

# CS 305: Computer Networks

## Fall 2022

### **Lecture 7: Transport Layer**

**Ming Tang**

Department of Computer Science and Engineering  
Southern University of Science and Technology (SUSTech)

# 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

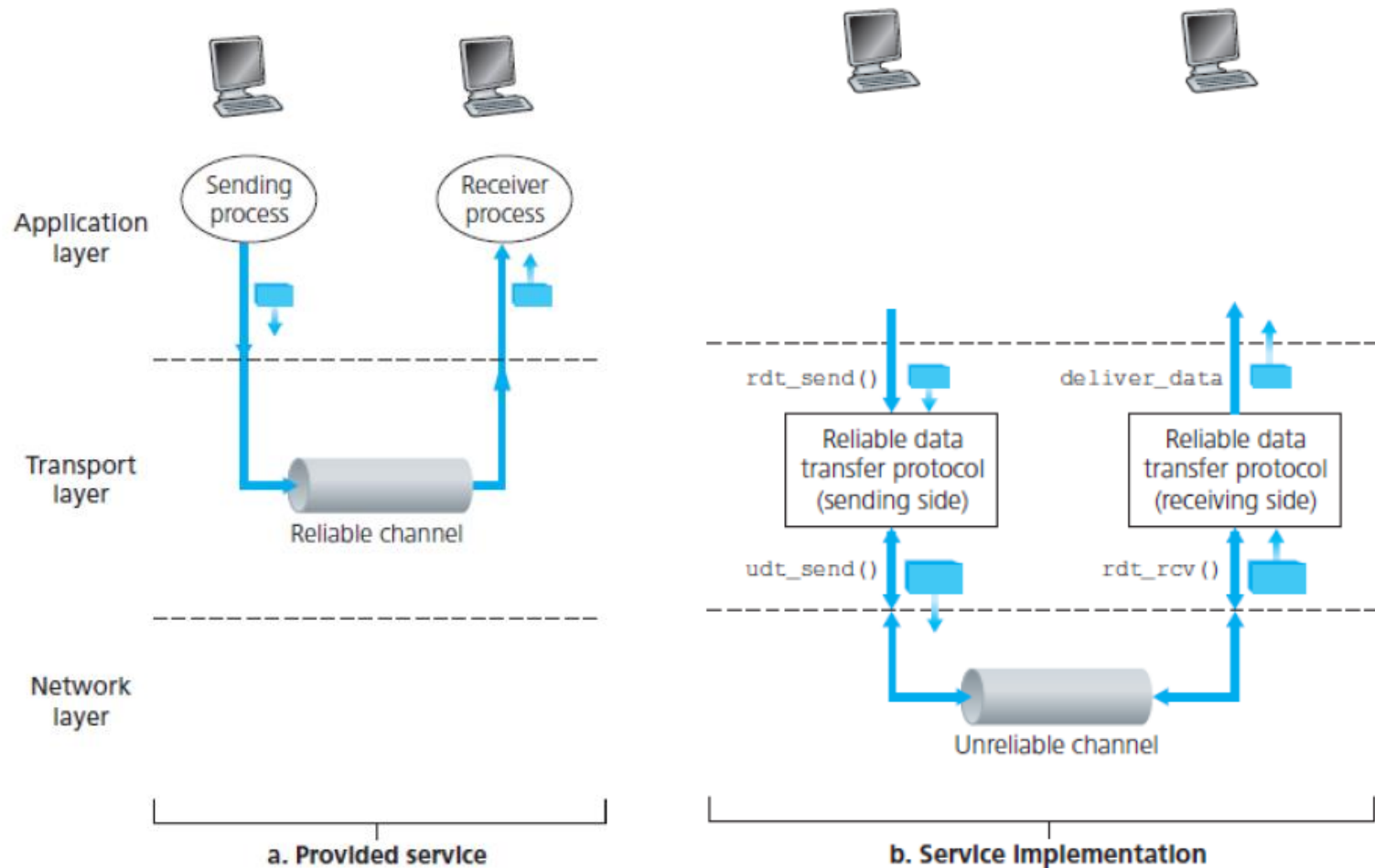
3.5 connection-oriented transport: TCP

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

3.6 principles of congestion control

3.7 TCP congestion control

# Reliable Data Transfer (rdt)



# Overview

---

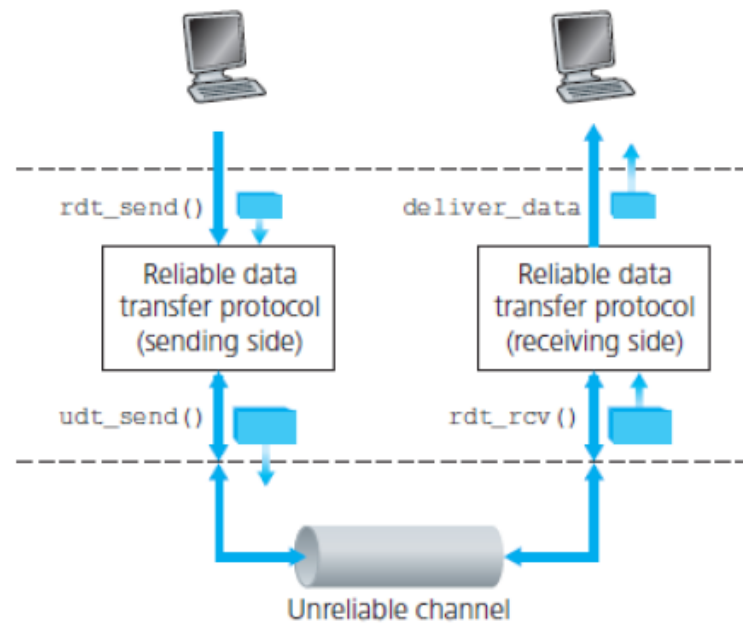
## Roadmap:

- ❖ Perfectly reliable channel: rdt1.0
- ❖ Channel with bit error:
  - bit error in packet: rdt 2.0
  - bit error in ACK: 2.1
  - NAK-free: 2.2
- ❖ Lossy channel: rdt 3.0

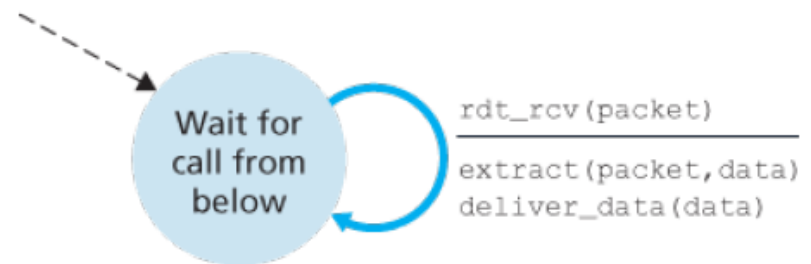
## Summary of Techniques

- Checksum
- Sequence number
- ACK packets
- Retransmission
- Timeout

# rdt1.0: reliable transfer over a reliable channel



sender

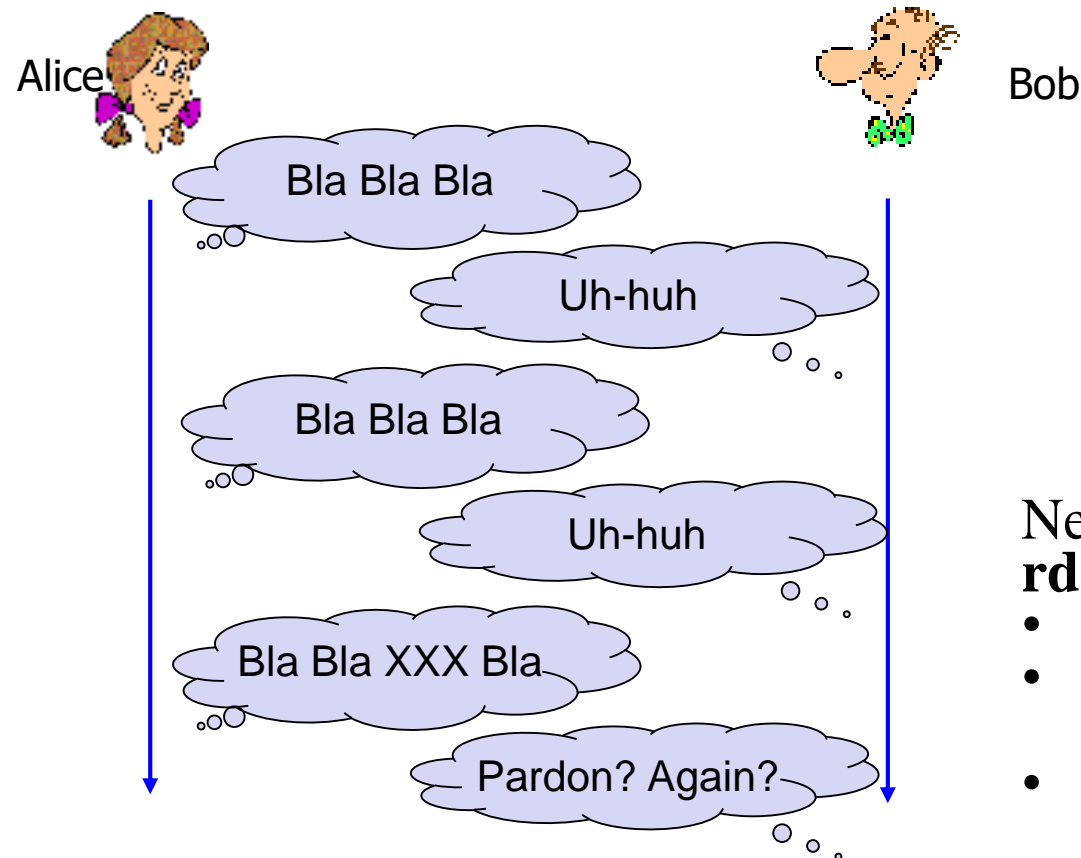


receiver

# rdt2.0: channel with bit errors

- ❖ Underlying channel may **flip bits** ( $0 \rightarrow 1$ ) in packet

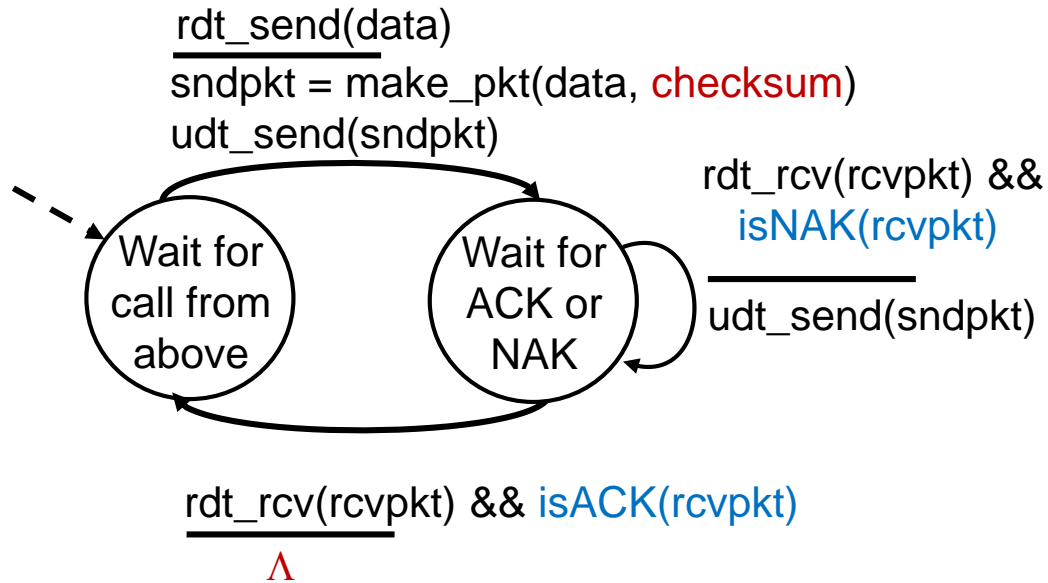
How do humans recover from “errors” during conversation?



