# Chapter 3. Transport Layer

# Chapter 3: Transport Layer

## Our goals:

❖ understand principles behind transport layer services:

- multiplexing/demultipl exing
- reliable data transfer
- flow control
- congestion control

❖ learn about transport layer protocols in the Internet:

- UDP: connectionless transport
- TCP: connection-oriented transport
- TCP congestion control

# 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
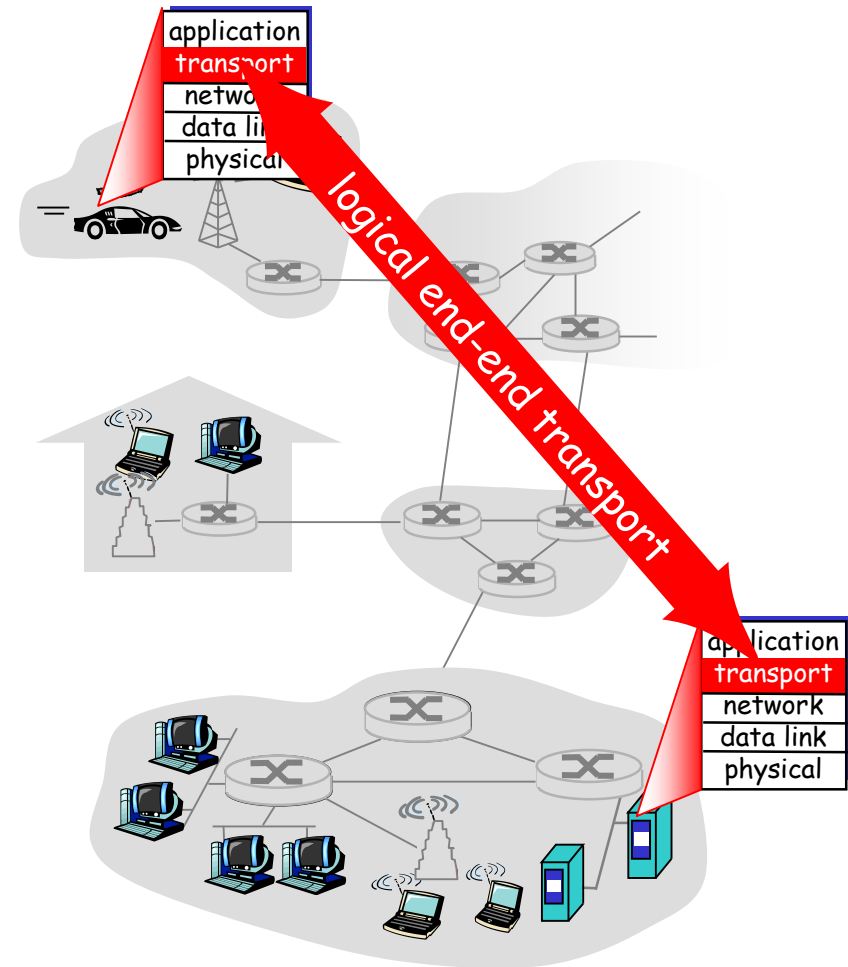
3.5 Connection-oriented transport: TCP

- 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

❖ 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
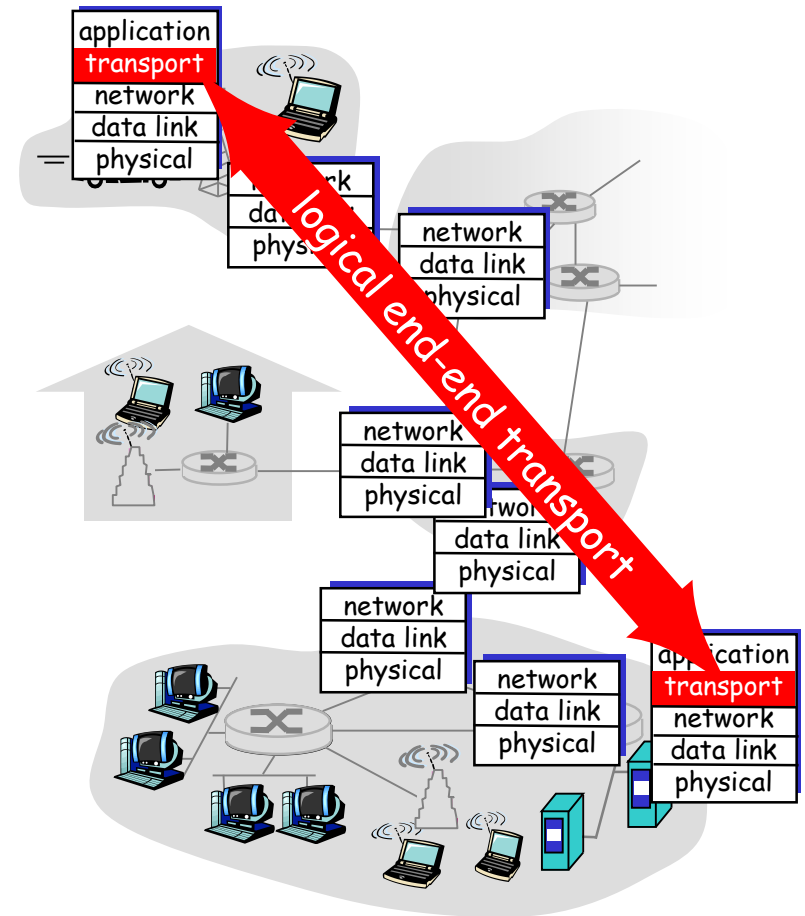- relies on, enhances, network layer services

