# COMS 2014A / 2020A

# Computer Networks

Mr. Gift Khangamwa

Office: TWK MSB UG05

# Lecture 5: Transport layer

*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
  - We will briefly outline reliable transfer protocol requirements
  - TCP: connection-oriented reliable transport
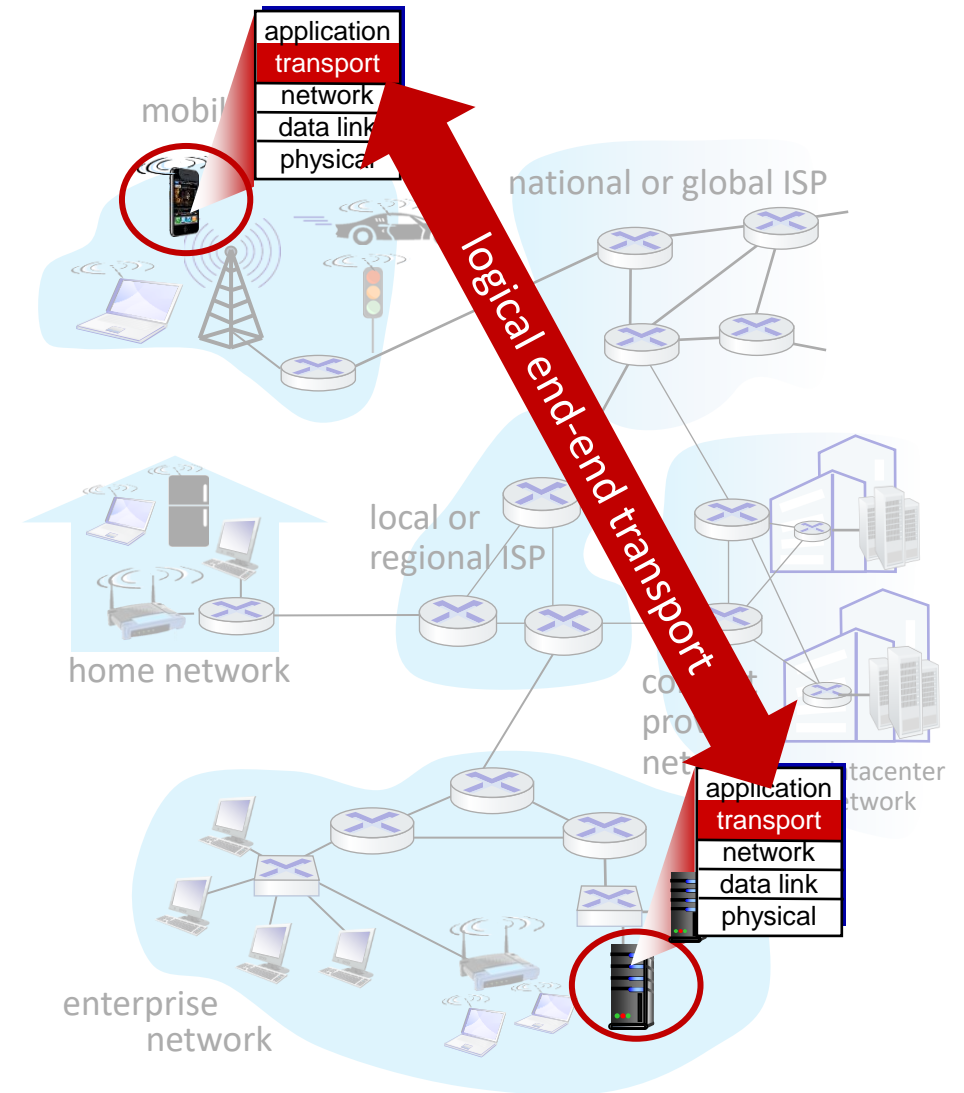  - TCP congestion control

# 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



© 1922. WM. EVANS.

THERE was an old woman who lived in a shoe,
She had so many children, she didn't know what to do.
She gave them some milk and nice butter bread,
She kissed them all round and put them to bed.

## *household analogy:*

*12 kids in Ann's house sending letters to 12 kids in Bill's house:*

- hosts = houses
- processes = kids
- app messages = letters in envelopes

# Transport vs. network layer services and protocols

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

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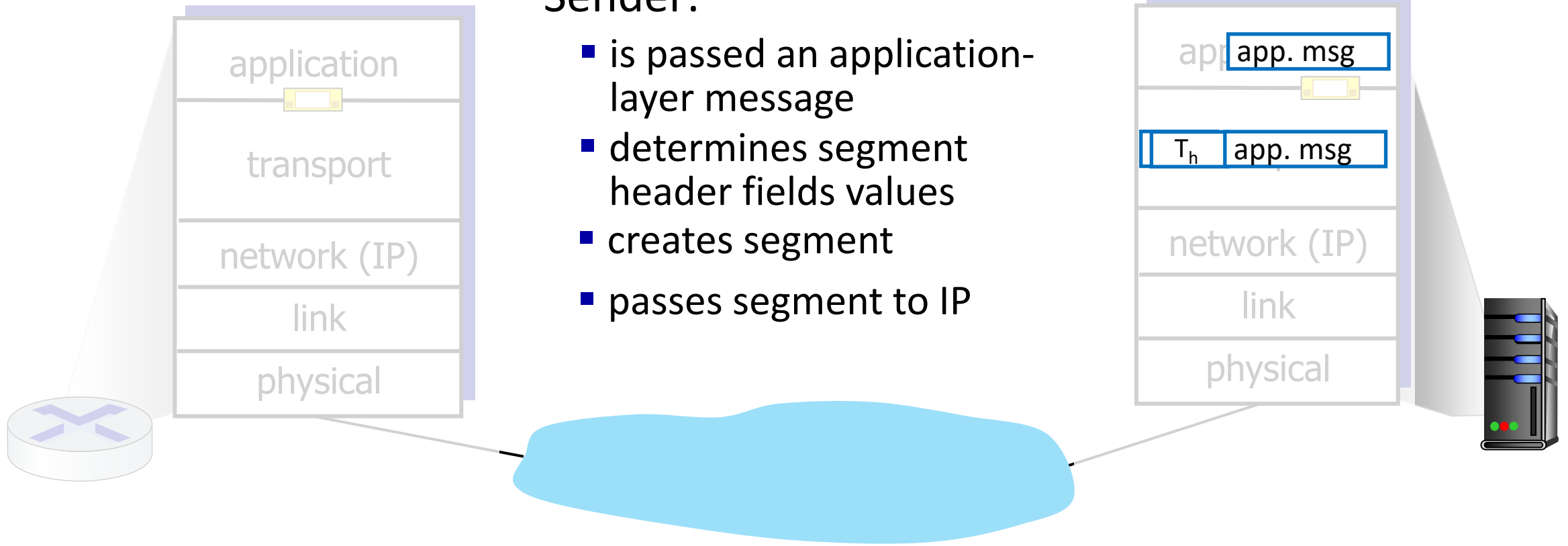  - relies on, enhances, network layer services

