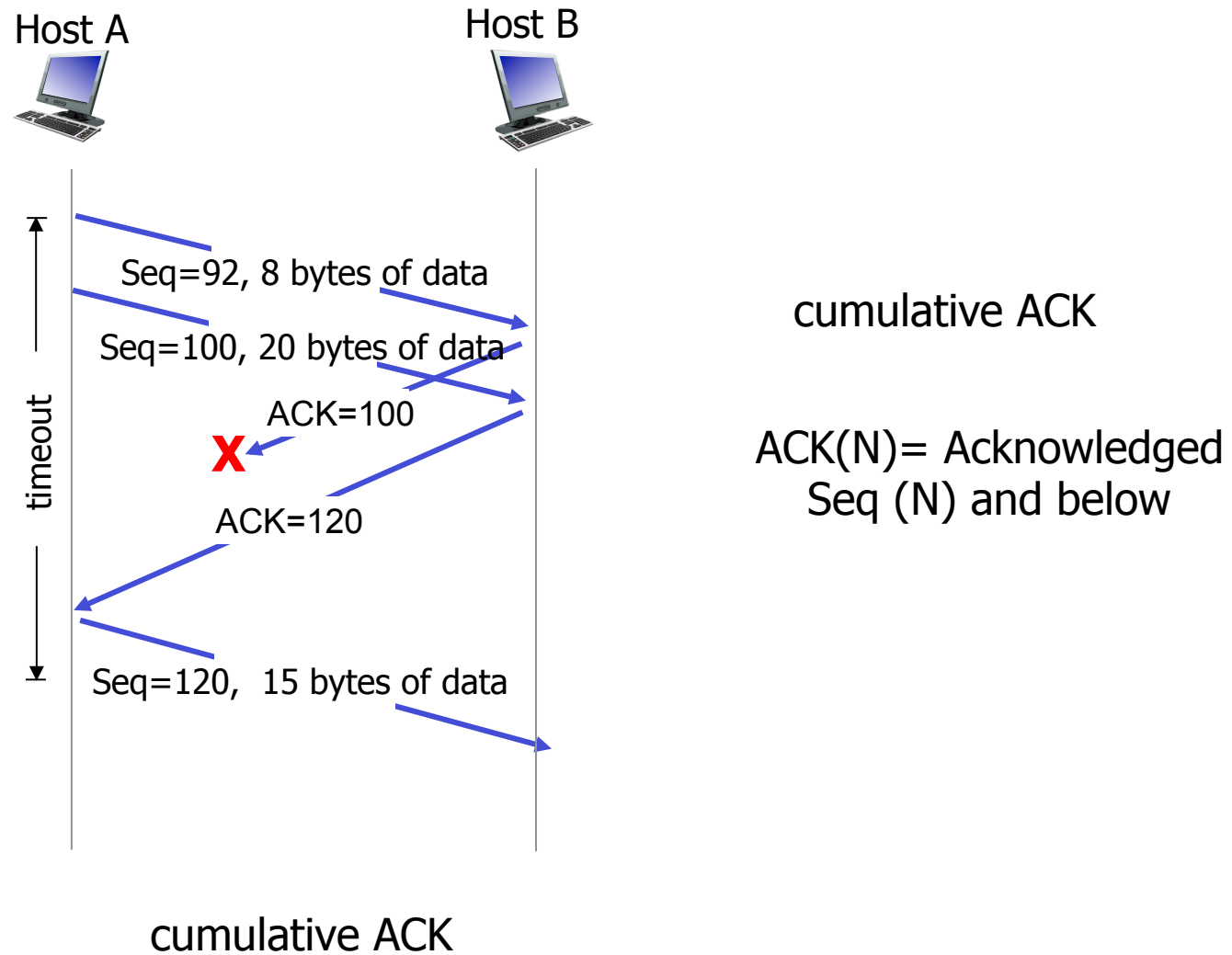- ❖ retransmit segment that caused timeout
- ❖ restart timer

**ack rcvd:**

- ❖ if ack acknowledges previously unacked segments
  - ▪ update what is known to be ACKed
  - ▪ start timer if there are still unacked segments

# TCP: retransmission scenarios

Host A                                    Host B

"single transmission"

Seq=92, 8 bytes of data

ACK=100

X

timeout:

Seq=92, 8 bytes of data

ACK=100

lost ACK scenario

Host A                                    Host B

"2-pipeline transmission"

SendBase=92

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

ACK=100

ACK=120

timeout:

SendBase=100

Seq=92, 8 bytes of data

SendBase=120

ACK=120

SendBase=120

premature timeout

# TCP: retransmission scenarios

Host A

Host B

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

ACK=100

**X**

ACK=120

Seq=120, 15 bytes of data

timeout

cumulative ACK

ACK(N)= Acknowledged
Seq (N) and below

cumulative ACK

# 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 *(retransmit – use oldest timer)* |
| arrival of out-of-order segment higher-than-expect seq. # . *Gap detected* | immediately send *duplicate ACK,* indicating seq. # of next expected byte *(TCP fast retransmit )* |
| arrival of segment that partially or completely fills gap *(between seq #)* | 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

