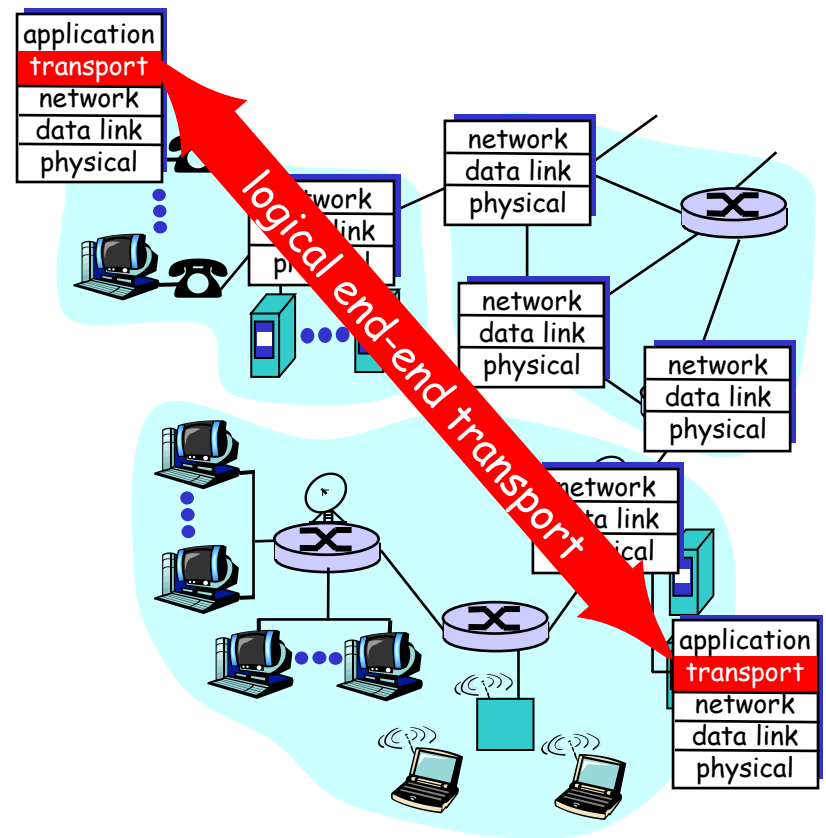# Chapter 3: Transport Layer

## Chapter goals:

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

□ instantiation and implementation in the Internet

## Chapter Overview:

□ transport layer services

□ multiplexing/demultiplexing

□ connectionless transport: UDP

□ principles of reliable data transfer

□ connection-oriented transport: TCP
- ○ reliable transfer
- ○ flow control
- ○ connection management

□ principles of congestion control

□ TCP congestion control

# Transport services and protocols

- provide *logical communication* between app' processes running on different hosts
- transport protocols run in end systems
- transport vs network layer services:
- *network layer:* data transfer between end systems
- *transport layer:* data transfer between processes
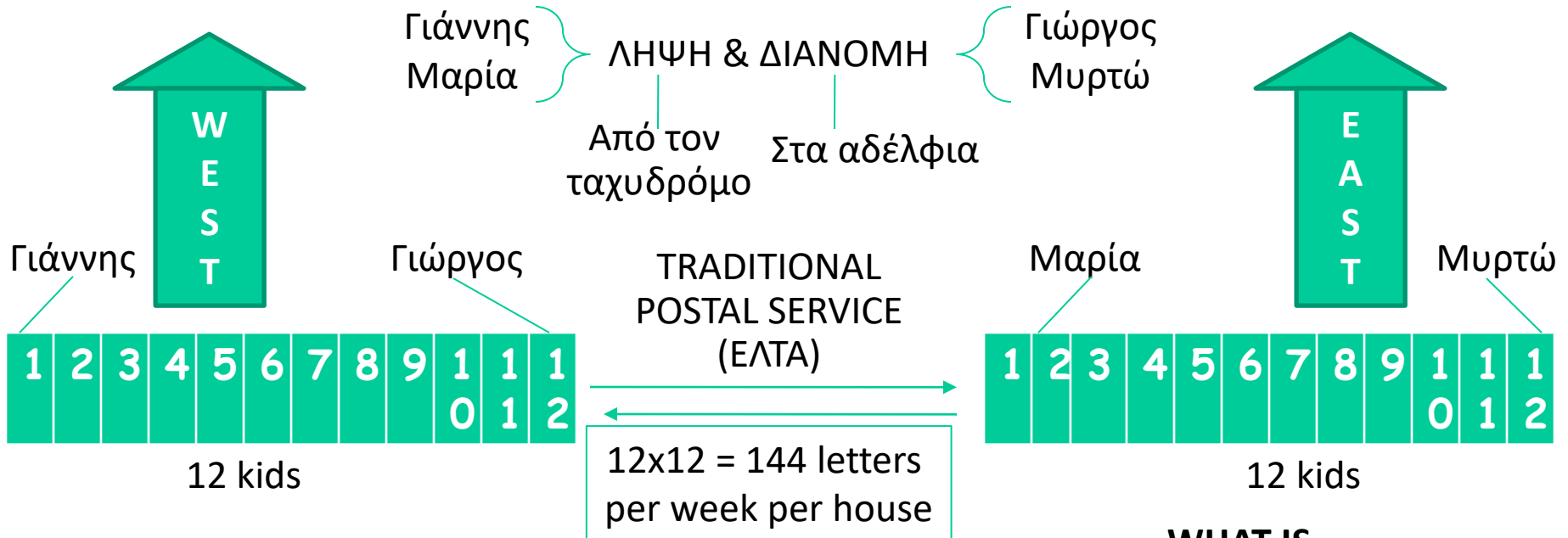  - relies on, enhances, network layer services

# Introduction to Transport Layer Protocols

□ Transport Layer Protocols provide logical communication between processes running on different hosts.

