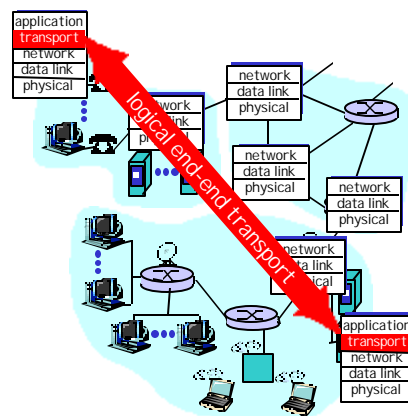# Chapter 3: Transport Layer

## Chapter goals:

r understand principles behind transport layer services:
  - m multiplexing/demultiplexing
  - m reliable data transfer
  - m flow control
  - m congestion control

r instantiation and implementation in the Internet

## Chapter Overview:

r transport layer services

r multiplexing/demultiplexing

r connectionless transport: UDP

r principles of reliable data transfer

r connection-oriented transport: TCP
  - m reliable transfer
  - m flow control
  - m connection management

r principles of congestion control

r TCP congestion control

3: Transport Layer    3a-1

---

# Transport services and protocols
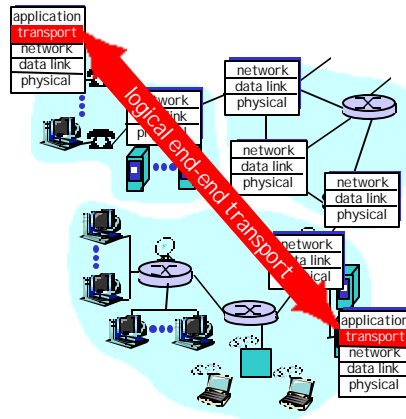
r provide *logical communication* between app' processes running on different hosts

r transport protocols run in end systems

r transport vs network layer services:

r *network layer:* data transfer between end systems

r *transport layer:* data transfer between processes
  - m relies on, enhances, network layer services

3: Transport Layer    3a-2

## Transport-layer protocols

Internet transport services:

r  reliable, in-order unicast delivery (TCP)
  m  congestion
  m  flow control
  m  connection setup
r  unreliable ("best-effort"), unordered unicast or multicast delivery: UDP
r  services not available:
  m  real-time
  m  bandwidth guarantees
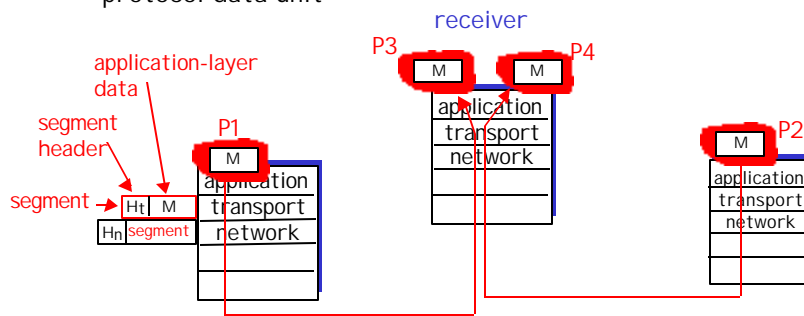  m  reliable multicast

3: Transport Layer    3a-3

# Multiplexing/demultiplexing

Recall: *segment* - unit of data exchanged between transport layer entities
  m  aka TPDU: transport protocol data unit

