# Network Transport Layer: TCP

Qiao Xiang

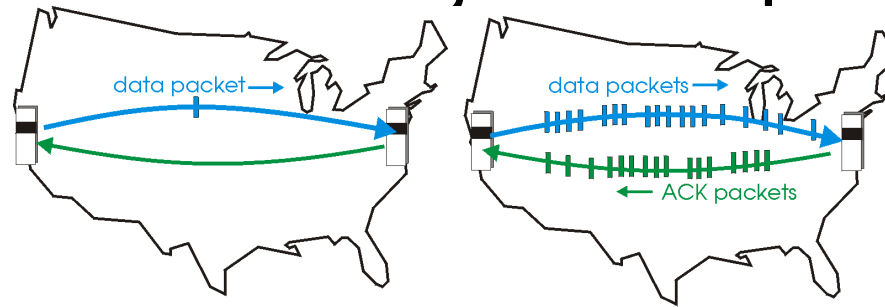https://qiaoxiang.me/courses/cnns-xmuf21/index.shtml

11/16/2021

# Admin

❑ Lab assignment 2 to be returned on Nov. 18
❑ Class Project
  ○ 15% of your final score
  ○ Please start ASAP
  ○ Talk to the instructor or TA to get feedback

# Recap: Reliable Transport

❑ Basic structure: sliding window protocols



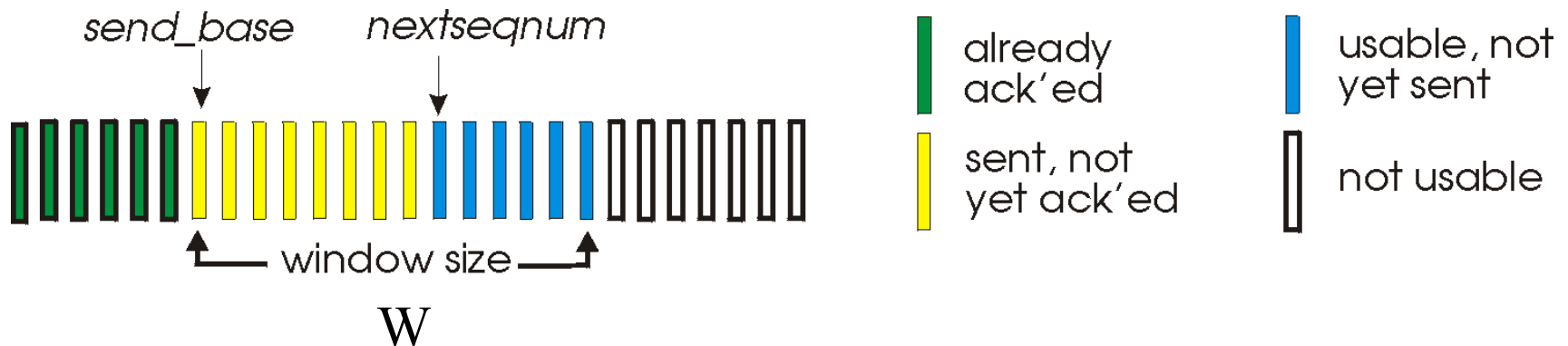(a) a stop-and-wait protocol in operation  (b) a pipelined protocol in operation

General technique: pipelining.
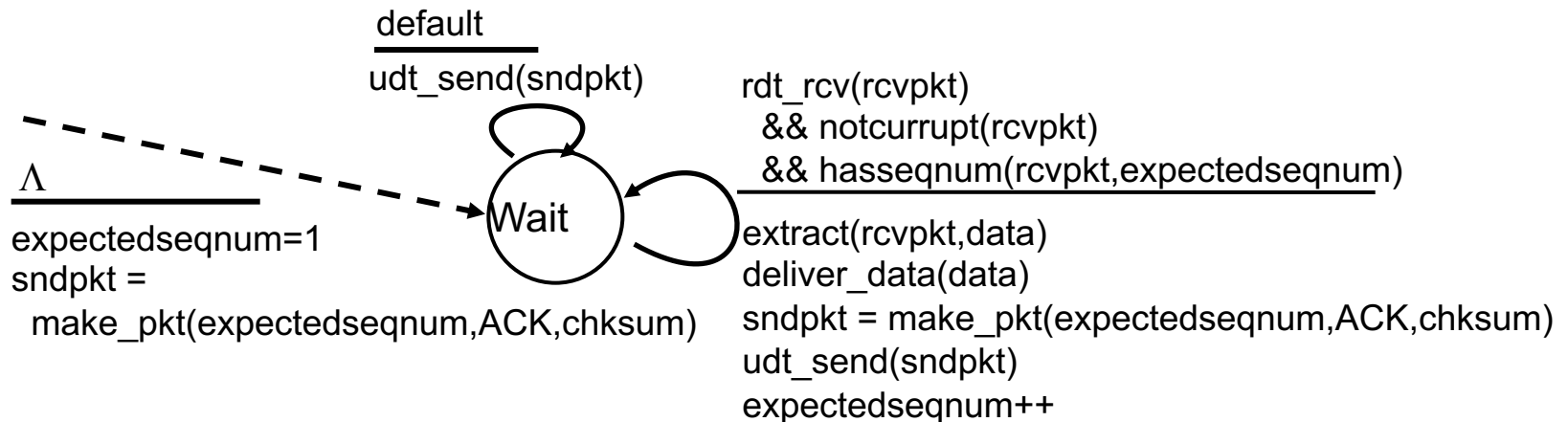
❑ Realization: GBN or SR

# Recap: Go-Back-N

## Sender:

- k-bit *seq #* in pkt header
- "window" of up to W, consecutive unack'ed pkts allowed



- ACK(n): ACKs all pkts up to, including seq # n - "cumulative ACK"
  - note: ACK(n) could mean two things: I have received upto and include n, or I am waiting for n
- timer for the packet at base
- *timeout(n):* retransmit pkt n and all higher seq # pkts in window

# Recap: Go-Back-N

default
udt_send(sndpkt)

rdt_rcv(rcvpkt)
  && notcurrupt(rcvpkt)
  && hasseqnum(rcvpkt,expectedseqnum)

Λ
expectedseqnum=1
sndpkt =
  make_pkt(expectedseqnum,ACK,chksum)