**Household analogy:**

*12 kids sending letters to 12 kids*

❖ processes = kids

❖ app messages = letters in envelopes

❖ hosts = houses

❖ 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

application
transport
network
data link
physical

network
data link
physical

network
data link
physical

network
data link
physical

network
data link
physical

network
data link
physical

application
transport
network
data link
physical

logical end-end transport

# 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
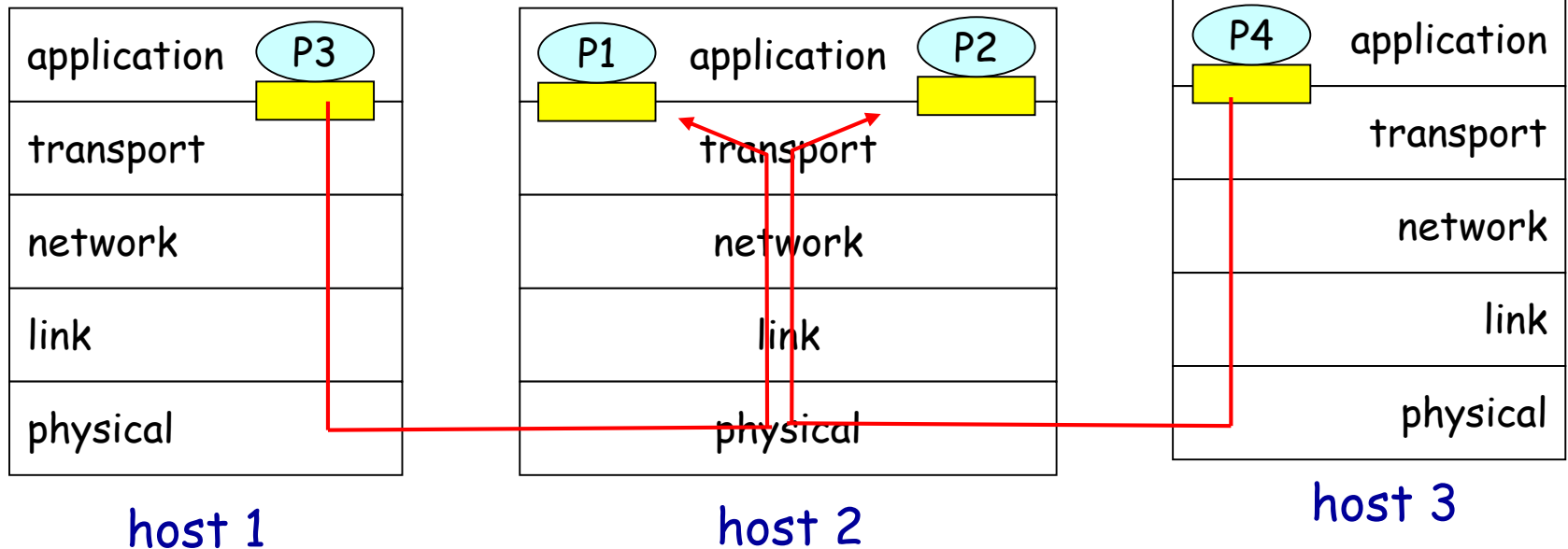
# Multiplexing/demultiplexing

delivering received segments to correct socket

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

☐ = socket   ⬭ = process
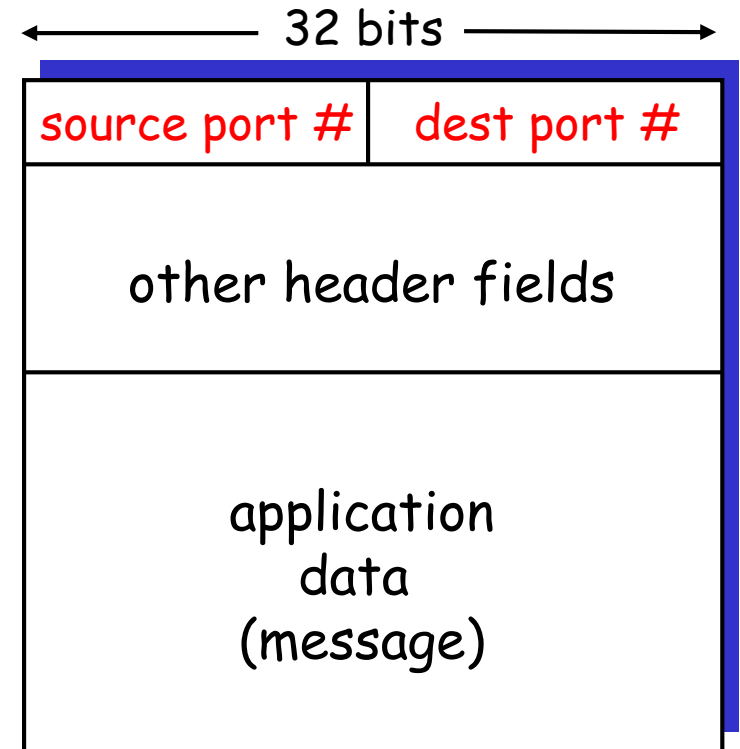


host 1        host 2        host 3

# How demultiplexing works

❖ **host receives IP datagrams**

- each datagram has source IP address, destination IP address
- each datagram carries 1 transport-layer segment
- each segment has source, destination port number

❖ **host uses IP addresses & port numbers to direct segment to appropriate socket**

32 bits

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

TCP/UDP segment format

# Connectionless demultiplexing

❖ *recall:* create sockets with host-local port numbers:

