# Last Lecture

- Peer-to-Peer (P2P) Applications

# This Lecture

- Overview of the transport layer

- Principles of Reliable Data Transfers

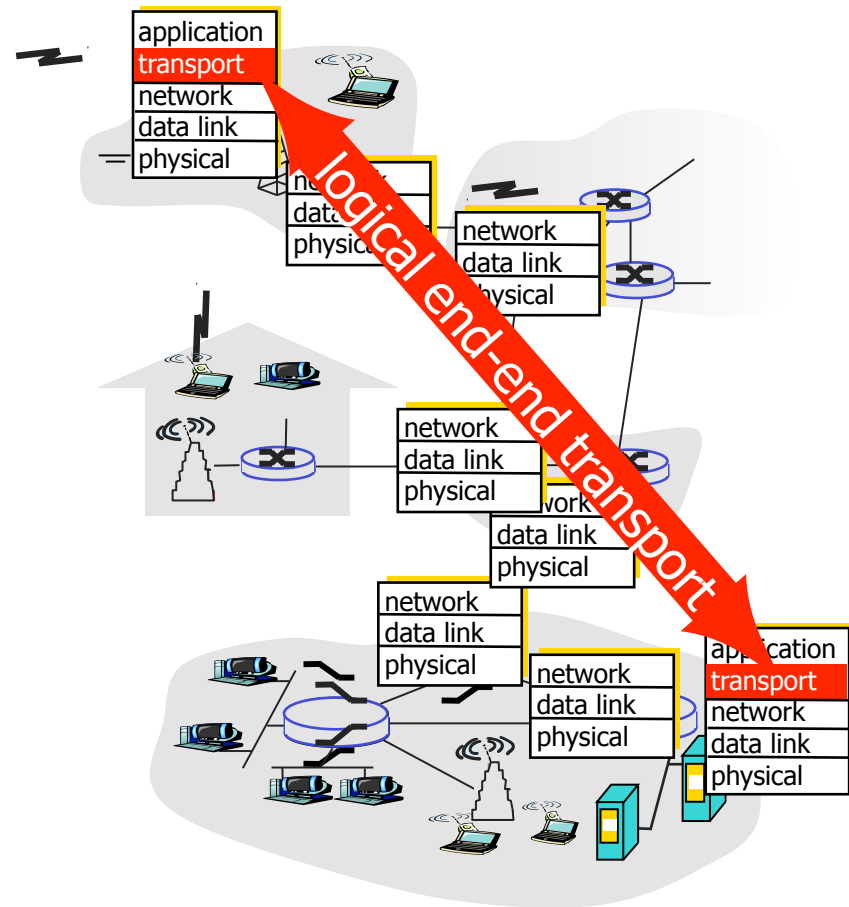# The Transport Layer

- ## Provide services to applications
  - ### What kind of services?
  - ### How to implement them?

- ## Make use of services provided by the network layer
  - ### Network gives *best-effort* packet delivery service

- ## Help networks out too
  - ### Don't pump too much data in if networks can't handle, i.e. *congestion control*

# What Services to Provide to Applications?

- Difficult to decide, because
  - Can't envision all future applications
  - Even current applications are too diverse in requirements
  - Can't provide services which can't be implemented
- Currently, two main services are
  - TCP: reliable, connection-oriented
  - UDP: unreliable, connectionless
- There are *many* other proposals & implementations, but not widespread
  - RTP, RSTP for real-time streaming
  - SCTP, DCCP: somewhere between UDP and TCP
- Many applications have transport functionalities built-in

# Services Provided by the Network Layer

- Depend on the network
  - *Datagram network*: service sucks! Just best-effort
  - *ATM network*: connection-oriented, virtual-circuit, some QoS guarantee
  - …

- A general transport protocol can only assume that network service is best effort
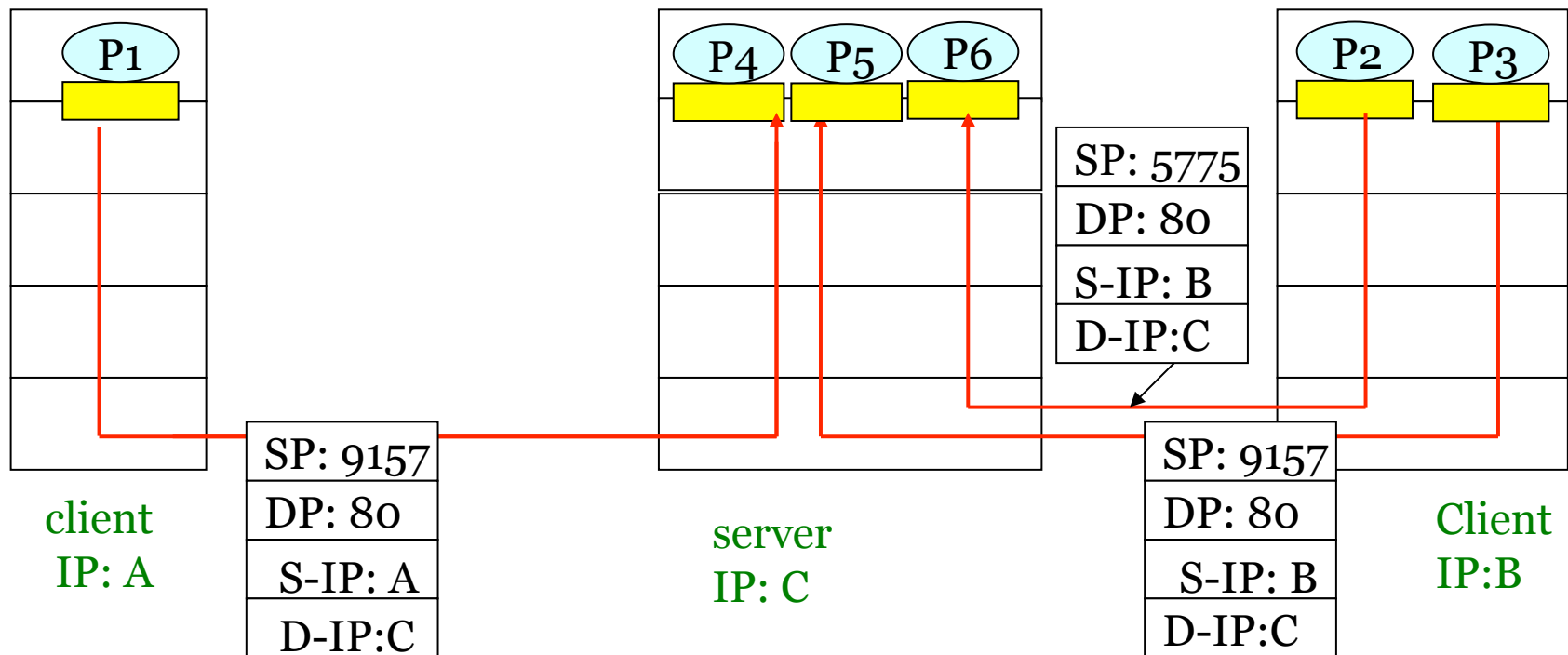
# What Is Best-Effort Again?

