# Network Transport Layer: Reliable Data Transfer; TCP

Qiao Xiang

https://qiaoxiang.me/courses/cnns-xmuf22/index.shtml

10/25/2022

This deck of slides are heavily based on CPSC 433/533 at Yale University, by courtesy of Dr. Y. Richard Yang.
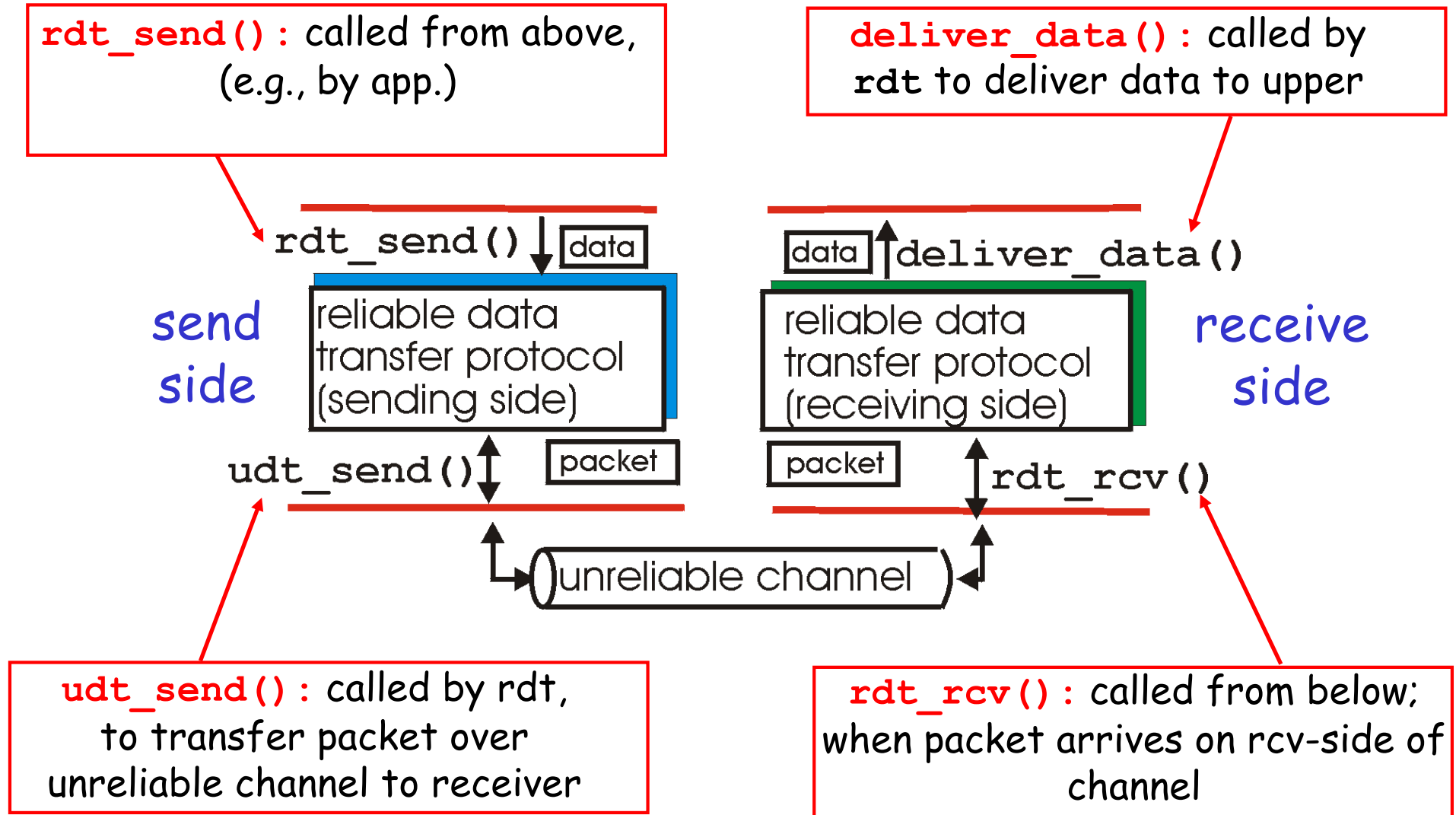
# Outline

- ❑ Admin and recap
- ❑ Reliable data transfer

# Admin

- Guest lecture on Oct. 27 class (morning)
  - Guest lecturer: Ennan Zhai (Alibaba)
  - Lecture title: Intent-Based Network in Alibaba

# Recap: Reliable Data Transfer Context

`rdt_send():` called from above, (e.g., by app.)

`deliver_data():` called by `rdt` to deliver data to upper

send side

`rdt_send()` | data

reliable data transfer protocol (sending side)

receive side

data | `deliver_data()`

reliable data transfer protocol (receiving side)

`udt_send()` | packet

packet | `rdt_rcv()`

unreliable channel

`udt_send():` called by rdt, to transfer packet over unreliable channel to receiver

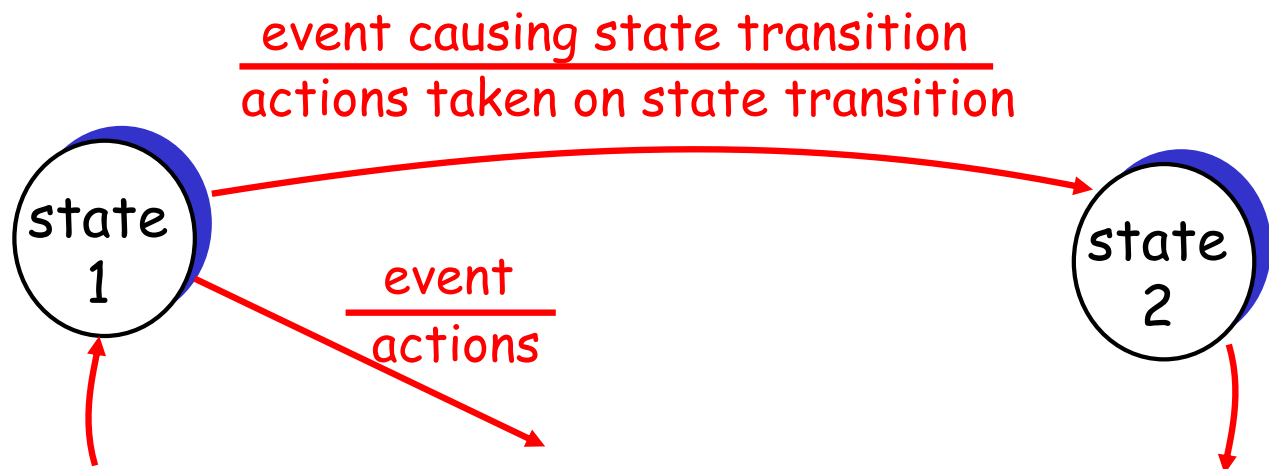`rdt_rcv():` called from below; when packet arrives on rcv-side of channel
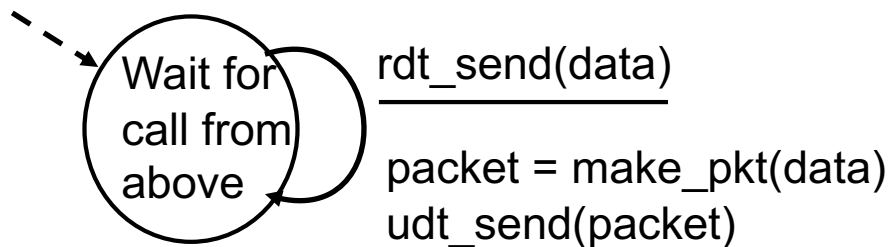
# Reliable Data Transfer: Getting Started

We'll:

❑ incrementally develop sender, receiver sides of reliable data transfer protocol (rdt)

❑ consider only unidirectional data transfer
  o but control info will flow on both directions !