Wait

extract(rcvpkt,data)
deliver_data(data)
sndpkt = make_pkt(expectedseqnum,ACK,chksum)
udt_send(sndpkt)
expectedseqnum++
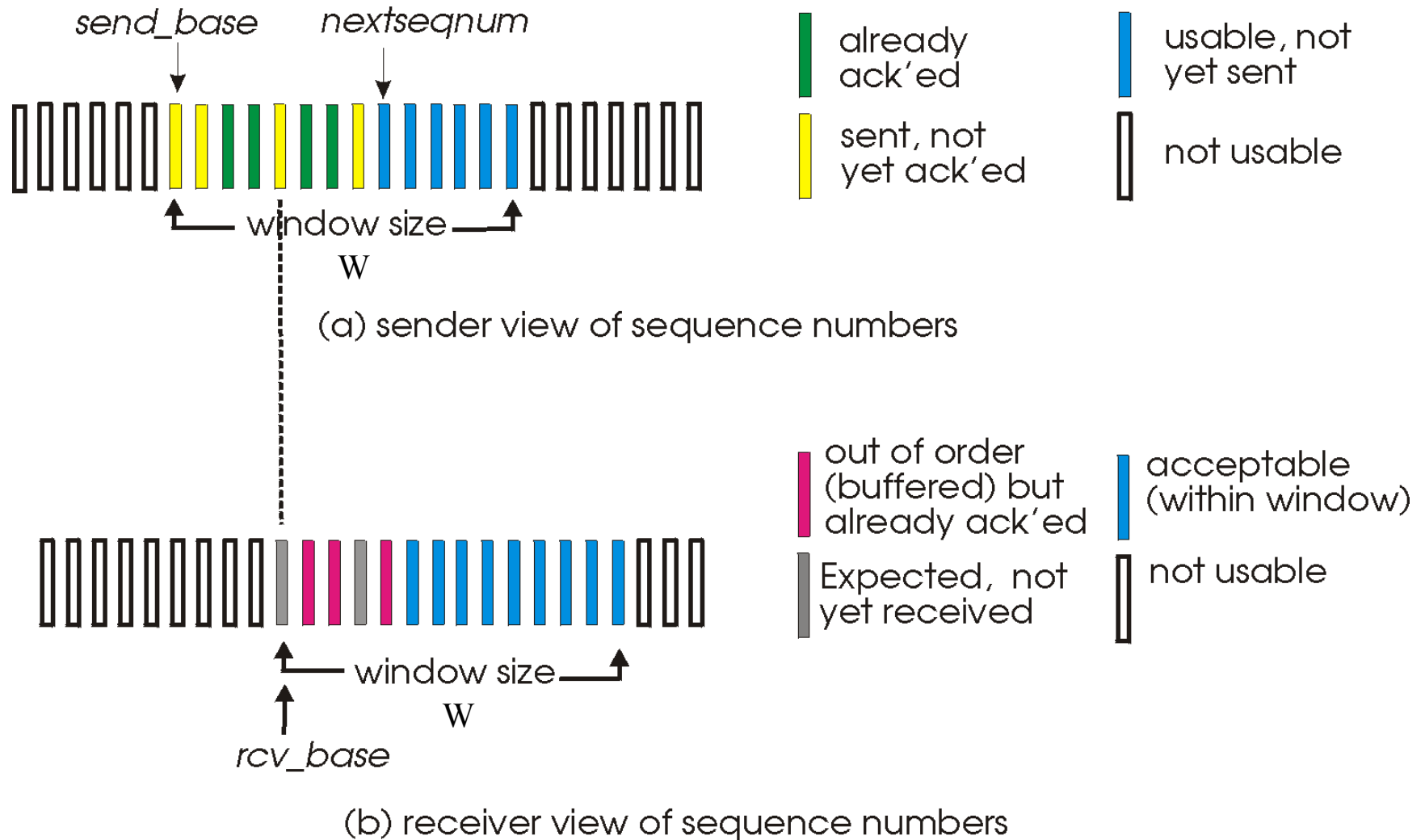
Only state: `expectedseqnum`

❑ out-of-order pkt:
- o  discard (don't buffer) -> no receiver buffering!
- o  re-ACK pkt with highest in-order seq #
- o  may generate duplicate ACKs

# Selective Repeat

- ❑ Sender window
  - ○ Window size W: W consecutive unACKed seq #'s
- ❑ Receiver *individually* acknowledges correctly received pkts
  - ○ buffers out-of-order pkts, for eventual in-order delivery to upper layer
  - ○ ACK(n) means received packet with seq# n only
  - ○ buffer size at receiver: window size
- ❑ Sender only resends pkts for which ACK not received
  - ○ sender timer for each unACKed pkt

# Selective Repeat: Sender, Receiver Windows



(a) sender view of sequence numbers

(b) receiver view of sequence numbers

# Selective Repeat

**sender**

**data from above :**
- unACKed packets is less than window size W, send; otherwise block app.

**timeout(n):**
- resend pkt n, restart timer

**ACK(n)** in [sendbase,sendbase+W-1]:
- mark pkt n as received
- update sendbase to the first packet unACKed

**receiver**

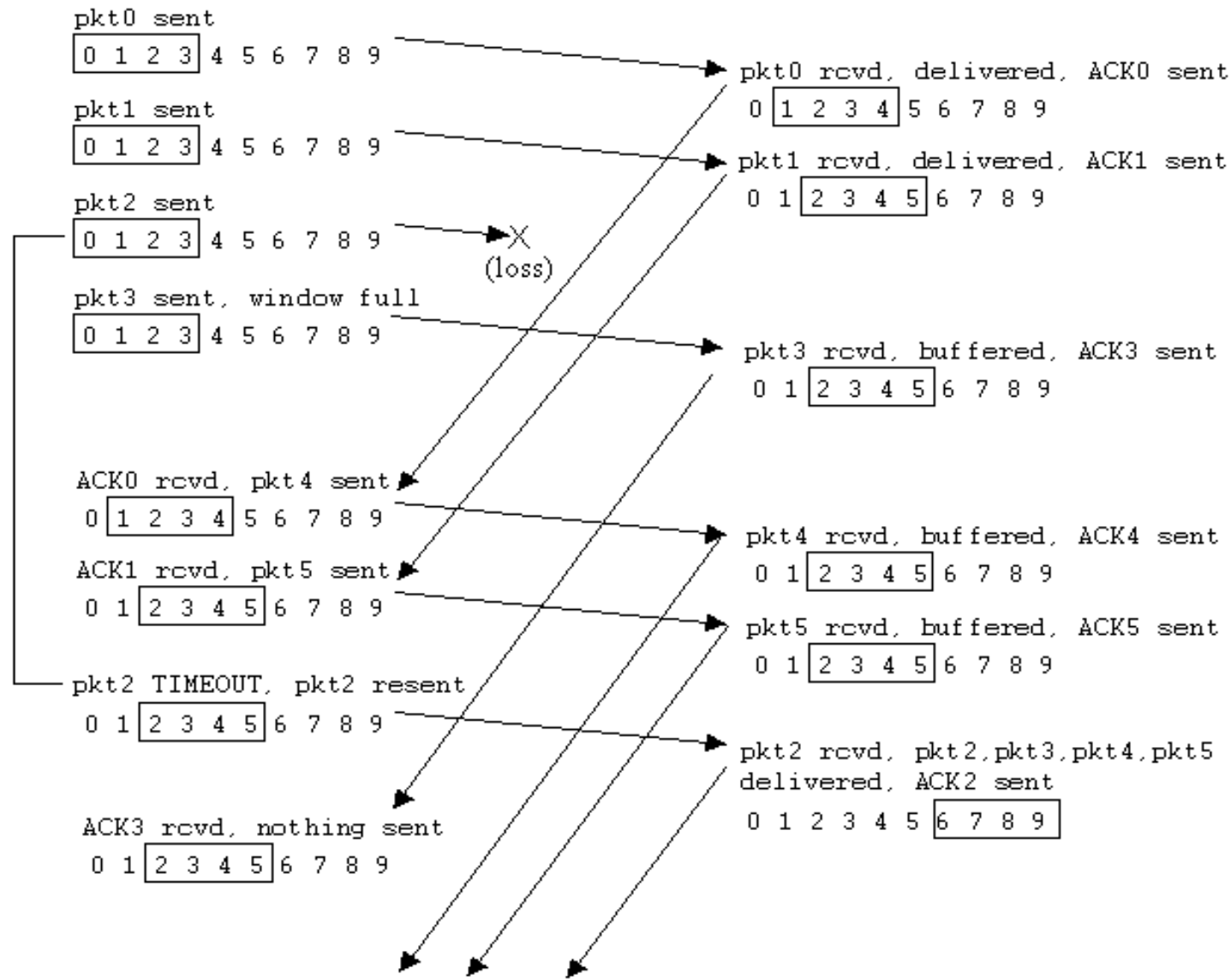**pkt n in** [rcvbase, rcvbase+W-1]
- send ACK(n)
- if (out-of-order)
  mark and buffer pkt n
  else /*in-order*/
  deliver any in-order packets

