# Chapter 3
# Transport Layer

A note on the use of these PowerPoint slides:
We're making these slides freely available to all (faculty, students, readers). They're in PowerPoint form so you see the animations; and can add, modify, and delete slides  (including this one) and slide content to suit your needs. They obviously represent a *lot* of work on our part. In return for use, we only ask the following:

- If you use these slides (e.g., in a class) that you mention their source (after all, we'd like people to use our book!)
- If you post any slides on a www site, that you note that they are adapted from (or perhaps identical to) our slides, and note our copyright of this material.

For a revision history, see the slide note for this page.

Thanks and enjoy!  JFK/KWR

*Computer Networking: A Top-Down Approach*
8th edition
Jim Kurose, Keith Ross
Pearson, 2020

# Transport layer: overview
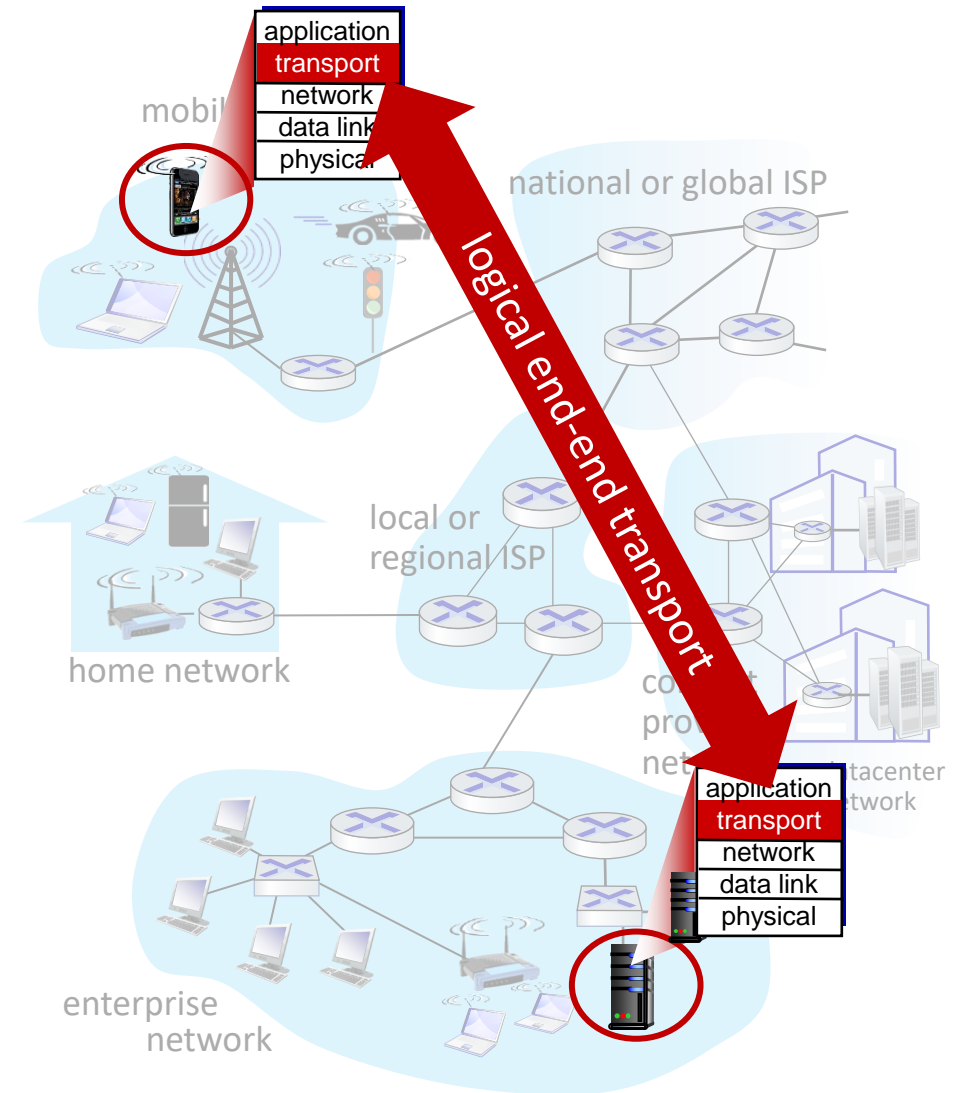
*Our goal:*

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

- learn about Internet transport layer protocols:
  - UDP: connectionless transport
  - TCP: connection-oriented reliable transport
  - TCP congestion control

# Transport layer: roadmap

- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- Connection-oriented transport: TCP
- Principles of congestion control
- TCP congestion control
- Evolution of transport-layer functionality

# Transport services and protocols

- **provide** *logical communication* between application processes running on different hosts

- **transport protocols actions in end systems:**
  - sender: breaks application messages into *segments*, passes to network layer
  - receiver: reassembles segments into messages, passes to application layer

