# CS 455/555
# Intro to Networks and Communications

## The Transport Layer
Multiplexing, UDP, & Reliable Transport

*Dr. Michele Weigle*
Department of Computer Science
Old Dominion University
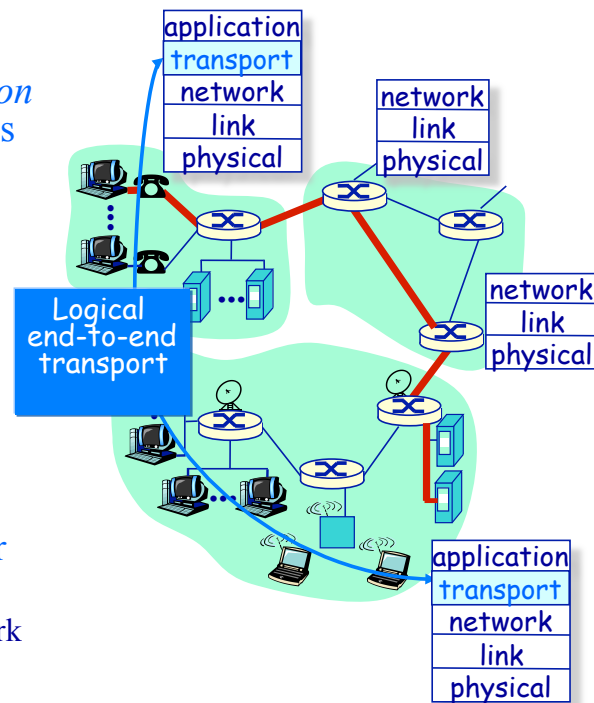*mweigle@cs.odu.edu*

*http://www.cs.odu.edu/~mweigle/CS455-S13*

1

---

# The Transport Layer
## Transport services and protocols

◆ Transport protocols:
  » Provide *logical communication* between application processes running on different hosts
  » Execute on the end systems (and *not* in the network)

◆ Transport *v.* network layer services:
  » *Network layer:* data transfer between end systems
  » *Transport layer:* data transfer between processes
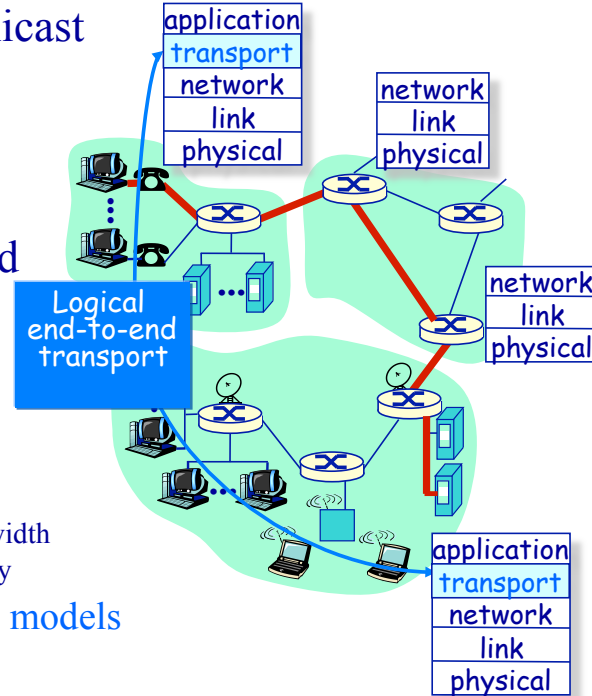    ❖ Relies on, and enhances, network layer services

2

# Transport Layer Protocols
## Internet transport services

- TCP: Reliable, in-order, unicast delivery
  - » Congestion control
  - » Flow control
  - » Connection setup
- UDP: Unreliable, unordered ("best-effort"), unicast or multicast delivery
  - » (Minimal) error detection
- Services not available:
  - » Performance guarantees
    - ❖ No guarantees of available bandwidth
    - ❖ No guarantees of end-to-end delay
  - » Other (non-unicast) delivery models
    - ❖ Multicast (reliable $v.$ unreliable)
    - ❖ Anycast



application
transport
network
link
physical

network
link
physical

Logical
end-to-end
transport

network
link
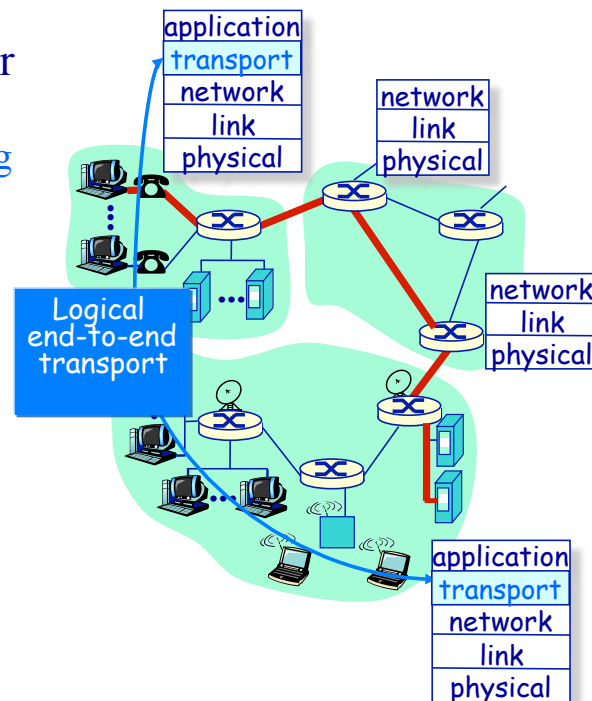physical

application
transport
network
link
physical

3

# Transport Layer Protocols & Services
## Outline

- Fundamental transport layer services
  - » Multiplexing/Demultiplexing
  - » Error detection
  - » Reliable data delivery
  - » Pipelining
  - » Flow control
  - » Congestion control

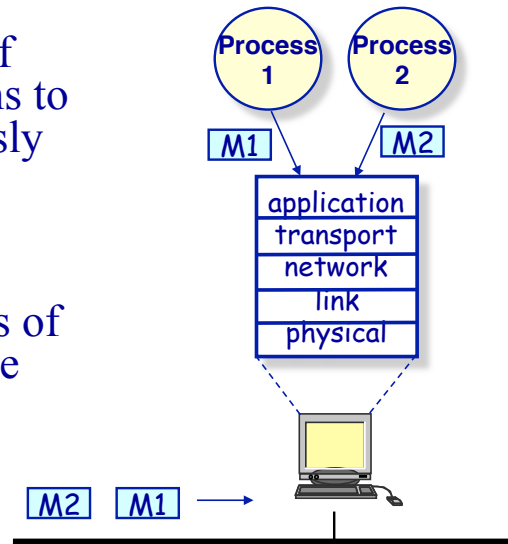- Service implementation in Internet transport protocols
  - » UDP
  - » TCP



application
transport
network
link
physical

network
link
physical

Logical
end-to-end
transport

network
link
physical

application
transport
network
link
physical

4

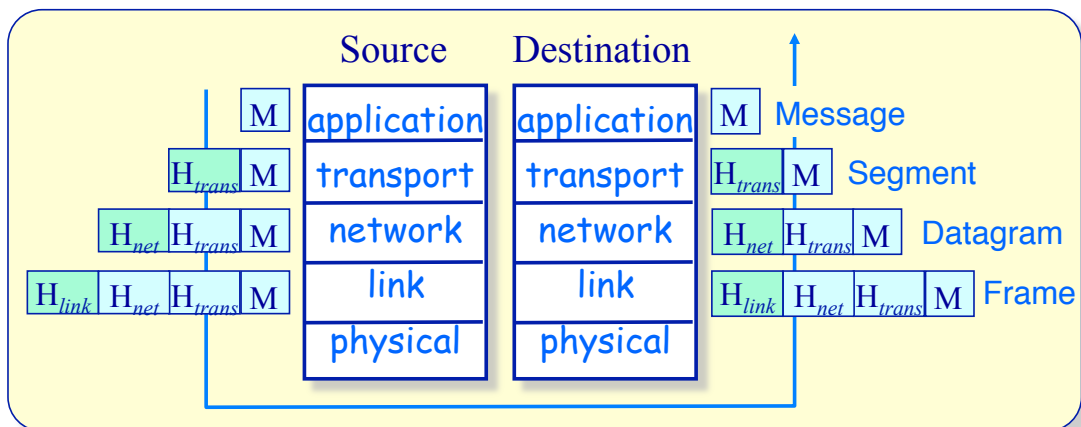# Fundamental Transport Layer Services
## Multiplexing/Demultiplexing

◆ Each end-system has a single protocol "stack"
  » The stack is shared between all applications using the network

◆ Multiplexing is the process of allowing multiple applications to use the network simultaneously
  » (To send data into the network concurrently)

◆ Demultiplexing is the process of delivering received data to the appropriate application

# Multiplexing/Demultiplexing
## Review: Protocol layering in the Internet

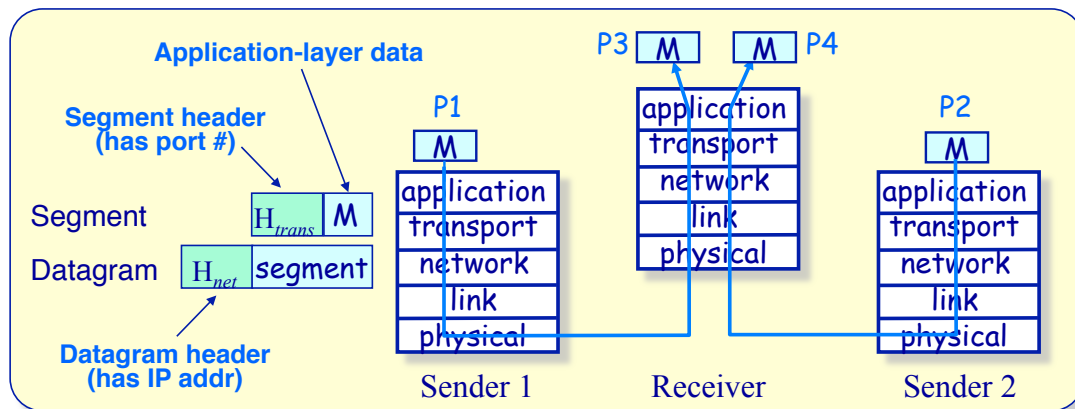◆ At the sender, each layer takes data from above
  » May subdivide into multiple data units at sending layer
  » Adds header information to create new data unit
  » Passes new data unit to layer below
◆ The process is reversed at the receiver

# Multiplexing/Demultiplexing
## Demultiplexing



- ◆ Demultiplexing is the process of delivering received segments to the correct application-layer process
  - » IP address (in network-layer datagram header) identifies the receiving machine
  - » Port number (in transport-layer segment header) identifies the receiving process

# Multiplexing/Demultiplexing
## Transport protocol specific demultiplexing

- ◆ Demultiplexing actions depend on whether the transport layer is connectionless (UDP) or connection-oriented (TCP)

- ◆ UDP demultiplexes segments to the *socket*
  - » UDP uses 2-tuple
    *<destination IP address, destination port number>*
    to identify the socket
  - » Socket is "owned" by some process (allocated by OS).