❑ use finite state machines (FSM) to specify sender, receiver

state: when in this "state" next state uniquely determined by next event

event causing state transition
actions taken on state transition

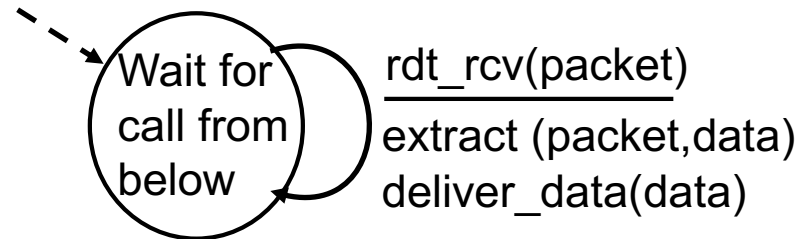event
actions

state 1

state 2

# Rdt1.0: reliable transfer over a reliable channel

❑ separate FSMs for sender, receiver:
  ○ sender sends data into underlying channel
  ○ receiver reads data from underlying channel

Wait for call from above

rdt_send(data)
_____

packet = make_pkt(data)
udt_send(packet)

Wait for call from below

rdt_rcv(packet)
_____
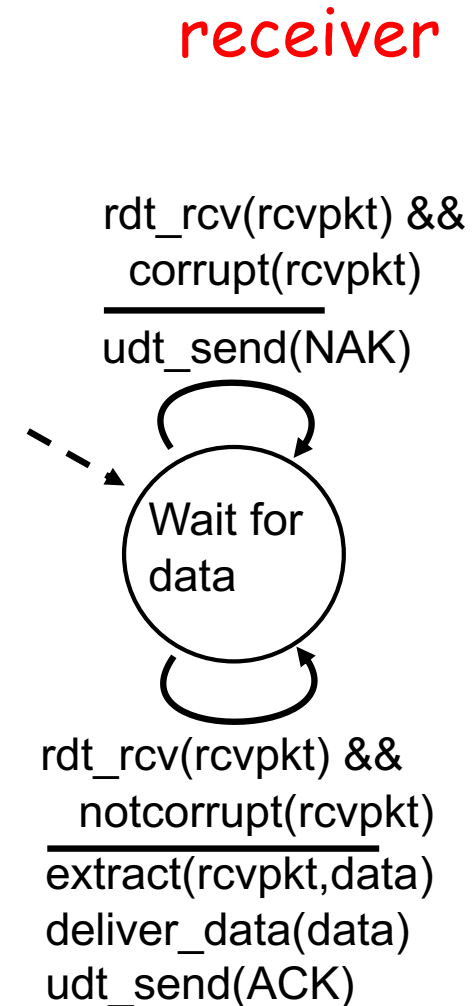
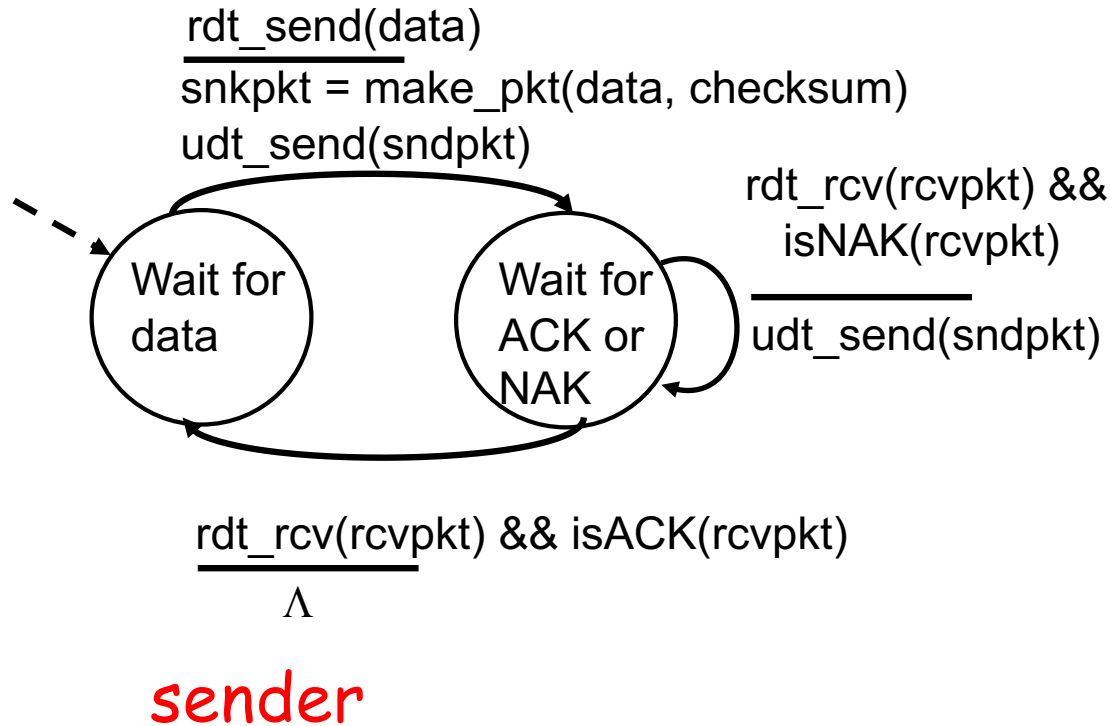extract (packet,data)
deliver_data(data)

sender

receiver

Exercise: Prove correctness of Rdt1.0.