- **two transport protocols available to Internet applications**
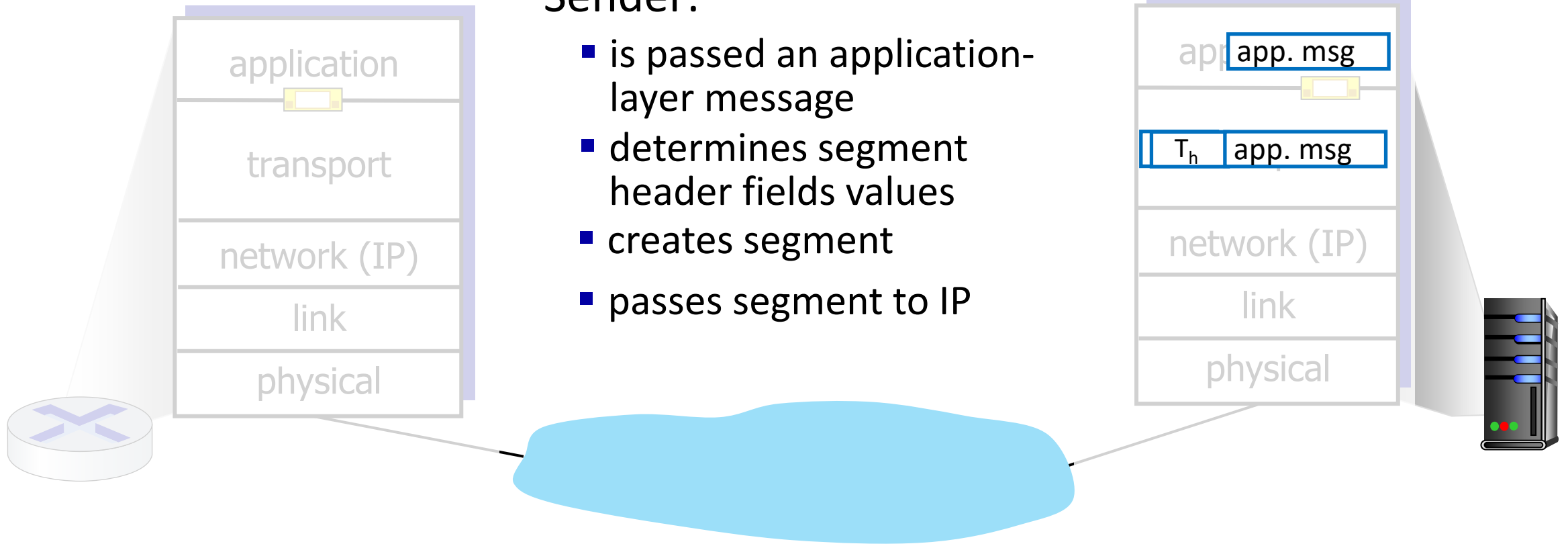  - TCP, UDP

# Transport vs. network layer services and protocols

- **transport layer**: logical communication between *processes*
  - relies on, enhances, network layer services
  - PDU: Segment
  - extends "host-to-host" communication to "process-to-process" communication

- **network layer**: logical communication between *hosts*

  - PDU: Datagram/Packet

  - Datagram's may be lost, duplicated, reordered in the Internet – "best effort" service

# Transport Layer Actions

**Sender:**

- is passed an application-layer message
- determines segment header fields values
- creates segment
- passes segment to IP

application

transport

network (IP)

link

physical



app. msg

$T_h$ | app. msg

app

network (IP)

link

physical

# Transport Layer Actions



**Receiver:**

- receives segment from IP
- checks header values
- extracts application-layer message
- **demultiplexes** message up to application via socket

# Two principal Internet transport protocols

- **TCP:** Transmission Control Protocol
  - reliable, in-order delivery
  - congestion control
  - flow control
  - connection setup

- **UDP:** User Datagram Protocol
  - unreliable, unordered delivery
  - no-frills extension of "best-effort" IP

- services *not* available:
  - delay guarantees
  - bandwidth guarantees

# Chapter 3: roadmap

# Multiplexing/demultiplexing

*multiplexing as sender:*

handle data from multiple sockets, add transport header (later used for demultiplexing)

*demultiplexing as receiver:*

use header info to deliver received segments to correct socket

# HTTP server

client

application

transport

$H_n\,H_t$ HTTP msg

link

physical

HTTP msg

$H_t$ HTTP msg

$H_n\,H_t$ HTTP msg

link

physical

application

transport

network

link

physical

$H_n\,H_t$ HTTP msg

**Q: how did transport layer know to deliver message to Firefox browser process rather than Netflix process or Skype process?**

client

application

HTTP msg

NETFLIX

transport

network

link

physical

APACHE
HTTP SERVER

HTTP msg

$H_t$ HTTP msg

network

link

physical

application

transport

network

link

physical

de-multiplexing

application

transport

de-multiplexing

multiplexing

application

transport

multiplexing

# How demultiplexing works

- host receives IP datagrams
  - each datagram has source IP address, destination IP address
  - each datagram carries one 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
(payload)

TCP/UDP segment format

# Connectionless demultiplexing

- UDP socket identified by two-tuple:
  - destination IP address
  - destination port #

when receiving host receives *UDP* segment:

- checks destination port # in segment
- directs UDP segment to socket with that port #

IP/UDP datagrams with *same dest. port #,* but different source IP addresses and/or source port numbers will be directed to *same socket* at receiving host

# Connectionless demultiplexing: an example



application

P3

transport

network

link

physical

application

P1

transport

network

link

physical

application

P4

transport

network

link

physical

B
source port: 6428
dest port: 9157

D
source port: ?
dest port: ?

A
source port: 9157
dest port: 6428

C
source port: ?
dest port: ?

# Connection-oriented demultiplexing

- TCP socket identified by 4-tuple:
  - source IP address
  - source port number
  - dest IP address
  - dest port number
- demux: receiver uses *all four values (4-tuple)* to direct segment to appropriate socket

- server may support many simultaneous TCP sockets:
  - each socket identified by its own 4-tuple
  - each socket associated with a different connecting client
- Web servers have different sockets for each connecting client.
  - non-persistent HTTP will have different socket for each request

# Connection-oriented demultiplexing: example



**Three segments, all destined to IP address: B,**
**dest port: 80 are demultiplexed to *different* sockets**

# Summary

- Multiplexing, demultiplexing: based on segment, datagram header field values

- **UDP:** demultiplexing using destination port number (only)

- **TCP:** demultiplexing using 4-tuple: source and destination IP addresses, and port numbers

- Multiplexing/demultiplexing happen at *all* layers

# Chapter 3: roadmap

- Transport-layer services
- Multiplexing and demultiplexing
- **Connectionless transport: UDP**
- Principles of reliable data transfer
- Connection-oriented transport: TCP
- Principles of congestion control
- TCP congestion control
- Evolution of transport-layer functionality