- ◆ TCP demultiplexes segments to the *connection*
  - » TCP uses 4-tuple
    *<source IP addr, source port nbr, destination IP addr, destination port nbr>*
    to the identify connection
  - » Connection (and its socket) is owned by some process

# Multiplexing/Demultiplexing
## Examples

| dest. IP: B |
|---|
| dest. port: 53 |
|  |

Host A → DNS Server B

| dest. IP: A |
|---|
| dest. port: x |
|  |

**DNS server port use**

Web client Host C

| source IP: C |
|---|
| dest IP: B |
| source port: y |
| dest. port: 80 |
|  |

| source IP: C |
|---|
| dest IP: B |
| source port: x |
| dest. port: 80 |
|  |

Web client Host A

| source IP: A |
|---|
| dest IP: B |
| source port: x |
| dest. port: 80 |
|  |

→ Web Server B

**Web server port use**

# Internet Transport Protocols
## User Datagram Protocol (UDP) [RFC 768]

- ◆ No frills, "bare bones" Internet transport protocol
- ◆ Best effort service — UDP segments may be:
  - » Lost
  - » Delivered out of order to the application
  - » Delivered multiple times to the application
- ◆ "Connectionless"
  - » No handshaking between UDP sender, receiver
  - » Each UDP segment handled independently of others
- ◆ Error Detection
  - » Based on checksum
  - » Make sure received packets haven't been corrupted

←——— 32 bits ———→

| source port # | dest. port # |
|---|---|
| length | checksum |
| application data (message payload) | |

**UDP segment format**

Length field is length in bytes, of UDP segment (including header)

# User Datagram Protocol (UDP)
## Is unreliable, unordered communications useful?

- ◆ Who uses UDP?
  - » Often used for streaming multimedia applications
  - » Loss tolerant
  - » Rate sensitive

- ◆ Other UDP uses (why?):
  - » DNS
  - » SNMP
  - » Routing protocols

Why use UDP?
- ◆ No connection establishment (which can add delay)
- ◆ Simple: no connection state at sender, receiver
- ◆ Small segment header
- ◆ No congestion control: UDP can blast away as fast as desired

- ◆ Reliable transfer over UDP still possible
  - » Reliability can always be added at the application layer
  - » (Application-specific error recovery)

11

# Fundamental Transport Layer Services
## Principles of reliable data transfer



- ◆ Goal: Provide a reliable channel abstraction
  - » The characteristics of the underlying channel will determine the complexity of providing reliable communications
- ◆ Issues: State required at sender and receiver and number of control messages exchanged

12

# Reliable Data Transfer
## Programming interfaces

Called "from above" by the application. Application passes in data to be delivered to receiving application

Called by **rdt** to deliver data to application

Application Layer

data      data      `deliver_data()`

`rdt_send()`

Transport Layer

Reliable Data Transfer Protocol (Sending Side)

Reliable Data Transfer Protocol (Receiving Side)

packet      packet      `rdt_rcv()`

Network Layer      `udt_send()`

Unreliable Channel

Called by **rdt** to transfer packet over unreliable channel to receiver

Called when packet arrives on receive-side of channel

13

# Reliable Data Transfer
## Protocol specification method

◆ Use finite state machines to specify sender and receiver algorithms

» When in a given state, the next state (and actions) are uniquely determined by the next event

event causing state transition
actions taken on state transition

State 1

event
actions

State 2

14

# Reliable Data Transfer Protocol 1.0

## Reliable transfer over a reliable channel

◆ The underlying channel is assumed to be perfectly reliable
  » No bit errors
  » No loss of packets

◆ Sender state machine                    ◆ Receiver state machine
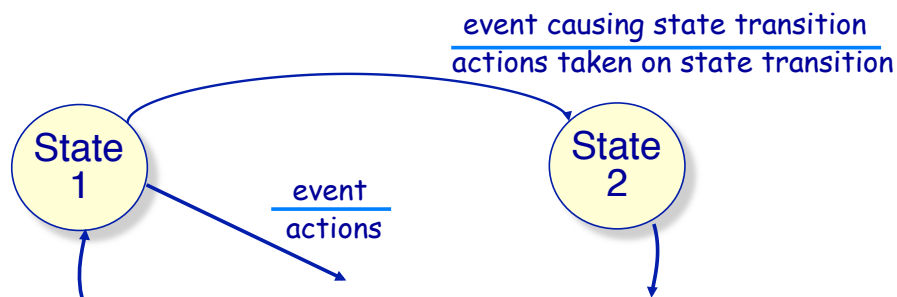
wait for call from above

```
rdt_send(data)
make_pkt(pkt,data)
udt_send(pkt)
```

wait for call from below

```
rdt_rcv(pkt)
extract(pkt,data)
deliver_data(data)
```

# Reliable Data Transfer Protocol 1.0

## Programming interfaces

Application Layer

data   data   `deliver_data()`

Transport Layer

`rdt_send()`

```
rdt_rcv(pkt)
extract(pkt,data)
deliver_data(data)
```

```
rdt_send(data)
make_pkt(pkt,data)
udt_send(pkt)
```

wait for call from above

wait for call from below

packet   packet   `rdt_rcv()`

Network Layer

`udt_send()`

Reliable Channel

◆ This is the complete protocol under the assumption of a reliable network channel

# Reliable Data Transfer Protocol 2.0
## Reliable transfer over a channel with bit errors

◆ Now assume the underlying channel may "flip" random bits in a packet
◆ How to detect errors?
◆ How to recover from errors:
  » *acknowledgements* (*ACKs*) — the receiver explicitly tells the sender that a packet was received OK
  » *negative acknowledgements* (*NAKs*) — the receiver explicitly tells the sender that a packet had errors
  » Sender retransmits packet on receipt of NAK
◆ New mechanisms to deal with bit errors:
  » Error detection
  » Control messages (ACK, NAK) from a receiver to the sender
  » Retransmission

# Reliable Data Transfer Protocol 2.0
## Reliable transfer over a channel with bit errors only

**Receiver FSM**

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

**Sender FSM**

rdt_send(data)
compute chksum
make_pkt(sndpkt,data,chksum)
udt_send(sndpkt)

wait for
call from
above

wait for
ACK or
NAK

wait for
call from
below

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

rdt_rcv(rcvpkt) &&
isACK(rcvpkt)

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

# Reliable Data Transfer Protocol 2.0
## Example 1: No Errors Occur

**Receiver FSM**

rdt_rcv(rcvpkt) &&
corrupt(rcvpkt)

udt_send(NAK)

**Sender FSM**

rdt_send(data)
compute chksum
make_pkt(sndpkt,data,chksum)
udt_send(sndpkt)

wait for
call from
above

wait for
ACK or
NAK

rdt_rcv(rcvpkt) &&
isNAK(rcvpkt)

udt_send(sndpkt)

wait for
call from
below

rdt_rcv(rcvpkt) &&
isACK(rcvpkt)

rdt_rcv(rcvpkt) &&
notcorrupt(rcvpkt)

extract(rcvpkt,data)
deliver_data(data)
udt_send(ACK)

19

# Reliable Data Transfer Protocol 2.0
## Example 2: A corrupted packet arrives at the receiver

**Receiver FSM**

rdt_rcv(rcvpkt) &&
corrupt(rcvpkt)

udt_send(NAK)

**Sender FSM**

rdt_send(data)
compute chksum
make_pkt(sndpkt,data,chksum)
udt_send(sndpkt)

wait for
call from
above

wait for
ACK or
NAK

rdt_rcv(rcvpkt) &&
isNAK(rcvpkt)

udt_send(sndpkt)

wait for
call from
below

rdt_rcv(rcvpkt) &&
isACK(rcvpkt)

rdt_rcv(rcvpkt) &&
notcorrupt(rcvpkt)

extract(rcvpkt,data)
deliver_data(data)
udt_send(ACK)

20

# Reliable Data Transfer Protocol 2.0
## Simple… but wrong!

wait for call from above

wait for ACK or NAK

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

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

wait for call from below

rdt_rcv(rcvpkt) &&
notcorrupt(rcvpkt)
_____
udt_send(ACK)

◆ What happens if an ACK/NAK is corrupted?
  » Sender doesn't know what happened at the receiver!

◆ What to do?
  » Sender ACKs/NAKs the receiver's ACK/NAK?
  » Retransmit last data packet?

# Reliable Data Transfer Protocol 2.0
## Simple… but wrong!

wait for call from above

wait for ACK or NAK

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

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

wait for call from below

rdt_rcv(rcvpkt) &&
notcorrupt(rcvpkt)
_____
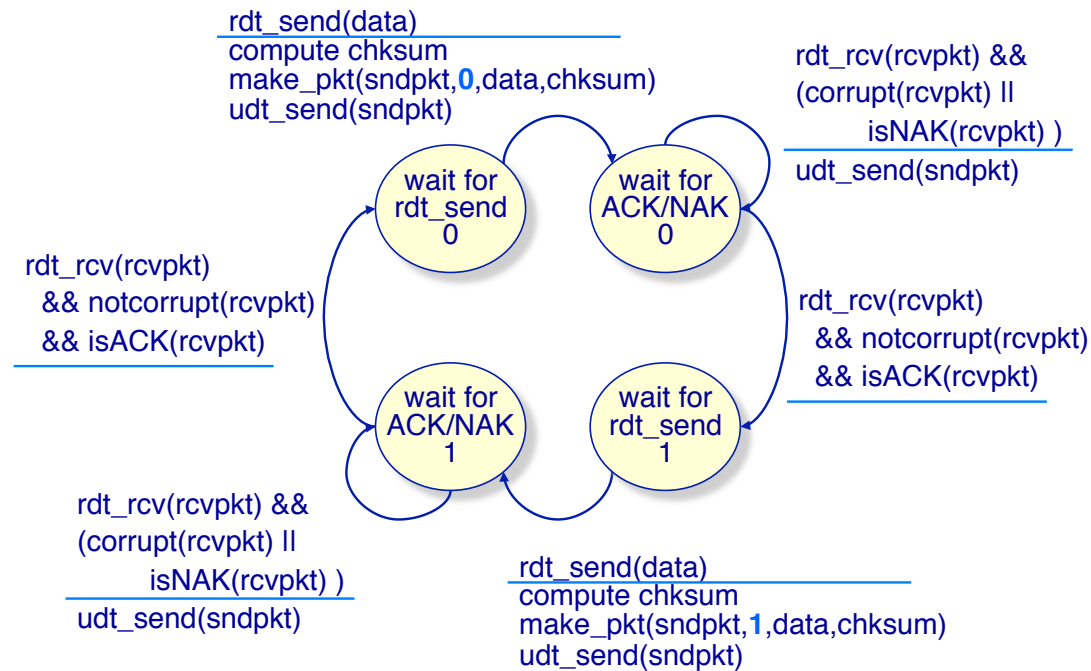udt_send(ACK)

◆ Deal with corrupted ACKs/NAKs by retransmission of data packets

◆ Sender will add a sequence number to each packet to allow the receiver to detect duplicate packets
  » Receiver's transport layer discards duplicate packets

