# Transport Layer

**Our goals:**

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

- learn about transport layer protocols in the Internet:
  - ❖ UDP: connectionless transport
  - ❖ TCP: connection-oriented transport
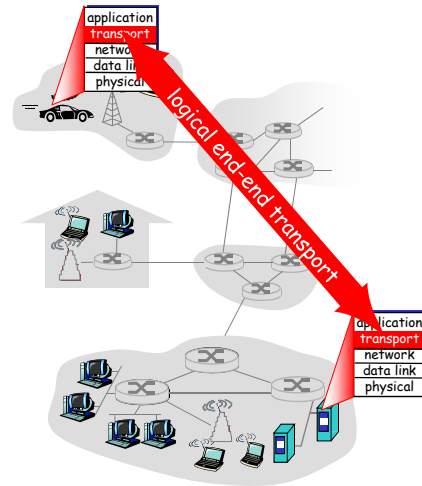  - ❖ TCP congestion control

# Outline

1

# Transport services and protocols

- provide *logical communication* between app processes running on different hosts
- transport protocols run in end systems
  - send side: breaks app messages into segments, passes to network layer
  - rcv side: reassembles segments into messages, passes to app layer
- more than one transport protocol available to apps
  - Internet: TCP and UDP

---

# Transport vs. network layer

- *network layer:* logical communication between hosts
- *transport layer:* logical communication between processes
  - relies on, enhances, network layer services

Household analogy:

*3 kids sending letters to 3 other kids*

- processes = kids
- app messages = letters in envelopes
- hosts = houses
- transport protocol = parents
- network-layer protocol = postal service

# Internet transport-layer protocols

- reliable, in-order delivery (TCP)
  - congestion control
  - flow control
  - connection setup
- unreliable, unordered delivery: UDP
  - no-frills extension of "best-effort" IP
- services not available:
  - delay guarantees
  - bandwidth guarantees

---

# Outline

# Multiplexing/demultiplexing

delivering received segments to correct socket

Multiplexing at send host:

gathering data from multiple sockets, enveloping data with header (later used for demultiplexing)

▮ = socket    ⬭ = process

| application P3 | P1 application P2 | P4 application |
|---|---|---|
| transport | transport | transport |
| network | network | network |
| link | link | link |
| physical | physical | physical |

host 1          host 2          host 3

# How demultiplexing works

❑ host receives IP datagrams
  ❖ each datagram has source IP address, destination IP address
  ❖ each datagram carries 1 transport-layer segment
  ❖ each segment has source, destination port number
❑ host uses IP addresses & port numbers to direct segment to appropriate socket

← 32 bits →

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

TCP/UDP segment format

# Connectionless demultiplexing

❑ Create sockets with port numbers:

```
DatagramSocket mySocket1 = new
    DatagramSocket(12534);
DatagramSocket mySocket2 = new
    DatagramSocket(12535);
```

❑ UDP socket identified by two-tuple:

(dest IP address, dest port number)

❑ When host receives UDP segment:
  ❖ checks destination port number in segment
  ❖ directs UDP segment to socket with that port number

❑ IP datagrams with different source IP addresses and/or source port numbers can be directed to same socket

# Connectionless demux (cont)

```
DatagramSocket serverSocket = new DatagramSocket(6428);
```



| SP: 6428 | | SP: 6428 |
| DP: 9157 | | DP: 5775 |

| SP: 9157 | | SP: 5775 |
| DP: 6428 | | DP: 6428 |

client
IP: A

server
IP: C

Client
IP:B

SP provides "return address"

# Connection-oriented demux

- ❑ TCP socket identified by 4-tuple:
  - ❖ source IP address
  - ❖ source port number
  - ❖ dest IP address
  - ❖ dest port number
- ❑ receiving host uses all four values to direct segment to appropriate socket

- ❑ Server host may support many simultaneous TCP sockets:
  - ❖ each socket identified by its own 4-tuple
- ❑ Web servers have different sockets for each connecting client
  - ❖ non-persistent HTTP will have different socket for each request

