# Module M6

## CPSC 317

## November 2, 2020

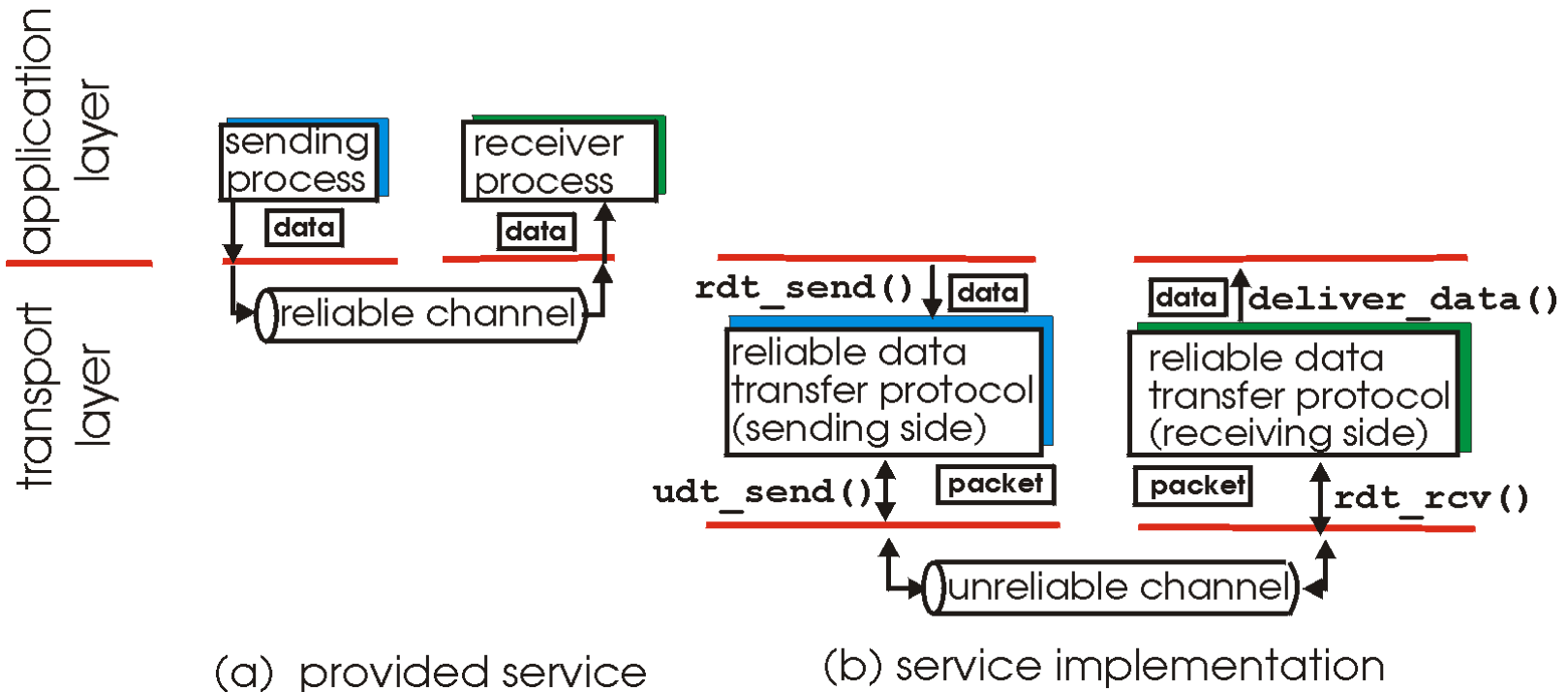(updated November 15, 2020 – slide 74 fix, reorganized TCP sliding window)

# RELIABLE DATA TRANSFER

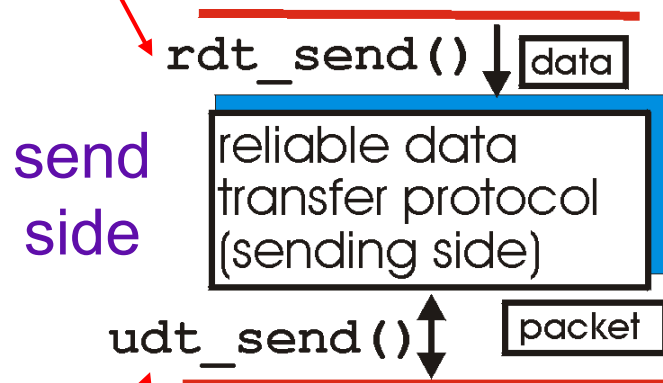# Principles of Reliable data transfer

❑ important in app., transport, link layers
❑ top-10 list of important networking topics!



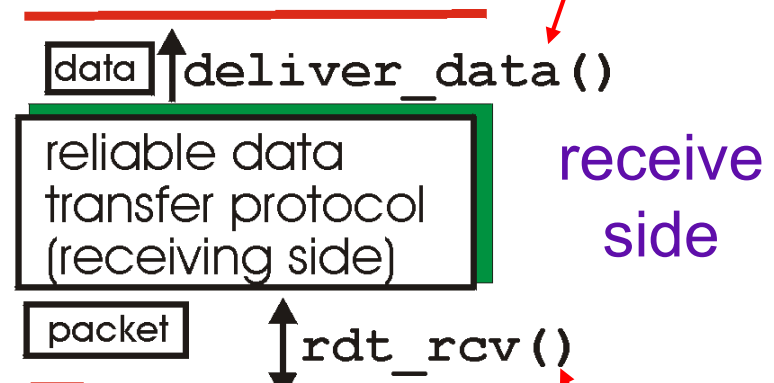(a) provided service       (b) service implementation

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

# Reliable data transfer: getting started

**rdt_send():** called from above, (e.g., by app.). Passed data to deliver to receiver upper layer

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

**send side**

**receive side**

rdt_send()  data

reliable data transfer protocol (sending 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 when packet arrives on rcv-side of channel

# The plan

❑ Reliable channel

❑ Channel that can corrupt messages

❑ Channel that can corrupt and lose messages

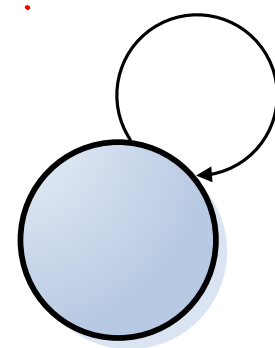❑ What if we can re-order messages???

# Programming State Machines

❑ What does the software look like?

```
Switch( event ):
    event:
          action()
    event:
          action()
End_switch
```
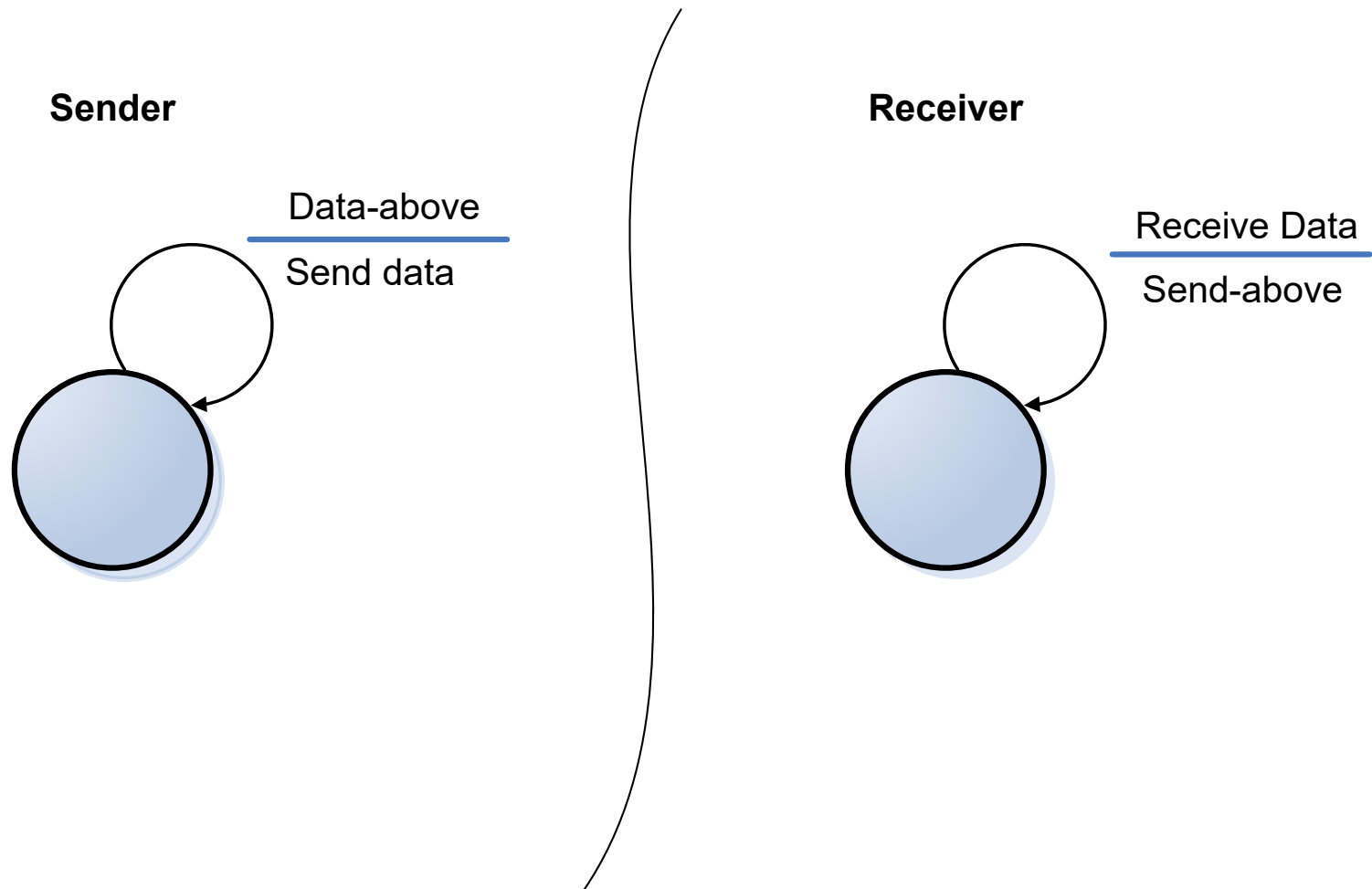
THE UNIVERSITY OF BRITISH COLUMBIA
Modified from Kurose-Ross

# State Machines: Events and Actions

❑ Events:




❑ Actions:

# Reliable Channel
# Communicating State Machines

**Sender**

Data-above
_____
Send data

**Receiver**

Receive Data
_____
Send-above

# Unreliable -- Bit Errors

❑ Messages contents may be garbled.

❑ What do we do?

# Scenario (trace)

# Solution rdt 2.0

**SENDER:**

❑ Events
- App message ready
- NAK recv'ed
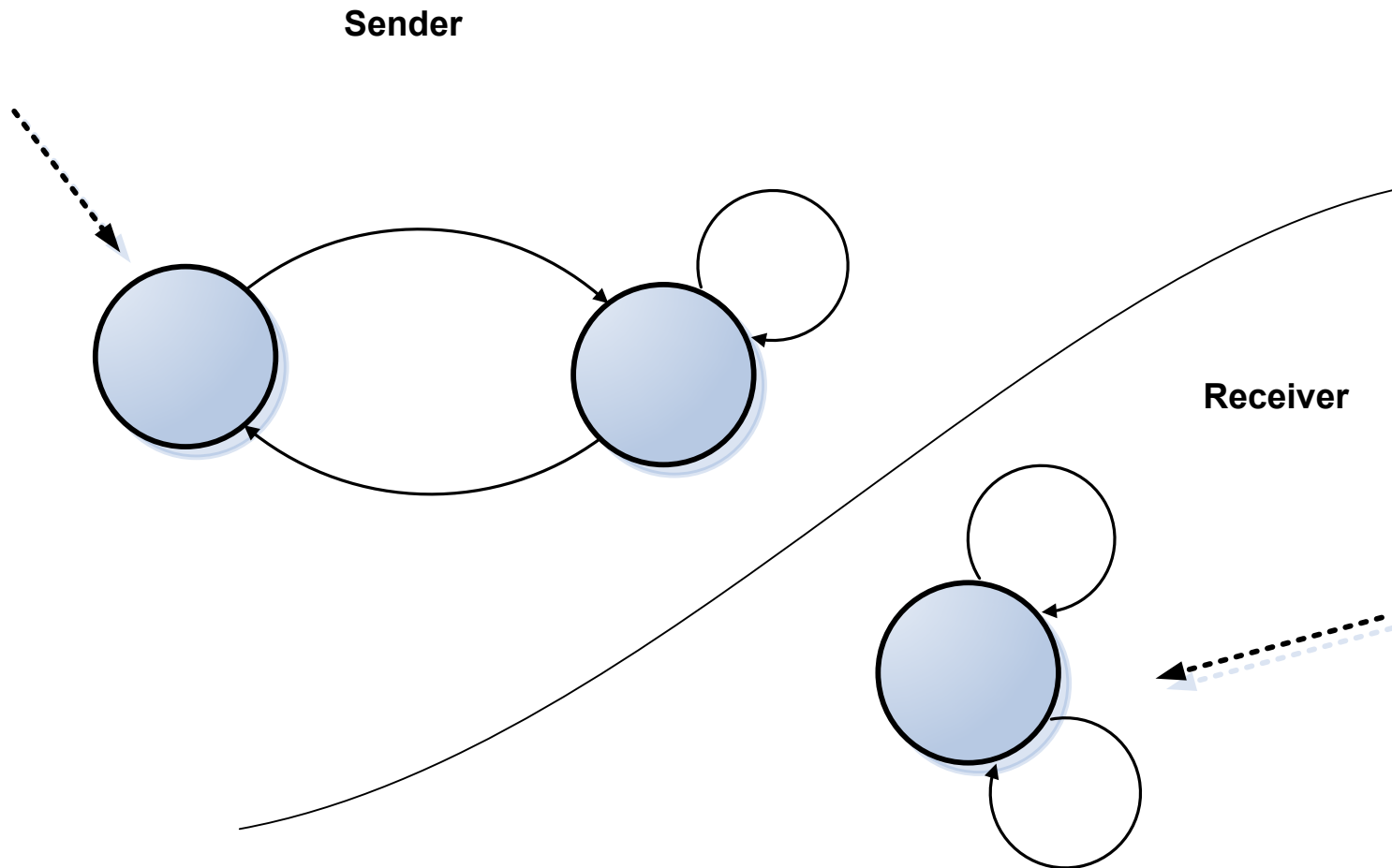- ACK recv'ed

❑ Actions
- Recv from app
- Send to link

**RECEIVER:**

❑ Events
- Link packet ready
- Corrupt packet

❑ Actions
- Send message to app
- Discard, send NAK
- Send ACK

THE UNIVERSITY OF BRITISH COLUMBIA
Modified from Kurose-Ross

# rdt 2.0 -- State Diagrams

**Sender**

**Receiver**

# Scenario (corrupt ptk)

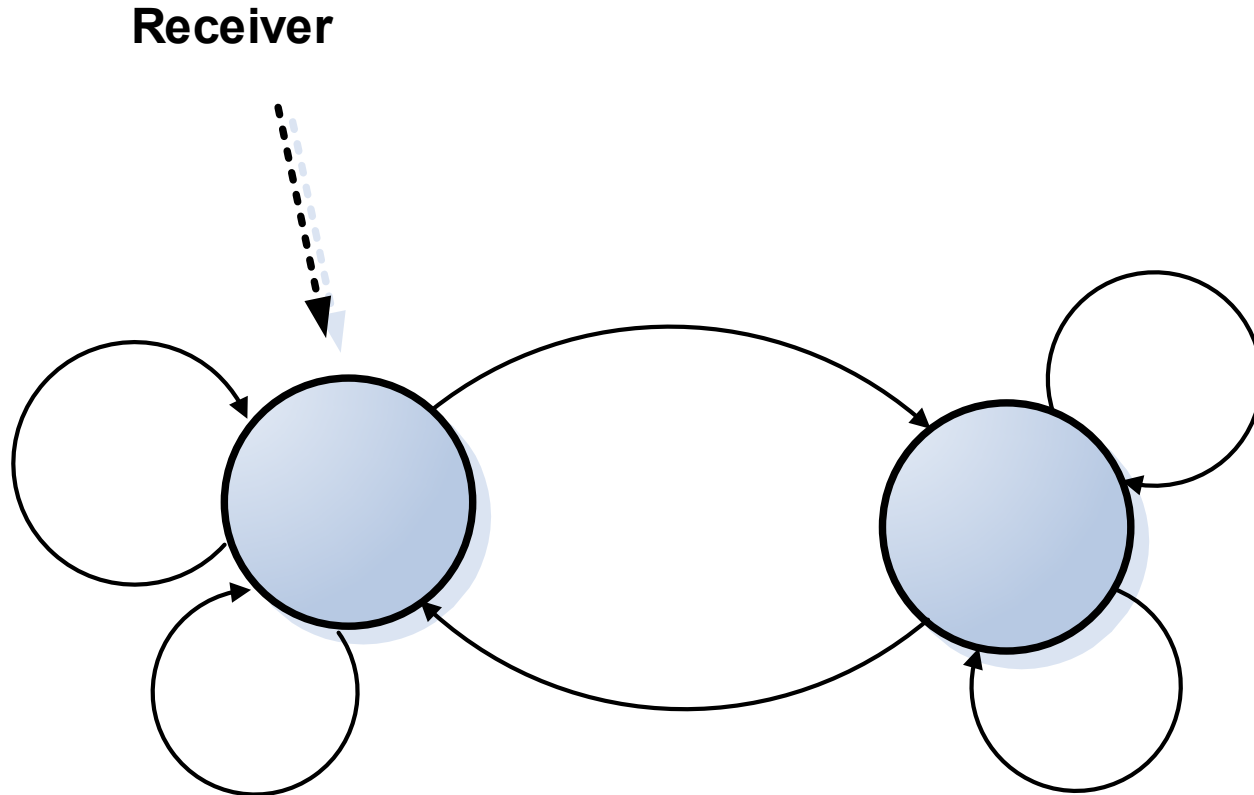THE UNIVERSITY OF BRITISH COLUMBIA
Modified from Kurose-Ross