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

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

❖ *recall:* when creating datagram to send into UDP socket, must specify
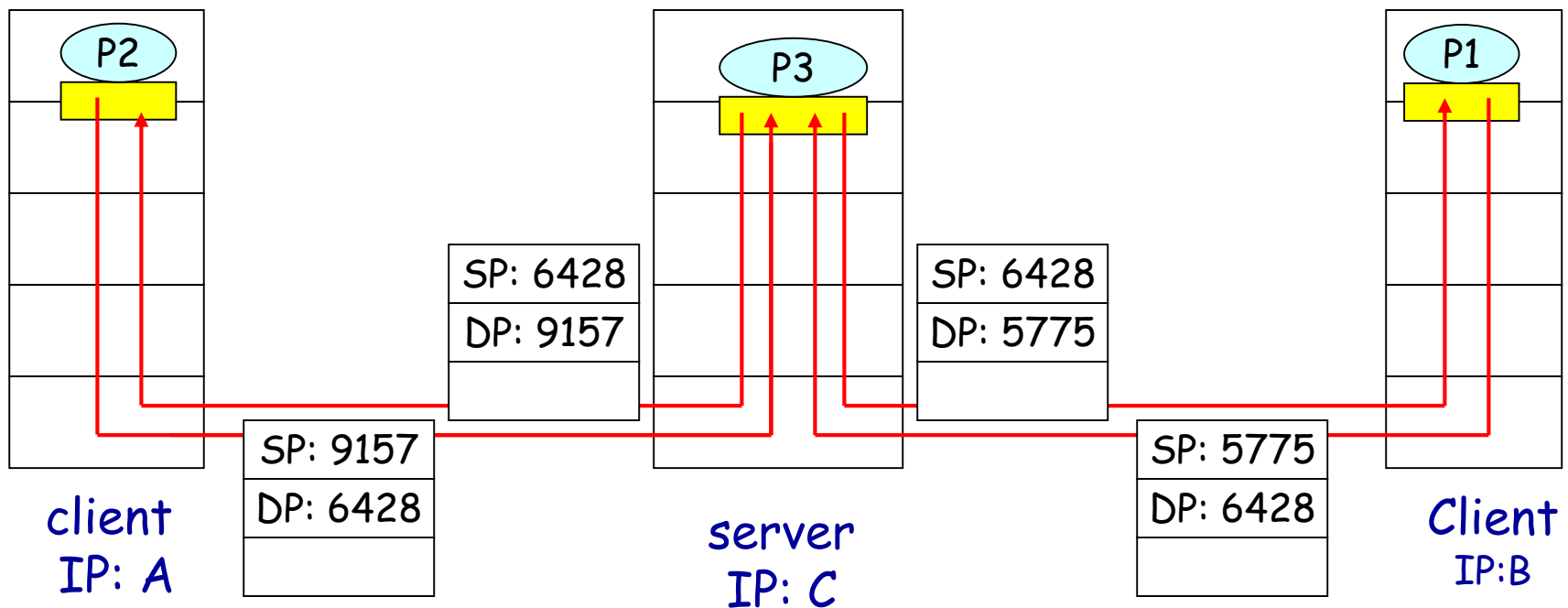
(dest IP address, dest port number)

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

❖ IP datagrams with different source IP addresses and/or source port numbers directed to same socket

# Connectionless demux (cont)

`DatagramSocket serverSocket = new DatagramSocket(6428);`



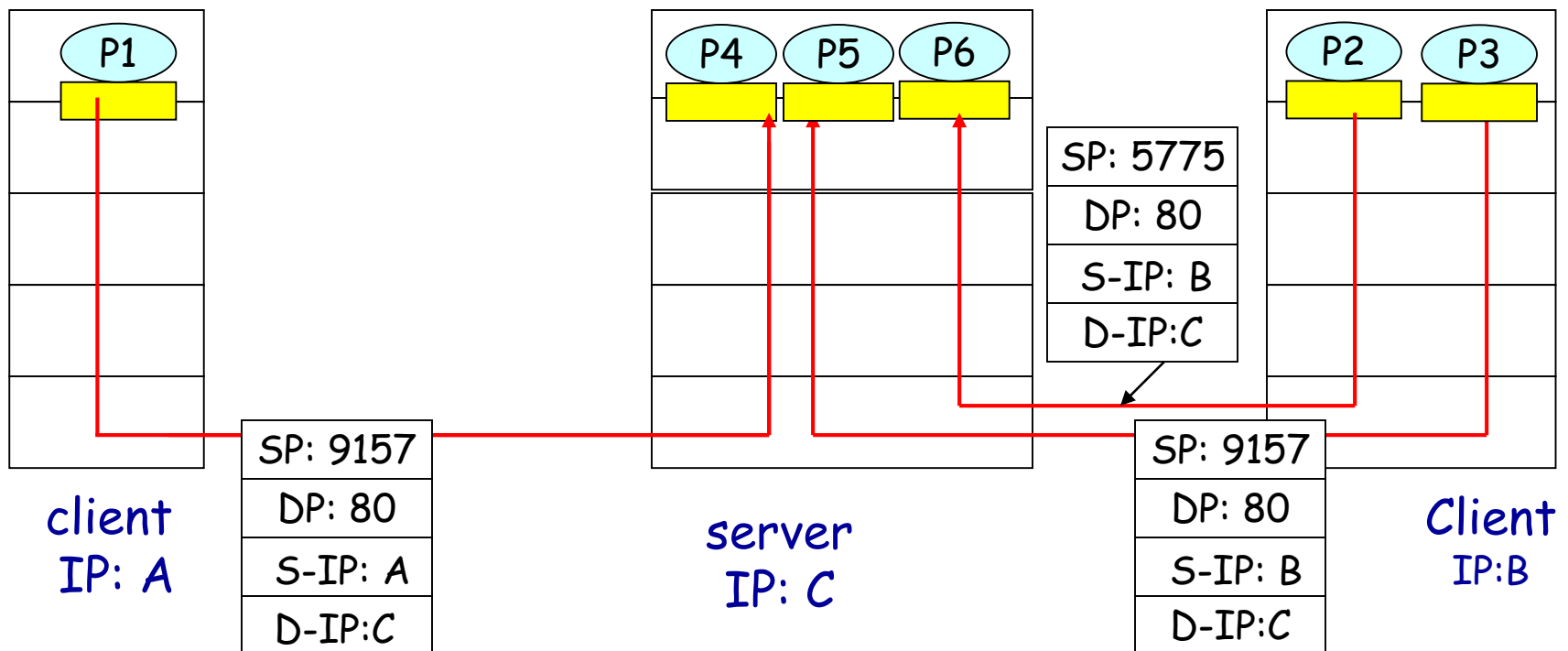client
IP: A

server
IP: C

Client
IP:B

SP provides "return address"