□ Όταν σε έναν host καταφθάνουν πακέτα από διάφορες εφαρμογές (application processes), το πρωτόκολλο του Στρώματος Μεταφοράς είναι αυτό που ξεδιαλύνει σε ποια εφαρμογή θα οδηγηθεί κάθε πακέτο.

□ Στο Διαδίκτυο έχουμε δύο πρωτόκολλα στο Στρώμα Μεταφοράς, το TCP (με αξιόπιστη λειτουργία) και το UDP (που θυσιάζει την αξιόπιστη λειτουργία χάριν της ταχύτητας – μεταφορά πληροφορίας χωρίς καθυστερήσεις).

□ Ακολούθως:

Ανθρώπινο ανάλογο που εξηγεί την λειτουργία των πρωτοκόλλων του Στρώματος Μεταφοράς και την σχέση τους με τα πρωτόκολλα του Στρώματος Δικτύου.

# Ανθρώπινο ανάλογο

Δυό σπίτια, WEST – EAST, σε μια μεγαλούπολη με πολυμελείς οικογένειες - εξαδέλφια

Γιάννης
Μαρία

ΛΗΨΗ & ΔΙΑΝΟΜΗ

Γιώργος
Μυρτώ

Από τον ταχυδρόμο

Στα αδέλφια

**W E S T**

**E A S T**

Γιάννης

Γιώργος

TRADITIONAL POSTAL SERVICE (ΕΛΤΑ)

Μαρία

Μυρτώ

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |

12 kids

12x12 = 144 letters per week per house

12 kids

**WHAT IS**

Houses – West & East

Kids

Letters (in the envelops)

ΕΛΤΑ - Ταχυδρόμοι

Γιάννης – Μαρία

Γιώργος – Μυρτώ

**WHAT IS**

Hosts – End systems

Application Processes

Messages

Network Layer Protocol

Transport Layer Protocol – TCP

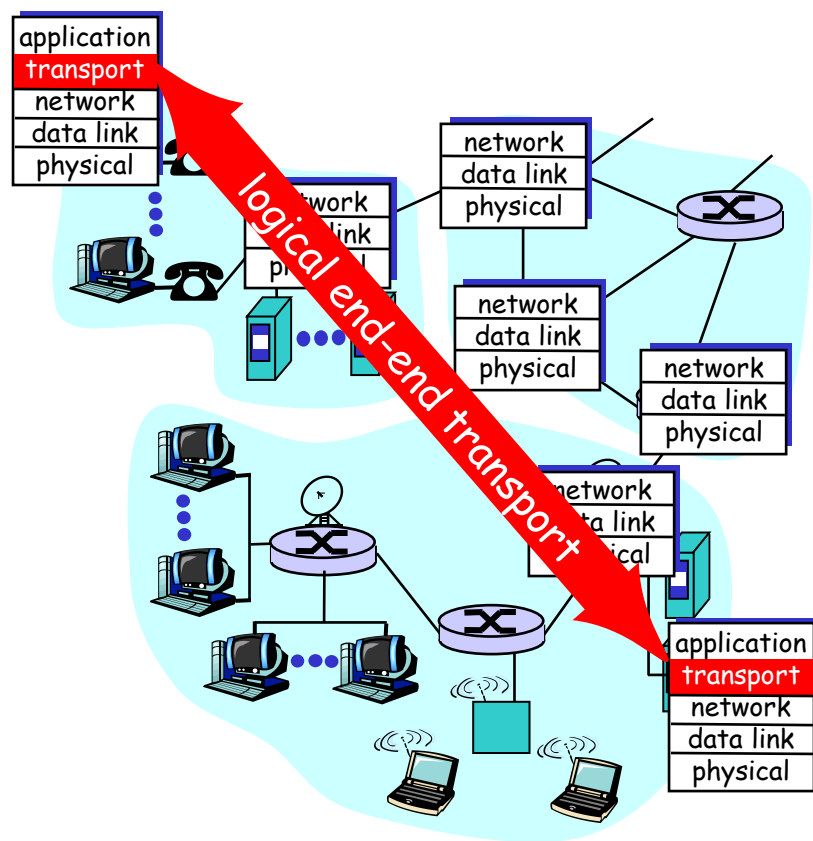Transport Layer Protocol – UDP

# Ανθρώπινο ανάλογο    (συνέχεια)

□ The possible services that Γιάννης & Μαρία can provide are clearly constrained by the possible services that the ΤΑΧΥΔΡΟΜΕΙΟ provides.