**Correctness**: for every single packet, one and only one copy is received by receiver correctly (no error) and in-order
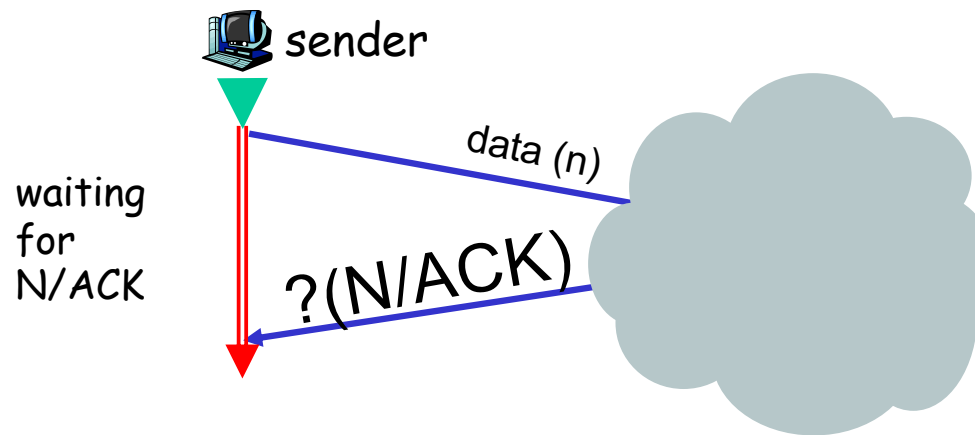
# Recap: rdt2.0: Reliability allowing only Data Msg Corruption

**sender**

rdt_send(data)
_____
snkpkt = make_pkt(data, checksum)
udt_send(sndpkt)

Wait for data

Wait for ACK or NAK

rdt_rcv(rcvpkt) && isNAK(rcvpkt)
_____
udt_send(sndpkt)

rdt_rcv(rcvpkt) && isACK(rcvpkt)
_____
Λ

**receiver**

rdt_rcv(rcvpkt) && corrupt(rcvpkt)
_____
udt_send(NAK)

Wait for data

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
_____
extract(rcvpkt,data)
deliver_data(data)
udt_send(ACK)

# rdt2.0 is Incomplete!

## What happens if ACK/NAK corrupted?

❑ Although sender receives feedback, but doesn't know what happened at receiver!

# rdt2.1b: Summary

<u>Sender:</u>

❑ seq # added to pkt

❑ must check if received ACK/NAK corrupted


<u>Receiver:</u>

❑ must check if received packet is duplicate

  o by checking if the packet has the expected pkt seq #

# rdt2.1c: Summary

Sender:

❑ state must "remember" whether "current" pkt has 0 or 1 seq. #

Receiver:

❑ must check if received packet is duplicate
   o state indicates whether 0 or 1 is expected pkt seq #

# rdt2.2: a NAK-free protocol

❑ Same functionality as rdt2.1c, using ACKs only

❑ Instead of NAK, receiver sends ACK for last pkt received OK
  ○ receiver must *explicitly* include seq # of pkt being ACKed

❑ Duplicate ACK at sender results in same action as NAK: *retransmit current pkt*