*"triple duplicate ACKs"*

timeout

**Event:**

arrival of out-of-order segment
higher-than-expect seq. # .
Gap detected

**Action:**

immediately send *duplicate ACK*,
indicating seq. # of next expected byte

fast retransmit after sender
receipt of triple duplicate ACK

# Chapter 3 outline

# TCP flow control

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

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

*flow control*

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

application
process

application
-------
OS

TCP socket
receiver buffers

TCP
code

IP
code

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

rwnd

buffered data

free buffer space

*TCP segment payloads*

*receiver-side buffering*

```
RcvBuffer = received buffer data

rwnd = received window
       free buffer space
```

# Chapter 3 outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer

3.5 connection-oriented transport: TCP
- segment structure
- reliable data transfer
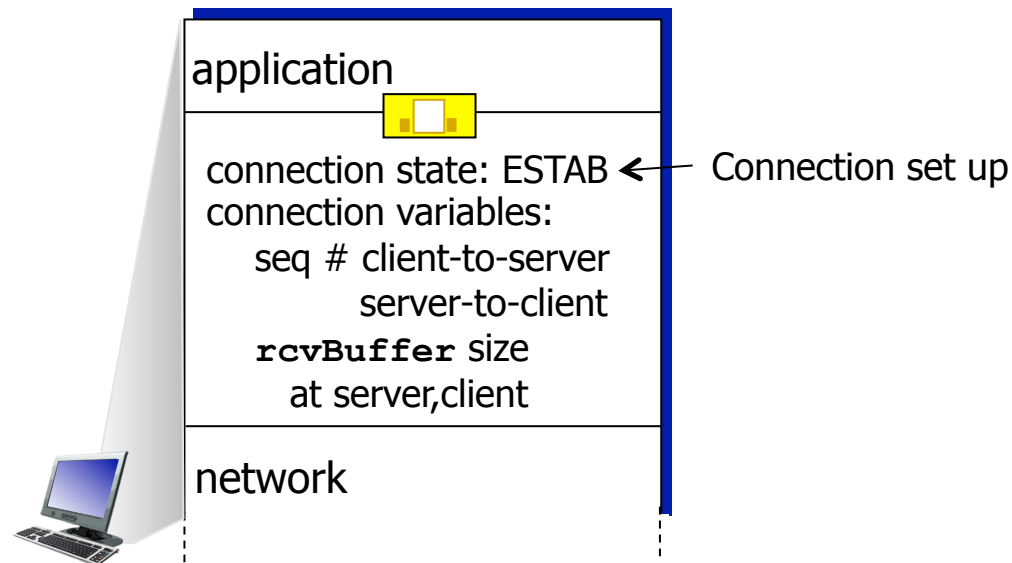- flow control
- connection management

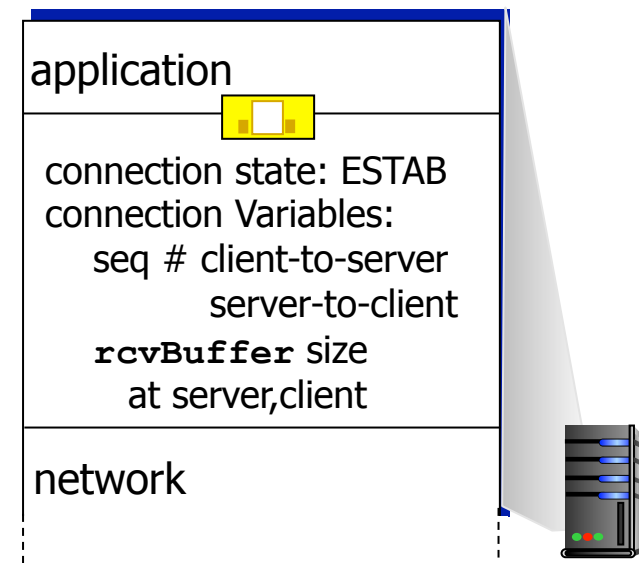3.6 principles of congestion control

3.7 TCP congestion control

# Connection Management (TCP)

before exchanging data, sender/receiver "handshake":

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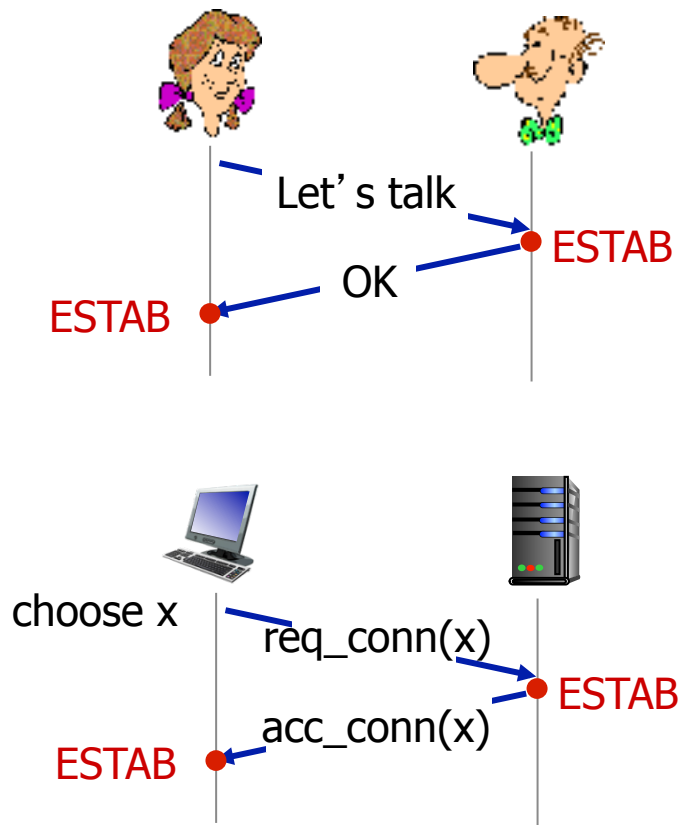- ❖ agree on connection parameters

