# CSCI 6760/4760 Computer Networks: Topology and applications

Chapter 3

Manijeh Keshtgari

Computer Science
UGA
M.keshtgari@uga.edu

# Chapter 3: Transport Layer

## our goals:

- ❖ 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

# Chapter 3 outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

3.4 principles of reliable data transfer

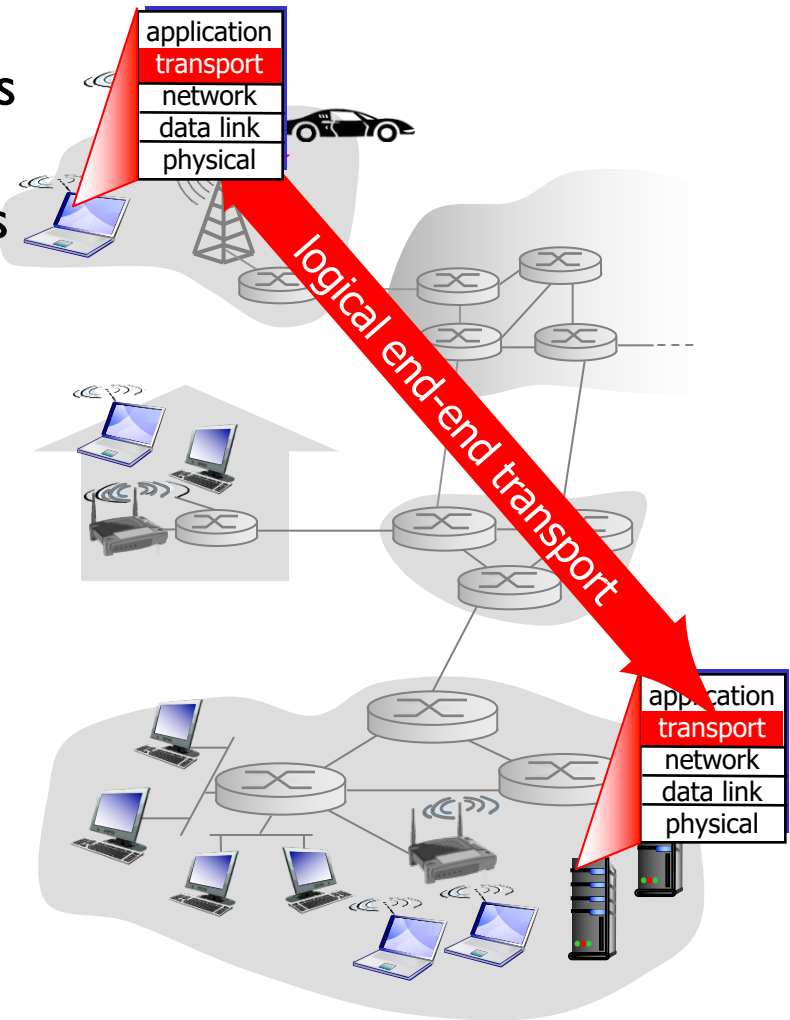3.5 connection-oriented transport: TCP
- segment structure
- reliable data transfer
- flow control
- connection management

3.6 principles of congestion control

3.7 TCP congestion control

# Transport services and protocols

* ❖ provide *logical communication* between app processes running on different hosts
* ❖ By *logical communication*, we mean that from an application's perspective, it is as if the hosts running the processes were directly connected;
* ❖ in reality, the hosts may be on opposite sides of the planet, 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

application
transport
network
data link
physical

logical end-end transport

application
transport
network
data link
physical

# Transport vs. network layer

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

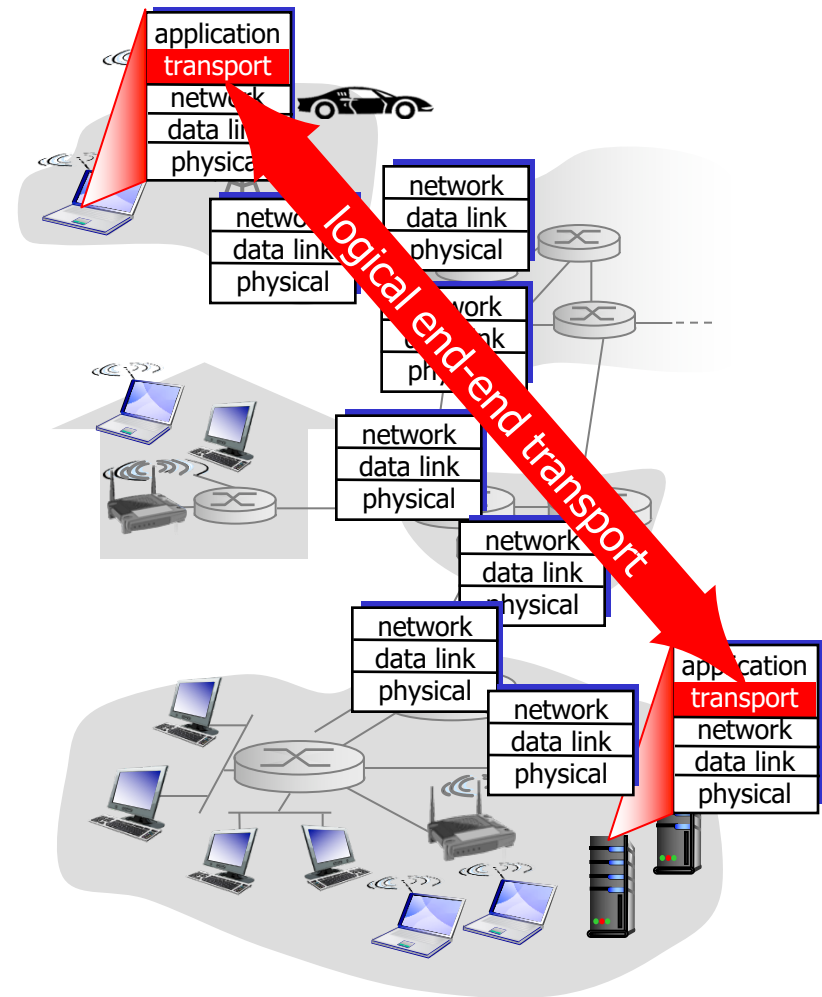*household analogy:*

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 protocol = Ann and Bill who demux to in-house siblings
* 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

# Chapter 3 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 principles of congestion control

3.7 TCP congestion control
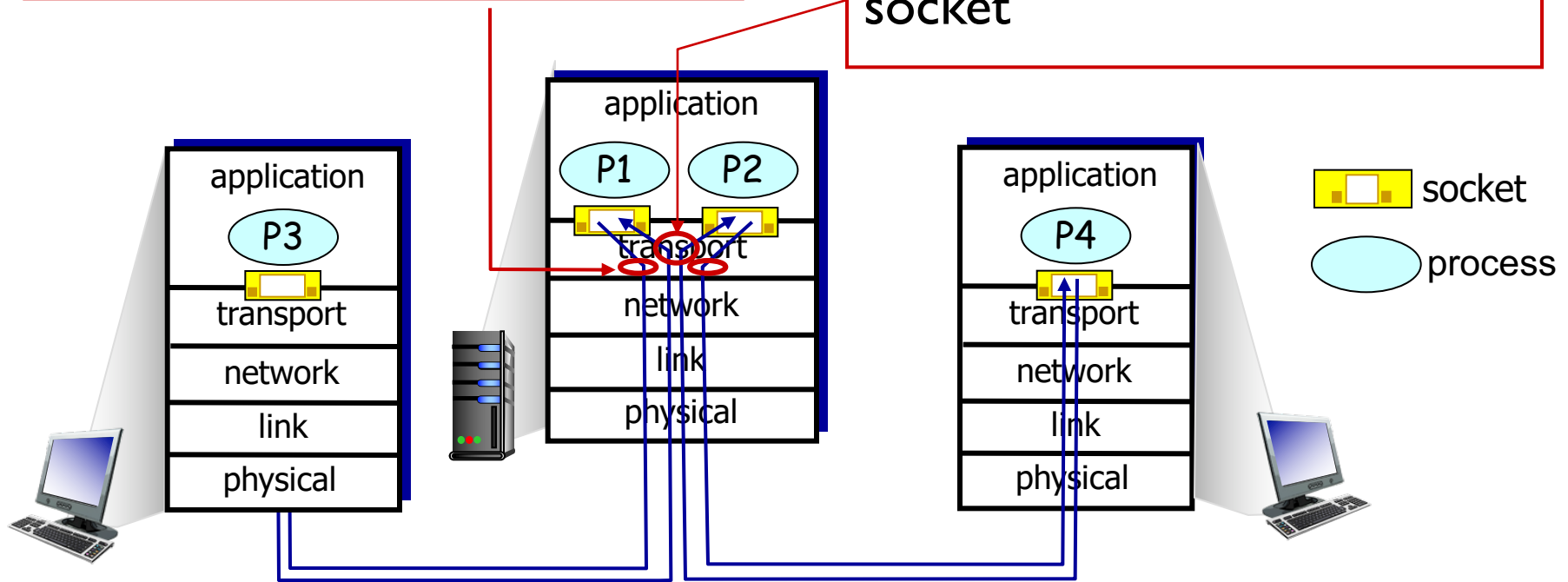
# Multiplexing/demultiplexing

*multiplexing at sender:*

handle data from multiple sockets, add transport header (later used for demultiplexing)
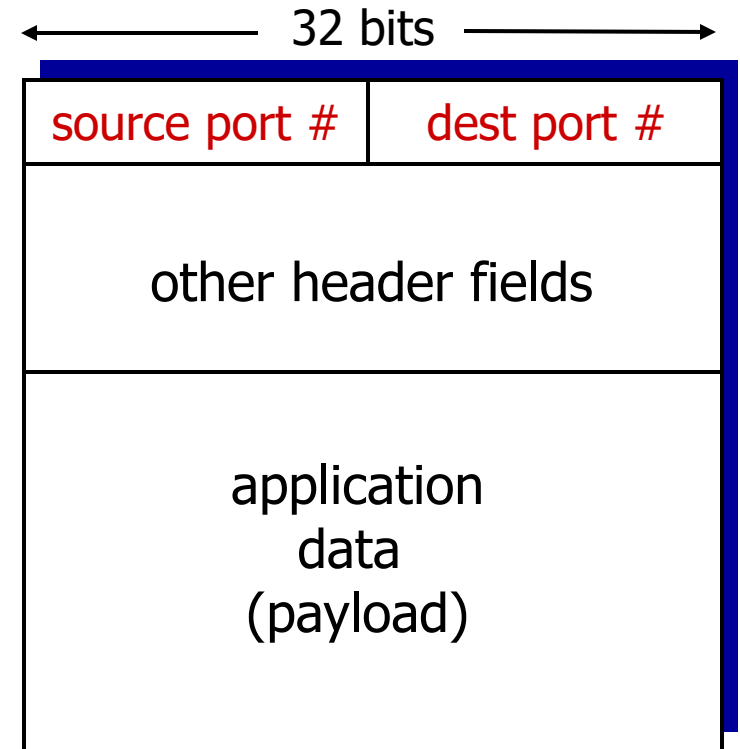
*demultiplexing at receiver:*