# rdt3.0: Channels with Errors _and_ Loss

**New assumption:**
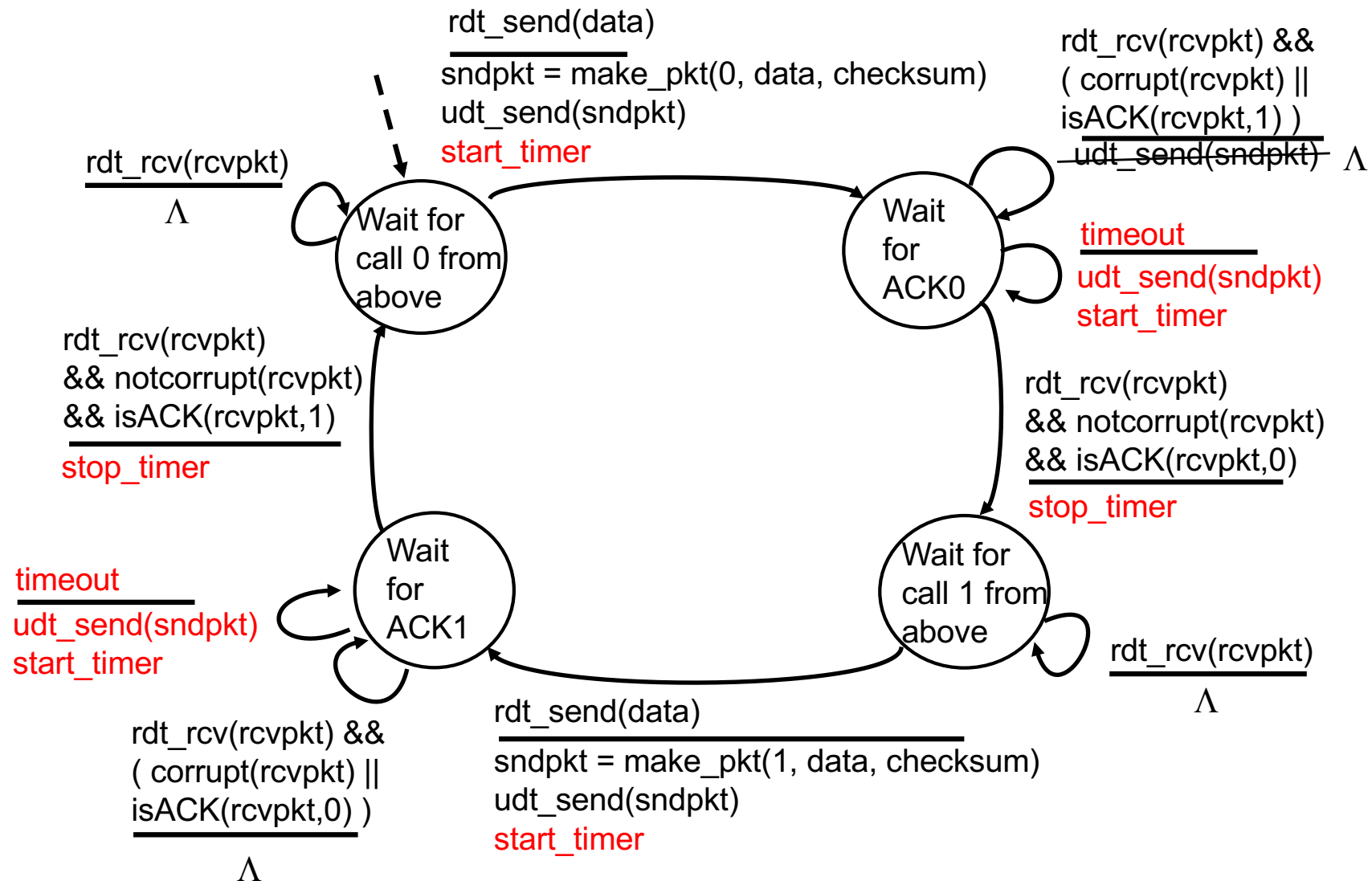underlying channel can also lose packets (data or ACKs)

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

**Q:** Does rdt2.2 work under losses?

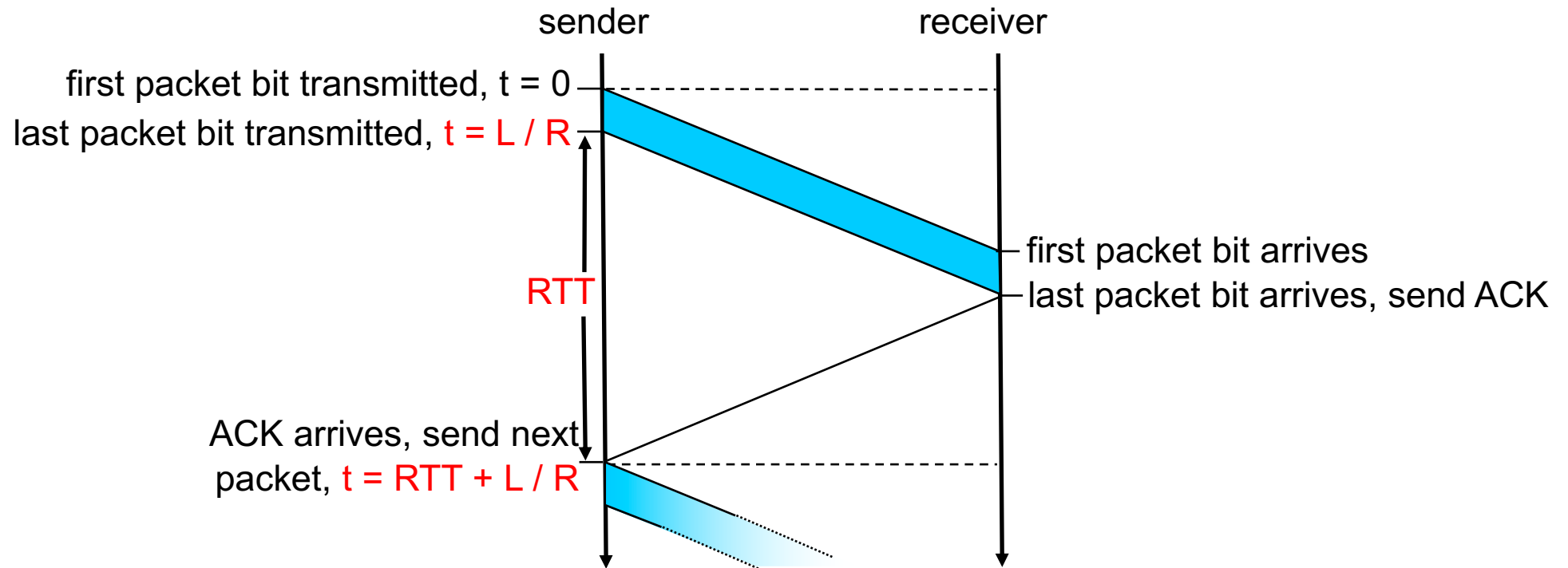**Approach:** sender waits "reasonable" amount of time for ACK

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

# rdt3.0 Sender



rdt_send(data)
sndpkt = make_pkt(0, data, checksum)
udt_send(sndpkt)
start_timer

rdt_rcv(rcvpkt) &&
( corrupt(rcvpkt) ||
isACK(rcvpkt,1) )
udt_send(sndpkt) Λ

rdt_rcv(rcvpkt)
Λ

timeout
udt_send(sndpkt)
start_timer

Wait for
call 0 from
above

Wait
for
ACK0

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& isACK(rcvpkt,1)
stop_timer

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& isACK(rcvpkt,0)
stop_timer

timeout
udt_send(sndpkt)
start_timer

Wait
for
ACK1

Wait for
call 1 from
above

rdt_rcv(rcvpkt)
Λ

rdt_rcv(rcvpkt) &&
( corrupt(rcvpkt) ||
isACK(rcvpkt,0) )
Λ

rdt_send(data)
sndpkt = make_pkt(1, data, checksum)
udt_send(sndpkt)
start_timer

13

# rdt3.0: Stop-and-Wait Performance

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

What is $U_{sender}$: utilization – fraction of time link busy sending?

Assume: 1 Gbps link, 15 ms e-e prop. delay, 1KB packet

# Performance of rdt3.0

❑ rdt3.0 works, but performance stinks

❑ Example: 1 Gbps link, 15 ms e-e prop. delay, 1KB packet:

$$T_{transmit} = \frac{L \text{ (packet length in bits)}}{R \text{ (transmission rate, bps)}} = \frac{8kb/pkt}{10^{**}9 \text{ b/sec}} = 8 \text{ microsec}$$

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

o   1KB pkt every 30 msec -> 33kB/sec throughput over 1 Gbps link

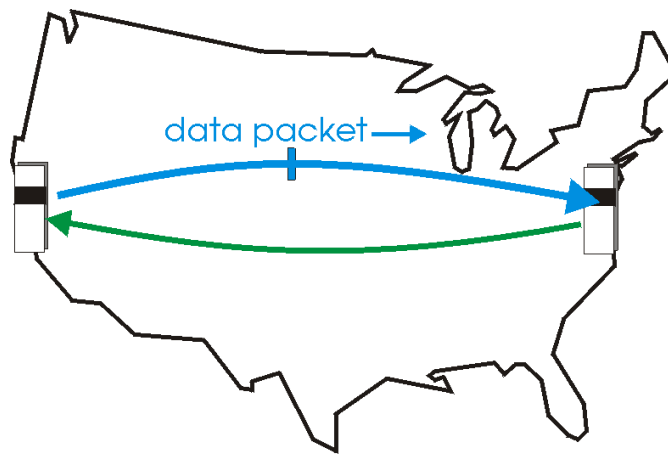o   network protocol limits use of physical resources !

# A Summary of Questions

- How to improve the performance of rdt3.0?

- What if there are reordering and duplication?

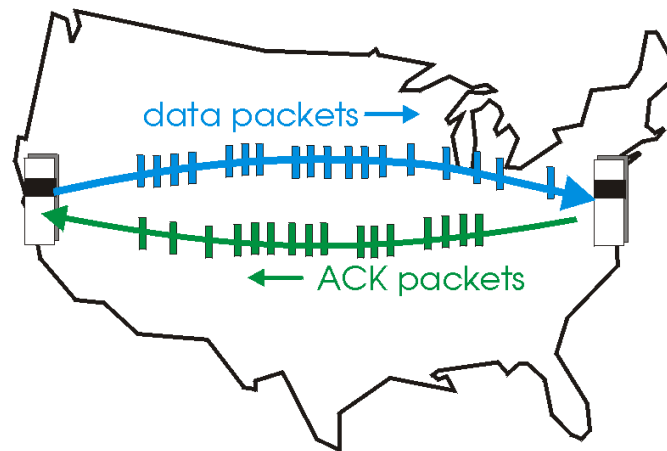- How to determine the "right" timeout value?