*household analogy:*

*12 kids in Ann's house sending letters to 12 kids in Bill's house:*

- hosts = houses
- processes = kids
- app messages = letters in envelopes
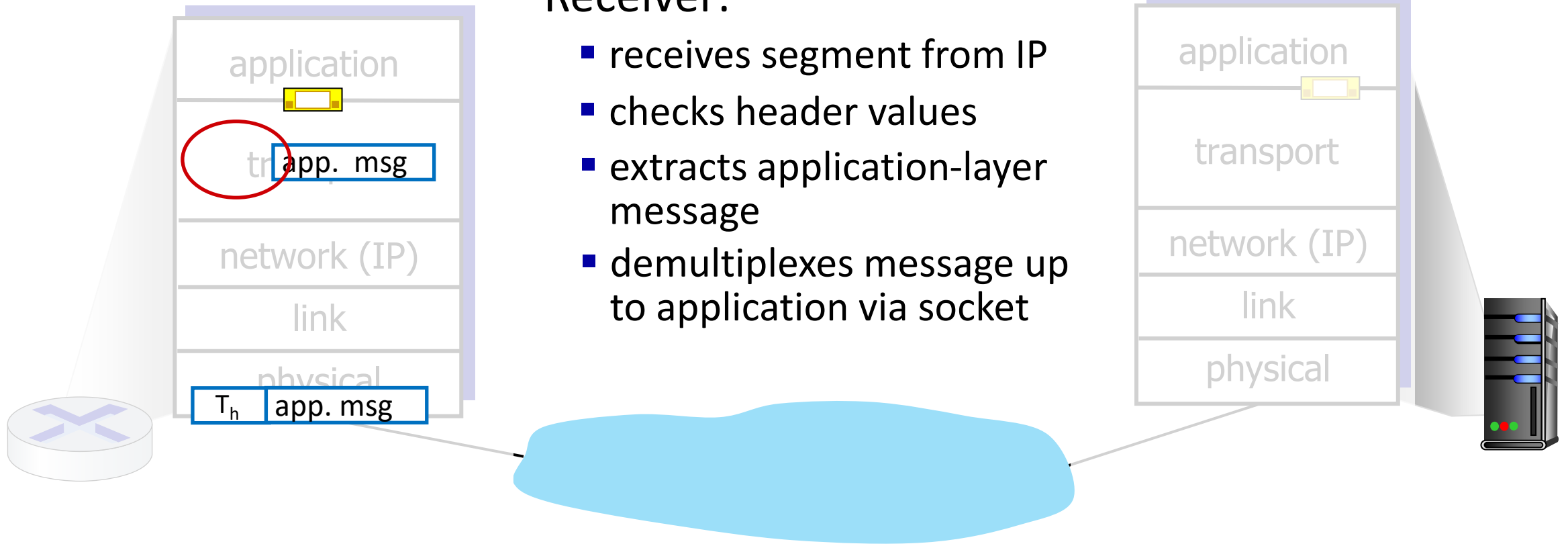
# Transport Layer Actions

**Sender:**

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

app. msg

| $T_h$ | app. msg |

application

transport

network (IP)

link

physical

app...

network (IP)