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

❖ created socket has host-local port #:

❖ when creating datagram to send into UDP socket, must specify
  ▪ destination IP address
  ▪ destination port #

❖ when host receives UDP segment:
  ▪ checks destination port # in segment
  ▪ directs UDP segment to socket with that port #

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

# Connection-oriented demux

❖ TCP socket identified by 4-tuple:
  ▪ source IP address
  ▪ source port number
  ▪ dest IP address
  ▪ dest port number
❖ demux: receiver 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: example



application

P4   P5   P6

application
P3

transport
network
link
physical

application
P2   P3

transport
network
link
physical

server: IP
address B

transport
network
link
physical

host: IP
address A

host: IP
address C

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

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

# Connection-oriented demux: example

threaded server

application

P4

transport

network

link

physical

server: IP address B

application

P3

transport

network

link

physical

host: IP address A

application

P2      P3

transport

network

link

physical

host: IP address C

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

# Chapter 3 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
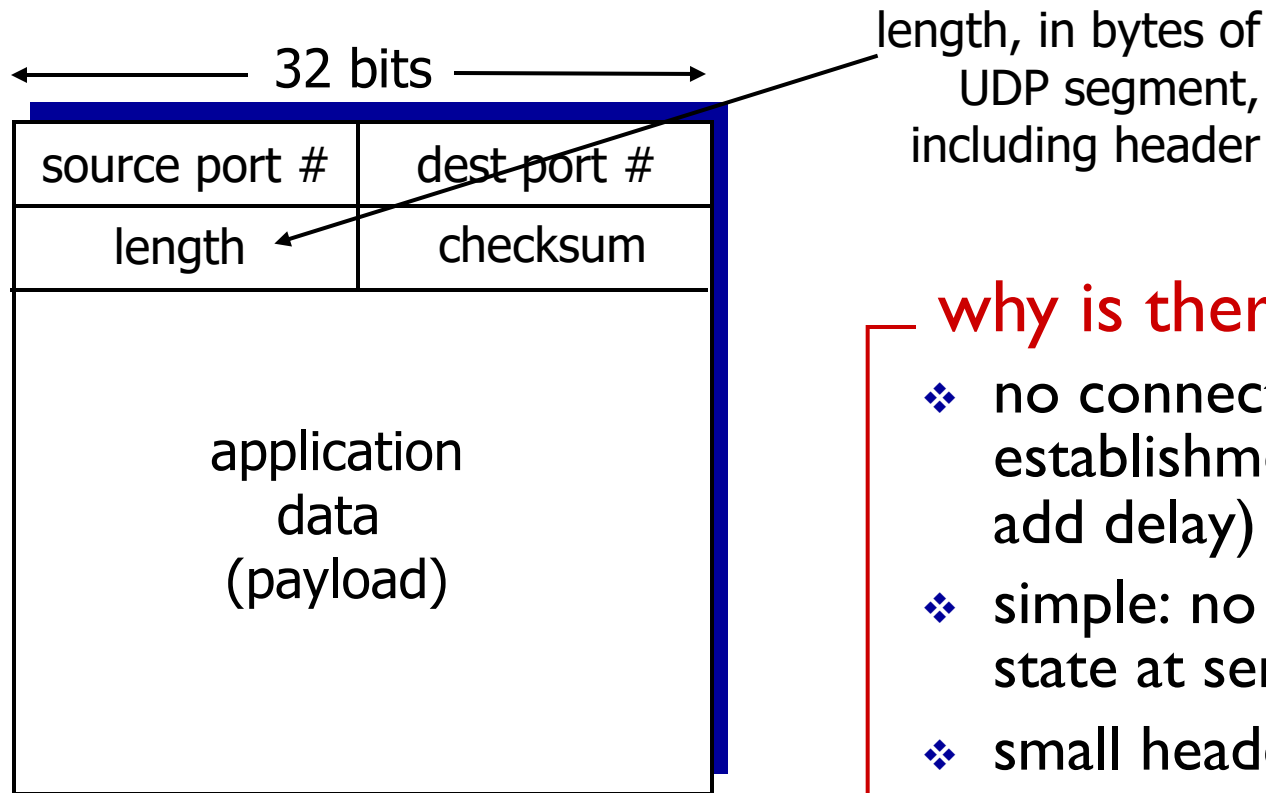  - connection management

3.6 principles of congestion control

3.7 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

- ❖ UDP use:
  - ▪ streaming multimedia apps (loss tolerant, rate sensitive)
  - ▪ DNS
- ❖ reliable transfer over UDP:
  - ▪ add reliability at application layer
  - ▪ application-specific error recovery!

# UDP: segment header

length, in bytes of
UDP segment,
including header

<------ 32 bits ------>

| source port # | dest port # |
|---|---|
| length | checksum |
| application data (payload) | |

UDP segment format

## why is there a UDP?

- ❖ no connection establishment (which can add delay)
- ❖ simple: no connection state at sender, receiver
- ❖ small header size
- ❖ no congestion control: UDP can blast away as fast as desired

# why is there a UDP?

❖ Finer control over what data is sent and when

  ▪ As soon as an application process writes into the socket

  ▪ … UDP will package the data and send the packet

❖ no delay for connection establishment

❖ no connection state at sender, receiver

  ▪ No allocation of buffers, parameters, sequence #s, etc.

  ▪ … making it easier to handle many active clients at once

❖ small header size

# Popular Applications That Use UDP

❖ **Multimedia streaming**
  - By the time the packet is retransmitted, it's too late
  - E.g., telephone calls, video conferencing, gaming

❖ **Simple query protocols like Domain Name System**

# UDP checksum

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

## sender:

- ❖ treat segment contents, including header fields, as sequence of 16-bit integers
- ❖ checksum: addition (one'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?* More later ….

# Internet checksum: 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  ⟨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

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

# Chapter 3 outline

3.1 transport-layer services

3.2 multiplexing and demultiplexing

3.3 connectionless transport: UDP

<span style="color:red">3.4 principles of reliable data transfer</span>

3.5 connection-oriented transport: TCP
- segment structure
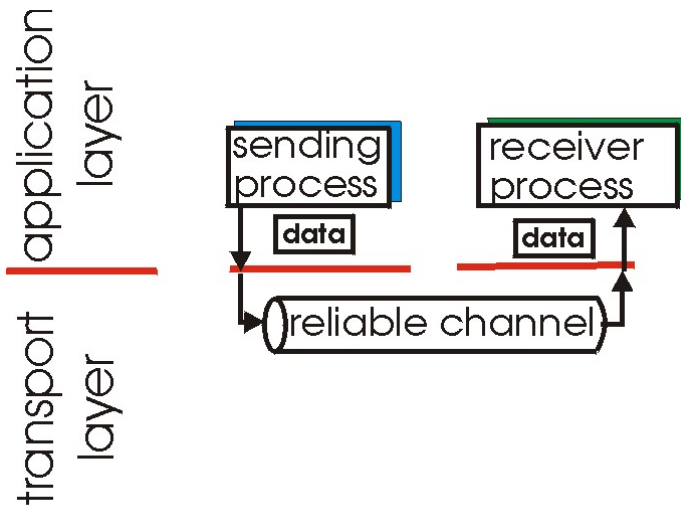- reliable data transfer
- flow control
- connection management

3.6 principles of congestion control

3.7 TCP congestion control

# Principles of reliable data transfer

❖ **important in application, transport, link layers**
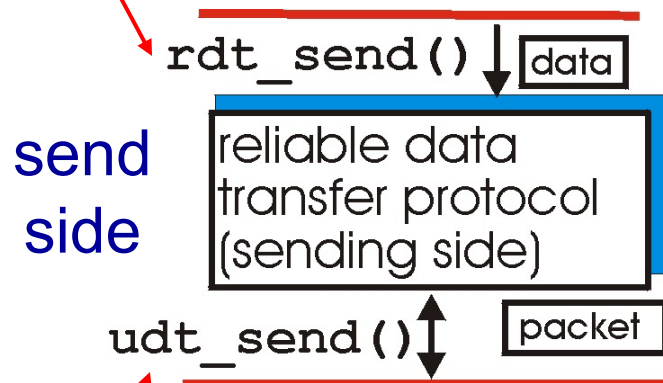  ▪ top-10 list of important networking topics!



(a) provided service

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

# 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

**we'll:**

❖ 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 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 errors (no flipped from 0 to 1, or vice versa)
  ▪ no loss of packets
  ▪ All the packets are delivered in the order that were sent
❖ **separate FSMs for sender, receiver:**
  ▪ sender sends data into underlying channel
  ▪ receiver reads 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: channel with bit errors

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

# rdt2.0: FSM specification

rdt_send(data)
_____
sndpkt = 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)

**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**

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

# 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!
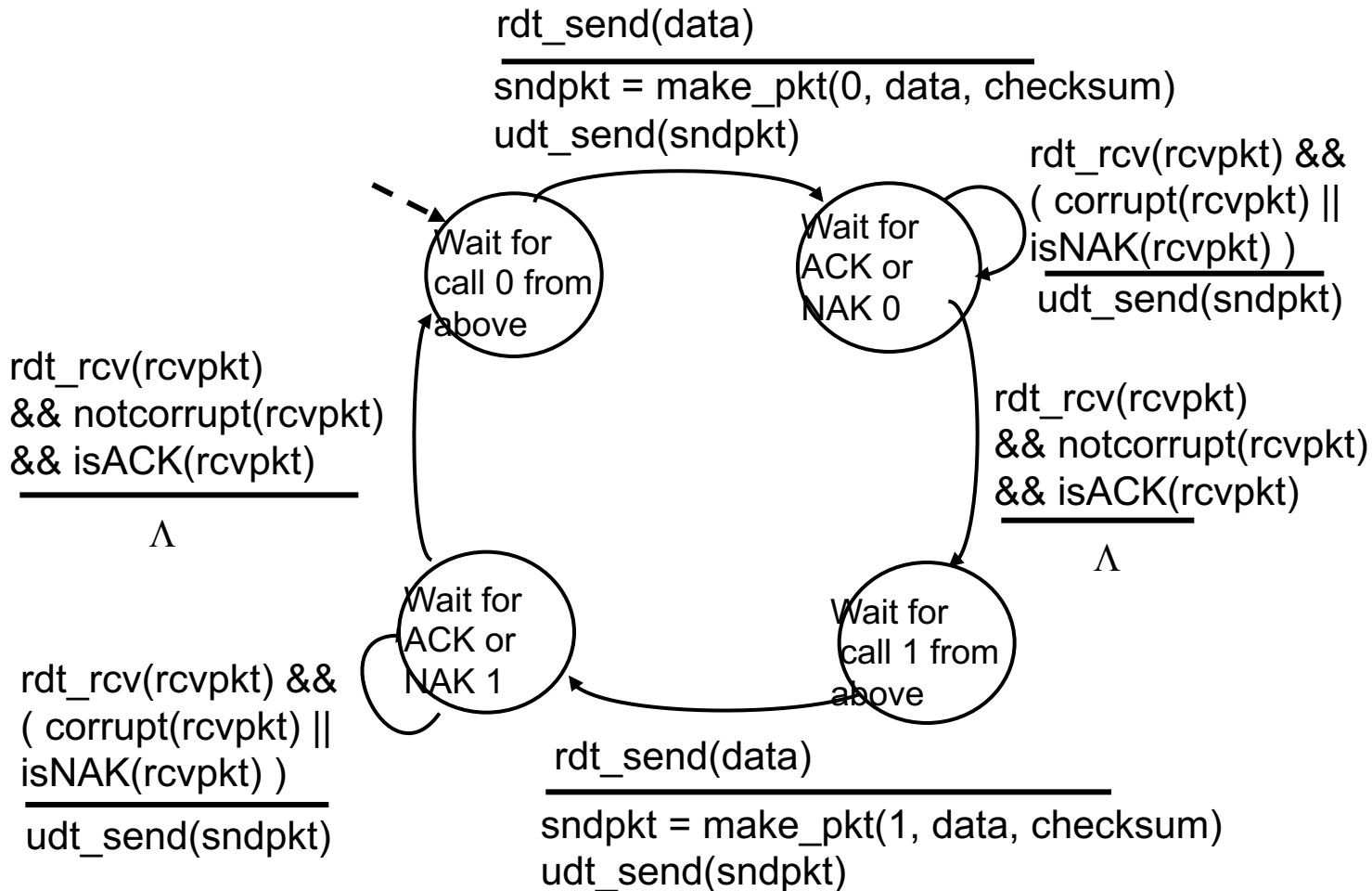
❖ can't just retransmit: possible duplicate