- Packets may be *corrupted*
- Packets may be *lost*
- Packets may be *duplicated*
- Packets may be delivered *out of order*
- Inter-arrival times can vary wildly
- End-to-end delay may vary wildly

# We Will Focus on TCP Alone

- *Multiplexing & de-multiplexing*
- *Reliable data transfer (& try to be efficient too)*
- *Connection-oriented*
- *Flow control*
- Help network with *congestion control & avoidance*
- **No** guarantee on timing (delay, jitter, bandwidth)

- TCP is sufficiently complex to illustrate fundamental ideas
- Services suitable for media streaming like RTP, RSTP, etc. are still active research topics!

# Multiplexing & De-multiplexing



client
IP: A

| SP: 9157 |
| DP: 80 |
| S-IP: A |
| D-IP:C |

server
IP: C

| SP: 5775 |
| DP: 80 |
| S-IP: B |
| D-IP:C |

| SP: 9157 |
| DP: 80 |
| S-IP: B |
| D-IP:C |

Client
IP:B

# Principles of Reliable Data Transfer
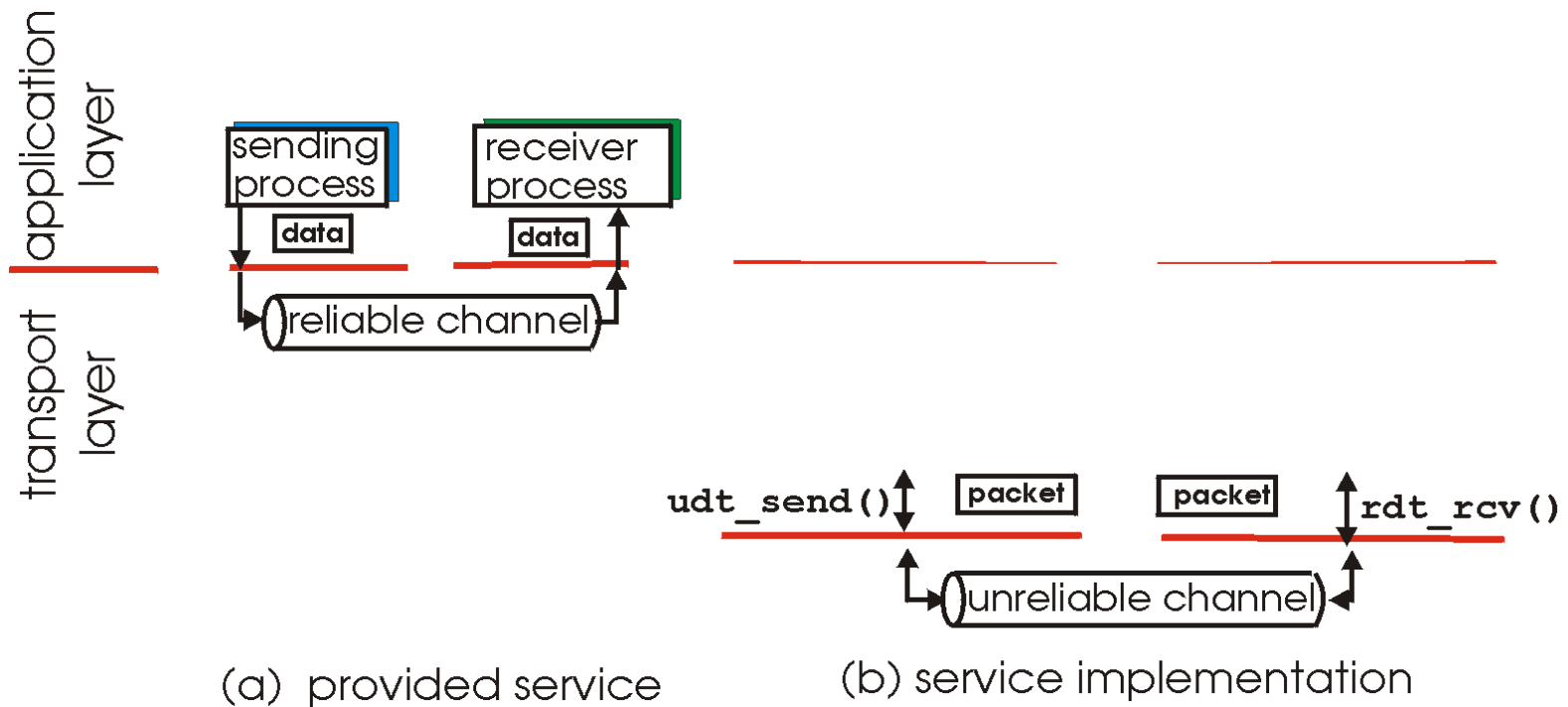
- Before looking at how TCP does it, let's try to design a *reliable data transfer protocol* (RDT) ourselves

- *Main question*: *how to reliably transfer data* when
  1. Network can *corrupt* packets (*bit error*)
  2. Network can *lose* packets
  3. Network can deliver *duplicates*
  4. Network can deliver packets *out of order*

- We address the 4 problems one by one, in that order

# The Bird-Eye View

application layer

transport layer

sending process

receiver process

data

data

reliable channel

udt_send() packet packet rdt_rcv()

unreliable channel

(a) provided service

(b) service implementation

# The Bird-Eye View



(a) provided service
(b) service implementation

- ■ Characteristics of unreliable channel will determine complexity of the RDT protocol

# RDT: Getting Started

rdt_send() ↓ data

data ↑ deliver_data()

send side

reliable data transfer protocol (sending side)

reliable data transfer protocol (receiving side)

receive side

udt_send() ↕ packet

packet ↕ rdt_rcv()

unreliable channel

# Finite State Machines

- FSMs are convenient for specifying protocol's behaviors

- FSM notations:

event causing state transition

actions taken on state transition

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

state
1

event

actions

state
2

# RDT 1.0: Perfectly Reliable Channel

Wait for call from above → rdt_send(data)

packet = make_pkt(data)
udt_send(packet)

Wait for call from below → rdt_rcv(packet)

extract (packet,data)
deliver_data(data)

Sender

Receiver

- *Next*: suppose network can corrupt packets (*i.e., bit errors* may occur)

- *But still:* no packet loss, no out of order packets, no duplicate packets

# RDT 2.0: Dealing with Bit-Errors

1. ## How to detect that a packet has been corrupted?
   - *Error-detecting* code (e.g. checksuming)

2. ## What to do when corrupted packet received?
   - *Error-correcting code*
     - May not always work, depend on how much error
     - Too much (time/space) overhead if error is rare
     - Decoding time might be too long
   - Tell sender to *retransmit*
     - ACK: acknowledgement of a good packet
     - NACK: acknowledgement of a bad packet

- ## New mechanisms in RDT 2.0:
  - Error detection
  - ARQ – *automatic repeat request*

# Error Detection

- *Problem*: detect bit errors in packets (frames)
- *Solution*: add *extra* bits to each packet
- *Goals*:
    - Reduce overhead (number of redundancy bits)
    - Increase the number and the type of bit error patterns that can be detected
- *Examples*:
    - Two-dimensional parity
    - (Internet) Checksum
    - Later (when we discuss data link layer)
        - Cyclic Redundancy Check (CRC)
        - Hamming Codes