**otherwise:**
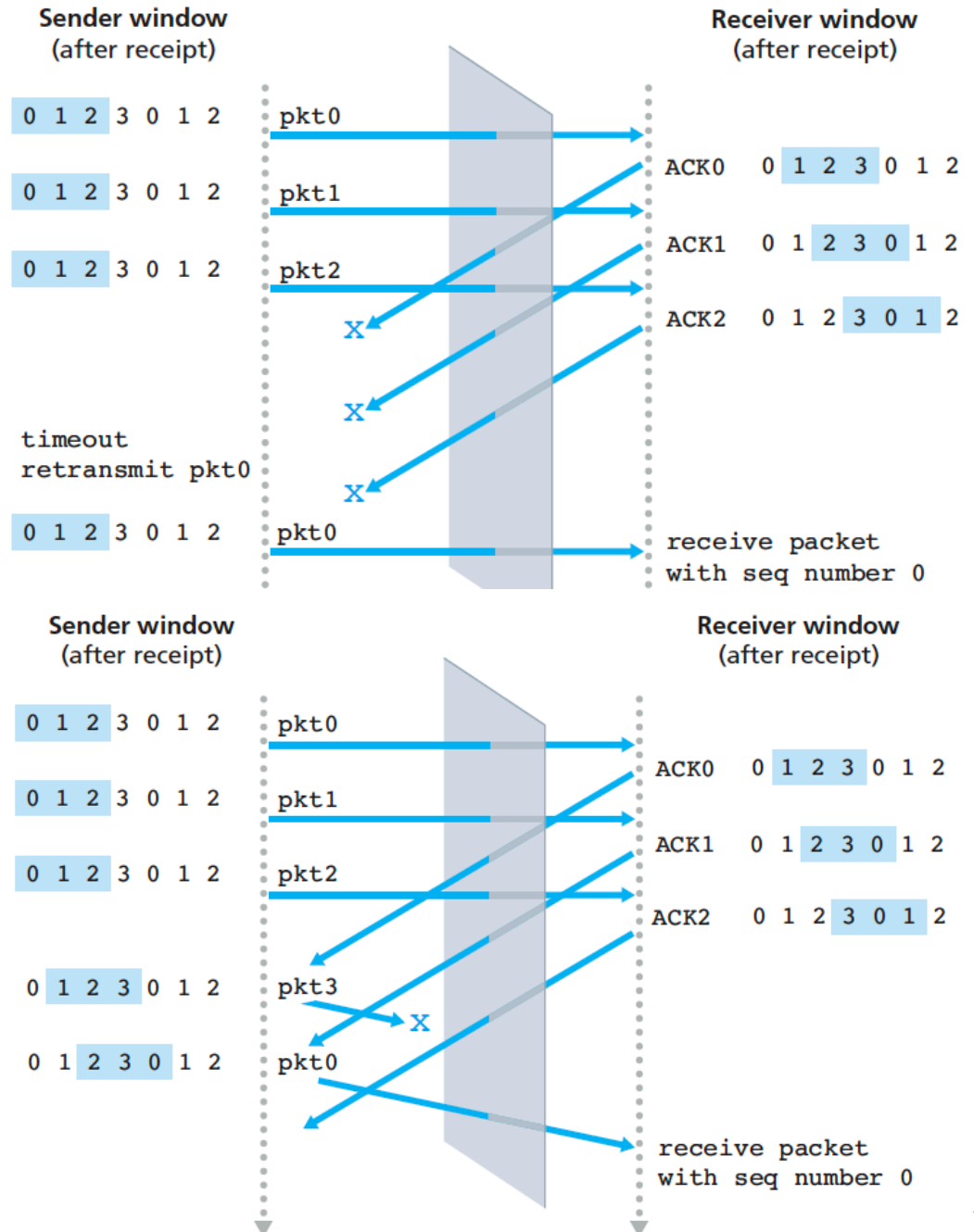- ignore

# Selective Repeat in Action



pkt0 sent
`0 1 2 3` 4 5 6 7 8 9

pkt1 sent
`0 1 2 3` 4 5 6 7 8 9

pkt2 sent
`0 1 2 3` 4 5 6 7 8 9

pkt3 sent, window full
`0 1 2 3` 4 5 6 7 8 9

ACK0 rcvd, pkt4 sent
0 `1 2 3 4` 5 6 7 8 9

ACK1 rcvd, pkt5 sent
0 1 `2 3 4 5` 6 7 8 9

pkt2 TIMEOUT, pkt2 resent
0 1 `2 3 4 5` 6 7 8 9

ACK3 rcvd, nothing sent
0 1 `2 3 4 5` 6 7 8 9

X
(loss)

pkt0 rcvd, delivered, ACK0 sent
0 `1 2 3 4` 5 6 7 8 9

pkt1 rcvd, delivered, ACK1 sent
0 1 `2 3 4 5` 6 7 8 9

pkt3 rcvd, buffered, ACK3 sent
0 1 `2 3 4 5` 6 7 8 9

pkt4 rcvd, buffered, ACK4 sent
0 1 `2 3 4 5` 6 7 8 9

pkt5 rcvd, buffered, ACK5 sent
0 1 `2 3 4 5` 6 7 8 9

pkt2 rcvd, pkt2,pkt3,pkt4,pkt5
delivered, ACK2 sent
0 1 2 3 4 5 `6 7 8 9`

# Discussion: Efficiency of Selective Repeat

❑ Assume window size W

❑ Assume each packet is lost with probability p

❑ On average, how many packets do we send for each data packet received?

# Selective Repeat: Seq# Ambiguity

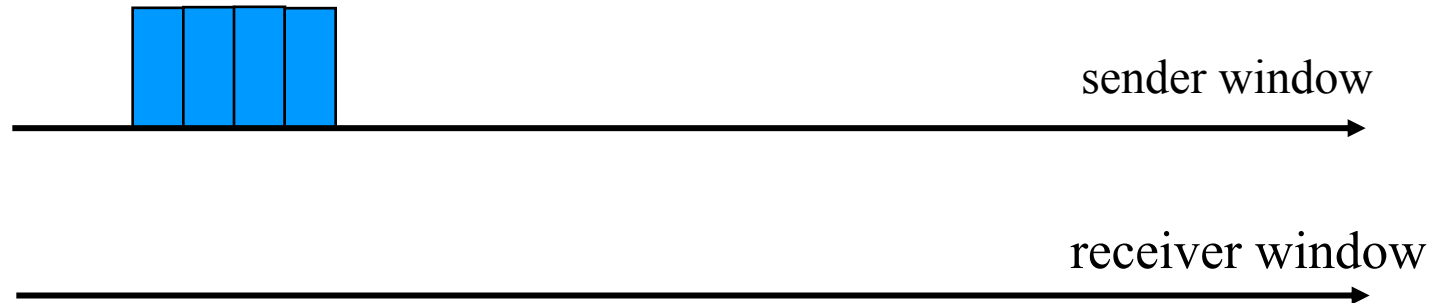Example:
- ❑ seq #'s: 0, 1, 2, 3
- ❑ window size=3

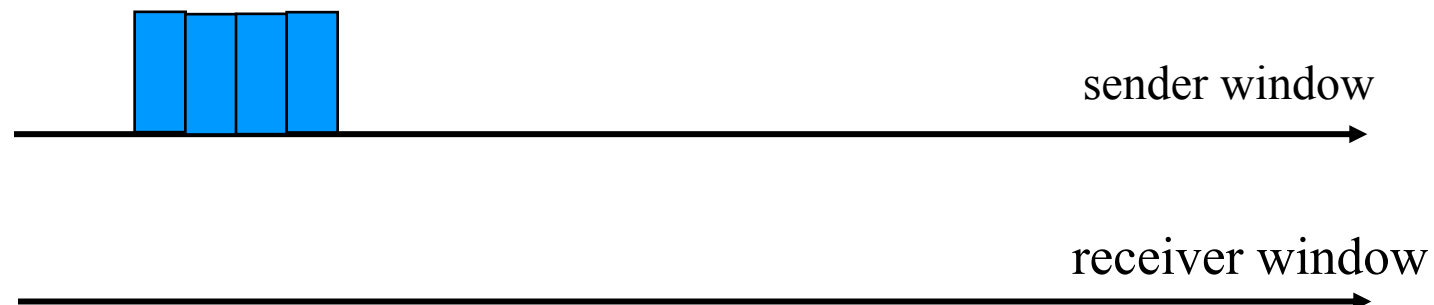- ❑ Error: incorrectly passes duplicate data as new.