New mechanisms in **rdt2.0** (beyond **rdt1.0**):

- error detection
- receiver feedback: control msgs (ACK, NAK) rcvr- $\rightarrow$ sender
- retransmission

# rdt2.0: FSM specification

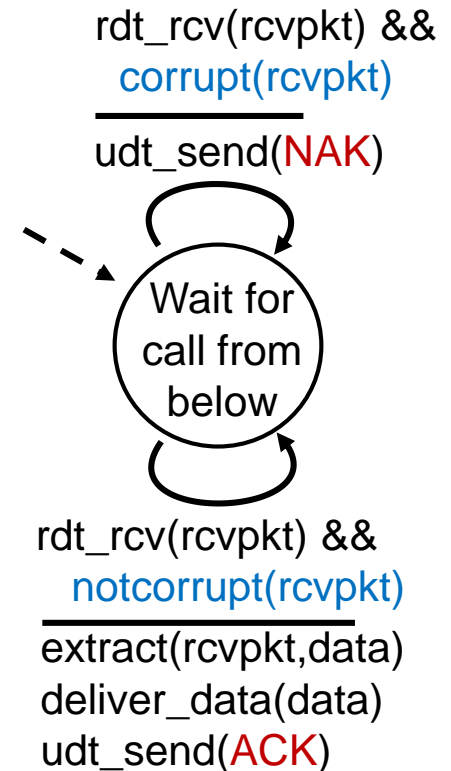


sender

## Stop and wait

Sender sends one packet,  
then waits for receiver  
response

receiver



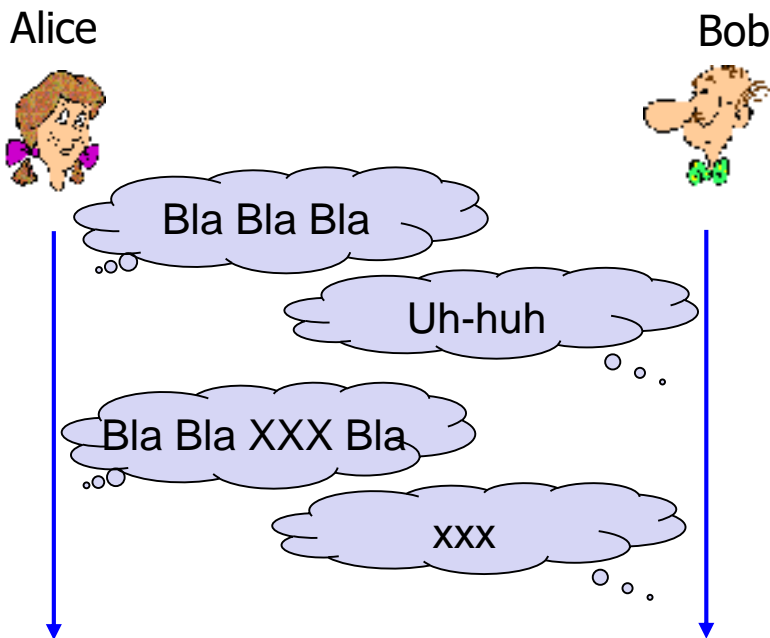
# rdt2.0 has a fatal flaw!

The possibility that ACK or NAK packet could be corrupted:

- ❖ Checksum bits

Handling corrupted ACKs or NAKs:

- ❖ Option 3: when garbled ACK or NAK, retransmit



**Problem:** can't just retransmit: new data or retransmission? possible duplicate

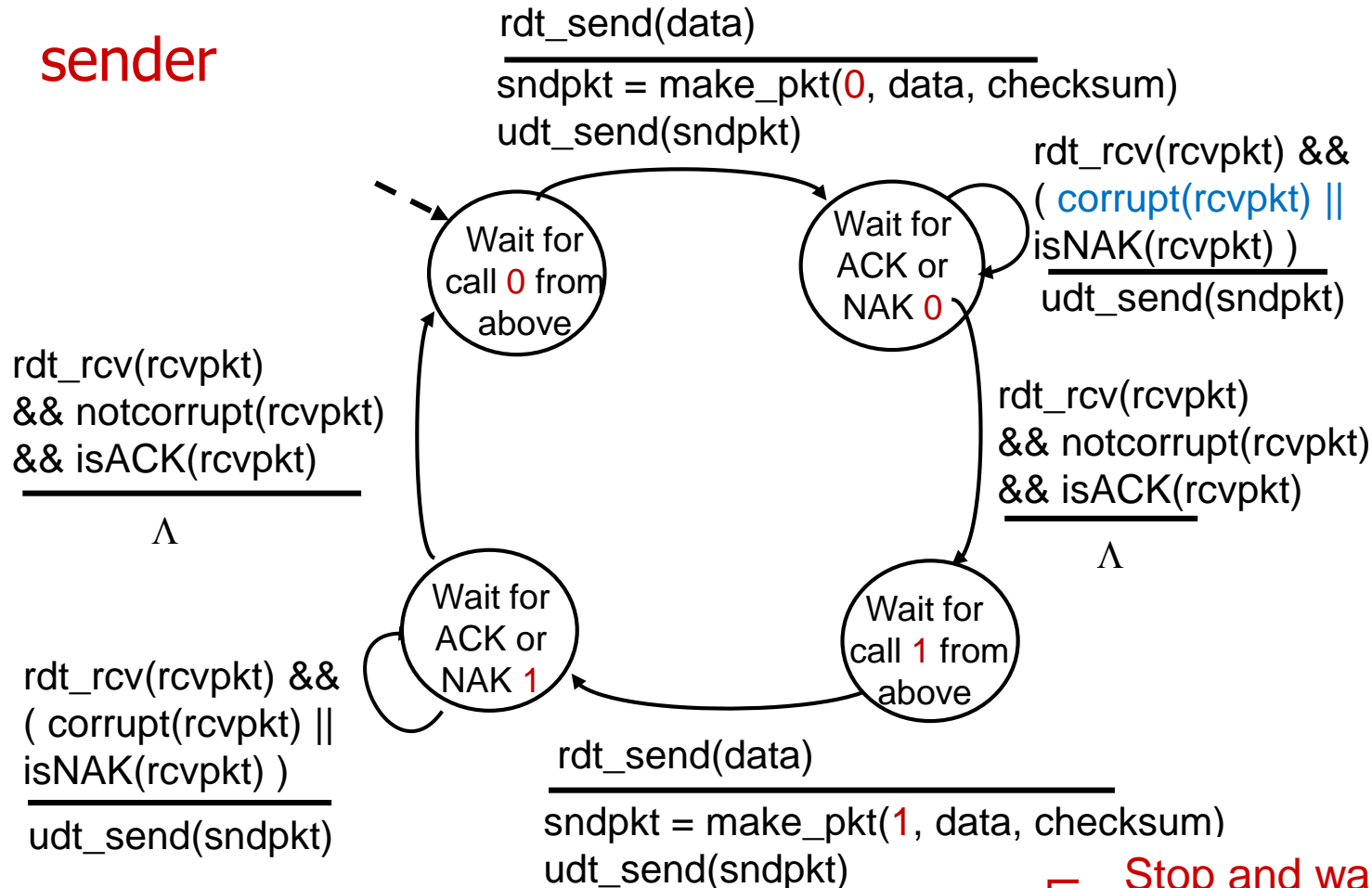
**Handling duplicates:**

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



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

sender



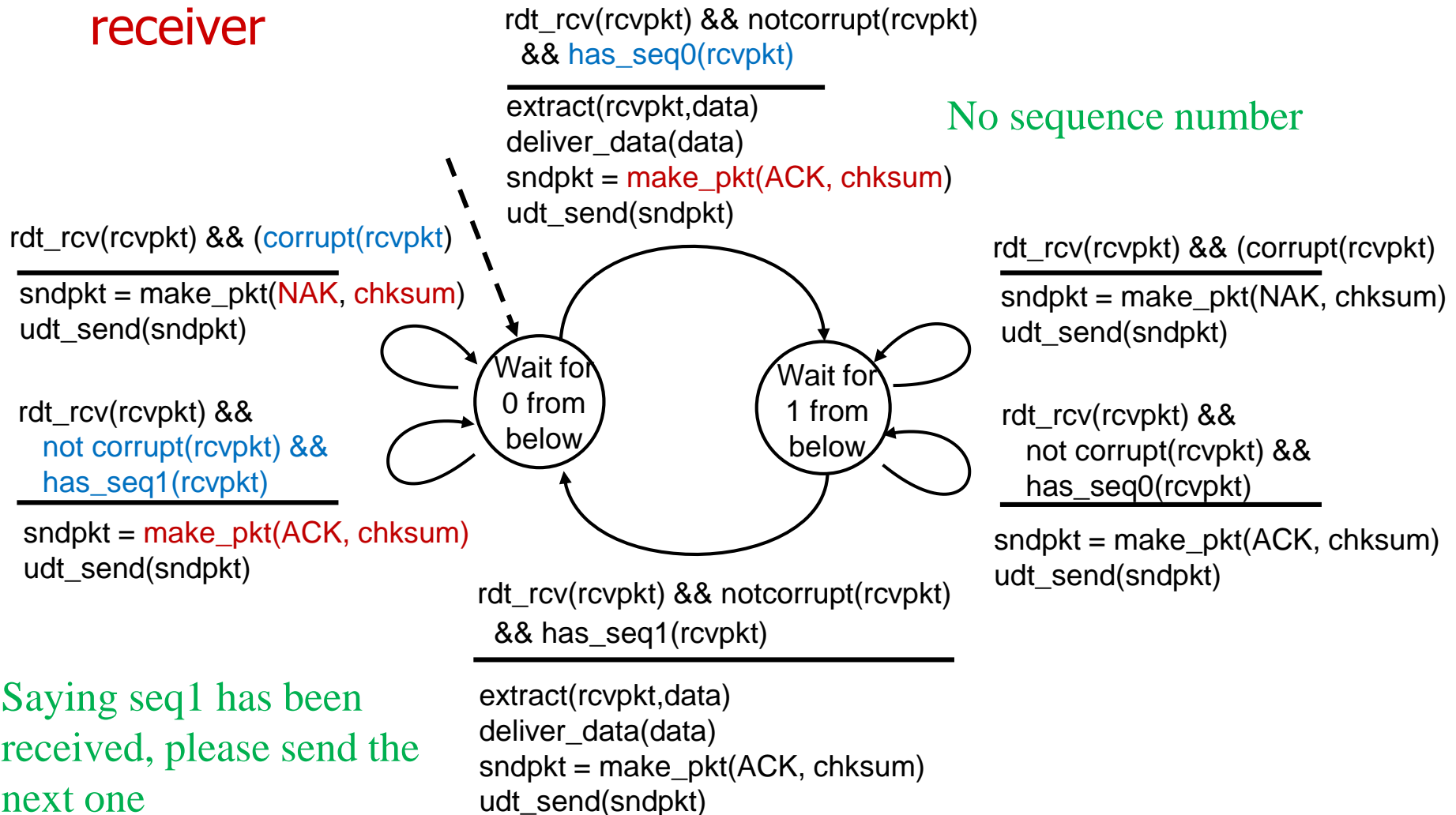
**Stop and wait**

Sender sends one packet, then waits for receiver response

Two sequence number would be sufficient!

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

## receiver



# rdt2.2: a NAK-free protocol

- ❖ same functionality as rdt2.1, **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*

# rdt2.2: sender, receiver fragments

sender

rdt\_send(data)

sndpkt = make\_pkt(0, data, checksum)

udt\_send(sndpkt)

Wait for  
call 0 from  
above

Wait for  
ACK  
0

sender FSM  
fragment

rdt\_rcv(rcvpkt) &&  
( corrupt(rcvpkt) ||  
isACK(rcvpkt,1) )  
udt\_send(sndpkt)

rdt\_rcv(rcvpkt)  
&& notcorrupt(rcvpkt)  
&& isACK(rcvpkt,0)