◆ How much space to reserve in a header field for sequence numbers?
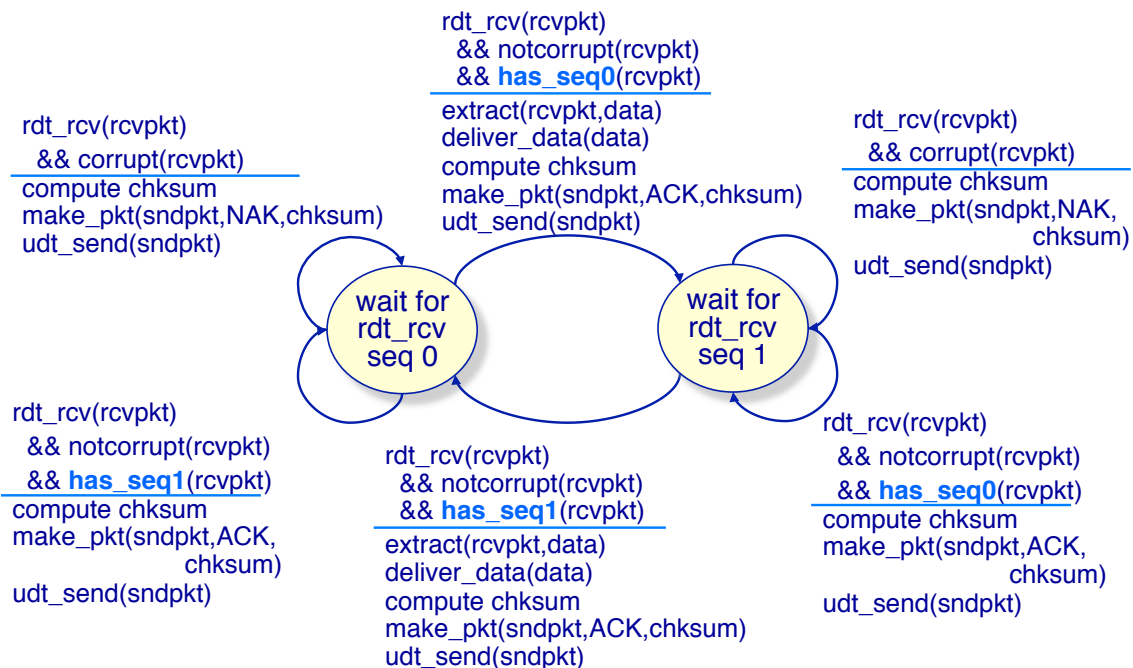
# Reliable Data Transfer Protocol 2.1
## Sender state machine to handle garbled ACKs/NAKs

rdt_send(data)
compute chksum
make_pkt(sndpkt,**0**,data,chksum)
udt_send(sndpkt)

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

( wait for rdt_send 0 )

( wait for ACK/NAK 0 )

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& isACK(rcvpkt)

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& isACK(rcvpkt)

( wait for ACK/NAK 1 )

( wait for rdt_send 1 )

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

rdt_send(data)
compute chksum
make_pkt(sndpkt,**1**,data,chksum)
udt_send(sndpkt)

# Reliable Data Transfer Protocol 2.1
## Receiver state machine to handle garbled ACKs/NAKs

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

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

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

( wait for rdt_rcv seq 0 )

( wait for rdt_rcv seq 1 )

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& **has_seq1**(rcvpkt)
compute chksum
make_pkt(sndpkt,ACK,
      chksum)
udt_send(sndpkt)

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

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& **has_seq0**(rcvpkt)
compute chksum
make_pkt(sndpkt,ACK,
      chksum)
udt_send(sndpkt)

# Reliable Data Transfer Protocol 2.1
## Discussion (Handling garbled ACKs/NAKs)

◆ Sender issues

Sequence number added to header
» Two sequence numbers suffice

Must check if received ACK/ NAK is corrupted

Number of states doubles
» State encodes whether current packet has sequence number 0 or 1

◆ Receiver issues

Must check if received packet is duplicate
» State encodes whether expected packet sequence number is 0 or 1

Note: receiver can *not* know if its last ACK/NAK received OK at sender
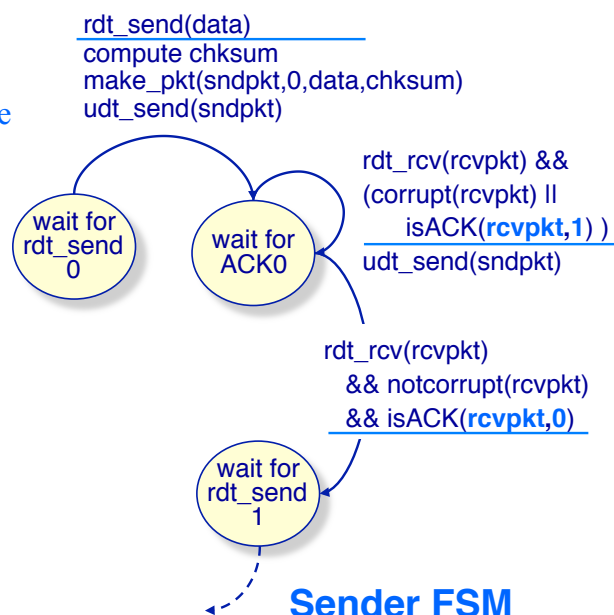
# Reliable Data Transfer Protocol 2.2
## A NAK-free protocol

◆ Instead of NAKing, receiver sends ACK for last packet received OK
» Receiver must include the sequence number of packet being ACKed in ACK

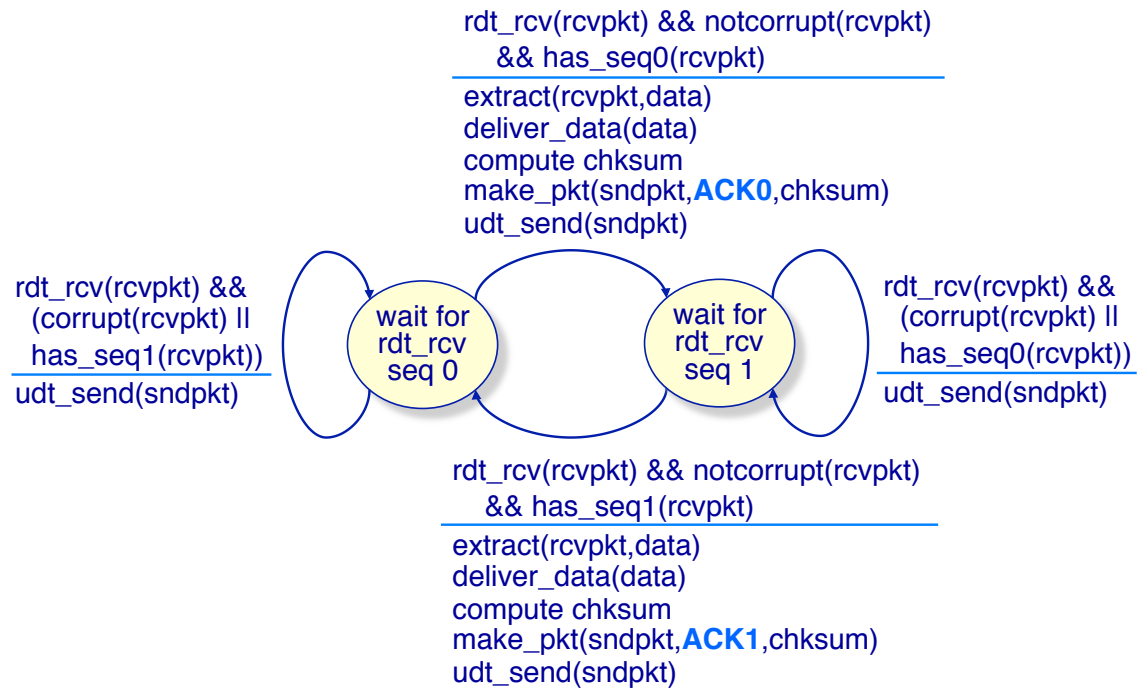◆ Receipt of duplicate ACKs at sender is equivalent to a NAK
» Sender retransmits current packet

rdt_send(data)
compute chksum
make_pkt(sndpkt,0,data,chksum)
udt_send(sndpkt)

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

wait for rdt_send 0

wait for ACK0

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

wait for rdt_send 1

**Sender FSM**

# Reliable Data Transfer Protocol 2.2

**Receiver state machine to eliminate NAKs**

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
&& has_seq0(rcvpkt)

extract(rcvpkt,data)
deliver_data(data)
compute chksum
make_pkt(sndpkt,**ACK0**,chksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) &&
(corrupt(rcvpkt) ||
has_seq1(rcvpkt))

udt_send(sndpkt)

**wait for rdt_rcv seq 0**

**wait for rdt_rcv seq 1**

rdt_rcv(rcvpkt) &&
(corrupt(rcvpkt) ||
has_seq0(rcvpkt))

udt_send(sndpkt)

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
&& has_seq1(rcvpkt)

extract(rcvpkt,data)
deliver_data(data)
compute chksum
make_pkt(sndpkt,**ACK1**,chksum)
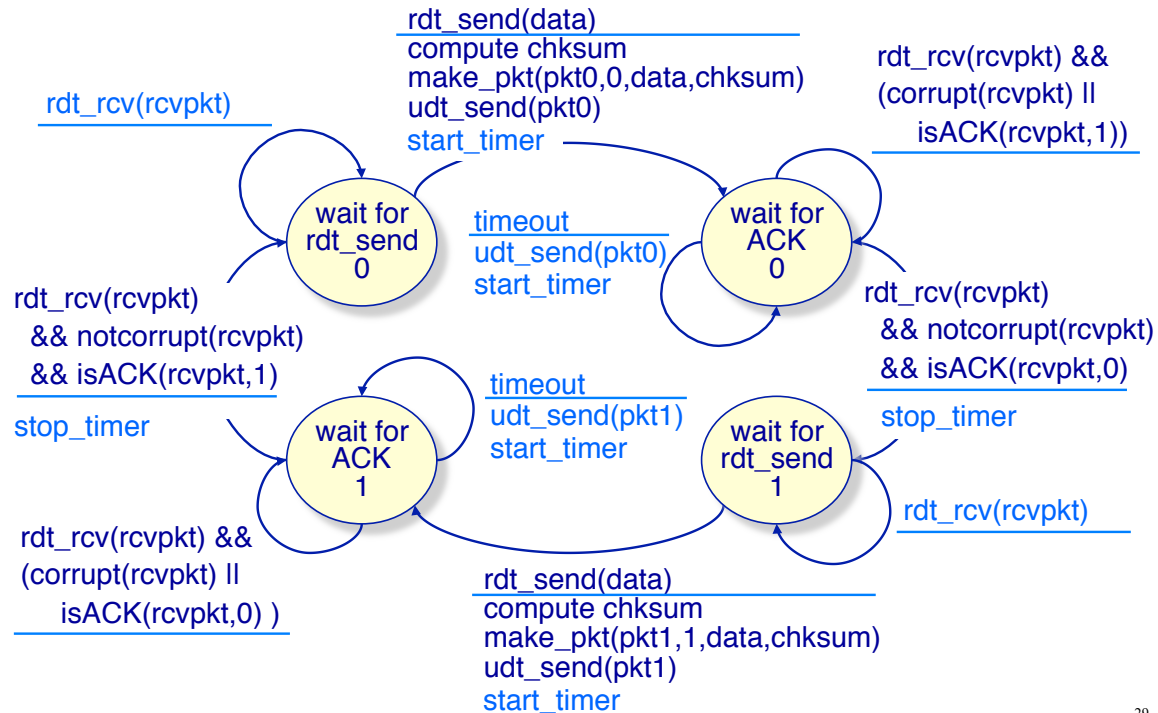udt_send(sndpkt)

# Reliable Data Transfer Protocol 3.0

**Dealing with channels with errors *and* loss**

- ◆ Now assume the underlying channel can also lose packets

- ◆ New problem: How to detect loss?
  - » Are checksums, ACKs, sequence numbers, retransmissions enough?

- ◆ Approach: sender waits "reasonable" amount of time and retransmits if no ACK received in this time
  - » Requires the use of a countdown timer

- ◆ What if packet (or ACK) just delayed beyond its timer?
  - » Retransmission will be  duplicate…
  - » But use of sequence numbers already handles this!