# Parity

- **Even parity**
  - Add a parity bit to 7 bits of data to make an even number of 1's

| 0110100 | **1** |
|---------|-------|
| 1011010 | **0** |

- **How many bits of error can be detected by a parity bit?**

- **What's the overhead?**

# Two-dimensional parity

o Add one extra bit to a 7-bit code such that the number of 1's in the resulting 8 bits is even (for even parity, and odd for odd parity)
o Add a parity byte for the packet
o Example: five 7-bit character packet, even parity

| 0110100 | 1 |
| 1011010 | 0 |
| 0010110 | 1 |
| 1110101 | 1 |
| 1001011 | 0 |
| 1000110 | 1 |

# Two-dimensional parity detection capability

- All 1-bit errors
- Example:

| | |
|---|---|
| 0110100 | 1 |
| 1011010 | 0 |
| 0000110 | 1 | ← odd number of 1's |
| 1110101 | 1 |
| 1001011 | 0 |
| 1000110 | 1 |

error bit → (pointing to 0000110)

# Two-dimensional parity detection capability

- All 2-bit errors
- Example:

| 0110100 | **1** |
|---------|-------|
| 1011010 | **0** |
| 0000111 | **1** |
| 1110101 | **1** |
| 1001011 | **0** |
| 1000110 | **1** |

error bits →

odd number of 1's on columns

# Two-dimensional parity detection capability

- All 3-bit errors
- Example:

| | |
|---|---|
| 0110100 | 1 |
| 1011010 | 0 |
| 0000111 | 1 |
| 1100101 | 1 |
| 1001011 | 0 |
| 1000110 | 1 |

error bits

odd number of 1's on column

# Two-dimensional parity detection capability

- Most 4-bit errors
- Example of 4-bit error that is not detected:

| 0110100 | 1 |
|---------|---|
| 1011010 | 0 |
| 0000111 | 1 |
| 1100100 | 1 |
| 1001011 | 0 |
| 1000110 | 1 |

error bits

*How many errors can this code correct?*

# Internet Checksum [RFC1071]

- ## Used in TCP, UDP, IP

- ## *The Internet checksum algorithm:*

  - Adjacent octets to be checksummed are paired to form 16-bit integers, and the *1's complement sum* of these 16-bit integers is formed.

  - To generate a checksum, the checksum field itself is cleared, the 16-bit 1's complement sum is computed over the octets concerned, and the 1's complement of this sum is placed in the checksum field.

  - To check a checksum, the 1's complement sum is computed over the same set of octets, including the checksum field. If the result is all 1 bits (-0 in 1's complement arithmetic), the check succeeds.



*16-bit units of data*     **message to be checksummed**     *zeroes appended to make a multiple of 16 bits*

0

# Example of Internet Checksum

```
Computing the Checksum
     1000 0110 0101 1110     First 16-bit value
+    1010 1100 0110 0000     Second 16-bit value
     --------------------
   1 0011 0010 1011 1110     Carry-out "loops"
+  \---------------> 1       back into LBb
     --------------------
     0011 0010 1011 1111
+    0111 0001 0010 1010     Third 16-bit value
     --------------------
   0 1010 0011 1110 1001     No carry to swing around
+    1000 0001 1011 0101     Fourth 16-bit value
     --------------------
   1 0010 0101 1001 1110     Carry-out "loops"
+  \---------------> 1       back into LBb
     --------------------
     0010 0101 1001 1111     "One's complement sum"
     --------------------
     1101 1010 0110 0000     Take 1's complement
                             again, that's the
                             checksum of the data
```

```
Input data:
1000 0110 0101 1110
1010 1100 0110 0000
0111 0001 0010 1010
1000 0001 1011 0101
```

# RDT2.0: the FSM Specification

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

receiver

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

Wait for call from above

Wait for ACK or NAK

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

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

sender

Wait for call from below

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

# RDT2.0: Operation with No Errors



rdt_send(data)

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

**Wait for call from above**

**Wait for ACK or NAK**

rdt_rcv(rcvpkt) &&
isNAK(rcvpkt)

udt_send(sndpkt)

rdt_rcv(rcvpkt) && isACK(rcvpkt)

Λ

rdt_rcv(rcvpkt) &&
corrupt(rcvpkt)

udt_send(NAK)

**Wait for call from below**

rdt_rcv(rcvpkt) &&
notcorrupt(rcvpkt)

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

# RDT2.0: Error Scenario

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

**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)

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

**Wait for call from below**

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

# RDT2.0: A Fatal Flaw

- *What if the ACK/NACK is corrupted?*

- Let's try a few options
    1. Ask receiver to resend
        - Lead to a sort of infinite loop + too much state information ("*I'm sending the ACK of the ACK of the NACK of the ACK that I sent an hour ago*")
    2. Just retransmit the packet (i.e. assume the worst)
        - Potentially create duplicate packets

- We'll pick option # 2
    - How to solve the duplication problem?
        - Use *sequence number* (i.e. packet ID)

- New mechanism in RDT2.1:
    - Sequence number

# Stop-and-Wait Protocol's Sequence Numbers

- RDT so far is a *stop and wait* protocol:
  - Sends a packet
  - Wait for ACK/NACK
  - Resend if ACK/NACK corrupted or NACK received
  - Only then, next packet can be sent

- We only need two sequence numbers: 0 and 1
  - 1-bit sequence number space

# RDT2.1: Sender Side

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

**Wait for call 0 from above**

**Wait for ACK or NAK 0**

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

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)

# RDT2.1: Receiver Side

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)

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)

**Wait for 0 from below**

**Wait for 1 from below**

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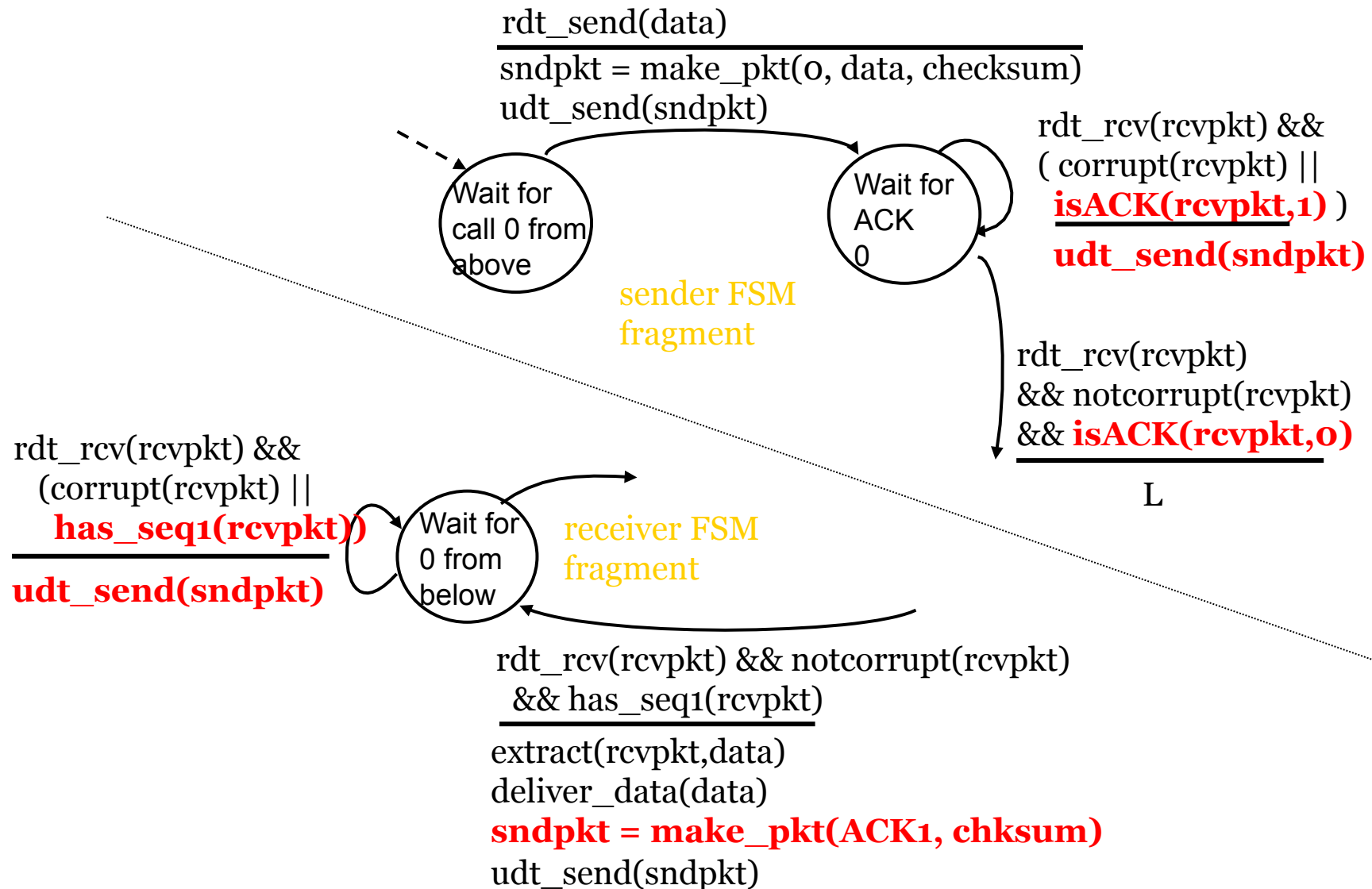deliver_data(data)
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