$\Lambda$

rdt\_rcv(rcvpkt) &&  
( corrupt(rcvpkt) ||  
has\_seq1(rcvpkt) )  
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(ACK, 1, chksum)

udt\_send(sndpkt)

receiver

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

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

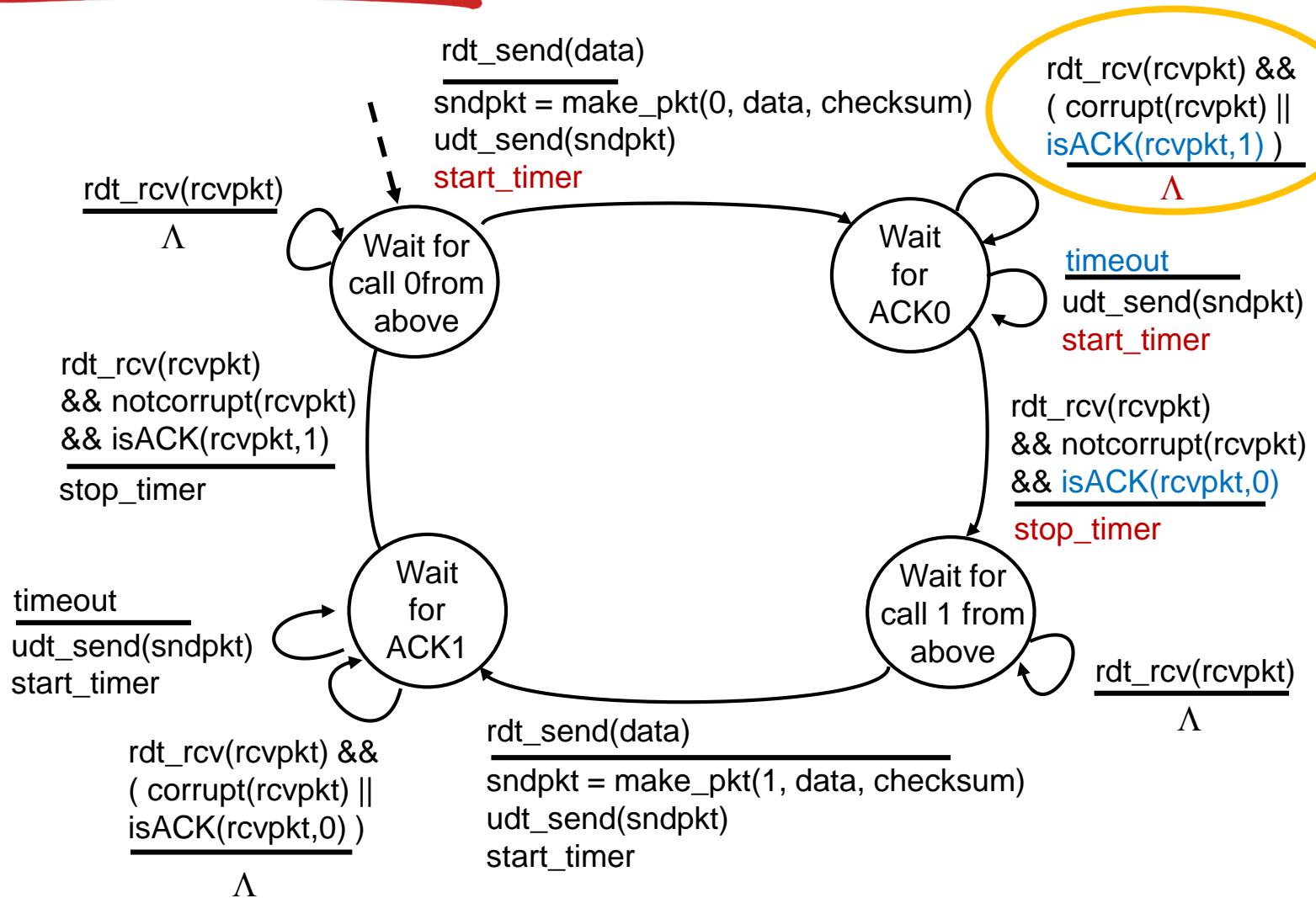
- checksum, seq. #, ACKs, retransmissions will be of help ... but not enough

Approach: sender waits “reasonable” amount of time for ACK

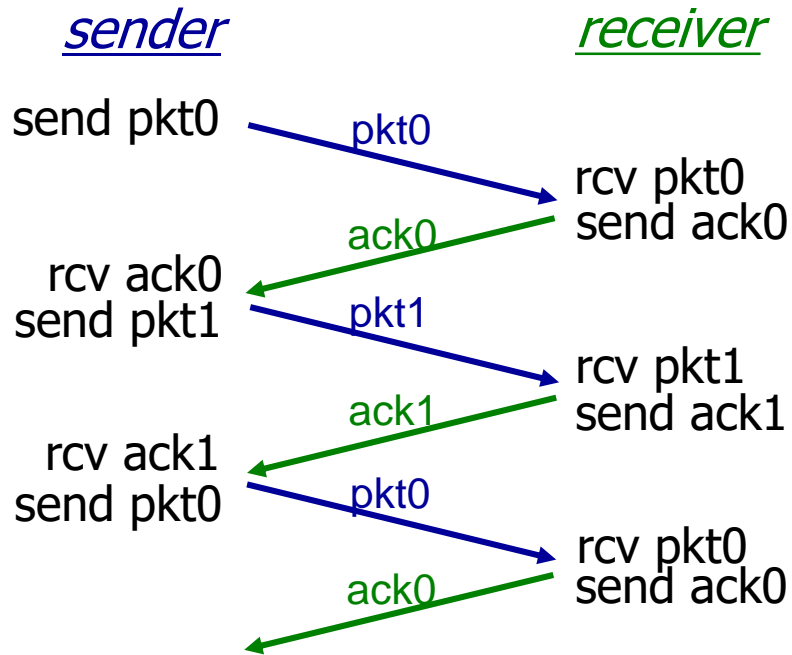
- ❖ retransmits if no ACK received in this time
- ❖ if pkt (or ACK) just delayed (not lost):
  - retransmission will be duplicate, but seq. #'s already handles this
  - receiver specifies seq # of pkt being ACKed
- ❖ requires countdown timer
  - start timer, timer interrupt, stop timer

How long should the sender wait?

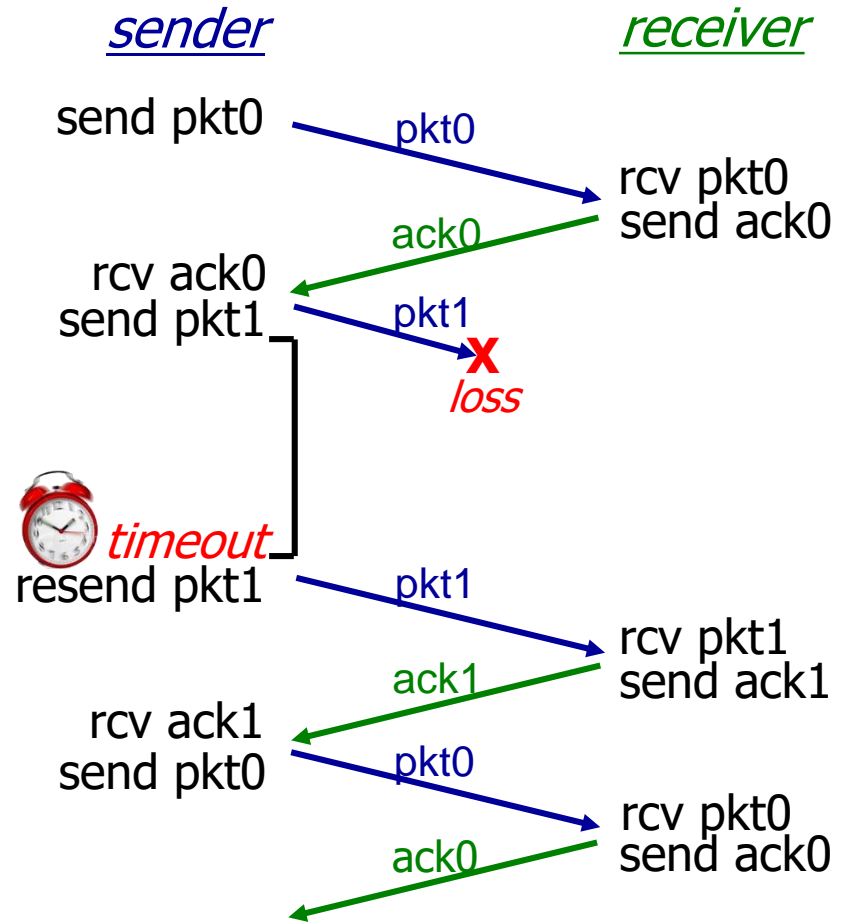
# rdt3.0 sender



# rdt3.0 in action

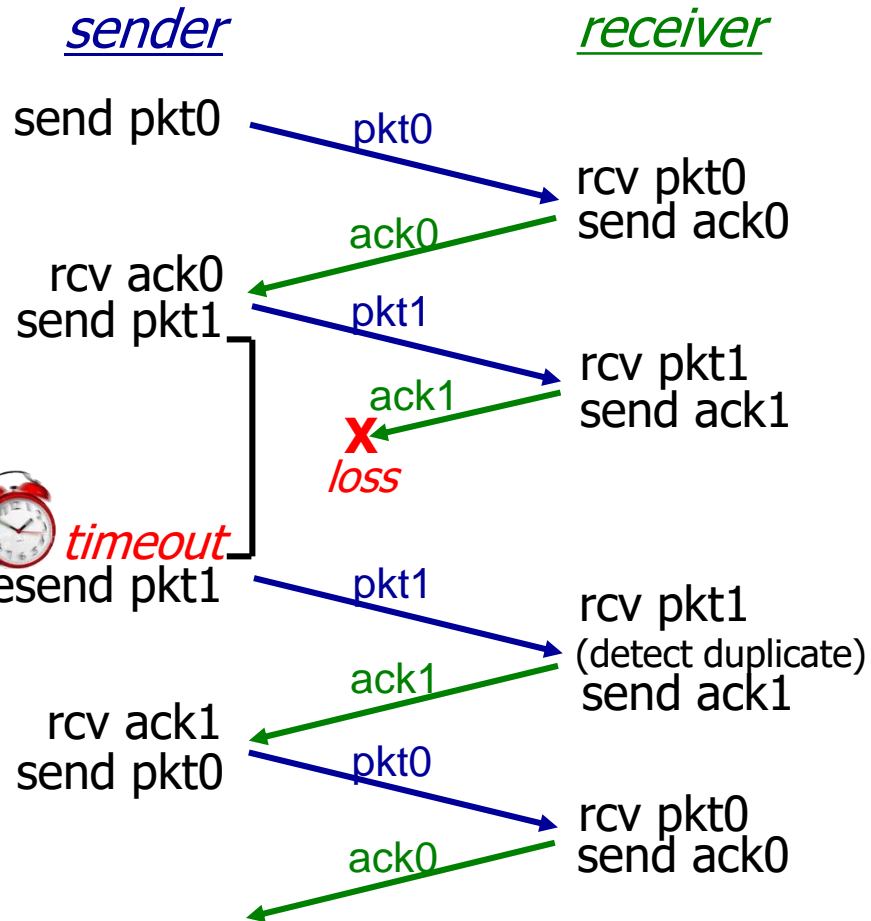


(a) no loss

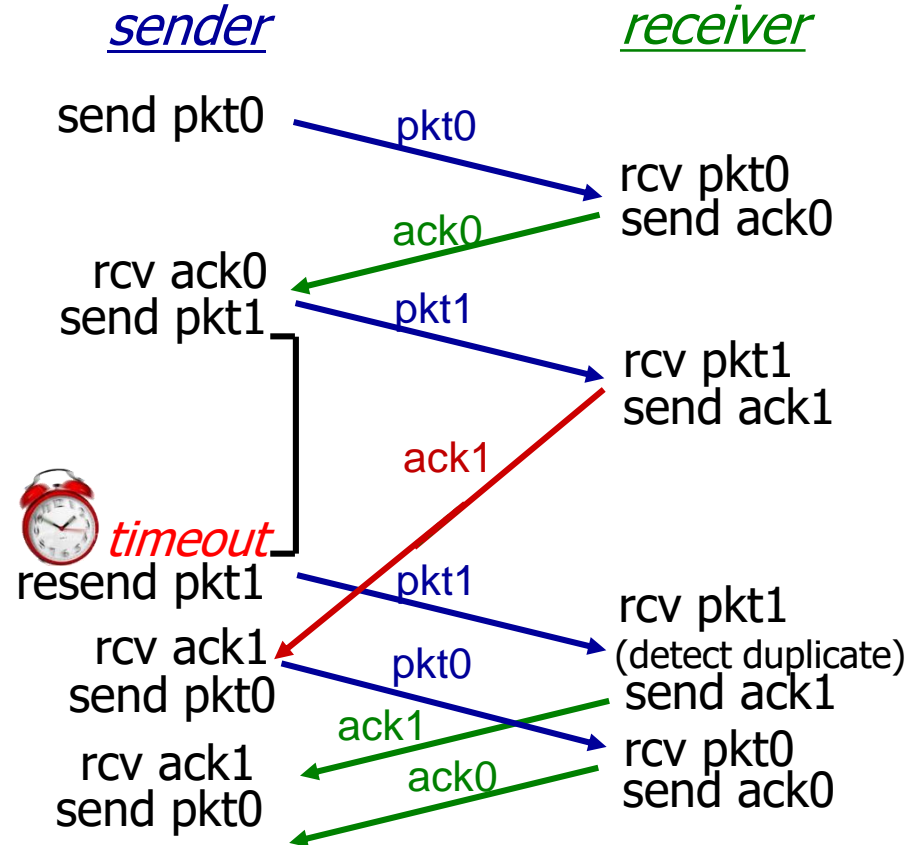


(b) packet loss

# rdt3.0 in action



(c) ACK loss



(d) premature timeout/ delayed ACK



# Summary

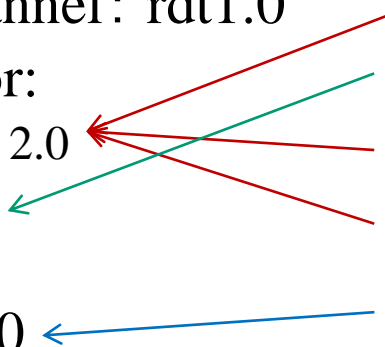
---

## Roadmap:

- ❖ Perfectly reliable channel: rdt1.0
- ❖ Channel with bit error:
  - bit error in packet: rdt 2.0
  - bit error in ACK: 2.1
  - NAK-free: 2.2
- ❖ Lossy channel: rdt 3.0