link

physical

# Transport Layer Actions

application

transport

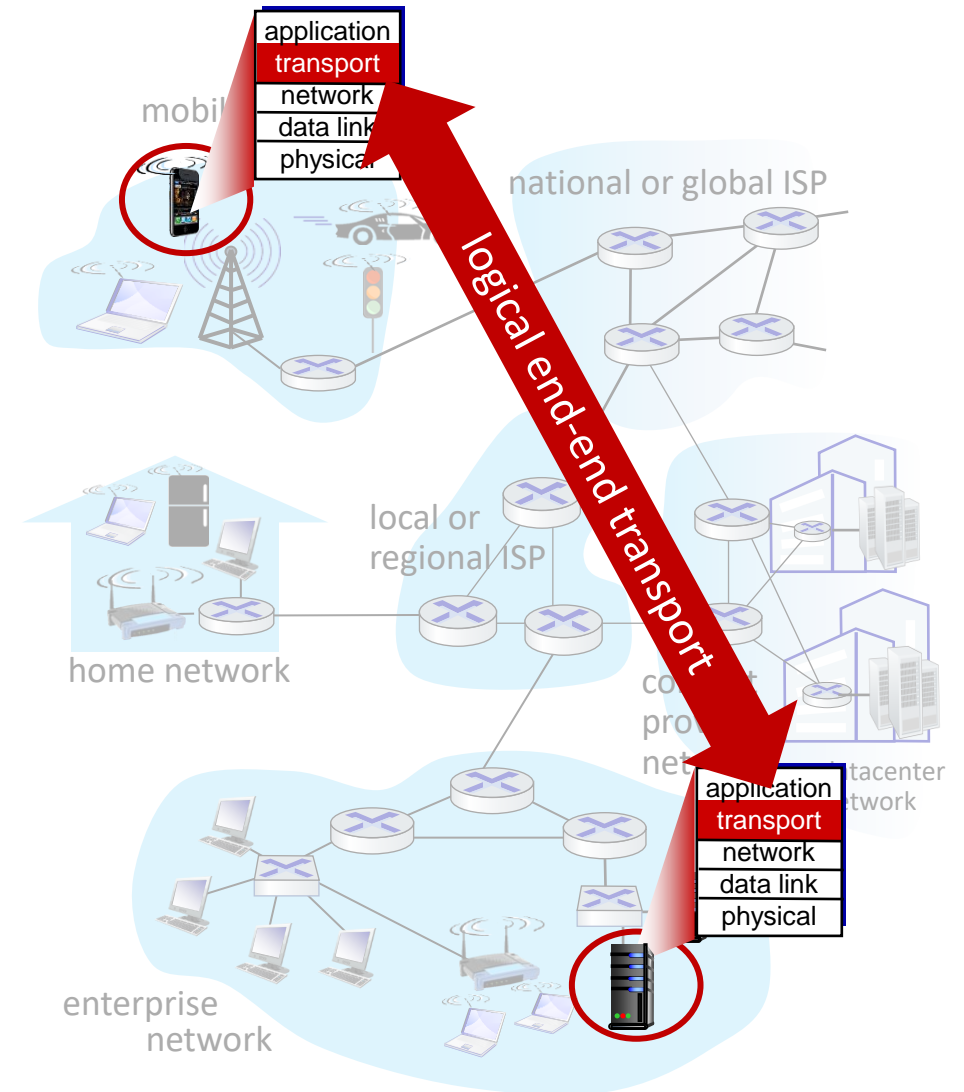app. msg

network (IP)

link

physical

T_h  app. msg

Receiver:

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

application

transport

network (IP)

link

physical

# 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

# TCP vrs UDP

# 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

client

HTTP server

application

NETFLIX

transport

network

link

physical

HTTP msg

transport

network

link

physical

application

transport

network

link

physical

client

HTTP server

application

transport

network

link

physical

HTTP msg

H_t HTTP msg

network

link

physical

application

transport

network

link

physical

# client

## HTTP server

application

transport

network

link

physical

HTTP msg

$H_t$ HTTP msg

$H_n H_t$ HTTP msg

link

physical

application

transport

network

link

physical

# HTTP server

## client



application

NETFLIX

transport

network

link

physical

APACHE
HTTP SERVER

transport

network

link

physical

application

transport

network

link

physical

$H_n H_t$ HTTP msg

transport-layer multiplexing requires (1) that sockets have unique identifiers, and (2) that each segment have special fields that indicate the socket to which the segment is to be delivered

# HTTP server

client₁

application

transport

network

link

physical

P-client₁    P-client₂

transport

network

link

physical

client₂

application

transport

network

link

physical

# Multiplexing/demultiplexing

*multiplexing at sender:*

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

*demultiplexing at receiver:*

use header info to deliver received segments to correct socket

# 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

*Recall:*

- when creating socket, must specify *host-local* port #:

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

- when creating datagram to send into UDP socket, must specify
  - destination IP address
  - destination port #
  - Range **1024 - 65535**

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

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

`DatagramSocket mySocket2 = new DatagramSocket (9157);`

`DatagramSocket mySocket1 = new DatagramSocket (5775);`

application

P3

transport

network

link

physical

application

P1

transport

network

link

physical

application

P4

transport

network

link

physical

source port: 6428
dest port: 9157

source port: ?
dest port: ?

source port: 9157
dest port: 6428

source port: ?
dest port: ?

# Quick overview UDP Socket programming

Two socket types for two transport services:

- *UDP:* unreliable datagram
- *TCP:* reliable, byte stream-oriented

Application Example  Tutorial 1:

1. client reads a line of characters (data) from its keyboard and sends data to server
2. server receives the data and converts characters to uppercase
3. server sends modified data to client
4. client receives modified data and displays line on its screen

# Socket programming with UDP

UDP: no "connection" between client and server:

- no handshaking before sending data
- sender explicitly attaches IP destination address and port # to each packet
- receiver extracts sender IP address and port# from received packet

UDP: transmitted data may be lost or received out-of-order

Application viewpoint:

- UDP provides *unreliable* transfer of groups of bytes ("datagrams") between client and server processes

# Client/server socket interaction: UDP

**server** (running on serverIP)

**client**

create socket, port= x:
serverSocket =
socket(AF_INET,SOCK_DGRAM)

create socket:
clientSocket =
socket(AF_INET,SOCK_DGRAM)

read datagram from
serverSocket

Create datagram with serverIP address
And port=x; send datagram via
clientSocket

write reply to
serverSocket
specifying
client address,
port number

read datagram from
clientSocket

close
clientSocket

# Example app: UDP server

*Python UDPServer*

```python
from socket import *
serverPort = 12000
serverSocket = socket(AF_INET, SOCK_DGRAM)
serverSocket.bind(('', serverPort))
print ("The server is ready to receive")
while True:
    message, clientAddress = serverSocket.recvfrom(2048)
    modifiedMessage = message.decode().upper()
    serverSocket.sendto(modifiedMessage.encode(),
                                        clientAddress)
```

create UDP socket →

bind socket to local port number 12000 →

loop forever →

Read from UDP socket into message, getting → client's address (client IP and port)

send upper case string back to this client →

# Example app: UDP client

*Python UDPClient*

include Python's socket library ⟶ from socket import *

serverName = 'hostname'

serverPort = 12000

create UDP socket for server ⟶ clientSocket = socket(AF_INET,
                                            SOCK_DGRAM)

get user keyboard input ⟶ message = raw_input('Input lowercase sentence:')

attach server name, port to message; send into socket ⟶ clientSocket.sendto(message.encode(),
                                            (serverName, serverPort))

read reply characters from socket into string ⟶ modifiedMessage, serverAddress =
                                            clientSocket.recvfrom(2048)

print out received string and close socket ⟶ print modifiedMessage.decode()

clientSocket.close()

# 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

# Connection-oriented demultiplexing: example



server: IP address B

host: IP address A

source IP,port: B,80
dest IP,port: A,9157

source IP,port: A,9157
dest IP,port: B,80

source IP,port: C,5775
dest IP,port: B,80

source IP,port: C,9157
dest IP,port: B,80

host: IP address C

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

Questions ?