# Scenario (corrupt ack/nack)

# FIXING ACK/NACK PROBLEM

# Receiver rdt2.1

**Receiver**

# rdt2.1: receiver, handles garbled ACK/NAKs

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
&& has_seq0(rcvpkt)
_____
extract(rcvpkt,data)
deliver_data(data)
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) && corrupt(rcvpkt)
_____
sndpkt = make_pkt(NAK, chksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) &&
not corrupt(rcvpkt) &&
has_seq1(rcvpkt)
_____
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)
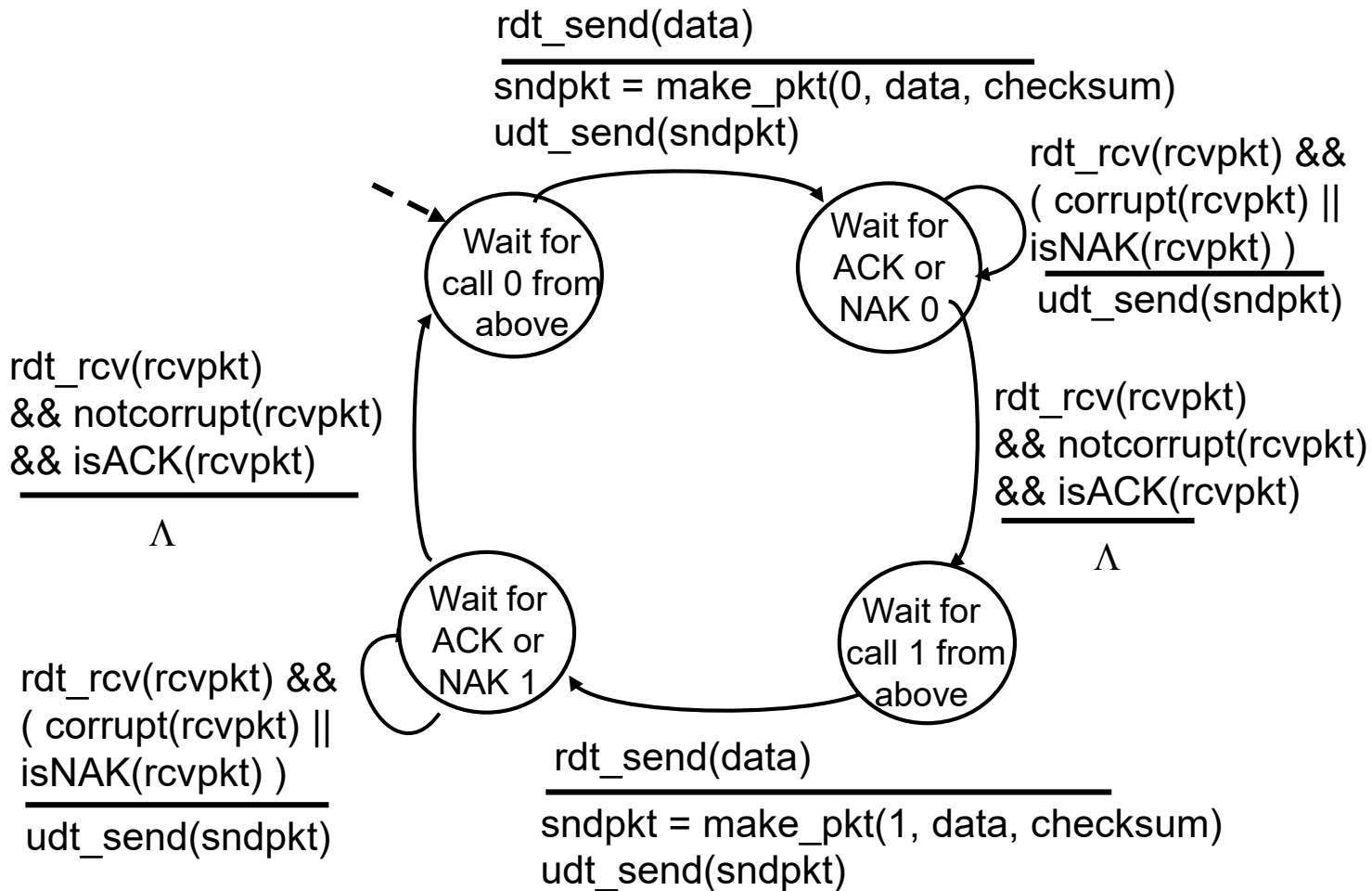
**Wait for 0 from below**

**Wait for 1 from below**

rdt_rcv(rcvpkt) && corrupt(rcvpkt)
_____
sndpkt = make_pkt(NAK, chksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) &&
not corrupt(rcvpkt) &&
has_seq0(rcvpkt)
_____
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
&& has_seq1(rcvpkt)
_____
extract(rcvpkt,data)
deliver_data(data)
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

THE UNIVERSITY OF BRITISH COLUMBIA
Modified from Kurose-Ross

# Sender rdt2.1

# rdt2.1: sender, handles garbled ACK/NAKs

rdt_send(data)
_____

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

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

Wait for
call 0 from
above

Wait for
ACK or
NAK 0

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

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

Wait for
ACK or
NAK 1

Wait for
call 1 from
above

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

rdt_send(data)
_____

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

# Solution rdt2.1

**Sender**

**Receiver**

# Scenario (corrupt nackless)

# Nakless Sender rdt2.2

# Nakless Receiver rdt2.2

THE UNIVERSITY OF BRITISH COLUMBIA

# rdt2.2: sender, receiver fragments: sequence numbers

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

**Wait for call 0 from above**

**Wait for ACK 0**

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

**sender FSM fragment**

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& **isACK(rcvpkt,0)**
_____
$\Lambda$

rdt_rcv(rcvpkt) &&
(corrupt(rcvpkt) ||
**has_seq1(rcvpkt))**
_____
sndptk=make_pkt(ACK,
1,checksum)
**udt_send(sndpkt)**

**Wait for 0 from below**

**receiver FSM fragment**

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
&& has_seq1(rcvpkt)
_____
extract(rcvpkt,data)
deliver_data(data)
**sndpkt = make_pkt(ACK1, chksum)**
udt_send(sndpkt)

# LOSS

# The plan

❏ Reliable channel

❏ Channel that can corrupt messages

❏ <span style="color:red">Channel that can corrupt and lose messages</span>

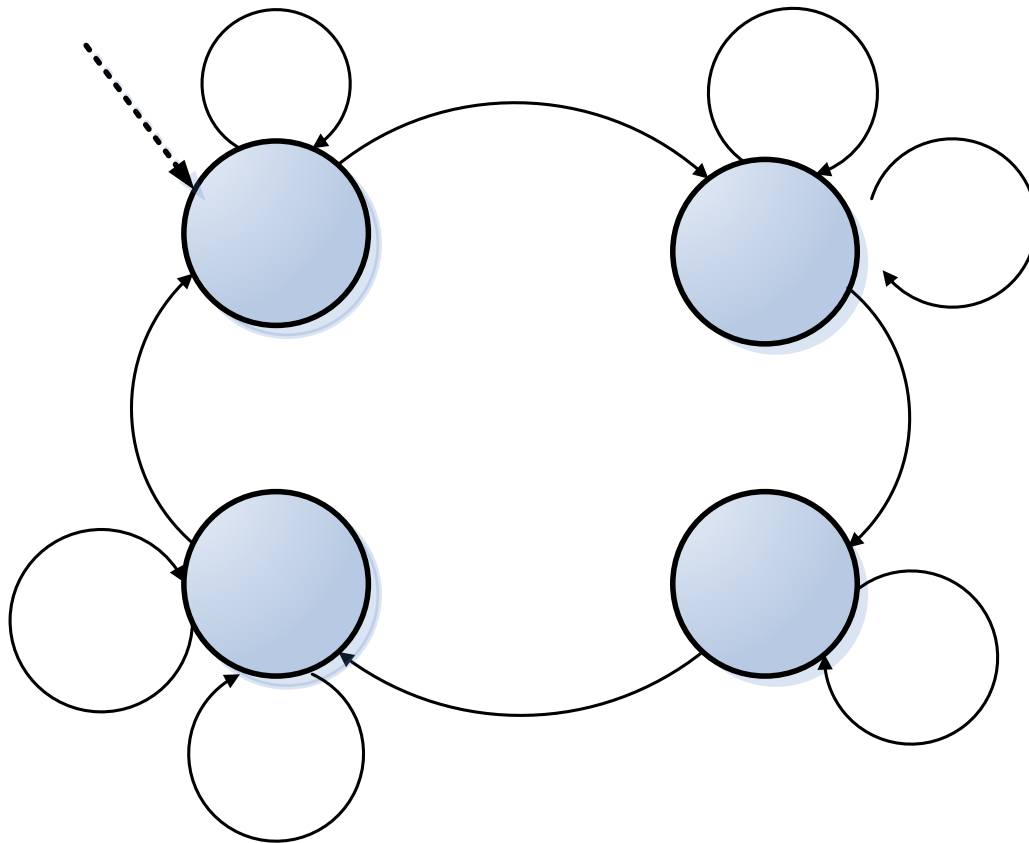❏ What if we can re-order messages???
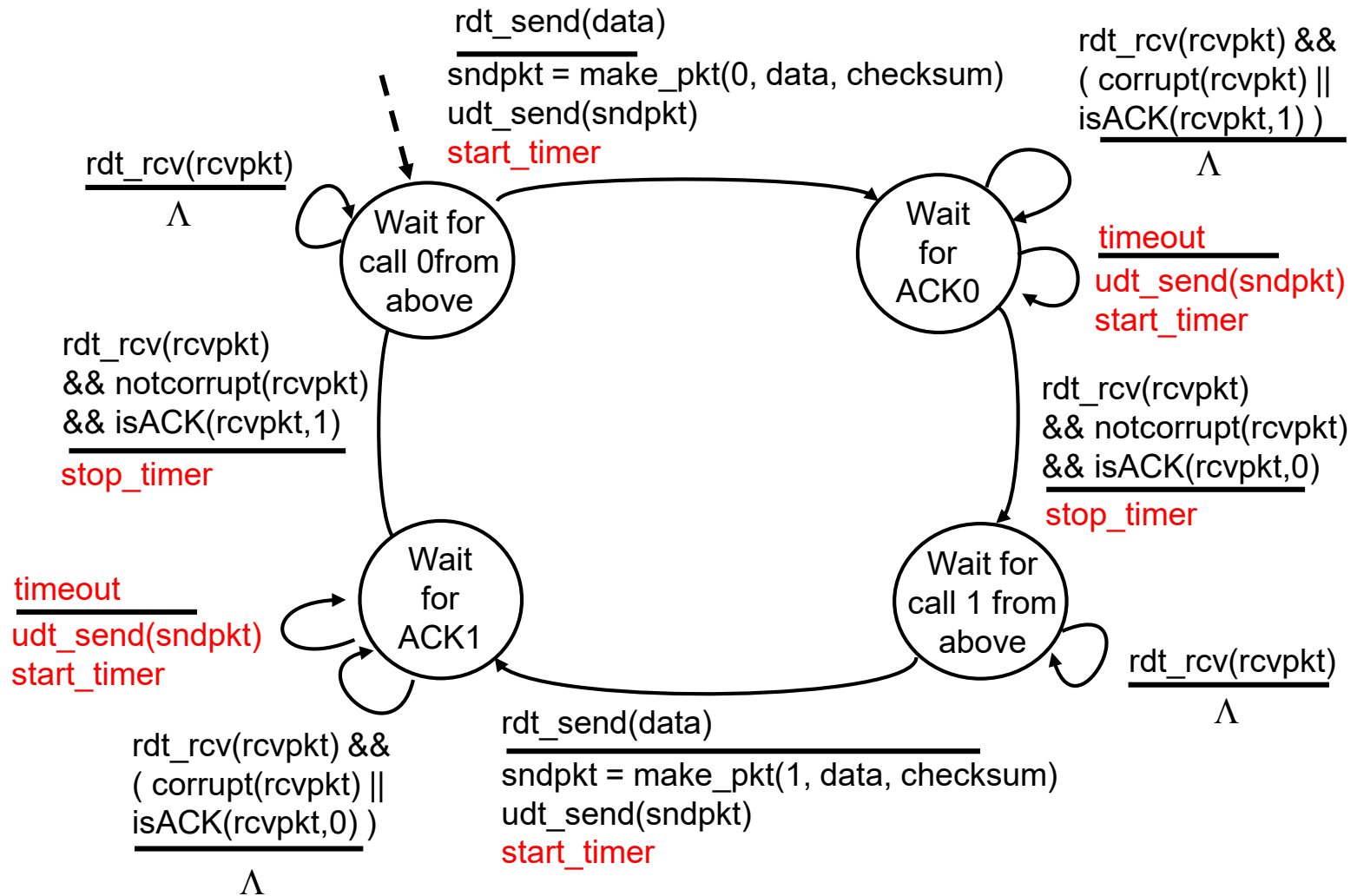
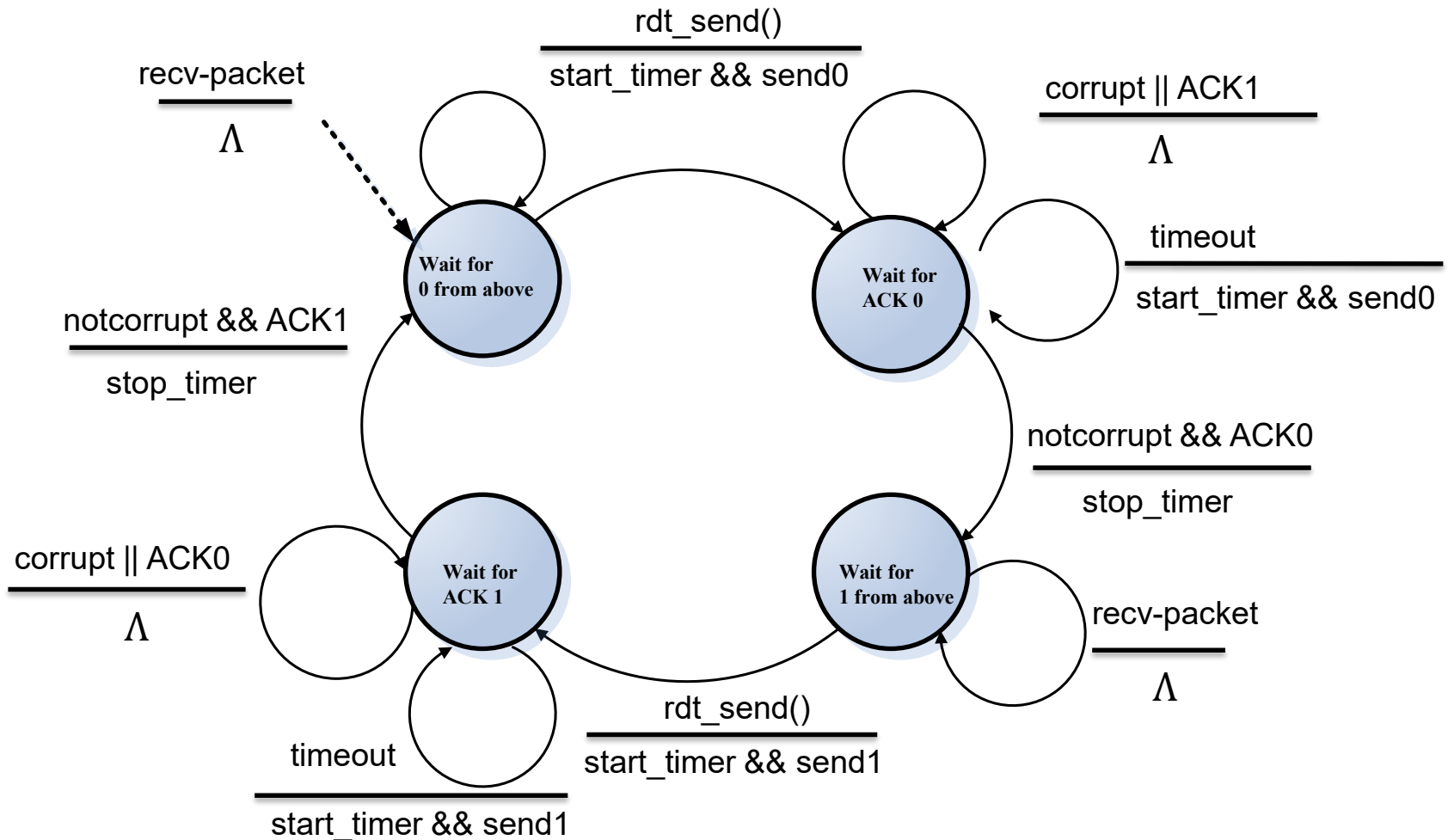# Need

What to do about loss?