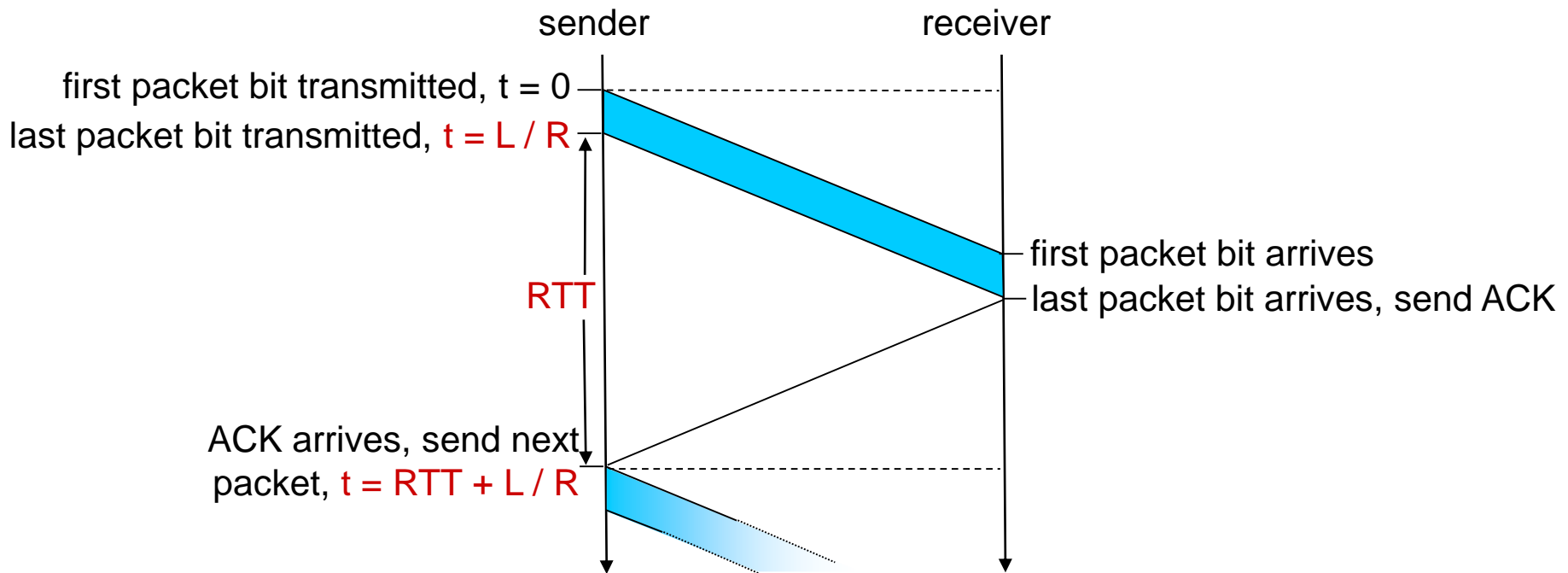
## Summary of Techniques

- Checksum
- Sequence number
- ACK packets
- Retransmission
- Timeout



# Performance of rdt3.0

- ❖ rdt3.0 is correct, but performance is bad
- ❖ e.g.: link rate  $R=1$  Gbps, prop. delay  $T_{pd}=15$  ms, packet length  $L=8000$  bit



- Calculate **utilization**  $U_{\text{sender}}$ : fraction of time sender busy sending

# Performance of rdt3.0

- ❖ link rate  $R=1$  Gbps, prop. delay  $T_{pd}=15$  ms, packet length  $L=8000$  bit

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

- **utilization**  $U_{sender}$ : fraction of time sender busy sending

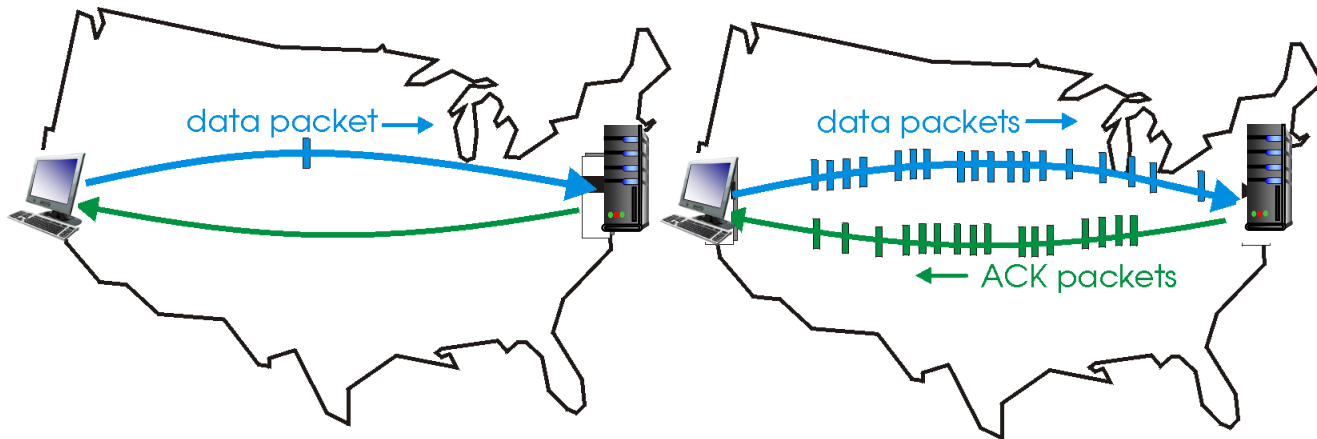
$$U_{sender} = \frac{L/R}{RTT + L/R} = \frac{.008}{30.008} = 0.00027$$

- RTT=30 msec, 1KB pkt every 30 msec:  
33kB/sec throughput over 1 Gbps link
- ❖ network protocol limits use of physical resources!

# 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

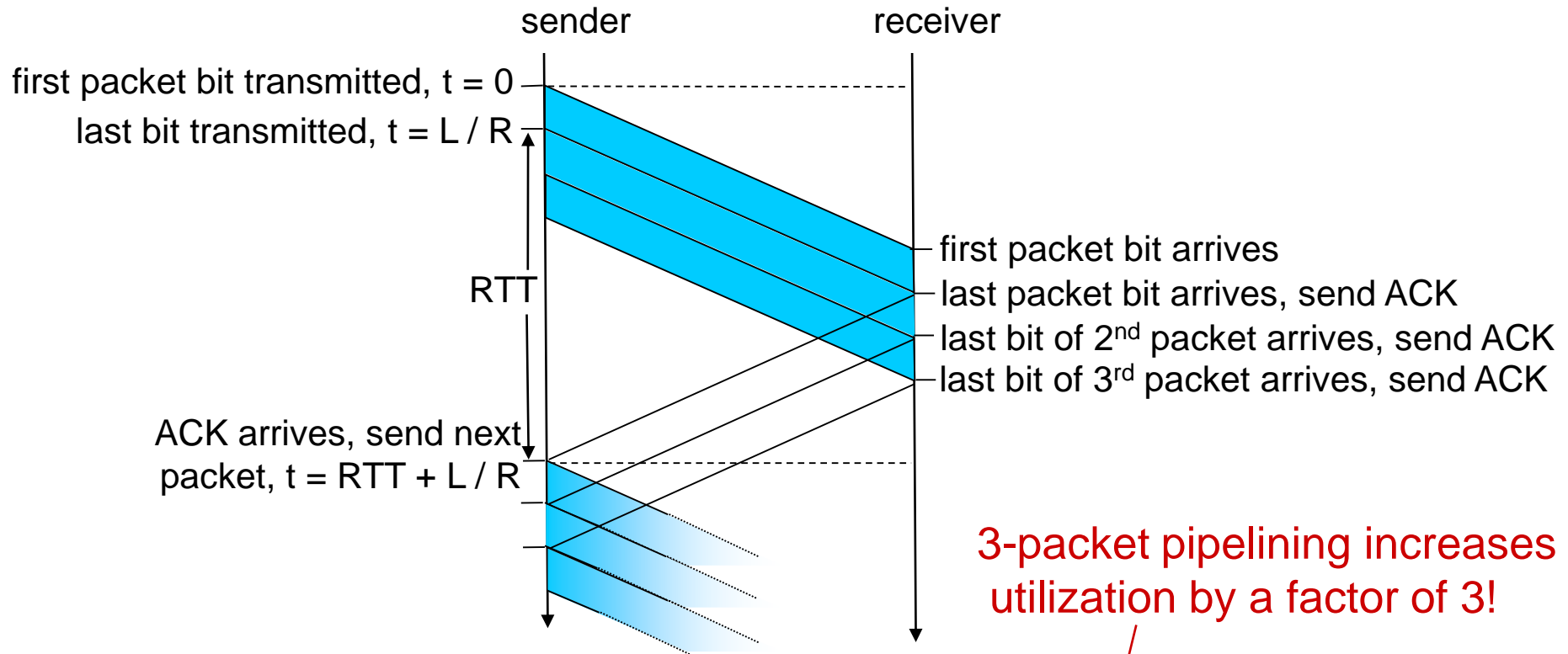


(a) a stop-and-wait protocol in operation

(b) a pipelined protocol in operation

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

# Pipelining: increased utilization



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

$$U_{\text{sender}} = \frac{3L / R}{RTT + L / R} = \frac{0.024}{30.008} = 0.00081$$

# Road Map

## ❖ Go-Back-N

- Timer for the oldest unACKed packet
- Cumulative ACK
- Retransmit all packets in the window

## ❖ Selective repeat

- Timer for each packet in window
- Individual ACK for each correctly received packets
- Retransmit only those packets that might be lost or corrupted

# Go-Back-N overview

sender window (N=4)

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

sender

send pkt0

send pkt1

send pkt2

send pkt3

(wait)