# Reliable Data Transfer Protocol 3.0
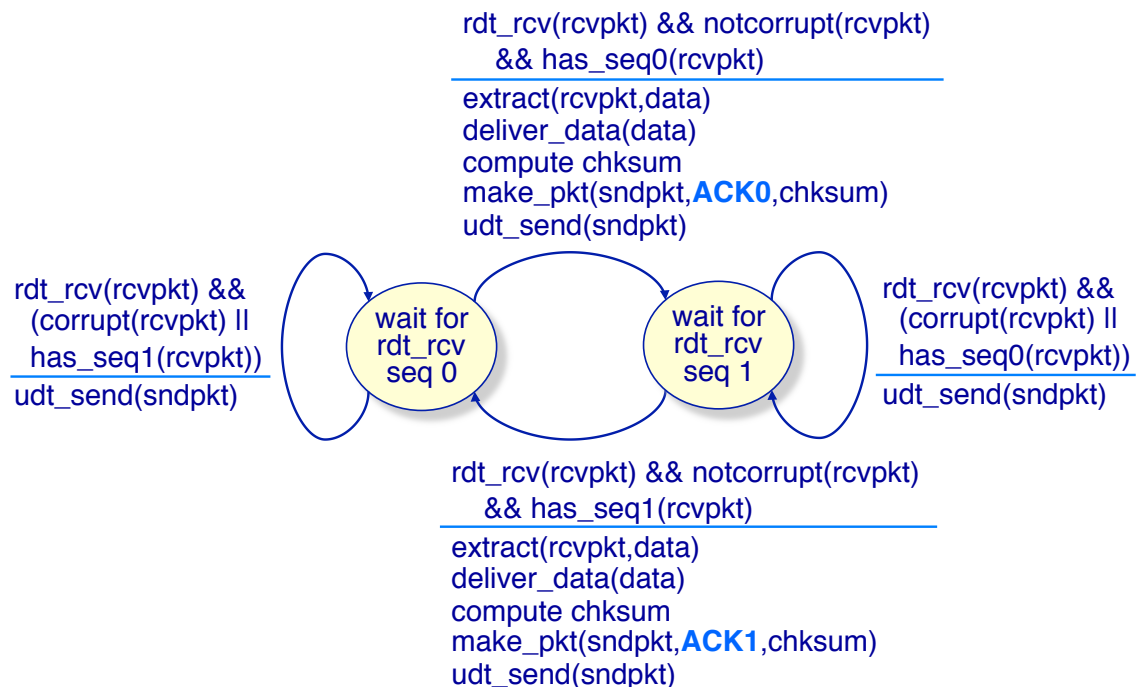## Sender state machine to handle lost/garbled packets

rdt_send(data)
compute chksum
make_pkt(pkt0,0,data,chksum)
udt_send(pkt0)
start_timer

rdt_rcv(rcvpkt)

rdt_rcv(rcvpkt) &&
(corrupt(rcvpkt) ||
  isACK(rcvpkt,1))

**wait for rdt_send 0**

**wait for ACK 0**

timeout
udt_send(pkt0)
start_timer

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

stop_timer

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

stop_timer

timeout
udt_send(pkt1)
start_timer

**wait for ACK 1**

**wait for rdt_send 1**

rdt_rcv(rcvpkt)

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

rdt_send(data)
compute chksum
make_pkt(pkt1,1,data,chksum)
udt_send(pkt1)
start_timer

29

# Receiver State Machine for RDT 2.2
## What changes are needed to handle lost/garbled packets?

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
  && has_seq0(rcvpkt)

extract(rcvpkt,data)
deliver_data(data)
compute chksum
make_pkt(sndpkt,**ACK0**,chksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) &&
(corrupt(rcvpkt) ||
has_seq1(rcvpkt))
udt_send(sndpkt)

**wait for rdt_rcv seq 0**

**wait for rdt_rcv seq 1**

rdt_rcv(rcvpkt) &&
(corrupt(rcvpkt) ||
has_seq0(rcvpkt))
udt_send(sndpkt)

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
  && has_seq1(rcvpkt)

extract(rcvpkt,data)
deliver_data(data)
compute chksum
make_pkt(sndpkt,**ACK1**,chksum)
udt_send(sndpkt)

30

# Fundamental Transport Layer Services
## Principles of reliable data transfer

◆ Use acknowledgements (ACKs) to indicate that a packet has been received

◆ Simple protocol:
- » *stop-and-wait* - can't send a new packet until the previous packet has been acknowledged
- » packet loss - sender sets a timer and re-sends the packet if no ACK received when timer expires
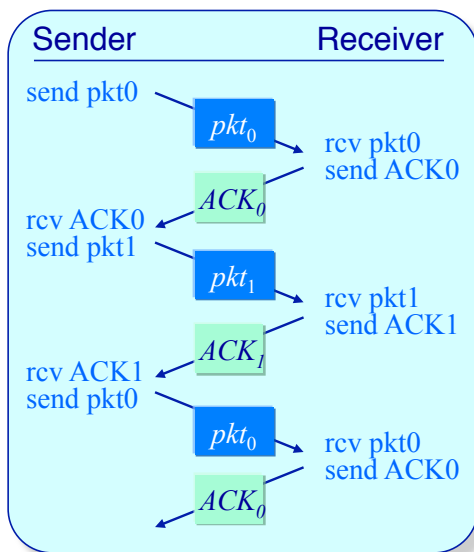- » ACK loss - ACKs are not retransmitted

# RDT 3.0
## Overview

◆ Sender
- » put a sequence number (0 or 1) on each packet
- » when receive an non-duplicate ACK
  - ❖ advance seqno
  - ❖ reset the timer
  - ❖ send the next packet
- » when receive a duplicate ACK
  - ❖ wait for a non-duplicate ACK
- » if timer expires before ACK received
  - ❖ re-send the previous packet

◆ Receiver
- » keep track of which seqno expected next (0 or 1)
- » when receive the next seqno expected
  - ❖ send an ACK for this seqno
  - ❖ advance next seqno expected
- » when receive a duplicate packet (packet isn't the next expected)
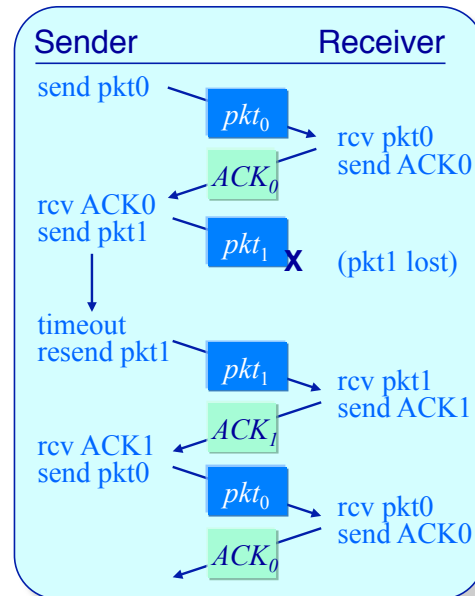  - ❖ re-send last ACK (for last seqno)

# Reliable Data Transfer
## Simple Protocol Examples



**Sender**      **Receiver**

send pkt0
$pkt_0$
rcv pkt0
send ACK0
$ACK_0$
rcv ACK0
send pkt1
$pkt_1$
rcv pkt1
send ACK1
$ACK_1$
rcv ACK1
send pkt0
$pkt_0$
rcv pkt0
send ACK0
$ACK_0$

- ◆ Protocol operation with no loss



**Sender**      **Receiver**

send pkt0
$pkt_0$
rcv pkt0
send ACK0
$ACK_0$
rcv ACK0
send pkt1
$pkt_1$ **X** (pkt1 lost)
timeout
resend pkt1
$pkt_1$
rcv pkt1
send ACK1
$ACK_1$
rcv ACK1
send pkt0
$pkt_0$
rcv pkt0
send ACK0
$ACK_0$

- ◆ Protocol operation with a lost packet

33

# Reliable Data Transfer
## Simple Protocol Examples



**Sender**      **Receiver**

send pkt0
$pkt_0$
rcv pkt0
send ACK0
$ACK_0$
rcv ACK0
send pkt1
$pkt_1$
rcv pkt1
send ACK1
$ACK_1$ **X**
timeout
resend pkt1
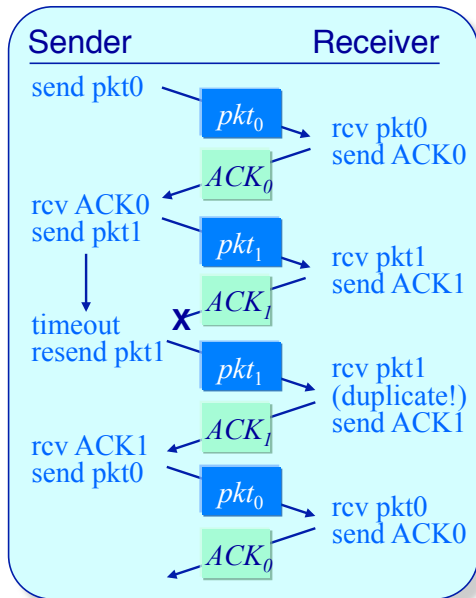$pkt_1$
rcv pkt1
(duplicate!)

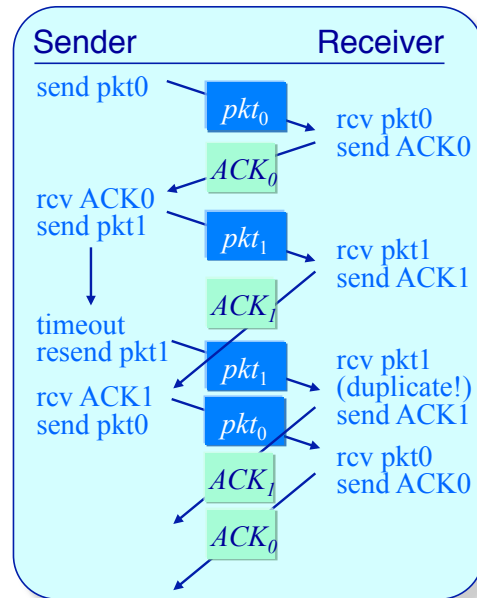What should happen next?