# State Invariant: Window Location

❑ Go-back-n (GBN)
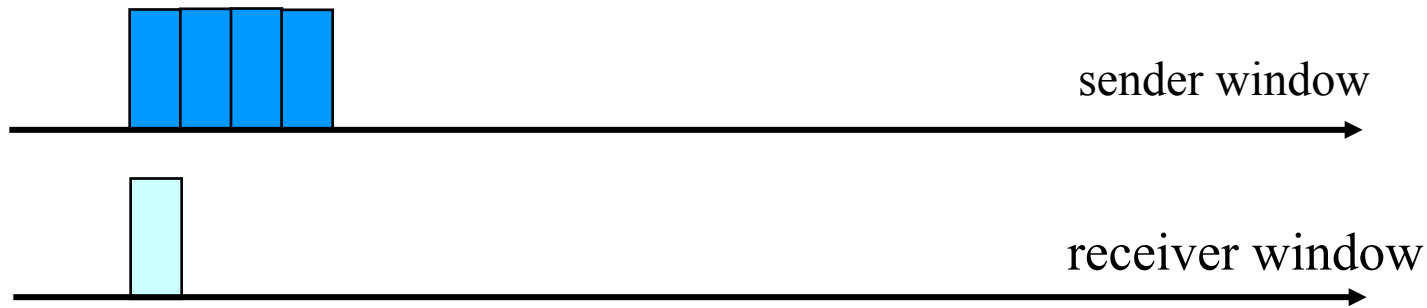
sender window

receiver window
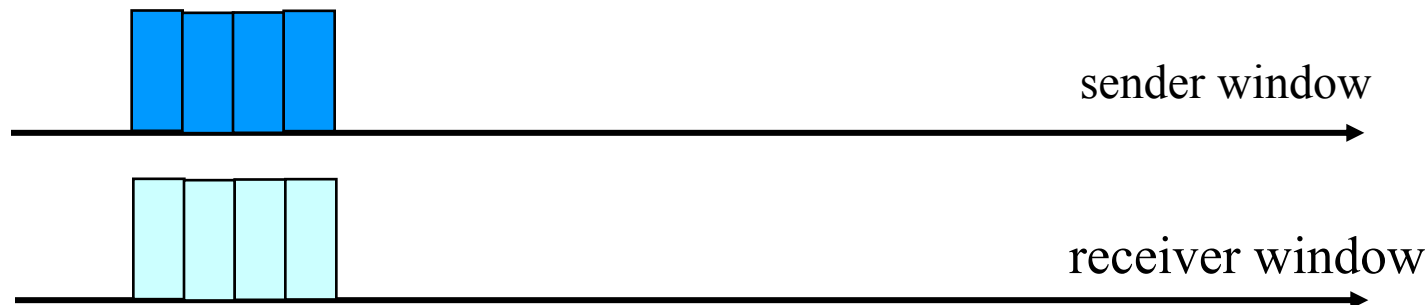
❑ Selective repeat (SR)

sender window

receiver window

12

# Window Location

❑ Go-back-n (GBN)

sender window

receiver window

❑ Selective repeat (SR)

sender window

receiver window

13

# Selective Repeat

## sender

**data from above :**

- ❑ unACKed packets is less than window size W, send; otherwise block app.

**timeout(n):**

- ❑ resend pkt n, restart timer

**ACK(n)** in [sendbase,sendbase+W-1]:

- ❑ mark pkt n as received
- ❑ update sendbase to the first packet unACKed

## receiver

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

- ❑ send ACK(n)
- ❑ if (out-of-order)
    mark and buffer pkt n
  else /*in-order*/
    deliver any in-order packets

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

- ❑ send ACK(n)

**otherwise:**

- ❑ ignore

14

# Sliding Window Protocols: Go-back-n and Selective Repeat

|  | Go-back-n | Selective Repeat |
|---|---|---|
| data bandwidth: sender to receiver (avg. number of times a pkt is transmitted) | Less efficient $\frac{1-p+pw}{1-p}$ | More efficient $\frac{1}{1-p}$ |
| ACK bandwidth (receiver to sender) | More efficient | Less efficient |
| Relationship between M (the number of seq#) and W (window size) | M > W | M ≥ 2W |
| Buffer size at receiver | 1 | W |
| Complexity | Simpler | More complex |

p: the loss rate of a packet; M: number of seq# (e.g., 3 bit M = 8); W: window size

# Outline

❑ Admin and Recap

❑ Reliable data transfer
- o perfect channel
- o channel with bit errors
- o channel with bit errors and losses
- o sliding window: reliability with throughput