rcv ack0, send pkt4

rcv ack1, send pkt5

ignore duplicate ACK



*pkt 2 timeout*

send pkt2

send pkt3

send pkt4

send pkt5

receiver

No buffer,  
Cumulative ACK

receive pkt0, send ack0

receive pkt1, send ack1

receive pkt3, discard,  
(re)send **ack1**

receive pkt4, discard,  
(re)send **ack1**

receive pkt5, discard,  
(re)send **ack1**

rcv pkt2, deliver, send ack2

rcv pkt3, deliver, send ack3

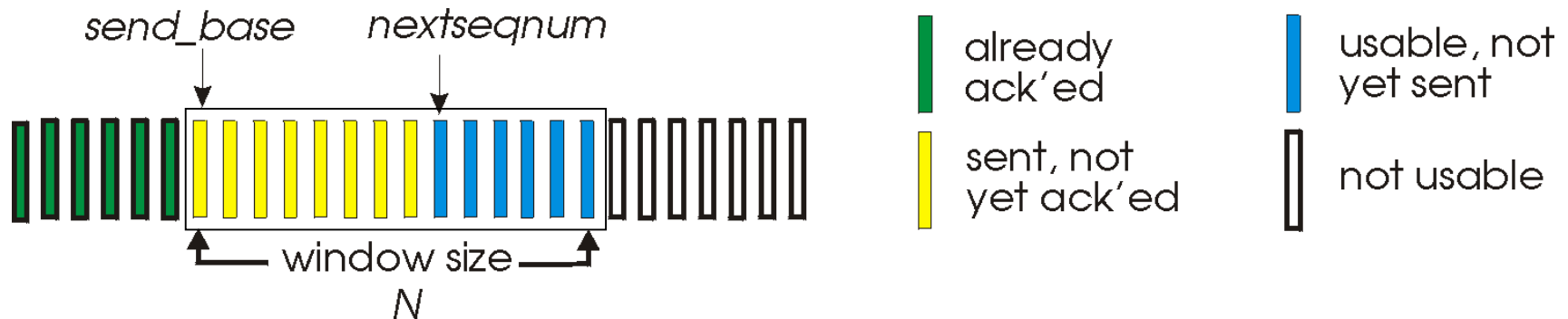
rcv pkt4, deliver, send ack4

rcv pkt5, deliver, send ack5

Retransmit all pkts  
upon pkt loss or

# Go-Back-N: Sequence #, ACK

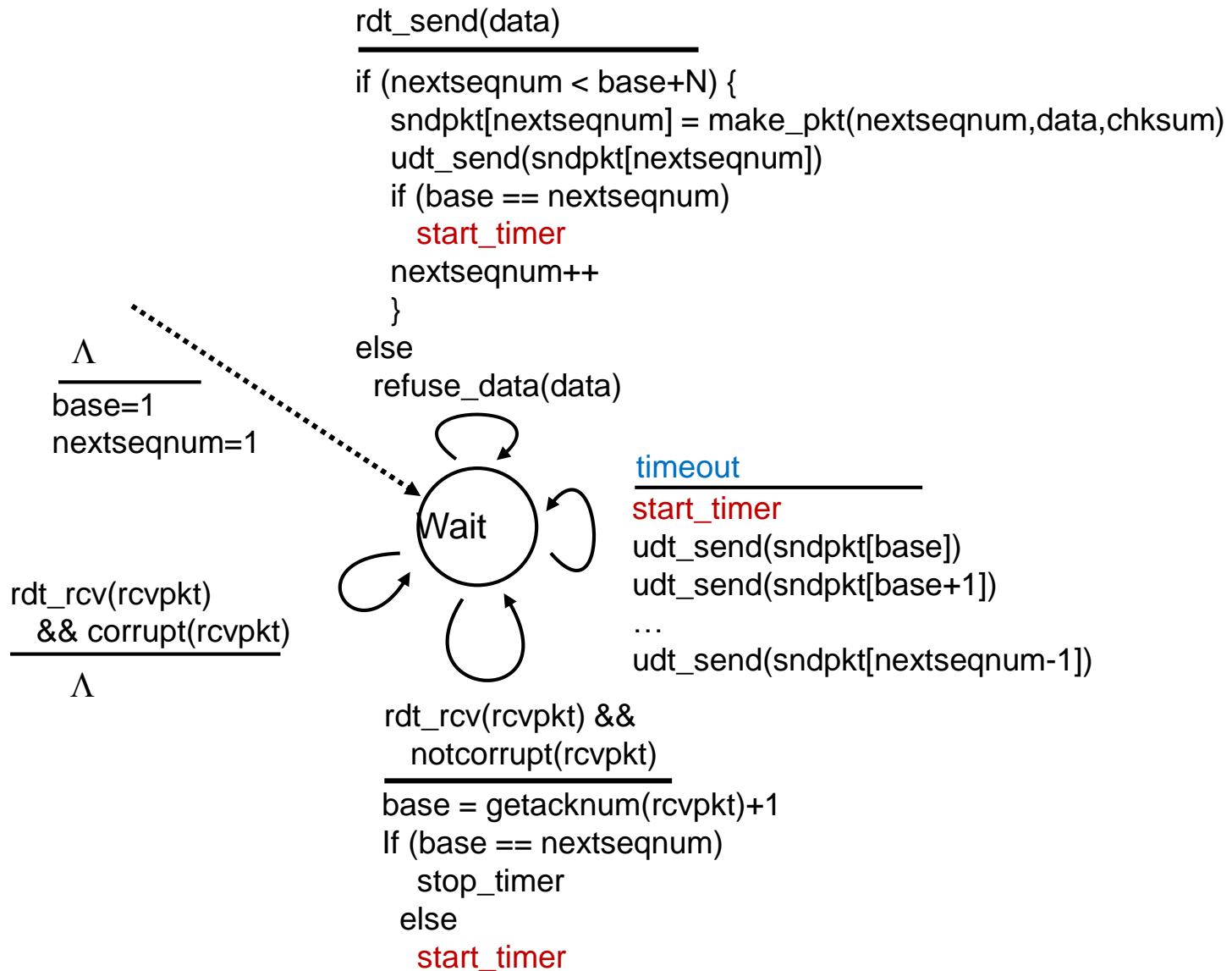
- ❖  $k$ -bit seq # in pkt header (not 0 or 1):  $[0, 2^k - 1]$
- ❖ At most  $N$  pkts in flight: window size =  $N$ , ( $N$  consecutive unacked pkts allowed)



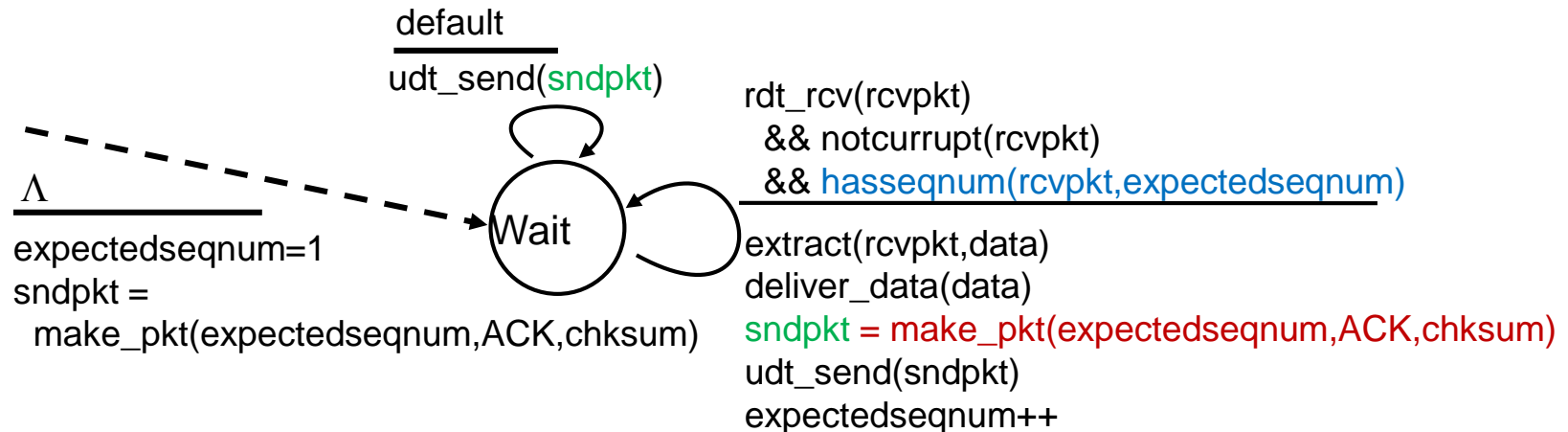
- ❖ ACK( $n$ ) means all pkts with a sequence # up to and including  $n$  are correctly received - *“cumulative ACK”*
  - Sender may receive duplicate ACKs (see receiver)
- ❖ Thought of as a timer for oldest in-flight pkt
- ❖  $timeout(n)$ : retransmit packet  $n$  and all higher seq # pkts in window



# GBN: sender extended FSM



# GBN: receiver extended FSM



ACK-only: always send ACK for correctly-received pkt with highest *in-order* seq #

- may generate duplicate ACKs
- need only remember **expectedseqnum**
- ❖ out-of-order pkt:
  - discard (don't buffer): *no receiver buffering!*
  - re-ACK pkt with highest in-order seq #

# Go-Back-N Recall

sender window (N=4)

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

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

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

sender

send pkt0  
send pkt1  
send pkt2  
send pkt3  
(wait)

rcv ack0, send pkt4  
rcv ack1, send pkt5

ignore duplicate ACK



*pkt 2 timeout*

send pkt2  
send pkt3  
send pkt4  
send pkt5

receiver

No buffer,  
Cumulative ACK

receive pkt0, send ack0  
receive pkt1, send ack1

receive pkt3, discard,  
(re)send ack1

receive pkt4, discard,  
(re)send ack1

receive pkt5, discard,  
(re)send ack1

rcv pkt2, deliver, send ack2  
rcv pkt3, deliver, send ack3  
rcv pkt4, deliver, send ack4  
rcv pkt5, deliver, send ack5

Retransmit all pkts  
upon pkt loss or

# Road Map

- ❖ Go-Back-N
  - Timer for the oldest unACKed packet
  - Cumulative ACK
  - Retransmit all packets in the window
- ❖ Selective repeat
  - Timer for each packet in window
  - Individually ACK correctly received packets
  - Retransmit only those packets that it suspects were lost or corrupted

# Selective repeat

sender window (N=4)

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

Only retransmit the  
unacked pkt (SR)

sender

send pkt0

send pkt1

send pkt2

send pkt3

(wait)

rcv ack0, send pkt4

rcv ack1, send pkt5

record ack3 arrived



*pkt 2 timeout*

send pkt2

record ack4 arrived

record ack5 arrived

receiver have buffer,  
individual ACK

receive pkt0, send ack0

receive pkt1, send ack1

receive pkt3, buffer,  
send ack3

receive pkt4, buffer,  
send ack4

receive pkt5, buffer,  
send ack5

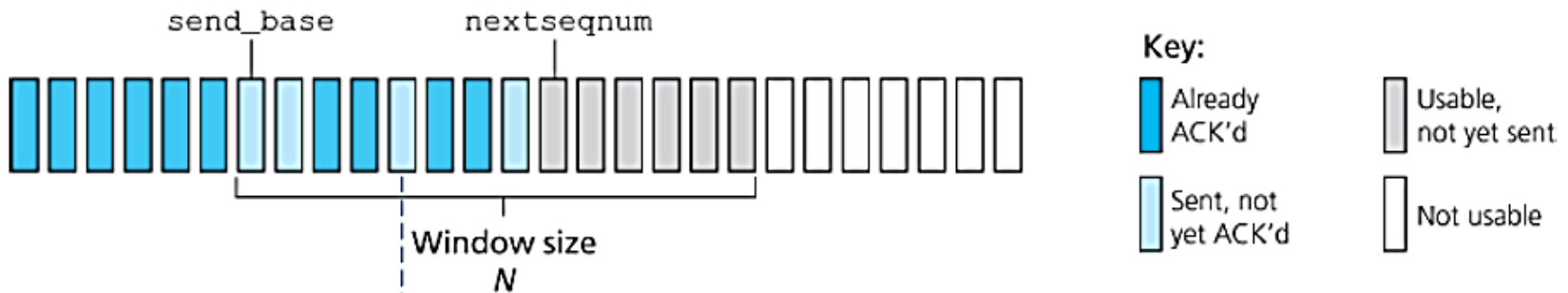
rcv pkt2; deliver pkt2,  
pkt3, pkt4, pkt5; send ack2

*what happens when ack2 arrives?*

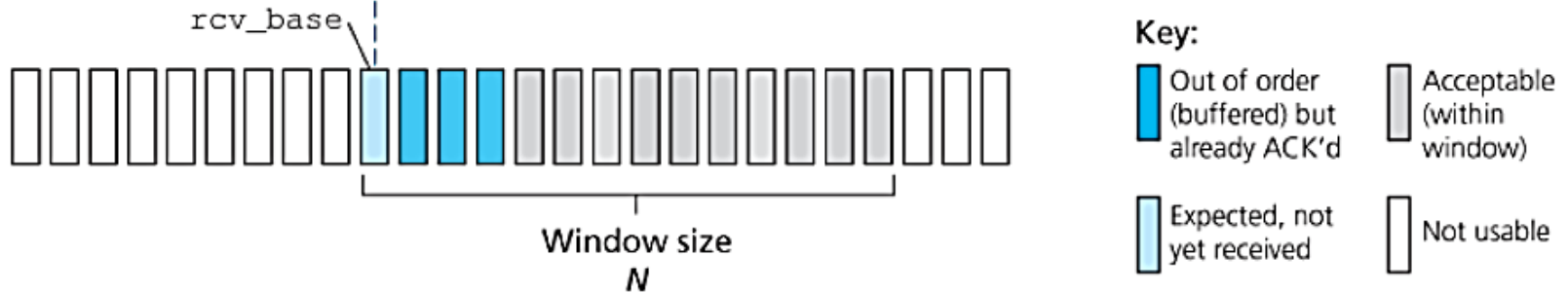
# Selective repeat

- ❖ 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
  - $N$  consecutive seq #'s
  - limits seq #s of sent, unACKed pkts

# Selective repeat: sender, receiver windows



a. Sender view of sequence numbers



b. Receiver view of sequence numbers

# Selective repeat

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: dilemma

Example:

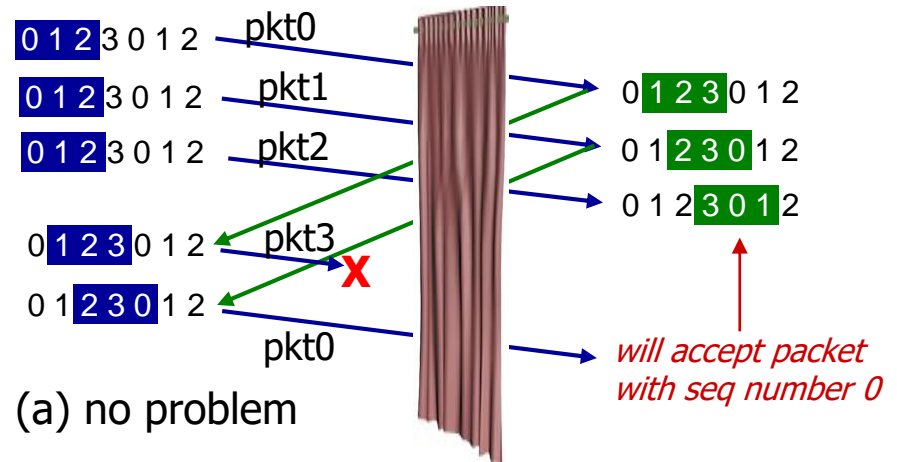
- ❖ seq #'s: 0, 1, 2, 3
- ❖ window size=3
- ❖ receiver sees no difference in two scenarios!
- ❖ duplicate data accepted as new in (b)

**Q:** what relationship between seq # size and window size to avoid problem in (b)?

The window size must be less than or equal to half the size of the sequence number space for SR protocols.