# RDT2.1: Observations

- **Senders and receivers need twice as many states!**
- **Receiver must check if a packet is a duplicate**

- **Can make it *NACK-free***
    - Specify in the ACK the sequence # it is acknowledging
    - If a packet is corrupted, send ACK of the previous sequence # again
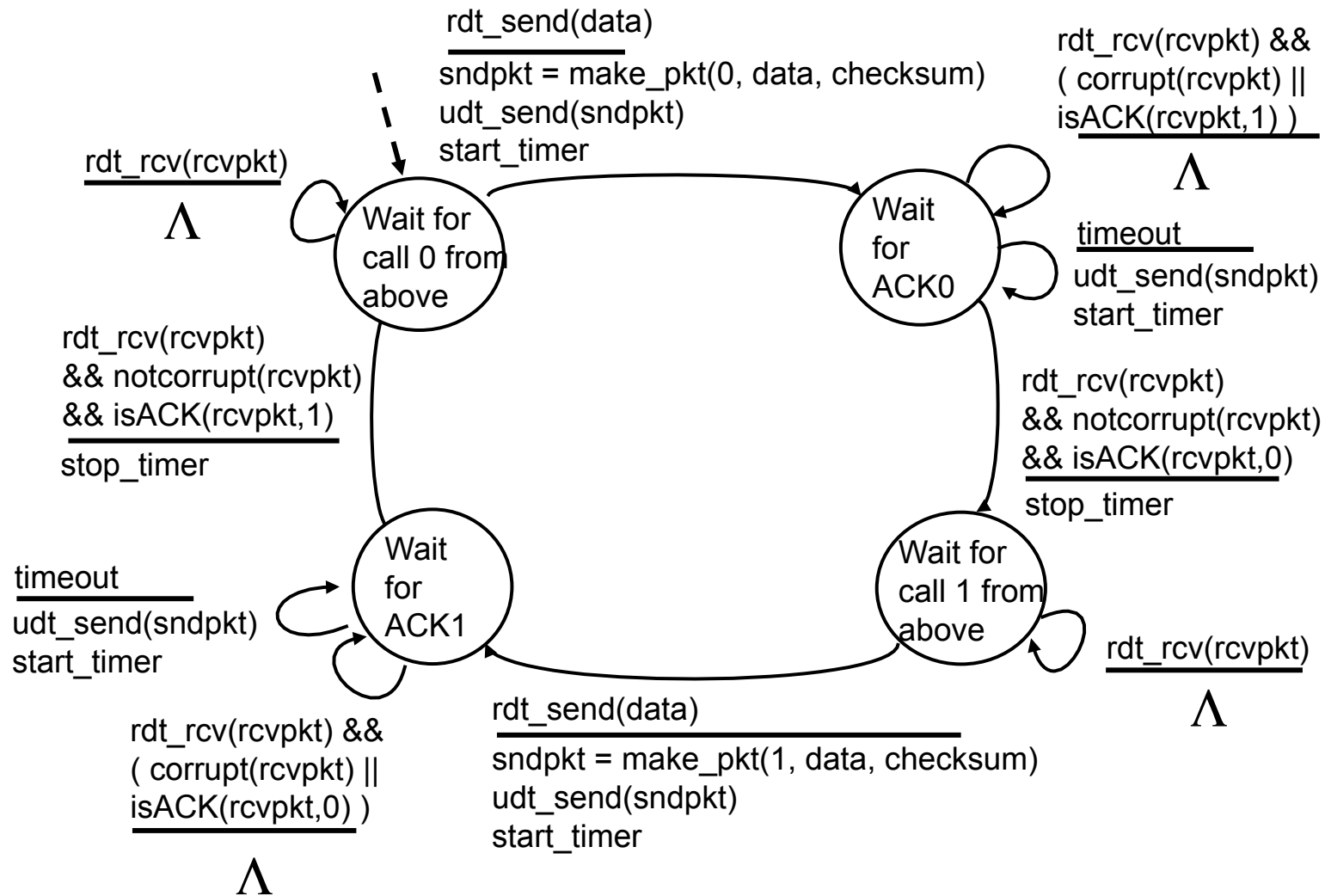    - This idea seems to appear "out of no where", will have utility later

# RDT2.2: the NACK-free Version of 2.1

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

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

**Wait for call 0 from above**

**Wait for ACK 0**

*sender FSM fragment*

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

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)
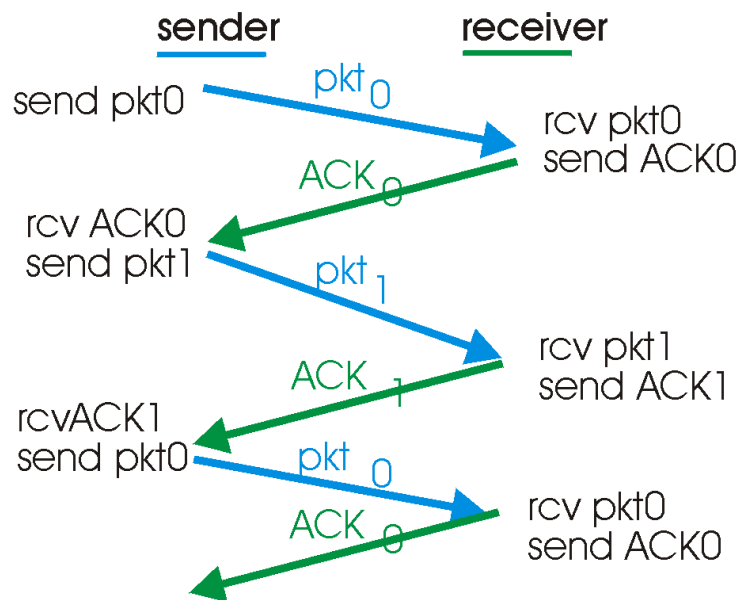**sndpkt = make_pkt(ACK1, chksum)**
udt_send(sndpkt)