---

# Connection-oriented demux (cont)

P1

P4  P5  P6

P2  P3

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

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

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

client
IP: A

server
IP: C

Client
IP:B

# Connection-oriented demux: Threaded Web Server



| P1 | | | P4 | | | P2 | P3 |

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

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

client
IP: A

server
IP: C

Client
IP:B

---

# Outline

- ❑ 3.1 Transport-layer services
- ❑ 3.2 Multiplexing and demultiplexing
- ❑ 3.3 Connectionless transport: UDP
- ❑ 3.4 Principles of reliable data transfer

- ❑ 3.5 Connection-oriented transport: TCP
  - ❖ segment structure
  - ❖ reliable data transfer
  - ❖ flow control
  - ❖ connection management
- ❑ 3.6 TCP congestion control

# 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)
- no (delay for) recovering lost segments as in TCP
- 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
  - DNS
  - SNMP
- reliable transfer over UDP: add reliability at application layer
  - application-specific error recovery!

Length, in bytes of UDP segment, including header

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

UDP segment format

8

# 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 nonetheless? ….*

# Internet Checksum Example

- Note
  - When adding numbers, a carryout from the most significant bit needs to be added to the result
- Example: add two 16-bit integers

```
              1 1 1 0 0 1 1 0 0 1 1 0 0 1 1 0
              1 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1
             _____
wraparound (1) 1 0 1 1 1 0 1 1 1 0 1 1 1 0 1 1

       sum    1 0 1 1 1 0 1 1 1 0 1 1 1 1 0 0
  checksum    0 1 0 0 0 1 0 0 0 1 0 0 0 0 1 1
```

# Outline

- 3.1 Transport-layer services
- 3.2 Multiplexing and demultiplexing
- 3.3 Connectionless transport: UDP
- 3.4 Principles of reliable data transfer

- 3.5 Connection-oriented transport: TCP
  - ❖ segment structure
  - ❖ reliable data transfer
  - ❖ flow control
  - ❖ connection management
- 3.6 TCP congestion control

---

# Principles of Reliable data transfer

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



(a) provided service

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

# 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)

---

# 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

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

**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

In this section we will:

❑ 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

**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

# Rdt1.0: <u>reliable transfer over a reliable channel</u>

❑ underlying channel perfectly reliable
  ❖ no bit errors
  ❖ 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)
_____
packet = make_pkt(data)
udt_send(packet)

Wait for call from below    rdt_rcv(packet)
_____
extract (packet,data)
deliver_data(data)

**sender**       **receiver**

# Rdt2.0: <u>channel with bit errors</u>

❑ underlying channel may flip bits in packet
  ❖ 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
❑ new mechanisms in `rdt2.0` (beyond `rdt1.0`):
  ❖ error detection
  ❖ receiver feedback: control msgs (ACK,NAK) rcvr->sender

# rdt2.0: FSM specification

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

rdt_rcv(rcvpkt) &&
isNAK(rcvpkt)

udt_send(sndpkt)

Wait for call from above

Wait for ACK or NAK

rdt_rcv(rcvpkt) && isACK(rcvpkt)

Λ

**sender**

**receiver**

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: operation with no errors

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

rdt_rcv(rcvpkt) &&
isNAK(rcvpkt)

Wait for call from above

Wait for ACK or NAK

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)

14

# 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 has a fatal flaw!

## What happens if ACK/NAK corrupted?

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

## Handling duplicates:

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

**stop and wait**
Sender sends one packet, then waits for receiver Response before sending anything

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

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, handles garbled ACK/NAKs

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)

Wait for 0 from below

Wait for 1 from below

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)