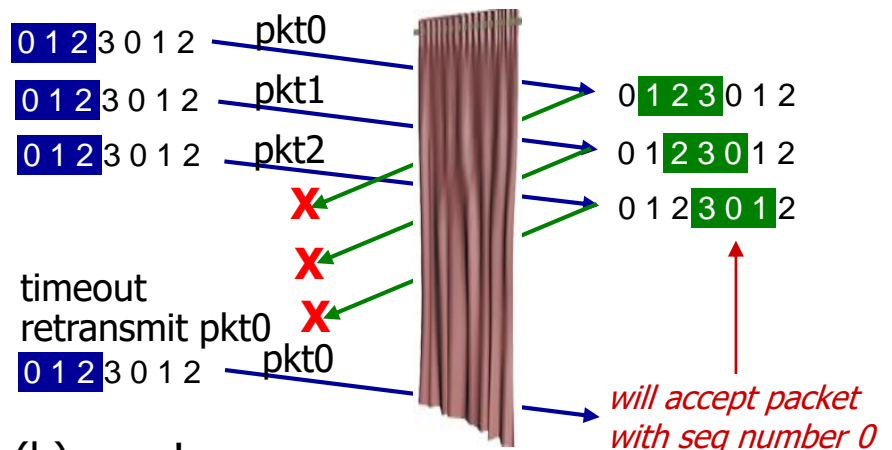
sender window  
(after receipt)

receiver window  
(after receipt)



(a) no problem

*receiver can't see sender side.  
receiver behavior identical in both cases!  
something's (very) wrong!*



(b) oops!

# GBN and SR comparison

## Go-back-N:

- ❖ sender can have up to  $N$  unacked 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

# Reliable Data Transfer Summary

Checksum	Used to detect bit errors in a transmitted packet.
Timer	Used to timeout/retransmit a packet, possibly because the packet (or its ACK) was lost within the channel. Because timeouts can occur when a packet is delayed but not lost (premature timeout), or when a packet has been received by the receiver but the receiver-to-sender ACK has been lost, duplicate copies of a packet may be received by a receiver.
Sequence number	Used for sequential numbering of packets of data flowing from sender to receiver. Gaps in the sequence numbers of received packets allow the receiver to detect a lost packet. Packets with duplicate sequence numbers allow the receiver to detect duplicate copies of a packet.

# Reliable Data Transfer Summary

Acknowledgment	Used by the receiver to tell the sender that a packet or set of packets has been received correctly. Acknowledgments will typically carry the sequence number of the packet or packets being acknowledged. Acknowledgments may be individual or cumulative, depending on the protocol.
Negative acknowledgment	Used by the receiver to tell the sender that a packet has not been received correctly. Negative acknowledgments will typically carry the sequence number of the packet that was not received correctly.
Window, pipelining	The sender may be restricted to sending only packets with sequence numbers that fall within a given range. By allowing multiple packets to be transmitted but not yet acknowledged, sender utilization can be increased over a stop-and-wait mode of operation. We'll see shortly that the window size may be set on the basis of the receiver's ability to receive and buffer messages, or the level of congestion in the network, or both.

# 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, RTT measurement
- 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
- No buffers or variables are allocated to network elements between hosts

## ❖ reliable, in-order byte stream:

- no “message boundaries”
- Seq # and Ack # are in unit of byte, rather than pkt

## ❖ pipelined:

- TCP congestion and flow control set window size

## ❖ full duplex data:

- bi-directional data flow in same connection
- MSS: maximum segment size

## ❖ connection-oriented:

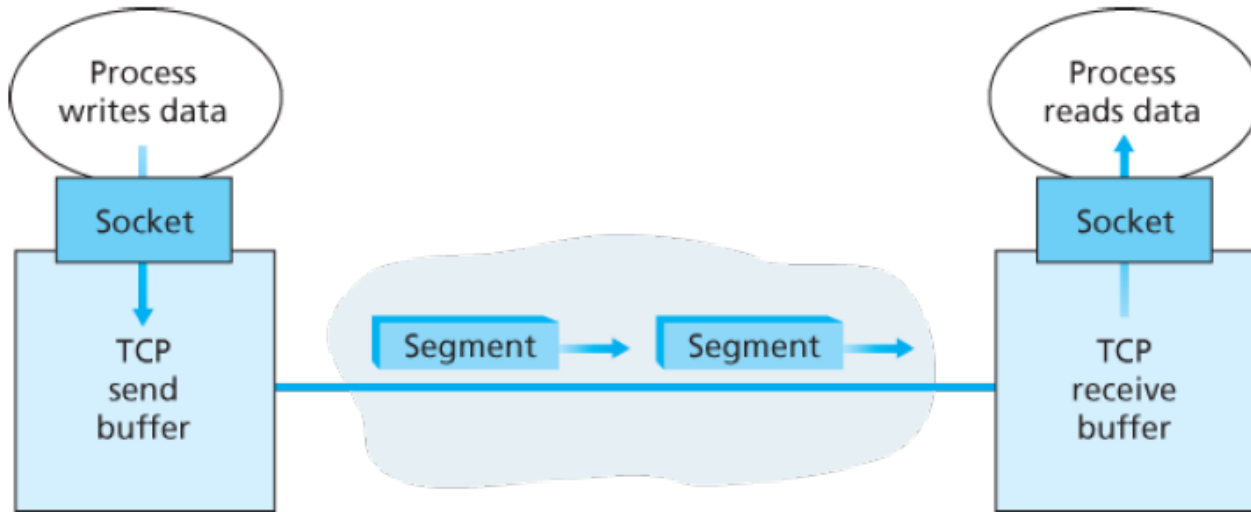
- handshaking (exchange of control msgs) initiates sender and receiver state before data exchange

## ❖ flow controlled:

- sender will not overwhelm receiver

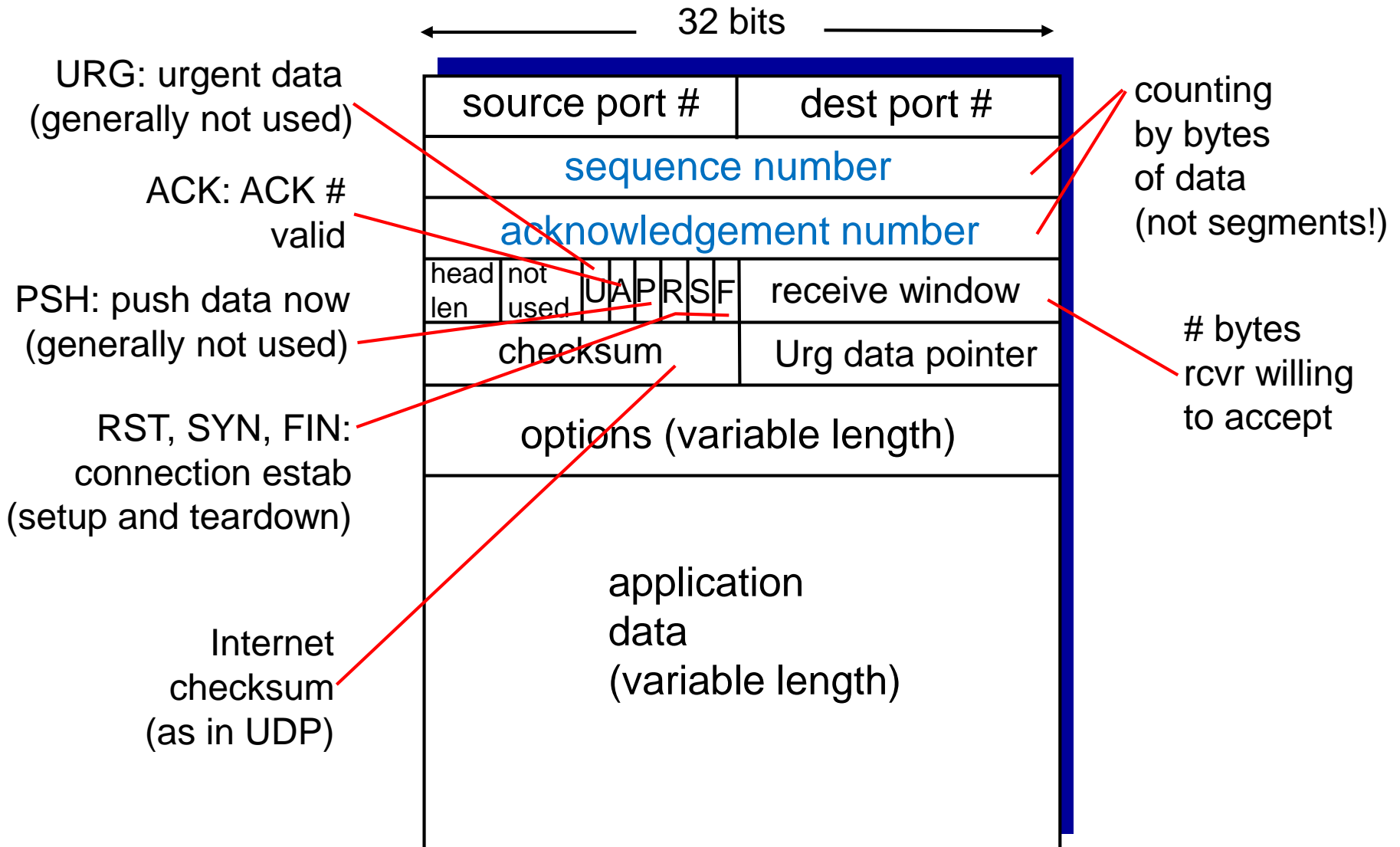
# TCP: Overview

RFCs: 793, 1122, 1323, 2018, 2581



- TCP connection
- TCP grab chunks of data from the sender buffer
  - MSS: maximum segment size, typically 1460 bytes
  - MTU: maximum transmission unit (link-layer frame), typically 1500 bytes
    - Application data + TCP/IP header (typically 40 bytes)
- TCP receives a segment at the other end, place it in receiver buffer
- application reads the stream from the receive buffer

# TCP segment structure





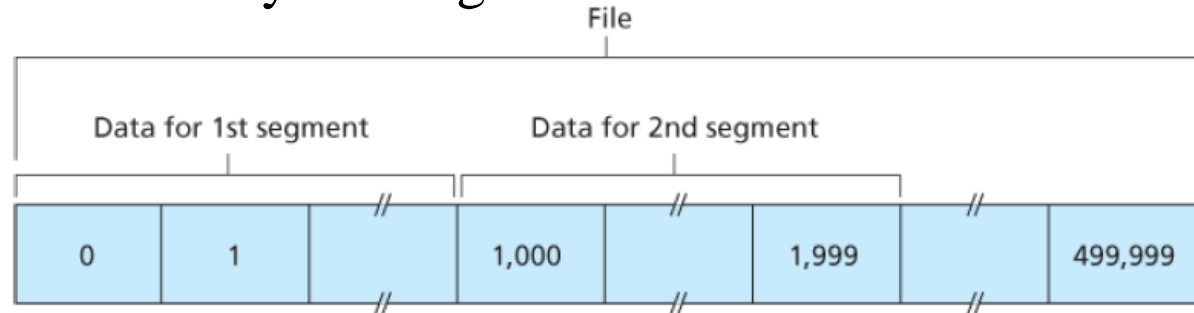
# TCP seq. numbers, ACKs

TCP views data as an unstructured, but ordered, stream of bytes.