# RDT3.0: Dealing with Packet Losses Too

- **Both data packets and ACK can be corrupted or lost**

- **How to detect that there's a lost packet?**
  - Use a timer
  - *But how long?*
    - Too long → inefficiency
    - Too short → duplicate packet (seq. no. took care of this!)
    - We'll get back to this issue later

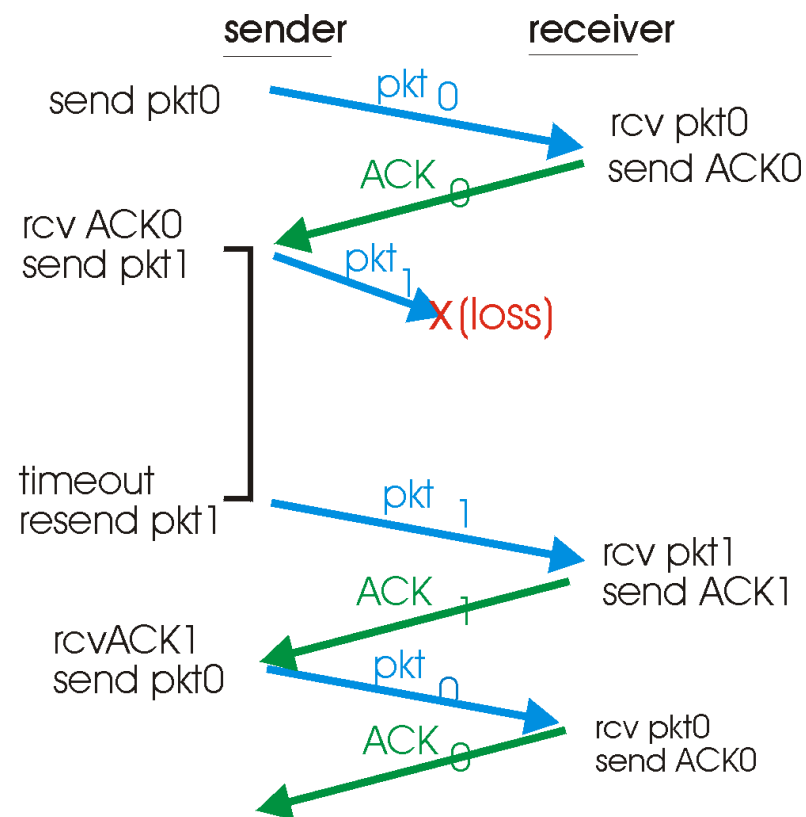- **What to do when a packet is lost? Well, retransmit!**

# RDT3.0: Sender Side

# RDT3.0 In Action



**sender**     **receiver**

send pkt0    pkt$_0$    rcv pkt0
send ACK0

ACK$_0$

rcv ACK0
send pkt1    pkt$_1$

rcv pkt1
send ACK1

ACK$_1$

rcvACK1
send pkt0    pkt$_0$

ACK$_0$    rcv pkt0
send ACK0

(a) operation with no loss

**sender**     **receiver**

send pkt0    pkt$_0$    rcv pkt0
send ACK0

ACK$_0$

rcv ACK0
send pkt1    pkt$_1$

X (loss)

timeout
resend pkt1    pkt$_1$

rcv pkt1
send ACK1

ACK$_1$

rcvACK1
send pkt0    pkt$_0$

ACK$_0$    rcv pkt0
send ACK0
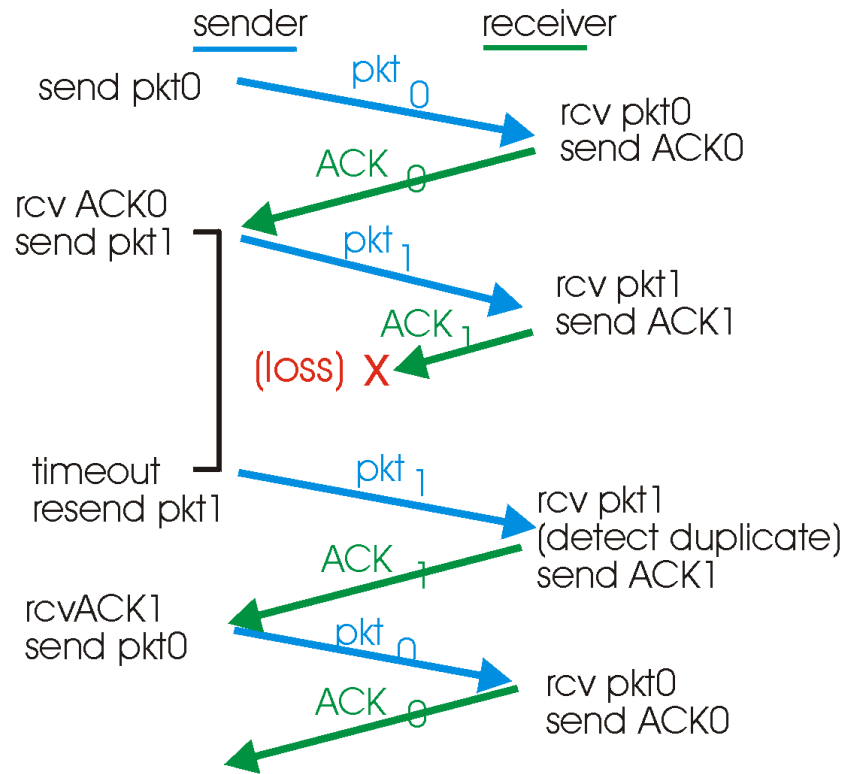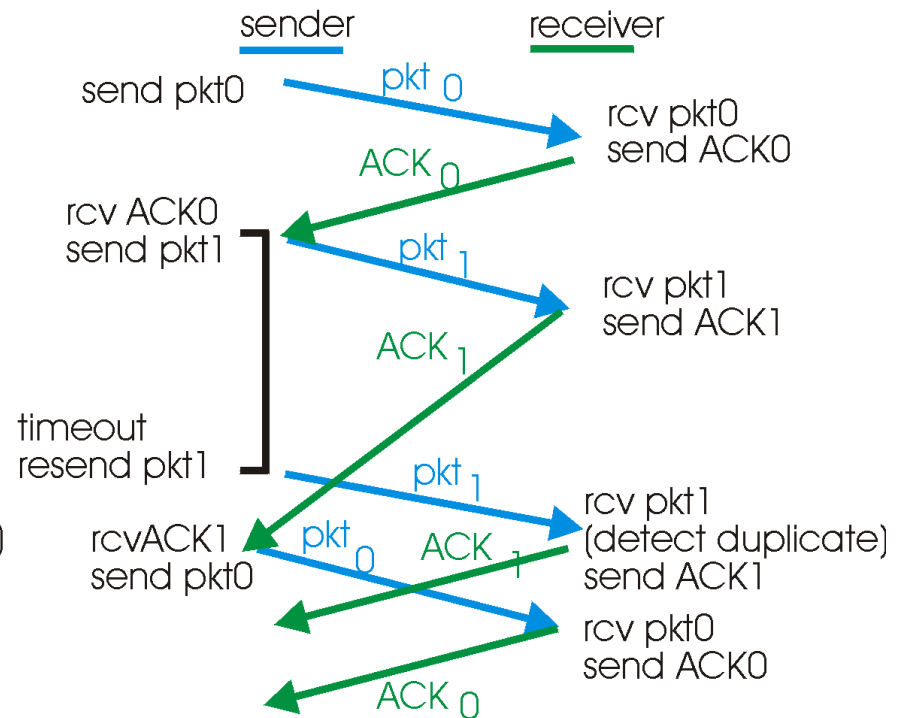
(b) lost packet