How to detect it?

THE UNIVERSITY OF BRITISH COLUMBIA
Modified from Kurose-Ross

# Scenario (loss?)

# Sender rdt3.0

**Sender**

# rdt3.0 sender

# Simplified Sender rdt3.0
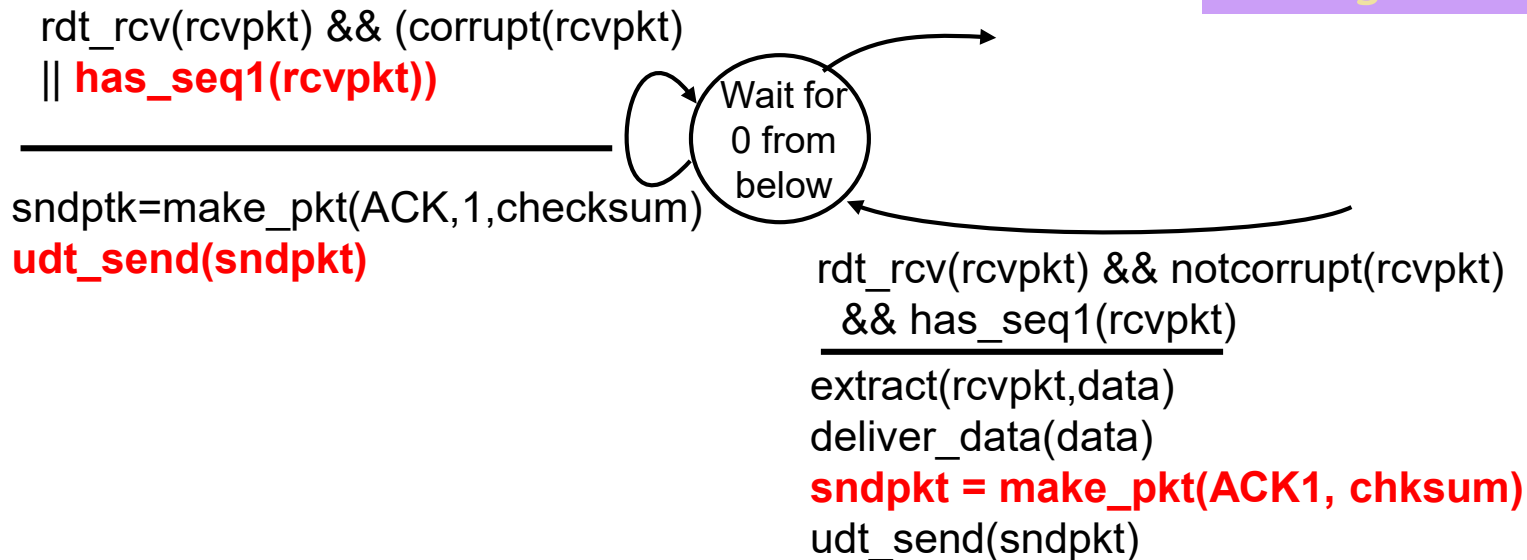
THE UNIVERSITY OF BRITISH COLUMBIA
Modified from Kurose-Ross

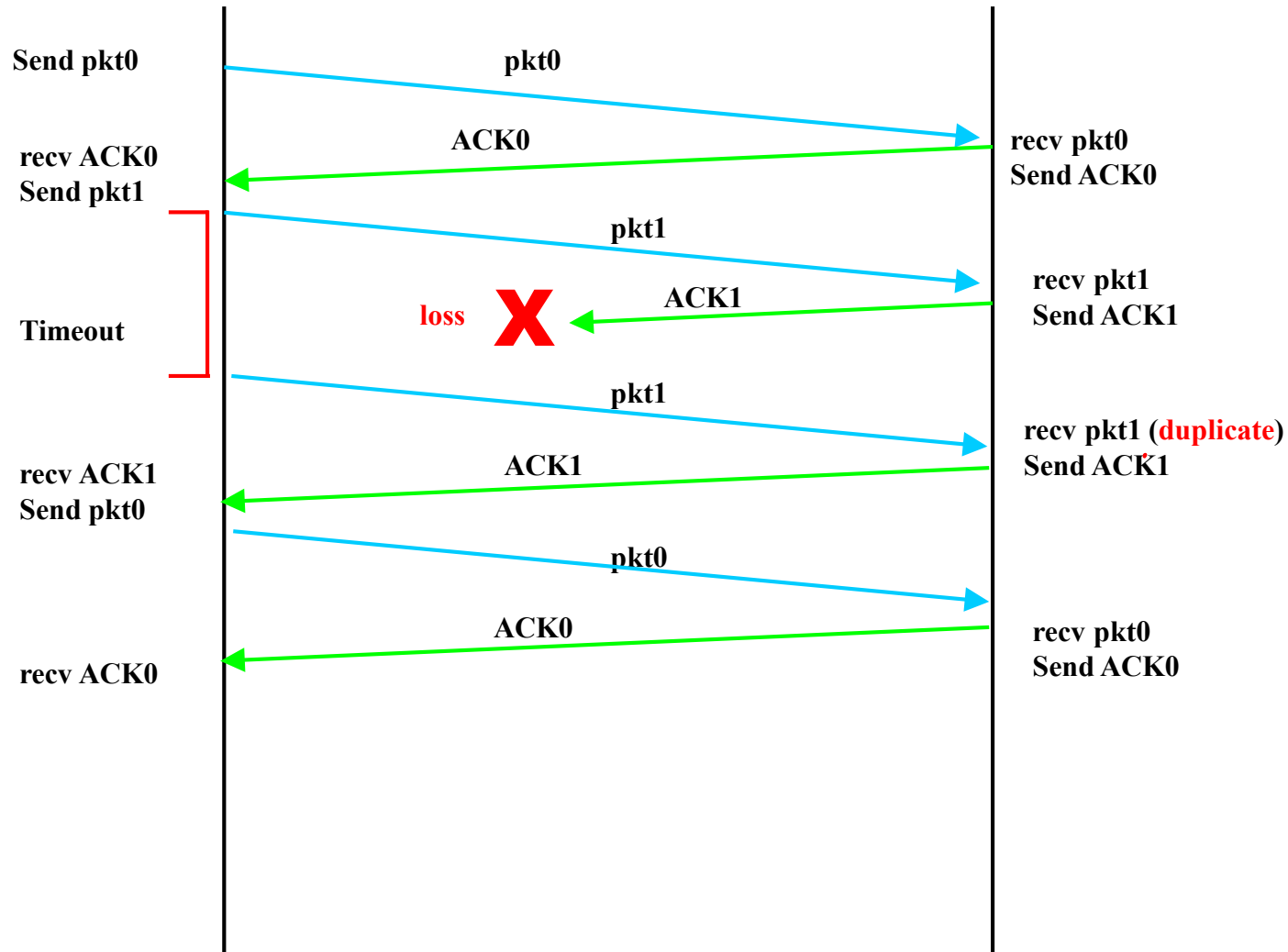# Receiver rdt3.0

# rdt3.0: receiver fragments

receiver FSM fragment

rdt_rcv(rcvpkt) && (corrupt(rcvpkt)
|| **has_seq1(rcvpkt))**
_____
sndptk=make_pkt(ACK,1,checksum)
**udt_send(sndpkt)**

Wait for
0 from
below

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
 && has_seq1(rcvpkt)
_____
extract(rcvpkt,data)
deliver_data(data)
**sndpkt = make_pkt(ACK1, chksum)**
udt_send(sndpkt)

# Normal Operation, no loss

**Send pkt0**  →  pkt0  →  **recv pkt0**
**Send ACK0**

**recv ACK0**  ←  ACK0
**Send pkt1**  →  pkt1  →  **recv pkt1**
**Send ACK1**

**recv ACK1**  ←  ACK1
**Send pkt0**  →  pkt0  →  **recv pkt0**
**Send ACK0**

**recv ACK0**  ←  ACK0

# Lost Packet

**Send pkt0** → pkt0 → **recv pkt0**
**Send ACK0**

← ACK0 ←

**recv ACK0**
**Send pkt1** → pkt1 → **X**

**Timeout**

**loss**

**Send pkt1** → pkt1 → **recv pkt1**
**Send ACK1**

← ACK1 ←

**recv ACK0**
→ pkt0 →

# Duplicate Packet at Receiver

**Send pkt0**

pkt0

**recv ACK0**
**Send pkt1**

ACK0

**recv pkt0**
**Send ACK0**

pkt1

**recv pkt1**
**Send ACK1**

**Timeout**

**loss** X ACK1

pkt1

**recv pkt1 (duplicate)**
**Send ACK1**

**recv ACK1**
**Send pkt0**

ACK1

pkt0

**recv ACK0**

ACK0

**recv pkt0**
**Send ACK0**

# Duplicate Acknowlegement(s)



Send pkt0

pkt0

recv pkt0
Send ACK0

ACK0

recv ACK0
Send pkt1

pkt1

recv pkt1
Send ACK1

Timeout

ACK1

pkt1

recv pkt1 (**duplicate**)
Send ACK1

recv ACK1
Send pkt0

pkt0

ACK1

Duplicate ACK1 **do nothing**

recv pkt0
Send ACK0

ACK0

recv ACK0

- ❑ second ACK1 is ignored by sender
- ❑ sender has sent pkt0; now expecting ACK0, ignores all else.

THE UNIVERSITY OF BRITISH COLUMBIA
Modified from Kurose-Ross

# Clarify Time-out situation

❑ Second ACK1 is ignored.

❑ Sender has sent pkt0 so is now expecting a ACK0, ignores everything else.



(d) premature timeout

# rdt3.0 in action



(a) operation with no loss

(b) lost packet

# Reliable Data Transfer Summary

❑ Acknowledgements (Negative ACKS)

❑ Re-transmissions

❑ Checksum (for detecting corrupt packets)

❑ Sequence Numbers

❑ Timer  (needed when there is loss)


❑ No solution for OUT-OF-ORDER

# PERFORMANCE STOP&WAIT

# Performance of rdt3.0

❑ rdt3.0 works, but performance is TERRIBLE
❑ example: 1 Gbps link, 15 millisecond propagation delay, 8000 bit packet:

$$d_{trans} = \frac{L}{R} = \frac{8000 \text{bits}}{10^9 \text{bps}} = 8 \text{ microseconds}$$

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

$$U_{sender} = \frac{L \,/\, R}{RTT + L \,/\, R} = \frac{.008}{30.008} = 0.00026$$

○ 1 pkt every 30 msec -> 0.26 Mbps throughput over 1 Gbps link
○ network protocol limits use of physical resources!

# rdt3.0: stop-and-wait operation



first packet bit transmitted, t = 0

last packet bit transmitted, t = L / R

RTT

first packet bit arrives

last packet bit arrives, send ACK

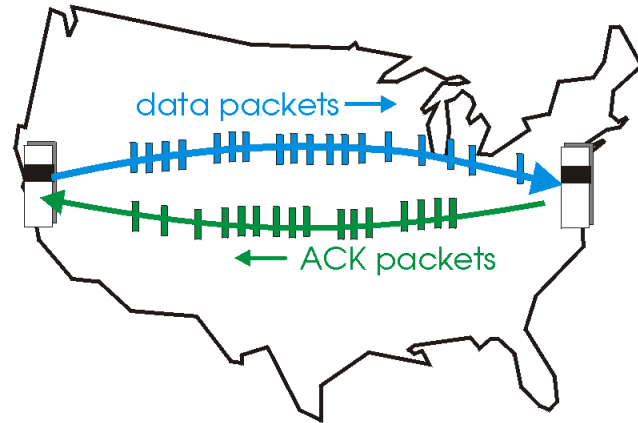ACK arrives, send next packet, t = RTT + L / R

sender

receiver

# Pipelined protocols

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

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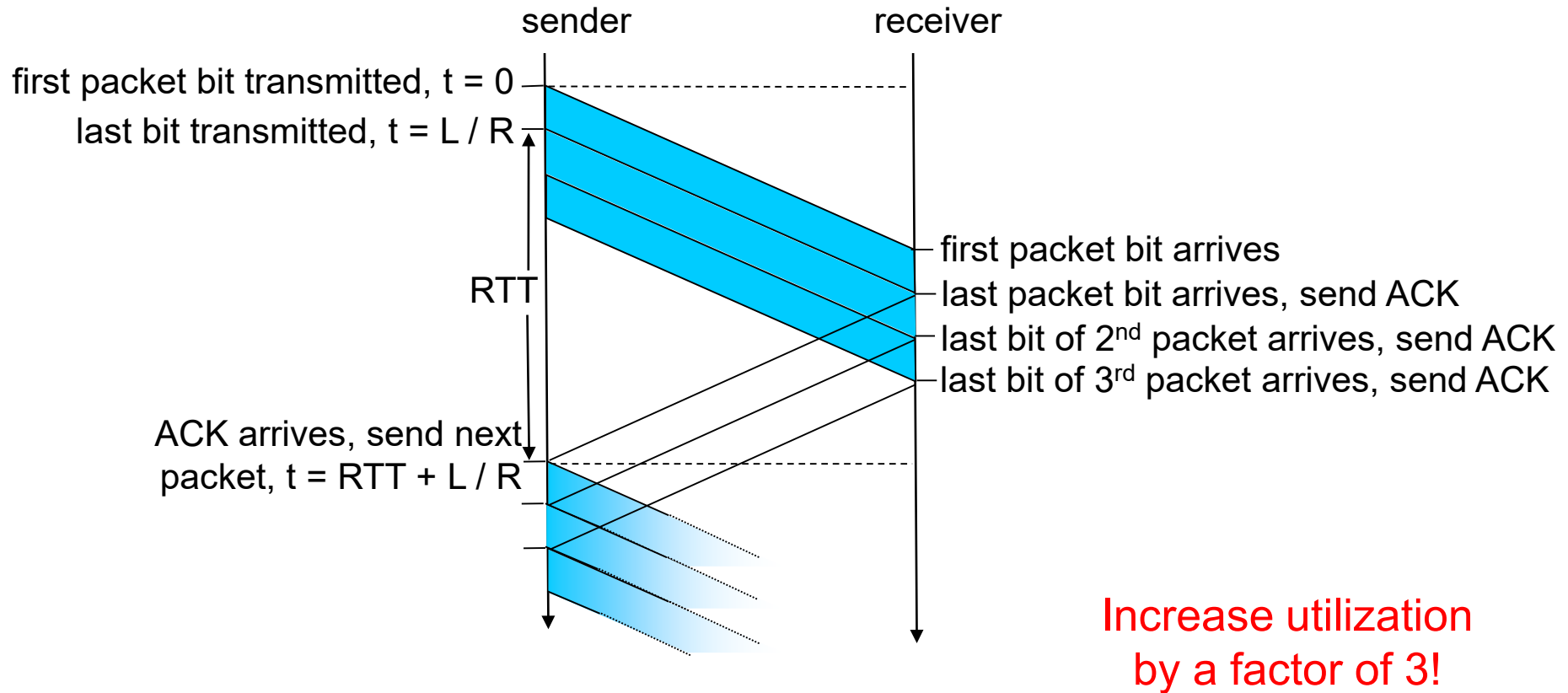- buffering at sender and/or receiver



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

❑ Two generic forms of pipelined protocols: *various TCP ones, 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

Increase utilization
by a factor of 3!