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

extract(rcvpkt,data)
deliver_data(data)
sndpkt = make_pkt(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 ACKs only
- instead of NAK, receiver sends ACK for last pkt received OK
  - receiver must *explicitly* include seq # of pkt being ACKed
  - sender then knows that the current packet was not received correctly
- duplicate ACK at sender results in same action as NAK: *retransmit current pkt*
- This is a simpler protocol because it does away with NAKs

# rdt2.2: sender, receiver fragments

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

Wait for call 0 from above

Wait for ACK 0

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

*sender FSM fragment*

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

Λ

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: channels with errors *and* loss

New assumption:
underlying channel can also lose packets (data or ACKs)
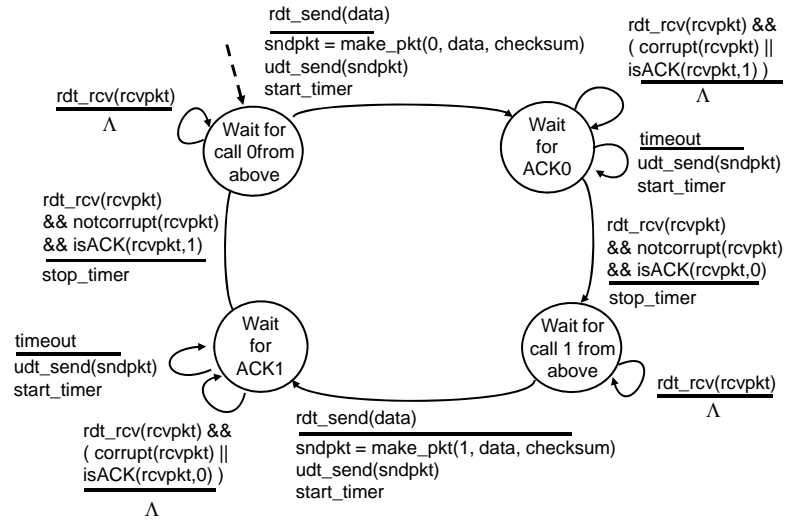
- ❖ checksum, seq. #, ACKs, retransmissions will be of help, but not enough

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)
sndpkt = make_pkt(0, data, checksum)
udt_send(sndpkt)
start_timer

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

rdt_rcv(rcvpkt)
Λ

Wait for
call 0from
above

Wait
for
ACK0

timeout
udt_send(sndpkt)
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(sndpkt)
start_timer

Wait
for
ACK1

Wait for
call 1 from
above

rdt_rcv(rcvpkt)
Λ

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

rdt_send(data)
sndpkt = make_pkt(1, data, checksum)
udt_send(sndpkt)
start_timer

# rdt3.0 in action



(a) operation with no loss

(b) lost packet

# rdt3.0 in action



(c) lost ACK

(d) premature timeout

# Performance of rdt3.0

❑ rdt3.0 works, but performance stinks
❑ eg: 1 Gb/s link, 15 ms propagation delay, 8000 bit packet:

$$d_{trans} = \frac{L}{R} = \frac{8000\text{bits}}{10^9\,\text{b/s}} = 8\,\text{microseconds}$$

❖ U $_{sender}$: utilization – fraction of time sender busy sending

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

❖ 1KB pkt every 30 msec -> 33KB/sec throughput over 1 Gb/s link
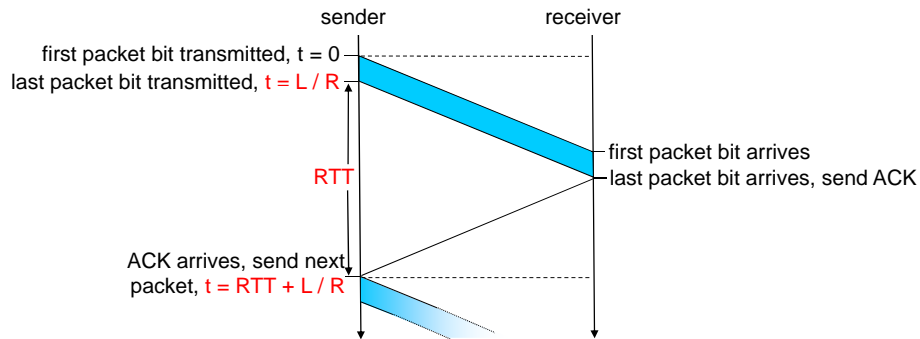❖ network protocol limits use of physical resources!