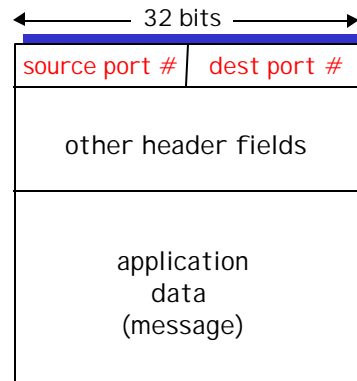Demultiplexing: delivering received segments to correct app layer processes

receiver

application-layer data

segment header

segment

3: Transport Layer    3a-4
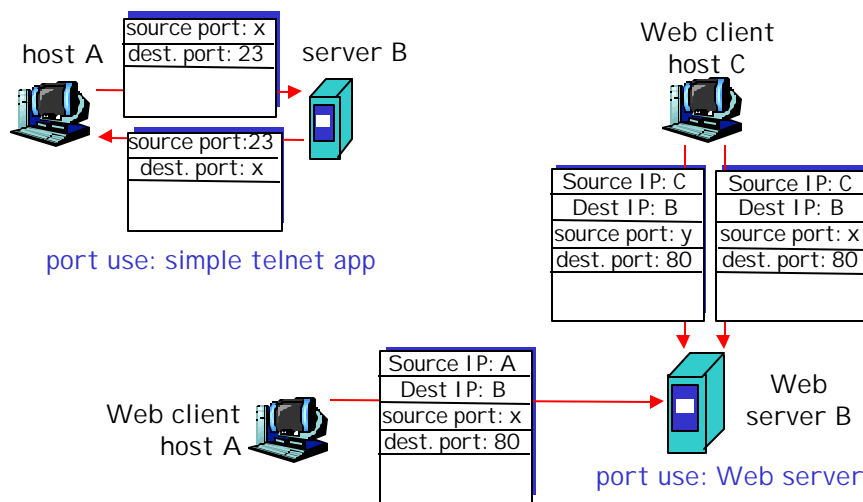
# Multiplexing/demultiplexing

**Multiplexing:**
gathering data from multiple app processes, enveloping data with header (later used for demultiplexing)

multiplexing/demultiplexing:
- r based on sender, receiver port numbers, IP addresses
  - m source, dest port #s in each segment
  - m recall: well-known port numbers for specific applications

← 32 bits →

| source port # | dest port # |
|---|---|

other header fields

application data (message)

TCP/UDP segment format

3: Transport Layer    3a-5

---

# Multiplexing/demultiplexing: examples

host A    source port: x / dest. port: 23    server B

source port:23 / dest. port: x

port use: simple telnet app

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

port use: Web server

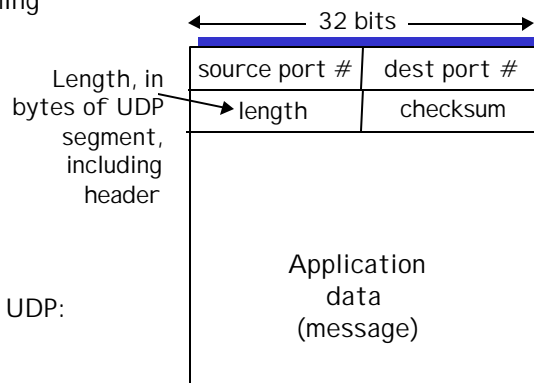3: Transport Layer    3a-6

# UDP: User Datagram Protocol [RFC 768]

r "no frills," "bare bones" Internet transport protocol

r "best effort" service, UDP segments may be:
- m lost
- m delivered out of order to app

r *connectionless:*
- m no handshaking between UDP sender, receiver
- m each UDP segment handled independently of others

## Why is there a UDP?

r no connection establishment (which can add delay)

r simple: no connection state at sender, receiver

r small segment header

r no congestion control: UDP can blast away as fast as desired

3: Transport Layer    3a-7

---

# UDP: more

r often used for streaming multimedia apps
- m loss tolerant
- m rate sensitive

r other UDP uses (why?):
- m DNS
- m SNMP

r reliable transfer over UDP: add reliability at application layer
- m application-specific error recover!