# UDP: User Datagram Protocol

- "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 RTT delay)
- simple: no connection state at sender, receiver
- small header size
- no congestion control
  - UDP can blast away as fast as desired!
  - can function in the face of congestion

# UDP: User Datagram Protocol

- UDP use:
  - streaming multimedia apps (loss tolerant, rate sensitive)
  - DNS
  - SNMP
  - HTTP/3
- if reliable transfer needed over UDP (e.g., HTTP/3):
  - add needed reliability at application layer
  - add congestion control at application layer

# UDP: User Datagram Protocol [RFC 768]

```
                                                          INTERNET STANDARD

RFC 768                                                        J. Postel
                                                                     ISI
                                                          28 August 1980


                         User Datagram Protocol
                         ----------------------

Introduction
------------

This User Datagram  Protocol  (UDP)  is  defined  to  make  available  a
datagram   mode  of  packet-switched   computer   communication  in  the
environment  of  an  interconnected  set  of  computer  networks.    This
protocol  assumes   that  the  Internet   Protocol  (IP)  [1] is used as the
underlying protocol.

This protocol  provides  a procedure  for application  programs  to send
messages  to other programs  with a minimum  of protocol mechanism.  The
protocol  is transaction oriented, and delivery and duplicate protection
are not guaranteed.  Applications requiring ordered reliable delivery of
streams of data should use the Transmission Control Protocol (TCP) [2].

Format
------


                0      7 8     15 16    23 24     31
               +--------+--------+--------+--------+
               |     Source      |   Destination   |
               |      Port       |      Port       |
               +--------+--------+--------+--------+
               |                 |                 |
               |     Length      |    Checksum     |
               +--------+--------+--------+--------+
               |
               |          data octets ...
               +--------------- ...
```

# UDP: Transport Layer Actions

### SNMP client

application

transport
(UDP)

network (IP)

link

physical

### SNMP server

application

transport
(UDP)

network (IP)

link

physical

# UDP: Transport Layer Actions

SNMP client

SNMP server

application

transport
(UDP)

network (IP)

link

physical

UDP sender actions:

- is passed an application-layer message
- determines UDP segment header fields values
- creates UDP segment
- passes segment to IP

application

SNMP msg

transport
(UDP)

$UDP_h$ | SNMP msg

network (IP)

link

physical

# UDP: Transport Layer Actions

### SNMP client

application

transport
SNMP msg
(UDP)

UDP_h | SNMP msg

link

physical

### SNMP server

application

transport
(UDP)

network (IP)

link

physical

UDP receiver actions:

- receives segment from IP
- checks UDP checksum header value
- extracts application-layer message
- demultiplexes message up to application via socket

# UDP segment header



UDP segment format

# UDP checksum

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

|  | 1st number | 2nd number | sum |
|---|---|---|---|
| Transmitted: | 5 | 6 | 11 |
| Received: | 4 | 6 | 11 |

receiver-computed checksum ≠ sender-computed checksum (as received)

# Internet checksum

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

### sender:

- treat contents of UDP segment (including UDP header fields and IP addresses) as sequence of 16-bit integers

- checksum: addition (one's complement sum) of segment content

- checksum value put into UDP checksum field

### receiver:

- compute checksum of received segment

- check if computed checksum equals checksum field value:
  - not equal - error detected
  - equal - no error detected. *But maybe errors nonetheless?* More later ....

# Internet checksum: an example

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 ⓵ 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

*Note:* when adding numbers, a carryout from the most significant bit needs to be added to the result

* Check out the online interactive exercises for more examples: http://gaia.cs.umass.edu/kurose_ross/interactive/

# Internet checksum: weak protection!

example: add two 16-bit integers

```
        1 1 1 0 0 1 1 0 0 1 1 0 0 1 1 0        0 1
        1 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1        1 0
```

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

Even though numbers have changed (bit flips), *no* change in checksum!

# Summary: UDP

- "no frills" protocol:
  - segments may be lost, delivered out of order
  - best effort service: "send and hope for the best"
- UDP has its plusses:
  - no setup/handshaking needed (no RTT incurred)
  - can function when network service is compromised
  - helps with reliability (checksum)
- build additional functionality on top of UDP in application layer (e.g., HTTP/3)

# Chapter 3: roadmap

- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- **Principles of reliable data transfer**
- Connection-oriented transport: TCP
- Principles of congestion control
- TCP congestion control
- Evolution of transport-layer functionality

# Principles of Reliable Data Transfer

- important in application, transport, and link layers



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

# Reliable Data Transfer with FSMs

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

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

38

# rdt1.0: Reliable Transfer over a Reliable (perfect) Channel

- 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



event

Wait for call from above

rdt_send(data)
_____

packet = make_pkt(data)
udt_send(packet)

actions

**sender**

event

Wait for call from below

rdt_rcv(packet)
_____
extract (packet,data)
deliver_data(data)

actions

**receiver**

# rdt2.0: channel with bit errors
## [stop & wait protocol]

- Assumptions
  - All packets are received
  - Packets may be corrupted (i.e., bits may be flipped)
  - Checksum to detect bit errors

*How do humans recover from "errors" during conversation?*

# rdt2.0: channel with bit errors
## [stop & wait protocol]

- Assumptions
  - All packets are received
  - Packets may be corrupted (i.e., bits may be flipped)
  - Checksum to detect bit errors

- How to recover from errors? ::: Use ARQ mechanism
  - *acknowledgements (ACKs):* receiver explicitly tells sender that packet received correctly
  - *negative acknowledgements (NAKs):* receiver explicitly tells sender that packet had errors
  - sender retransmits <u>packet</u> on receipt of NAK

- in rdt2.0 (beyond rdt1.0):
  - error detection
  - feedback: control msgs (ACK,NAK) from receiver to sender

# rdt2.0: channel with bit errors
## [stop & wait protocol]

- Three additional protocol capabilities are required in ARQ protocols:

  (1) Error detection
  - Extra bits are required

  (2) Receiver feedback
  - Using ACK / NAK

  (3) Retransmission

# rdt2.0: FSM specification
## [stop & wait protocol]

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

**receiver**

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

**sender**

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

# rdt2.0: Observations

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

**sender**

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

1. A stop-and-Wait protocol

2. What happens when ACK or NAK has bit errors?

Approach 1: resend the current data packet?