# 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

SNMP server

application

transport
(UDP)

network (IP)

link

physical

application

transport
(UDP)

network (IP)

link

physical

# UDP: Transport Layer Actions

## SNMP client

| |
|---|
| application |
| transport (UDP) |
| network (IP) |
| link |
| physical |

## SNMP server

| |
|---|
| application   SNMP msg |
| transport (UDP)   UDP$_h$ SNMP msg |
| 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

# UDP: Transport Layer Actions

## SNMP client



## SNMP server

**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)

# UDP 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

# 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)

Questions ?

# 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



reliable service *abstraction*

# Principles of reliable data transfer



reliable service *abstraction*

reliable service *implementation*

sending process

data

application
transport

reliable channel

receiving process

data

sending process

data

application
transport

sender-side of reliable data transfer protocol

receiving process

data

receiver-side of reliable data transfer protocol

transport
network

unreliable channel

# Principles of reliable data transfer

Complexity of reliable data transfer protocol will depend (strongly) on characteristics of unreliable channel (lose, corrupt, reorder data?)

sending process

receiving process

application
transport

data

data

sender-side of reliable data transfer protocol

receiver-side of reliable data transfer protocol

transport
network

unreliable channel

reliable service *implementation*

# Principles of reliable data transfer

Sender, receiver do *not* know the "state" of each other, e.g., was a message received?

- unless communicated via a message



reliable service *implementation*

# Reliable data transfer protocol (rdt): interfaces

**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 layer

sending process

receiving process

`rdt_send()`

data

data

`deliver_data()`

data

sender-side *implementation* of `rdt` reliable data transfer protocol

receiver-side *implementation* of `rdt` reliable data transfer protocol

packet

`udt_send()`

Header data

Header data

`rdt_rcv()`

unreliable channel

**udt_send():** called by `rdt` to transfer packet over unreliable channel to receiver

**rdt_rcv():** called when packet arrives on receiver side of channel

Bi-directional communication over unreliable channel

# Reliable data transfer: getting started

## We will:

- incrementally develop sender, receiver sides of <u>r</u>eliable <u>d</u>ata <u>t</u>ransfer protocol (`rdt`)

- consider only unidirectional data transfer
  - but control info will flow in 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: reliable transfer over a reliable channel

- underlying channel perfectly reliable
  - no bit errors
  - no loss of packets

- *separate* FSMs for sender, receiver:
  - sender sends data into underlying channel
  - receiver reads data from underlying channel

sender

Wait for
call from
above

rdt_send(data)
——————————
packet = make_pkt(data)
udt_send(packet)

receiver

Wait for
call from
below

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

# rdt2.0: channel with bit errors

- underlying channel may flip bits in packet
  - checksum (e.g., Internet checksum) to detect bit errors
- *the* question: how to recover from errors?

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

# rdt2.0: channel with bit errors

- 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

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

# rdt2.0: FSM specifications

sender

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

# rdt2.0: FSM specification

sender

rdt_send(data)

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

**Wait for call from above**

rdt_rcv(rcvpkt) && isNAK(rcvpkt)

udt_send(sndpkt)

**Wait for ACK or NAK**

rdt_rcv(rcvpkt) && isACK(rcvpkt)

Λ

Note: "state" of receiver (did the receiver get my message correctly?) isn't known to sender unless somehow communicated from receiver to sender
- that's why we need a protocol!

# 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

Questions ?