# Sliding Window Protocols: Pipelining

**Pipelining:** sender allows multiple, "in-flight", yet-to-be-acknowledged pkts

- o range of sequence numbers must be increased
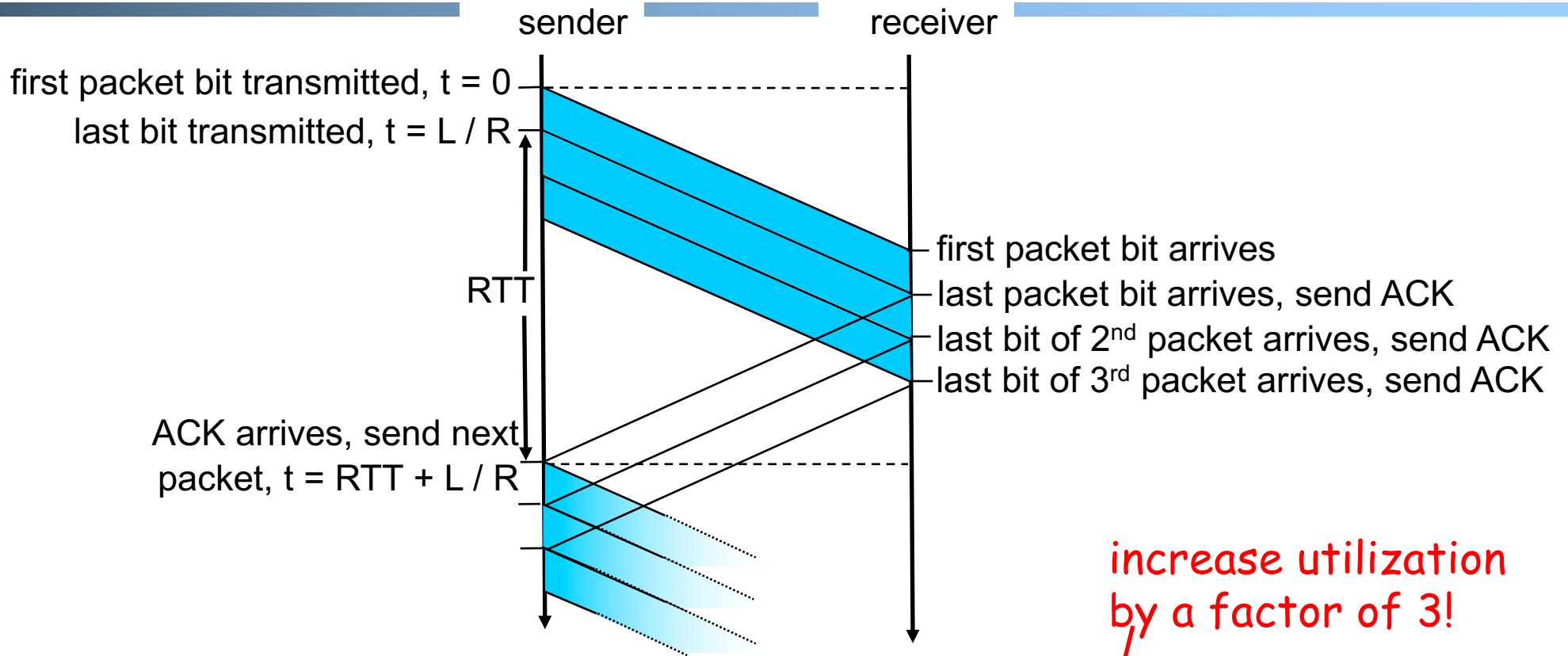- o buffering at sender and/or receiver



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

# 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
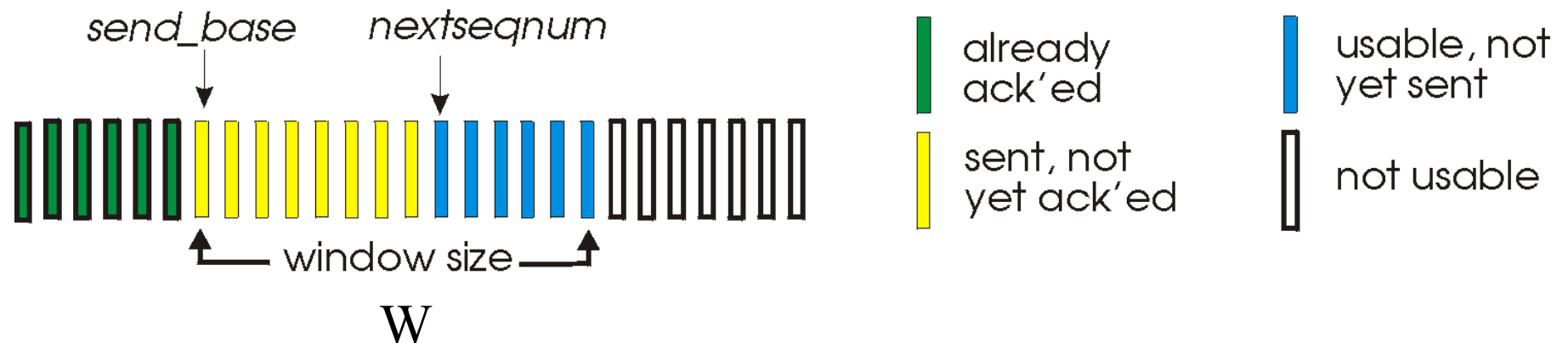packet, t = RTT + L / R

increase utilization
by a factor of 3!

$$U_{sender} = \frac{3 * L / R}{RTT + L / R} = \frac{.024}{30.008} = 0.0008$$

Question: a rule-of-thumb window size?

# Realizing Sliding Window: Go-Back-n

## Sender:

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



send_base  nextseqnum

already ack'ed

sent, not yet ack'ed

usable, not yet sent

not usable

window size

W

- 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

# GBN: Sender FSM

rdt_send(data)
_____
```
if (nextseqnum < base+W) {
    sndpkt[nextseqnum] = make_pkt(nextseqnum,data,chksum)
    udt_send(sndpkt[nextseqnum])
    if (base == nextseqnum) start_timer
    nextseqnum++
} else
    block sender
```
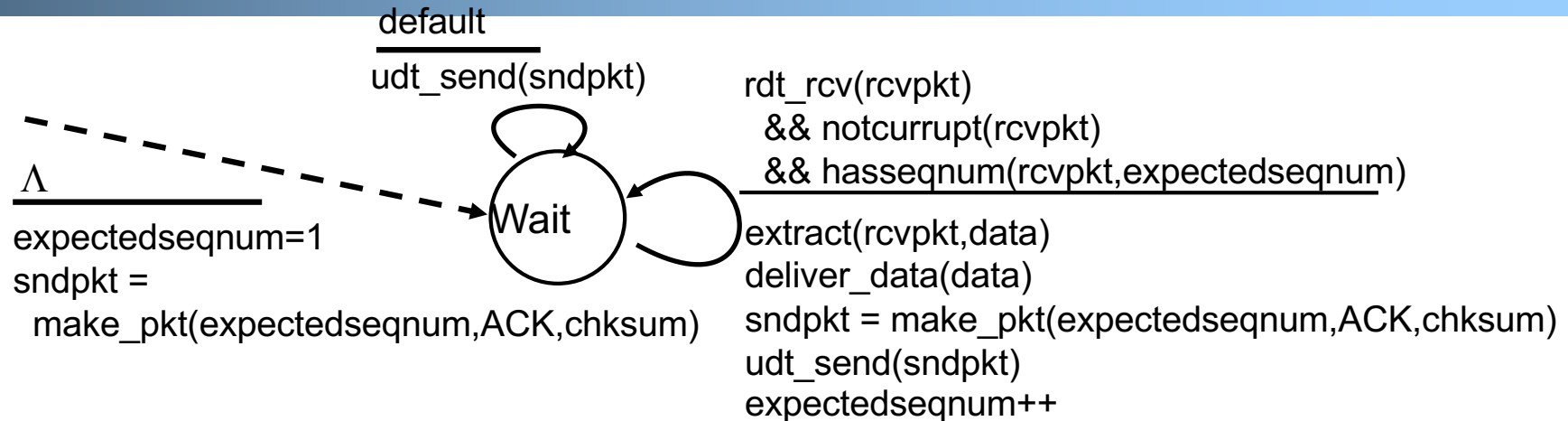
$\Lambda$
_____
base=1
nextseqnum=1