**handling duplicates:**

❖ sender retransmits current pkt if ACK/NAK 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, 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)
_____
$\Lambda$

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)

# 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
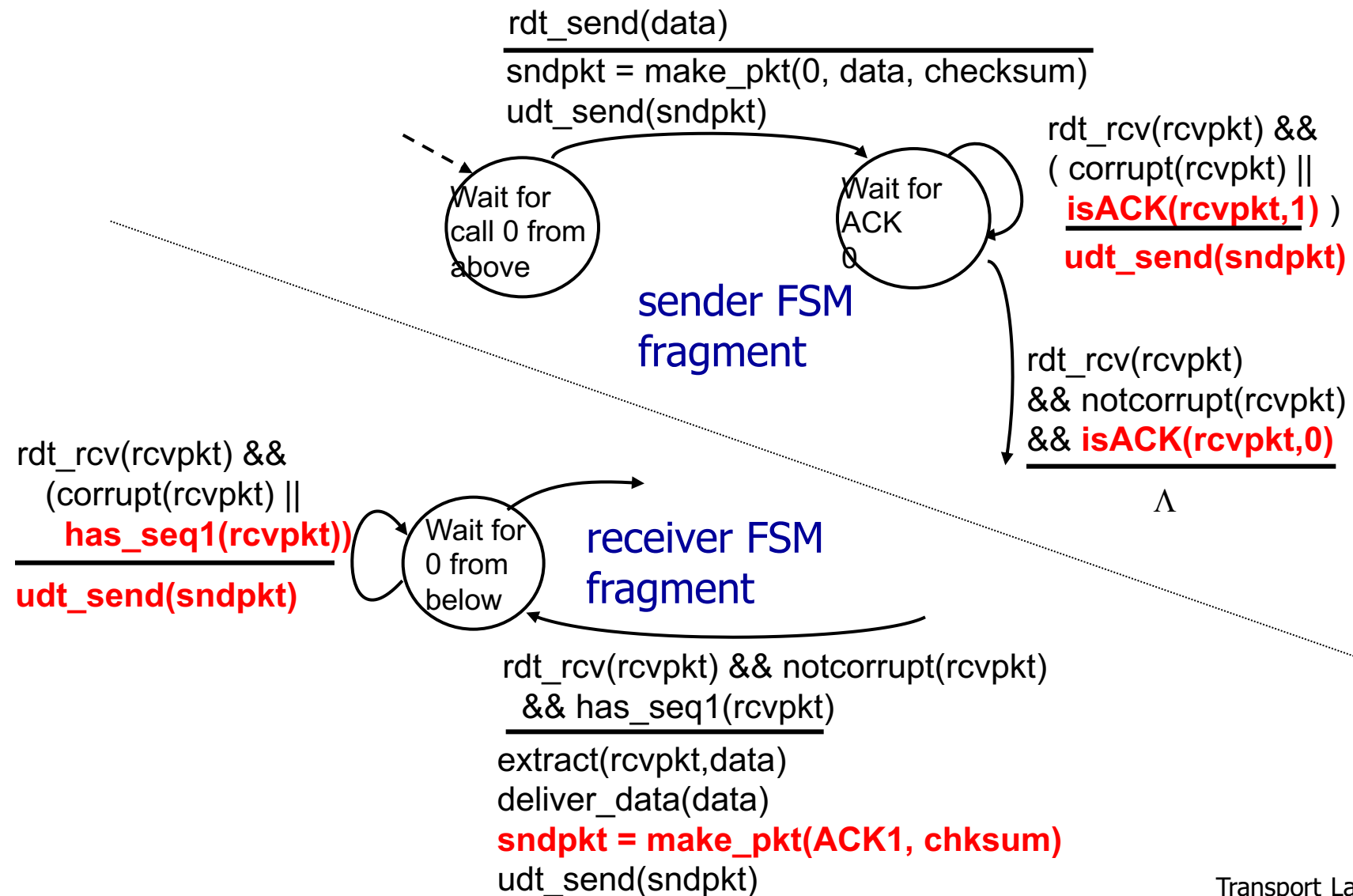  - ▪ state must "remember" whether "expected" pkt should have seq # of 0 or 1

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*

# 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)**
—————————————————————————————
$\Lambda$

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, ACKs)

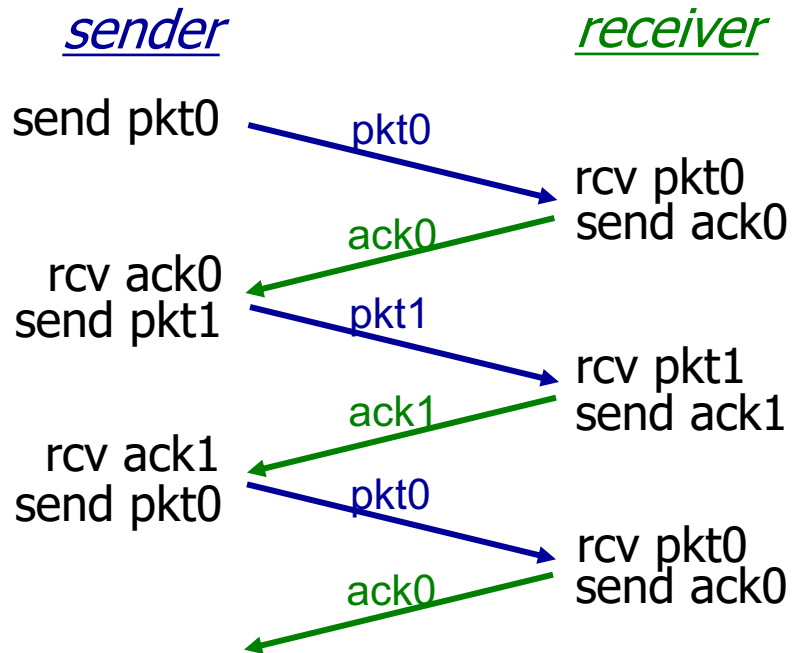- checksum, seq. #, ACKs, retransmissions will be of help … but not enough

**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 seq. #'s already handles this
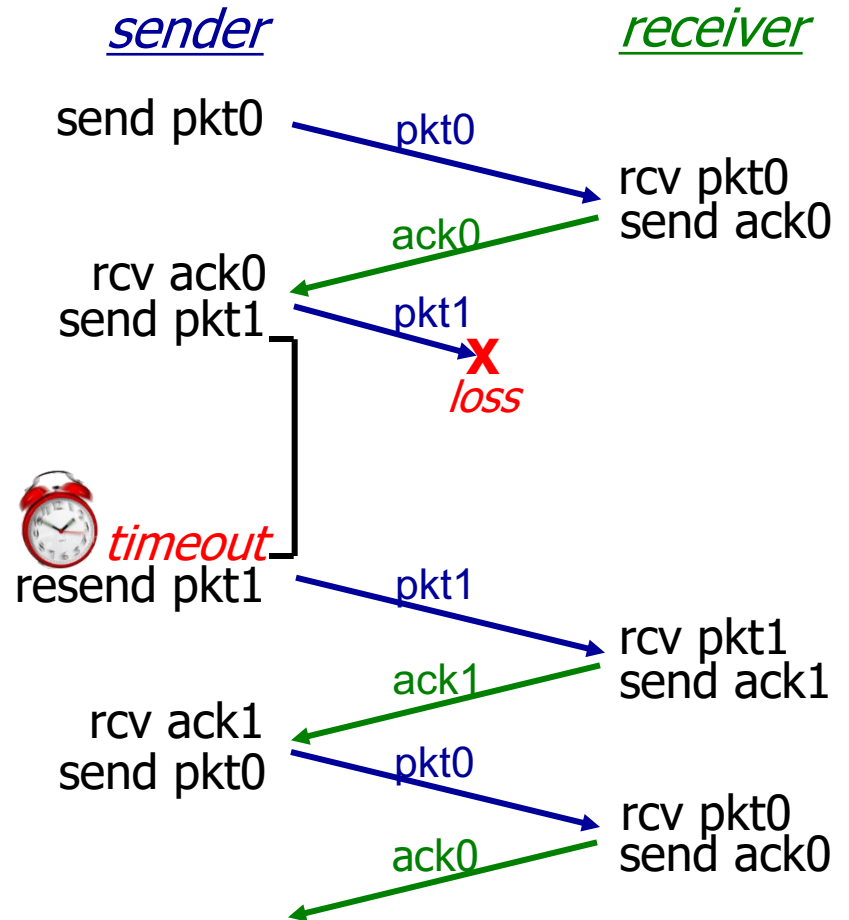  - receiver must specify seq # of pkt being ACKed