# Principles of reliable data transfer



reliable service *abstraction*
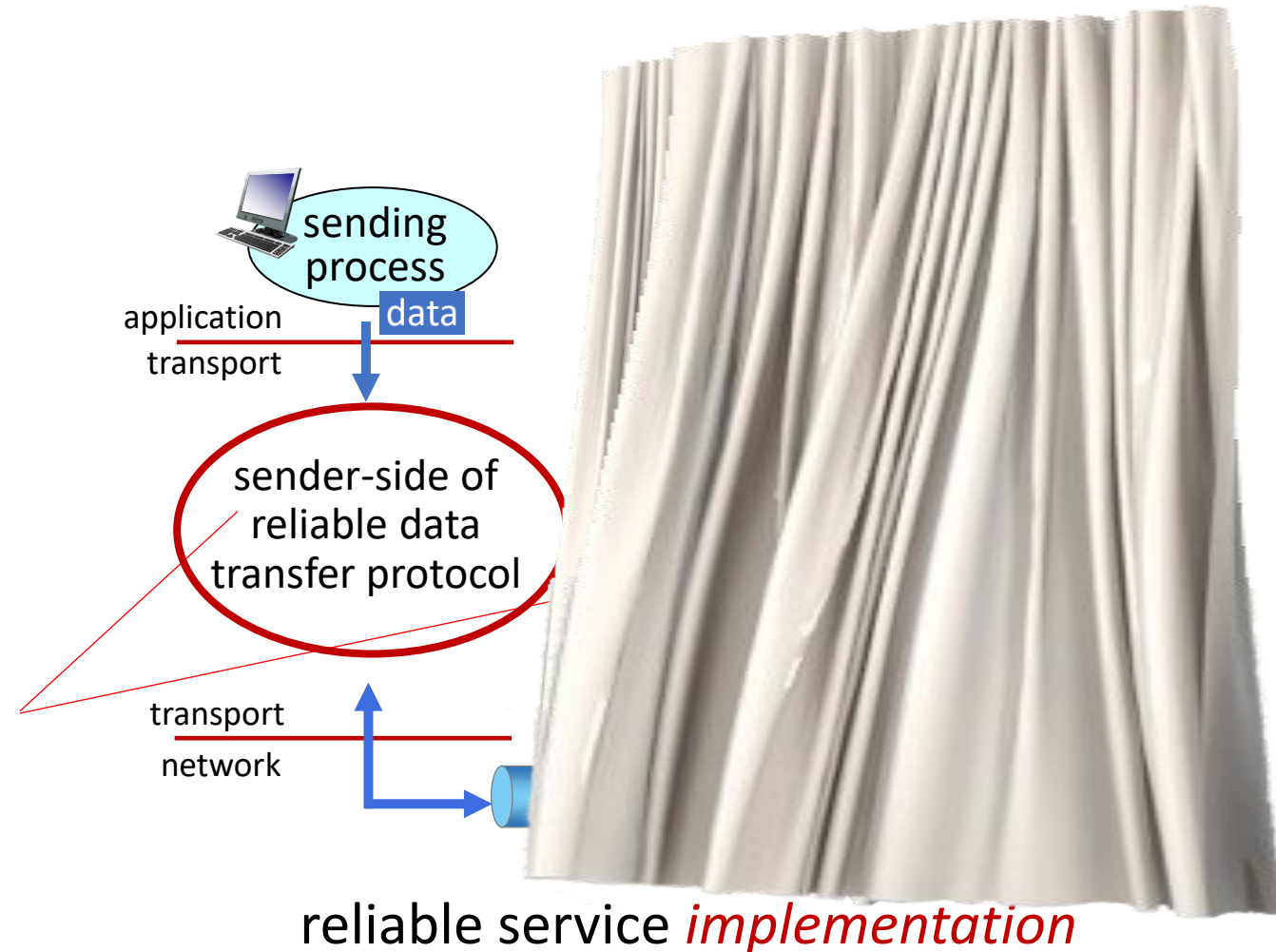
# Principles of reliable data transfer

reliable service *abstraction*

reliable service *implementation*

# Principles of reliable data transfer

Complexity of reliable data transfer protocol will depend (strongly) on characteristics of unreliable channel (lose, corrupt, reorder data?)



reliable service *implementation*

# Principles of reliable data transfer



Sender, receiver do *not* know the "state" of each other, e.g., was a message received?
- unless communicated via a message

reliable service *implementation*

# Reliable data transfer protocol (rdt): interfaces

**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 layer

sending process

receiving process

`rdt_send()`

data

data

`deliver_data()`

data

sender-side *implementation* of `rdt` reliable data transfer protocol

receiver-side *implementation* of `rdt` reliable data transfer protocol

packet

`udt_send()`

Header data

Header data

`rdt_rcv()`

unreliable channel

**udt_send():** called by `rdt` to transfer packet over unreliable channel to receiver

Bi-directional communication over unreliable channel

**rdt_rcv():** called when packet arrives on receiver side of channel

# Reliable data transfer: getting started

## We will:

- incrementally develop sender, receiver sides of <u>r</u>eliable <u>d</u>ata <u>t</u>ransfer protocol (`rdt`)

- consider only unidirectional data transfer
  - but control info will flow in 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: reliable transfer over a reliable channel

- underlying channel perfectly reliable
  - no bit errors
  - no loss of packets

- *separate* FSMs for sender, receiver:
  - sender sends data into underlying channel
  - receiver reads data from underlying channel

sender

Wait for call from above

rdt_send(data)
_____
packet = make_pkt(data)
udt_send(packet)

receiver

Wait for call from below

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

# rdt2.0: channel with bit errors

- underlying channel may flip bits in packet
  - checksum (e.g., Internet checksum) to detect bit errors
- *the* question: how to recover from errors?

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

# rdt2.0: channel with bit errors

- 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

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

# rdt2.0: FSM specifications

sender

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

# rdt2.0: FSM specification

sender

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

Note: "state" of receiver (did the receiver get my message correctly?) isn't known to sender unless somehow communicated from receiver to sender
- that's why we need a protocol!

# rdt2.0: operation with no errors

rdt_send(data)

---

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

**sender**

Wait for call from above

Wait for ACK or NAK

rdt_rcv(rcvpkt) && isNAK(rcvpkt)

---

udt_send(sndpkt)

rdt_rcv(rcvpkt) && isACK(rcvpkt)

---

$\Lambda$

rdt_rcv(rcvpkt) && corrupt(rcvpkt)

---

udt_send(NAK)

Wait for call from below

**receiver**

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)

---

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

# rdt2.0: corrupted packet scenario

sender

Wait for call from above

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

Wait for ACK or NAK

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

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

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 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 corrupted
- 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

# rdt2.1: sender, handling 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)
_____
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, handling 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)

# 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 "expected" pkt should have seq # of 0 or 1

## 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

- duplicate ACK at sender results in same action as NAK: *retransmit current pkt*

As we will see, TCP uses this approach to be NAK-free

# 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)

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

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

- checksum, sequence #s, ACKs, retransmissions will be of help … but not quite enough

*Q:* How do *humans* handle lost sender-to-receiver words in conversation?