➢ *TCP reliability*
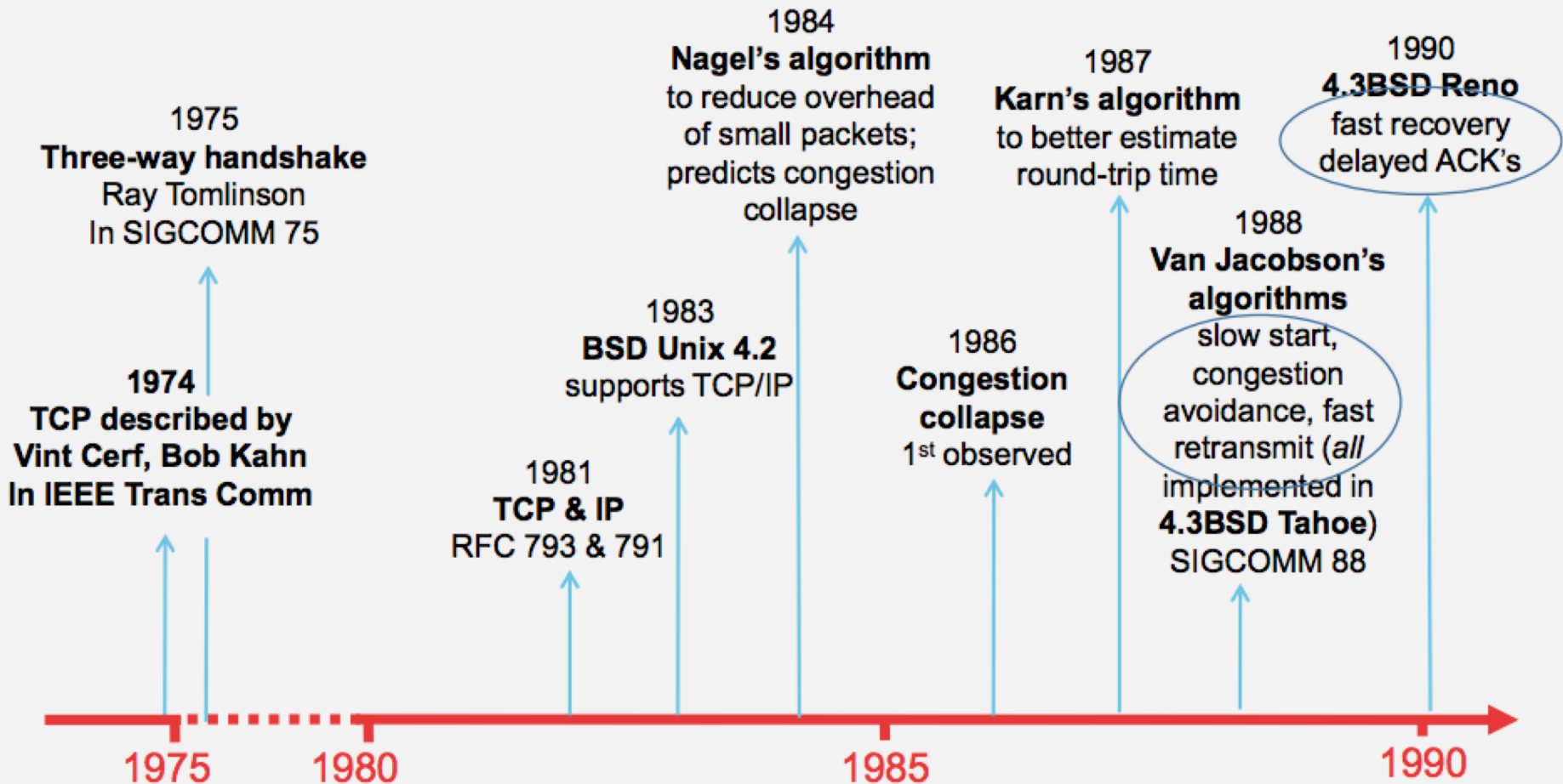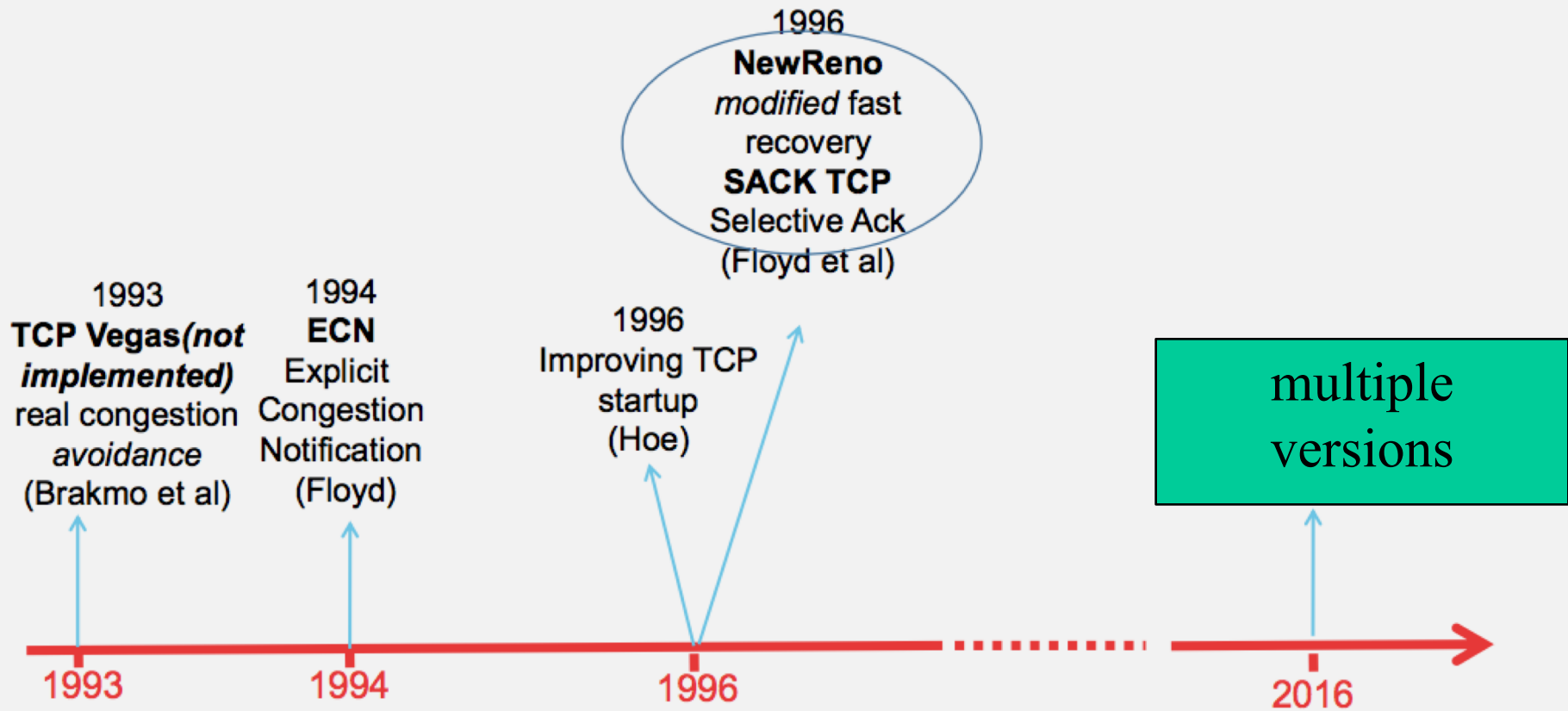
# TCP: Overview     RFCs: 793, 1122, 1323, 2018, 2581

❑ Point-to-point reliability: one sender, one receiver

❑ Flow controlled and congestion controlled

# Evolution of TCP



**1975**
**Three-way handshake**
Ray Tomlinson
In SIGCOMM 75

**1974**
**TCP described by**
**Vint Cerf, Bob Kahn**
**In IEEE Trans Comm**

**1981**
**TCP & IP**
RFC 793 & 791

**1983**
**BSD Unix 4.2**
supports TCP/IP

**1984**
**Nagel's algorithm**
to reduce overhead
of small packets;
predicts congestion
collapse

**1986**
**Congestion**
**collapse**
1st observed

**1987**
**Karn's algorithm**
to better estimate
round-trip time

**1988**
**Van Jacobson's**
**algorithms**
slow start,
congestion
avoidance, fast
retransmit (*all*
implemented in
**4.3BSD Tahoe)**
SIGCOMM 88

**1990**
**4.3BSD Reno**
fast recovery
delayed ACK's

1975     1980     1985     1990

*Source: http://webcourse.cs.technion.ac.il/236341/Winter2015-2016/ho/WCFiles/Tutorial10.pdf*

18

# Evolution of TCP

1996
**NewReno**
*modified* fast
recovery
**SACK TCP**
Selective Ack
(Floyd et al)

1993
**TCP Vegas** *(not implemented)*
real congestion
*avoidance*
(Brakmo et al)

1994
**ECN**
Explicit
Congestion
Notification
(Floyd)

1996
Improving TCP
startup
(Hoe)

multiple
versions

1993    1994    1996    2016

*Source: http://webcourse.cs.technion.ac.il/236341/Winter2015-2016/ho/WCFiles/Tutorial10.pdf*