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

# SLIDING WINDOW

# SLIDING WINDOW

https://media.pearsoncmg.com/aw/ecs_kurose_compnetwork_7/cw/content/interactiveanimations/go-back-n-protocol/index.html

https://media.pearsoncmg.com/aw/ecs_kurose_compnetwork_7/cw/content/interactiveanimations/selective-repeat-protocol/index.html

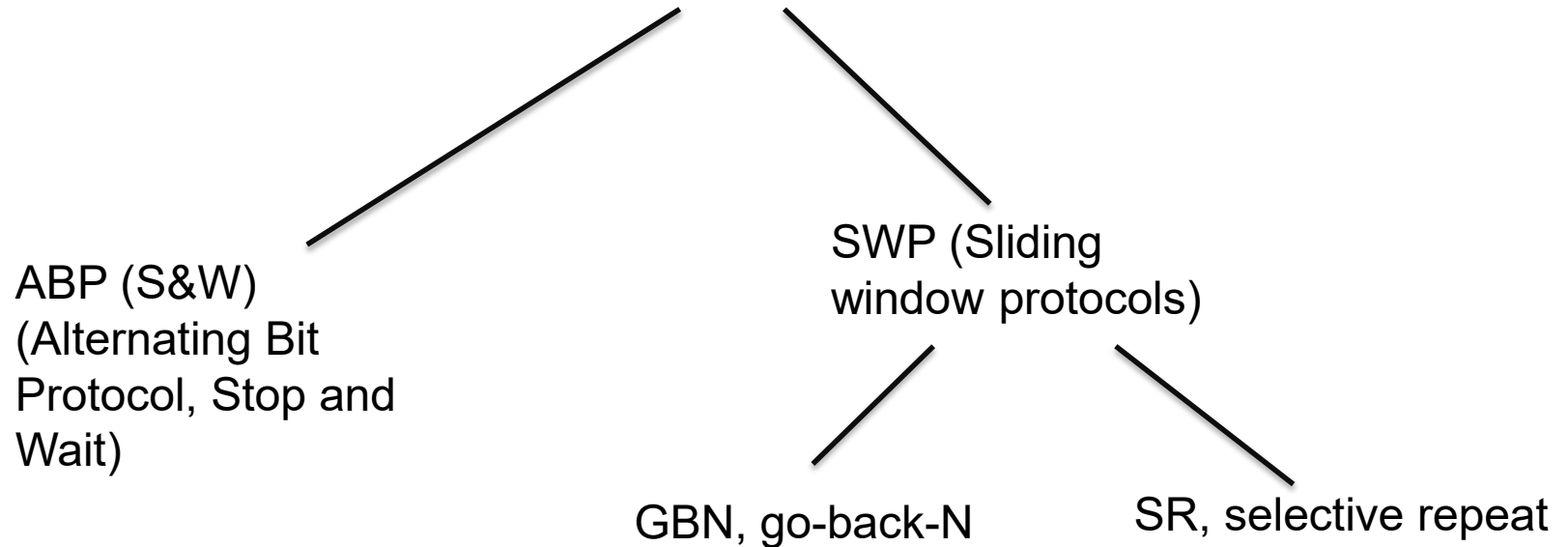http://www.ccs-labs.org/teaching/rn/animations/gbn_sr/

# SLIDING WINDOW

❑ Developed ARQ method called
  ○ Alternating Bit Protocol or
  ○ Stop and Wait


❑ Link utilization (throughput) is low and solution was pipelining (more packets in flight)

# ARQ
# (automatic repeat request)

ABP (S&W)
(Alternating Bit
Protocol, Stop and
Wait)

SWP (Sliding
window protocols)

GBN, go-back-N

SR, selective repeat

# Sliding Window in Action

THE UNIVERSITY OF BRITISH COLUMBIA

# Terminology

**Sender side:**
    **SWS:** send window size
    **LAR**: last ACK received
    **LFS:** last frame sent

**Receiver side:**
    **LFR:** last frame received
    **LAF**: largest acceptable frame

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|

# Stop and Wait

LAR

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |

LFS

**Sender side:**
**LAR**: last ACK received
**LFS:** last frame sent

LFR          LAF

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | |

**Receiver side:**
**LFR:** last frame received
**LAF**: largest acceptable frame

# Sliding Window



**Sender side:**
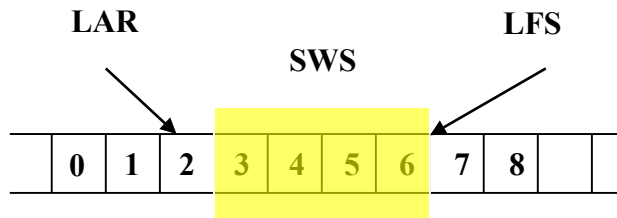  **SWS:** send window size
  **LAR**: last ACK received
  **LFS:** last frame sent

**Receiver side:**
  **RWS:** receive window size
  **LFR:** last frame received
  **LAF**: largest acceptable frame

THE UNIVERSITY OF BRITISH COLUMBIA
Modified from Kurose-Ross

# Sender



**Sender side:**
**SWS:** send window size
**LAR**: last ACK received
**LFS:** last frame sent

Sender:

if more data to send (LFS-LAR < SWS)

then send data, LFS++
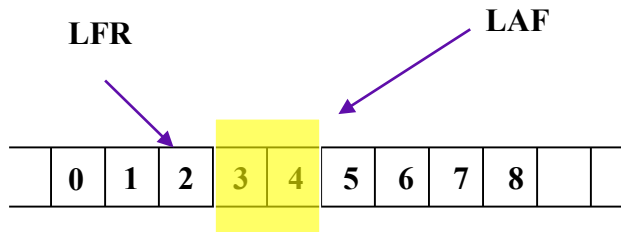
if recv'ed ACK for LAR+1

then LAR++

if timer expires

then send [3 or 3-4-5-6]

Two strategies:
(a) Go-Back-N
(b) Selective Repeat

THE UNIVERSITY OF BRITISH COLUMBIA
Modified from Kurose-Ross

# Receiver

LFR        LAF

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | | |

**Receiver side:**

**LFR:** last frame received

**LAF**: largest acceptable frame

Receiver:

if recv'ed K > LAF

then discard

else

if K == LFR+1 then

store

LFR++, LAF++ (slide window)

else store [or discard]

ACK, largest in-order received frame

Two strategies:
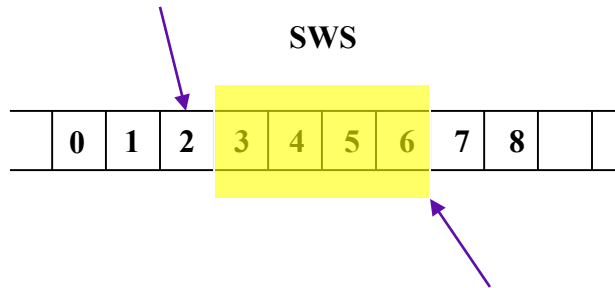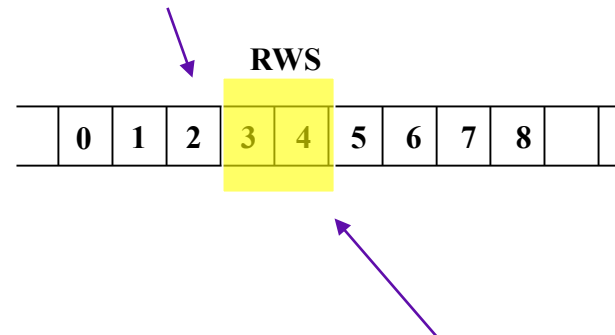(a) Go-Back-N
(b) Selective Repeat

# Sliding Window

**Receiver**:
    if recv'ed K > LAF
        then discard
        else
            if K == LFR+1 then
                store
                LFR++, LAF++ (slide window)
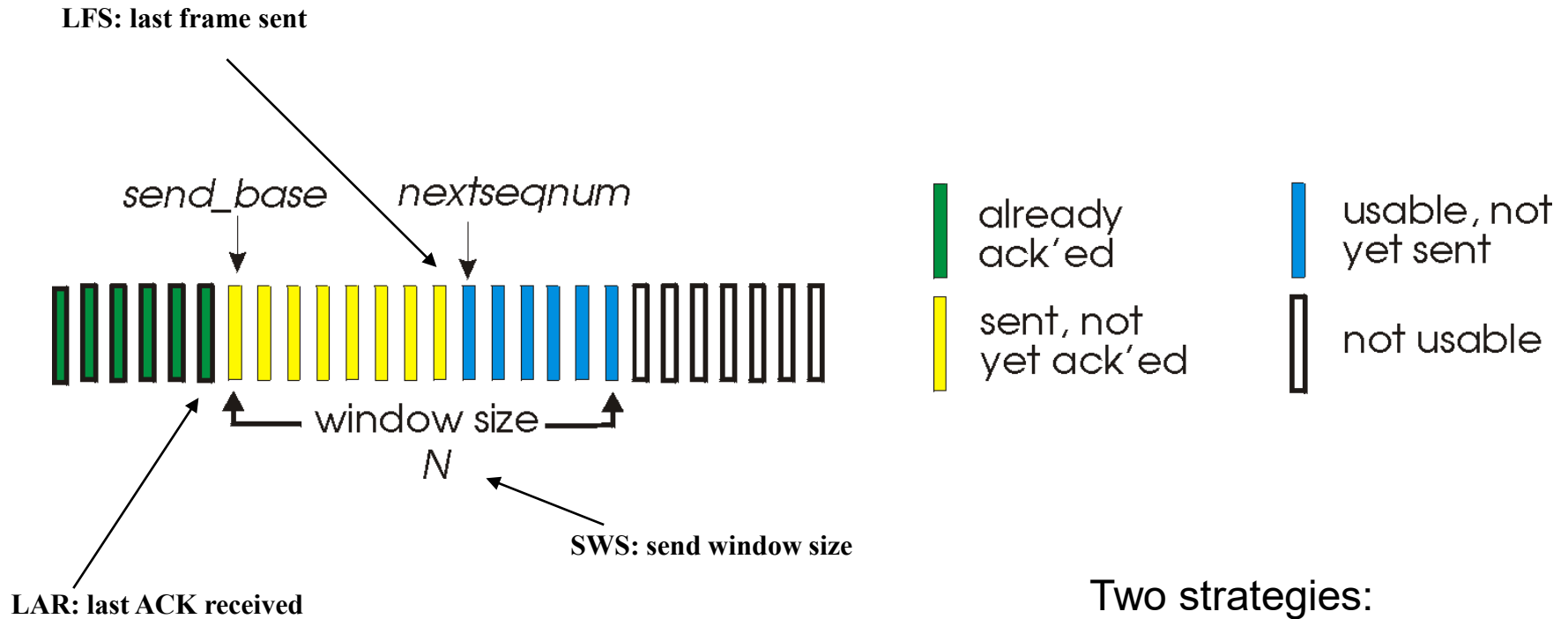            else
                discard
    ACK, largest in-order received frame

**Sender**:
    if more data to send (LFS-LAR < SWS)
        then send data, LFS++
    if recv'ed ACK for LAR+1
        then LAR++
    if timer expires
        then send LAR+1

# Book's Terminology



LFS: last frame sent

send_base

nextseqnum

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

window size N

SWS: send window size
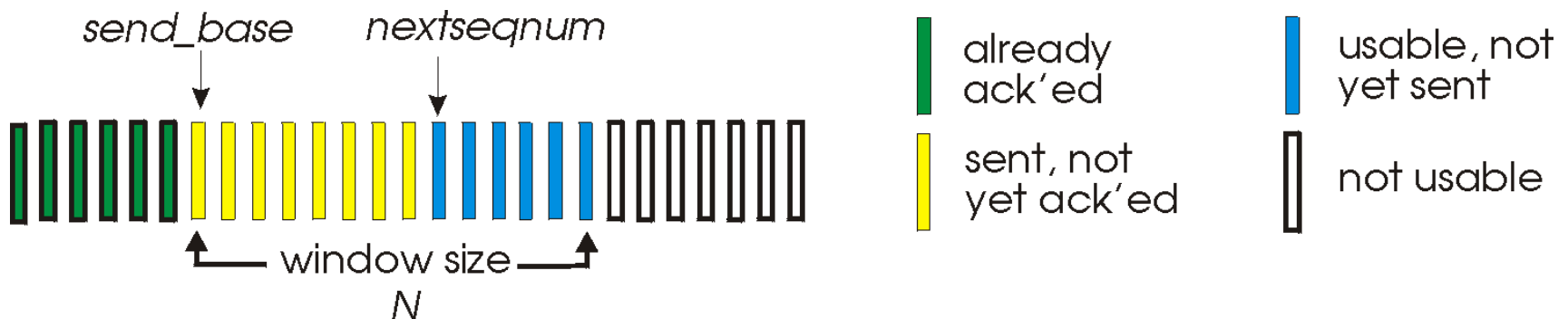
LAR: last ACK received

Two strategies:
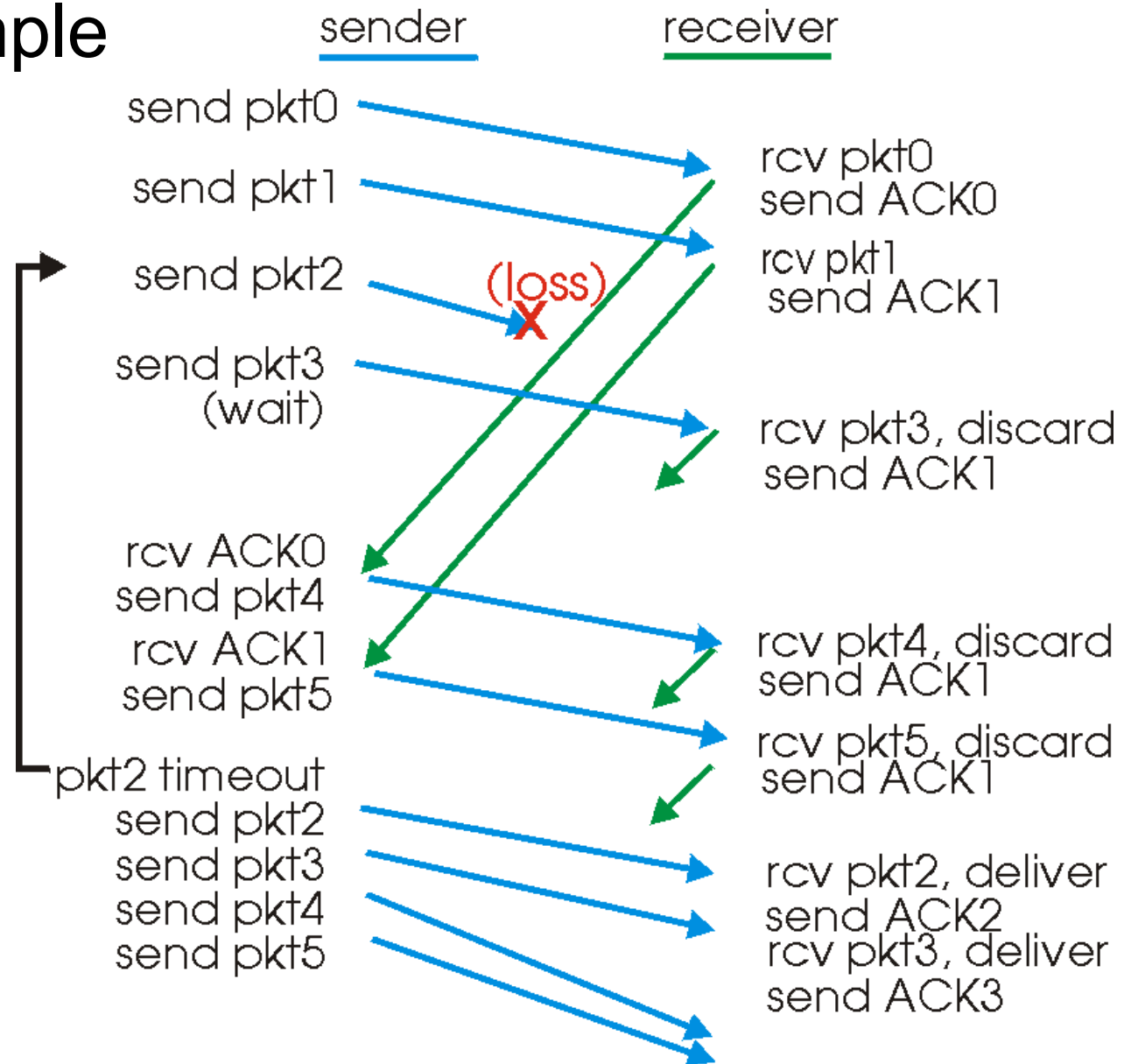(a) Go-Back-N
(b) Selective Repeat

# Go-Back-N

**Sender:**

❑ k-bit 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 only for smallest sequence number sent but not ack'ed

❑ *timeout(n):* re-transmit pkt n and all higher seq # pkts in window

# GBN Example

# Selective Repeat

❑ receiver *individually* acknowledges all correctly received packets

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

❑ sender only resends pkts for which ACK not received

   o sender timer for each unACKed pkt

❑ sender window

   o N consecutive seq #'s

   o again limits seq #s of sent, unACKed pkts

# Selective repeat

**sender**

**data from above :**

- ❑ if next available seq # in window, send pkt

**timeout(n):**

- ❑ resend pkt n, restart timer

**ACK(n)** in [send-window]:

- ❑ mark pkt n as received
- ❑ if n smallest unACKed pkt, advance window base to next unACKed seq #