Length, in bytes of UDP segment, including header

```
  |<------------- 32 bits ------------->|
  |  source port #  |  dest port #  |
  |  length         |  checksum     |
  |                                 |
  |        Application              |
  |           data                  |
  |        (message)                |
  |                                 |
```

UDP segment format

3: Transport Layer    3a-8

# UDP checksum

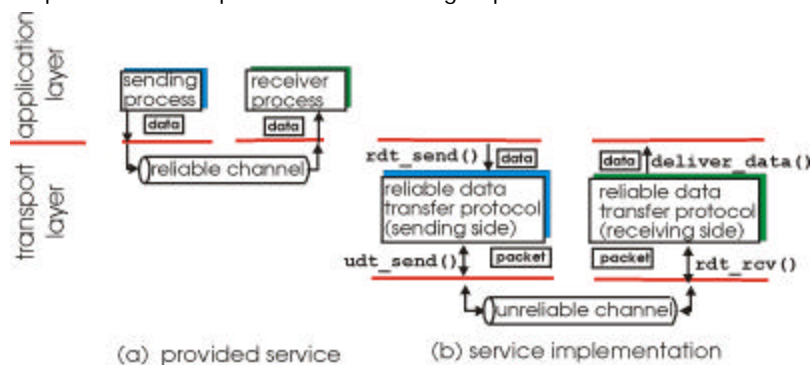Goal: detect "errors" (e.g., flipped bits) in transmitted segment

### Sender:

r  treat segment contents as sequence of 16-bit integers
r  checksum: addition (1's complement sum) of segment contents
r  sender puts checksum value into UDP checksum field

### Receiver:

r  compute checksum of received segment
r  check if computed checksum equals checksum field value:
  m  NO - error detected
  m  YES - no error detected. *But maybe errors nonethless?* More later … .

3: Transport Layer    3a-9

# Principles of Reliable data transfer

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



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

3: Transport Layer    3a-10

## 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

rdt_send() ↓ data

reliable data transfer protocol (sending side)

data ↑ deliver_data()

reliable data transfer protocol (receiving side)

receive 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

3: Transport Layer   3a-11

---

## Reliable data transfer: getting started

We'll:

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

r  consider only unidirectional data transfer
   m  but control info will flow on both directions!

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

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

3: Transport Layer   3a-12

## Rdt1.0: reliable transfer over a reliable channel

r  underlying channel perfectly reliable
  m  no bit erros
  m  no loss of packets
r  separate FSMs for sender, receiver:
  m  sender sends data into underlying channel
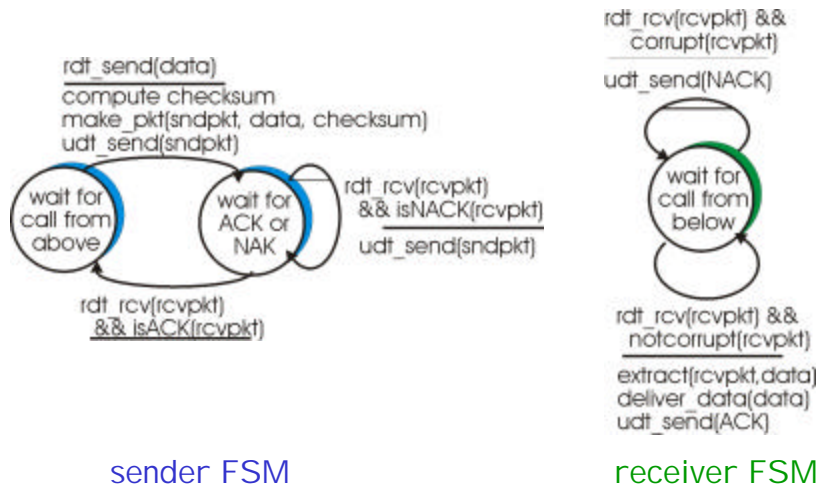  m  receiver read data from underlying channel

| wait for call from above | rdt_send(data) |
| --- | --- |
| | make_pkt(packet,data) |
| | udt_send(packet) |

| wait for call from below | rdt_rcv(packet) |
| --- | --- |
| | extract(packet,data) |
| | deliver_data(data) |