❖ requires countdown timer

# rdt3.0 sender

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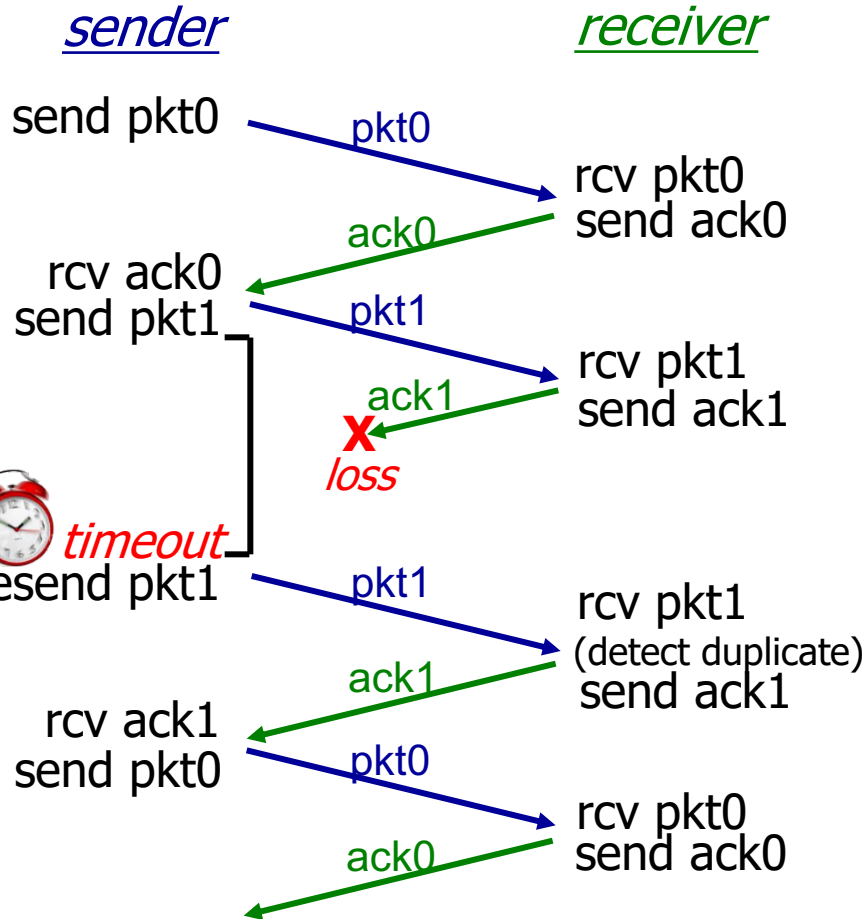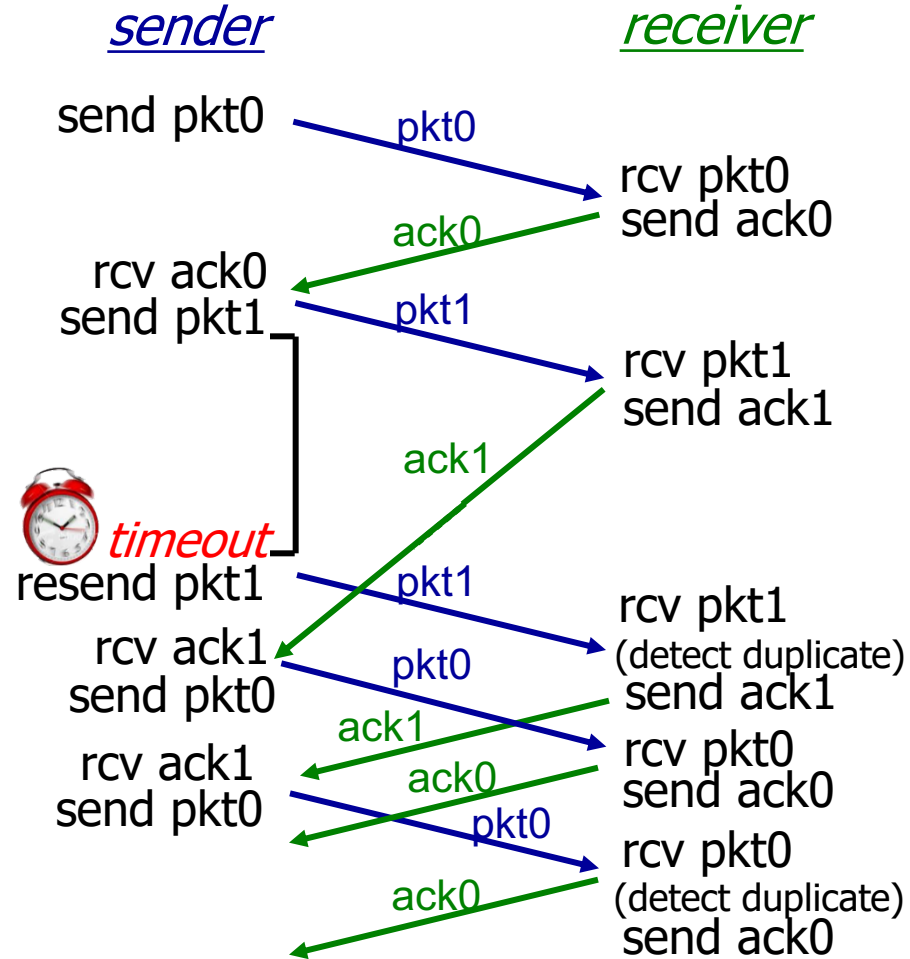
sndpkt = make_pkt(0, data, checksum)
udt_send(sndpkt)
start_timer

rdt_rcv(rcvpkt)

$\Lambda$

**Wait for call 0 from above**

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

$\Lambda$

**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)

$\Lambda$

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

$\Lambda$

rdt_send(data)

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

# rdt3.0 in action

**sender**          **receiver**

send pkt0 ——— pkt0 ———→ rcv pkt0
                             send ack0
rcv ack0 ←——— ack0 ———
send pkt1 ——— pkt1 ———→ rcv pkt1
                             send ack1
rcv ack1 ←——— ack1 ———
send pkt0 ——— pkt0 ———→ rcv pkt0
                             send ack0
        ←——— ack0 ———

(a) no loss

**sender**          **receiver**

send pkt0 ——— pkt0 ———→ rcv pkt0
                             send ack0
rcv ack0 ←——— ack0 ———
send pkt1 ——— pkt1 ———→ **X**
                             *loss*

*timeout*
resend pkt1 ——— pkt1 ———→ rcv pkt1
                             send ack1
rcv ack1 ←——— ack1 ———
send pkt0 ——— pkt0 ———→ rcv pkt0
                             send ack0
        ←——— ack0 ———

(b) packet loss

# rdt3.0 in action

**sender**      **receiver**

send pkt0
→ pkt0 →
    rcv pkt0
    send ack0

rcv ack0
← ack0 ←
send pkt1
→ pkt1 →
    rcv pkt1
    send ack1

ack1
**X**
*loss*

*timeout*
resend pkt1
→ pkt1 →
    rcv pkt1
    (detect duplicate)
← ack1 ←
    send ack1

rcv ack1
send pkt0
→ pkt0 →
    rcv pkt0
    send ack0
← ack0 ←

(c) ACK loss

**sender**      **receiver**

send pkt0
→ pkt0 →
    rcv pkt0
    send ack0

rcv ack0
← ack0 ←
send pkt1
→ pkt1 →
    rcv pkt1
    send ack1

ack1

*timeout*
resend pkt1
→ pkt1 →
    rcv pkt1
    (detect duplicate)
rcv ack1
send pkt0
→ pkt0 →
    send ack1
← ack1 ←
    rcv pkt0
rcv ack1
send pkt0
    send ack0
← ack0 ←
→ pkt0 →
    rcv pkt0
    (detect duplicate)
    send ack0
← ack0 ←

(d) premature timeout/ delayed ACK

# Performance of rdt3.0

❖ rdt3.0 is correct, but performance is not good

❖ e.g.: 1 Gbps link, 15 ms prop. Delay (RTT=30 ms), 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, 1KB pkt every 30 msec: 33kB/sec thruput over 1 Gbps 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
RTT                                                     last packet bit arrives, send ACK

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



(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

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 has timer for oldest unacked packet
  - when timer expires, retransmit *all* unacked packets

## Selective Repeat:

- sender can have up to N unack'ed packets in pipeline
- rcvr sends *individual ack* for each packet
- sender maintains timer for each unacked packet
  - when timer expires, retransmit only that unacked packet

# 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 receive duplicate ACKs (see receiver)
❖ timer for oldest in-flight pkt
❖ *timeout(n):* retransmit packet n and all higher seq # pkts in window

# 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

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)

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

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

Wait

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 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 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
- limits seq #s of sent, unACKed pkts

# Selective repeat: sender, receiver windows



(a) sender view of sequence numbers

- already ack'ed
- sent, not yet ack'ed
- usable, not yet sent
- not usable

(b) receiver view of sequence numbers

- out of order (buffered) but already ack'ed
- Expected, not yet received
- acceptable (within window)
- 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) (this pkt has previously acked)

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        **X** *loss*

`0 1 2 3` 4 5 6 7 8     send pkt3

                      (wait)

                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

                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,

0 1 `2 3 4 5` 6 7 8                pkt3, pkt4, pkt5; send ack2

     *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!

- duplicate data accepted as new in (b)

Q: what relationship between seq # size and window size to avoid problem in (b)?



sender window (after receipt)    receiver window (after receipt)

pkt0
pkt1
pkt2
pkt3
pkt0

(a) no problem

will accept packet with seq number 0

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

pkt0
pkt1
pkt2

timeout retransmit pkt0
pkt0

(b) oops!

will accept packet with seq number 0

# Chapter 3 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 principles of congestion control

3.7 TCP congestion control

# TCP: Overview RFCs: 793,1122,1323, 2018, 2581

- ❖ point-to-point:
  - one sender, one receiver
- ❖ reliable, in-order *byte steam:*
  - no "message boundaries"
- ❖ pipelined:
  - TCP congestion and flow control set window size

- ❖ full duplex data:
  - bi-directional data flow in same connection
  - MSS: maximum segment size
- ❖ connection-oriented:
  - handshaking (exchange of control msgs) inits sender, receiver state before data exchange
- ❖ flow controlled:
  - sender will not overwhelm receiver

# TCP segment structure

32 bits

URG: urgent data
(generally not used)

ACK: ACK #
valid

PSH: push data now
(generally not used)

RST, SYN, FIN:
connection estab
(setup, teardown
commands)

Internet
checksum
(as in UDP)

| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |

| head len | not used | U | A | P | R | S | F | receive window |
|---|---|---|---|---|---|---|---|---|

| checksum | Urg data pointer |
|---|---|

options (variable length)

application
data
(variable length)

counting
by bytes
of data
(not segments!)

# bytes
rcvr willing
to accept

# TCP seq. numbers

❖ Suppose that the data stream consists of a file consisting of 500,000 bytes, that the MSS is 1,000 bytes, and that the first byte of the data stream is numbered 0.

❖ TCP constructs 500 segments out of the data stream. The first segment gets assigned sequence number 0, the second segment gets 1,000 and so on.