# Connection-oriented demux

❖ TCP socket identified by 4-tuple:
  - source IP address
  - source port number
  - dest IP address
  - dest port number

❖ recv host uses all four values to direct segment to appropriate socket

❖ server host may support many simultaneous TCP sockets:
  - each socket identified by its own 4-tuple

❖ web servers have different sockets for each connecting client
  - non-persistent HTTP will have different socket for each request

# Connection-oriented demux (cont)



client
IP: A

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

server
IP: C

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

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

Client
IP:B

P1   P4   P5   P6   P2   P3

# Connection-oriented demux: Threaded Web Server

P1

P4

P2    P3

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

client
IP: A

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

server
IP: C

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

client
IP:B

# 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
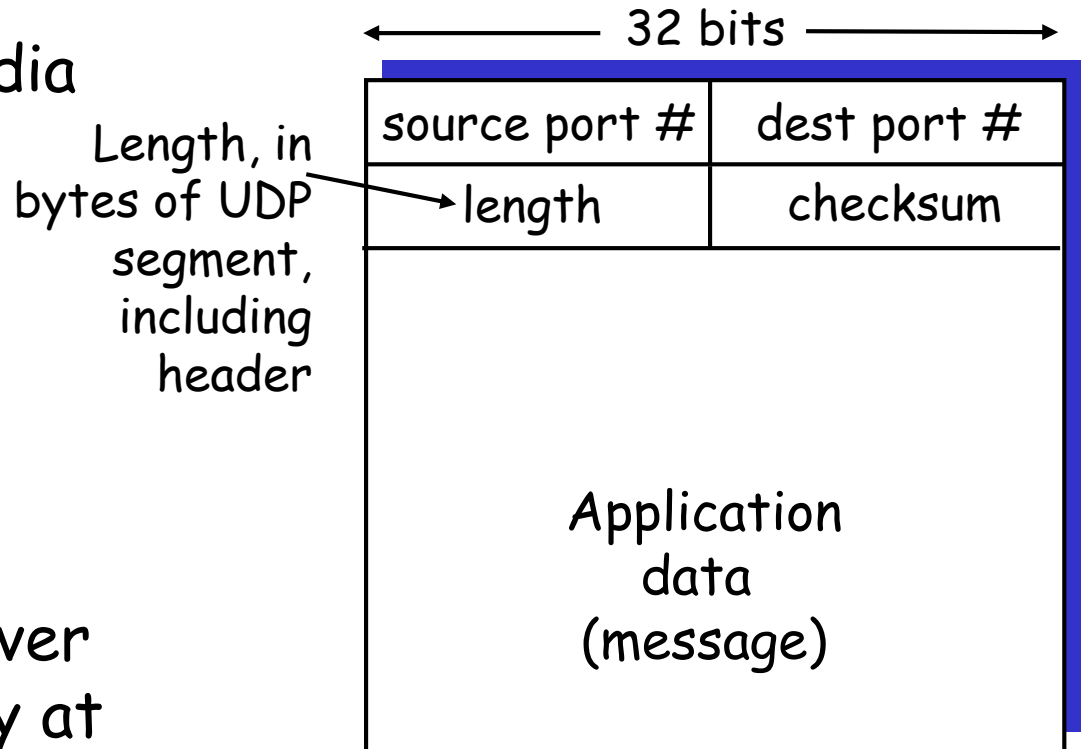
# UDP: User Datagram Protocol [RFC 768]

- ❖ "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

## Why is there a UDP?

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

# UDP: more

❖ **often used for streaming multimedia apps**
  ▪ loss tolerant
  ▪ rate sensitive

❖ **other UDP uses**
  ▪ DNS
  ▪ SNMP

❖ **reliable transfer over UDP: add reliability at application layer**
  ▪ application-specific error recovery!

Length, in bytes of UDP segment, including header

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

UDP segment format

# UDP checksum

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

## Sender:

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

## Receiver:

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

# Internet Checksum Example

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

❖ Example: add two 16-bit integers

```
              1 1 1 0 0 1 1 0 0 1 1 0 0 1 1 0
              1 1 0 1 0 1 0 1 0 1 0 1 0 1 0 1
            _____
wraparound  (1) 1 0 1 1 1 0 1 1 1 0 1 1 1 0 1 1
            _____
     sum      1 0 1 1 1 0 1 1 1 0 1 1 1 1 0 0
checksum      0 1 0 0 0 1 0 0 0 1 0 0 0 0 1 1
```