- ◆ Protocol operation with a lost ACK

34

# Reliable Data Transfer
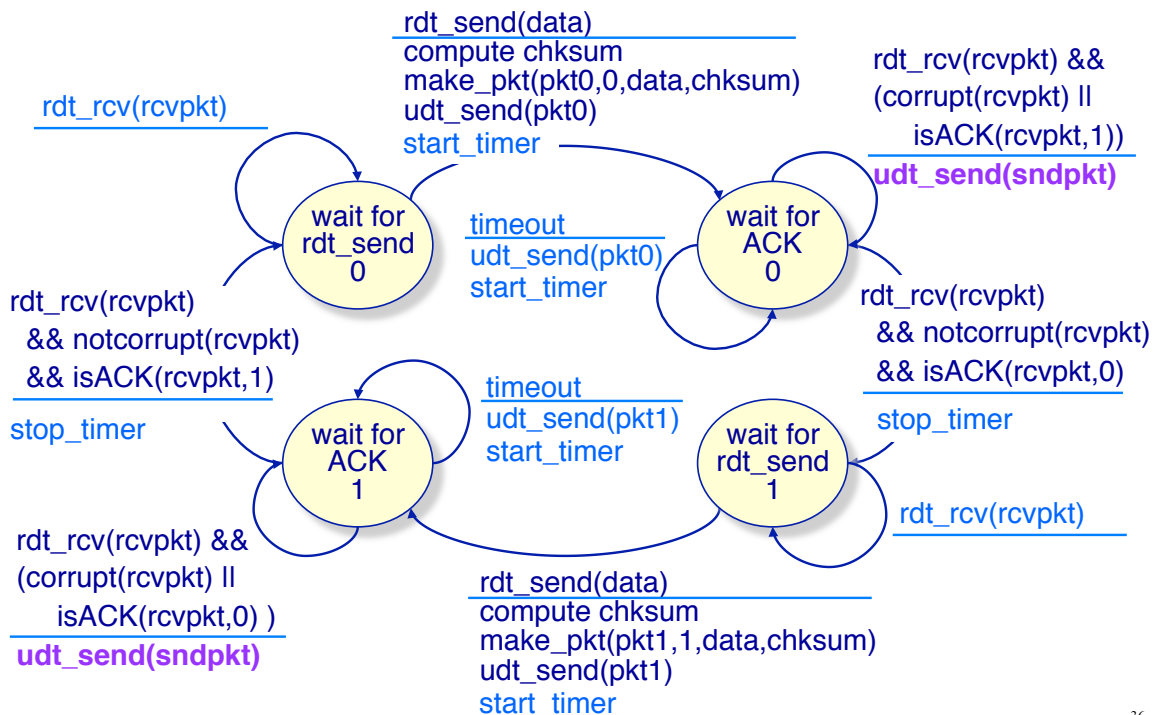## Simple Protocol Examples



◆ Protocol operation with a lost ACK

◆ Protocol operation with a poor timeout value
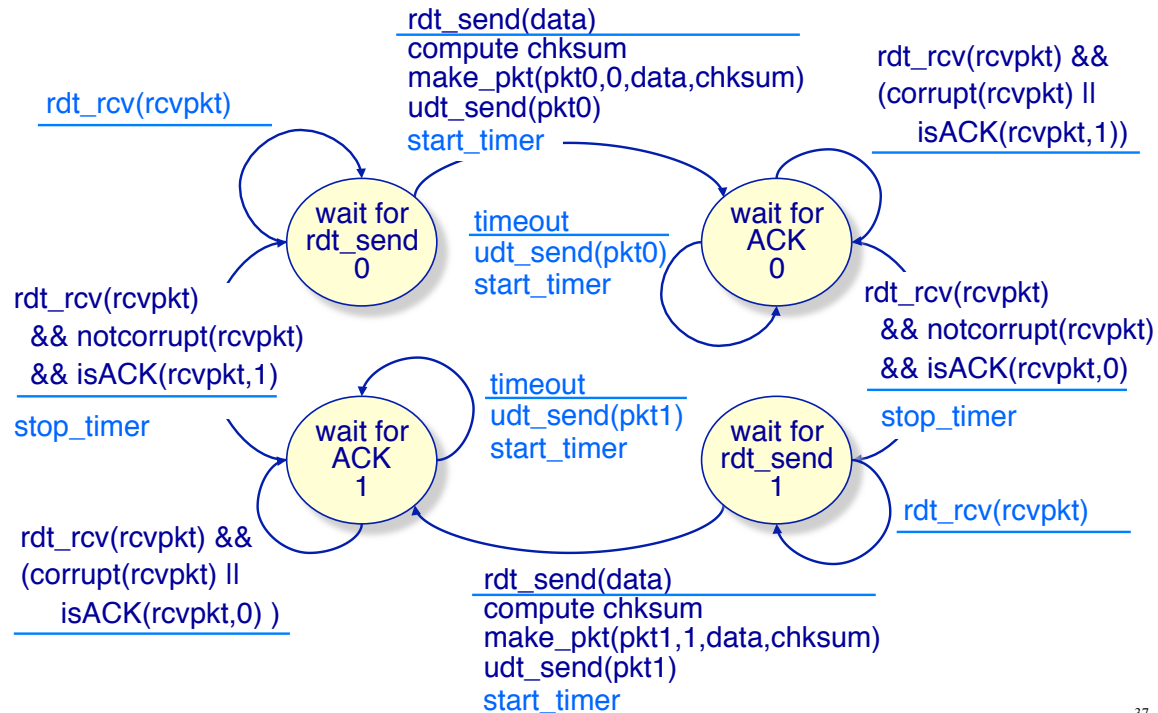
# Reliable Data Transfer Protocol 3.0
## Sender state machine to handle lost/garbled packets



rdt_send(data)
compute chksum
make_pkt(pkt0,0,data,chksum)
udt_send(pkt0)
start_timer

rdt_rcv(rcvpkt)

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

wait for rdt_send 0

timeout
udt_send(pkt0)
start_timer

wait for ACK 0

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

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

timeout
udt_send(pkt1)
start_timer

wait for ACK 1

wait for rdt_send 1

rdt_rcv(rcvpkt)

rdt_rcv(rcvpkt) &&
(corrupt(rcvpkt) ll
    isACK(rcvpkt,0) )
**udt_send(sndpkt)**

rdt_send(data)
compute chksum
make_pkt(pkt1,1,data,chksum)
udt_send(pkt1)
start_timer

# Reliable Data Transfer Protocol 3.0
## Sender state machine to handle lost/garbled packets

rdt_send(data)
compute chksum
make_pkt(pkt0,0,data,chksum)
udt_send(pkt0)
start_timer

rdt_rcv(rcvpkt)

rdt_rcv(rcvpkt) &&
(corrupt(rcvpkt) ||
isACK(rcvpkt,1))

**wait for rdt_send 0**

timeout
udt_send(pkt0)
start_timer

**wait for ACK 0**

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

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

timeout
udt_send(pkt1)
start_timer

**wait for ACK 1**

**wait for rdt_send 1**

rdt_rcv(rcvpkt)

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

rdt_send(data)
compute chksum
make_pkt(pkt1,1,data,chksum)
udt_send(pkt1)
start_timer

37

# Transport Protocol Performance
## Performance of RDT 3.0

◆ Can an end-system make efficient use of a network under this protocol?

◆ Consider a 1 Gbps link with 15 *ms* end-to-end propagation delay

◆ How busy is the network?

$$utilization = \frac{time\ network\ busy}{observation\ interval} = \frac{time\ to\ transmit\ a\ packet}{packet\ generation\ time}$$

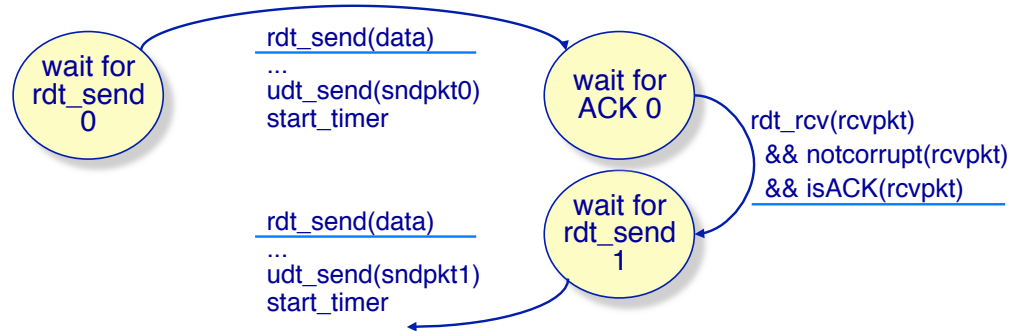◆ How long does it take to transmit a 1,000 byte packet?

$$transmission\ time = \frac{1\ kB\ packet \times 8\ bits/B}{10^9\ bps} = 8\ \mu s$$

◆ How fast can an end-system transmit packets?

38

# Transport Protocol Performance
## Performance of RDT 3.0



- ◆ How fast can an end-system transmit packets?
  - » Packet generation/transmission time = 8 $\mu$s (0.008 ms)
  - » Propagation delay to receiver = 15 ms
  - » ACK generation/transmission time ≈ 8 $\mu$s (0.008 ms)
  - » Propagation time for ACK to return to sender = 15 ms
- ◆ 1 packet every 30.016 ms

# Reliable Data Transfer
## Performance

- ◆ How busy is the network?

$$utilization = \frac{time\ network\ busy}{observation\ interval} = \frac{time\ to\ transmit\ a\ packet}{packet\ generation\ time}$$

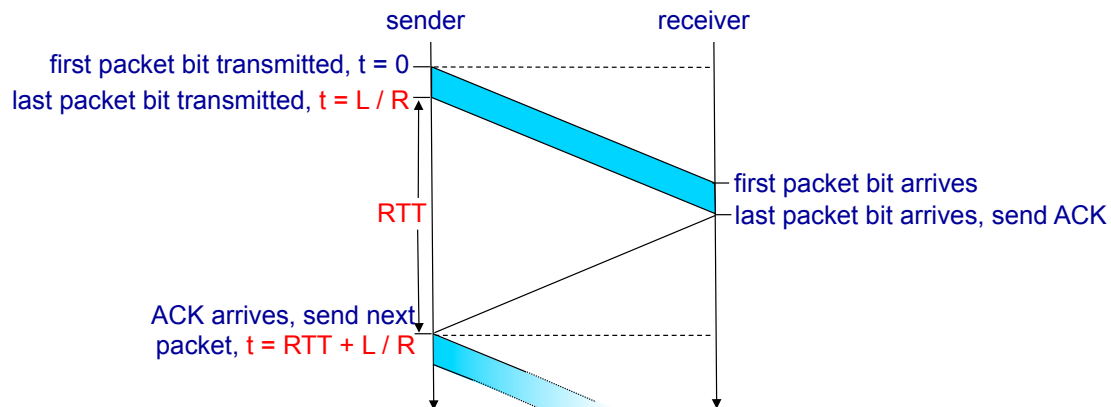$$= \frac{8\ \mu s}{30.016\ ms} = 0.027\%$$

- ◆ Is this good?
  - » 1,000 byte packet every 30 *ms* results in (maximum) throughput of 266 *kbps* over a 1 Gbps link!
  
    (266,000 *bps* over a 1,000,000,000 *bps* link)

*Network protocols limit the use of physical resources!*

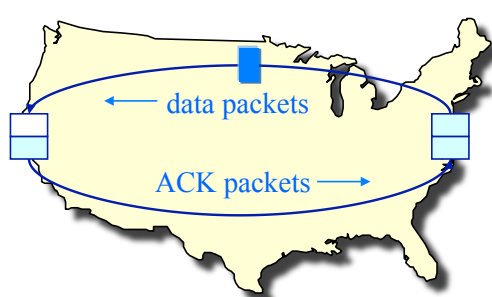# Reliable Data Transfer 3.0
## Stop and Wait

sender       receiver

first packet bit transmitted, t = 0

last packet bit transmitted, t = L / R

first packet bit arrives

last packet bit arrives, send ACK

RTT

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

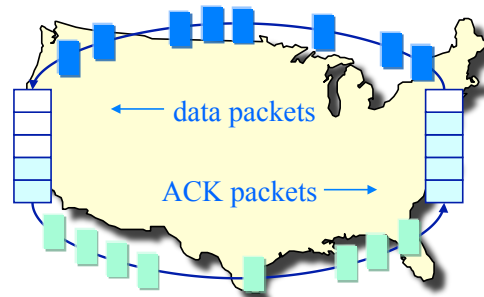$$U_{sender} = \frac{L / R}{RTT + L / R} = \frac{.008}{30.008} = 0.00027$$