Duplicate packets

# Handling Duplicate Packets

- sender adds *sequence number* to each packet

- sender retransmits current packet if ACK/NAK garbled

- receiver discards (doesn't deliver up) duplicate packet

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

rdt_send(data)

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

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

udt_send(sndpkt)

**Wait for call 0 from above**

**Wait for ACK or NAK 0**

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

$\Lambda$

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

$\Lambda$

**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)

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)

# rtd2.1: examples

# rdt2.1: summary

Sender:

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

Receiver:

- must check if received packet is duplicate
    - state indicates whether 0 or 1 is expected <u>packet</u> 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 <u>packet</u> received OK
  - receiver must *explicitly* include seq # of <u>packet</u> being ACKed

- **duplicate ACK** at sender results in same action as NAK: *retransmit current <u>packet</u>*

# rdt2.2: sender, receiver fragments

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)**
_____
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)

# rdt2.2: Receiver

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

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

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

udt_send(sndpkt)

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

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(ACK1, chksum)**
udt_send(sndpkt)

# rdt3.0:The case of "Lossy" Channels

- **<u>Assumption:</u>** underlying channel can also lose packets (data or ACKs)
  - checksum, seq. #, ACKs, retransmissions will be of help … but not enough

- **<u>Approach:</u>** sender waits "reasonable" amount of time for ACK (a Time-Out)
  - Time-out value?
  - Possibility of duplicate packets/ACKs?

- if packet (or ACK) just delayed (not lost):
  - retransmission will be  duplicate, but use of seq. #'s already handles this
  - receiver must specify seq # of packet being ACKed

- Countdown timer is required.

# rdt3.0:The case of "Lossy" Channels

- The sender will need to be able
    (1) Start the timer each time a packet is sent
    (2) Respond to a timer interrupt
    (3) Stop the timer

# rdt3.0 sender

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

rdt_rcv(rcvpkt)
_____
Λ

**Wait for call 0from above**

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

**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) &&
( corrupt(rcvpkt) ||
isACK(rcvpkt,0) )
_____
Λ

rdt_rcv(rcvpkt)
_____
Λ

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

# 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$_0$ → rcv pkt0
send ACK0
← ACK$_0$ ←

(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$_0$ → rcv pkt0
send ACK0
← ACK$_0$ ←

(b) lost packet

57

# rdt3.0 in action



(c) ACK loss

(d) premature timeout/ delayed ACK

# stop-and-wait operation



first packet bit transmitted, t = 0

last packet bit transmitted, t = L / R

D

first packet bit arrives

last packet bit arrives, send ACK

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

sender

receiver

**Protocol rdt3.0 is a stop-and-wait protocol.**

# Performance of rdt3.0

- rdt3.0 is correct, but performance stinks

- e.g.: 1 Gbps link, 15 ms prop. delay, 8000 bit packet:

$$D_{trans} = \frac{L}{R} = \frac{8000 \text{ bits}}{10^9 \text{ bits/sec}} = 8 \text{ microsecs}$$

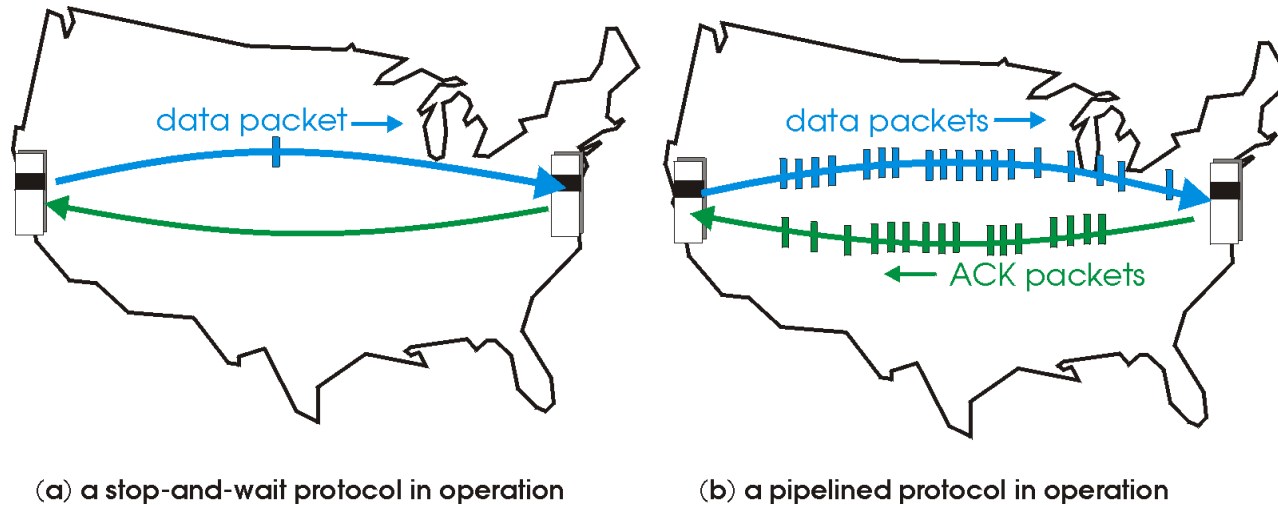  - U$_{sender}$: *utilization* – fraction of time sender busy sending