**receiver**

**pkt n in** [recv-window]

- ❑ send ACK(n)
- ❑ out-of-order: buffer
- ❑ in-order: deliver (also deliver buffered, in-order pkts), advance window to next not-yet-received pkt
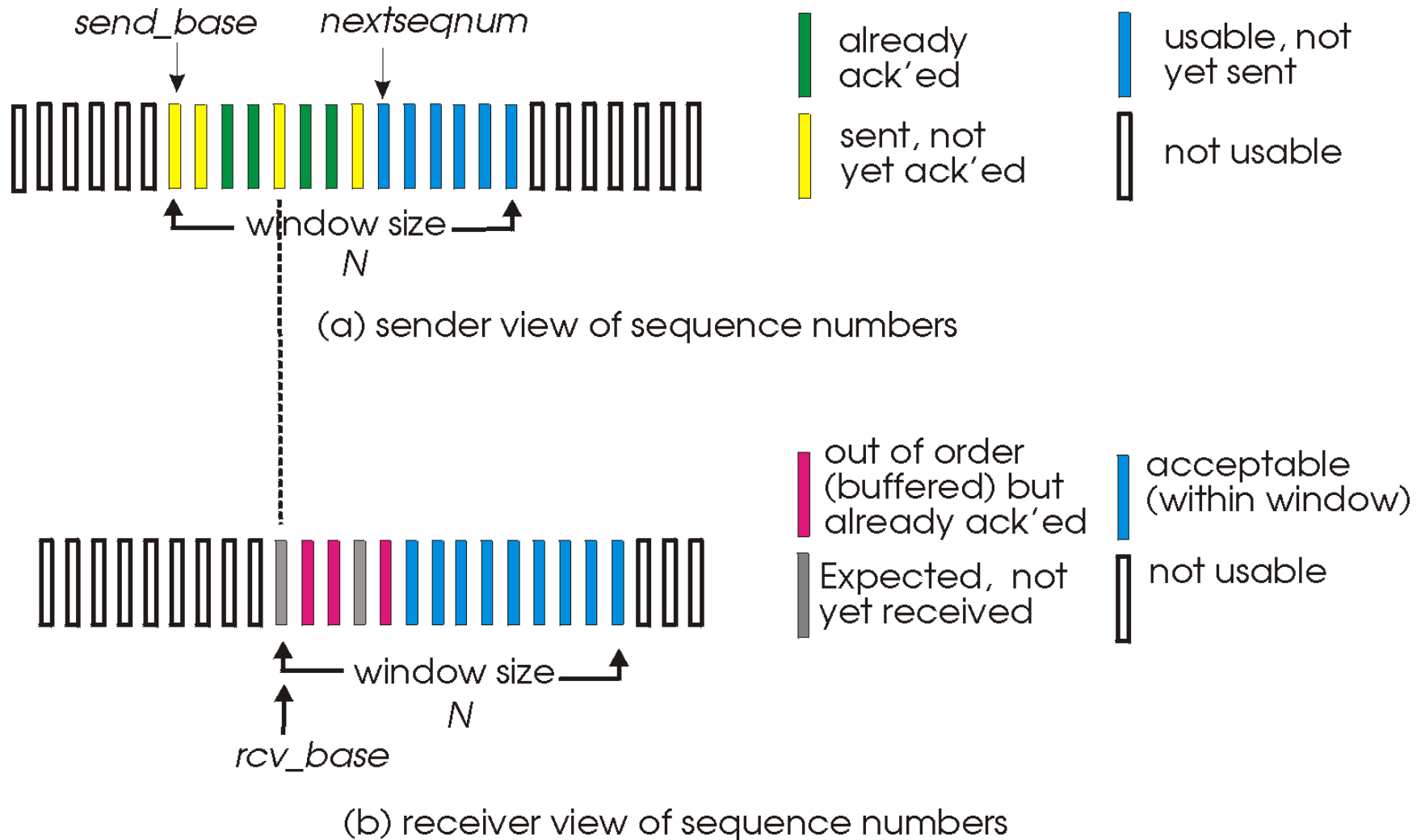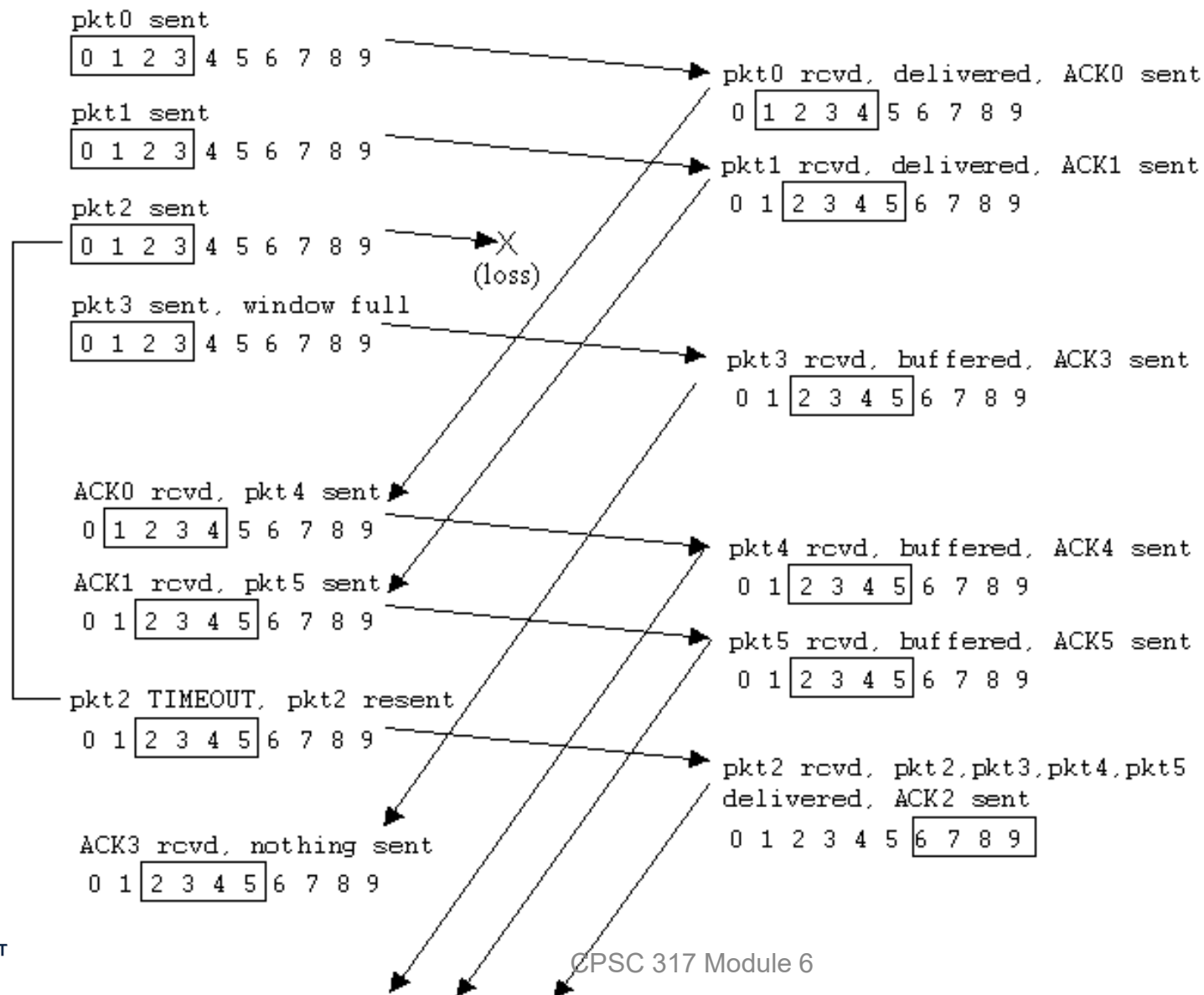- ❑ **pkt n in** [rcvbase-N,rcvbase-1]
- ❑ ACK(n)
- ❑ **otherwise:**
- ❑ ignore

# Selective repeat: sender, receiver windows

# 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 → X (loss)

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

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]

# SEQUENCE NUMBERS

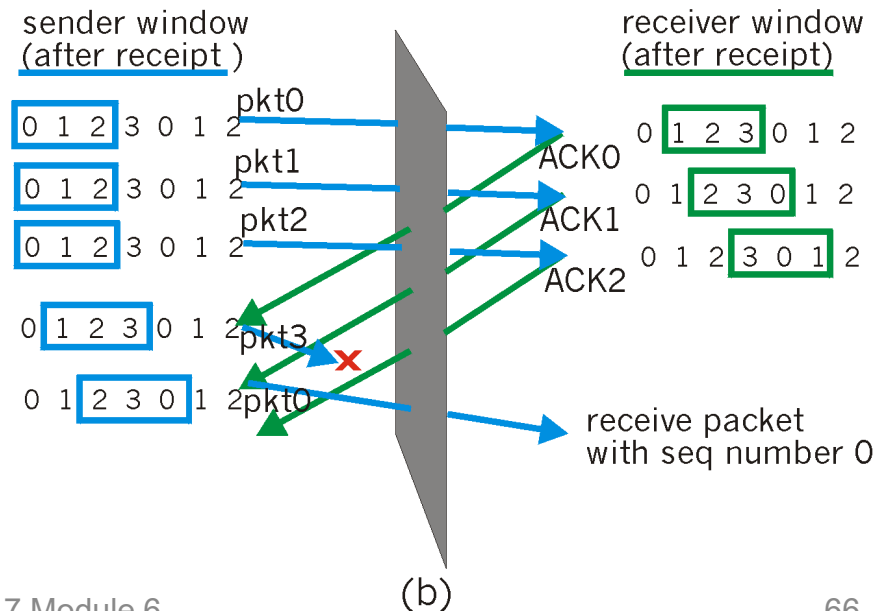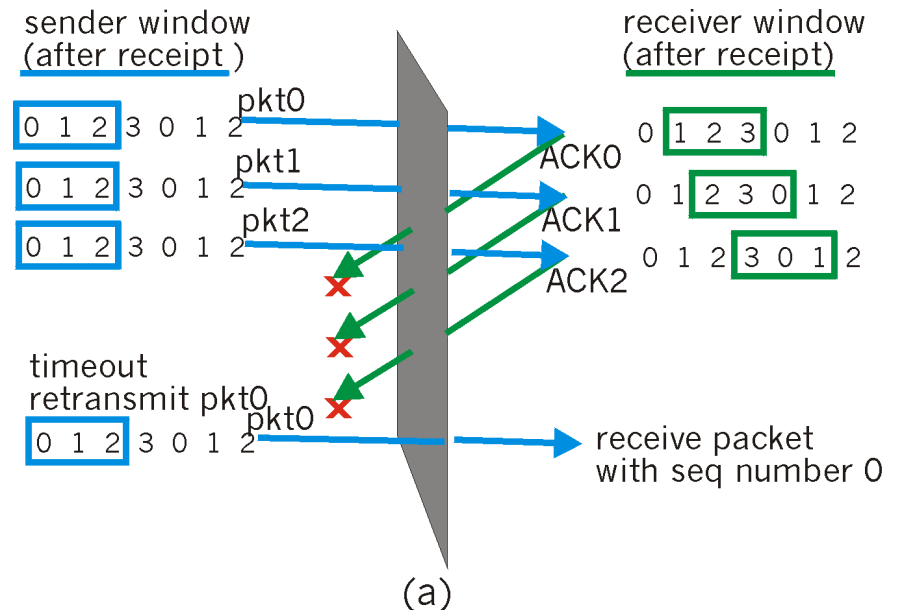# Sequence Number Range

❑ Must fit into K bits

❑ Finite

❑ Is there a limit on the ranges that work?

　o SWS = N, RWS = 1

　o SWS = N, RWS = N

❑ Does it make sense for RWS>SWS?

# Selective repeat: dilemma

## Example:
- ❏ seq #'s: 0, 1, 2, 3
- ❏ window size=3
- ❏ receiver sees no difference in two scenarios!
- ❏ incorrectly passes duplicate data as new in (a)

sender window (after receipt )

receiver window (after receipt)

pkt0
0 1 2 3 0 1 2

pkt1
0 1 2 3 0 1 2

pkt2
0 1 2 3 0 1 2

ACK0
0 1 2 3 0 1 2

ACK1
0 1 2 3 0 1 2

ACK2
0 1 2 3 0 1 2

timeout
retransmit pkt0
0 1 2 3 0 1 2

pkt0

receive packet with seq number 0

(a)

sender window (after receipt )

receiver window (after receipt)

pkt0
0 1 2 3 0 1 2

pkt1
0 1 2 3 0 1 2

pkt2
0 1 2 3 0 1 2

ACK0
0 1 2 3 0 1 2

ACK1
0 1 2 3 0 1 2

ACK2
0 1 2 3 0 1 2

0 1 2 3 0 1 2
pkt3

0 1 2 3 0 1 2 pkt0

receive packet with seq number 0

(b)

THE UNIVERSITY OF BRITISH COLUMBIA
Modified from Kurose-Ross

# Sequence Number Range

❑ Must fit into K bits

❑ Finite

❑ What is the relationship between RWS, SWS and the number of sequence numbers?

# of sequence numbers >= SWS + RWS

# GBN Sequence Space Example

3 sequence numbers
SWS = 3
RWS = 1

**SWS**

| | 0 | 1 | 2 | 0 | 1 | 2 | 0 | 1 | 2 | | |

In 1st case, receiver gets the original zero, fine.

**RWS**

2

**Expecting "0"**

In 2nd case, receiver gets 0, 1, and 2, and ACKS them all, slides and is now expecting the next zero, WRONG

Sequence space must be at least SWS+1

# SR Sequence Space Example



**SWS**

| 0 | 1 | 2 | 0 | 1 | 2 | 0 | 1 | 2 |
|---|---|---|---|---|---|---|---|---|

In 1$^{st}$ case, all packets loss, 1$^{st}$ zero is recv'ed,
In 2$^{nd}$ case, all acks loss, receiver is expecting the 2$^{nd}$ zero i

**RWS**

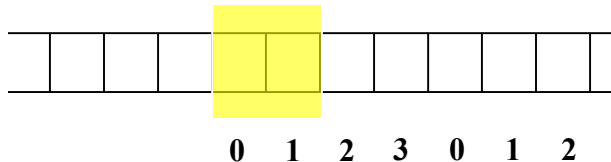0   1   2   0

# SR Example

**SWS**

| 0 | 1 | 2 | 3 | 0 | 1 | 2 | 3 | 0 | | |

Sender: Did 0,1,2 get lost and I need to resend orginal 0, or did 0,1,2 get received and the receiver is expecting the next 0
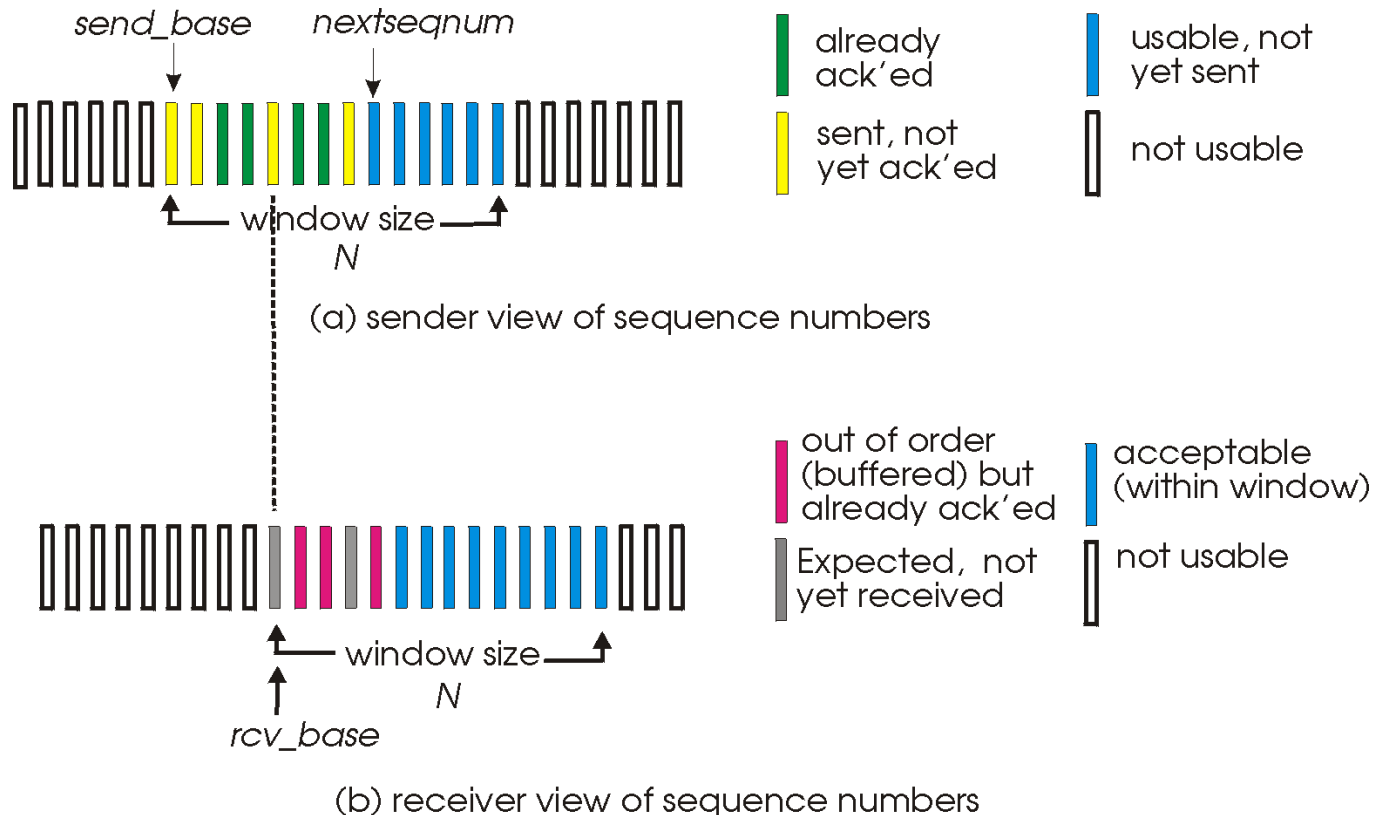
**RWS**