# Improving Transport Protocol Performance
## Pipelining data transmissions

◆ Performance can be improved by allowing the sender to have multiple unacknowledged packets "in flight"

data packets

ACK packets

data packets

ACK packets

Stop-and-Wait protocol      Pipelined protocol
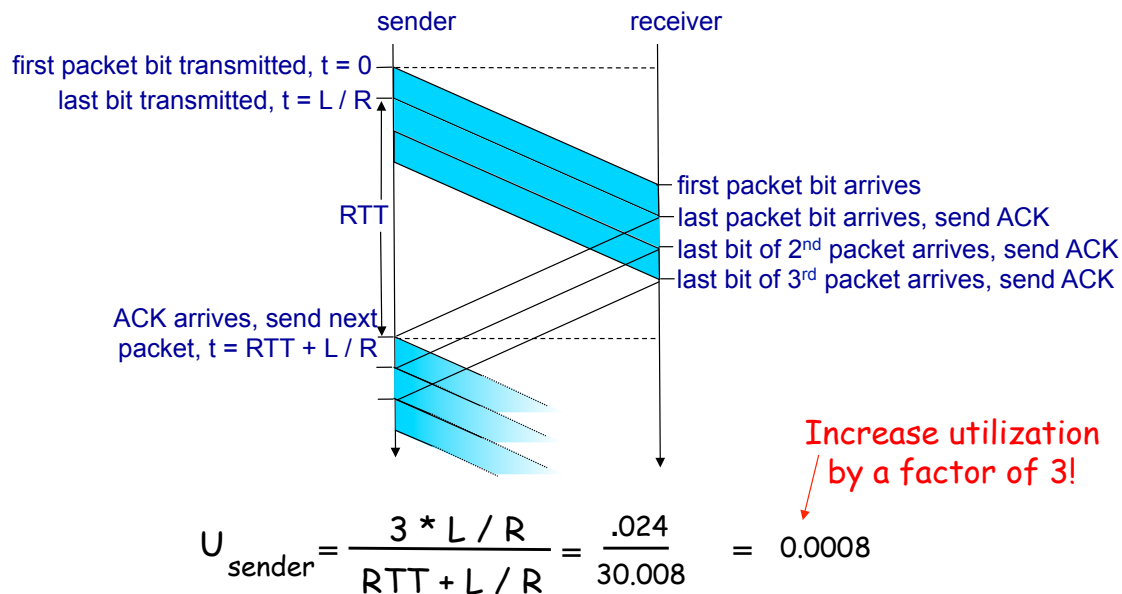
◆ Issues
  » The range of sequence numbers must be increased
  » ACKs need sequence numbers (what packet is being ACKed?)
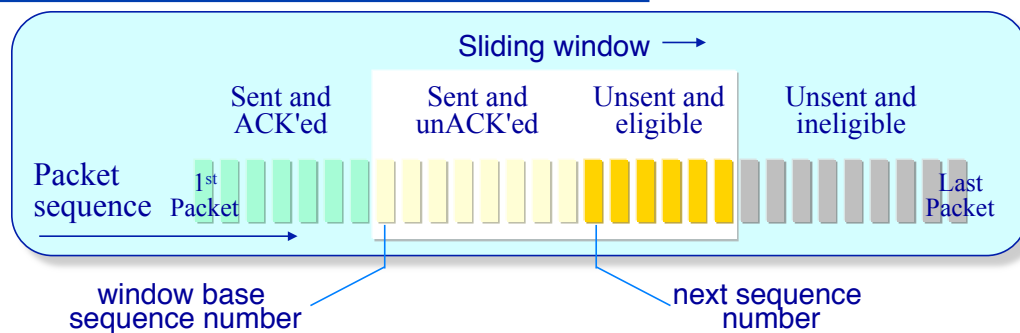  » More packets must be buffered at sender and receiver

# Reliable Data Transfer
## Pipelining



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 2nd packet arrives, send ACK

last bit of 3rd packet arrives, send ACK

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

Increase utilization by a factor of 3!

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

# Pipelined Protocols
## "Go-Back-*n*" protocols



Sliding window ⟶

Sent and ACK'ed

Sent and unACK'ed

Unsent and eligible

Unsent and ineligible

Packet sequence

1st Packet

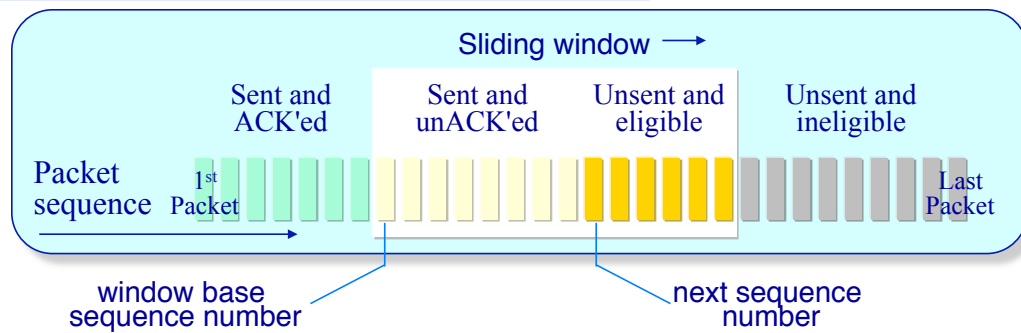Last Packet

window base sequence number

next sequence number

- ◆ Packet header contains a *k*-bit sequence number
- ◆ A "window" of up to $N \le 2^k$ consecutive, unacknowledged packets allowed to be in-flight
    - » Up to *N* packets may be buffered at the sender
    - » Window advances as ACKs are received
- ◆ Receiver generates "cumulative ACKs"
    - » ACKs contain the sequence number of the last in-order packet received

# Pipelined Protocols
## "Go-Back-*n*" protocols

Sliding window →

| Sent and ACK'ed | Sent and unACK'ed | Unsent and eligible | Unsent and ineligible |

Packet sequence — 1ˢᵗ Packet ... Last Packet

window base sequence number
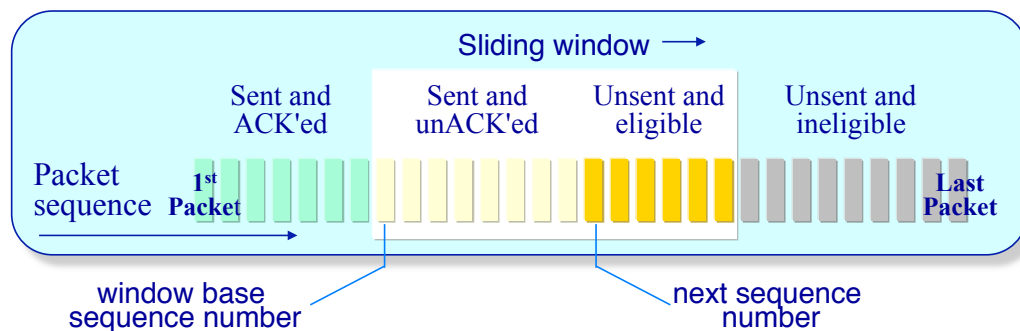
next sequence number

- ◆ Receiver protocol
    - » Use cumulative ACKs — ACK packet *n* only if all packets numbered less than *n* have been received
    - » If losses occur, sender may receive duplicate ACKs
- ◆ Sender protocol
    - » A timer is set for the each (or just the oldest) in-flight packet
    - » On timeout for packet *n*, retransmit packet *n* and all higher number packets in the current window

# Go-Back-n Protocol
## Sender

Sliding window →

| Sent and ACK'ed | Sent and unACK'ed | Unsent and eligible | Unsent and ineligible |

Packet sequence — 1ˢᵗ Packet ... Last Packet

window base sequence number

next sequence number

- ◆ Sender waits for an event:
    - » application has data to send
    - » timer goes off
    - » ACK is received

# Go-Back-*n* Protocol
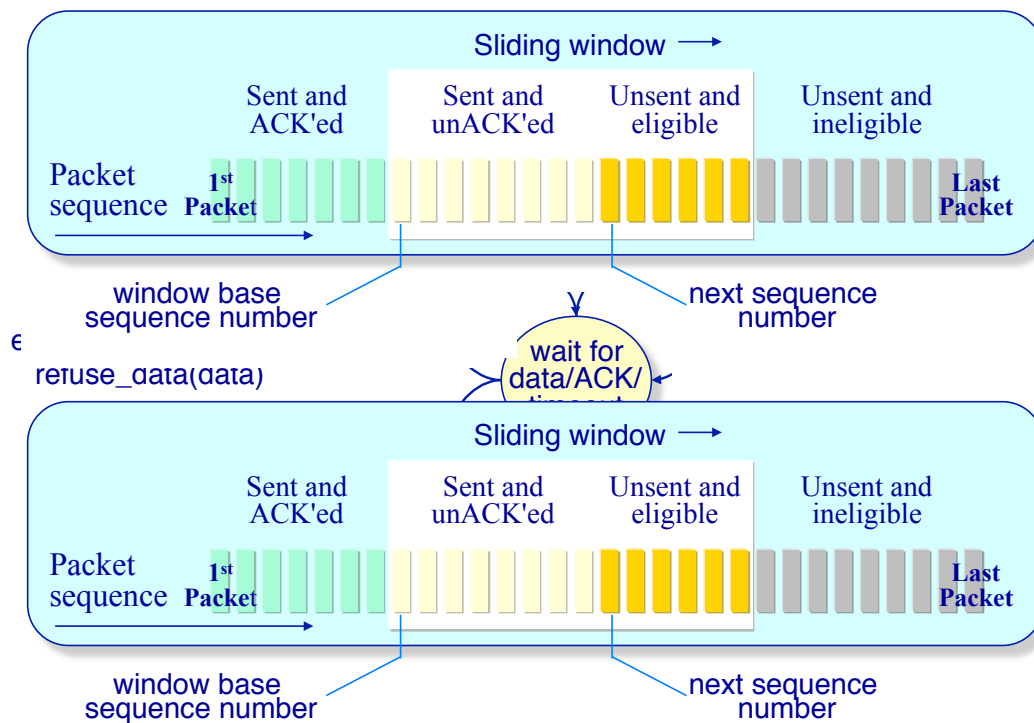## Sender extended FSM

◆ THIS SLIDE INTENTIONALLY LEFT BLANK!

# Go-Back-n Protocol
## Sender extended FSM

Sliding window ⟶

| Sent and ACK'ed | Sent and unACK'ed | Unsent and eligible | Unsent and ineligible |

Packet sequence

1st Packet

Last Packet

window base sequence number

next sequence number

wait for data/ACK/

refuse_data(data)

Sliding window ⟶

| Sent and ACK'ed | Sent and unACK'ed | Unsent and eligible | Unsent and ineligible |

Packet sequence

1st Packet

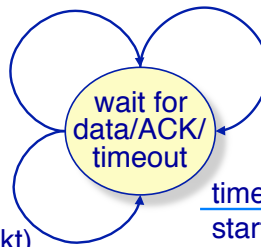Last Packet

window base sequence number

next sequence number

# Go-Back-n Protocol
## Sender extended FSM

rdt_send(data)
if (nextseqnum < base+N) {
   compute chksum
   make_pkt(sndpkt[nextseqnum],nextseqnum,data,chksum)
   udt_send(sndpkt[nextseqnum])
   if (base == nextseqnum) start_timer
   nextseqnum += 1
  }
else
   refuse_data(data)

wait for data/ACK/ timeout

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
base := getacknum(rcvpkt) + 1
if (base == nextseqnum)
  stop_timer
else
  start_timer

timeout
start_timer
udt_send(sndpkt[base])
udt_send(sndpkt[base+1])
...
udt_send(sndpkt[nextseqnum–1])

# Go-Back-*n* Protocol
## Sender extended FSM

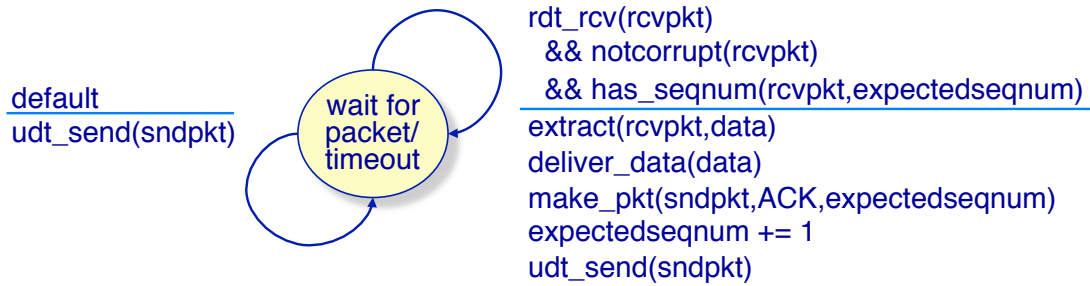◆ THIS SLIDE INTENTIONALLY LEFT BLANK!