(a) rdt1.0: sending side          (b) rdt1.0: receiving side

3: Transport Layer   3a-13

## Rdt2.0: channel with bit errors

r  underlying channel may flip bits in packet
  m  recall: UDP checksum to detect bit errors
r  *the* question: how to recover from errors:
  m  *acknowledgements (ACKs):* receiver explicitly tells sender that pkt received OK
  m  *negative acknowledgements (NAKs):* receiver explicitly tells sender that pkt had errors
  m  sender retransmits pkt on receipt of NAK
  m  human scenarios using ACKs, NAKs?
r  new mechanisms in `rdt2.0` (beyond `rdt1.0`):
  m  error detection
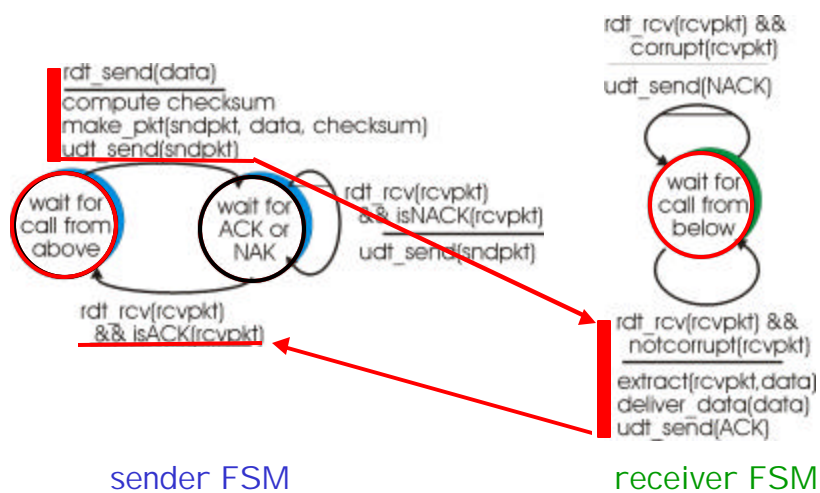  m  receiver feedback: control msgs (ACK,NAK) rcvr->sender

3: Transport Layer   3a-14

# rdt2.0: FSM specification
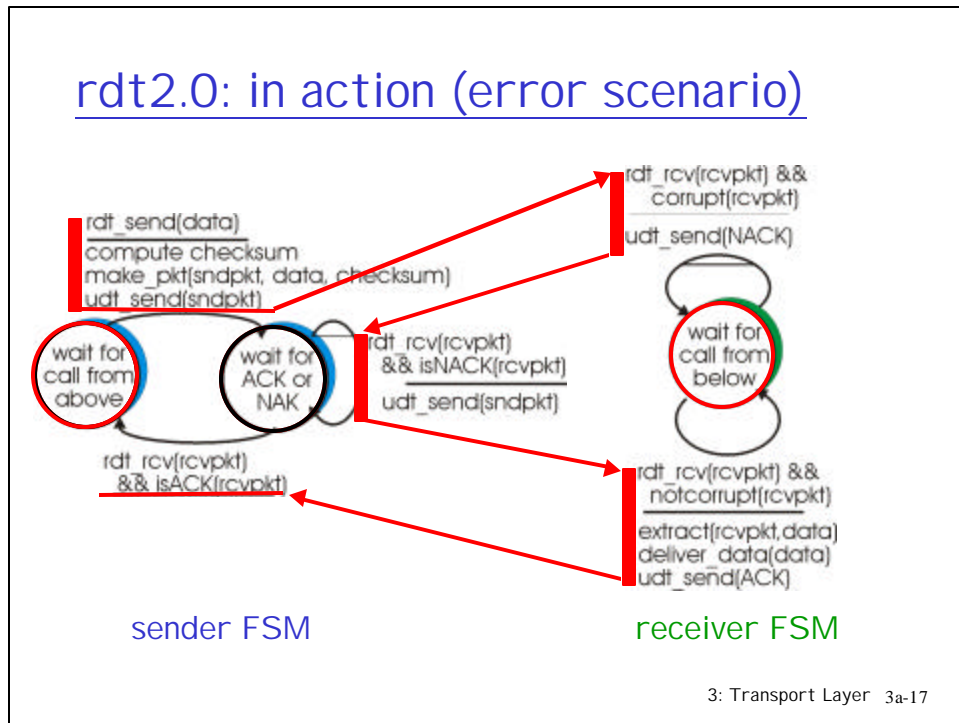
rdt_send(data)
compute checksum
make_pkt(sndpkt, data, checksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt)
&& isNACK(rcvpkt)