# RDT3.0 In Action



(c) lost ACK

(d) premature timeout

# Performance of RDT3.0

- RDT3.0 works, but performance stinks!
- Example: 1 Gbps link, 15 ms propagation delay, 8000 bit packet:
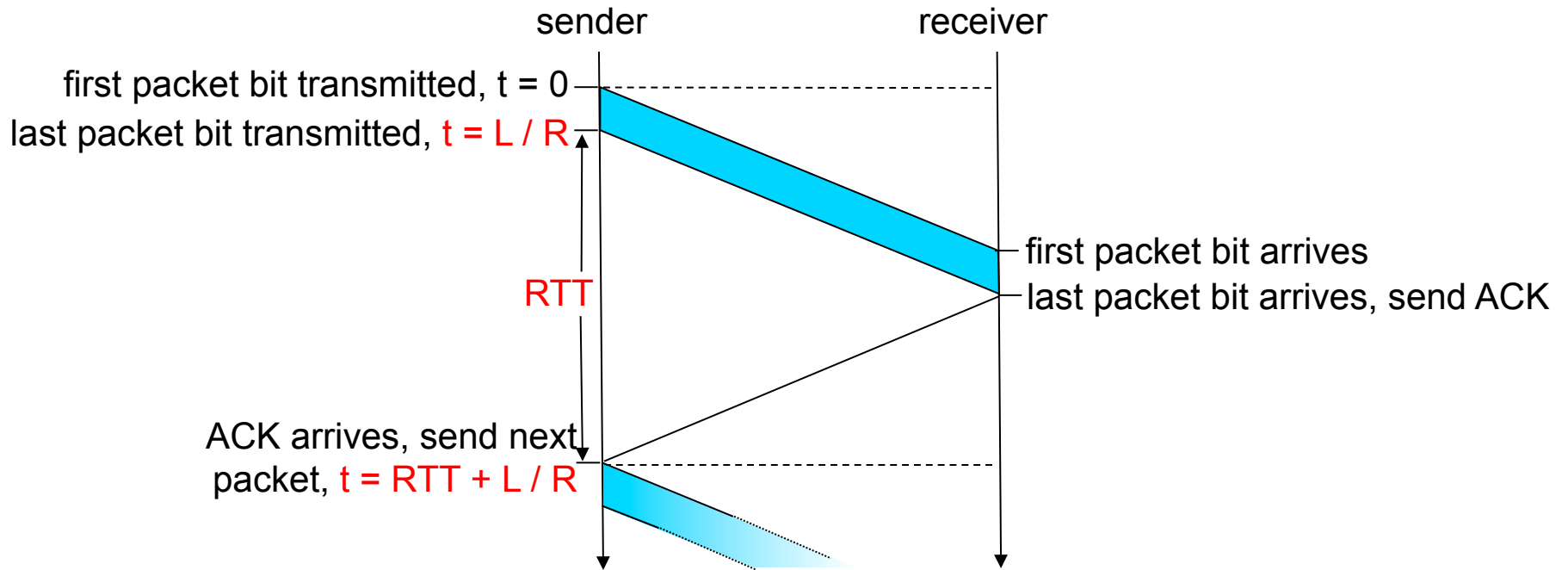
$$d_{trans} = \frac{L}{R} = \frac{8000\,\text{bits}}{10^9\,\text{bps}} = 8\,\text{microseconds}$$

- *Channel utilization* – fraction of time sender busy sending

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

- 1KB packet every 30 msec, i.e. 33kB/sec *throughput* over 1 Gbps link

- *Lesson: network protocol limits use of physical resources!*

# RDT3.0: Stop-and-Wait → Bad Utilization

sender                                    receiver

first packet bit transmitted, t = 0

last packet bit transmitted, t = L / R

first packet bit arrives

RTT                                       last packet bit arrives, send ACK

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

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

**Question**: *how to utilize the channel better?*

# Pipelined Protocols

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



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

Pipelining introduces new problems:

- 1-bit seq. # no longer works → enlarge seq. # space
- We have to deal with out-of-order packets now
- Larger buffers at senders and receivers

# Pipelining Helps Increase 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 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!

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

# Pipelining Protocols

| *Go-back-N: big picture* | *Selective Repeat: big picture* |
|---|---|
| ■ **Sender can have up to *N* unack'ed packets in pipeline**<br>　■ i.e. keep a *window* of size N<br><br>■ **Receiver only sends cumulative ACKs**<br>　■ *Drop any out-of-order packet*<br>　■ Re-ACK oldest received pkt<br><br>■ **Sender has a timer for *oldest* unacked packet**<br>　■ When timer expires, retransmit all unack'ed packets in the window | ■ **Sender can have up to *N* unack'ed packets in pipeline**<br>　■ i.e. keep a *window* of size N<br><br>■ **Receiver acknowledges individual packets**<br>　■ Including out-of-order packets within window<br><br>■ **Sender maintains a timer for *each* unacked packet**<br>　■ When any timer expires, retransmit that unacked packet |

**Question**: why limit the number of unack'ed packets to N?

# Go-Back-N

## Sender:

- Uses a $k$-bit sequence number in packet header
- Maintains a *window* of up to $N$, consecutive unack'ed packets allowed
- Also called *sliding window protocol*



- ACK(n): ACKs all pkts up to, including seq # n - "*cumulative ACK*"
  - may receive duplicate ACKs (see receiver)
- One timer for all in-flight packet
- *timeout(n):* retransmit pkt *n* and all higher seq # pkts in window

# Go-Back-N: Sender Extended FSM

rdt_send(data)
_____

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

$\Lambda$
_____
base=1
nextseqnum=1

**Wait**

rdt_rcv(rcvpkt)
  && corrupt(rcvpkt)
_____

timeout
_____
start_timer // and resend all unack'ed packets
udt_send(sndpkt[base])
udt_send(sndpkt[base+1])
…
udt_send(sndpkt[nextseqnum-1])

rdt_rcv(rcvpkt) &&
  notcorrupt(rcvpkt)
_____

base = getacknum(rcvpkt)+1
If (base == nextseqnum)
  stop_timer
 else
  start_timer // restart timer for new base

# Go-Back-N: Receiver Extended FSM

default
—————————
udt_send(sndpkt)

rdt_rcv(rcvpkt)
  && notcurrupt(rcvpkt)
  && hasseqnum(rcvpkt,expectedseqnum)
————————————————————————————
extract(rcvpkt,data)
deliver_data(data)
sndpkt = make_pkt(expectedseqnum,ACK,chksum)
udt_send(sndpkt)
expectedseqnum++