**Figure 3.30** ♦ Dividing file data into TCP segments

# TCP seq. numbers, ACKs

sequence numbers:
- byte stream "number" of first byte in segment's data

acknowledgements:
- seq # of next byte expected from other side
- cumulative ACK

Q: how receiver handles out-of-order segments
- A: TCP spec doesn't say, - up to implementor

outgoing segment from sender

| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |
| | rwnd |
| checksum | urg pointer |

window size
$N$

sender sequence number space

sent ACKed

sent, not-yet ACKed ("in-flight")

usable but not yet sent

not usable

incoming segment to sender

| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |
| | A | rwnd |
| checksum | urg pointer |

# TCP seq. numbers, ACKs

Host A                                    Host B

User
types
'C'
    Seq=42, ACK=79, data = 'C'

                                  host ACKs
                                  receipt of
    Seq=79, ACK=43, data = 'C'  'C', echoes
                                  back 'C'

host ACKs
receipt
of echoed
'C'
    Seq=43, ACK=80

simple telnet scenario

# TCP round trip time, timeout

Q: how to set TCP timeout value?

❖ longer than RTT
  ▪ but RTT varies
❖ *too short:* premature timeout, unnecessary retransmissions
❖ *too long:* slow reaction to segment loss

Q: how to estimate RTT?

❖ **SampleRTT**: measured time from segment transmission until ACK receipt
  ▪ ignore retransmissions
❖ **SampleRTT** will vary, want estimated RTT "smoother"
  ▪ average several *recent* measurements, not just current **SampleRTT**

# TCP round trip time, timeout

**EstimatedRTT = (1- α)\*EstimatedRTT + α\*SampleRTT**

- ❖ exponential weighted moving average
- ❖ influence of past sample decreases exponentially fast
- ❖ typical value: $\alpha = 0.125$



RTT: gaia.cs.umass.edu to fantasia.eurecom.fr

◆ sampleRTT
■ EstimatedRTT

RTT (milliseconds)

time (seconds)

# TCP round trip time, timeout

❖ **timeout interval:** `EstimatedRTT` plus "safety margin"
  ▪ large variation in `EstimatedRTT` `->` larger safety margin

❖ estimate SampleRTT deviation from EstimatedRTT:

```
DevRTT = (1-β)*DevRTT +
            β*|SampleRTT-EstimatedRTT|

         (typically, β = 0.25)
```

**`TimeoutInterval = EstimatedRTT + 4*DevRTT`**

estimated RTT          "safety margin"

# Chapter 3 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 principles of congestion control

3.7 TCP congestion control

# TCP reliable data transfer

❖ **TCP creates rdt service on top of IP's unreliable service**
  - pipelined segments
  - cumulative acks
  - single retransmission timer

❖ **retransmissions triggered by:**
  - timeout events
  - duplicate acks

**let's initially consider simplified TCP sender:**
  - ignore duplicate acks
  - ignore flow control, congestion control

# TCP sender events:

**Three events in sender:**

*1. data rcvd from app:*

- ❖ create segment with seq #
- ❖ seq # is byte-stream number of first data byte in segment
- ❖ start timer if not already running
  - ▪ think of timer as for oldest unacked segment
  - ▪ expiration interval: `TimeOutInterval`

*2. timeout:*

- ❖ retransmit segment that caused timeout
- ❖ restart timer

*3. ack rcvd:*

- ❖ if ack acknowledges previously unacked segments
  - ▪ update what is known to be ACKed
  - ▪ start timer if there are still unacked segments

```
NextSeqNum = InitialSeqNum
SendBase = InitialSeqNum

loop (forever) {
  switch(event)

  event: data received from application above
      create TCP segment with sequence number NextSeqNum
      if (timer currently not running)
          start timer
      pass segment to IP
      NextSeqNum = NextSeqNum + length(data)

  event: timer timeout
      retransmit not-yet-acknowledged segment with
          smallest sequence number
      start timer

  event: ACK received, with ACK field value of y
      if (y > SendBase) {
          SendBase = y
          if (there are currently not-yet-acknowledged segments)
              start timer
      }

} /* end of loop forever */
```

# TCP sender (simplified)

Comment:
• SendBase-1: last cumulatively ack'ed byte
Example:
• SendBase-1 = 71; y= 73, so the rcvr wants 73+ ; y > SendBase, so that new data is acked

# TCP: retransmission scenarios



lost ACK scenario

premature timeout

# TCP: retransmission scenarios



Host A                    Host B

timeout

Seq=92, 8 bytes of data

Seq=100, 20 bytes of data

ACK=100

X

ACK=120

Seq=120,  15 bytes of data

cumulative ACK

# TCP ACK generation [RFC 1122, RFC 2581]

| event at receiver | TCP receiver action |
|---|---|
| arrival of in-order segment with expected seq #. All data up to expected seq # already ACKed | delayed ACK. Wait up to 500ms for next segment. If no next segment, send ACK **(prevents creating extra traffic)** |
| arrival of in-order segment with expected seq #. One other segment has ACK pending | immediately send single cumulative ACK, ACKing both in-order segments **(do not let sender wait too long)** |
| arrival of out-of-order segment higher-than-expect seq. # . Gap detected | immediately send *duplicate ACK,* indicating seq. # of next expected byte **(enables fast retransmit)** |
| arrival of segment that partially or completely fills gap | immediate send ACK, provided that segment starts at lower end of gap |

# Doubling the Timeout Interval

❖ Each time TCP retransmits, it sets the next timeout interval to twice the previous value, rather than deriving it from the last EstimatedRTT and DevRTT

❖ If Time out Interval is .75 sec when the timer first expires. TCP will then retransmit this segment and set the new expiration time to 1.5 sec.

❖ If the timer expires again 1.5 sec later, TCP will again retransmit this segment, now setting the expiration time to 3.0 sec. Thus the intervals grow exponentially after each retransmission

# TCP fast retransmit

❖ time-out period  often relatively long:
- long delay before resending lost packet

❖ detect lost segments via duplicate ACKs.
- sender often sends many segments back-to-back
- if segment is lost, there will likely be many duplicate ACKs.

*TCP fast retransmit*

if sender receives 3 duplicate ACKs for same data

("triple duplicate ACKs"),

resend unacked segment with smallest seq #
- likely that unacked segment lost, so don't wait for timeout

# Is TCP a GBN or SR?

❖ TCP looks a lot like GBN but TCP will buffer correctly received but out-of-order segments

❖ Some version, has selective ack

❖ So it is hybrid GBN and SR

# TCP fast retransmit



Host A                          Host B

Seq=92, 8 bytes of data
Seq=100, 20 bytes of data

X

ACK=100

ACK=100

timeout

ACK=100

ACK=100

Seq=100, 20 bytes of data

fast retransmit after sender
receipt of triple duplicate ACK

# Chapter 3 outline

# TCP flow control

application may
remove data from
TCP socket buffers ....

... slower than TCP
receiver is delivering
(sender is sending)

application
process

TCP socket
receiver buffers

TCP
code

IP
code

application
---------
OS

*flow control*
receiver controls sender, so
sender won't overflow
receiver's buffer by transmitting
too much, too fast

from sender

receiver protocol stack

# TCP flow control



counting by bytes of data (not segments!)

# bytes rcvr willing to accept

❖ receiver "advertises" free buffer space by including **rwnd** value in TCP header of receiver-to-sender segments

- **RcvBuffer** size set via socket options (typical default is 4096 bytes)
- many operating systems autoadjust **RcvBuffer**

❖ sender limits amount of unacked ("in-flight") data to receiver's **rwnd** value

❖ guarantees receive buffer will not overflow

to application process



RcvBuffer    buffered data

rwnd    free buffer space

TCP segment payloads

*receiver-side buffering*

# TCP flow control

- ❖ Because TCPis not permitted to overflow the allocated buffer, we must have

- ❖ LastByteRcvd − LastByteRead less than or equal to RcvBuffer

- ❖ The receive window, denoted rwnd is set to the amount of spare room in the buffer:

- ❖ rwnd = RcvBuffer − [LastByteRcvd − LastByteRead]

*to application process*

`RcvBuffer`

buffered data

`rwnd`

free buffer space

*TCP segment payloads*

*receiver-side buffering*

# TCP flow control

❖ Host A is never informed that some space has opened up in Host B's receive buffer

❖ Host A is blocked and can transmit no more data!

❖ To solve this problem, the TCP specification requires Host A to continue to send segments with one data byte when B's receive window is zero.

❖ These segments will be acknowledged by the receiver.

❖ Eventually the buffer will begin to empty and the acknowledgments will contain a nonzero rwnd value.

# Chapter 3 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 principles of congestion control

3.7 TCP congestion control

# Connection Management

before exchanging data, sender/receiver "handshake":

❖ agree to establish connection (each knowing the other willing to establish connection)

❖ agree on connection parameters

application

connection state: ESTAB
connection variables:
    seq # client-to-server
        server-to-client
    **rcvBuffer** size
    at server,client

network

application

connection state: ESTAB
connection Variables:
    seq # client-to-server
        server-to-client
    **rcvBuffer** size
    at server,client

network

```
Socket clientSocket =
  newSocket("hostname","port
  number");
```

```
Socket connectionSocket =
  welcomeSocket.accept();
```

# Agreeing to establish a connection

2-way handshake:



*Q:* will 2-way handshake always work in network?

❖ variable delays
❖ retransmitted messages (e.g. req_conn(x)) due to message loss
❖ message reordering
❖ can't "see" other side

# Three way handshake:

Step 1: client host sends TCP SYN segment to server.
- SYN bit is set to 1
- specifies initial seq #
- no data

Step 2: server host receives SYN, replies with SYNACK segment
- SYN bit is set to 1
- server allocates buffers
- specifies server initial seq. #
- no data

Step 3: client receives SYNACK, replies with ACK segment, which may contain data. SYN bit is set to 0. Client also allocate buffers

# TCP 3-way handshake

client state

LISTEN

SYNSENT

ESTAB

choose init seq num, x
send TCP SYN msg

SYNbit=1, Seq=x

SYNbit=1, Seq=y
ACKbit=1; ACKnum=x+1

received SYNACK(x)
indicates server is live;
send ACK for SYNACK;
this segment may contain
client-to-server data

ACKbit=1, ACKnum=y+1

SYN bit is set to zero

choose init seq num, y
send TCP SYNACK
msg, acking SYN

received ACK(y)
indicates client is live

server state

LISTEN

SYN RCVD

ESTAB

# TCP: closing a connection

- ❖ client, server each close their side of connection
  - ▪ send TCP segment with FIN bit = 1
- ❖ respond to received FIN with ACK
  - ▪ on receiving FIN, ACK can be combined with own FIN
- ❖ simultaneous FIN exchanges can be handled

# TCP: closing a connection

*client state*                                                                    *server state*

ESTAB                                                                              ESTAB

`clientSocket.close()`

FIN_WAIT_1     can no longer
               send but can         FINbit=1, seq=x
               receive data                                                        CLOSE_WAIT

                                     ACKbit=1; ACKnum=x+1                           can still
FIN_WAIT_2     wait for server                                                      send data
               close

                                                                                   LAST_ACK

TIMED_WAIT                           FINbit=1, seq=y                                can no longer
                                                                                   send data

               timed wait           ACKbit=1; ACKnum=y+1
               for 2*max
               segment lifetime                                                    CLOSED

CLOSED

# RST flag

❖ Suppose a host receives a TCP SYN packet with destination port 80 but host does not accept connection on port 80.