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

  - if RTT=30 msec, 1 KB pkt every 30 msec: 33 KB/sec throughput over 1 Gbps link
- network protocol limits use of physical resources!
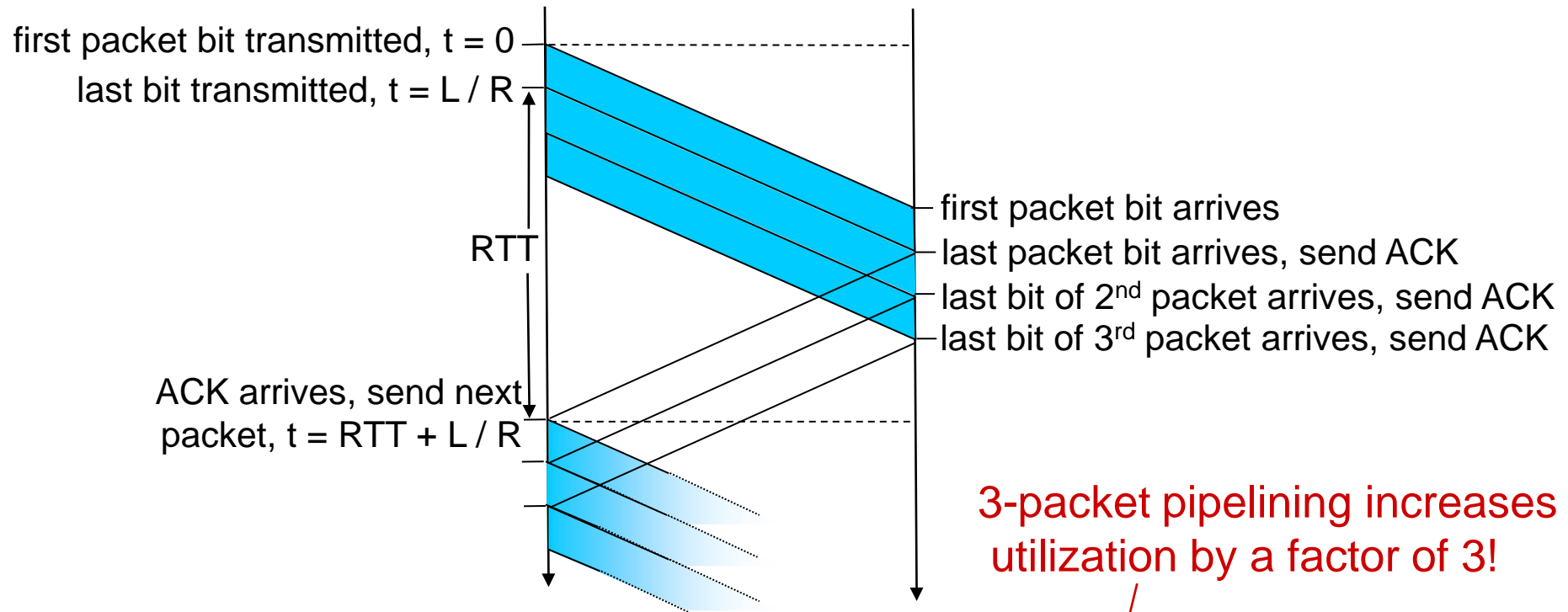
*** U = Utilization ***

# 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 and selective repeat*

# Pipelining: increased utilization

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 2$^{nd}$ packet arrives, send ACK

last bit of 3$^{rd}$ packet arrives, send ACK

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

3-packet pipelining increases utilization by a factor of 3!

$$U_{sender} = \frac{3L / R}{RTT + L / R} = \frac{.0024}{30.008} = 0.00081$$

# Pipelined protocols: overview

## Go-back-N:

- sender can have up to N unacked packets in pipeline
- receiver only sends *cumulative ack*
  - doesn't ack packet if there's a gap
- sender uses a single timer for oldest unacked packet
  - when timer expires, retransmit *all* unacked packets

## Selective Repeat:
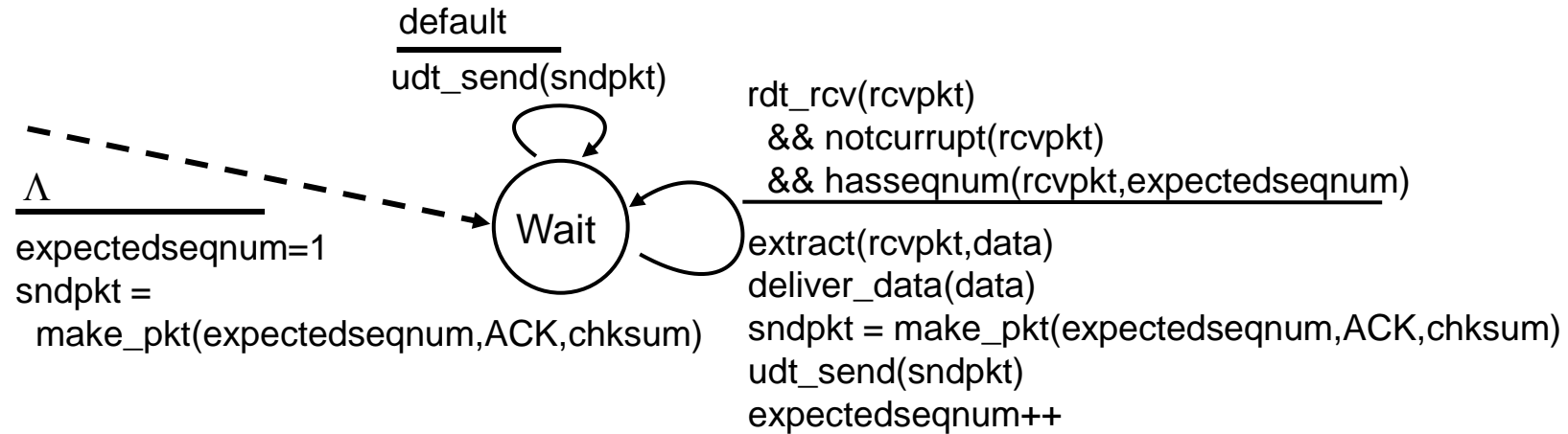
- sender can have up to N unacked packets in pipeline
- receiver sends *individual ack* for each packet