- Sequence numbers are over the **stream** of transmitted bytes and **not** over the series of transmitted **segments**

## sequence numbers:

- byte stream “number” of first byte in segment’s data



## acknowledgements:

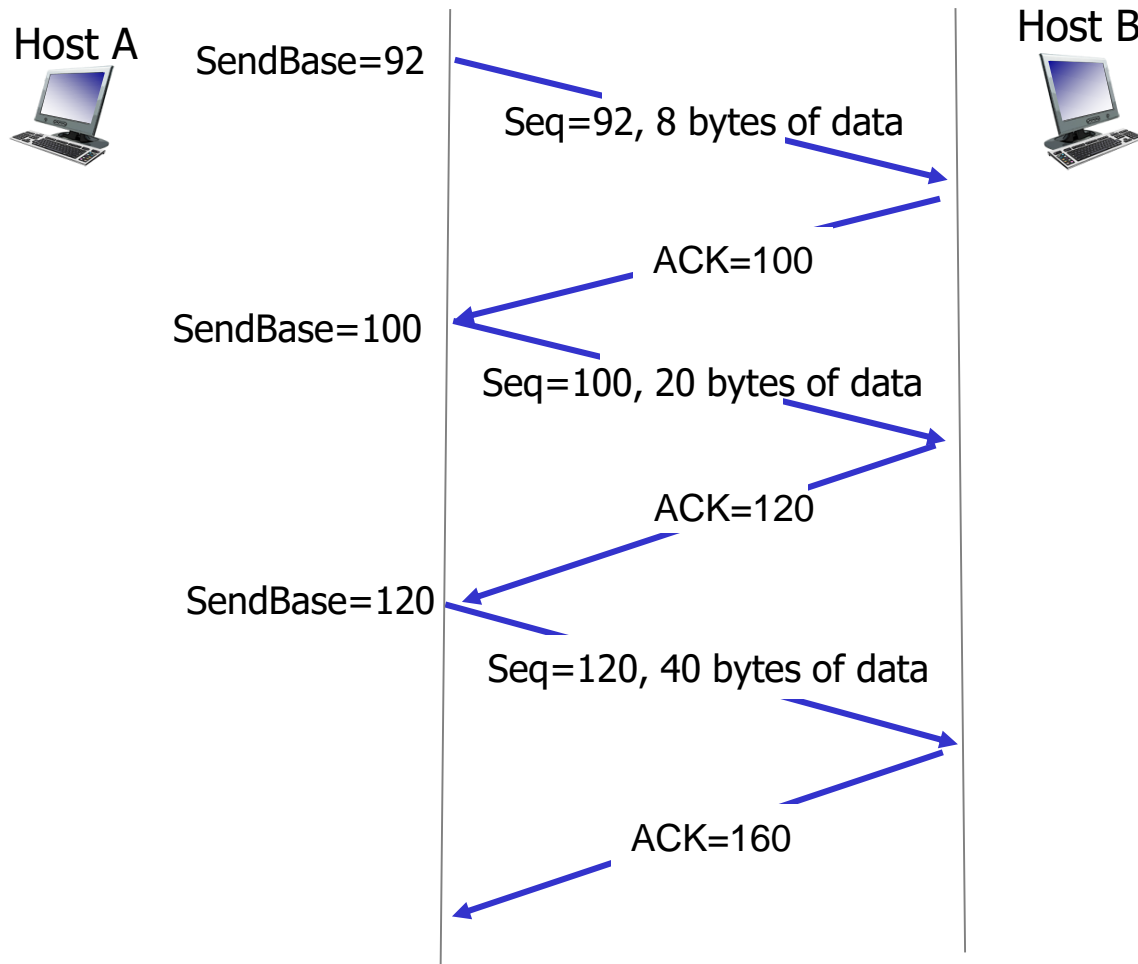
- seq # of next byte expected from other side
  - E.g., receiver has received bytes numbered 0 through 535 and 900 through 1000; then, acknowledgement number is 536.
- cumulative ACK

**Q:** how receiver handles out-of-order segments

- A:** TCP spec doesn't say, - up to implementor

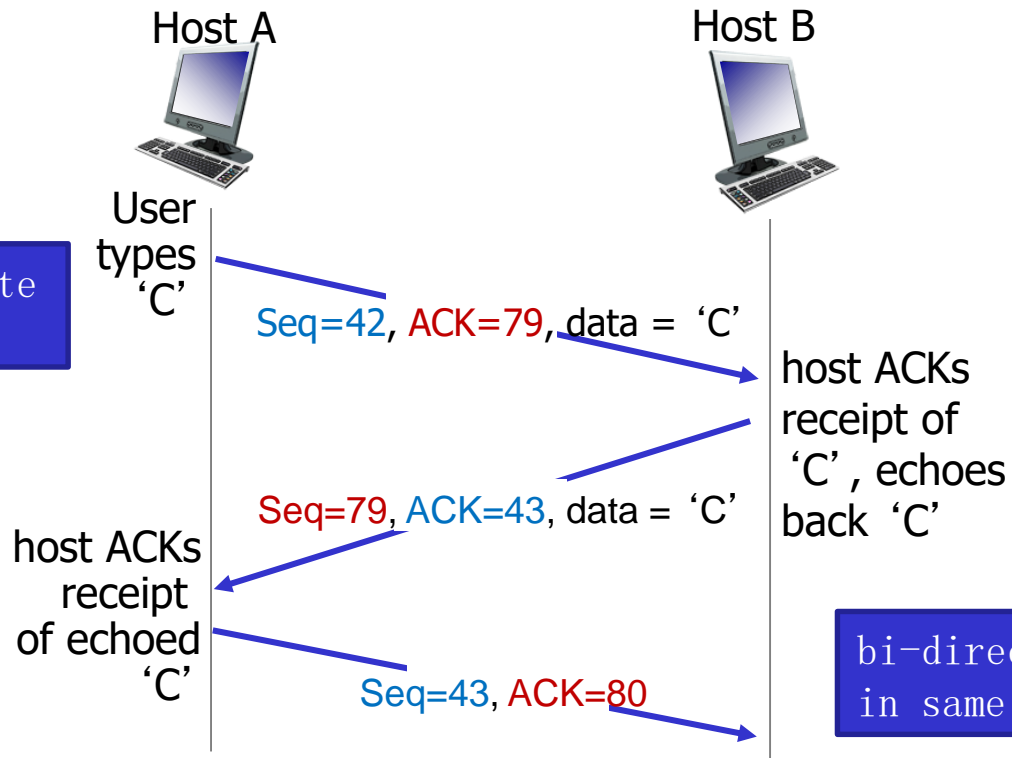
**Initial sequence number  
is randomly chosen**

# TCP Example



# Telnet Case Study

- User types a character at host A, and host A sends the character to host B
- Host B sends back a copy of the character
- Host A displays the character on user's screen



ACK # is the next byte expected from other

host ACKs receipt of 'C', echoes back 'C'

bi-directional data flow in same connection

# 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

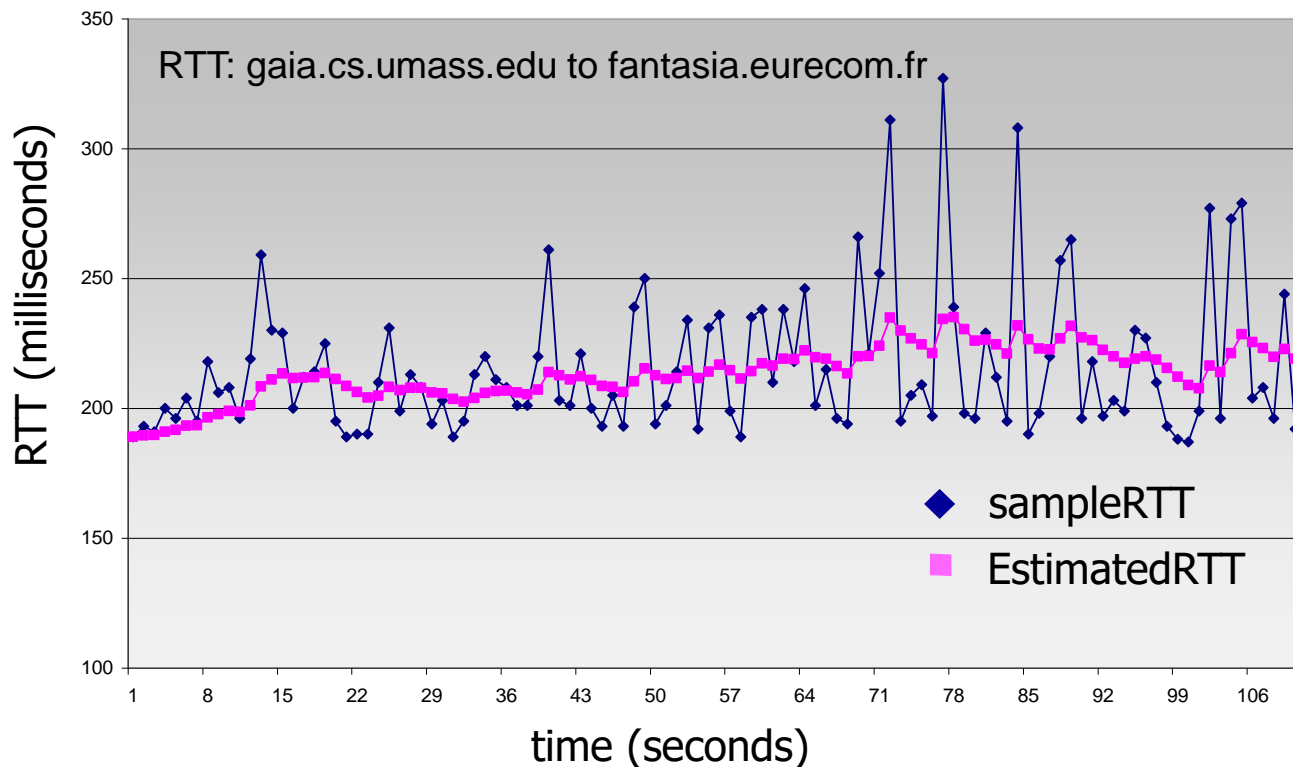
Q: How to estimate RTT?

- ❖ **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**

# TCP round trip time, timeout

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

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



# TCP round trip time, timeout


**Variability** of the RTT: how much **SampleRTT** deviates from **EstimatedRTT**:

$$\text{DevRTT} = (1-\beta) * \text{DevRTT} + \beta * |\text{SampleRTT} - \text{EstimatedRTT}|$$

(typically,  $\beta = 0.25$ )

**TCP timeout interval:** **EstimatedRTT** plus “safety margin”  
large variation in **EstimatedRTT** -> larger safety margin

$$\text{TimeoutInterval} = \text{EstimatedRTT} + 4 * \text{DevRTT}$$

 ↑  
estimated RTT ↑  
“safety margin”

# 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, RTT measurement
- **reliable data transfer**
- flow control
- connection management

3.6 principles of congestion control

3.7 TCP congestion control

# TCP reliable data transfer

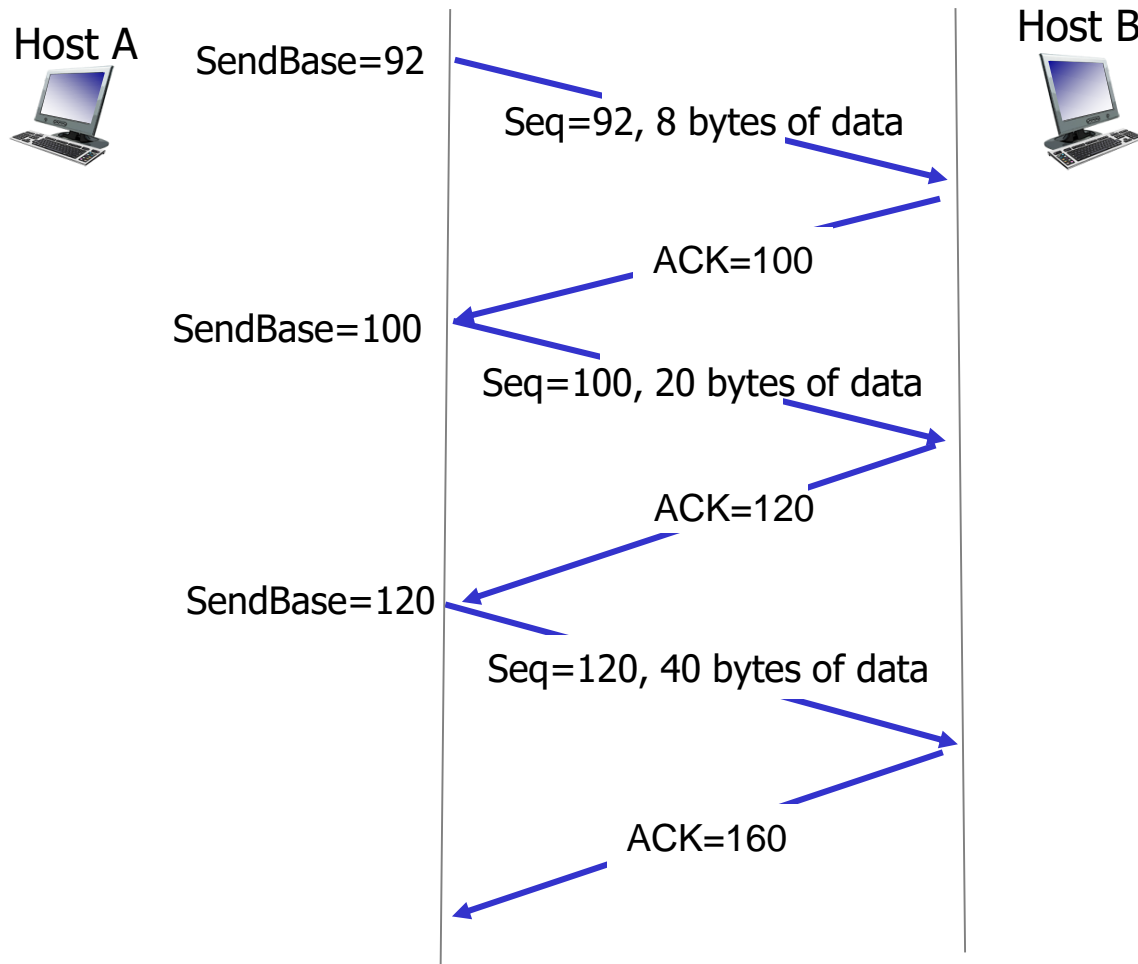
- ❖ TCP creates rdt service on top of IP' s unreliable service
  - pipelined segments: window size, SendBase
  - **cumulative acks**
  - single retransmission timer
- ❖ retransmissions triggered by:
  - timeout events
  - duplicate acks