Wait

timeout
_____
```
start_timer
udt_send(sndpkt[base])
udt_send(sndpkt[base+1])
…
udt_send(sndpkt[nextseqnum-1])
```

rdt_rcv(rcvpkt)
   && corrupt(rcvpkt)
_____

rdt_rcv(rcvpkt) &&
   notcorrupt(rcvpkt)
_____
```
if (new packets ACKed) {
    advance base;
    if (more packets waiting)
        send more packets
}
if (base == nextseqnum)
  stop_timer
else
  start_timer for the packet at new base
```

*send_base*        *nextseqnum*

window size
N

20

# GBN: Receiver FSM



**default**
udt_send(sndpkt)

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

**Wait**

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

extract(rcvpkt,data)
deliver_data(data)
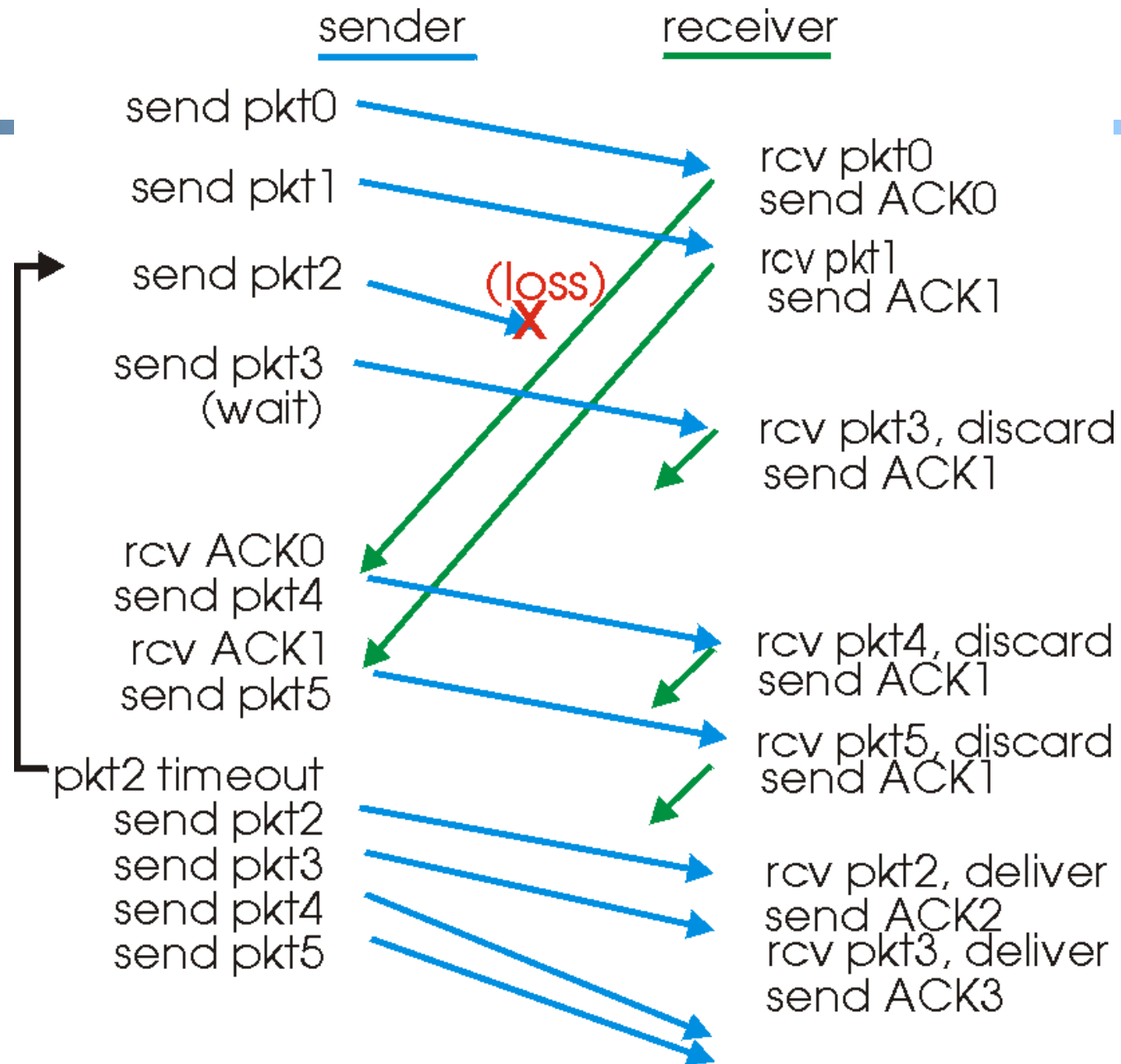sndpkt = make_pkt(expectedseqnum,ACK,chksum)
udt_send(sndpkt)
expectedseqnum++

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

# GBN in Action

window
size = 4



sender                                    receiver

send pkt0
send pkt1
                                          rcv pkt0
                                          send ACK0
send pkt2    (loss)                       rcv pkt1
                                          send ACK1
send pkt3
(wait)
                                          rcv pkt3, discard
                                          send ACK1

rcv ACK0
send pkt4
rcv ACK1                                  rcv pkt4, discard
send pkt5                                 send ACK1