udt_send(sndpkt)

wait for
call from
above

wait for
ACK or
NAK

rdt_rcv(rcvpkt)
&& isACK(rcvpkt)

rdt_rcv(rcvpkt) &&
corrupt(rcvpkt)

udt_send(NACK)

wait for
call from
below

rdt_rcv(rcvpkt) &&
notcorrupt(rcvpkt)

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

sender FSM

receiver FSM

3: Transport Layer   3a-15

# rdt2.0: in action (no errors)

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

rdt_rcv(rcvpkt)
&& isNACK(rcvpkt)

udt_send(sndpkt)

wait for
call from
above

wait for
ACK or
NAK

rdt_rcv(rcvpkt)
&& isACK(rcvpkt)

rdt_rcv(rcvpkt) &&
corrupt(rcvpkt)

udt_send(NACK)

wait for
call from
below

rdt_rcv(rcvpkt) &&
notcorrupt(rcvpkt)

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

sender FSM

receiver FSM

3: Transport Layer   3a-16

# rdt2.0: in action (error scenario)



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

rdt_rcv(rcvpkt) &&
  corrupt(rcvpkt)
udt_send(NACK)

rdt_rcv(rcvpkt)
&& isNACK(rcvpkt)
udt_send(sndpkt)

wait for
call from
above

wait for
ACK or
NAK

wait for
call from
below

rdt_rcv(rcvpkt)
&& isACK(rcvpkt)

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

sender FSM                                    receiver FSM

3: Transport Layer   3a-17

# rdt2.0 has a fatal flaw!

**What happens if ACK/NAK corrupted?**
r sender doesn't know what happened at receiver!
r san't just retransmit: possible duplicate

**What to do?**
r sender ACKs/NAKs receiver's ACK/NAK? What if sender ACK/NAK lost?
r retransmit, but this might cause retransmission of correctly received pkt!

**Handling duplicates:**
r sender adds *sequence number* to each pkt
r sender retransmits current pkt if ACK/NAK garbled
r receiver discards (doesn't deliver up) duplicate pkt