# Chapter 3 outline

# Principles of Reliable data transfer

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



(a) provided service

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

# Principles of Reliable data transfer

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

application layer

transport layer

sending process

receiver process

data

data

reliable channel

unreliable channel

(a) provided service

(b) service implementation

- ❖ characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)
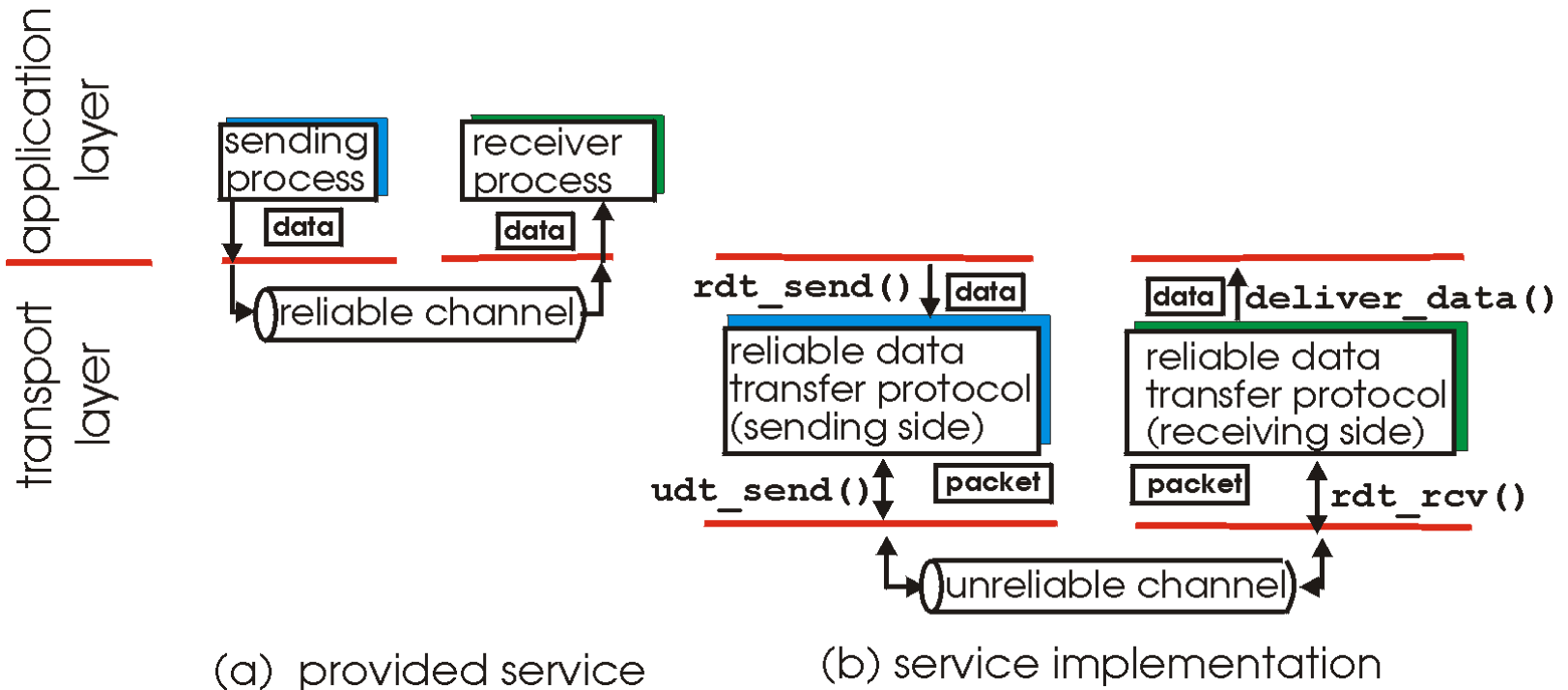
# Principles of Reliable data transfer

❖ **important in app., transport, link layers**

❖ **top-10 list of important networking topics!**



(a) provided service          (b) service implementation

❖ **characteristics of unreliable channel will determine complexity of reliable data transfer protocol (rdt)**
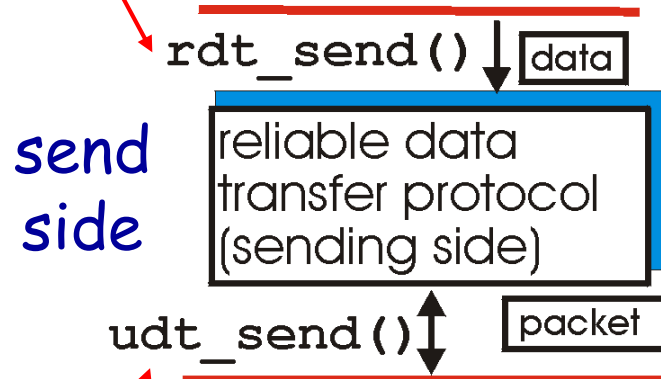
# Reliable data transfer: getting started

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

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

rdt_send() | data |

| data | deliver_data()

send side

reliable data transfer protocol (sending side)

reliable data transfer protocol (receiving side)

receive side

udt_send() | packet |

| packet | rdt_rcv()

( unreliable channel )

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

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

# Reliable data transfer: getting started

**We'll:**

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

❖ consider only unidirectional data transfer
  ▪ but control info will flow on both directions!

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

event causing state transition
actions taken on state transition

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

state 1

event
actions

state 2

# 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 read data from underlying channel

Wait for call from above

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

Wait for call from below

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

sender                    receiver

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

❖ underlying channel may flip bits in packet
  ▪ checksum to detect bit errors
❖ *the* question: how to recover from errors:

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

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

❖ **underlying channel may flip bits in packet**
  ▪ checksum to detect bit errors

❖ *the* **question: how to recover from errors:**
  ▪ *acknowledgements (ACKs):* receiver explicitly tells sender that pkt received OK
  ▪ *negative acknowledgements (NAKs):* receiver explicitly tells sender that pkt had errors
  ▪ sender retransmits pkt on receipt of NAK

❖ **new mechanisms in `rdt2.0` (beyond `rdt1.0`):**
  ▪ error detection
  ▪ receiver feedback: control `msgs` (ACK,NAK) rcvr->sender

# rdt2.0: FSM specification

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

**receiver**

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

Wait for call from above → Wait for ACK or NAK

rdt_rcv(rcvpkt) &&
corrupt(rcvpkt)
——————
udt_send(NAK)

rdt_rcv(rcvpkt) && isACK(rcvpkt)
——————
Λ

**sender**

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)
_____
Λ

rdt_rcv(rcvpkt) &&
corrupt(rcvpkt)
_____
udt_send(NAK)