19

# TCP Reliable Data Transfer
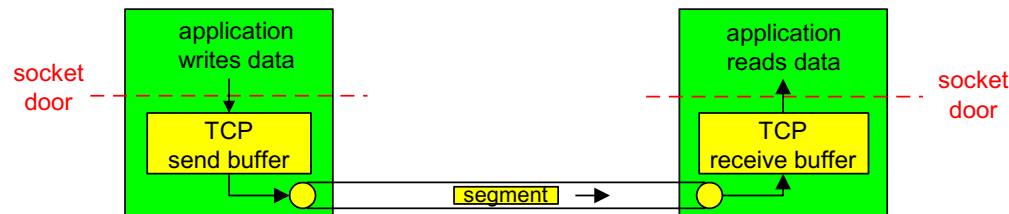
❑ **Connection-oriented:**
  - connection management
    - setup (exchange of control msgs) init's sender, receiver state before data exchange
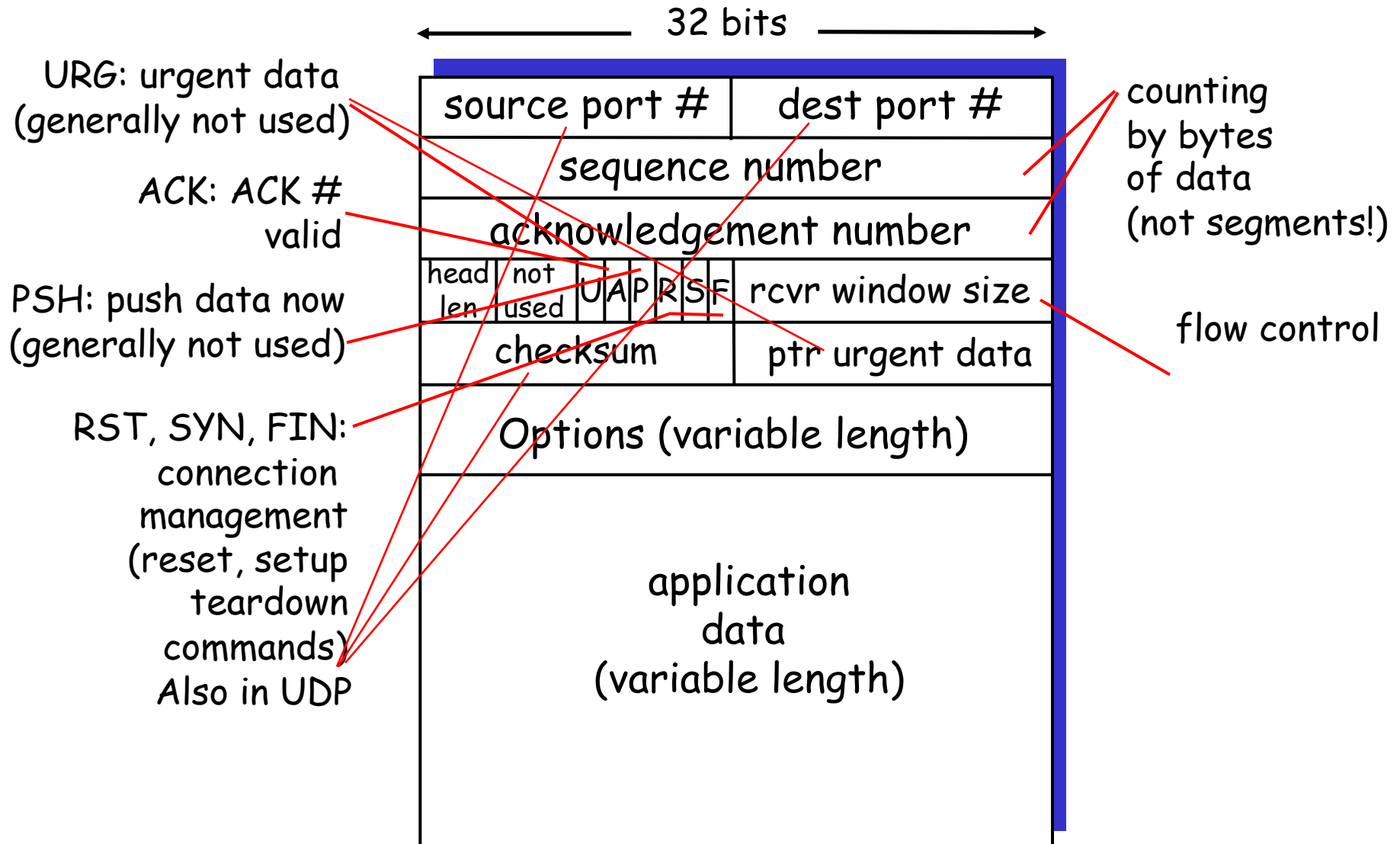    - close

❑ **Full duplex data:**
  - bi-directional data flow in same connection

❑ **A sliding window protocol**
  - a combination of go-back-n and selective repeat:
    - send & receive buffers
    - cumulative acks
    - TCP uses a single retransmission timer
    - do not retransmit all packets upon timeout

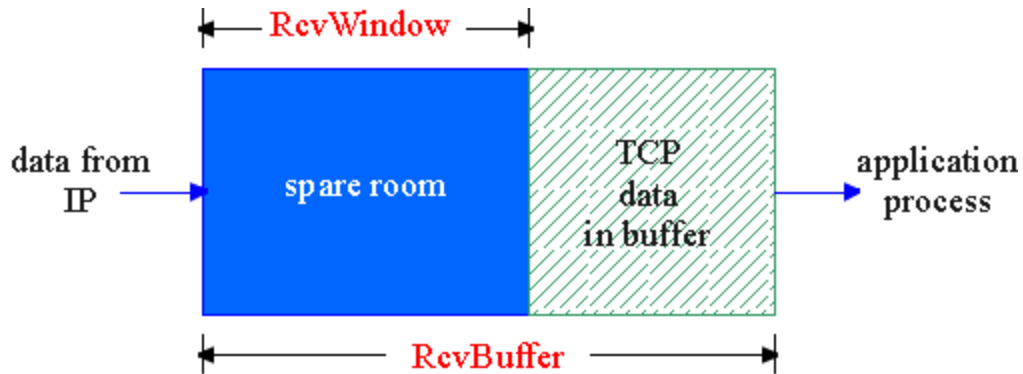application writes data

socket door

TCP send buffer

segment

socket door

application reads data

TCP receive buffer

# TCP Segment Structure



URG: urgent data
(generally not used)

ACK: ACK #
valid

PSH: push data now
(generally not used)

RST, SYN, FIN:
connection
management
(reset, setup
teardown
commands)
Also in UDP

32 bits

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

| head len | not used | U | A | P | R | S | F | rcvr window size |

| checksum | ptr urgent data |

Options (variable length)

application
data
(variable length)

counting
by bytes
of data
(not segments!)

flow control

# Outline

❑ Admin and Recap

❑ Reliable data transfer
  - ○ perfect channel
  - ○ channel with bit errors
  - ○ channel with bit errors and losses
  - ○ sliding window: reliability with throughput

❑ TCP reliability
  - ➢ *data seq#, ack, buffering*

# Flow Control
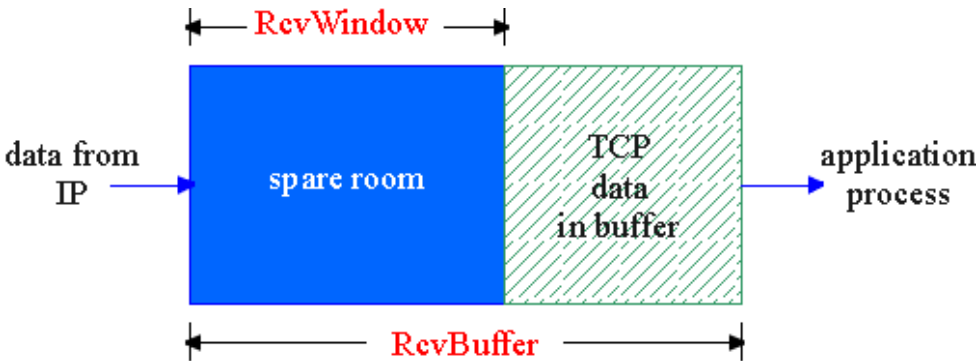
☐ receive side of a connection has a receive buffer:



☐ app process may be slow at reading from buffer

**flow control**
sender won't overflow receiver's buffer by transmitting too much, too fast

☐ speed-matching service: matching the send rate to the receiving app's drain rate

# TCP Flow Control: How it Works



❑ spare room in buffer
= **RcvWindow**

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

| head len | not used | U | A | P | R | S | F | rcvr window size |
|---|---|---|---|---|---|---|---|---|

| checksum | ptr urgent data |
|---|---|

Options (variable length)

application
data
(variable length)

# TCP Seq. #'s and ACKs

Seq. #'s:
- byte stream "number" of first byte in segment's data

ACKs:
- seq # of next byte expected from other side
- cumulative ACK in standard header
- selective ACK in options

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=42, ACK=80

time

simple telnet scenario

# TCP Send/Ack Optimizations

❑ TCP includes many tune/optimizations, e.g.,

  o the "small-packet problem": sender sends a lot of small packets (e.g., telnet one char at a time)

   • Nagle's algorithm: do not send data if there is small amount of data in send buffer and there is an unack'd segment

  o the "ack inefficiency" problem: receiver sends too many ACKs, no chance of combing ACK with data

   • Delayed ack to reduce # of ACKs/combine ACK with reply

# TCP Receiver 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 |

# Outline

❑ Admin and Recap

❑ Reliable data transfer
  - o   perfect channel
  - o   channel with bit errors
  - o   channel with bit errors and losses
  - o   sliding window: reliability with throughput

❑ TCP reliability
  - o   data seq#, ack, buffering
  - ➢ *timeout realization*

# TCP Reliable Data Transfer

❑ Basic structure: sliding window protocol

❑ Remaining issue: How to determine the "right" parameters?