application

connection state: ESTAB ← Connection set up
connection variables:
    seq # client-to-server
            server-to-client
    `rcvBuffer` size
        at server,client

network

application

connection state: ESTAB
connection Variables:
    seq # client-to-server
            server-to-client
    `rcvBuffer` size
        at server,client

network

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

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

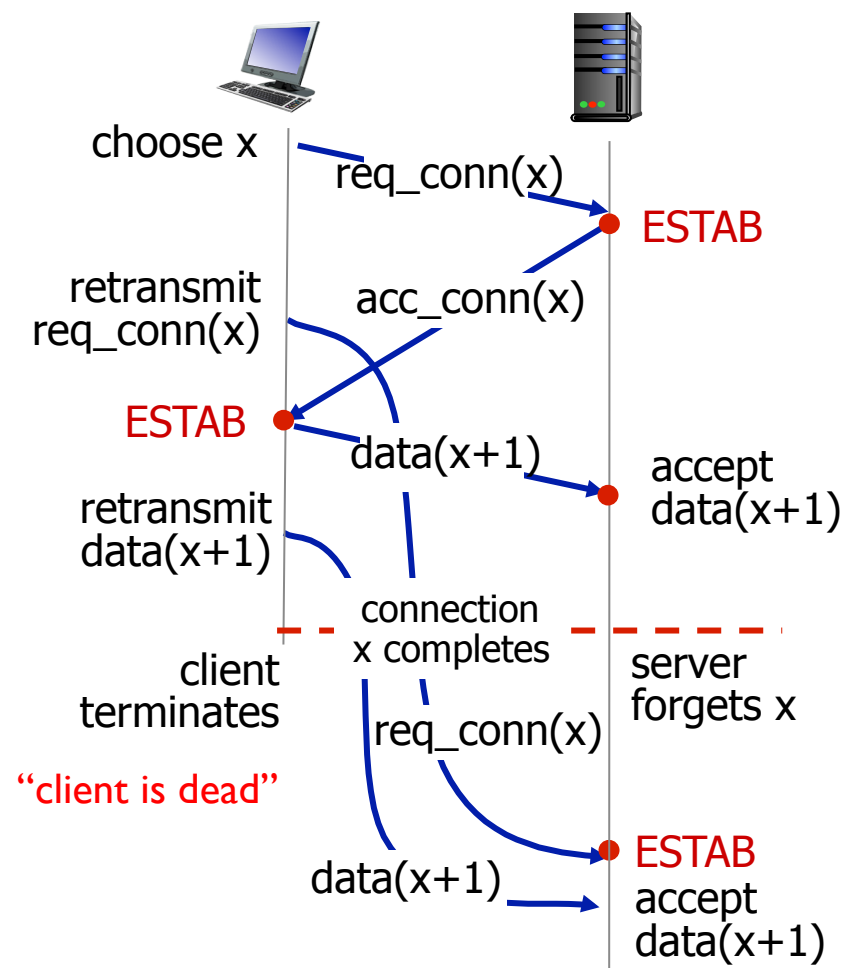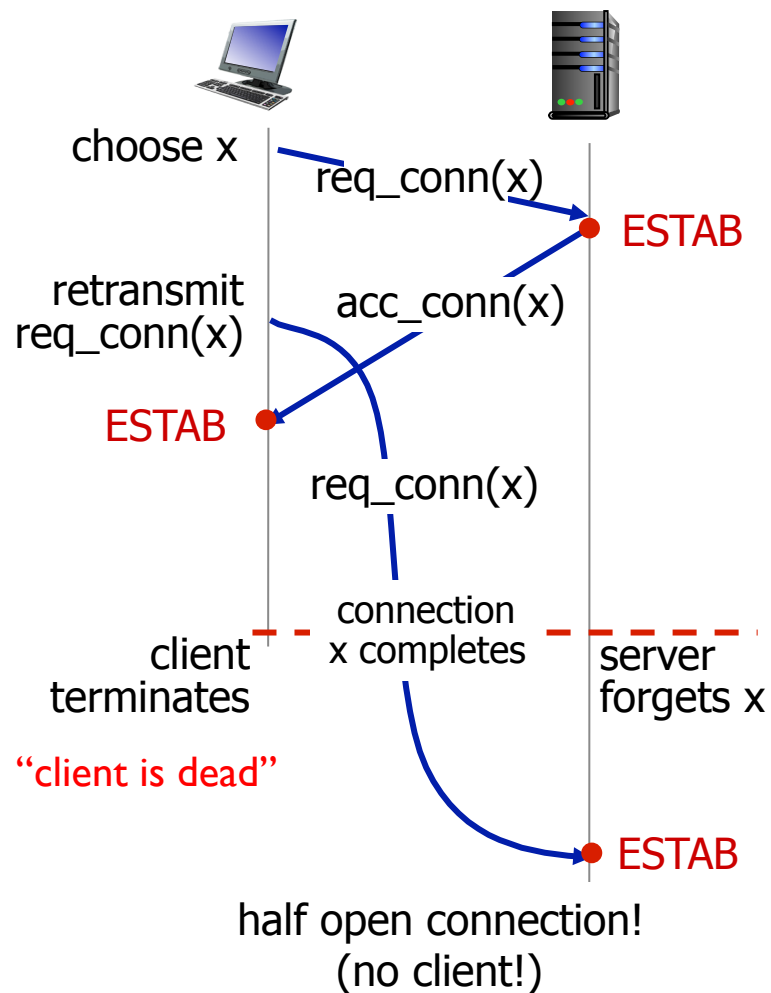# Agreeing to establish a connection