Let' s initially consider simplified TCP sender:

- ignore duplicate acks
- ignore flow control, congestion control



# TCP without retransmission



# 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 sender events:

```
NextSeqNum=InitialSeqNumber  
SendBase=InitialSeqNumber
```

```
loop (forever) {  
    switch(event)
```

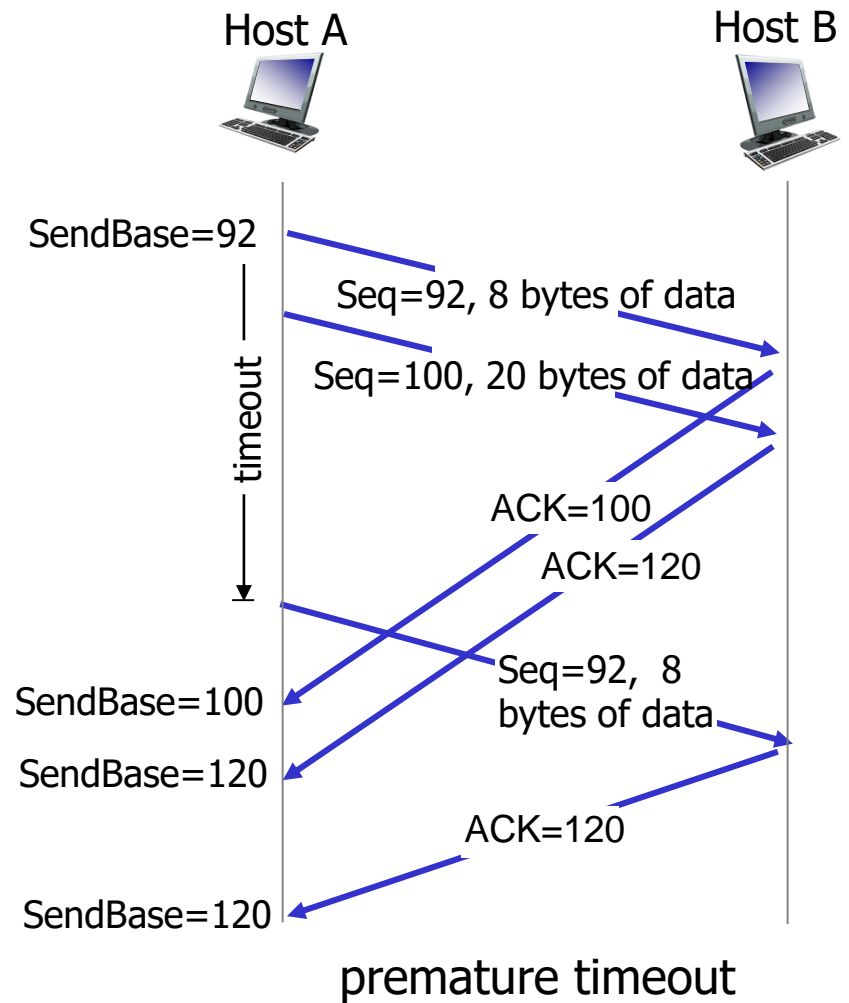
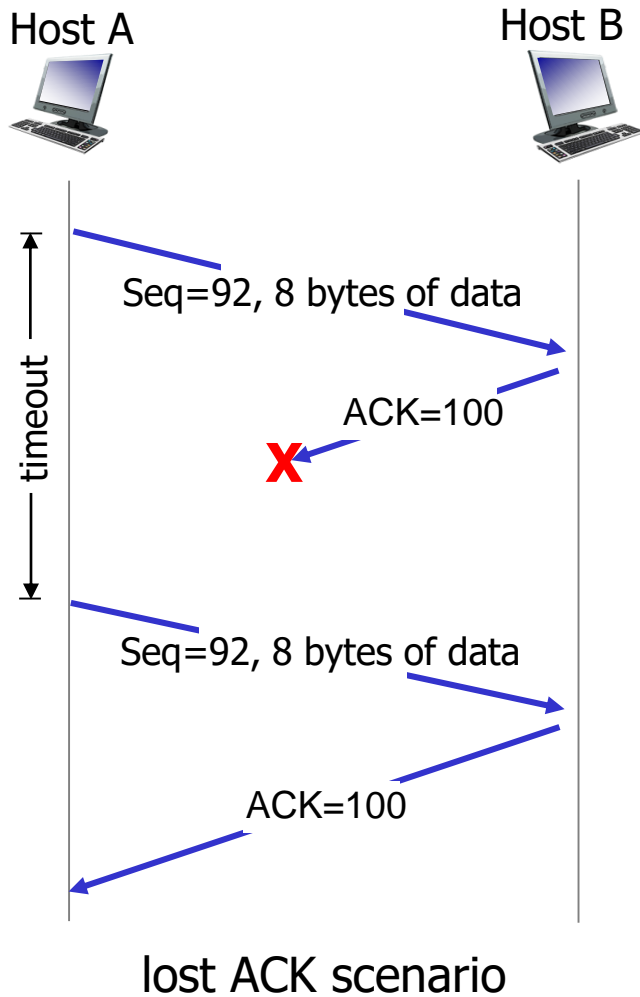
```
event: data received from application above  
    create TCP segment with sequence number NextSeqNum  
    if (timer currently not running)  
        start timer  
    pass segment to IP  
    NextSeqNum=NextSeqNum+length(data)  
    break;
```

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

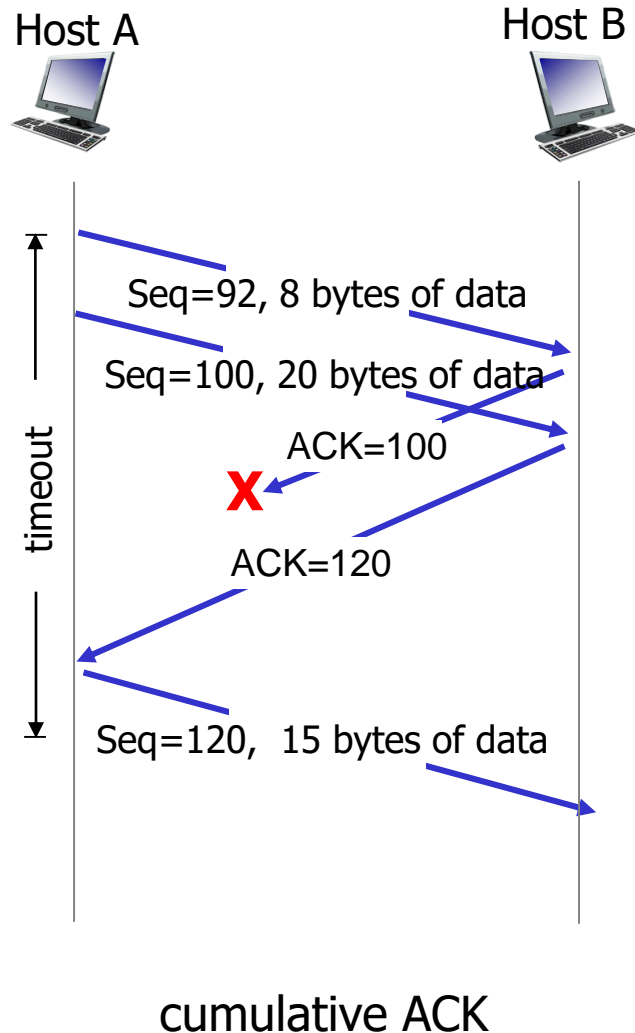
```
event: ACK received, with ACK field value of y  
    if (y > SendBase) {  
        SendBase=y  
        if (there are currently any not-yet-acknowledged segments)  
            start timer  
    }  
    break;
```

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

# TCP: retransmission scenarios



# TCP: retransmission scenarios



# TCP receiver [RFC 1122, RFC 2581]

<i>event at receiver</i>	<i>TCP receiver action</i>
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 <i>duplicate ACK</i> , indicating seq. # of next expected byte
arrival of segment that partially or completely fills gap	immediate send ACK, provided that segment starts at lower end of gap

# TCP fast retransmit

- ❖ time-out period often relatively long:
  - long delay before resending lost packet
- ❖ detect lost segments via duplicate ACKs.
  - sender often sends many segments back-to-back
  - if segment is lost, there will likely be many duplicate ACKs.

## *TCP fast retransmit*

if sender receives **3 ACKs** for same data (“triple duplicate ACKs”), resend unacked segment with smallest seq #

- likely that unacked segment lost, so don't wait for timeout

# TCP fast retransmit

```
NextSeqNum=InitialSeqNumber
SendBase=InitialSeqNumber

loop (forever) {
    switch(event)

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

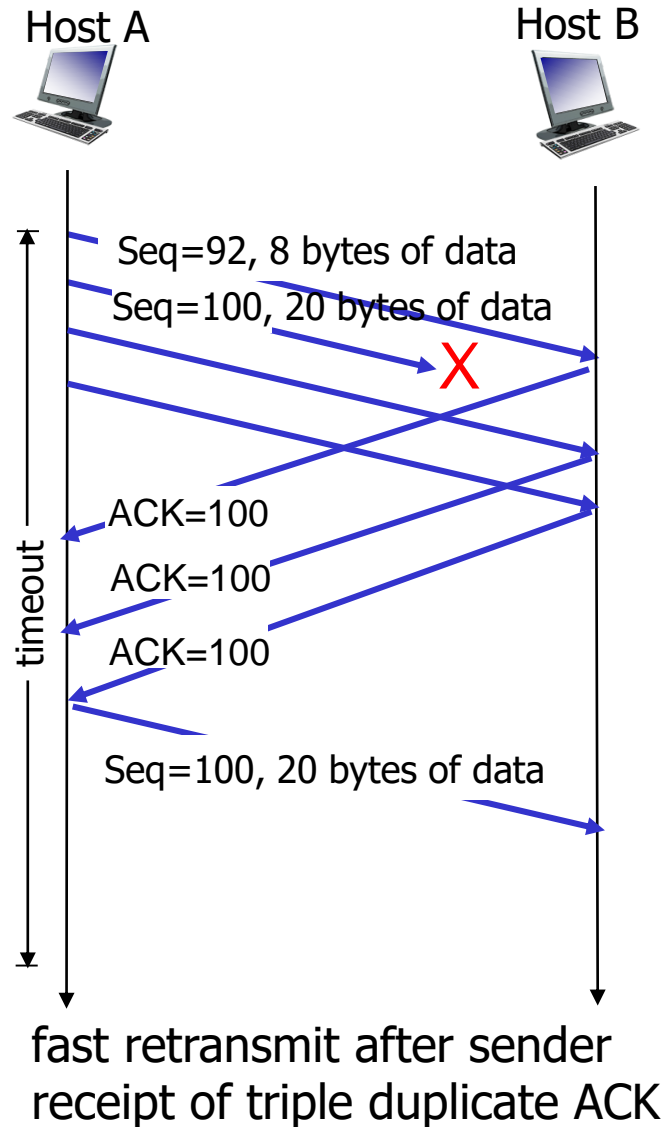
        event: timer timeout
            retransmit not-yet-acknowledged segment ,
                smallest sequence number
            start timer
            break;

        event: ACK received, with ACK field value of y
            if (y > SendBase) {
                SendBase=y
                if (there are currently any not yet
                    acknowledged segments)
                    start timer
            }
            else { /* a duplicate ACK for already ACKed
                segment */
                increment number of duplicate ACKs
                received for y
                if (number of duplicate ACKs received
                    for y==3)
                    /* TCP fast retransmit */
                    resend segment with sequence number y
            }
            break;

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



# TCP fast retransmit



# 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