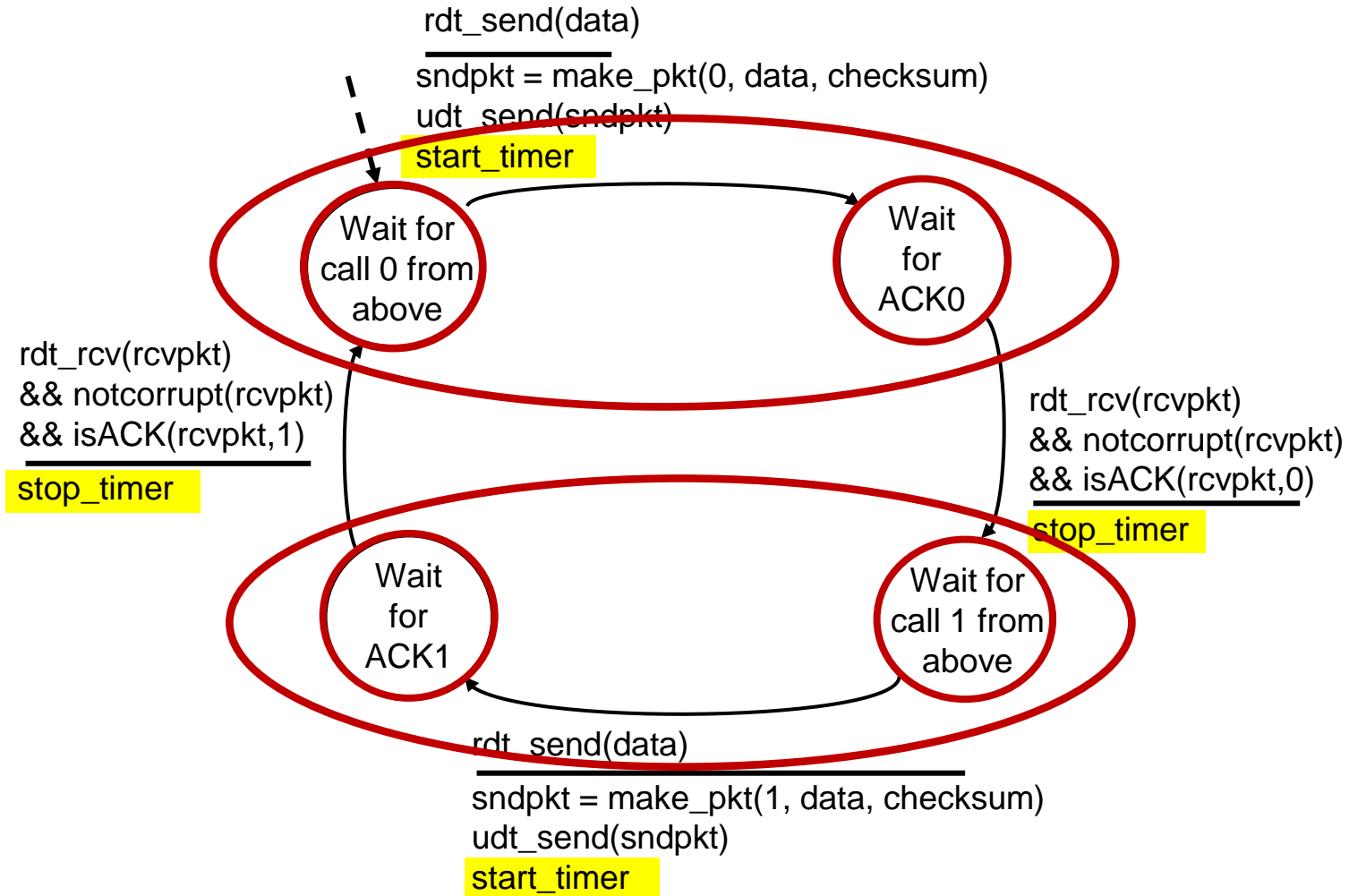# rdt3.0: channels with errors *and* loss

*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 seq #s already handles this!
  - receiver must specify seq # of packet being ACKed
- use countdown timer to interrupt after "reasonable" amount of time

*timeout*

# rdt3.0 sender

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

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& isACK(rcvpkt,1)
—————————————
stop_timer

Wait for call 0 from above

Wait for ACK0

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& isACK(rcvpkt,0)
—————————————
stop_timer

Wait for ACK1

Wait for call 1 from above

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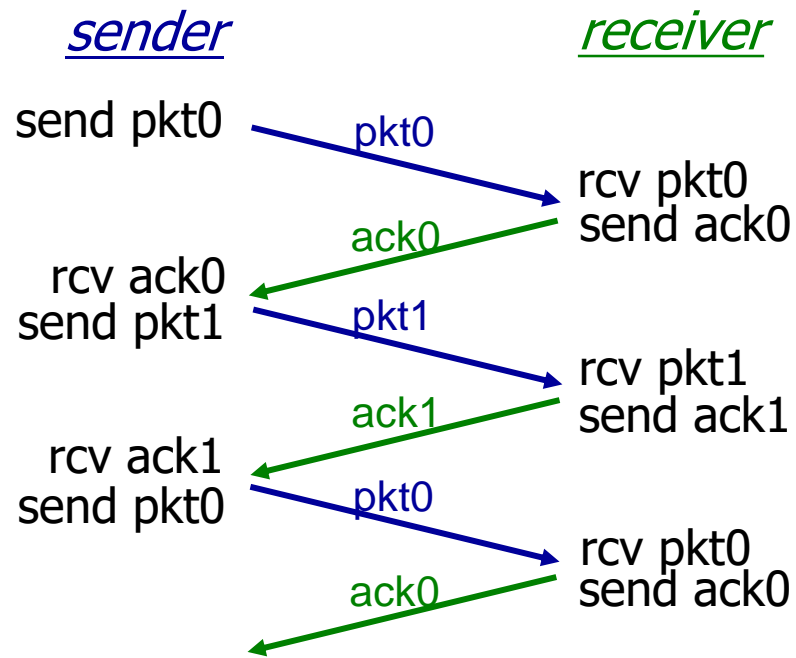sndpkt = make_pkt(1, data, checksum)
udt_send(sndpkt)
start_timer