2-way handshake:



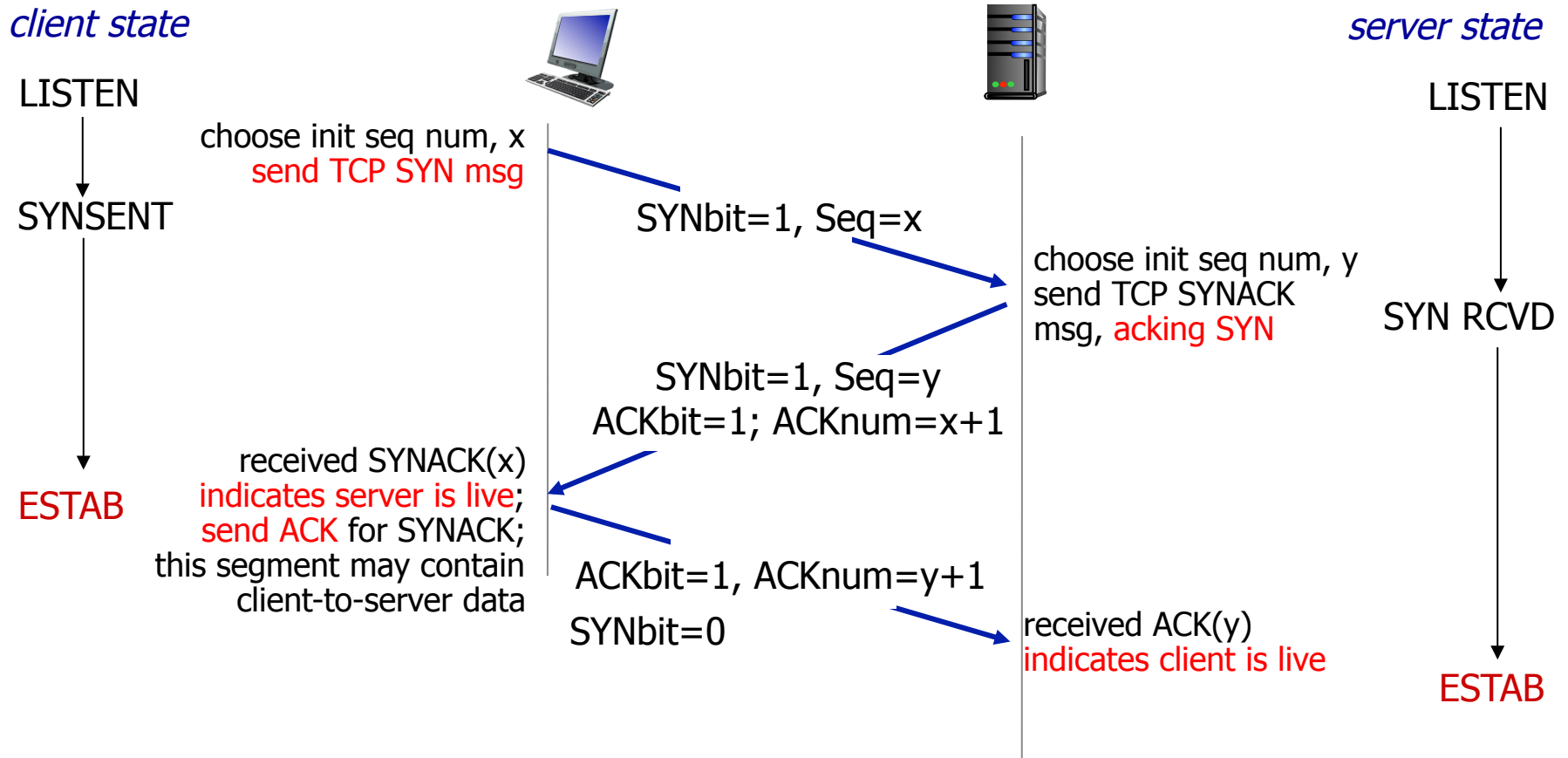*Q:* will 2-way handshake always work in network?

❖ 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)

client
terminates
connection
x completes
server
forgets x

"client is dead"

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)

client
terminates
connection
x completes
server
forgets x

"client is dead"
req_conn(x)

data(x+1)
ESTAB
accept
data(x+1)

# TCP 3-way handshake

*client state*

LISTEN

SYNSENT

ESTAB

choose init seq num, x
send TCP SYN msg

SYNbit=1, Seq=x

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
SYNbit=0

*server state*

LISTEN

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

SYN RCVD

received ACK(y)
indicates client is live

ESTAB

# 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

              

ESTAB

clientSocket.close()