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

**Wait**

- ACK-only: always send ACK for correctly-received packet with highest *in-order* sequence number
  - may generate duplicate ACKs
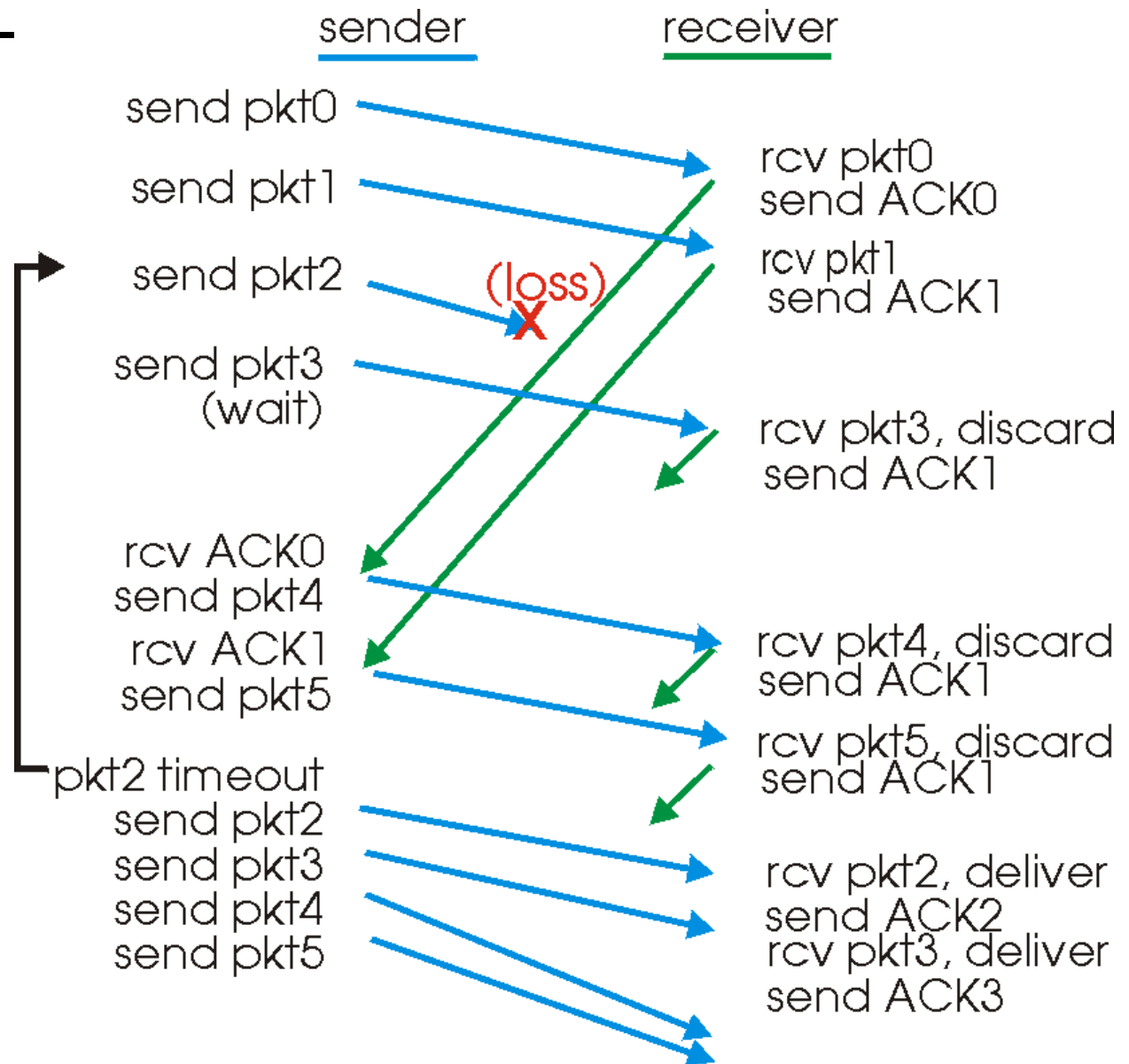  - need only remember `expectedseqnum`
- Out-of-order pkt:
  - discard (don't buffer) -> *no receiver buffering*!
  - Re-ACK packet with highest in-order sequence number

# Go-Back-N in Action



| sender | receiver |
|---|---|
| send pkt0 | rcv pkt0 send ACK0 |
| send pkt1 | rcv pkt1 send ACK1 |
| send pkt2 (loss) X | |
| send pkt3 (wait) | rcv pkt3, discard send ACK1 |
| rcv ACK0 send pkt4 | rcv pkt4, discard send ACK1 |
| rcv ACK1 send pkt5 | rcv pkt5, discard send ACK1 |
| pkt2 timeout send pkt2 send pkt3 send pkt4 send pkt5 | rcv pkt2, deliver send ACK2 rcv pkt3, deliver send ACK3 |

# Selective Repeat: Sender, Receiver Windows



send_base    nextseqnum

already ack'ed

sent, not yet ack'ed

usable, not yet sent

not usable

window size N

(a) sender view of sequence numbers

out of order (buffered) but already ack'ed

Expected, not yet received

acceptable (within window)

not usable

window size N

rcv_base

(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
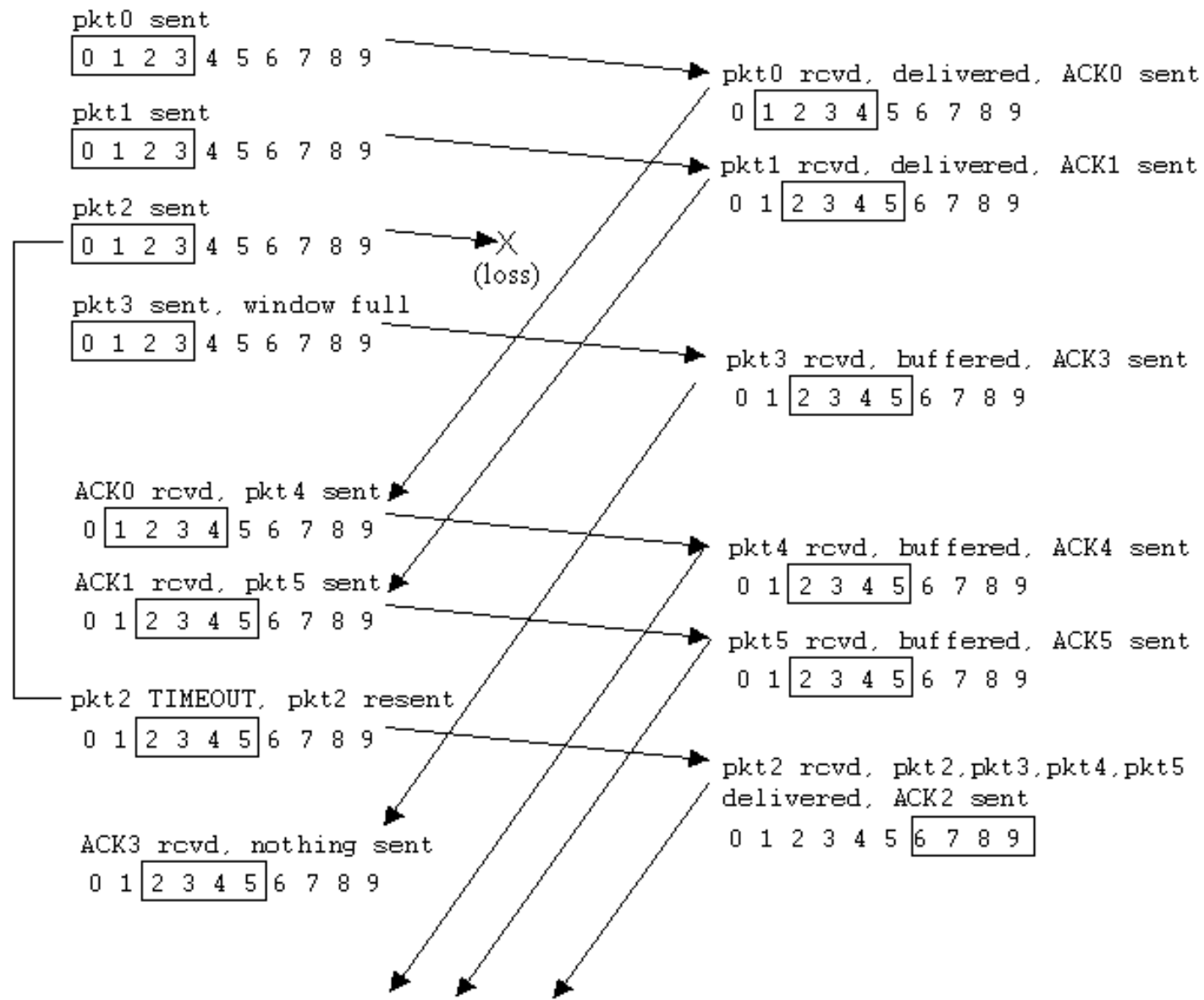
pkt n in [rcvbase-N,rcvbase-1]