0    1    2    3    0    1    2

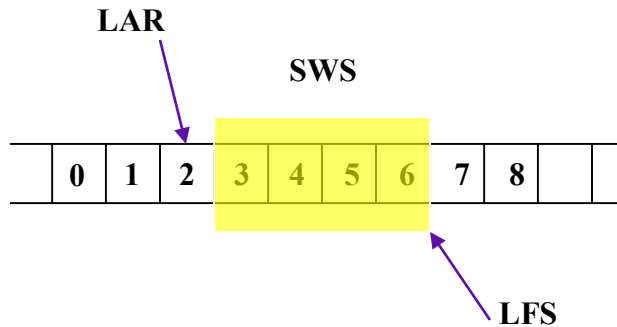**Sequence space SWS + RWS**

# SEQUENCE NUMBERS
# (sliding window)

# What do we ACK?

❑ The packet we just received (Kurose and Ross rdt3.0, and earlier ones)

❑ Sliding window (Kurose and Ross)

- o ACK the packet we just received
- o Cumulative ACK, ACK the largest in-order received packet

❑ Same as above but ACKing next expected packet rather than the one received. (TCP)

# Selective repeat (vs GBN, vs CUMULATIVE)

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

# Sliding Window (TCP like)

**LAR**

**SWS**

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | |
|---|---|---|---|---|---|---|---|---|---|

**LFS**

**Sender**:

    if more data to send (LFS-LAR < SWS)

        then send data, LFS++

    if recv'ed ACK for LAR+1

        then LAR++

    if timer expires

        then send LAR+1

LAR: last ACK received
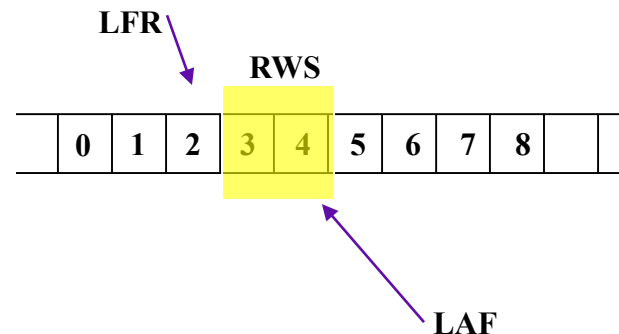
LFS: last frame sent

**Receiver**:

    if recv'ed K > LAF
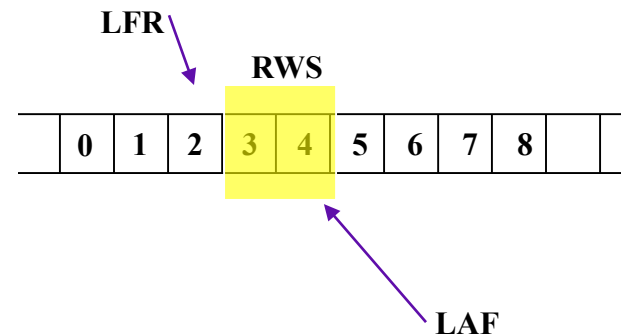
        then discard

    else

            if K == LFR+1

            then

                store

            LFR++, LAF++ (slide window)

ACK, LFR (after it was incremented)

**LFR**

**RWS**

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | |
|---|---|---|---|---|---|---|---|---|---|

**LAF**

LFR: last frame received

LAF: largest acceptable frame

# Sliding Window (TCP like)

LAF

SWS

| | | | 3 | 4 | 5 | 6 | | | | |
|0|1|2| | | | |7|8| | |

LFS

**Sender**:

if more data to send (LFS-LAR < SWS)

then send data, LFS++

if recv'ed ACK for LAR+1

then LAR++

if timer expires

then send LAR+1

LAR: last ACK received
LFS: last frame sent

**Receiver**:

if recv'ed K > LAF

then discard

else

if K == LFR+1

then

store

LFR++, LAF++ (slide window)

ACK, LFR (after it was incremented)

LFR

RWS

| | | | 3 | 4 | | | | | |
|0|1|2| | |5|6|7|8| |

LAF

LFR: last frame received
LAF: largest acceptable frame

# Sequence Numbers Cases

**# of sequence numbers >= SWS + RWS**

❑ RWS > SWS?

❑ SWS > # of sequence numbers

❑ SWS=RWS=1

❑ SWS=N, RWS=1

❑ SWS=RWS=N

# Summary

RDT:

❑ Added retransmit, checksum, sequence numbers, acknowledgments, and timers.

❑ Pipelined, needed to introduce sliding windows and more sequence number space.

❑ Still cannot handle out of order packets (i.e. packet A leaves before packet B, but packet B arrives before A, or to say that an earlier packet in transit can arrive later than a packet sent after the earlier packet)

<span style="color:red">STILL A PROBLEM</span>

Strategies for Sliding Window
     Go-back-N
     Selective Repeat

Slight variations of the above, cumulative ACK or next packet instead of last one.

Sliding window makes it possible to improve throughput and a mechanism for flow control.

# TCP

THE UNIVERSITY OF BRITISH COLUMBIA

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

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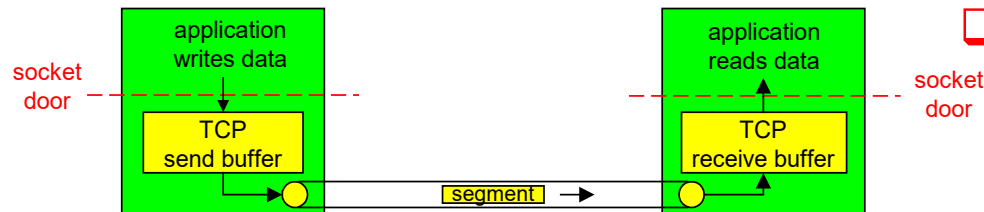- one sender, one receiver

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

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

❑ ***send & receive buffers***



❑ **full duplex data:**
- bi-directional data flow in same connection
- MSS: maximum segment size

❑ **connection-oriented:**
- handshaking (exchange of control msgs) init's sender, receiver state before data exchange

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

# TCP

❑ <span style="color:red">What's in the header?</span>

❑ Sliding window

❑ RTT estimation

❑ More on sliding window

❑ Flow control

❑ Connection management

❑ Congestion

# 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)

| source port # | dest port # |
| sequence number | |
| acknowledgement number | |
| head len | not used | U A P R S F | Receive window |
| checksum | Urgent data ptr |
| Options (variable length) | |
| application data (variable length) | |

counting
by bytes
of data
(not segments!)

# bytes
rcvr willing
to accept

# TCP

❑ What's in the header?

❑ Sliding window

❑ RTT estimation

❑ More on sliding window

❑ Flow control

❑ Connection management

❑ Congestion

THE UNIVERSITY OF BRITISH COLUMBIA
Modified from Kurose-Ross

# TCP Sliding Window

❑ Cumulative acknowledgements

❑ Store out of order frames that are within the size of the receive window

❑ ACK next expected byte

❑ Sequence number is of the first byte in segment
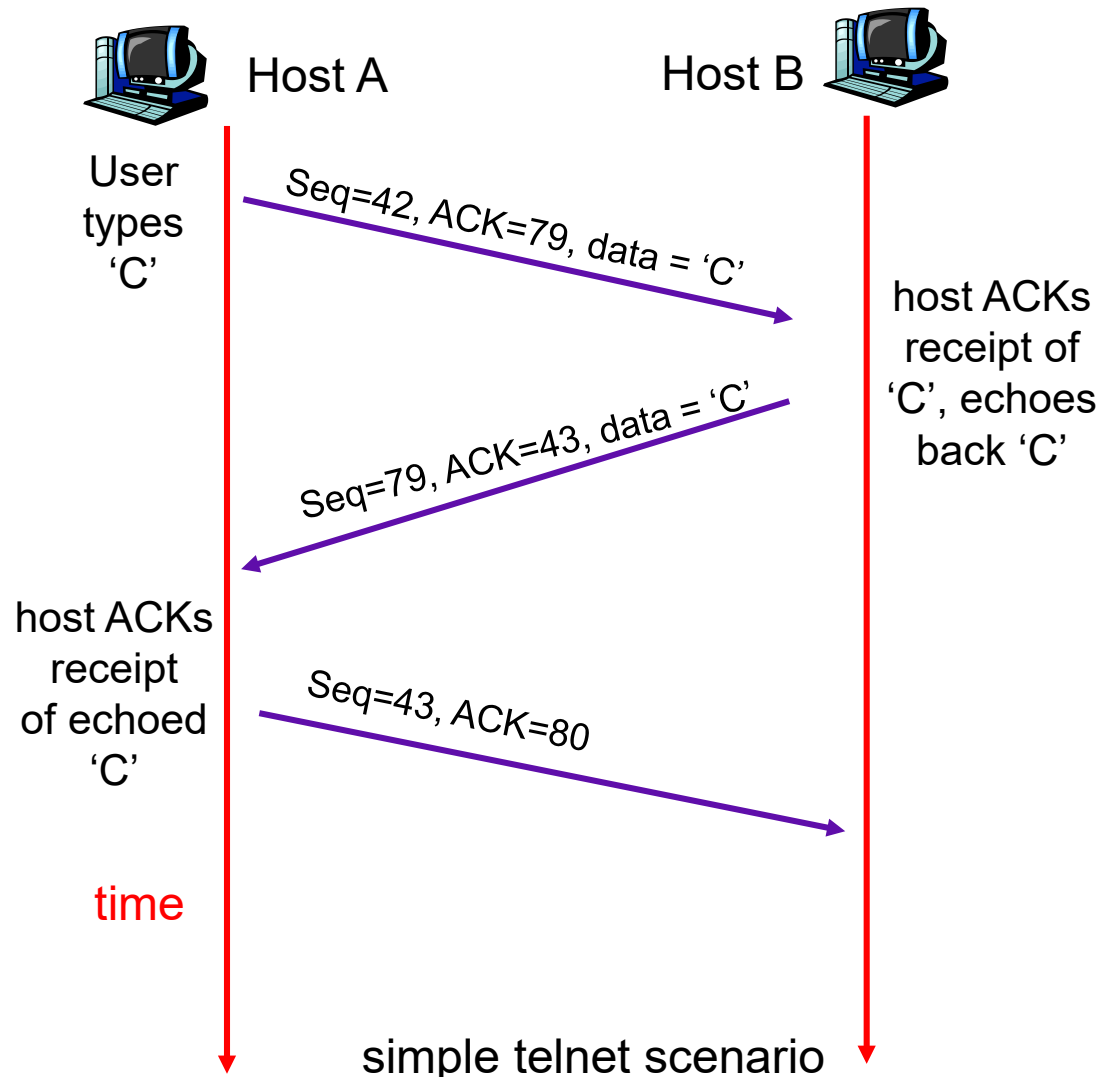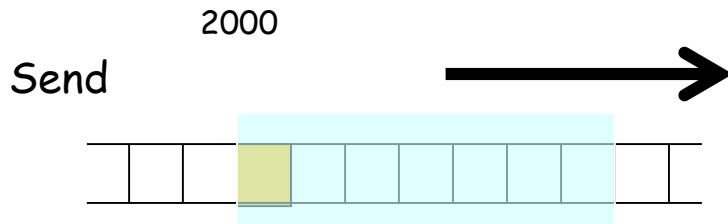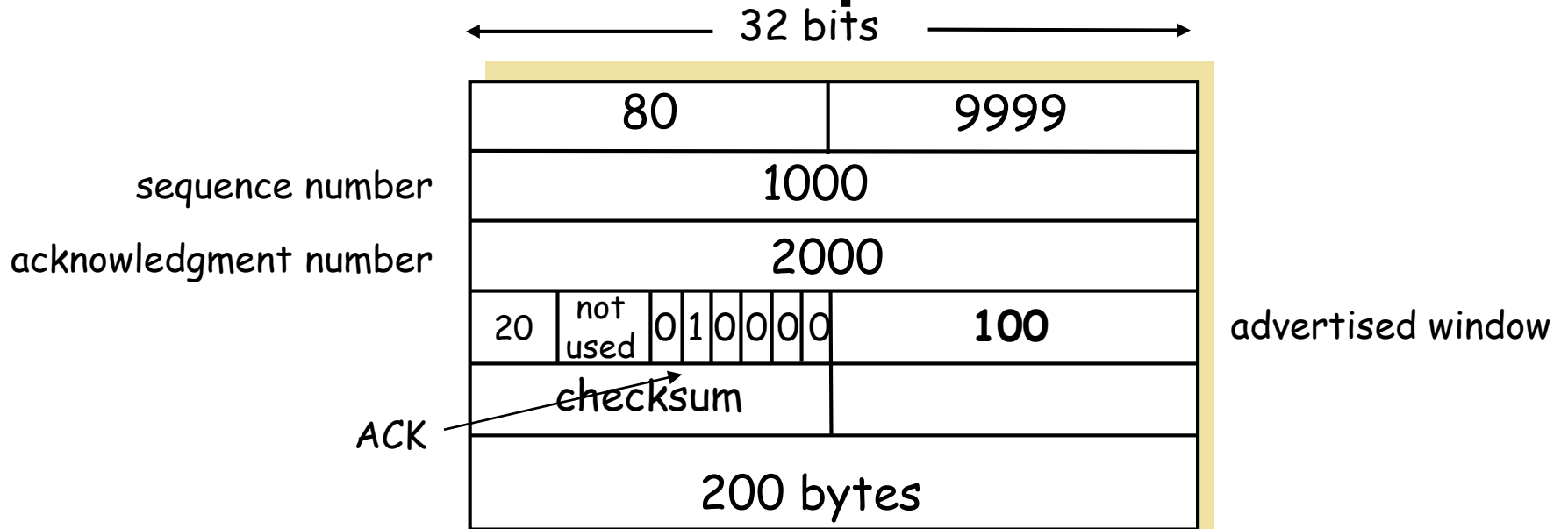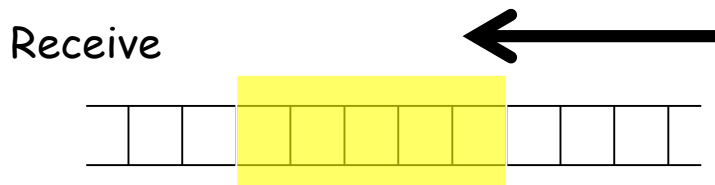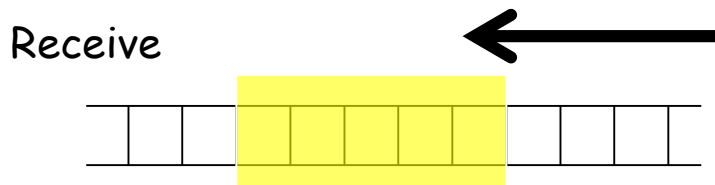
❑ Variations of TCP: TCP-vegas, TCP-reno, TCP-sack

# 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

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

time

simple telnet scenario

# TCP seq. #'s and 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

time

simple telnet scenario

# TCP

❑ What's in the header?