**Wait for call from below**

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

# rdt2.0: error scenario

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

Wait for call from above

Wait for ACK or NAK

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

rdt_rcv(rcvpkt) && corrupt(rcvpkt)
——————————
udt_send(NAK)

rdt_rcv(rcvpkt) && isACK(rcvpkt)
——————————
Λ

Wait for call from below

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

# rdt2.0 has a fatal flaw!

## What happens if ACK/NAK corrupted?

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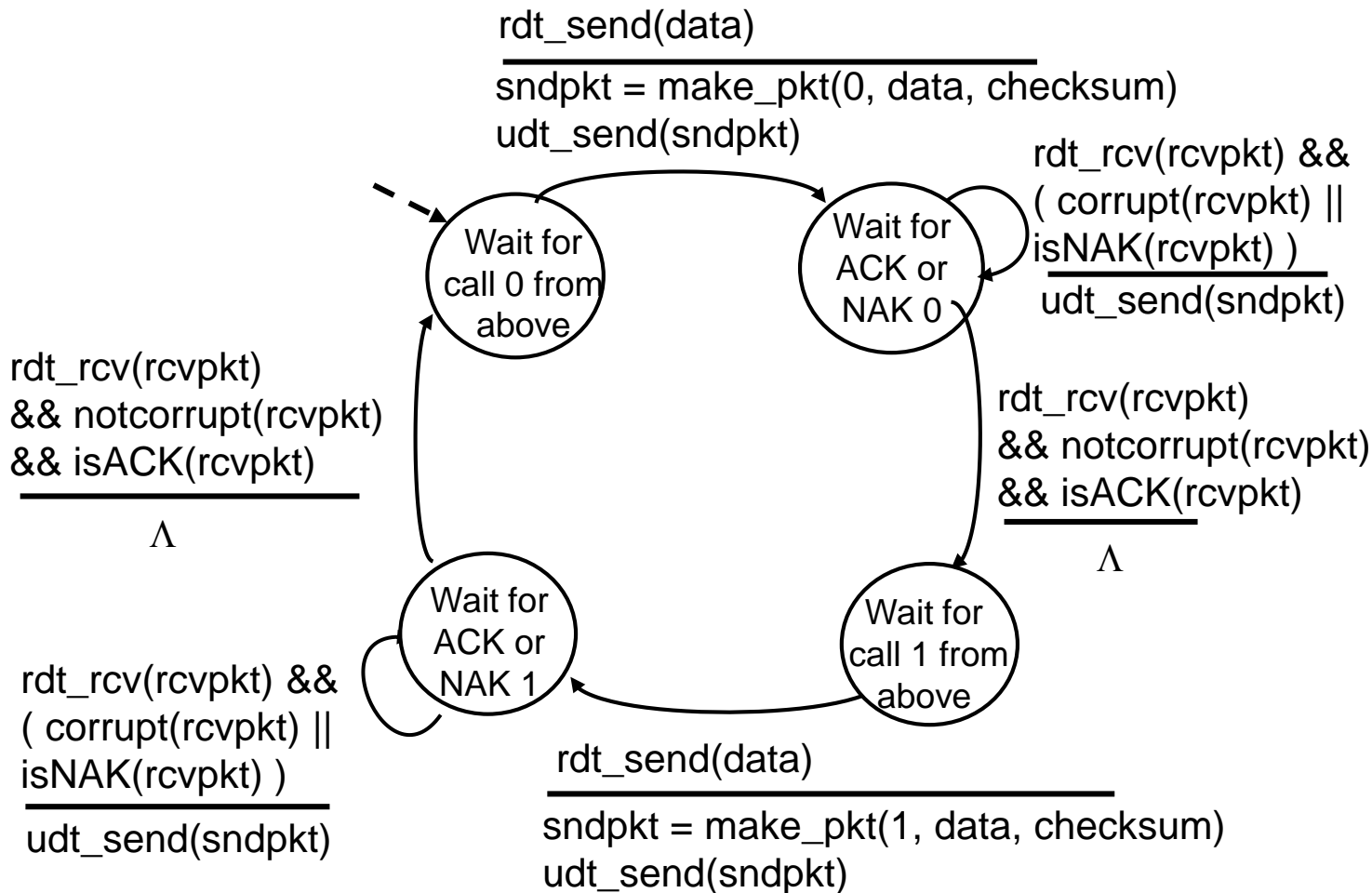- can't just retransmit: possible duplicate