FIN_WAIT_1    can no longer
send but can
receive data

        FINbit=1, seq=x

                          ESTAB

                          CLOSE_WAIT

      ACKbit=1; ACKnum=x+1

FIN_WAIT_2    wait for server
close

                         can still
send data

         FINbit=1, seq=y

                         LAST_ACK

TIMED_WAIT                        can no longer
send data

      ACKbit=1; ACKnum=y+1

   timed wait
for 2*max
segment lifetime

                         CLOSED

CLOSED

# Chapter 3 outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer

3.5 connection-oriented transport: TCP
- segment structure
- reliable data transfer
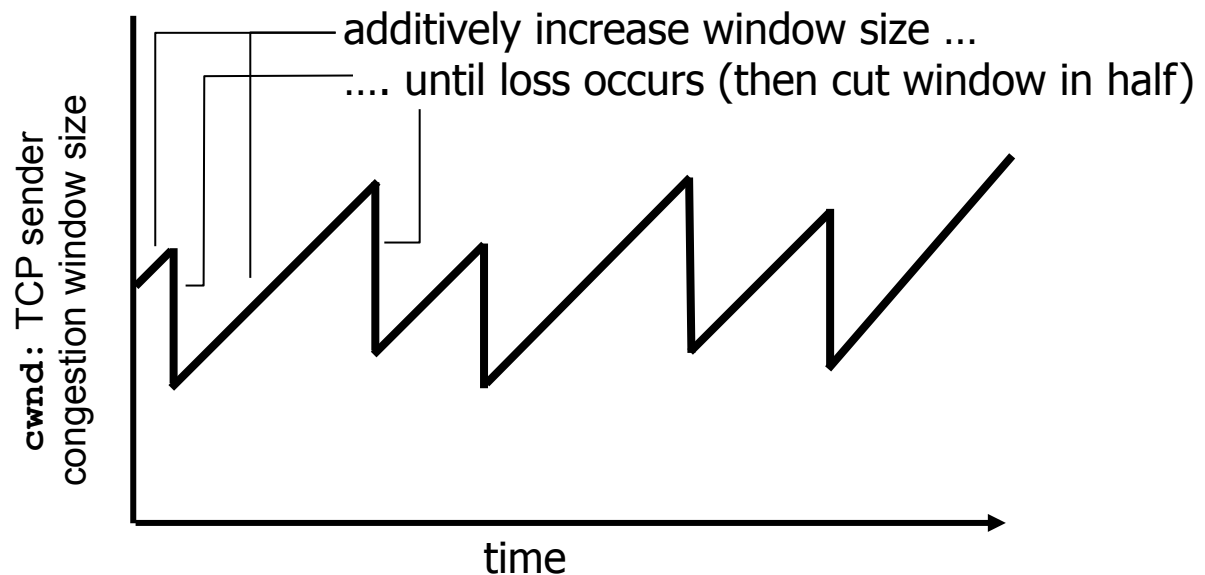- flow control
- connection management

3.6 principles of congestion control

3.7 TCP congestion control

# Principles of congestion control

*congestion:*

❖ informally: "too many sources sending too much data too fast for *network* to handle"

❖ different from flow control!

❖ manifestations:

- lost packets (buffer overflow at routers)
- long delays (queueing in router buffers)

❖ a top-10 problem!

# Approaches towards congestion control

two broad approaches towards congestion control:

## end-end congestion control:

- ❖ no explicit feedback from network
- ❖ congestion inferred from end-system observed loss, delay (e.g. from timeout, duplicate ACK)
- ❖ approach taken by TCP

## network-assisted congestion control:

- ❖ routers provide feedback to end systems
  - single bit indicating congestion (as implemented by SNA, DECbit, TCP/IP ECN, ATM)
  - explicit rate for sender to send at

# Chapter 3 outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer

3.5 connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management

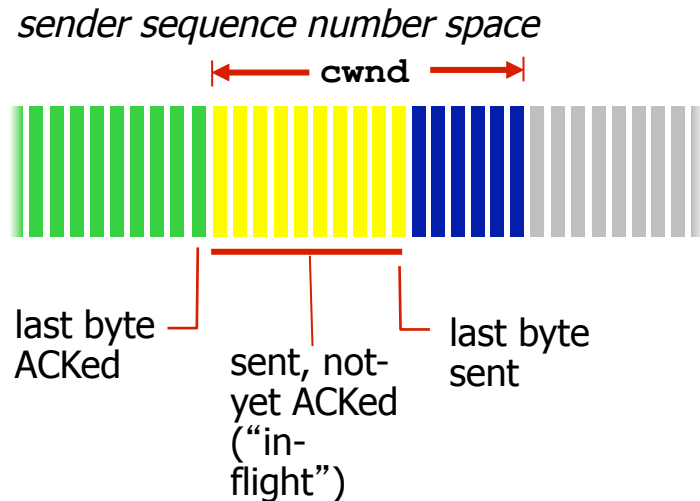3.6 principles of congestion control

3.7 TCP congestion control

# TCP congestion control: additive increase multiplicative decrease (AIMD)

❖ *approach:* sender increases transmission rate (window size), probing for usable bandwidth, until loss occurs

- *additive increase:* increase `cwnd (congestion window)` by 1 MSS (Maximum Segment Size) every RTT until loss detected

- *multiplicative decrease:* cut `cwnd` in half after loss

additively increase window size ...