- sender maintains timer for each unacked packet
  - when timer expires, retransmit only that unacked packet

# GBN: 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
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
```

64

# GBN: receiver extended FSM

default
udt_send(sndpkt)

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

Λ

Wait

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

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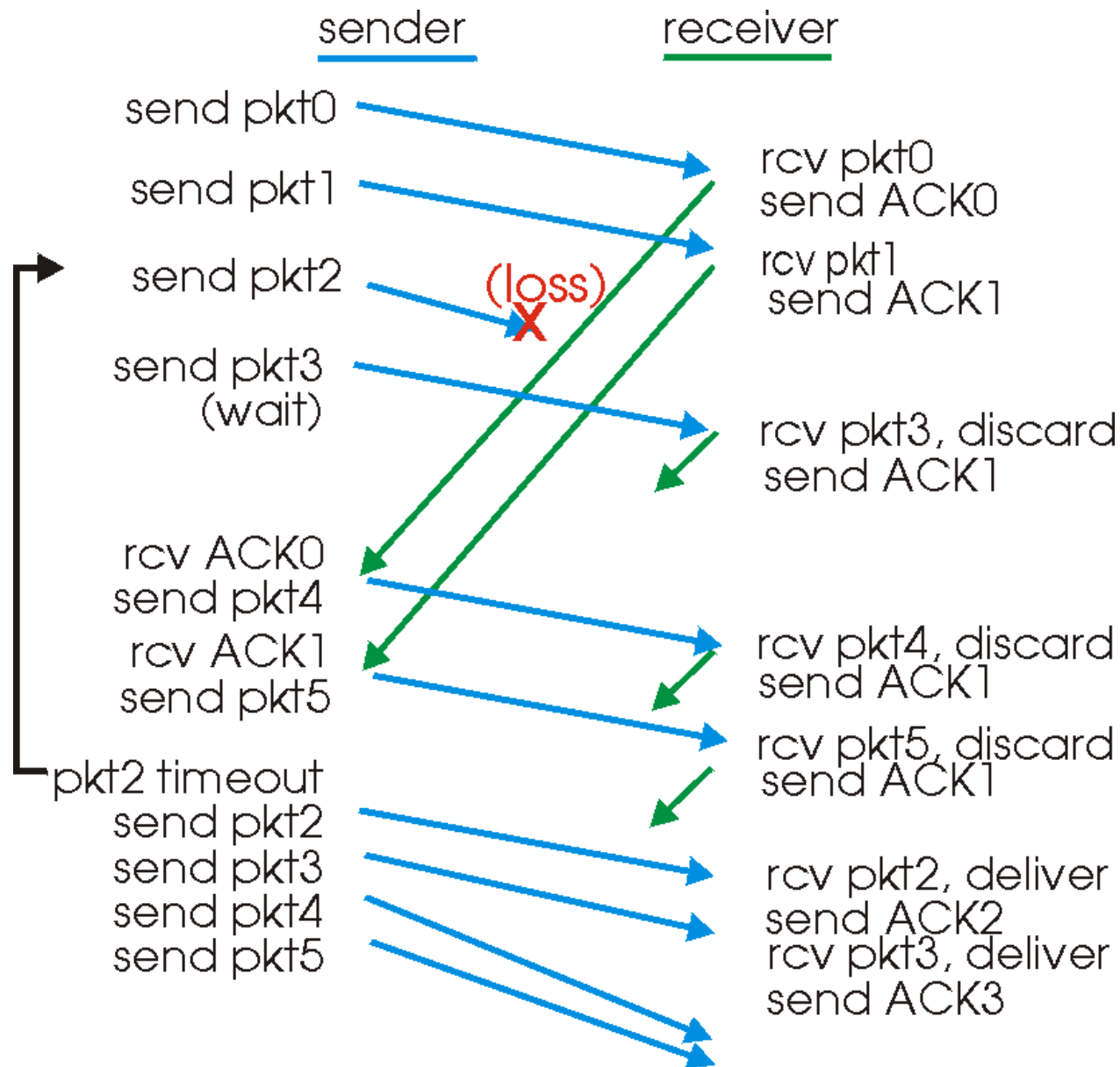expectedseqnum++

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!
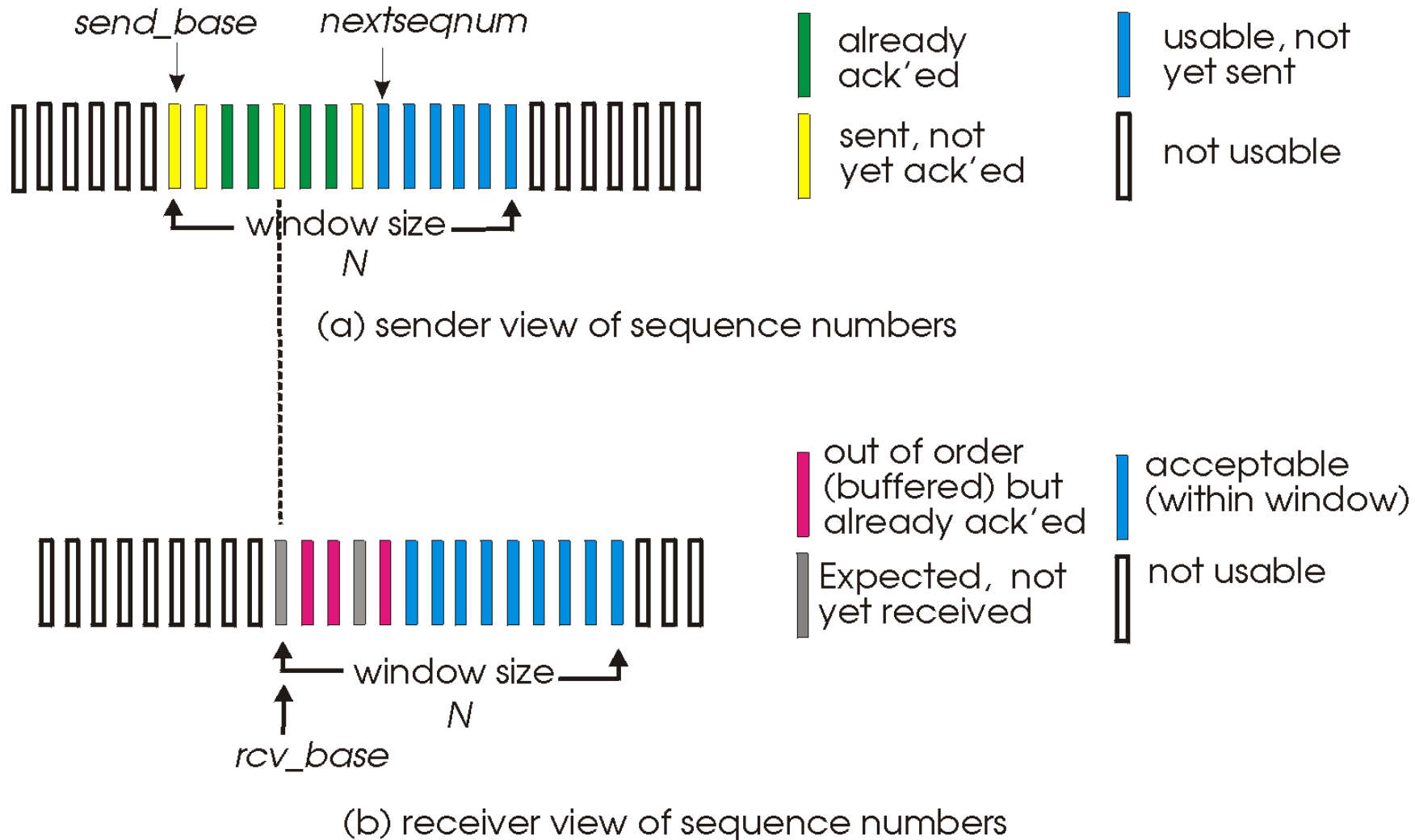- Re-ACK pkt with highest in-order seq #

# GBN in Action

# 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

*send_base*  *nextseqnum*

window size N

already ack'ed

sent, not yet ack'ed