□ Αν το ΤΑΧΥΔΡΟΜΕΙΟ δεν εγγυηθεί max delay of (π.χ.) 3 days, ούτε ο Γιάννης & η Μαρία μπορούν να εγγυηθούν πότε θα παραδώσουν γράμματα στα αδέλφια τους.

------------------------------------------------------------------------

□ If the network layer protocol cannot provide delay or bandwidth guarantees for the segments (received from the transport layer and) sent to the hosts, thus

□ the transport layer protocol cannot provide delay or bandwidth guarantees for the messages sent between application processes.

# Transport-layer protocols

**Internet transport services:**

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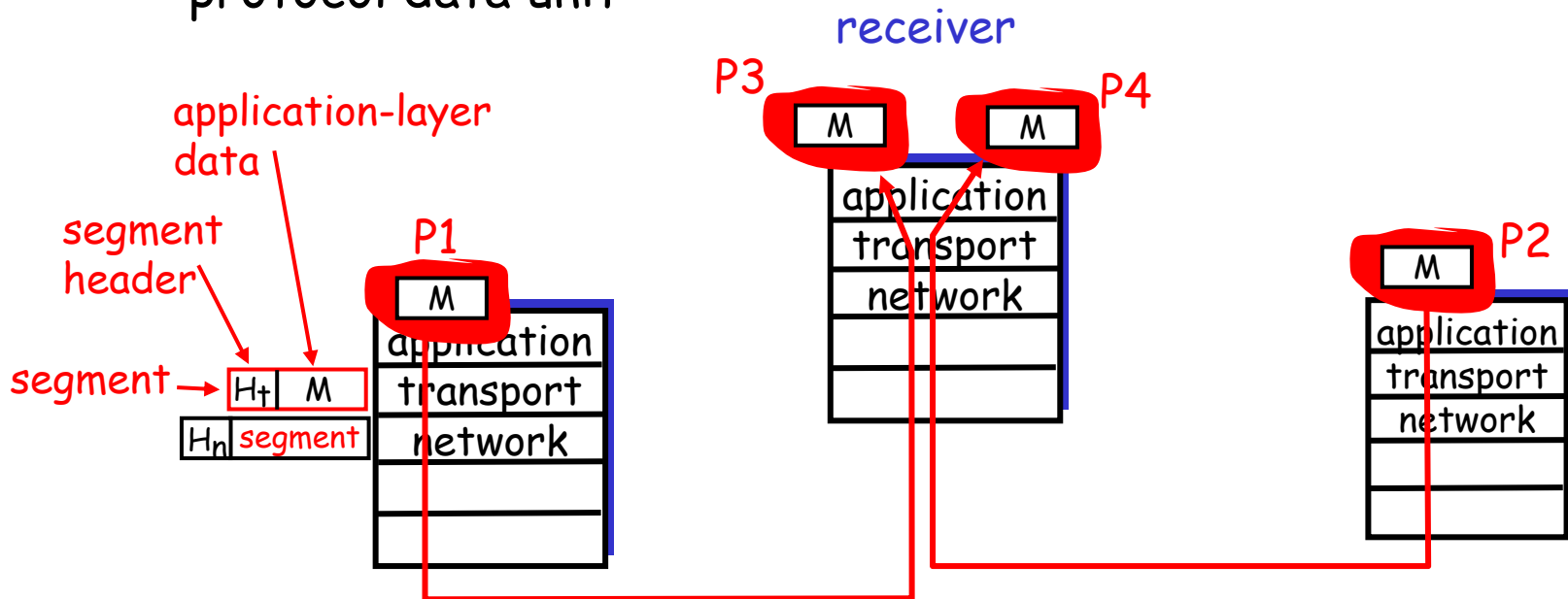  - reliable multicast



logical end-end transport

# Multiplexing/demultiplexing

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

Demultiplexing: delivering received segments to correct app layer processes

# Multiplexing/demultiplexing

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

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

← 32 bits →

| source port # | dest port # |
|---|---|
| other header fields | |
| application data (message) | |

TCP/UDP segment format

# 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

# UDP: User Datagram Protocol [RFC 768]

- "no frills," "bare bones" Internet transport protocol
- "best effort" service, UDP segments may be:
  - lost
  - delivered out of order to app
- *connectionless:*
  - no handshaking between UDP sender, receiver
  - each UDP segment handled independently of others

**Why is there a 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

# UDP: more

- often used for streaming multimedia apps
  - loss tolerant
  - rate sensitive
- **other UDP uses (why?):**
  - DNS
  - SNMP
- reliable transfer over UDP: add reliability at application layer
  - 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