.... until loss occurs (then cut window in half)

AIMD saw tooth behavior: probing for bandwidth

**cwnd:** TCP sender congestion window size

time

# TCP Congestion Control: details

*sender sequence number space*



last byte
ACKed

sent, not-
yet ACKed
("in-
flight")

last byte
sent

❖ sender limits transmission:

$$LastByteSent - LastByteAcked \le cwnd$$

❖ **cwnd** is dynamic, function of perceived (recognized) network congestion

*TCP sending rate:*

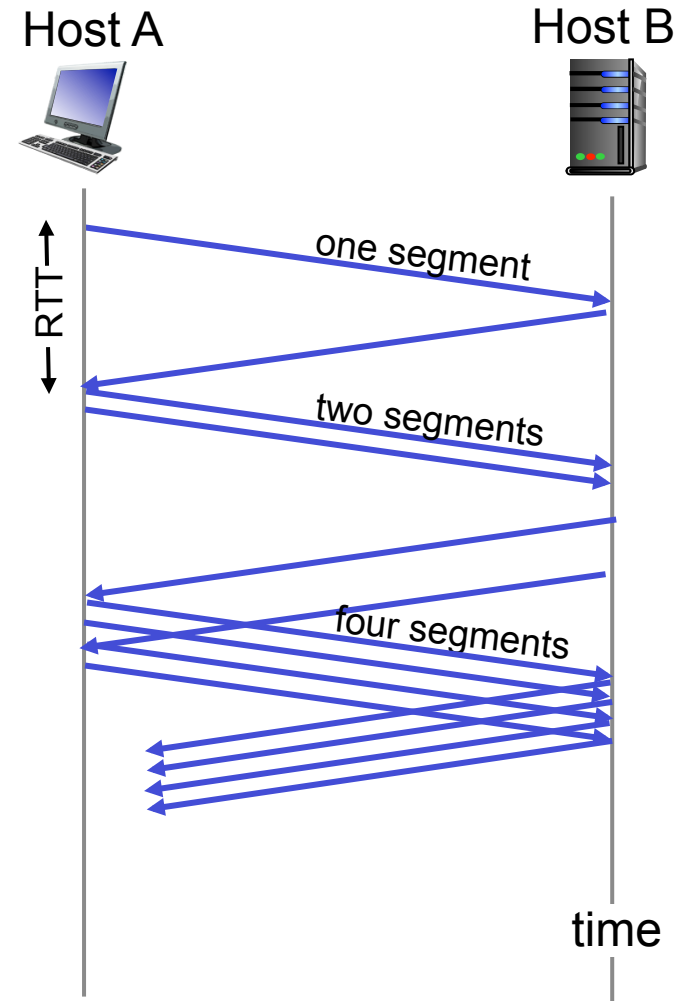❖ *roughly:* send cwnd bytes, wait RTT for ACKS, then send more bytes

$$rate \approx \frac{cwnd}{RTT} \text{ bytes/sec}$$

# TCP Slow Start

❖ **when connection begins, increase rate exponentially until first loss event:**
  - initially `cwnd` = 1 MSS
  - double `cwnd` every RTT
  - done by incrementing `cwnd` for every ACK received

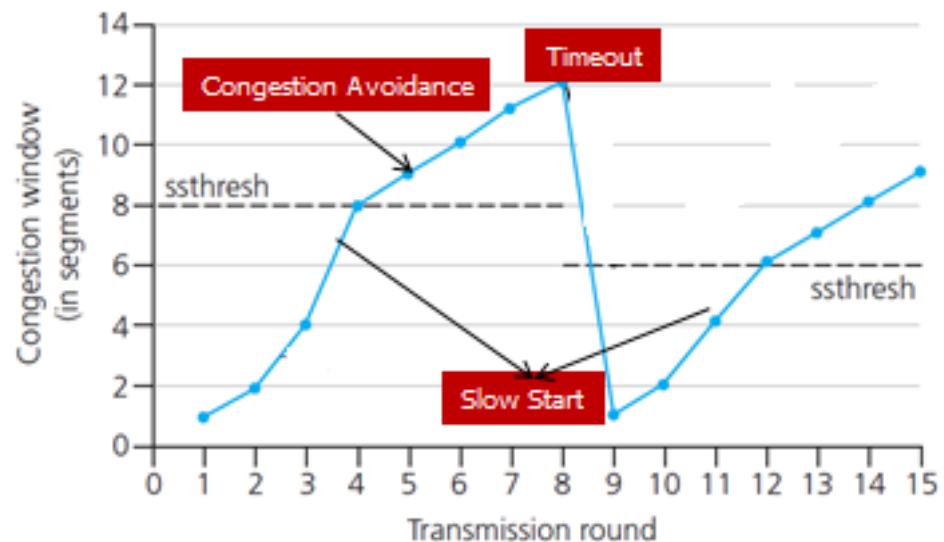❖ *summary:* initial rate is slow but ramps up <span style="color:red">exponentially</span> fast