usable, not yet sent

not usable

(b) receiver view of sequence numbers

*rcv_base*

window size N

out of order (buffered) but already ack'ed

Expected, not yet received

acceptable (within window)

not usable

# Selective Repeat

**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

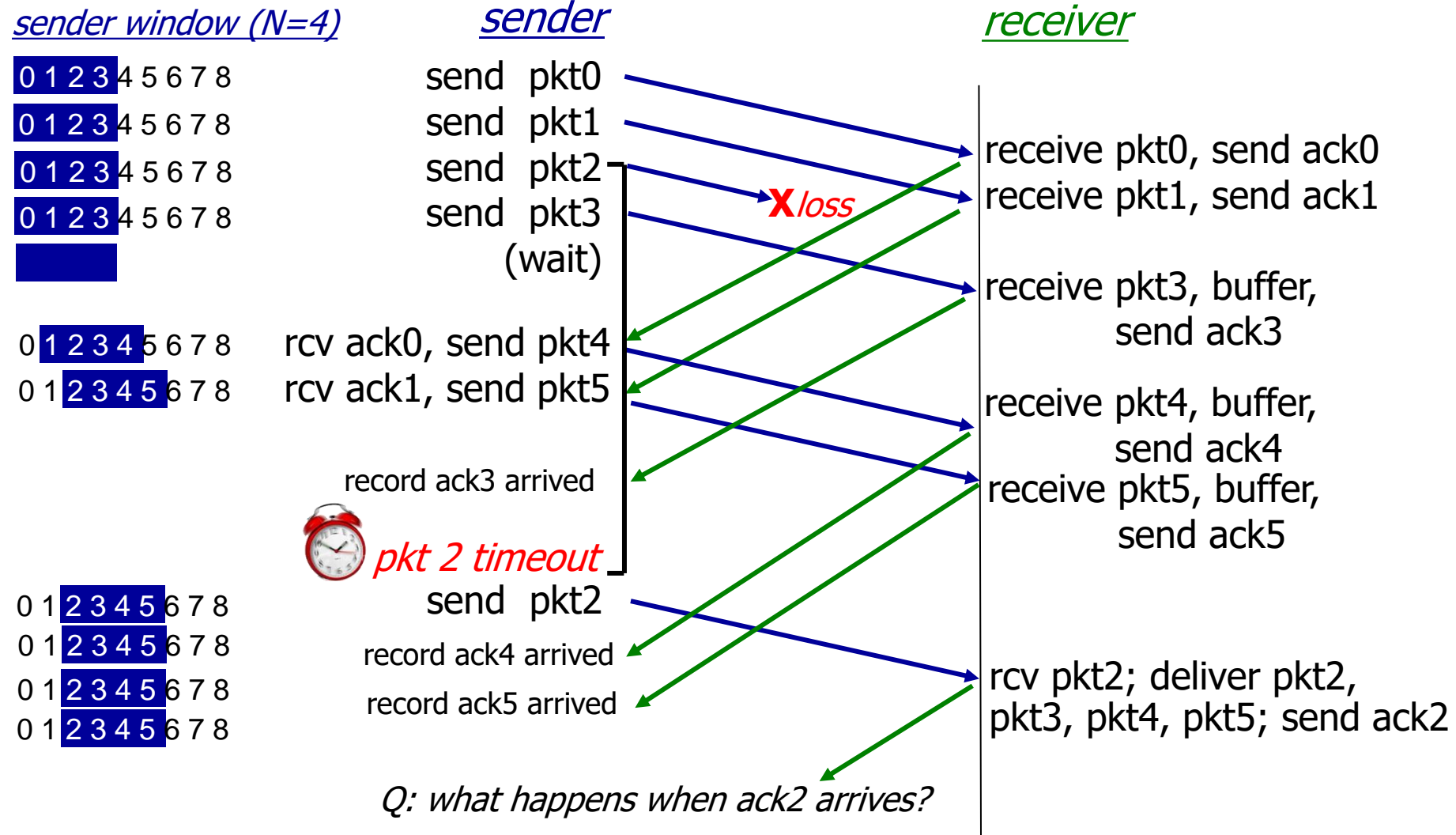- in-order: deliver (also deliver buffered, in-order pkts), advance window to next not-yet-received pkt

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

- ACK(n)

**otherwise:**

- ignore

# Selective Repeat in Action

**sender window (N=4)**   **sender**   **receiver**

`0 1 2 3` 4 5 6 7 8   send  pkt0

`0 1 2 3` 4 5 6 7 8   send  pkt1

`0 1 2 3` 4 5 6 7 8   send  pkt2 — receive pkt0, send ack0

`0 1 2 3` 4 5 6 7 8   send  pkt3   **X** *loss*   receive pkt1, send ack1