# rdt3.0 sender

# rdt3.0 in action



(a) no loss



(b) packet loss

# rdt3.0 in action



(c) ACK loss

(d) premature timeout/ delayed ACK

# Performance of rdt3.0 (stop-and-wait)

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

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

  - time to transmit packet into channel:

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

# rdt3.0: stop-and-wait operation

sender                          receiver

first packet bit transmitted, t = 0

RTT

first packet bit arrives
last packet bit arrives, send ACK

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

# rdt3.0: stop-and-wait operation

$$U_{sender} = \frac{L \: / \: R}{RTT + L \: / \: R}$$

$$= \frac{.008}{30.008}$$

$$= 0.00027$$

sender       receiver

L/R

RTT

- rdt 3.0 protocol performance stinks!
- Protocol limits performance of underlying infrastructure (channel)

# rdt3.0: pipelined protocols operation

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

- range of sequence numbers must be increased
- buffering at sender and/or receiver



data packet

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

# Pipelining: increased utilization

sender                    receiver

first packet bit transmitted, $t = 0$

last bit transmitted, $t = L / R$

RTT

first packet bit arrives

last packet bit arrives, send ACK

last bit of 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!

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

# Go-Back-N: sender

- sender: "window" of up to N, consecutive transmitted but unACKed pkts
  - k-bit seq # in pkt header



- *cumulative ACK:* ACK($n$): ACKs all packets up to, including seq # $n$
  - on receiving ACK($n$): move window forward to begin at $n+1$
- timer for oldest in-flight packet
- *timeout(n):* retransmit packet n and all higher seq # packets in window

# Go-Back-N: receiver

- ACK-only: always send ACK for correctly-received packet so far, with highest *in-order* seq #
  - may generate duplicate ACKs
  - need only remember `rcv_base`

- on receipt of out-of-order packet:
  - can discard (don't buffer) or buffer: an implementation decision
  - re-ACK pkt with highest in-order seq #

Receiver view of sequence number space:



... ...

rcv_base

received and ACKed

Out-of-order: received but not ACKed

Not received

# Go-Back-N 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, discard,
                                          (re)send ack1

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, discard,
                                          (re)send ack1

            ignore duplicate ACK    receive pkt5, discard,
                                          (re)send ack1

*pkt 2 timeout*

0 1 2 3 4 5 6 7 8     send  pkt2

0 1 2 3 4 5 6 7 8     send  pkt3

0 1 2 3 4 5 6 7 8     send  pkt4         rcv pkt2, deliver, send ack2

0 1 2 3 4 5 6 7 8     send  pkt5         rcv pkt3, deliver, send ack3

                                          rcv pkt4, deliver, send ack4

                                          rcv pkt5, deliver, send ack5

# Selective repeat

- receiver *individually* acknowledges all correctly received packets
  - buffers packets, as needed, for eventual in-order delivery to upper layer

- sender times-out/retransmits individually for unACKed packets
  - sender maintains timer for each unACKed pkt

- sender window
  - *N* consecutive seq #s
  - limits seq #s of sent, unACKed packets

# Selective repeat: sender, receiver windows



(a) sender view of sequence numbers

# Selective repeat: sender and receiver

**sender**

data from above:

- if next available seq # in window, send packet

timeout(*n*):

- resend packet *n*, restart timer

ACK(*n*) in [sendbase,sendbase+N]:

- mark packet *n* as received
- if n smallest unACKed packet, advance window base to next unACKed seq #

**receiver**

packet *n* in [rcvbase, rcvbase+N-1]

- send ACK(*n*)
- out-of-order: buffer
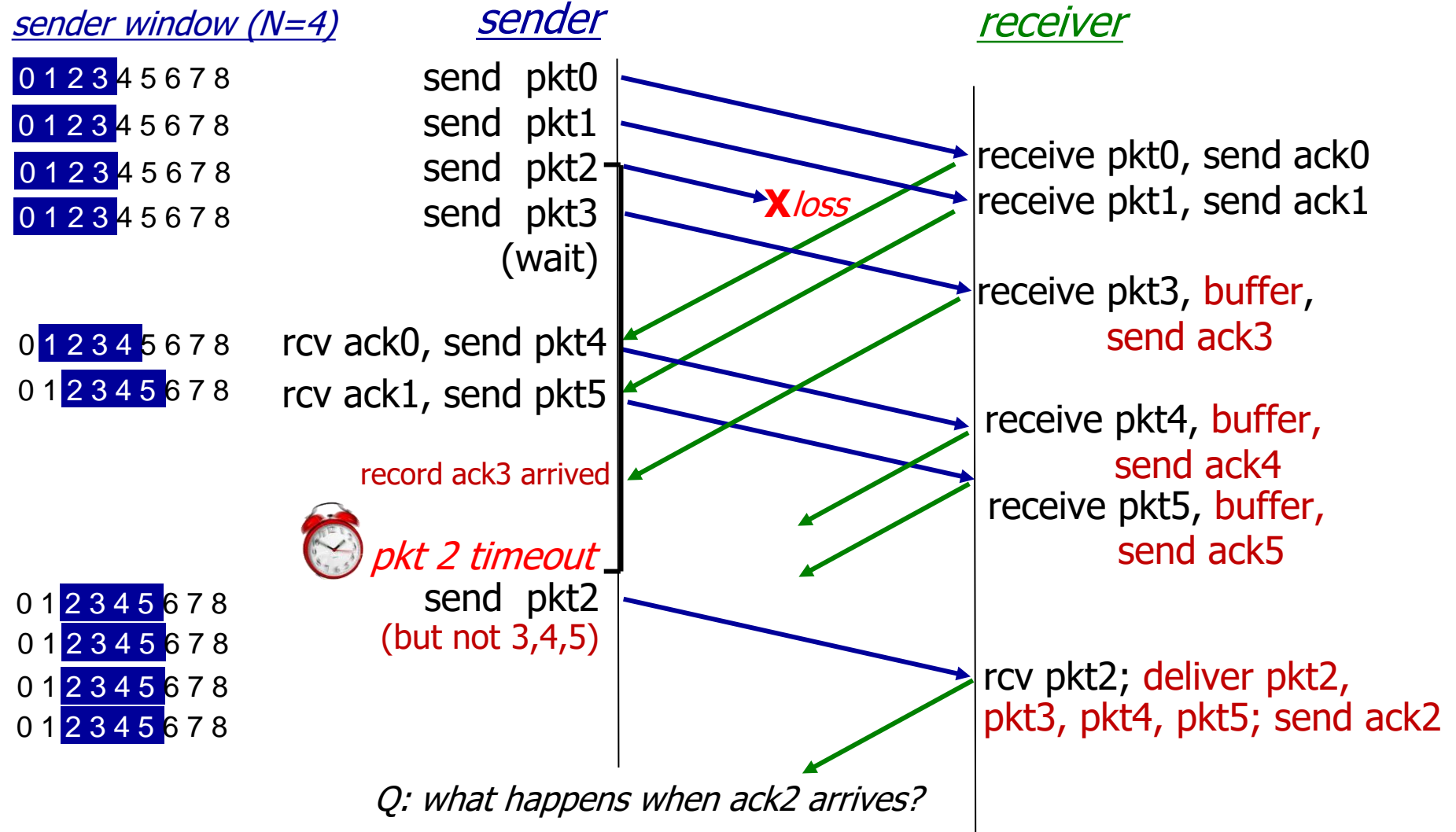- in-order: deliver (also deliver buffered, in-order packets), advance window to next not-yet-received packet