Host A
Host B

RTT

one segment
two segments
four segments

time

# TCP: Slow Start & Congestion Avoidance (CA) (Loss because of Timeout)

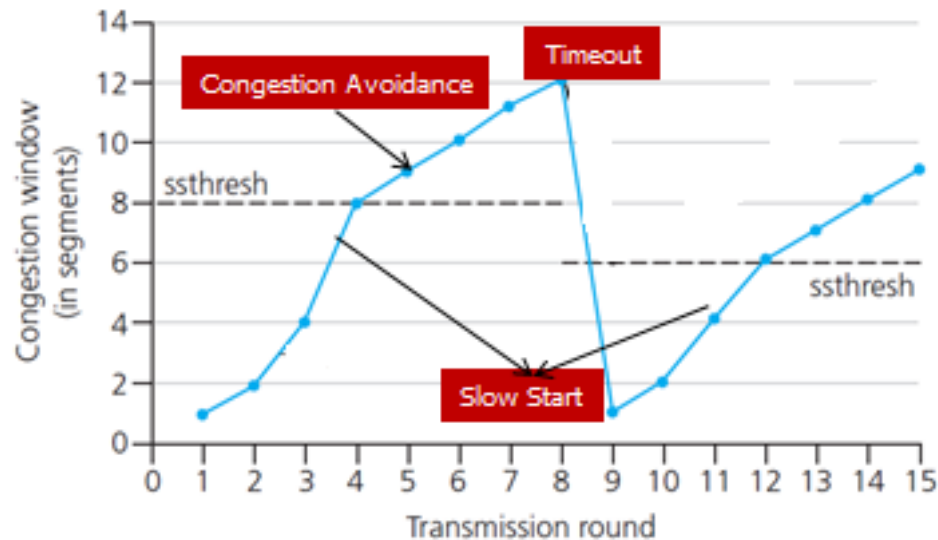Q: when should the exponential increase switch to linear?

A: when **cwnd** gets to 1/2 of its value before timeout. (Congestion Avoidance)

## Implementation:

❖ variable **ssthresh (slow-start threshold)**

❖ on loss event:
  - ❖ **ssthresh** is set to 1/2 of **cwnd** just before loss event
  - ❖ Value of cwnd is set to 1 MSS (slow start)

# Switching from slow start to CA



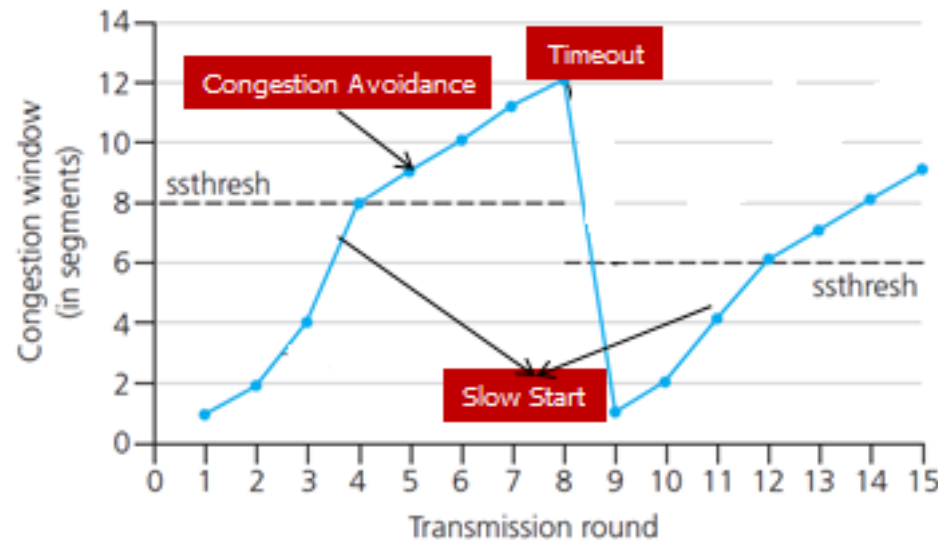| Phase | TR | CW | SS | ssth |
|-------|----|----|----|------|
| Slow start | 1 | 1 | 1 | 8 |
| | 2 | 2 | 3 | 8 |
| | 3 | 4 | 7 | 8 |
| | 4 | 8 | 15 | 8 |
| CA | 5 | 9 | 24 | 8 |
| | 6 | 10 | 34 | 8 |
| | 7 | 11 | 44 | 8 |
| | 8 | 12 | 56 | 12/2 = 6 |

TR=Transmission round
CW=Congestion Window
SS=Segment Send
ssthreshold=slow start threshold

TR 1 to 4
- Slow Start, Exponential growth, ssth=8

TR 4 = ssth is detected and Congestion Avoidance (CA) starts

TR 5 to 8
- Operate at CA, Linear growth

# LOSS because of TIMEOUT



| TR | CW | SS | ssth |
|----|------|------|------|
| 9 | 1 | 57 | 6 |
| 10 | 2 | 59 | |
| 11 | 4 | 63 | |
| 12 | 6 (8) | 69 | 6 |
| 13 | 7 | 78 | |
| 14 | 8 | 86 | |
| 15 | 9 | 95 | |
| | | | |

TR=Transmission round
CW=Congestion Window
SS=Segment Send
ssthreshold=slow start threshold