# Go-Back-n Protocol
## Receiver extended FSM

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& has_seqnum(rcvpkt,expectedseqnum)

extract(rcvpkt,data)
deliver_data(data)
make_pkt(sndpkt,ACK,expectedseqnum)
expectedseqnum += 1
udt_send(sndpkt)

wait for packet/ timeout
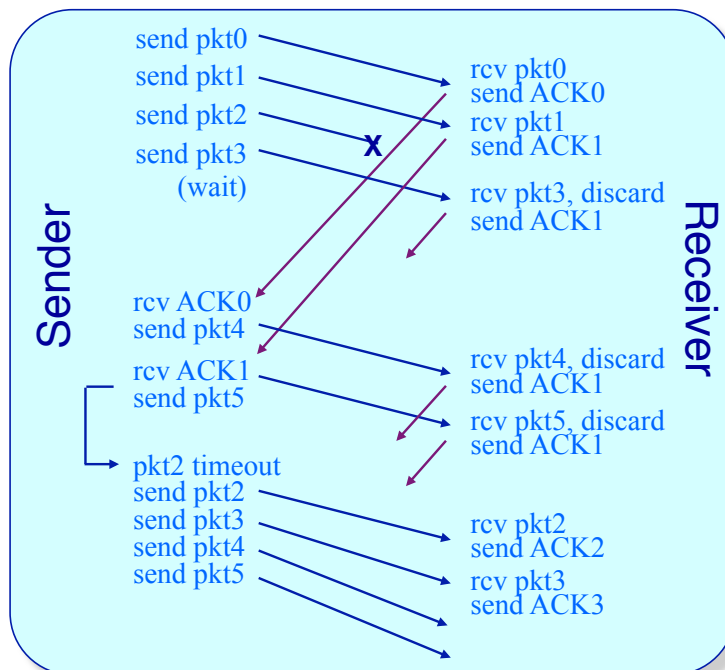
default
udt_send(sndpkt)

- ◆ In-order packets processed, out-of-order packets discarded
  - » Sender will eventually timeout and retransmit out-of-order packets
  - » Thus the receiver need not buffer any packets

- ◆ Always send ACK for correctly-received packet with highest in-order sequence number
  - » May generate duplicate ACKs
  - » But minimal state — need only remember *expectedseqnum*

51

# Go-Back-*n* Protocol
## Execution example

Sender

send pkt0
send pkt1
send pkt2
send pkt3
(wait)

rcv ACK0
send pkt4

rcv ACK1
send pkt5

pkt2 timeout
send pkt2
send pkt3
send pkt4
send pkt5

Receiver

rcv pkt0
send ACK0
rcv pkt1
send ACK1

rcv pkt3, discard
send ACK1

rcv pkt4, discard
send ACK1

rcv pkt5, discard
send ACK1

rcv pkt2
send ACK2
rcv pkt3
send ACK3

- ◆ Assume a window size of 4 packets

- ◆ Receiver ignores out-of-order packets

- ◆ Sender retransmits only on timeout

52

# Transport Protocol Performance
## Performance of Go-Back-*n* protocols

- ◆ Can an end-system make more efficient use of a network under a Go-Back-*n* protocol?

- ◆ Consider again transmitting 1,000 byte packets on a 1 Gbps link with 15 *ms* end-to-end propagation delay

$$utilization \quad = \quad \frac{time\ to\ transmit\ a\ packet}{packet\ generation\ time}$$

$$\begin{matrix} transmission \\ time \end{matrix} \quad = \quad \frac{1\ kB\ packet \times 8\ bits/B}{10^9\ bps} \quad = \quad 8\ \mu s$$

- ◆ How fast can an end-system transmit packets?
  - » Depends on the window size!

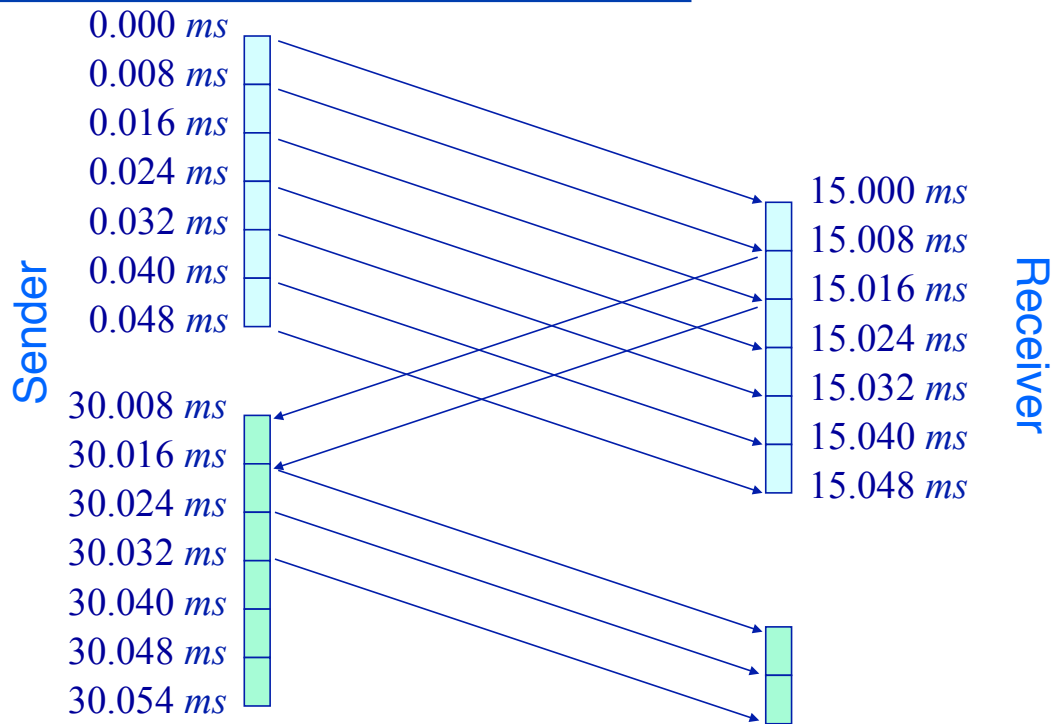# Transport Protocol Performance
## Performance of Go-Back-n protocols

- ◆ How fast can an end-system transmit packets?
  - » N packets can be sent before the sender must wait for an ACK

- ◆ N packets sent every 30.016 ms
  - » Packet generation/transmission time = 8 *μ*s
  - » Round-trip-time to receiver = 30 ms
  - » ACK generation/transmission time ≈ 8 *μ*s

# Transport Protocol Performance

## Performance of Go-Back-*n* protocols
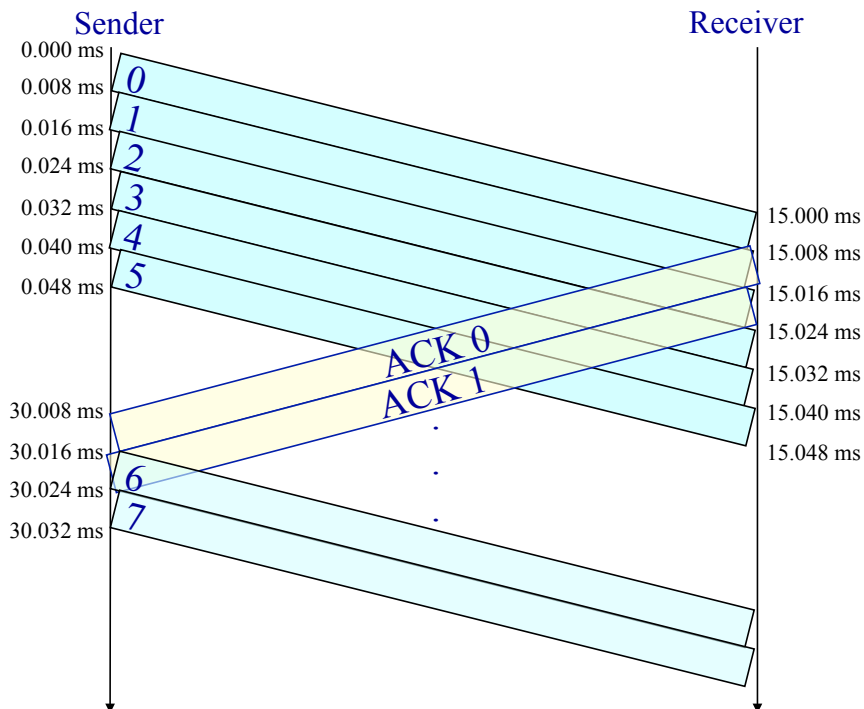
| Sender | | Receiver |
|---|---|---|
| 0.000 *ms* | | |
| 0.008 *ms* | | 15.000 *ms* |
| 0.016 *ms* | | 15.008 *ms* |
| 0.024 *ms* | | 15.016 *ms* |
| 0.032 *ms* | | 15.024 *ms* |
| 0.040 *ms* | | 15.032 *ms* |
| 0.048 *ms* | | 15.040 *ms* |
| | | 15.048 *ms* |
| 30.008 *ms* | | |
| 30.016 *ms* | | |
| 30.024 *ms* | | |
| 30.032 *ms* | | |
| 30.040 *ms* | | |
| 30.048 *ms* | | |
| 30.054 *ms* | | |

# Transport Protocol Performance

## Performance of Go-Back-*n* protocols

Sender                                      Receiver

| | |
|---|---|
| 0.000 ms | |
| 0.008 ms | 0 |
| 0.016 ms | 1 |
| 0.024 ms | 2 |
| 0.032 ms | 3 |
| 0.040 ms | 4 |
| 0.048 ms | 5 |

15.000 ms
15.008 ms
15.016 ms
15.024 ms
15.032 ms
15.040 ms
15.048 ms

ACK 0
ACK 1

| | |
|---|---|
| 30.008 ms | |
| 30.016 ms | |
| 30.024 ms | 6 |
| 30.032 ms | 7 |

# Transport Protocol Performance
## Performance of Go-Back-*n* protocols

◆ Performance with a window size of $N = 64$ packets:

$$utilization = \frac{\text{time to transmit N packets}}{\text{time to receipt of first ACK}}$$

$$= \frac{512 \ \mu s}{30.016 \ ms} = 1.7\%$$

A 64**x** improvement!

◆ Is this good?
  - » 64 1,000 byte packets every 30 *ms* results in (maximum) throughput of 17 *Mbps* over a 1 Gbps link!
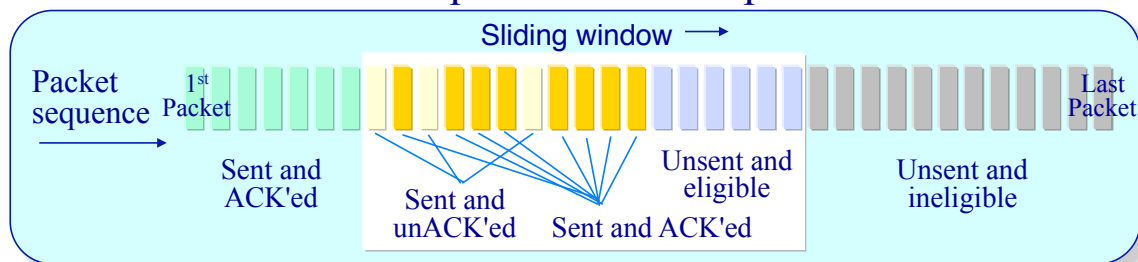
# Pipelined Protocols
## "Selective Repeat" protocols