- ACK(n)

otherwise:

- ignore

# Selective Repeat 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

✗ (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
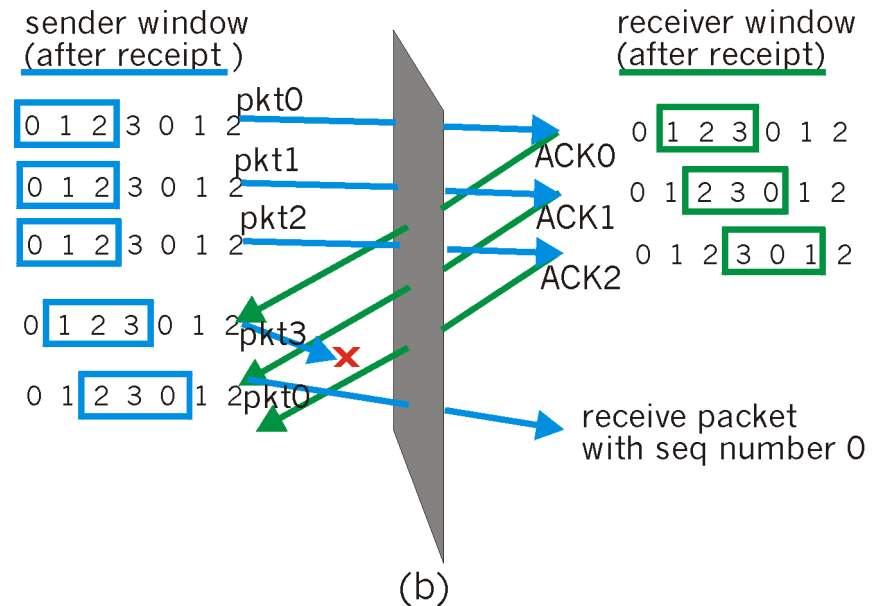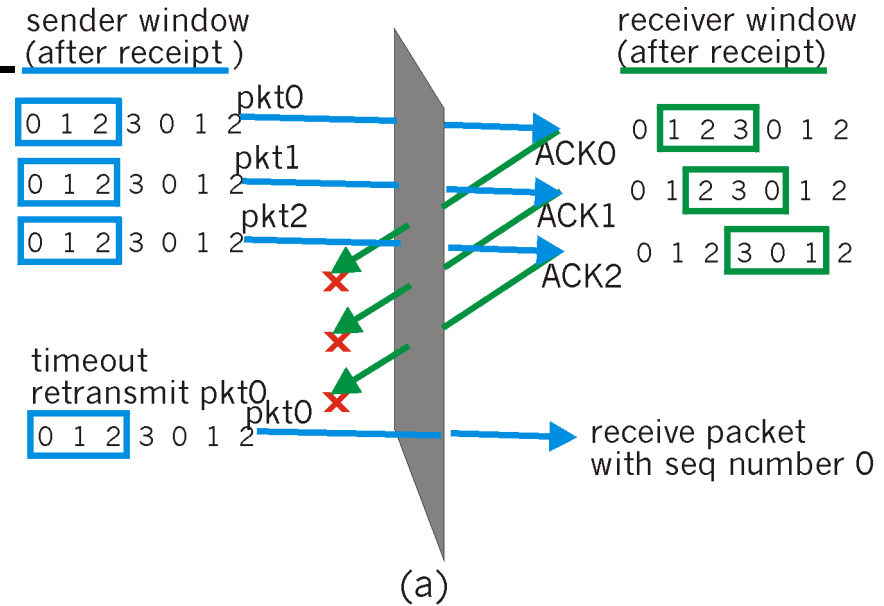0 1 2 3 4 5 | 6 7 8 9 |

# 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)

*Question:* what relationship between seq # size and window size must hold?

sender window (after receipt)    receiver window (after receipt)

0 1 2 3 0 1 2    pkt0    0 1 2 3 0 1 2
0 1 2 3 0 1 2    pkt1    ACK0    0 1 2 3 0 1 2
0 1 2 3 0 1 2    pkt2    ACK1    0 1 2 3 0 1 2
                          ACK2

timeout
retransmit pkt0    pkt0
0 1 2 3 0 1 2            receive packet
                        with seq number 0

(a)

sender window (after receipt)    receiver window (after receipt)

0 1 2 3 0 1 2    pkt0    0 1 2 3 0 1 2
0 1 2 3 0 1 2    pkt1    ACK0    0 1 2 3 0 1 2
0 1 2 3 0 1 2    pkt2    ACK1    0 1 2 3 0 1 2
                          ACK2
0 1 2 3 0 1 2    pkt3
0 1 2 3 0 1 2    pkt0    receive packet
                        with seq number 0

(b)

# Summary of Ideas We Have Learned

- Channel *bit errors* require
  - Error detecting codes
  - Receiver feedback (ACK/NAK)
  - Retransmissions (when NAK received or ACK/NAK corrupted)
- Retransmissions introduce *duplicates*
  - Need sequence numbers
- *Packet loss* requires
  - Timeout + retransmission (again introduce duplicates)
  - Estimating the "right" timeout is a fundamental problem!
- *Pipelining* improves utilization + throughput
  - Needs to enlarge sequence number space
  - Needs more buffer space at both sender & receiver
  - ACK + retransmission strategies: Go-Back-N & Selective Repeat
  - Window size & sequence number range strongly related
- *We have not discussed out-of-order, **long**-delayed packets*
- *And how long the timeout should be before retransmission*