                                          rcv pkt5, discard
pkt2 timeout                              send ACK1
send pkt2
send pkt3                                 rcv pkt2, deliver
send pkt4                                 send ACK2
send pkt5                                 rcv pkt3, deliver
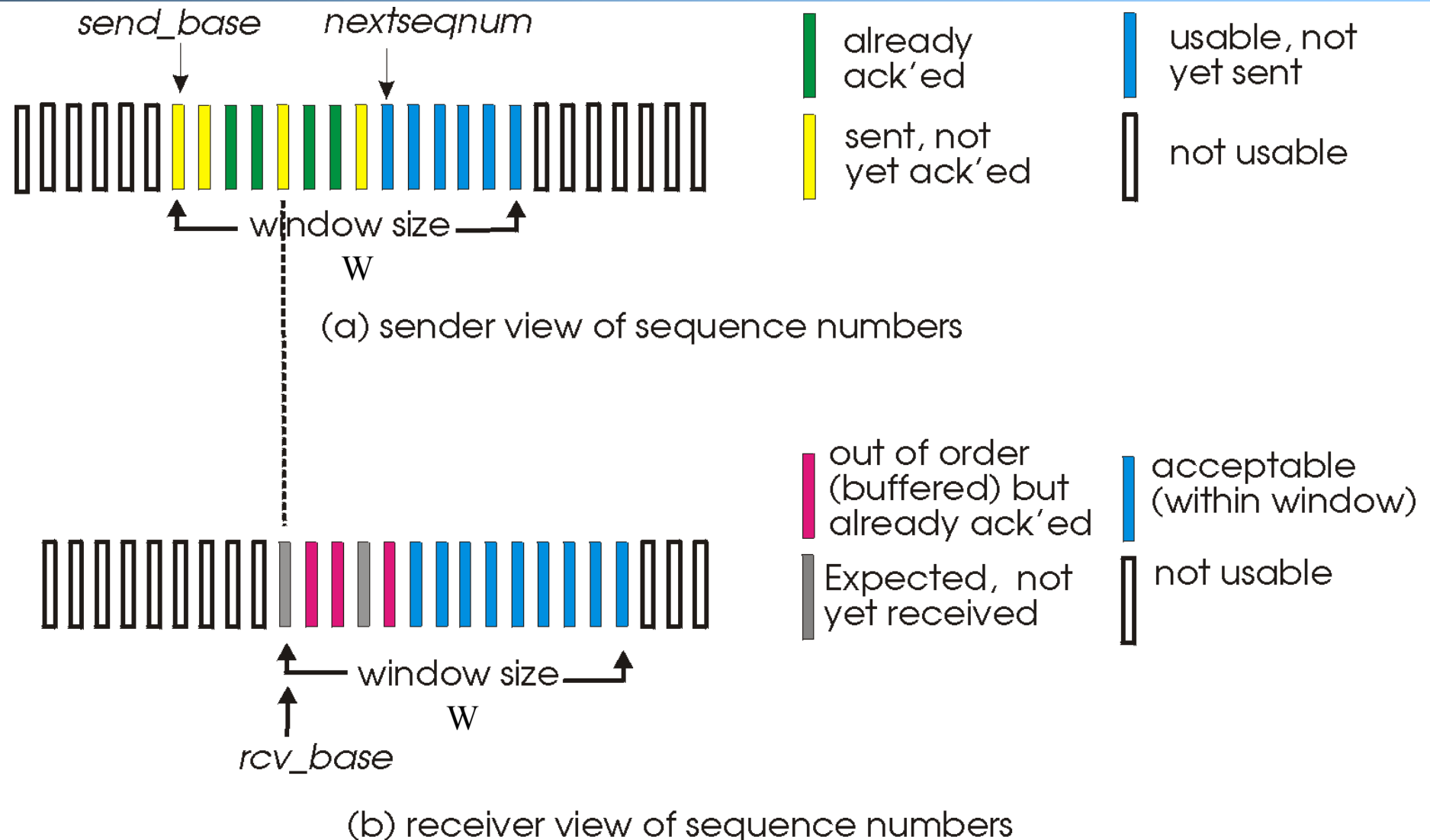                                          send ACK3

# Analysis: Efficiency of Go-Back-n

❑ 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

- **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



send_base    nextseqnum

already ack'ed

sent, not yet ack'ed

usable, not yet sent

not usable

window size W

(a) sender view of sequence numbers

out of order (buffered) but already ack'ed

Expected, not yet received

acceptable (within window)