**stop and wait**
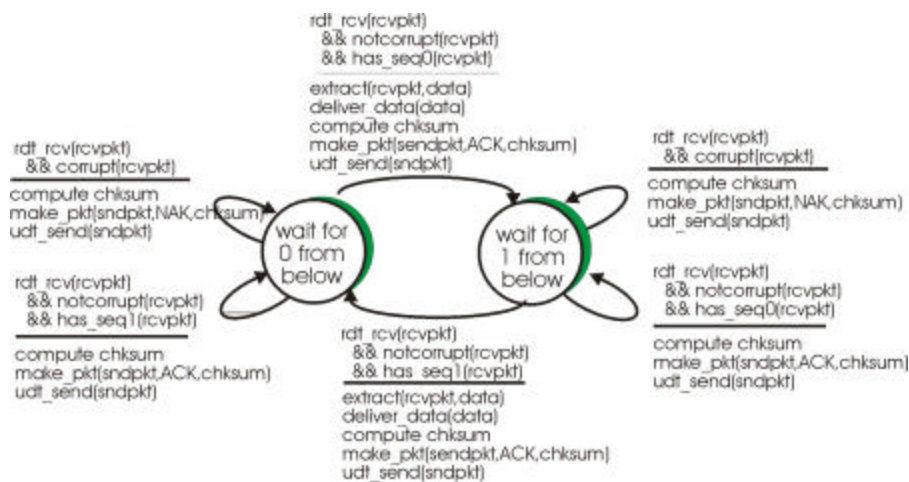Sender sends one packet, then waits for receiver response

3: Transport Layer   3a-18

## rdt2.1: sender, handles garbled ACK/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
call0 from
above

wait
ACK or
NAK
0

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

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

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

udt_send(sndpkt)

wait
ACK or
NAK
1

wait for
call1 from
above

rdt_send(data)

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

3: Transport Layer   3a-19

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

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

extract(rcvpkt,data)
deliver_data(data)
compute chksum
make_pkt(sendpkt,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
0 from
below

wait for
1 from
below

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

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)

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

extract(rcvpkt,data)
deliver_data(data)
compute chksum
make_pkt(sendpkt,ACK,chksum)
udt_send(sndpkt)

3: Transport Layer   3a-20

# rdt2.1: discussion

**Sender:**

- r seq # added to pkt
- r two seq. #'s (0,1) will suffice. Why?
- r must check if received ACK/NAK corrupted
- r twice as many states
  - m state must "remember" whether "current" pkt has 0 or 1 seq. #

**Receiver:**

- r must check if received packet is duplicate
  - m state indicates whether 0 or 1 is expected pkt seq #
- r note: receiver can *not* know if its last ACK/NAK received OK at sender

3: Transport Layer  3a-21

---

# rdt2.2: a NAK-free protocol

- r same functionality as rdt2.1, using NAKs only
- r instead of NAK, receiver sends ACK for last pkt received OK
  - m receiver must *explicitly* include seq # of pkt being ACKed
- r duplicate ACK at sender results in same action as NAK: *retransmit current pkt*

sender FSM



3: Transport Layer  3a-22

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

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

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

Q: how to deal with loss?

   m sender waits until certain data or ACK lost, then retransmits

   m yuck: drawbacks?

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

r retransmits if no ACK received in this time

r if pkt (or ACK) just delayed (not lost):

   m retransmission will be duplicate, but use of seq. #'s already handles this

   m receiver must specify seq # of pkt being ACKed

r requires countdown timer

3: Transport Layer   3a-23

---

# rdt3.0 sender



3: Transport Layer   3a-24

# rdt3.0 in action



(a) operation with no loss

(b) lost packet

3: Transport Layer   3a-25

# rdt3.0 in action



(c) lost ACK

(d) premature timeout

3: Transport Layer   3a-26

# Performance of rdt3.0

r  rdt3.0 works, but performance stinks

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

$$T_{transmit} = \frac{8kb/pkt}{10^{**}9 \ b/sec} = 8 \ microsec$$

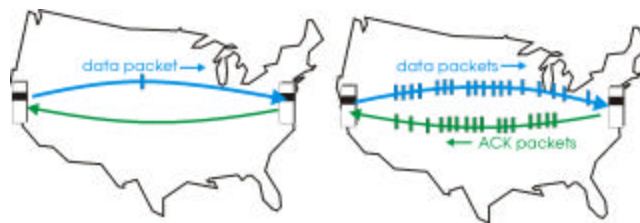$$\text{Utilization} = U = \frac{\text{fraction of time}}{\text{sender busy sending}} = \frac{8 \ microsec}{30.016 \ msec} = 0.00015$$

   m  1KB pkt every 30 msec -> 33kB/sec thruput over 1 Gbps link
   m  network protocol limits use of physical resources!