# rdt3.0: stop-and-wait operation

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

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

# 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

data packet →

data packets →

← ACK packets

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

(b) a pipelined protocol in operation

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

# Pipelining: increased utilization

sender                     receiver

first packet bit transmitted, t = 0
last bit transmitted, t = L / R

first packet bit arrives
last packet bit arrives, send ACK

RTT

last bit of 2$^{nd}$ packet arrives, send ACK
last bit of 3$^{rd}$ packet arrives, send ACK
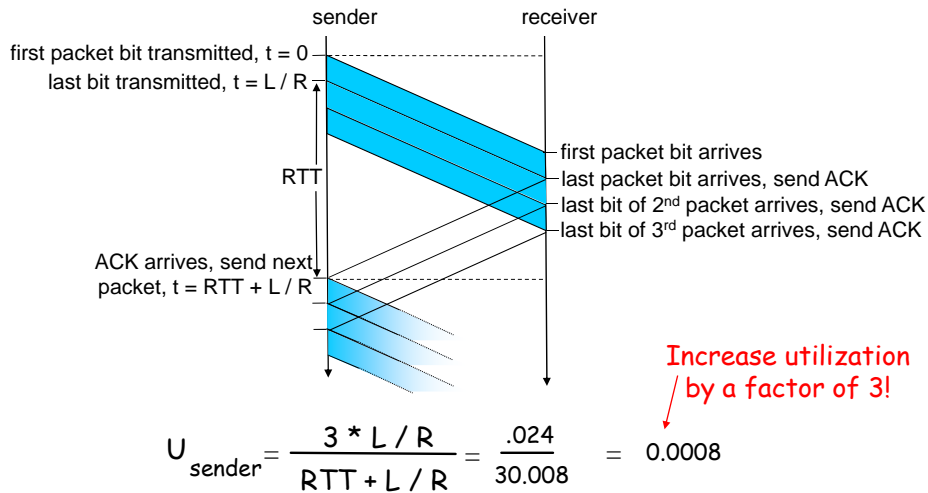
ACK arrives, send next
packet, t = RTT + L / R

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:  overview
- *sender:* up to N unACKed pkts in pipeline
- *receiver:* only sends cumulative ACKs
  - ❖ doesn't ACK pkt if there's a gap
- *sender:* has timer for oldest unACKed pkt
  - ❖ if timer expires: retransmit all unACKed packets
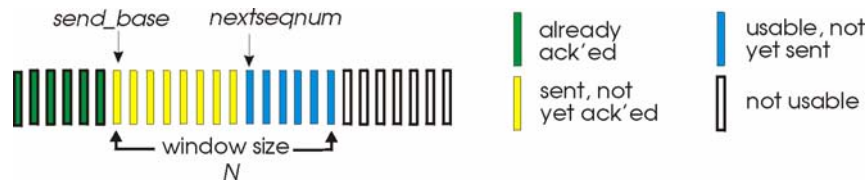
### Selective Repeat:  overview
- *sender:* up to N unACKed packets in pipeline
- *receiver:* ACKs individual pkts
- *sender:* maintains timer for each unACKed pkt
  - ❖ if timer expires: retransmit only unACKed packet

# Go-Back-N

- ❑ k-bit seq # in pkt header
- ❑ "sliding window" of up to N, consecutive unACKed pkts allowed