After TR 8
- Timeout is detected

TR 9 to 12 (refer table)
- CW=1, and  ssth=1/2*CW(current)=6
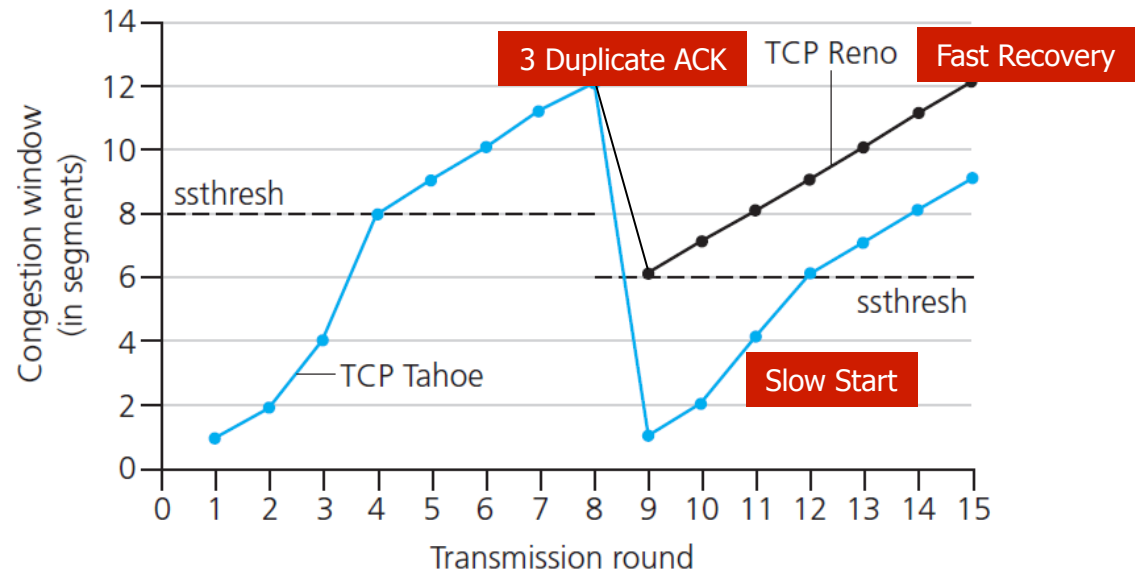- Start Slow, Exponential Growth and ssthreshold = 6

TR 12 to 15
- Operate at CA, Linear Growth

# TCP: Fast Recovery
## (Loss because of 3 Duplicate ACK)

Earlier version of TCP
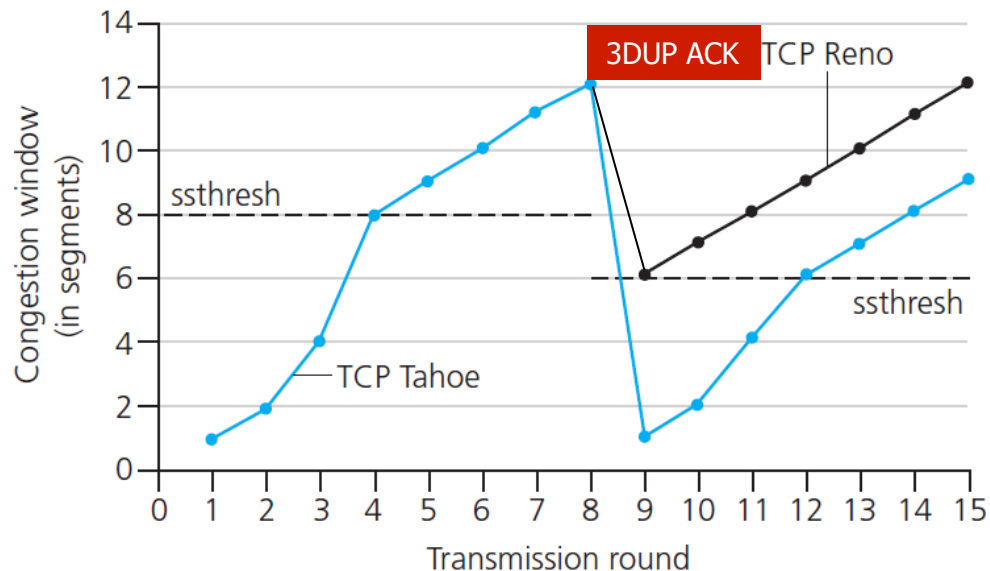(TCP Tahoe)
entered Slow start

Newer version of
TCP (TCP Renoe)
incorporated fast
recovery

## Implementation:

❖ on loss event, **ssthresh** is
set to 1/2 of **cwnd** just before
loss event

❖ **cwnd** is cut in half window then
grows linearly

# LOSS because of 3 DUPLICATE ACK (TCP RENO)



| TR | CW | SS | ssth |
|----|-----|-----|------|
| 9 | 6 | 62 | 6 |
| 10 | 7 | 69 | |
| 11 | 8 | 77 | |
| 12 | 9 | 86 | |
| 13 | 10 | 96 | |
| 14 | 11 | 107 | |
| 15 | 12 | 119 | |

TR=Transmission round
CW=Congestion Window
SS=Segment Send
ssthreshold=slow start threshold

After TR 8 3 DUP ACK is detected

TR 9
-CW=1/2*CW(current)=12/2=6

TR 9,10 to 15
- Enters Fast Recovery
- Operate at Congestion Avoidance (CA)
- Linear growth

# TCP: detecting, reacting to loss

## TCP RENO

❖ loss indicated by timeout: (Slow Start)
  ▪ `cwnd` set to 1 MSS;
  ▪ window then grows exponentially (as in slow start) to threshold, then grows linearly

❖ loss indicated by 3 duplicate ACKs (Fast Recovery)
  ▪ dup ACKs indicate network capable of delivering some segments
  ▪ `cwnd` is cut in half window then grows linearly

## TCP Tahoe

❖ loss indicated by timeout or 3 duplicate ACKs : (Slow Start)
  ▪ `cwnd` set to 1 MSS;
  ▪ window then grows exponentially (as in slow start) to threshold, then grows linearly

# Chapter 3: summary

❖ principles behind
transport layer services:

- multiplexing,
demultiplexing
- reliable data transfer
- flow control
- congestion control

next:

❖ leaving the
network
"edge" (application
, transport layers)
❖ into the network
"core"