3: Transport Layer  3a-27

# Pipelined protocols

Pipelining: sender allows multiple, "in-flight", yet-to-be-acknowledged pkts
   m  range of sequence numbers must be increased
   m  buffering at sender and/or receiver



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

r  Two generic forms of pipelined protocols: *go-Back-N, selective repeat*

3: Transport Layer  3a-28

# Go-Back-N

Sender:

r  k-bit seq # in pkt header

r  "window" of up to N, consecutive unack'ed pkts allowed

send_base        nextseqnum

already ack'ed          usable, not yet sent

sent, not yet ack'ed          not usable

window size
N

r  ACK(n): ACKs all pkts up to, including seq # n - "cumulative ACK"

   m  may deceive duplicate ACKs (see receiver)

r  timer for each in-flight pkt

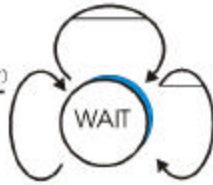r  *timeout(n):* retransmit pkt n and all higher seq # pkts in window

3: Transport Layer  3a-29

# GBN: 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 = nextseqnum + 1
    }
else
    refuse_data(data)
```

rdt_rcv(rcv_pkt) && notcorrupt(rcvpkt)

base = getacknum(rvcpkt)+1
if (base == nextseqnum)
    stop_timer
else
    start_timer

WAIT

timeout

start_timer
udt_send(sndpkt(base))
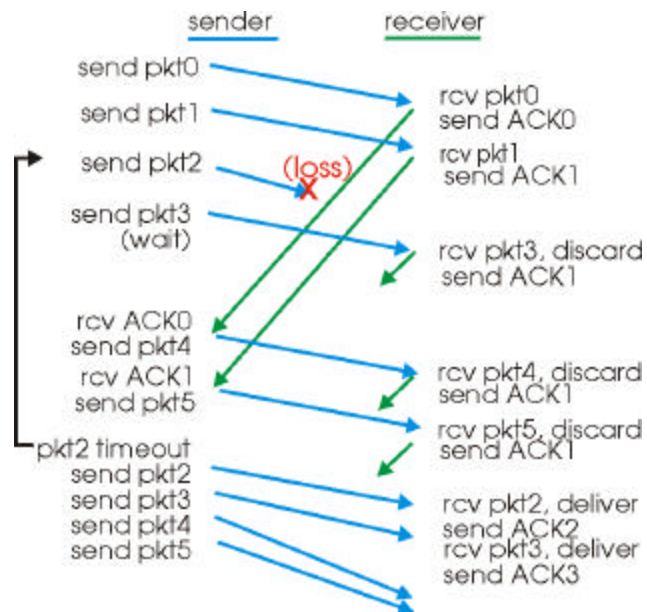udt_send(sndpkt(base+1))
. . . . .
udt_send(sndpkt(nextseqnum-1))

3: Transport Layer  3a-30

## GBN: receiver extended FSM



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

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

default
‾‾‾‾‾‾‾‾‾‾‾‾‾‾
udt_send(sndpkt)

WAIT

receiver simple:

r  ACK-only: always send ACK for correctly-received pkt with highest *in-order* seq #
  m  may generate duplicate ACKs
  m  need only remember `expectedseqnum`