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

---

# GBN: sender extended FSM

rt_send(data)

```
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)
```

Λ
base=1
nextseqnum=1

**Wait**

```
timeout
start_timer
udt_send(sndpkt[base])
udt_send(sndpkt[base+1])
…
udt_send(sndpkt[nextseqnum-1])
```

rdt_rcv(rcvpkt)
  && corrupt(rcvpkt)

rdt_rcv(rcvpkt) &&
  notcorrupt(rcvpkt)

```
base = getacknum(rcvpkt)+1
If (base == nextseqnum)
   stop_timer
 else
   start_timer
```

# GBN: receiver extended FSM

default
udt_send(sndpkt)

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

Wait

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

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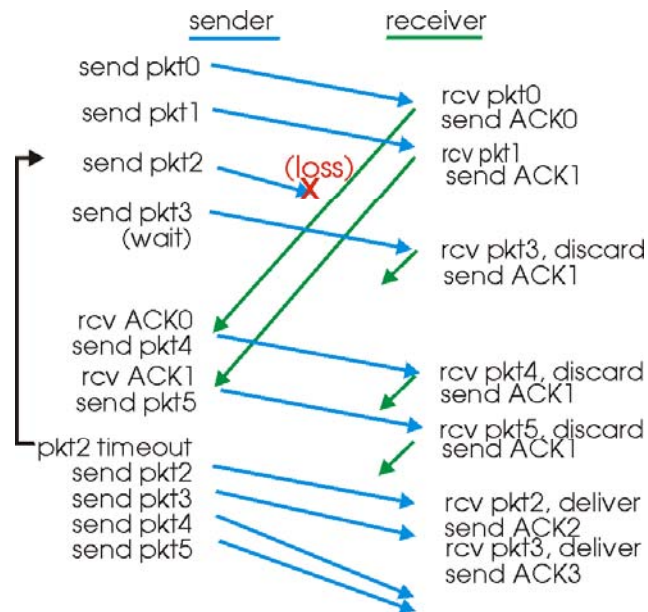  ❖ Re-ACK pkt with highest in-order seq #

---

# GBN in action

sender                    receiver

send pkt0 ————————→ rcv pkt0
send pkt1 ————————→ send ACK0
                         rcv pkt1
send pkt2  (loss)         send ACK1
             X
send pkt3 ————————→ rcv pkt3, discard
 (wait)                   send ACK1

rcv ACK0
send pkt4 ————————→ rcv pkt4, discard
rcv ACK1                  send ACK1
send pkt5 ————————→ rcv pkt5, discard
pkt2 timeout              send ACK1
send pkt2 ————————→ 
send pkt3 ————————→ rcv pkt2, deliver
send pkt4                 send ACK2
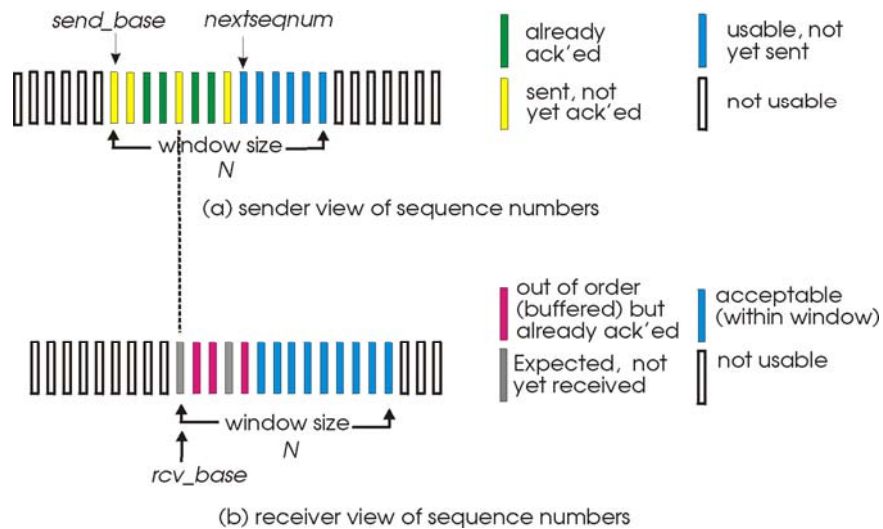send pkt5 ————————→ rcv pkt3, deliver
                         send ACK3