# UDP checksum

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

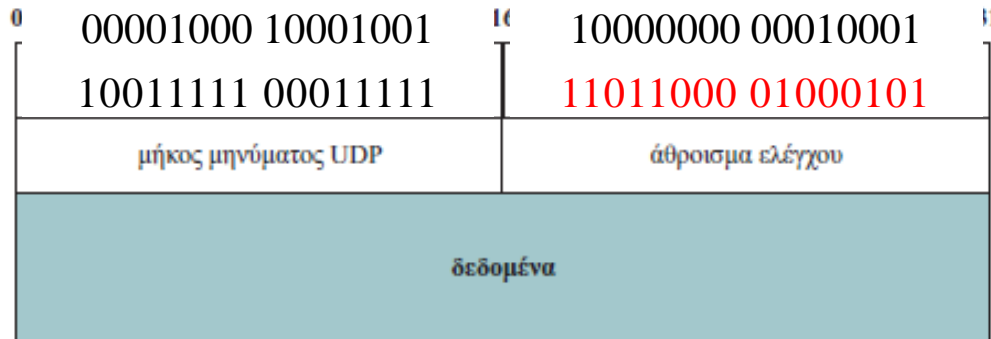**Sender:**
- treat segment contents as sequence of 16-bit integers
- checksum: addition (1's complement sum) of segment contents
- sender puts checksum value into UDP checksum field

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

# Checksum Example

00001000 10001001          10000000 00010001

10011111 00011111          11011000 01000101

μήκος μηνύματος UDP          άθροισμα ελέγχου

δεδομένα

  0000 1000 1000 1001
+ 1000 0000 0001 0001
  ───────────────────────
  10001000 1001 1010
+ 1001 1111 0001 1111
  ───────────────────────
  0010 0111 1011 1001
  + κρατούμενο      1
  0010 0111 1011 1010

Στον πομπό παίρνουμε το συμπλήρωμα  του 1:

**1101 1000 0100 0101**

Στον δέκτη:

  0000 1000 1000 1001
+ 1000 0000 0001 0001
  ───────────────────────
  10001000 1001 1010
+ 1001 1111 0001 1111
  ───────────────────────
  0010 0111 1011 1001
  + κρατούμενο      1
  0010 0111 1011 1010
+ 1101 1000 0100 0101

  **1111 1111 1111 1111**          Κανένα λάθος στη λήψη!

# Principles of Reliable data transfer

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



(a) provided service    (b) service implementation

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

# Reliable data transfer: getting started

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

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

send side

reliable data transfer protocol (sending side)

rdt_send() ↓ data

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

# Reliable data transfer: getting started

We'll:

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

☐ consider only unidirectional data transfer
  ○ but control info will flow on both directions!

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

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

event causing state transition
actions taken on state transition

state 1

event
actions

state 2

# Rdt1.0: reliable transfer over a reliable channel

□ **underlying channel perfectly reliable**
  ○ no bit erros
  ○ no loss of packets
□ **separate FSMs for sender, receiver:**
  ○ sender sends data into underlying channel
  ○ 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

# Rdt2.0: channel with bit errors

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

# rdt2.0: FSM specification

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

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

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

# rdt2.0: in action (no errors)

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

wait for
call from
above

wait for
ACK or
NAK

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

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

# rdt2.0: in action (error scenario)

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

wait for call from above

wait for ACK or NAK

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

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

# rdt2.0 has a fatal flaw!

## What happens if ACK/NAK corrupted?

- sender doesn't know what happened at receiver!
- can't just retransmit: possible duplicate

## What to do?

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

## Handling duplicates:

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

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

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

rdt_send(data)
_____

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

( wait for call0 from above )

( wait ACK or NAK 0 )

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

udt_send(sndpkt)

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

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

( wait ACK or NAK 1 )

( wait for call1 from above )

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

udt_send(sndpkt)

rdt_send(data)
_____

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

# 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)

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

compute chksum
make_pkt(sndpkt,ACK,chksum)
udt_send(sndpkt)