packet *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

0 1 2 3 4 5 6 7 8     send  pkt3

(wait)

X *loss*

receive pkt0, send ack0
receive pkt1, send ack1

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

record ack3 arrived

*pkt 2 timeout*

receive pkt4, buffer,
      send ack4
receive pkt5, buffer,
      send ack5

0 1 2 3 4 5 6 7 8     send  pkt2
(but not 3,4,5)

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

0 1 2 3 4 5 6 7 8

rcv pkt2; deliver pkt2,
pkt3, pkt4, pkt5; send ack2

*Q: what happens when ack2 arrives?*

# Selective repeat: a dilemma!

example:

- seq #s: 0, 1, 2, 3 (base 4 counting)
- window size=3



sender window (after receipt)

receiver 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

0 1 2 3 0 1 2   pkt3   X
0 1 2 3 0 1 2
                pkt0

0 1 2 3 0 1 2
0 1 2 3 0 1 2
0 1 2 3 0 1 2

*will accept packet with seq number 0*

(a) no problem

0 1 2 3 0 1 2   pkt0
0 1 2 3 0 1 2   pkt1
0 1 2 3 0 1 2   pkt2   X
                      X
                      X

timeout
retransmit pkt0
0 1 2 3 0 1 2   pkt0

0 1 2 3 0 1 2
0 1 2 3 0 1 2
0 1 2 3 0 1 2

*will accept packet with seq number 0*
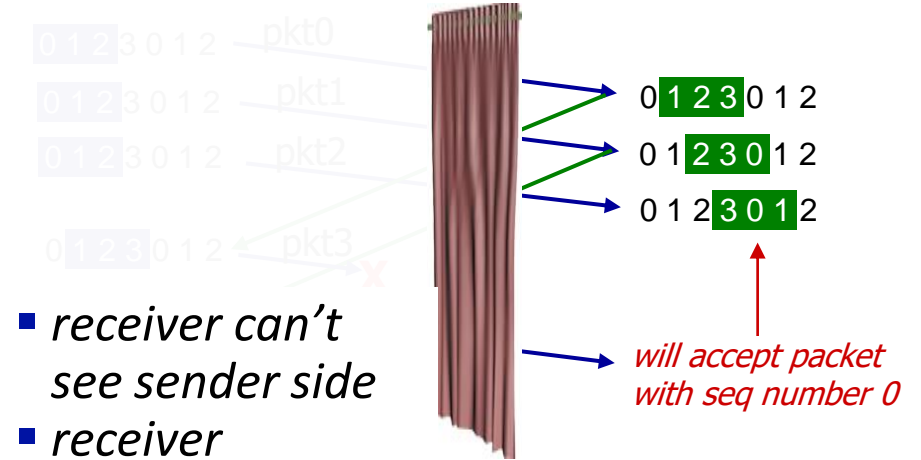
(b) oops!

# Selective repeat: a dilemma!

example:

- seq #s: 0, 1, 2, 3 (base 4 counting)
- window size=3

Q: what relationship is needed between sequence # size and window size to avoid problem in scenario (b)?

pkt0

0 1 2 3 0 1 2   pkt1

0 1 2 3 0 1 2   pkt2

0 1 2 3 0 1 2

0 1 2 3 0 1 2   pkt3

0 1 2 3 0 1 2

0 1 2 3 0 1 2

0 1 2 3 0 1 2

- *receiver can't see sender side*
- *receiver behavior identical in both cases!*
- *something's (very) wrong!*

*will accept packet with seq number 0*

0 1 2 3 0 1 2

0 1 2 3 0 1 2

0 1 2 3 0 1 2

0 1 2 3 0 1 2   pkt2

timeout
retransmit pkt0
0 1 2 3 0 1 2   pkt0

(b) oops!

*will accept packet with seq number 0*

# Chapter 3: roadmap

- Transport-layer services
- Multiplexing and demultiplexing
- Connectionless transport: UDP
- Principles of reliable data transfer
- **Connection-oriented transport: TCP**
  - segment structure
  - reliable data transfer
  - flow control
  - connection management
- Principles of congestion control
- TCP congestion control