## Handling duplicates:

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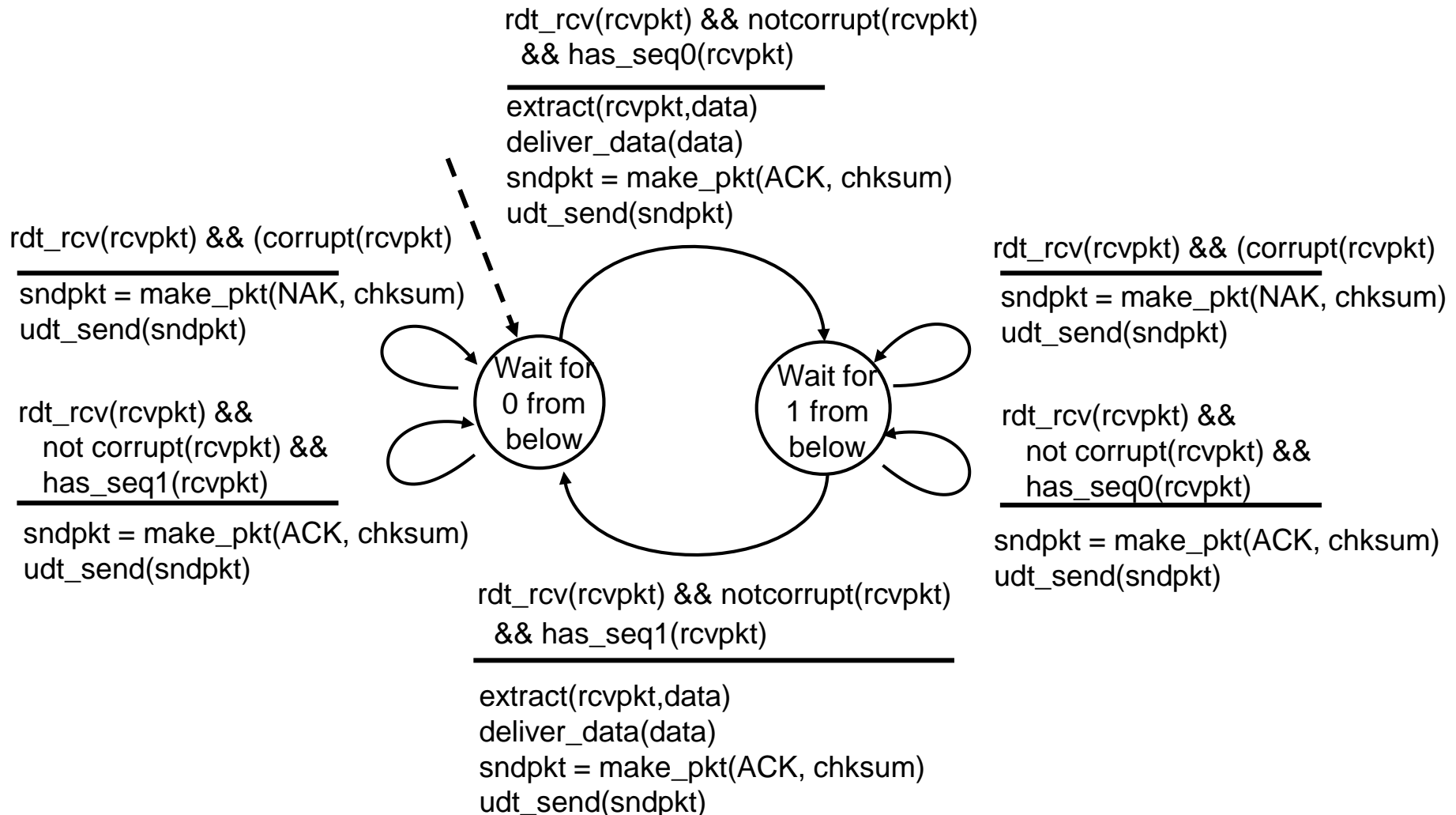- receiver discards (doesn't deliver up) duplicate pkt

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

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

rdt_send(data)
_____

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

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

Wait for
call 0 from
above

Wait for
ACK or
NAK 0

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& isACK(rcvpkt)
_____
$\Lambda$

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& isACK(rcvpkt)
_____
$\Lambda$

Wait for
ACK or
NAK 1

Wait for
call 1 from
above

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

rdt_send(data)
_____

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

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

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
  && has_seq0(rcvpkt)
_____
extract(rcvpkt,data)
deliver_data(data)
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) && (corrupt(rcvpkt)
_____
sndpkt = make_pkt(NAK, chksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) &&
  not corrupt(rcvpkt) &&
  has_seq1(rcvpkt)
_____
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) && (corrupt(rcvpkt)
_____
sndpkt = make_pkt(NAK, chksum)
udt_send(sndpkt)

rdt_rcv(rcvpkt) &&
  not corrupt(rcvpkt) &&
  has_seq0(rcvpkt)
_____
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

( Wait for 0 from below )   ( Wait for 1 from below )

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
  && has_seq1(rcvpkt)
_____
extract(rcvpkt,data)
deliver_data(data)
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

# rdt2.1: discussion

## Sender:

❖ seq # added to pkt

❖ two seq. #'s (0,1) will suffice.  Why?

❖ must check if received ACK/NAK corrupted

❖ twice as many states
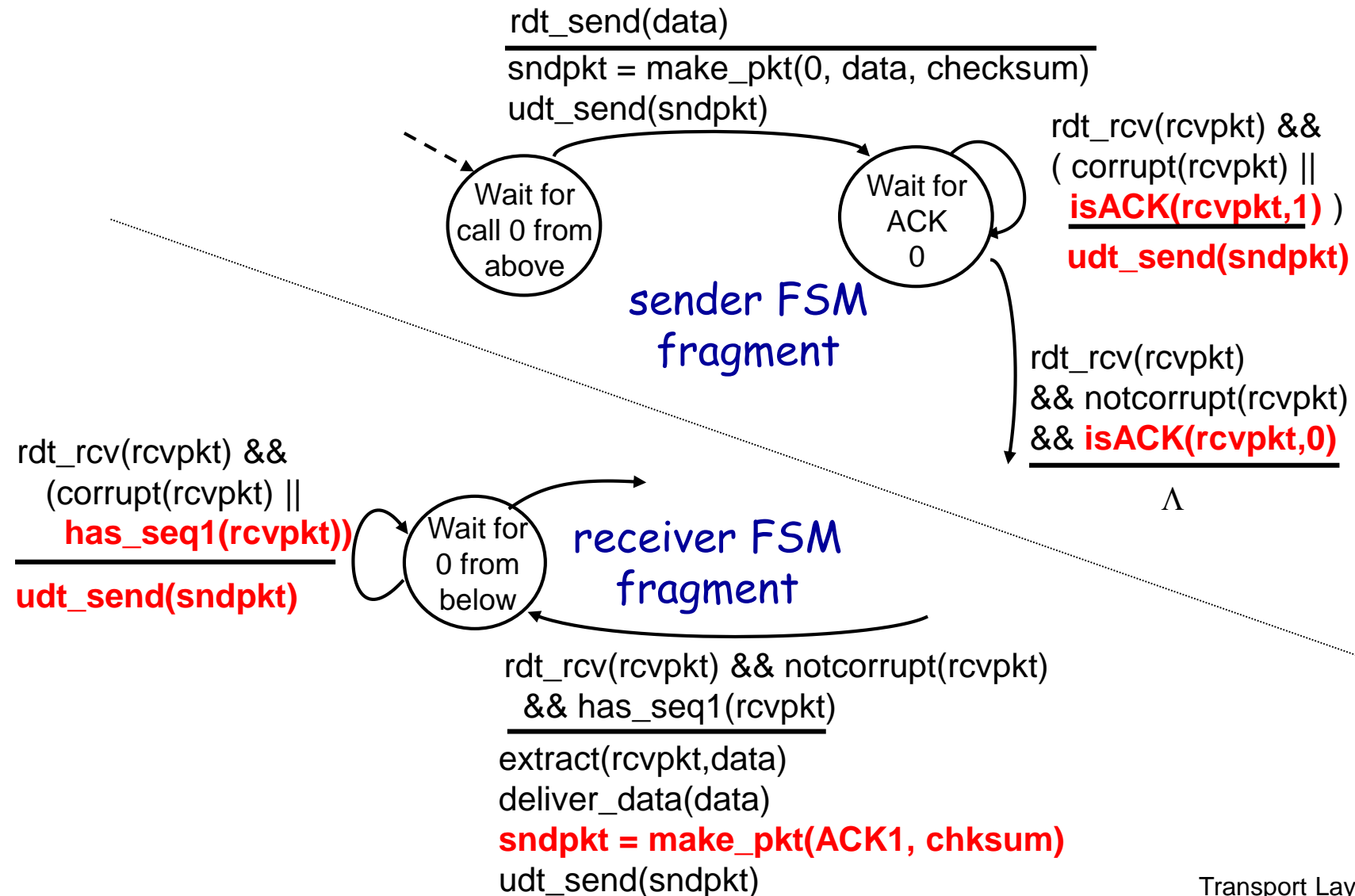  ▪ state must "remember" whether "current" pkt has 0 or 1 seq. #