# Selective Repeat

❑ Go-back-N can be inefficient if there can be many pkts in pipeline and an error occurs
  ❖ All these packets will be retransmitted unnecessarily
❑ With 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

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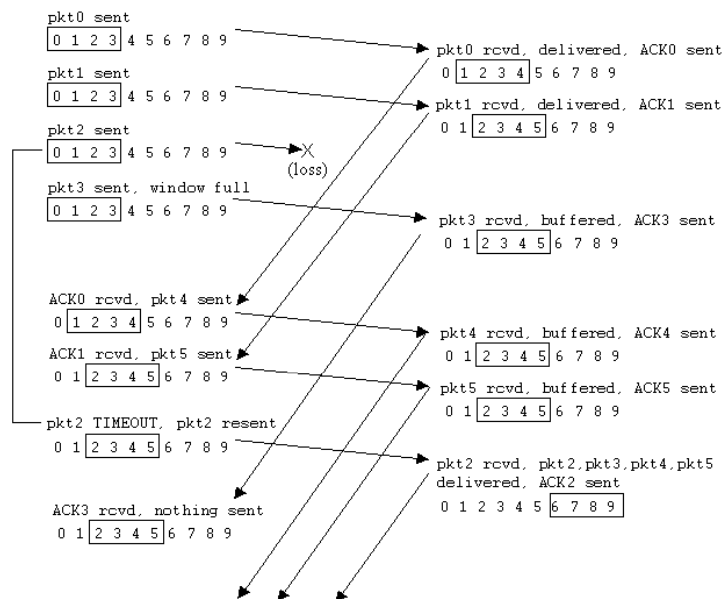
**pkt n in** [rcvbase-N,rcvbase-1]
- ACK(n)

**otherwise:**
- ignore

# Selective repeat in action

```
pkt0 sent
0 1 2 3 4 5 6 7 8 9                          pkt0 rcvd, delivered, ACK0 sent
pkt1 sent                                    0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9                          pkt1 rcvd, delivered, ACK1 sent
pkt2 sent                                    0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9          X
pkt3 sent, window full      (loss)
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

ACK0 rcvd, pkt4 sent
0 1 2 3 4 5 6 7 8 9                          pkt4 rcvd, buffered, ACK4 sent
ACK1 rcvd, pkt5 sent                         0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9                          pkt5 rcvd, buffered, ACK5 sent
pkt2 TIMEOUT, pkt2 resent                    0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9                          pkt2 rcvd, pkt2,pkt3,pkt4,pkt5
                                             delivered, ACK2 sent
ACK3 rcvd, nothing sent                      0 1 2 3 4 5 6 7 8 9
0 1 2 3 4 5 6 7 8 9
```
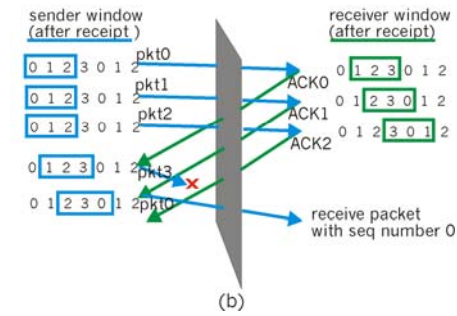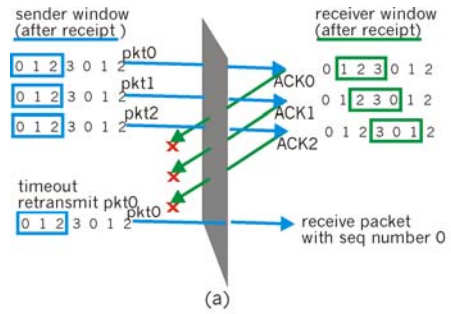
# Selective repeat: dilemma

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?