wait for
0 from
below

wait for
1 from
below

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)
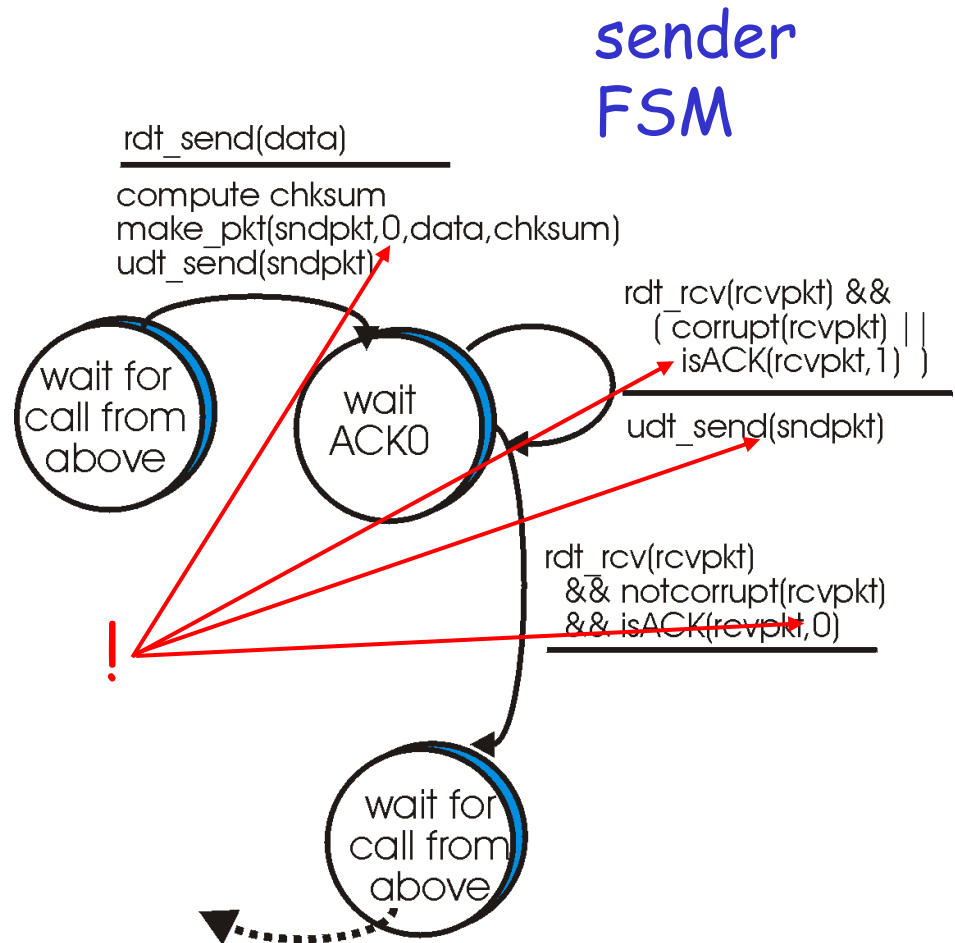udt_send(sndpkt)

# rdt2.1: discussion

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

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

# rdt2.2: a NAK-free protocol

- [ ] same functionality as rdt2.1, using NAKs 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*



sender FSM

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

wait for call from above

wait ACK0

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

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

wait for call from above

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

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

- checksum, seq. #, ACKs, retransmissions will be of help, but not enough
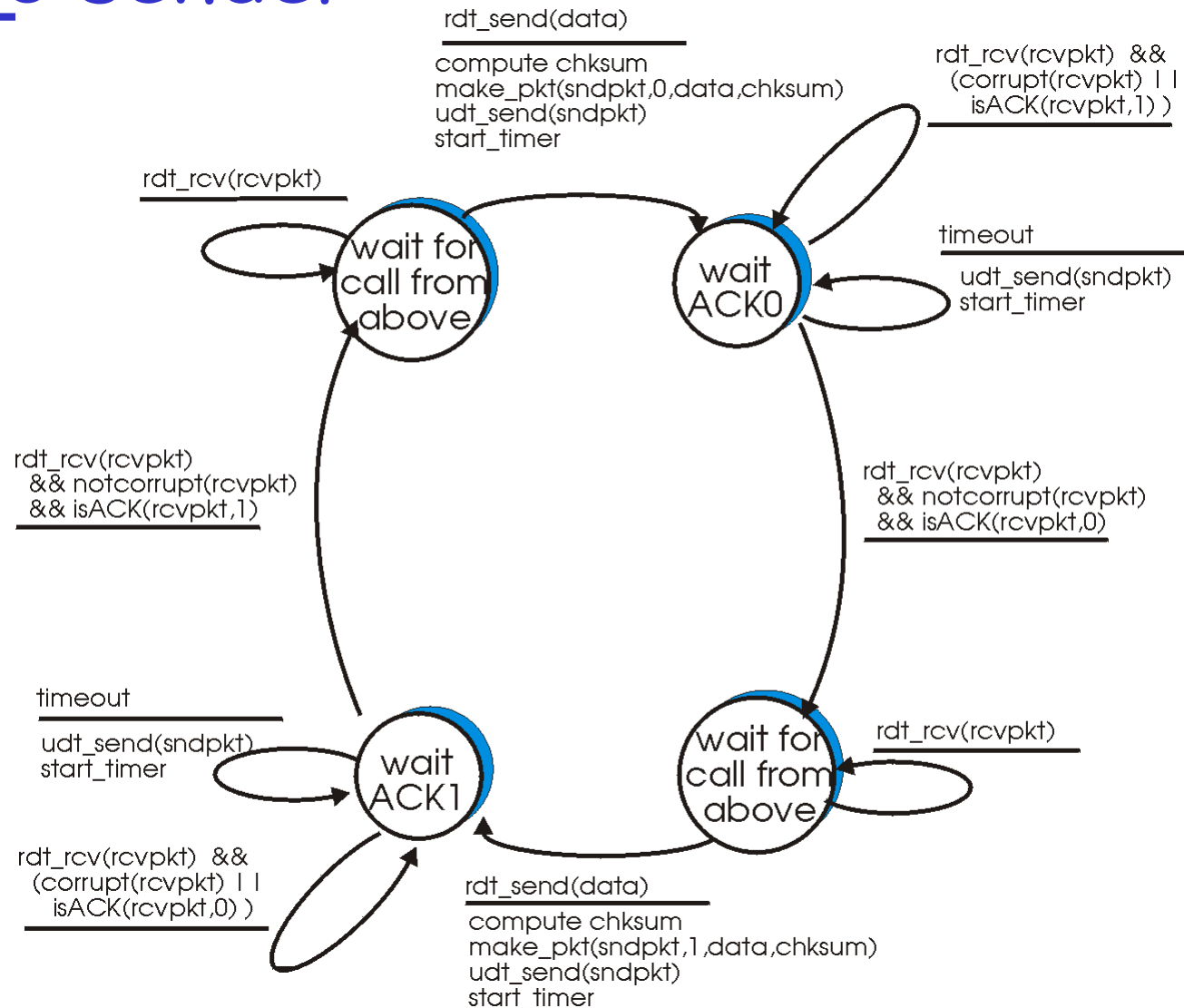
Q: how to deal with loss?