◆ Receiver *individually* acknowledges all correctly received packets
  - » Buffers packets as needed for eventual in-order delivery to upper layer
◆ Sender only resends packets for which an ACK has not been received
  - » Sender maintains a timer for each unACK'ed packet
◆ Sender window is the same as before
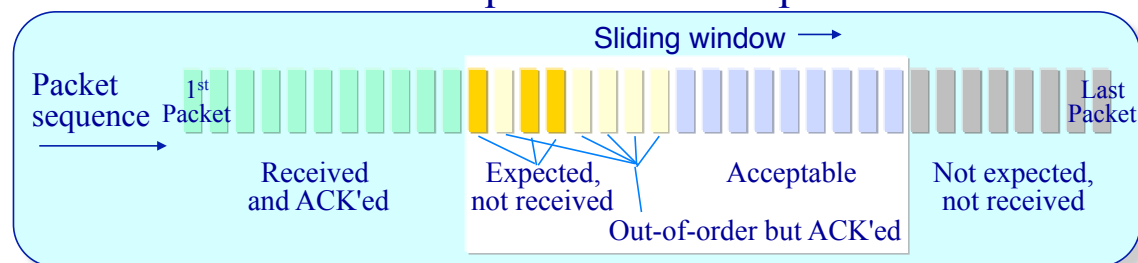  - » *N* consecutive sequence numbers (Limits the sequence numbers of sent, unACK'ed packets)

# Selective Repeat Protocols
## Sender and receiver windows

◆ Sender's view of sequence number space

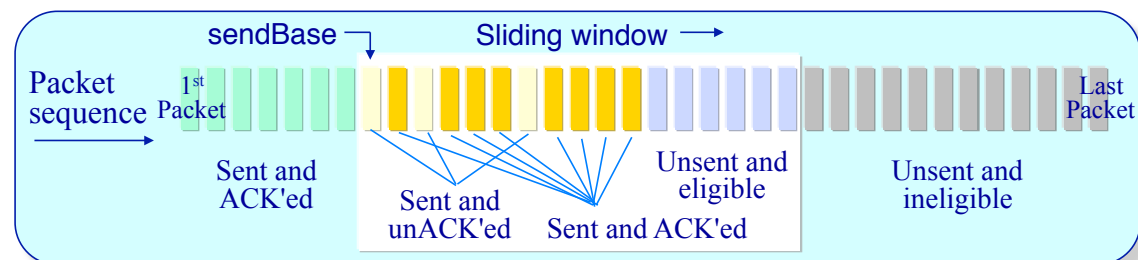Sliding window →

Packet sequence

1st Packet

Sent and ACK'ed

Sent and unACK'ed

Sent and ACK'ed

Unsent and eligible

Unsent and ineligible

Last Packet

◆ Receiver's view of sequence number space

Sliding window →

Packet sequence

1st Packet

Received and ACK'ed

Expected, not received

Out-of-order but ACK'ed

Acceptable

Not expected, not received

Last Packet

# Selective Repeat Protocols
## Sender state machine

sendBase

Sliding window →

Packet sequence

1st Packet

Sent and ACK'ed

Sent and unACK'ed

Sent and ACK'ed
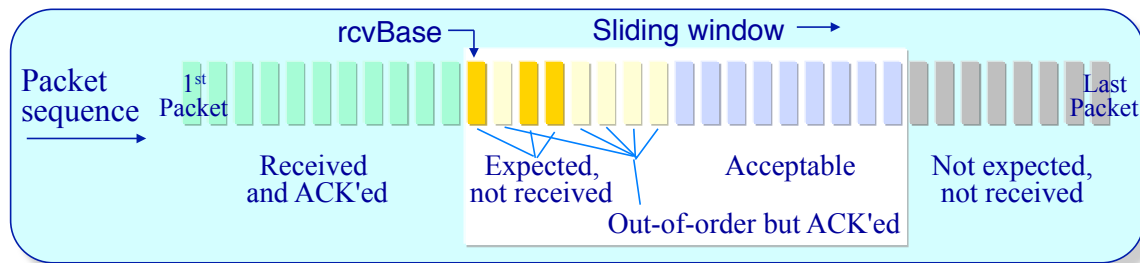
Unsent and eligible

Unsent and ineligible

Last Packet

◆ Call from above:
  » If next available sequence number is within window, send the packet and start a timer for it

◆ Timeout for packet *n*:
  » Resend packet *n*, restart timer for packet *n*

◆ ACK received for packet with sequence number *n*:
  » If *n* in [*sendBase*, *sendBase*+*N*–1] then mark packet *n* as received
  » If *n* == *sendBase*, advance *sendBase* to next highest unACKed sequence number and move the window forward by that amount
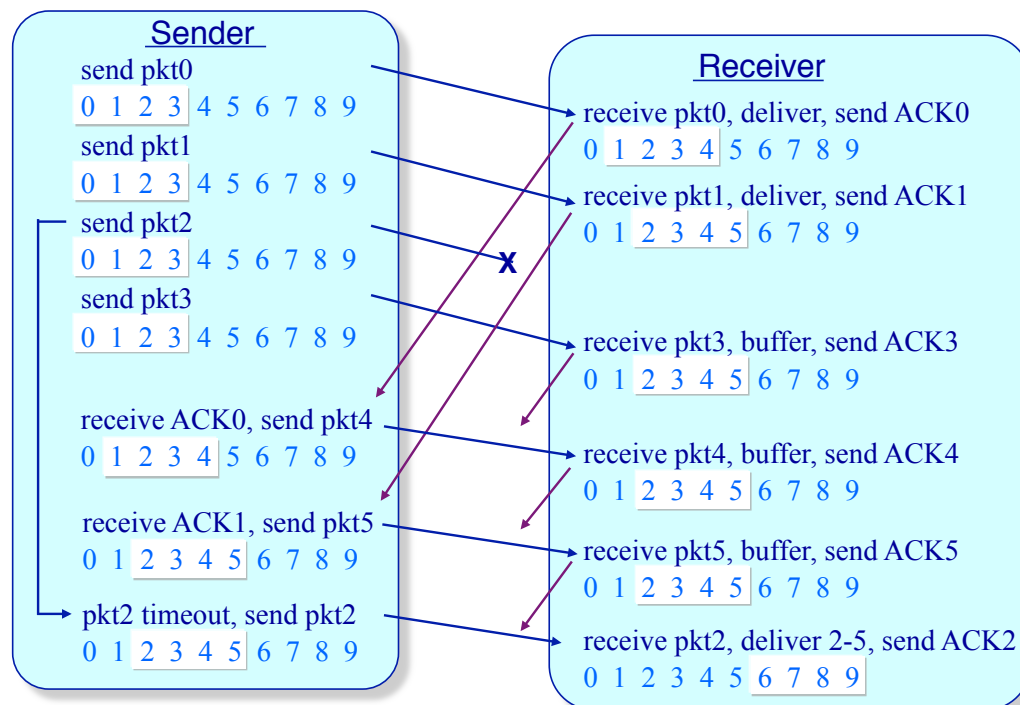
# Selective Repeat Protocols
## Receiver state machine



- ◆ Packet *n* in [*rcvbase*, *rcvbase+N*–1] correctly received:
  - » Send an ACK for packet *n*
  - » If packet *n* is out-of-order then buffer
  - » If *n* == *rcvBase*, deliver packet *n*, and all other buffered consecutive in-order packets, to application, and advance the window by the number of delivered packets
- ◆ Packet *n* in [*rcvbase*–*N*, *rcvbase*–1] received:
  - » Send an ACK for packet *n*
- ◆ Otherwise discard packet (without ACK'ing)
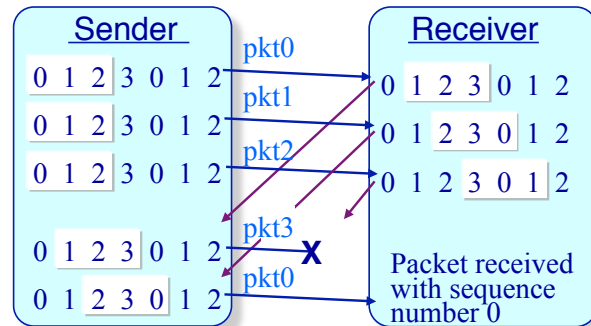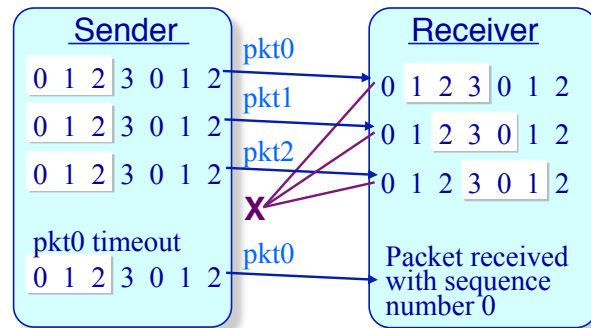
61

# Selective Repeat Protocols
## Execution example



62

# Selective Repeat Protocols
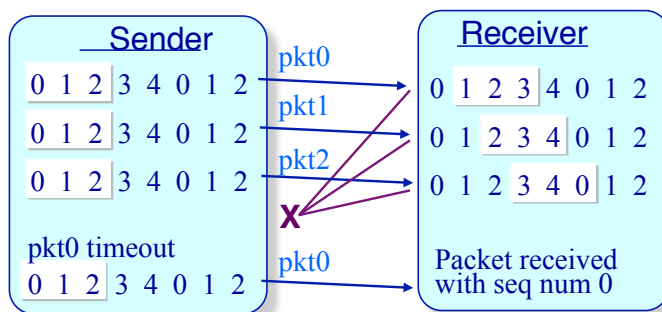## Window state ambiguity

- How many sequence numbers do we need?
  - » As many as the largest number of packets that can be in flight?

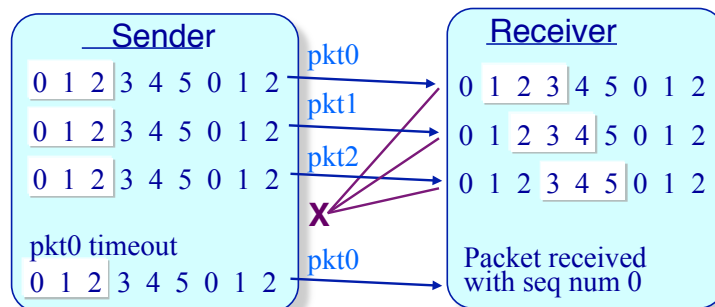- If the sequence number space is close to the window size then the receiver can get confused

**Sender**

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

pkt0 timeout

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

pkt0
pkt1
pkt2
X
pkt0

**Receiver**

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

Packet received with sequence number 0

**Sender**

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

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

pkt0
pkt1
pkt2
pkt3
X
pkt0

**Receiver**

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

Packet received with sequence number 0

63

# Selective Repeat Protocols
## Window state ambiguity

**Sender**

| 0 | 1 2 | 3 | 4 | 0 | 1 | 2 |
| 0 | 1 2 | 3 | 4 | 0 | 1 | 2 |
| 0 | 1 2 | 3 | 4 | 0 | 1 | 2 |

pkt0 timeout

| 0 | 1 2 | 3 | 4 | 0 | 1 | 2 |

pkt0
pkt1
pkt2
X
pkt0

**Receiver**

| 0 | 1 2 3 | 4 | 0 | 1 | 2 |
| 0 | 1 | 2 3 4 | 0 | 1 | 2 |
| 0 | 1 | 2 | 3 4 0 | 1 | 2 |

Packet received with seq num 0

**5 seq nums**

**Sender**

| 0 | 1 2 | 3 | 4 | 5 | 0 | 1 | 2 |
| 0 | 1 2 | 3 | 4 | 5 | 0 | 1 | 2 |
| 0 | 1 2 | 3 | 4 | 5 | 0 | 1 | 2 |

pkt0 timeout

| 0 | 1 2 | 3 | 4 | 5 | 0 | 1 | 2 |

pkt0
pkt1
pkt2
X
pkt0

**Receiver**

| 0 | 1 2 3 | 4 | 5 | 0 | 1 | 2 |
| 0 | 1 2 3 4 | 5 | 0 | 1 | 2 |
| 0 | 1 2 3 4 5 | 0 | 1 | 2 |

Packet received with seq num 0

**6 seq nums**

64