❖ Host will send a reset segment to client (set RST to 1)

❖ By this, host says " I don't have a socket for that segment"

# Sending SYN could have three possible outcomes

- ❖ SYNACK is received
- ❖ RST is received
- ❖ Nothing received which means that SYN segment was blocked by firewall and never reached the target host

# Chapter 3 outline

# Principles of congestion control

*congestion*:

❖ What is congestion?

❖ –Load is higher than capacity

❖ What do IP routers do? –Drop the excess packets

❖ informally: "too many sources sending too much data too fast for *network* to handle"

❖ different from flow control!

❖ manifestations:

  ▪ lost packets (buffer overflow at routers)
  ▪ long delays (queueing in router buffers)

❖ a top-10 problem!

# Causes/costs of congestion: scenario 1

❖ two senders, two receivers

❖ one router, infinite buffers

❖ output link capacity: R

❖ no retransmission

original data: $\lambda_{in}$

throughput: $\lambda_{out}$

Host A

unlimited shared output link buffers

Host B

❖ maximum per-connection throughput: R/2

❖ large delays as arrival rate, $\lambda_{in}$, approaches capacity

# Causes/costs of congestion: scenario 2

❖ one router, *finite* buffers

❖ sender retransmission of timed-out packet
  - application-layer input = application-layer output: $\lambda_{in} = \lambda_{out}$
  - transport-layer input includes *retransmissions* : $\lambda'_{in} \geq \lambda_{in}$

$\lambda_{in}$ : original data

$\lambda'_{in}$: original data, *plus* retransmitted data

$\lambda_{out}$

Host A

Host B

finite shared output link buffers

# Causes/costs of congestion: scenario 2

idealization: perfect knowledge

❖ sender sends only when router buffers available



$\lambda_{in}$ : original data

$\lambda'_{in}$: original data, *plus* retransmitted data

copy

$\lambda_{out}$

A

*free buffer space!*

Host B

finite shared output link buffers

# Causes/costs of congestion: scenario 2

*Idealization: known loss*
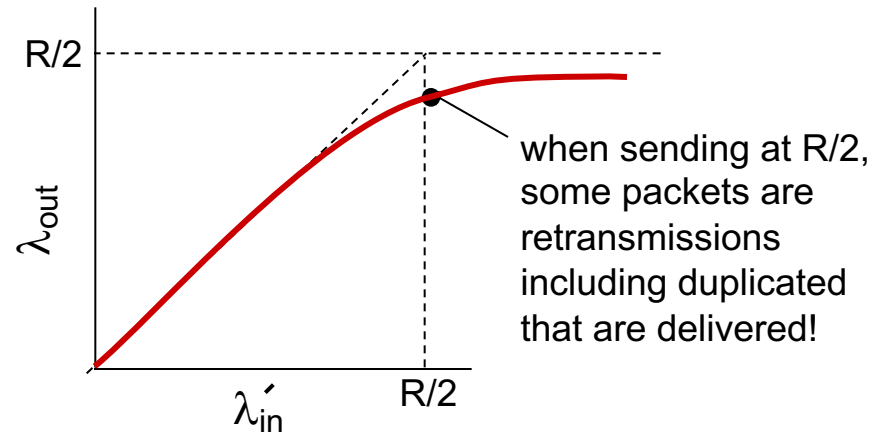    packets can be lost, dropped at router due to full buffers

❖ sender only resends if packet *known* to be lost



$\lambda_{in}$ : original data
$\lambda'_{in}$: original data, *plus* retransmitted data

$\lambda_{out}$

copy

no buffer space!

A

Host B

# Causes/costs of congestion: scenario 2

*Idealization: known loss*
  packets can be lost, dropped at router due to full buffers

❖ sender only resends if packet *known* to be lost

$\lambda_{out}$ vs $\lambda'_{in}$ graph with R/2 on both axes.

when sending at R/2, some packets are retransmissions but asymptotic goodput is still R/2 (why?)

$\lambda_{in}$ : original data

$\lambda'_{in}$: original data, *plus* retransmitted data

$\lambda_{out}$

A

*free buffer space!*

Host B

# Causes/costs of congestion: scenario 2

## *Realistic: duplicates*

❖ packets can be lost, dropped at router due to full buffers

❖ sender times out prematurely, sending *two* copies, both of which are delivered

when sending at R/2, some packets are retransmissions including duplicated that are delivered!

$\lambda_{out}$

$\lambda'_{in}$

R/2

timeout

$\lambda_{in}$

$\lambda'_{in}$

A

$\lambda_{out}$

*free buffer space!*

Host B

# Causes/costs of congestion: scenario 2

*Realistic: duplicates*

❖ packets can be lost, dropped at router due to full buffers

❖ sender times out prematurely, sending *two* copies, both of which are delivered

when sending at R/2, some packets are retransmissions including duplicated that are delivered!

## "costs" of congestion:

❖ more work (retrans) for given "goodput"

❖ unneeded retransmissions: link carries multiple copies of pkt

   ▪ decreasing goodput

# Causes/costs of congestion: scenario 3

* four senders
* multihop paths
* timeout/retransmit

Q: what happens as $\lambda_{in}$ and $\lambda_{in}'$ increase?

A: as red $\lambda_{in}'$ increases, all arriving blue pkts at upper queue are dropped, blue throughput $\rightarrow 0$

Host A

$\lambda_{in}$: original data

$\lambda_{in}'$: original data, *plus* retransmitted data

$\lambda_{out}$

Host B

finite shared output link buffers

Host D

Host C

Transport Layer 3-98

# Causes/costs of congestion: scenario 3



another "cost" of congestion:

❖ when packet dropped, any "upstream transmission capacity used for that packet was wasted!

# Approaches towards congestion control

two broad approaches towards congestion control:

## end-end congestion control:

❖ no explicit feedback from network layer to transport layer

❖ congestion inferred from end-system observed loss, delay

❖ approach taken by TCP

## network-assisted congestion control:

❖ routers provide feedback to end systems

  ▪ Feedback could be a single bit indicating congestion Examples: TCP/IP ECN, ATM ABR)

  ▪ explicit rate for sender to send at

# Two ways of feedback from the network to the sender

1. Direct feedback: sent from a network router to the sender, essentially saying, "I'm congested!".

2. Router marks/updates a field in a packet flowing from sender to receiver to indicate congestion.

❖ Upon receipt of a marked packet, the receiver notifies the sender of the congestion indication.

❖ Note: this form of notification takes at least a full round-trip time.

# Two ways of feedback from the network to the sender

# ATM (**Asynchronous Transfer Mode**)

❖ ATM takes virtual circuit (VC) approach toward packet switching. This means that each switch on the source-to-destination path will maintain state about VC

❖ Cell used instead of Packet

❖ Switch used instead of Router

❖ Resource-management (RM) cells are used to send congestion-related information among the hosts and switches

# Case study: ATM ABR congestion control

## ABR: Available Bit Rate:

- ❖ "elastic service"
- ❖ if sender's path "underloaded":
  - sender should use available bandwidth
- ❖ if sender's path congested:
  - sender limited to minimum guaranteed rate

## RM (resource management) cells:

- ❖ sent by sender, between data cells
- ❖ RM cells can be used to provide both direct network feedback and network feedback via the receiver
- ❖ bits in RM cell set by switches ("*network-assisted*")
- ❖ RM cells returned to sender by receiver, with bits intact

# ABR provides three mechanisms for signaling congestion

1. EFCI bit. Each data cell contains an explicit forward congestion indication (EFCI) bit. A congested network switch can set EFCI bit in a data cell to 1. The destination host checks the EFCI bit. When an RM cell arrives at the destination, then the destination sets the congestion indication bit (the CI bit) of the RM cell to 1 and sends the RM cell back to the sender

2. CI and NI bits: RM cells have a congestion indication (CI) bit (severe congestion) and a no increase (NI) bit ( mild congestion ) in RM cells

3. ER setting. Each RM cell also contains a 2-byte explicit rate (ER) field. A congested switch lower the value contained in the ER field in a passing RM cell

# Case study: ATM ABR congestion control



RM cell　　data cell

- ❖ **two-byte ER (explicit rate) field in RM cell**
  - congested switch may lower ER value in cell
  - senders' send rate thus max supportable rate on path
- ❖ **EFCI bit in data cells: set to 1 in congested switch**
  - if data cell preceding RM cell has EFCI set, receiver sets CI bit in returned RM cell
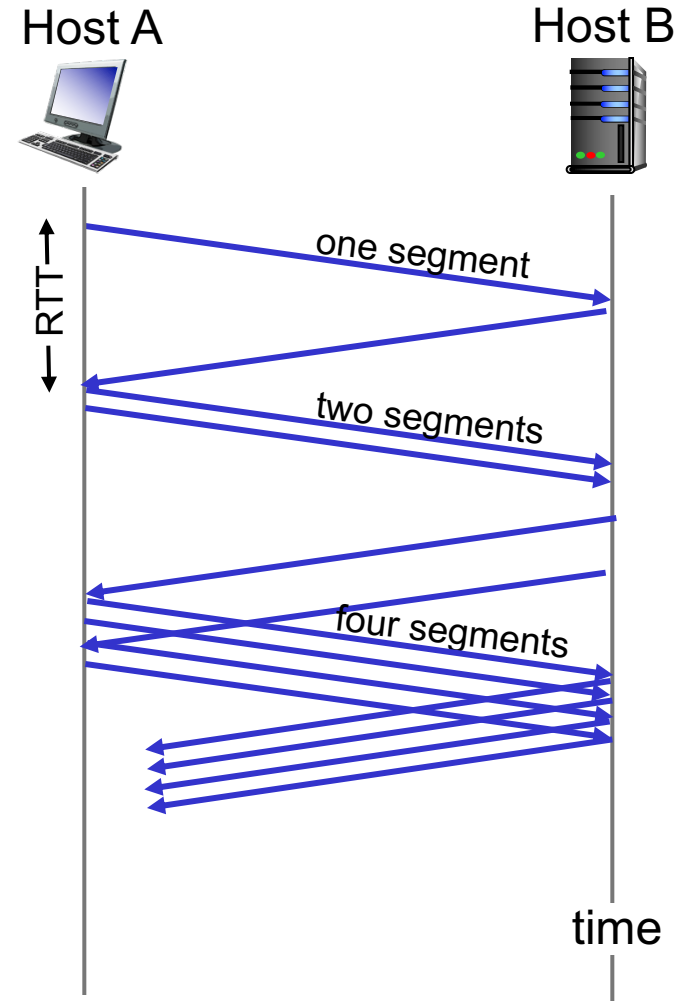
# Chapter 3 outline

# TCP congestion control

❖ sender keeps track of an, congestion window (cwnd) the amount of unacknowledged data at a sender may not exceed the minimum of cwnd and rwnd, that is:

$$\text{LastByteSent} - \text{LastByteAcked} \leq \min\{\text{cwnd, rwnd}\}$$

❖ Assuming TCP receive buffer is very large, then the sender's send rate is roughly cwnd/RTT bytes/sec

❖ TCP detects congestion if "lost event". Lost is detected using timeout and Duplicate ACK.