- sender waits until certain data or ACK lost, then retransmits
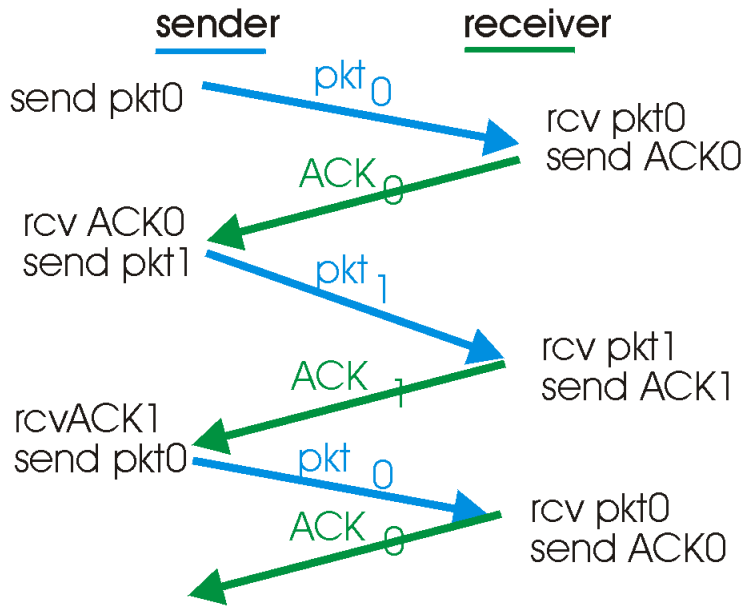- yuck: drawbacks?

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 use of seq. #'s already handles this
  - receiver must specify seq # of pkt being ACKed
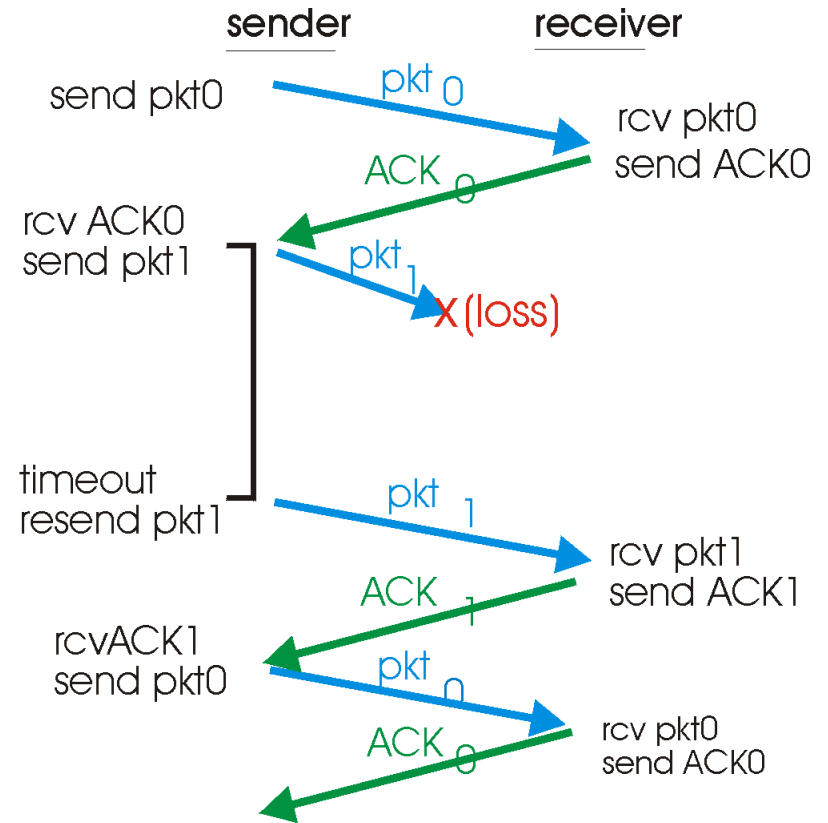- requires countdown timer

# rdt3.0 sender

rdt_send(data)
_____

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

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

rdt_rcv(rcvpkt)
_____

**wait for call from above**

**wait ACK0**

timeout
_____
udt_send(sndpkt)
start_timer

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

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

timeout
_____
udt_send(sndpkt)
start_timer

**wait ACK1**

**wait for call from above**

rdt_rcv(rcvpkt)
_____

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

rdt_send(data)
_____
compute chksum
make_pkt(sndpkt,1,data,chksum)
udt_send(sndpkt)
start_timer