❑ Sliding window

❑ RTT estimation

❑ More on sliding window

❑ Flow control

❑ Connection management

❑ Congestion

CPSC 317  Module 6
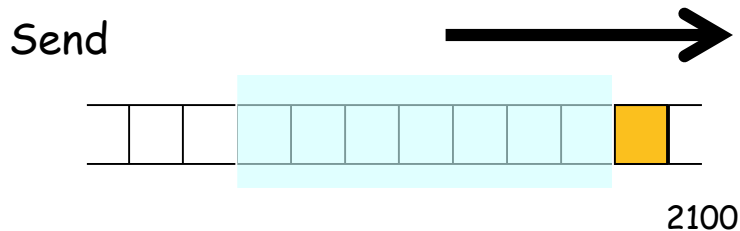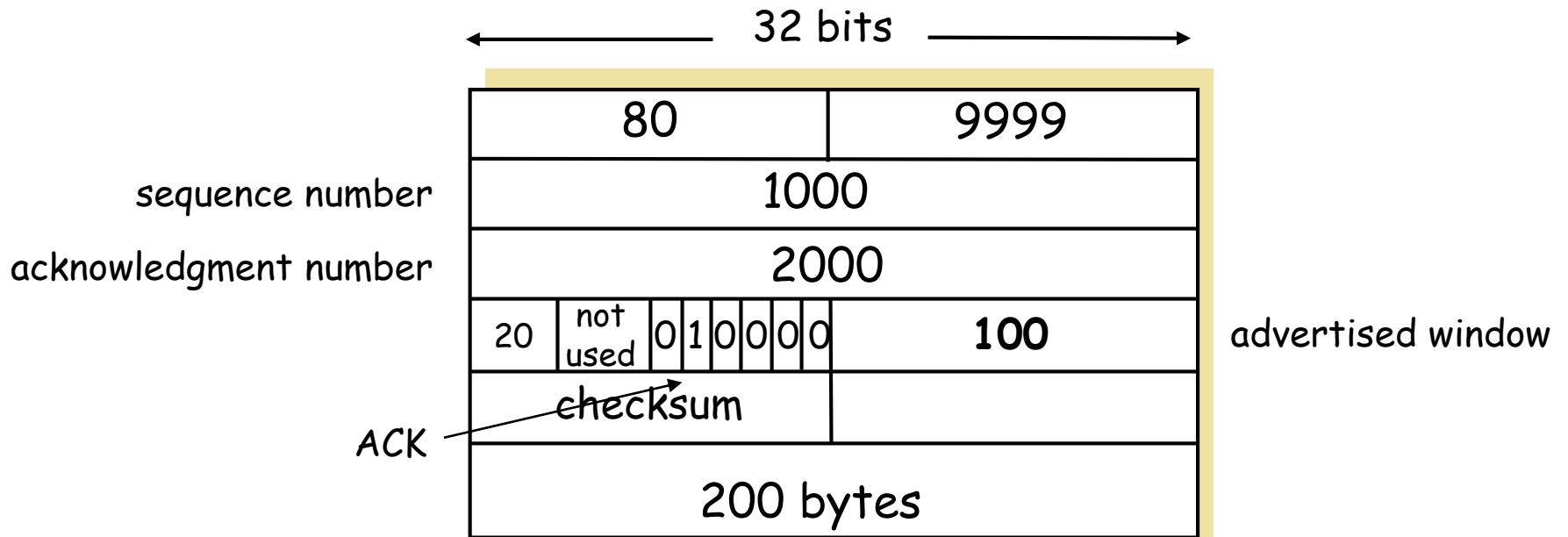
# Example



32 bits

| 80 | 9999 |
|---|---|
| sequence number | 1000 |
| acknowledgment number | 2000 |
| 20 | not used | 0 | 1 | 0 | 0 | 0 | 0 | **100** |
| ACK | checksum | |
| 200 bytes | |

advertised window

2000

Send

**What is the minimum that A's SendBase value can be when A is recv'ing this segment?**

**Assume the segment IS a duplicate ACK**

Receive

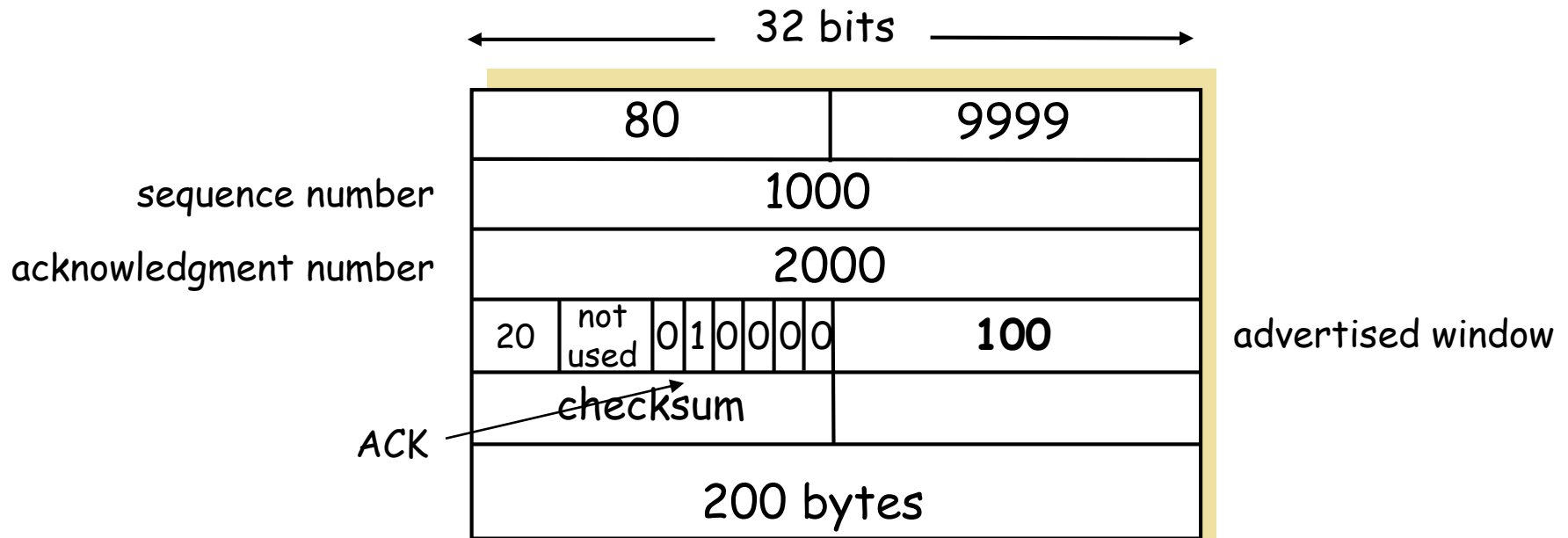THE UNIVERSITY OF BRITISH COLUMBIA
Modified from Kurose-Ross

# Example



32 bits

| 80 | 9999 |
|---|---|

sequence number — 1000

acknowledgment number — 2000

| 20 | not used | 0 | 1 | 0 | 0 | 0 | 0 | **100** |

advertised window

checksum

ACK

200 bytes

Send → 2100

**What is the maximum that A's NextSeqNum value can be when A is recv'ing this segment.**

Receive ←

**Assume the segment is NOT a duplicate ACK**
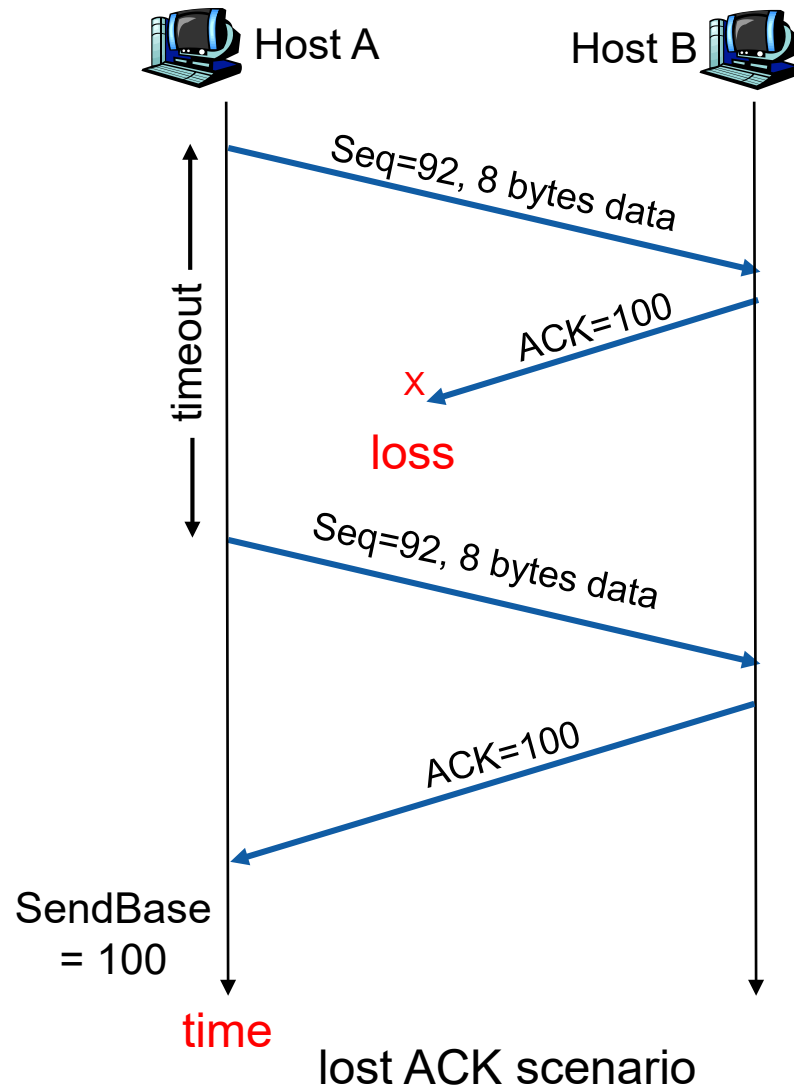
# Example

32 bits

| 80 | 9999 |
|---|---|
| sequence number | 1000 |
| acknowledgment number | 2000 |
| 20 \| not used \| 0 1 0 0 0 0 | 100 |
| checksum | |
| 200 bytes | |

advertised window

ACK

Send

249

Receive

950

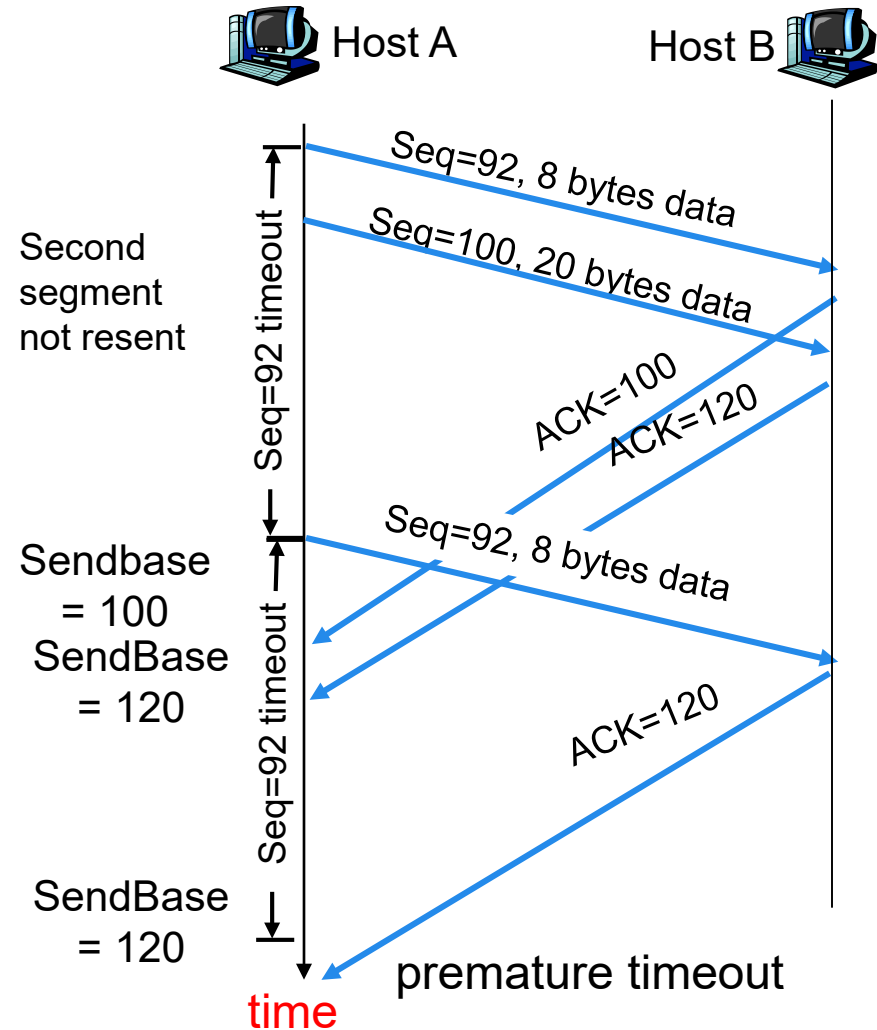**Suppose that B's data byte 950 is in A's recv'ing window.**

**What is the minimum amount of space in bytes, that A's receiving window must possess <u>beyond</u> byte 950.**
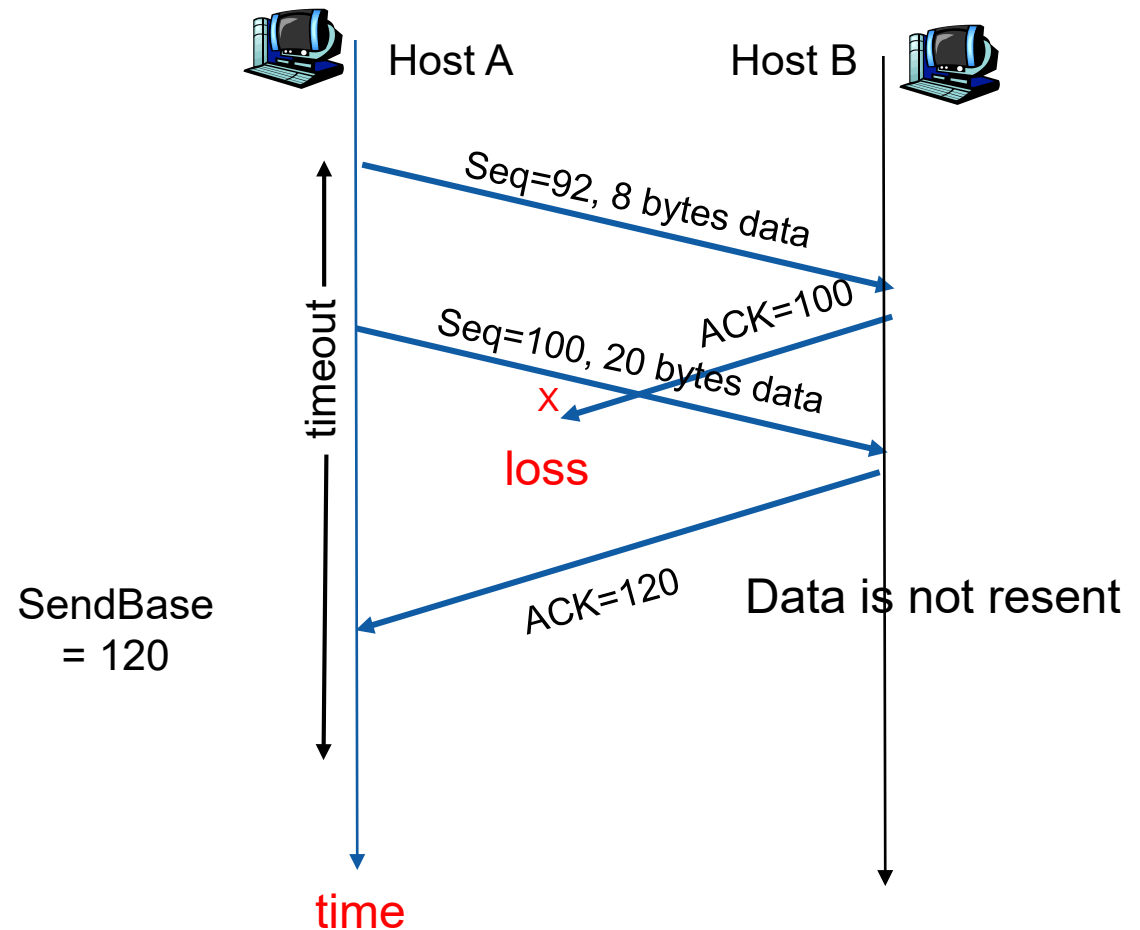
# TCP: retransmission scenarios



lost ACK scenario

# TCP: retransmission scenarios



Host A

Host B

Seq=92, 8 bytes data

Seq=100, 20 bytes data

Second segment not resent

Seq=92 timeout

ACK=100

ACK=120

Sendbase = 100

SendBase = 120

Seq=92 timeout

Seq=92, 8 bytes data

ACK=120

SendBase = 120

premature timeout

time

# TCP retransmission scenarios (more)



Host A                    Host B

timeout

Seq=92, 8 bytes data

ACK=100

Seq=100, 20 bytes data

x

loss

SendBase
= 120

ACK=120

Data is not resent

time

Cumulative ACK scenario

# TCP

❑ What's in the header?

❑ Sliding window

❑ <span style="color:red">RTT estimation</span>

❑ More on sliding window

❑ Flow control

❑ Connection management

❑ Congestion

# Time-out Problem?

❑ Rate depends on the window sizes, loss rate and round-trip time in acknowledging data.

❑ BUT
- ○ Congestion in the Internet
- ○ Conditions at the end-stations
- ○ Properties of the network
- ○ Size and timing of data segments

❑ Through-put rate is going to vary dynamically

# Time-out and Retransmission

➢ Purpose: sets timer for each segment, retransmit earliest unacknowledged segment when timer goes off. (one timer, adds segment to retransmission queue)

➢ Problem: In the Internet we don't a priori know the RTT of a segment.  It is going to vary depending on the traffic.

## TIMER MANAGEMENT

# Setting the Time-out value?

❑ too short: premature timeout
   o unnecessary retransmissions
   o add to network congestion
   o Retransmission (hurts everyone)

❑ too long: slow reaction to segment loss
   o sluggish performance
   o slow
   o delayed transmission (hurts you, helps everyone)

# How to estimate RTT?

❑ Static time-out?   No!

❑ Adaptive time-out Yes!

 o Estimating time-out is difficult because

  • Peer TCP entity may accumulate acknowledgements and not acknowledge immediately

  • For retransmitted segments, can't tell whether acknowledgement is response to original transmission or retransmission

  • Network conditions may change suddenly

 o Better over-estimate than under-estimate!

# RTT variance (Comer)



**Figure 13.10** A plot of Internet round trip times as measured for 100 successive IP datagrams. Although the Internet now operates with much lower delay, the delays still vary over time.

# How to estimate RTT?

❑ **SampleRTT**:

     ○ Moving average

**(smoothed estimate SRTT)**

| | | 2 | 1 | 8 | 5 | | |
|---|---|---|---|---|---|---|---|
| 2 | 5 | 2 | 1 | 8 | 5 | 3 | 4 |

**4**

     ○ Exponential weighted moving average

# Initially TCP used:

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

- ❑ Exponential weighted moving average. (A dial we can adjust to change the sensitivity of RTT-estimate to history)
- ❑ typical value: $\alpha$ **=** 0.125
- ❑ Time-out is a constant times EstimatedRTT

- ❑ **Time-out = β × EstimatedRTT**
- ❑ Recommended setting for β was 2

# Example RTT estimation:

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

THE UNIVERSITY OF BRITISH COLUMBIA
Modified from Kurose-Ross

# A better Estimate

❑ Research (Jacobson) showed that this estimate did not respond quickly in high variance situations.