# TCP Slow Start

- ❖ when connection begins, increase rate exponentially until first loss event:
  - initially `cwnd` = 1 MSS
  - double `cwnd` every RTT
  - done by incrementing `cwnd` for every ACK received

- ❖ *summary:* initial rate is slow but ramps up exponentially fast

Host A          Host B

RTT

one segment

two segments

four segments

time

# Slow start after loss event

❖ if there is a loss event (i.e., congestion) indicated by a timeout, the TCP sender sets the value of cwnd to 1 and sets ssthresh to cwnd/2

❖ when the value of cwnd equals ssthresh, slow start ends and TCP transitions into congestion avoidance mode

❖ slow start can end if three duplicate ACKs are detected, in which case TCP performs a fast retransmit and enters fast recovery state.

# Congestion Avoidance

- ❖ Rather than doubling cwnd every RTT, TCP increases cwnd by just a single MSS every RTT

- ❖ TCP's congestion-avoidance algorithm behaves the same as in the case of slow start when a timeout occurs.

- ❖ But for the case of loss detection using a triple duplicate ACK, TCP halves the value of cwnd and enter fast-recovery
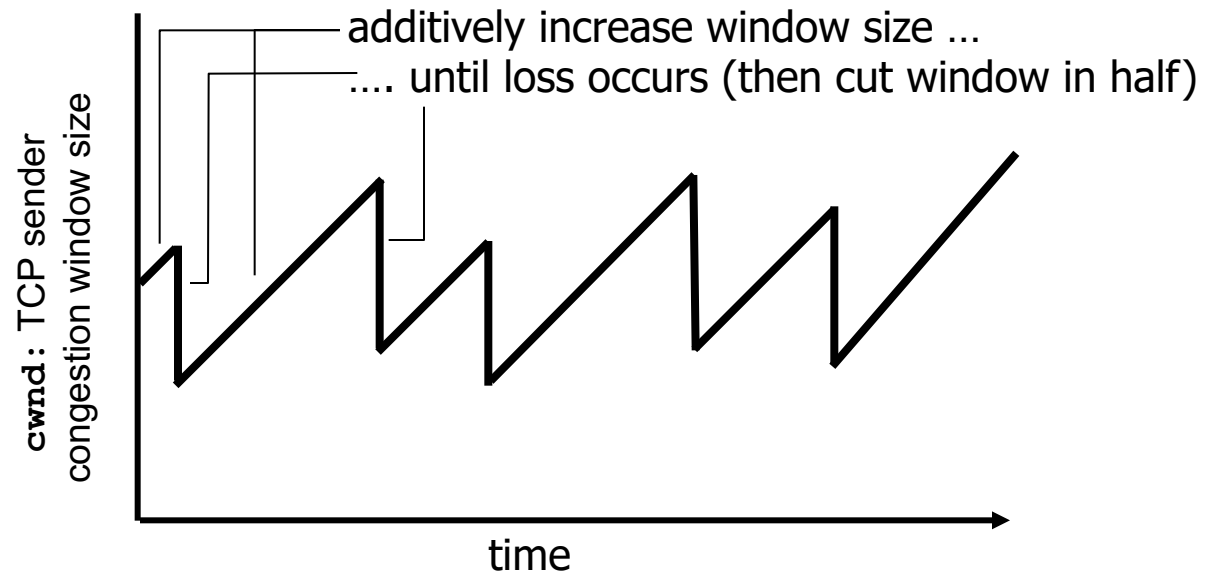
# Fast Recovery

❖ Set SSTHRESH=CWND/2

❖ Set CWND=SSTHRESH+3 MSS

❖ Every subsequent duplicate ack: CWND=CWND+1MSS

❖ When a new ACK arrives (not a duplicate ack), TCP enters the congestion-avoidance state

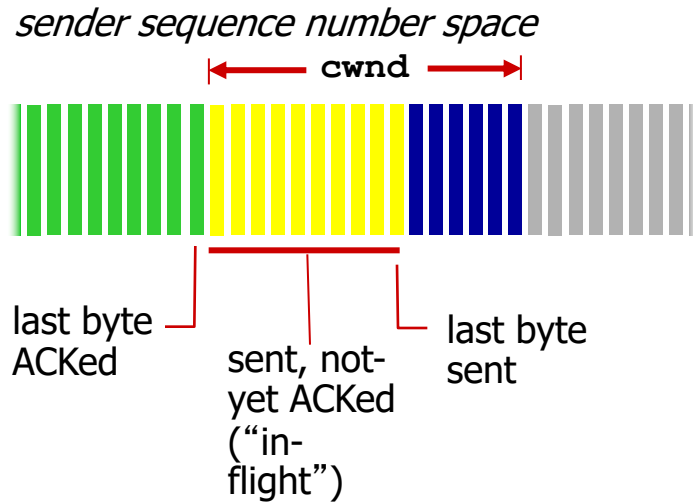❖ If a timeout event occurs, fast recovery transitions to the slow-start state

# TCP congestion control: additive increase multiplicative decrease

❖ *approach:* sender increases transmission rate (window size), probing for usable bandwidth, until loss occurs

▪ *additive increase:* increase `cwnd` by 1 MSS every RTT until loss detected

▪ *multiplicative decrease:* cut `cwnd` in half after loss

AIMD saw tooth behavior: probing for bandwidth

additively increase window size ...

.... until loss occurs (then cut window in half)

**cwnd**: TCP sender congestion window size

time

# TCP Congestion Control: details

*sender sequence number space*

$\longleftarrow$ `cwnd` $\longrightarrow$

last byte ACKed

sent, not-yet ACKed ("in-flight")

last byte sent

❖ sender limits transmission:

$$\text{LastByteSent} - \text{LastByteAcked} \leq \text{cwnd}$$

❖ **cwnd** is dynamic, function of perceived network congestion

*TCP sending rate:*

❖ *roughly:* send cwnd bytes, wait RTT for ACKS, then send more bytes

$$\text{rate} \approx \frac{\text{cwnd}}{\text{RTT}} \text{ bytes/sec}$$
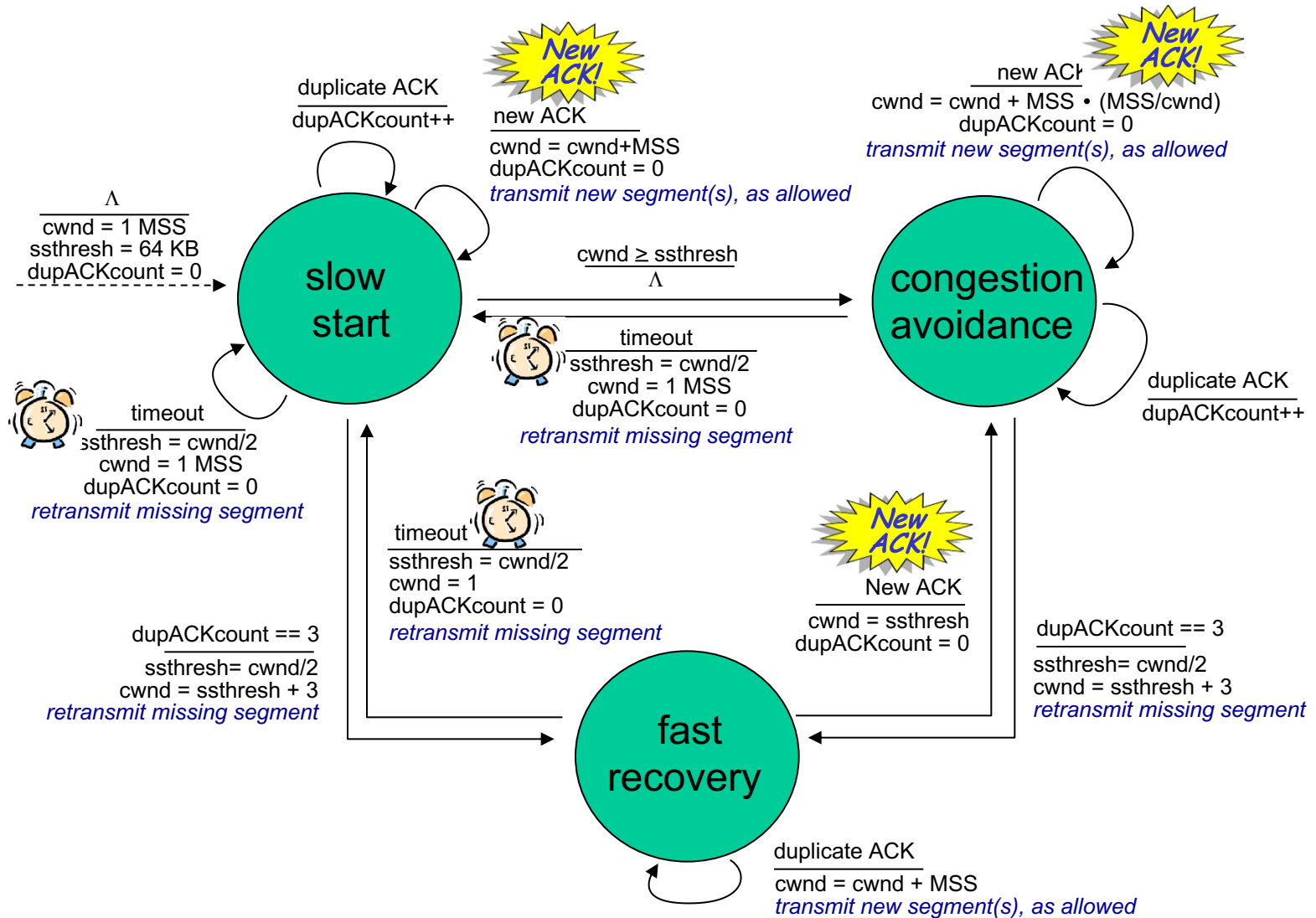
# TCP Taho and TCP Reno

TCP Tahoe, unconditionally cut its congestion window to 1 MSS and entered the slow-start phase after either a timeout or triple-duplicate-ACK.