# rdt3.0 in action

**sender**     **receiver**

send pkt0 → $pkt_0$ → rcv pkt0
send ACK0

rcv ACK0 ← $ACK_0$
send pkt1 → $pkt_1$ → rcv pkt1
send ACK1

rcvACK1 ← $ACK_1$
send pkt0 → $pkt_0$ → rcv pkt0
send ACK0

← $ACK_0$

## (a) operation with no loss

**sender**     **receiver**

send pkt0 → $pkt_0$ → rcv pkt0
send ACK0

rcv ACK0 ← $ACK_0$
send pkt1 → $pkt_1$ → X (loss)

timeout
resend pkt1 → $pkt_1$ → rcv pkt1
send ACK1

rcvACK1 ← $ACK_1$
send pkt0 → $pkt_0$ → rcv pkt0
send ACK0

← $ACK_0$

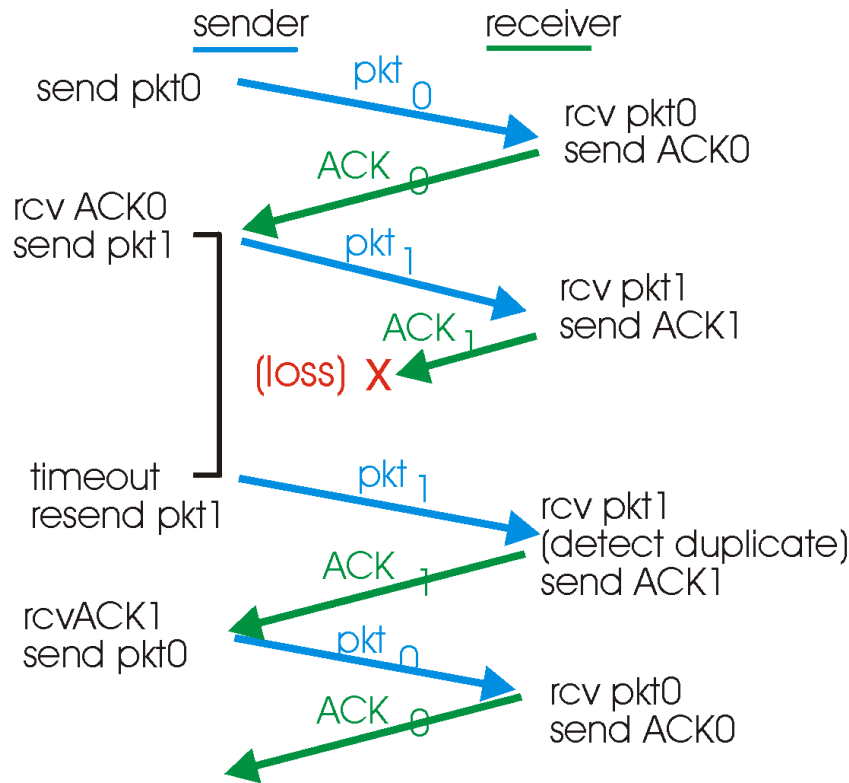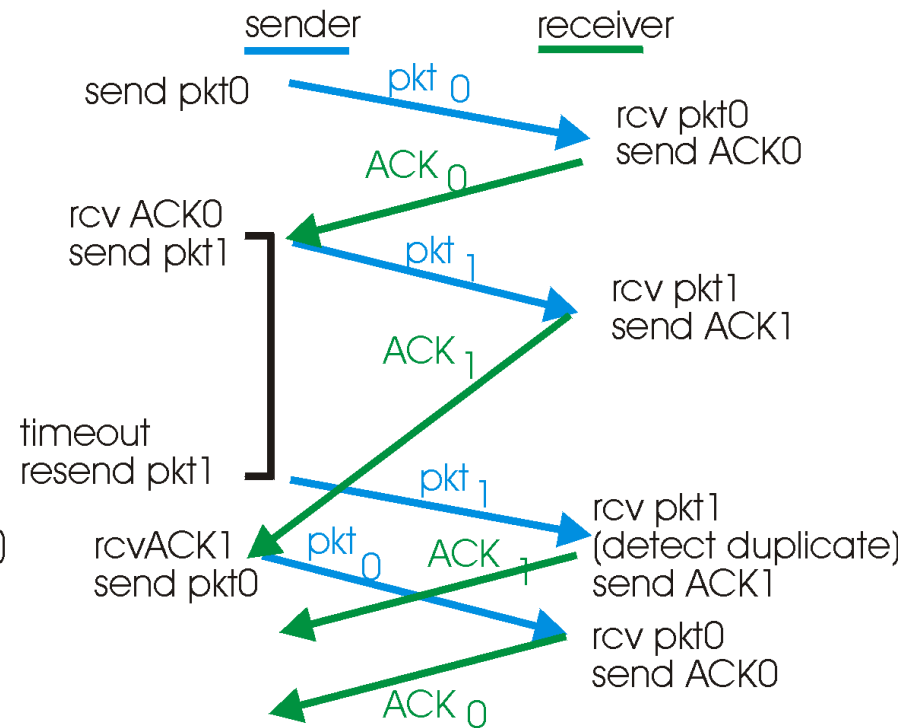## (b) lost packet