not usable

window size W

rcv_base

(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:**

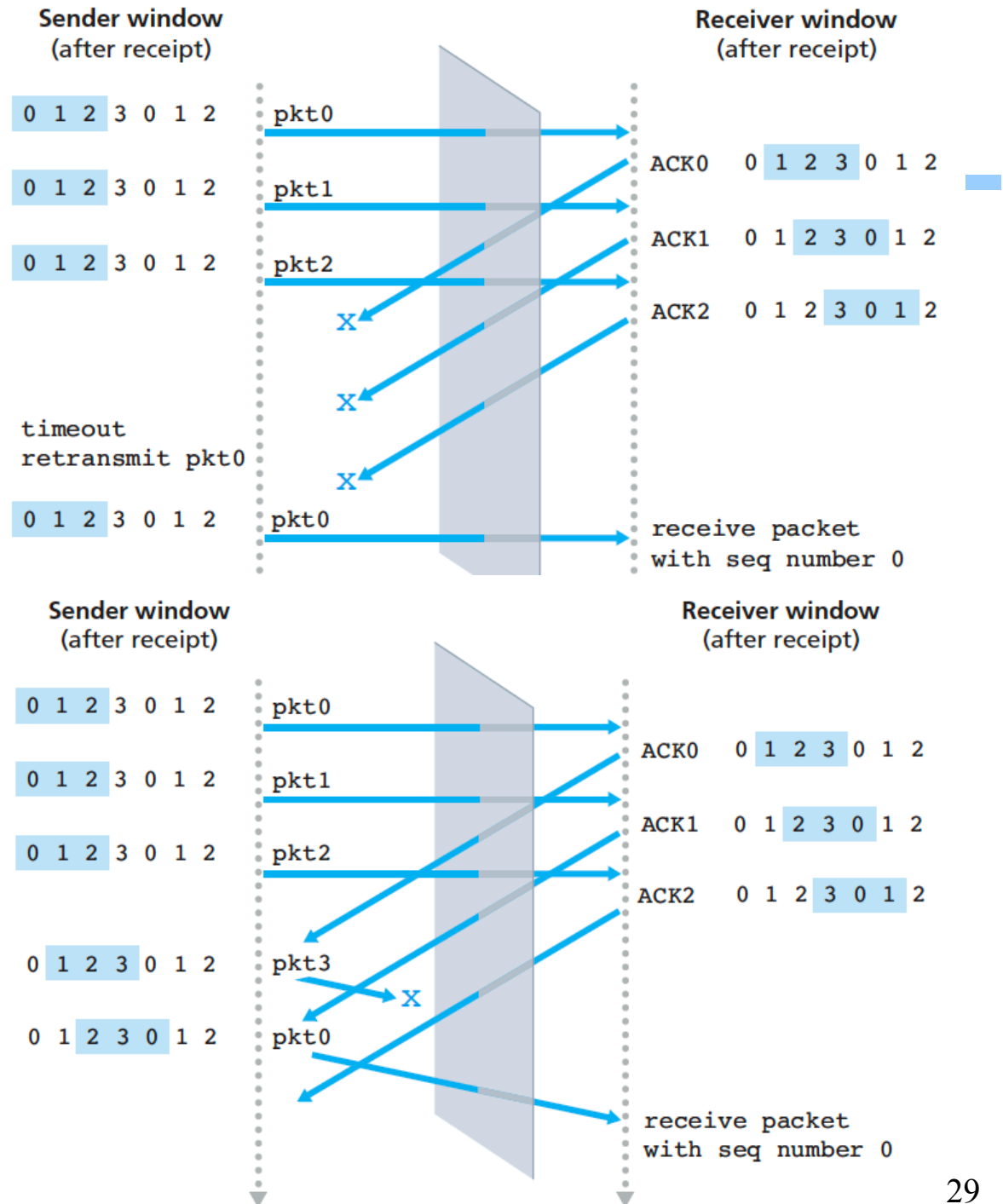- ignore

# 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

27

# 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?
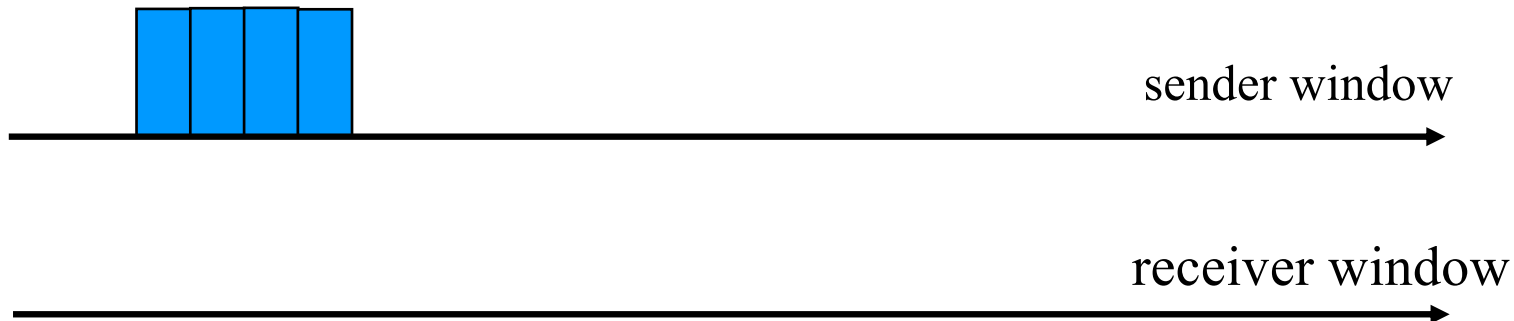
# Selective Repeat: Seq# Ambiguity

Example:

- seq #' s: 0, 1, 2, 3
- window size=3
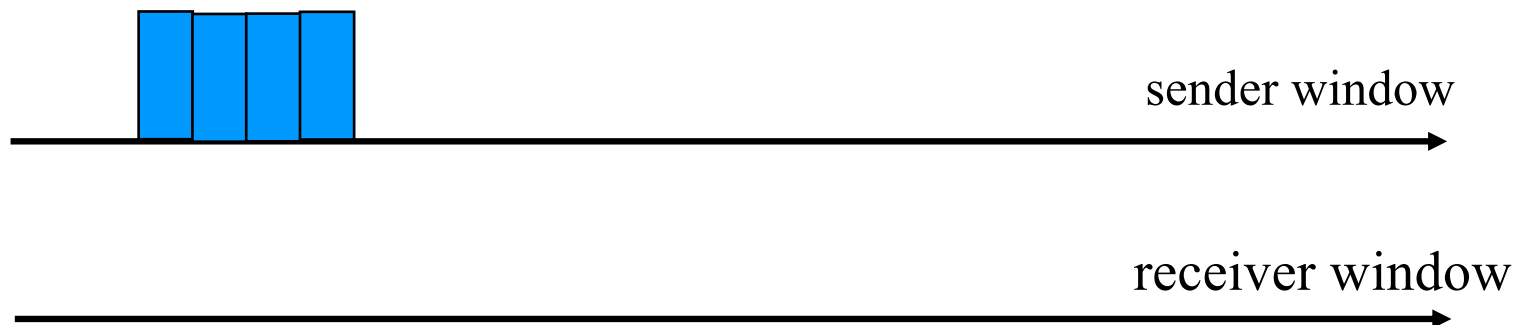
- Error: incorrectly passes duplicate data as new.
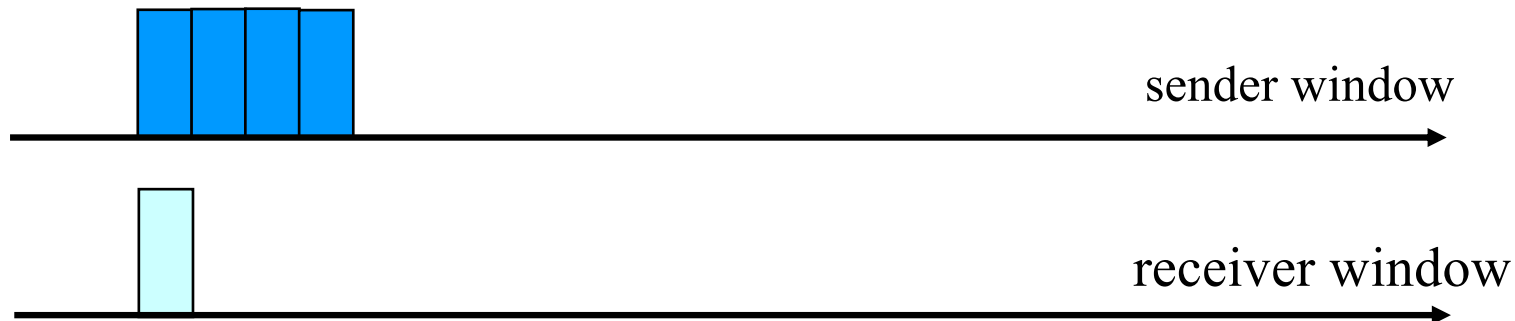
# State Invariant: Window Location

❑ Go-back-n (GBN)

sender window

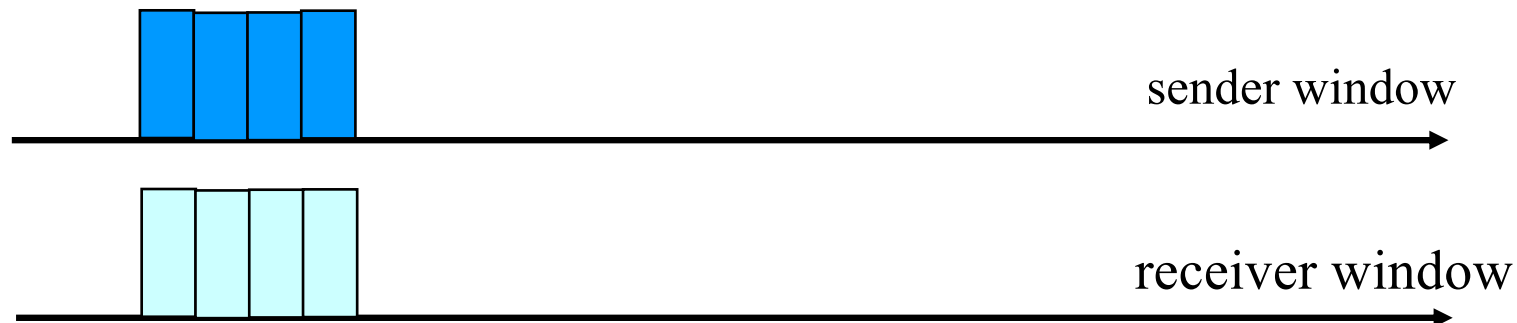receiver window

❑ Selective repeat (SR)

sender window

receiver window

# Window Location

- Go-back-n (GBN)

sender window

receiver window

- Selective repeat (SR)

sender window

receiver window

# 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

# 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