    ○ timeout value?

    ○ sliding window size?

# History

❑ Key parameters for TCP in mid-1980s
- ○ fixed window size W
- ○ timeout value = 2 RTT

❑ Network collapse in the mid-1980s
- ○ UCB ←→ LBL throughput dropped by 1000X !

❑ The intuition was that the collapse was caused by wrong parameters…

# Timeout: Cost of Timeout Param

Why is good timeout value important?
- ❑ too short
  - ○ premature timeout
  - ○ unnecessary retransmissions; many duplicates
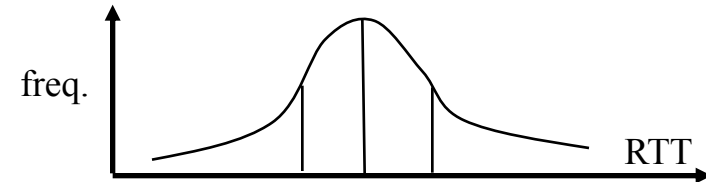
- ❑ too long
  - ○ slow reaction to segment loss

Q: Is it possible to set Timeout as a constant?

Q: Any problem w/ the early approach: Timeout = 2 RTT

# Setting Timeout

**Problem:**

❑ Ideally, we set timeout = RTT,
 but RTT is not a fixed value
 =>
 using the average of **RTT** will generate
 many timeouts due to network variations

❑ Possibility: using the average/median of RTT

❑ Issue: this will generate many timeouts due to network variations

freq.

RTT

**Solution:**

❑ **Set Timeout RTO = avg +** "safety margin" based on variation

TCP approach:

**Timeout = EstRTT + 4 * DevRTT**

32

# Compute EstRTT and DevRTT

❑ Exponential weighted moving average (EWMA)
 ○ influence of past sample decreases exponentially fast

```
EstRTT = (1-alpha)*EstRTT + alpha*SampleRTT
```

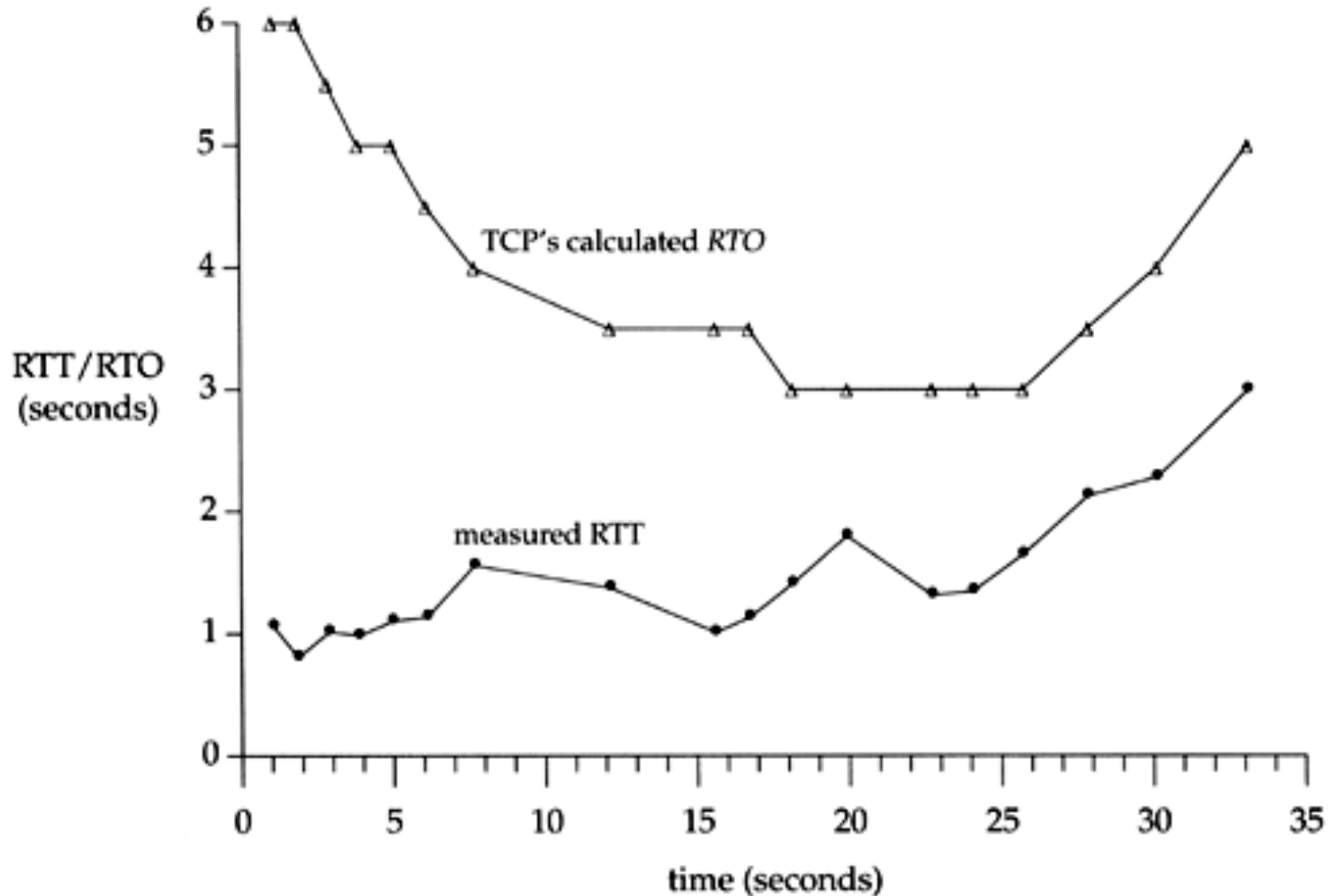– **SampleRTT:** measured time from segment transmission until ACK receipt

- typical value: `alpha` = 0.125

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



```
DevRTT = (1-beta)*DevRTT + beta|SampleRTT-EstRTT|

(typically, beta = 0.25)
```

33

# An Example TCP Session

# Fast Retransmit

❑ Issue: Timeout period often relatively long:
  o long delay before resending lost packet
❑ Question: Can we detect loss faster than RTT?