❑ 1989 TCP specification required estimates of both average and <span style="color:red">variance</span>

# High Variance (Comer)



Figure 13.12 The TCP retransmission timer for the data from Figure 13.10. Arrows mark two successive datagrams where the delay doubles.

CPSC 317 Module 6

# Summary

❑ Need both average and variance and being selective

❑ Want to avoid time-outs

❑ Existing techniques use selective sampling using estimates of the average variance of the RTT time

❑ Internet and TCP makes it difficult to predict

# TCP

❑ What's in the header?

❑ Sliding window

❑ RTT estimation

❑ More on sliding window

❑ Flow control

❑ Connection management

❑ Congestion  -- omit

THE UNIVERSITY OF BRITISH COLUMBIA
Modified from Kurose-Ross

# TCP Connection Management

**<u>Recall:</u>** TCP sender, receiver establish "connection" before exchanging data segments

❑ initialize TCP variables:
   - o seq. #s
   - o buffers, flow control info (e.g. `RcvWindow`)

❑ *client:* connection initiator

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

❑ *server:* contacted by client

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

**<u>Three way handshake:</u>**

<u>Step 1:</u> client host sends TCP SYN segment to server
   - o specifies initial seq #
   - o no data

<u>Step 2:</u> server host receives SYN, replies with SYNACK segment
   - o server allocates buffers
   - o specifies server initial seq. #

<u>Step 3:</u> client receives SYNACK, replies with ACK segment, which may contain data
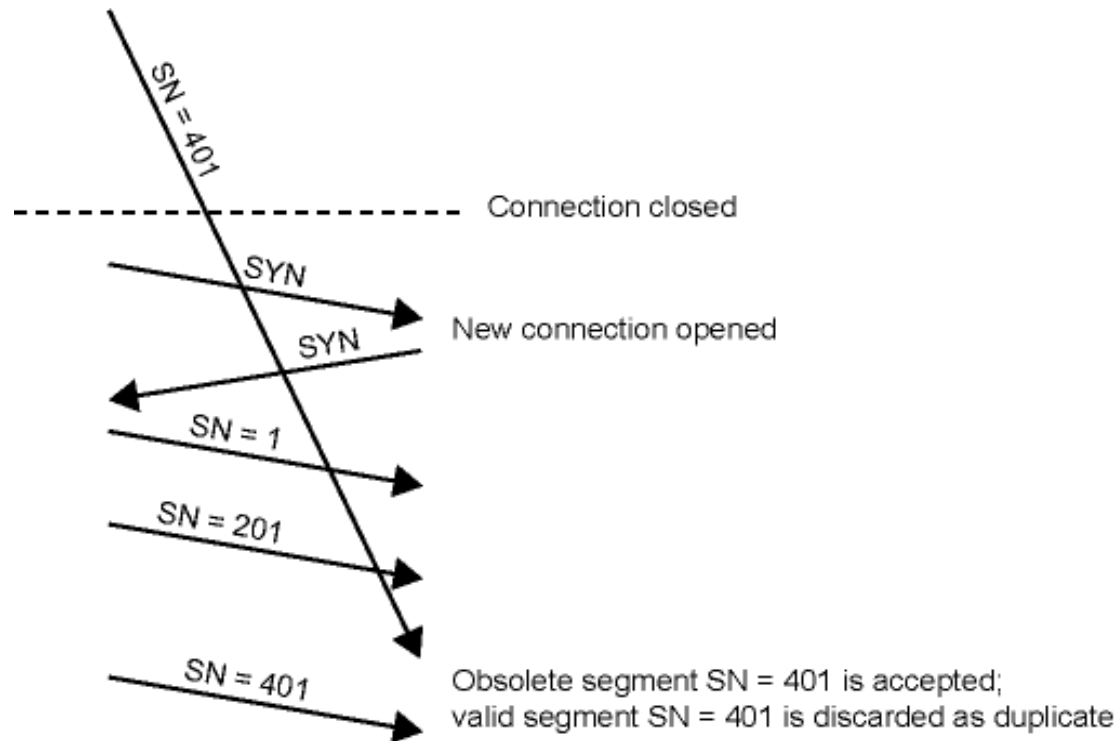
# Initial Sequence Number (ISN)

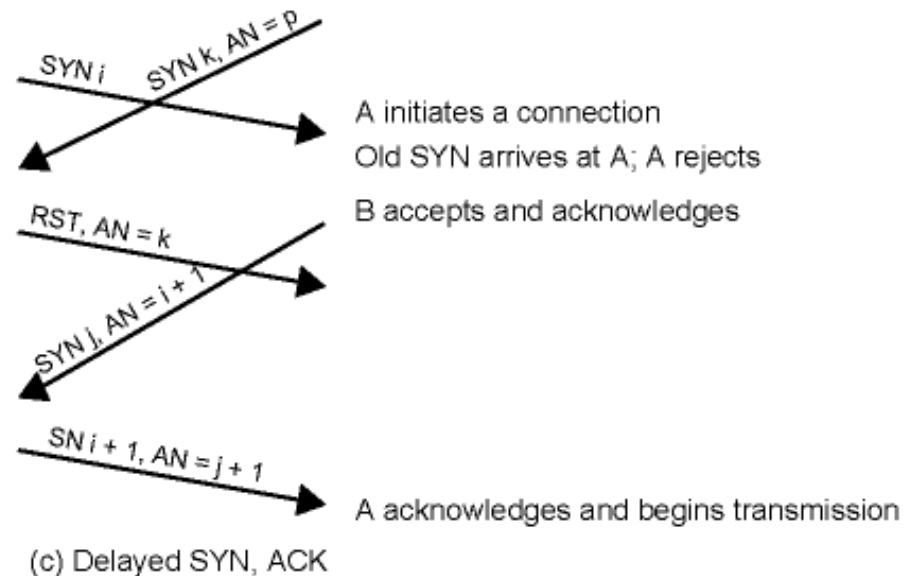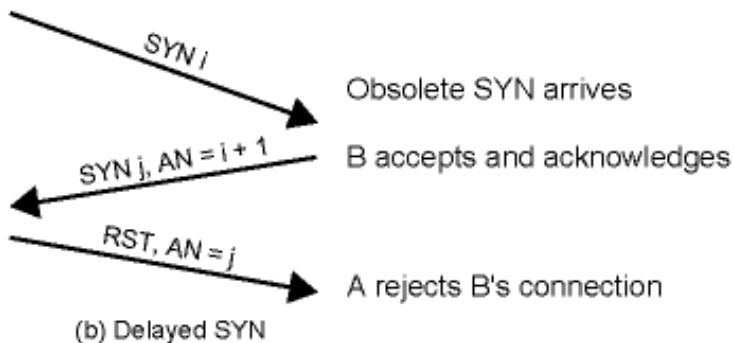Client                                    Server
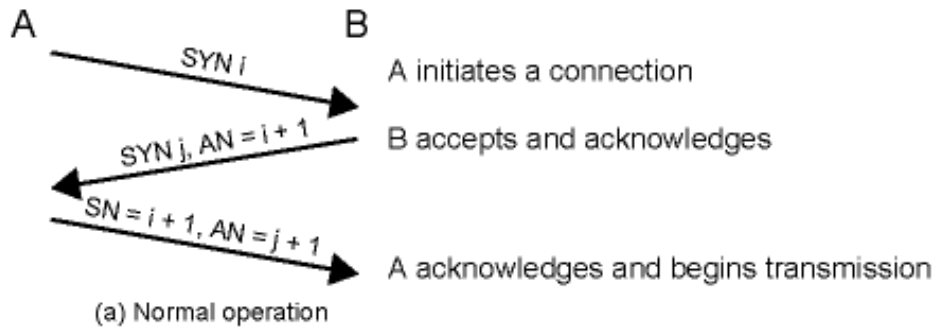
# Single ISN problems



- ❑ Add a unique session ID to each stream
- ❑ But what if machine re-boots, clocks?
  - ○ What if the machine re-boots?

# Some Possible Scenarios

**ISN initial sequence number**

**Assumption: MSL (maximum segment lifetime --- two minutes)**



(a) Normal operation

(b) Delayed SYN

(c) Delayed SYN, ACK

# Closing a connection

❑ Objective

  ○ Close without abruptly dropping the connection!

# Graceful Close

❑ Send FIN i and receive FINACK i

❑ Receive FIN j and send FINACK j

❑ Wait twice maximum expected segment lifetime
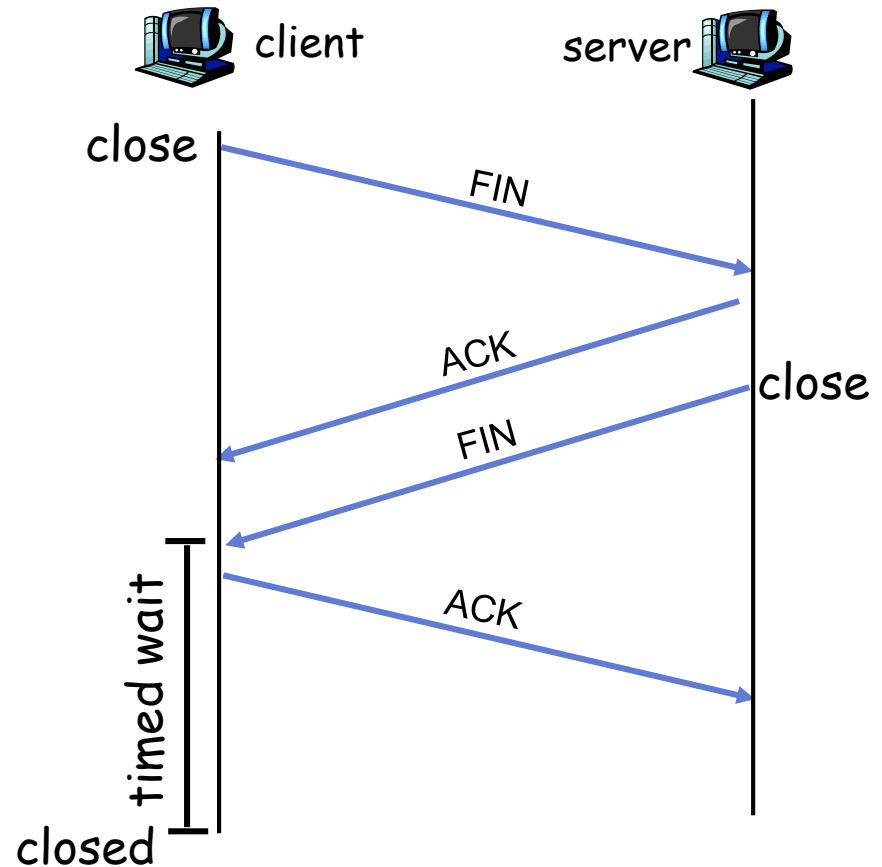
# TCP Connection Management (cont.)

**Closing a connection:**

client closes socket:
   **clientSocket.close();**

**Step 1:** client end system sends TCP FIN control segment to server

**Step 2:** server receives FIN, replies with ACK. Closes connection, sends FIN.
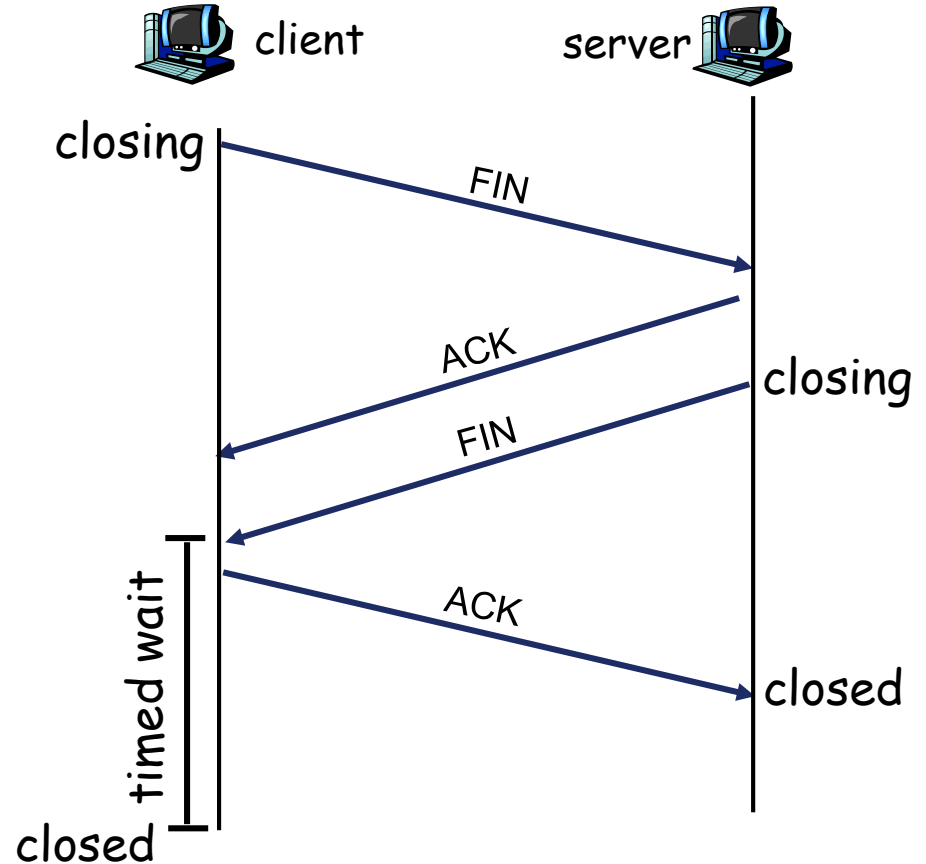
# TCP Connection Management (cont.)

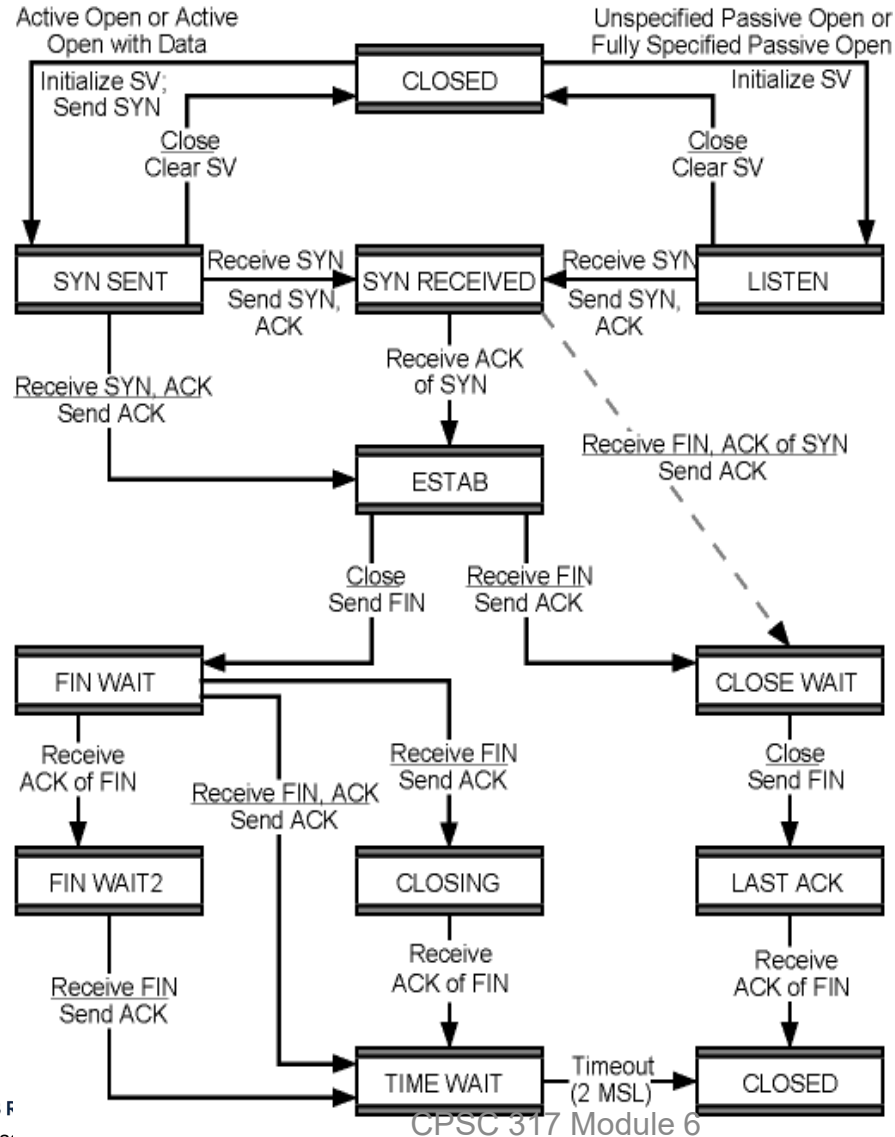Step 3: client receives FIN, replies with ACK.

- o Enters "timed wait" - will respond with ACK to received FINs

Step 4: server, receives ACK. Connection closed.

Note: with small modification, can handle simultaneous FINs.

# TCP State Machine

THE UNIVERSITY OF BR
Modified from Kurose-Ros

# Steven's TCP state machine