(wait)

receive pkt3, buffer,
send ack3

0 `1 2 3 4` 5 6 7 8   rcv ack0, send pkt4

0 1 `2 3 4 5` 6 7 8   rcv ack1, send pkt5

receive pkt4, buffer,
send ack4

record ack3 arrived

receive pkt5, buffer,
send ack5

*pkt 2 timeout*

0 1 `2 3 4 5` 6 7 8   send  pkt2

0 1 `2 3 4 5` 6 7 8   record ack4 arrived

0 1 `2 3 4 5` 6 7 8   record ack5 arrived   rcv pkt2; deliver pkt2,
pkt3, pkt4, pkt5; send ack2

0 1 `2 3 4 5` 6 7 8
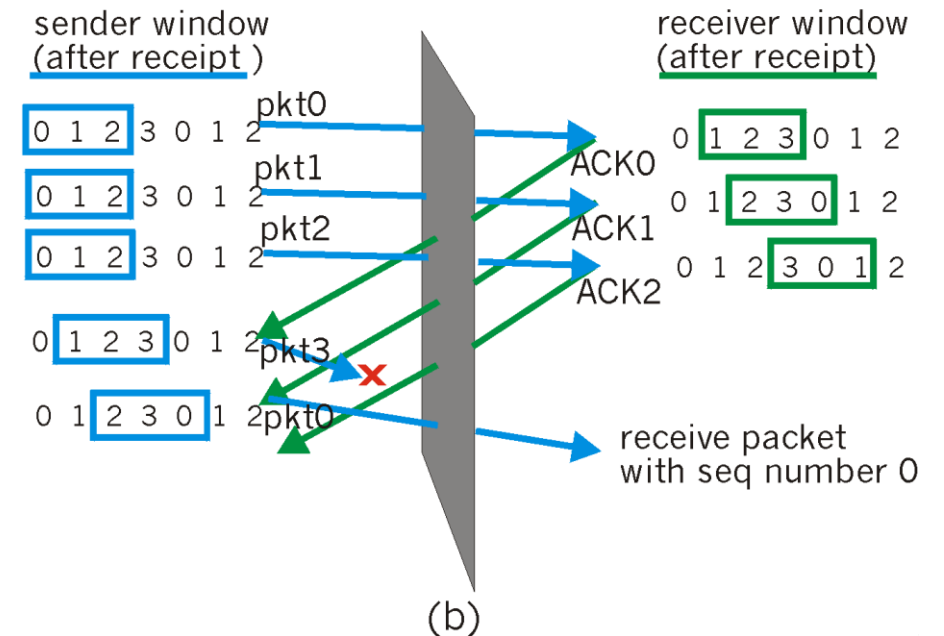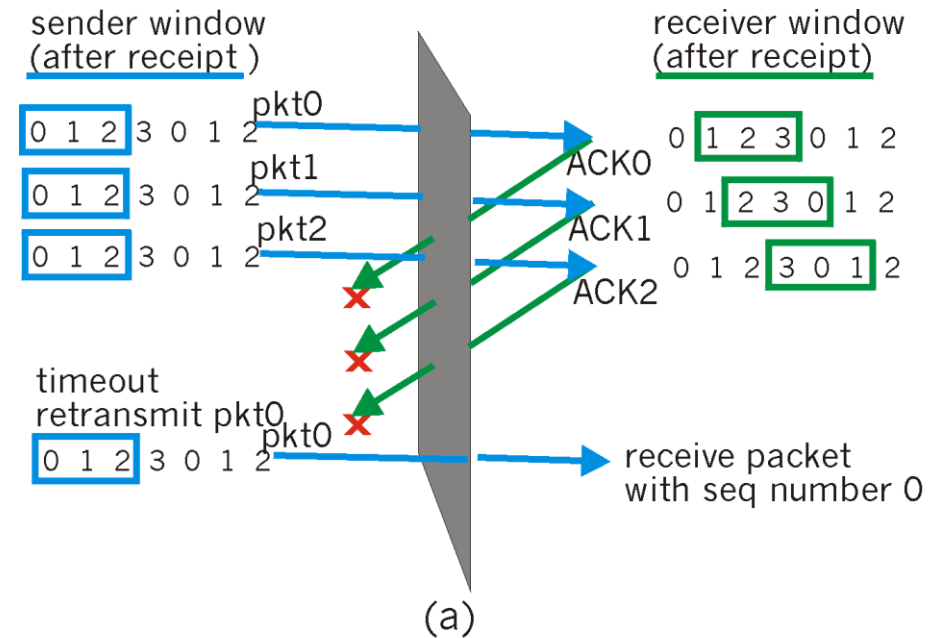
*Q: what happens when ack2 arrives?*

# 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?

# Chapter 3: Roadmap

- Connection-oriented transport: TCP
- Principles of congestion control
- TCP congestion control
- Evolution of transport-layer functionality

To Be Continued

つづく