❑ Detect lost segments via duplicate ACKs
  o sender often sends many segments back-to-back
  o if segment is lost, there will likely be many duplicate ACKs

❑ If sender receives 3 ACKs for the same data, it supposes that segment after ACKed data was lost:
  o resend segment before timer expires

# Triple Duplicate Ack

Packets

| 1 | 2 | 3 | 4 | 5 | 6 | 7 |

Acknowledgements (waiting seq#)

| 2 | 3 | 4 | 4 | 4 | 4 |

# Fast Retransmit:

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
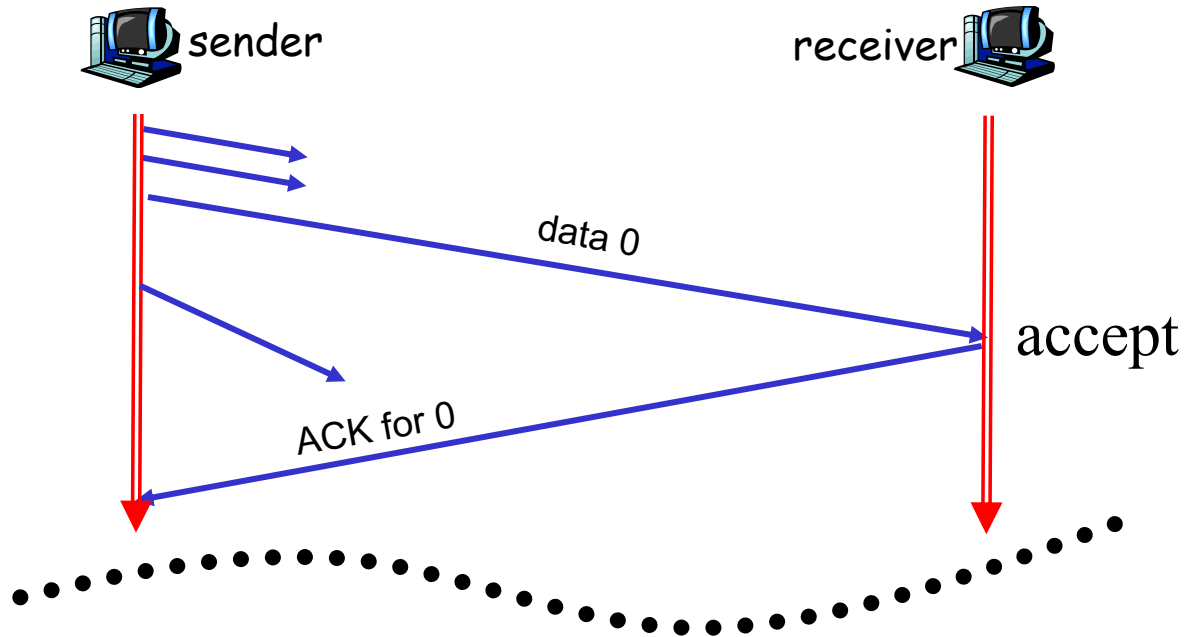
37

# TCP: reliable data transfer

Simplified TCP sender

```
00   sendbase = initial_sequence number agreed by TWH
01   nextseqnum = initial_sequence number by TWH
02   loop (forever) {
03     switch(event)
04     event: data received from application above
05           if (window allows send)
06              create TCP segment with sequence number nextseqnum
06              if (no timer) start timer
07              pass segment to IP
08              nextseqnum = nextseqnum + length(data)
           else put packet in buffer
09     event: timer timeout for sendbase
10        retransmit segment
11        compute new timeout interval
12        restart timer
13     event: ACK received, with ACK field value of y
14        if (y > sendbase) { /* cumulative ACK of all data up to y */
15            cancel the timer for sendbase
16            sendbase = y
17            if (no timer and packet pending) start timer for new sendbase
17            while (there are segments and window allow)
18               sent a segment;
18        }
19        else { /* y==sendbase, duplicate ACK for already ACKed segment */
20            increment number of duplicate ACKs received for y
21            if (number of duplicate ACKS received for y == 3) {
22                /* TCP fast retransmit */
23                resend segment with sequence number y
24                restart timer for segment y
25            }
26     } /* end of loop forever */
```
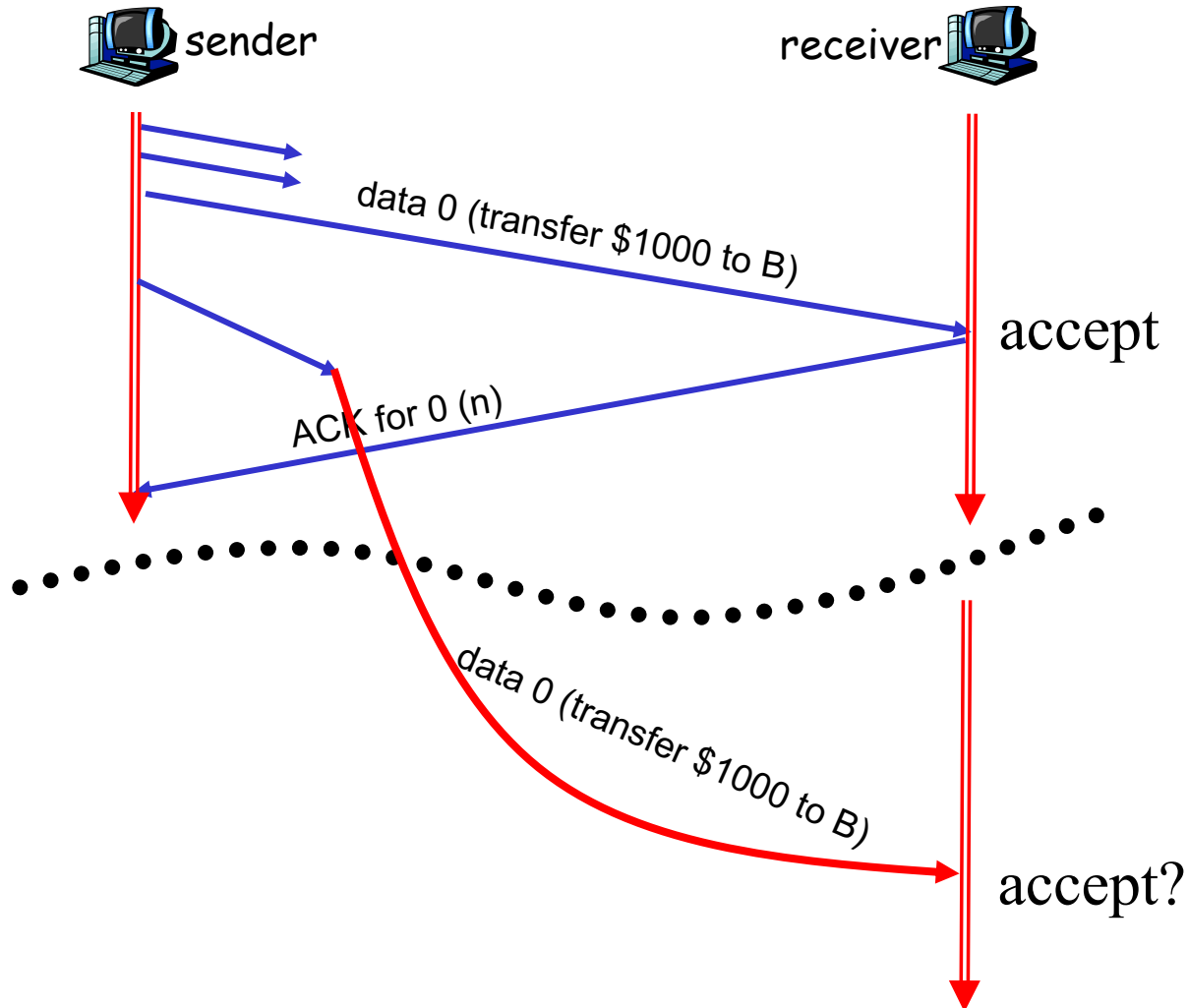
# Outline

❑ Admin and Recap

❑ Reliable data transfer
  o perfect channel
  o channel with bit errors
  o channel with bit errors and losses
  o sliding window: reliability with throughput

❑ TCP reliability
  o data seq#, ack, buffering
  o timeout realization
  ➢ *connection management*

sender

receiver

data 0

accept

ACK for 0

# Why Connection Setup/When to Accept (Safely Deliver) First Packet?

sender

receiver

data 0 (transfer $1000 to B)

accept

ACK for 0 (n)

data 0 (transfer $1000 to B)

accept?

# Transport "Safe-Setup" Principle

❑ A general safety principle for a receiver R to accept a message from a sender S is the general "authentication" principle, which consists of two conditions:

Transport authentication principle:
- [p1] Receiver can be sure that what Sender says is **fresh**
- [p2] Receiver receives something that *only* Sender can say

We first assume a secure setting: no malicious attacks.

Exercise: Techniques to allow a receiver to check for freshness (e.g., add a time stamp)?