# 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 e-e prop. delay, 1KB packet:

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

Utilization = U = $\frac{fraction \ of \ time}{sender \ busy \ sending} = \frac{8 \ microsec}{30.008 \ msec} = 0.00027$
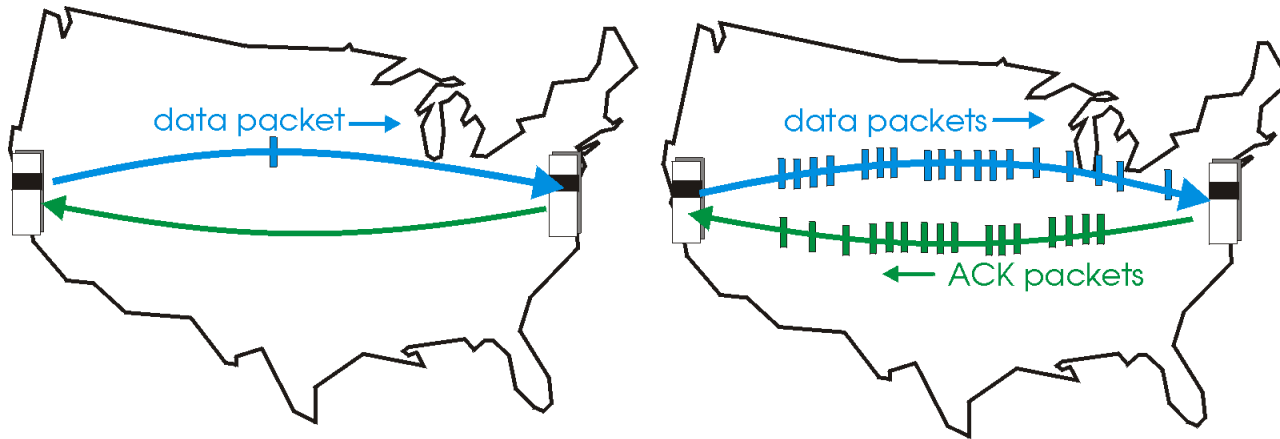
30.008 msec = total propagation delay +  transfer delay,
when ignoring the transmission delay of ACK

○ 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*

# Go-Back-N

Sender:

☐  k-bit *seq # in pkt header*

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



☐  ACK(n): ACKs all pkts up to, including seq # n - "cumulative ACK"
  ○  may deceive duplicate ACKs (see receiver)

☐  timer for each in-flight pkt

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

# 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))
```

# GBN: receiver extended FSM

```
                    rdt_rcv(rcvpkt) &&
                     notcorrupt(rcvpkt) &&
                     hasseqnum(rcvpkt,expectedseqnum)
default             _____
_____          extract(rcvpkt,data)
udt_send(sndpkt)     deliver_data(data)
                     make_pkt(sndpkt,ACK,expectedseqnum)
                     udt_send(sndpkt)
```

WAIT

## receiver simple:

☐ 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!
  ○ ACK pkt with highest in-order seq #

# GBN in action

# 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
  ○ again limits seq #s of sent, unACKed pkts

# Selective repeat: sender, receiver windows



(a) sender view of sequence numbers

- green: already ack'ed
- yellow: sent, not yet ack'ed
- blue: usable, not yet sent
- white: not usable

(b) receiver view of sequence numbers

- magenta: out of order (buffered) but already ack'ed
- gray: Expected, not yet received
- blue: acceptable (within window)
- white: not usable

# 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

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

pkt2 timeout, pkt2 resent
0 `1 2 3 4` 5 6 7 8 9

ACK1 rcvd, pkt5 sent
0 1 `2 3 4 5` 6 7 8 9

**X**

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

pkt2 rcvd, deliver pkts 2, 3, 4
ACK2 sent
0 1 2 3 4 `5 6 7 8` 9

pkt5 rcvd, delivered, ACK5 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)

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



sender window (after receipt)
0 1 2 3 0 1 2 — pkt0
0 1 2 3 0 1 2 — pkt1
0 1 2 3 0 1 2 — pkt2

receiver window (after receipt)
0 1 2 3 0 1 2 — ACK0
0 1 2 3 0 1 2 — ACK1
0 1 2 3 0 1 2 — ACK2

timeout
retransmit pkt0
0 1 2 3 0 1 2 — pkt0

receive packet with seq number 0

(a)

sender window (after receipt)
0 1 2 3 0 1 2 — pkt0
0 1 2 3 0 1 2 — pkt1
0 1 2 3 0 1 2 — pkt2

receiver window (after receipt)
0 1 2 3 0 1 2 — ACK0
0 1 2 3 0 1 2 — ACK1
0 1 2 3 0 1 2 — ACK2

0 1 2 3 0 1 2 — pkt3
0 1 2 3 0 1 2 — pkt0

receive packet with seq number 0

(b)