TCP Reno, incorporated fast recovery

# Summary: TCP Congestion Control



$\Lambda$
cwnd = 1 MSS
ssthresh = 64 KB
dupACKcount = 0

**slow start**

duplicate ACK
―――――――――
dupACKcount++

New ACK!
new ACK
―――――――――
cwnd = cwnd+MSS
dupACKcount = 0
*transmit new segment(s), as allowed*

New ACK!
new ACK
―――――――――
cwnd = cwnd + MSS · (MSS/cwnd)
dupACKcount = 0
*transmit new segment(s), as allowed*

cwnd ≥ ssthresh
―――――――――
$\Lambda$

**congestion avoidance**

timeout
―――――――――
ssthresh = cwnd/2
cwnd = 1 MSS
dupACKcount = 0
*retransmit missing segment*

duplicate ACK
―――――――――
dupACKcount++

timeout
―――――――――
ssthresh = cwnd/2
cwnd = 1 MSS
dupACKcount = 0
*retransmit missing segment*

timeout
―――――――――
ssthresh = cwnd/2
cwnd = 1
dupACKcount = 0
*retransmit missing segment*

New ACK!
New ACK
―――――――――
cwnd = ssthresh
dupACKcount = 0

dupACKcount == 3
―――――――――
ssthresh= cwnd/2
cwnd = ssthresh + 3
*retransmit missing segment*

dupACKcount == 3
―――――――――
ssthresh= cwnd/2
cwnd = ssthresh + 3
*retransmit missing segment*

**fast recovery**

duplicate ACK
―――――――――
cwnd = cwnd + MSS
*transmit new segment(s), as allowed*

# Summary: TCP Congestion Control

❖ When **CongWin** is below **Threshold**, sender in slow-start phase, window grows exponentially.

❖ When **CongWin** is above **Threshold**, sender is in congestion-avoidance phase, window grows linearly.

❖ When a triple duplicate ACK occurs, **Threshold** set to **CongWin/2** and **CongWin** set to **Threshold**.

❖ When timeout occurs, **Threshold** set to **CongWin/2** and **CongWin** is set to 1 MSS.
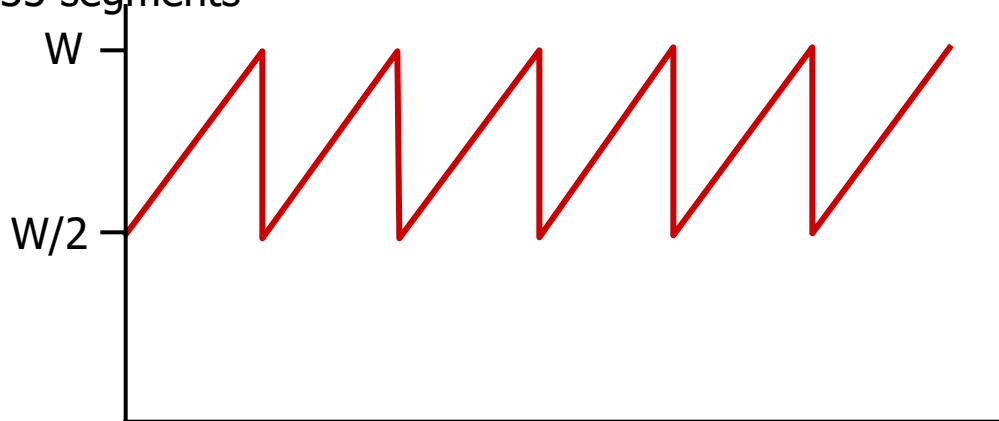
# TCP sender congestion control

| State | Event | TCP Sender Action | Commentary |
|---|---|---|---|
| Slow Start (SS) | ACK receipt for previously unacked data | CongWin = CongWin + MSS, If (CongWin > Threshold) set state to "Congestion Avoidance" | Resulting in a doubling of CongWin every RTT |
| Congestion Avoidance (CA) | ACK receipt for previously unacked data | CongWin = CongWin+MSS * (MSS/CongWin) | Additive increase, resulting in increase of CongWin by 1 MSS every RTT |
| SS or CA | Loss event detected by triple duplicate ACK | Threshold = CongWin/2, CongWin = Threshold, Set state to "Congestion Avoidance" | Fast recovery, implementing multiplicative decrease. CongWin will not drop below 1 MSS. |
| SS or CA | Timeout | Threshold = CongWin/2, CongWin = 1 MSS, Set state to "Slow Start" | Enter slow start |
| SS or CA | Duplicate ACK | Increment duplicate ACK count for segment being acked | CongWin and Threshold not changed |

# TCP throughput

❖ avg. TCP thruput as function of window size, RTT?
   ▪ ignore slow start, assume always data to send
❖ W: window size (measured in bytes)
   ▪ avg. window size (# in-flight bytes) is ¾ W (Average of W/2 and W)
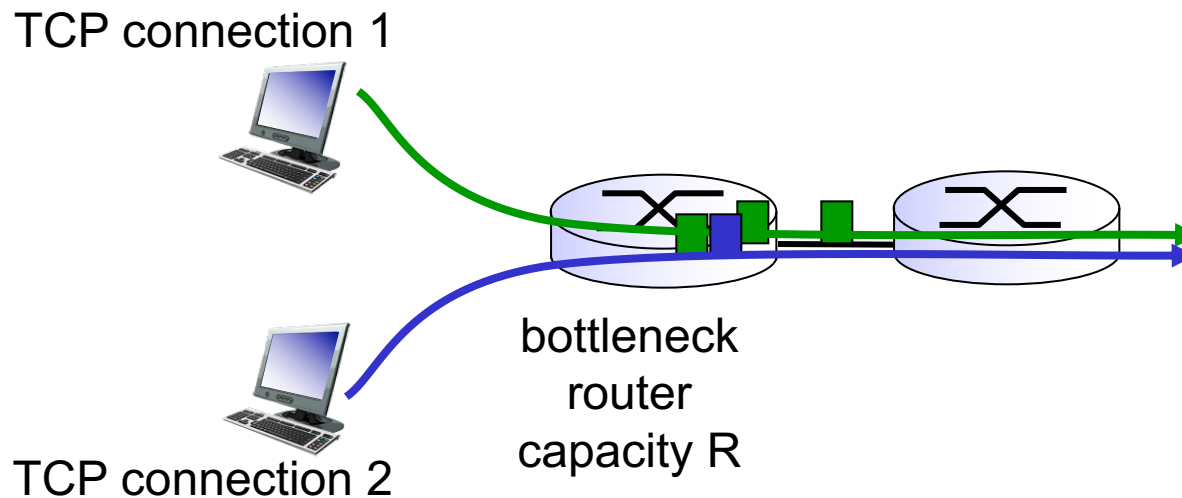   ▪ avg. thruput is 3/4W per RTT

$$\text{avg TCP thruput} = \frac{3}{4} \frac{W}{RTT} \text{ bytes/sec}$$

example: 1500 byte segments, 100ms RTT, want 10 Gbps throughput
requires W = 83,333 in-flight segments (Throughput is 1.25 GB, 1.25 GB
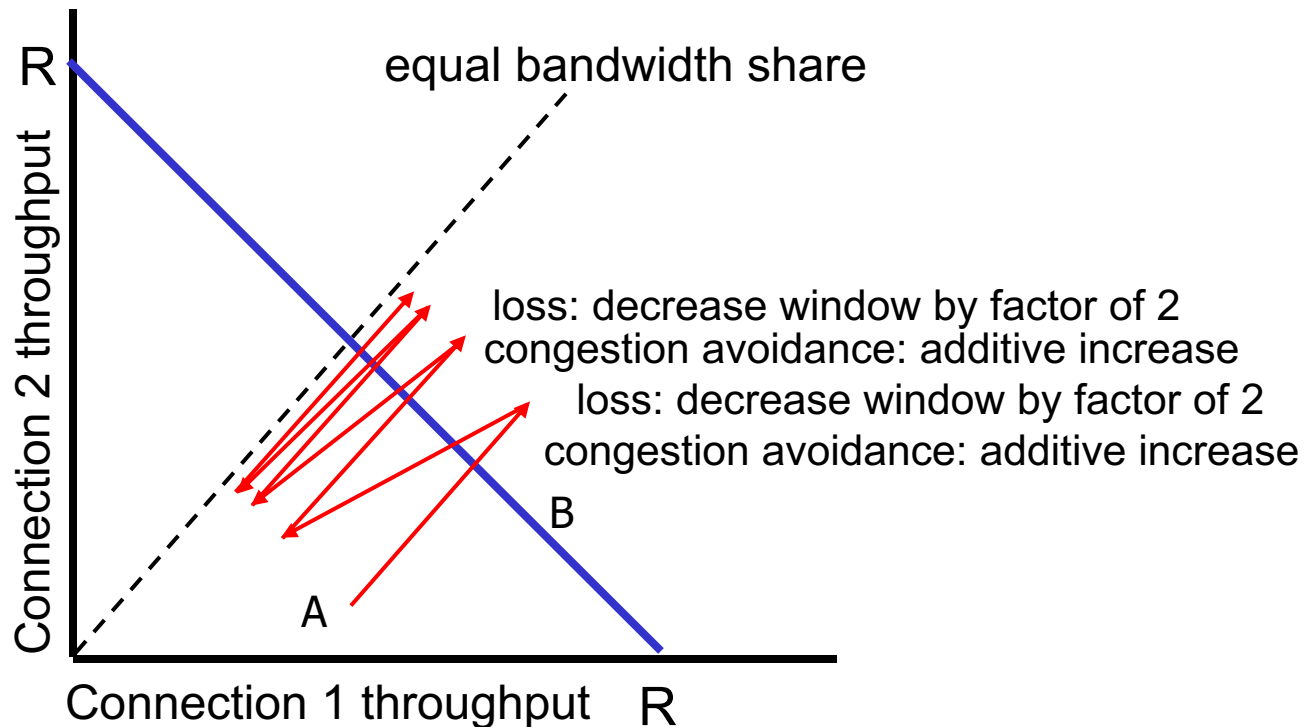/1500 =83,333 segments

# TCP Fairness

*fairness goal:* if K TCP sessions share same bottleneck link of bandwidth R, each should have average rate of R/K

TCP connection 1



bottleneck
router
capacity R

TCP connection 2

# Why is TCP fair?

two competing sessions:

❖ additive increase gives slope of 1, as throughout increases
❖ multiplicative decrease decreases throughput proportionally



equal bandwidth share

Connection 2 throughput

loss: decrease window by factor of 2
congestion avoidance: additive increase

loss: decrease window by factor of 2
congestion avoidance: additive increase

B

A

R

Connection 1 throughput    R

- Suppose that the TCP window sizes are at point *A*. Because the amount of link bandwidth jointly consumed by the two connections is less than *R*, both connections will increase their window by 1 MSS.

- Throughput of the two connections proceeds along a 45-degree line

- Both connections decrease their

- windows by a factor of two (Point C).

# Fairness (more)

## Fairness and UDP

* multimedia apps often do not use TCP
  * do not want rate throttled by congestion control
* instead use UDP:
  * send audio/video at constant rate, tolerate packet loss

## Fairness, parallel TCP connections

* application can open multiple parallel connections between two hosts
* web browsers do this
* e.g., link of rate R with 9 existing connections:
  * new app asks for 1 TCP, gets rate R/10
  * new app asks for 11 TCPs, gets more than R/2

# Chapter 3: summary

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

❖ instantiation, implementation in the Internet
- UDP
- TCP

next:

❖ leaving the network "edge" (application, transport layers)

❖ into the network "core"