r  out-of-order pkt:
  m  discard (don't buffer) -> no receiver buffering!
  m  ACK pkt with highest in-order seq #

3: Transport Layer   3a-31
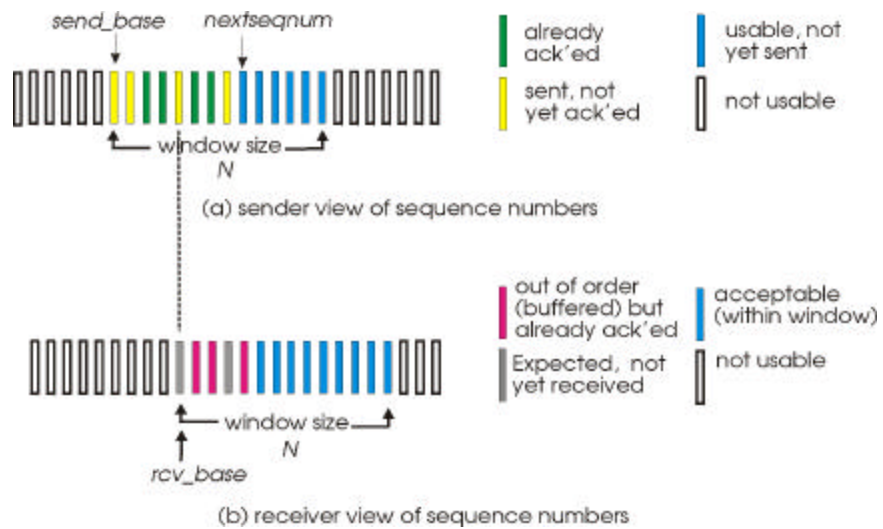
## GBN in action



3: Transport Layer   3a-32

# Selective Repeat

r receiver *individually* acknowledges all correctly received pkts
  m buffers pkts, as needed, for eventual in-order delivery to upper layer
r sender only resends pkts for which ACK not received
  m sender timer for each unACKed pkt
r sender window
  m N consecutive seq #'s
  m again limits seq #s of sent, unACKed pkts

3: Transport Layer 3a-33

# Selective repeat: sender, receiver windows



(a) sender view of sequence numbers

(b) receiver view of sequence numbers

3: Transport Layer 3a-34

# Selective repeat

### sender

**data from above :**
r if next available seq # in window, send pkt

**timeout(n):**
r resend pkt n, restart timer

**ACK(n)** in [sendbase,sendbase+N]:
r mark pkt n as received
r if n smallest unACKed pkt, advance window base to next unACKed seq #

### receiver

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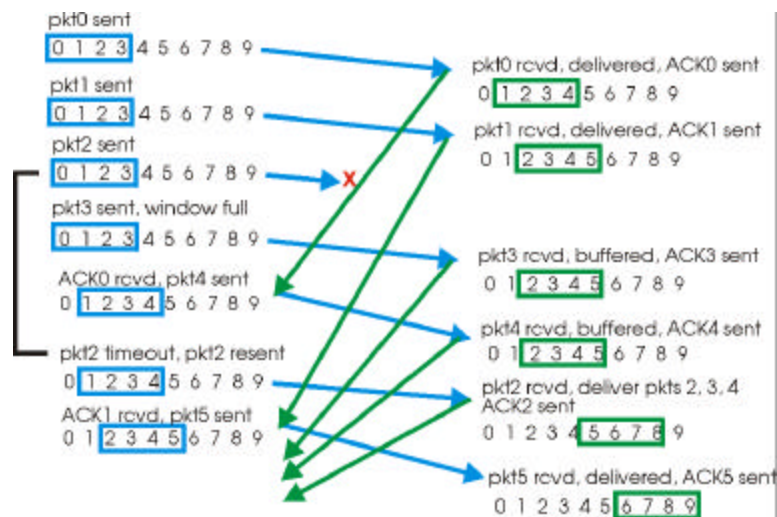r in-order: deliver (also deliver buffered, in-order pkts), advance window to next not-yet-received pkt

**pkt n in** [rcvbase-N,rcvbase-1]
r ACK(n)

**otherwise:**
r ignore

3: Transport Layer   3a-35

# Selective repeat in action



3: Transport Layer   3a-36

## Selective repeat: dilemma

Example:

r   seq #'s: 0, 1, 2, 3

r   window size=3


r   receiver sees no difference in two scenarios!

r   incorrectly passes duplicate data as new in (a)


Q: what relationship between seq # size and window size?



3: Transport Layer   3a-37