## Receiver:

❖ must check if received packet is duplicate
  ▪ state indicates whether 0 or 1 is expected pkt seq #

❖ note: receiver can *not* know if its last ACK/NAK received OK at sender

# rdt2.2: a NAK-free protocol

❖ same functionality as rdt2.1, using ACKs only

❖ instead of NAK, receiver sends ACK for last pkt received OK

- receiver must *explicitly* include seq # of pkt being ACKed

❖ 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 or ACKs)

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

Approach: sender waits "reasonable" amount of time for ACK

❖ retransmits if no ACK received in this time
❖ if pkt (or ACK) just delayed (not lost):

- retransmission will be duplicate, but use of seq. #'s already handles this
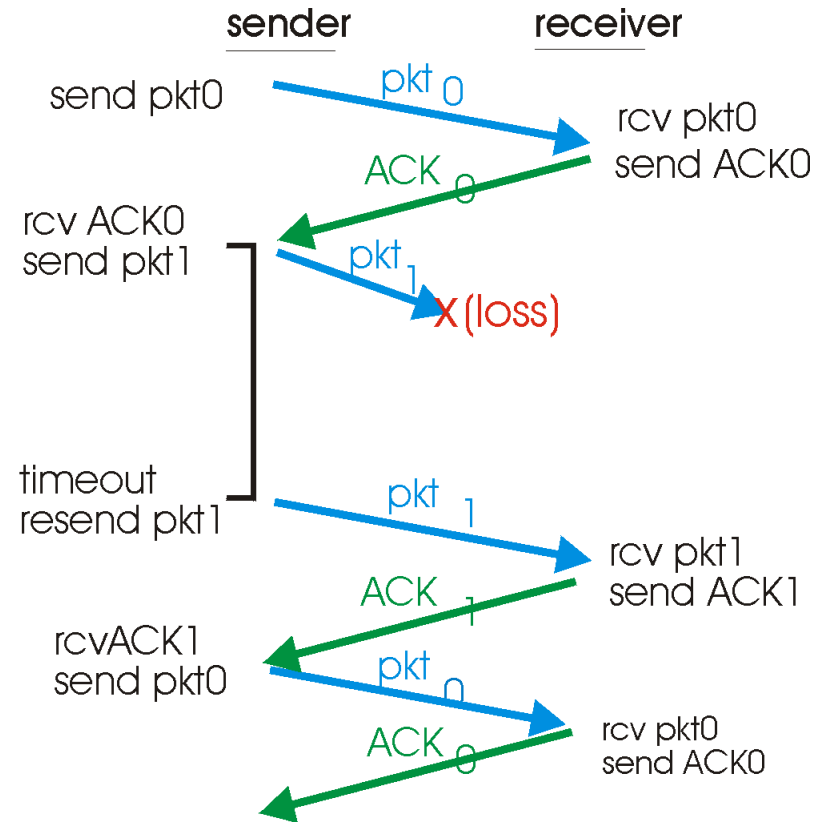- receiver must specify seq # of pkt being ACKed

❖ requires countdown timer

# rdt3.0 sender

rdt_send(data)

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

rdt_rcv(rcvpkt)
$\Lambda$

Wait for call 0from 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

Wait for ACK1

Wait for call 1 from above

timeout
udt_send(sndpkt)
start_timer

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    pkt$_0$    rcv pkt0
send ACK0

ACK$_0$

rcv ACK0
send pkt1    pkt$_1$    rcv pkt1
send ACK1

ACK$_1$

rcvACK1
send pkt0    pkt$_0$    rcv pkt0
send ACK0

ACK$_0$

(a) operation with no loss

sender     receiver

send pkt0    pkt$_0$    rcv pkt0
send ACK0

ACK$_0$

rcv ACK0
send pkt1    pkt$_1$    X (loss)

timeout
resend pkt1    pkt$_1$    rcv pkt1
send ACK1

ACK$_1$

rcvACK1
send pkt0    pkt$_0$    rcv pkt0
send ACK0

ACK$_0$

(b) lost packet

# rdt3.0 in action



(c) lost ACK

(d) premature timeout

# Performance of rdt3.0

❖ rdt3.0 works, but performance stinks
❖ ex: 1 Gbps link, 15 ms prop. delay, 8000 bit packet:

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

- $U_{sender}$: <span style="color:red">utilization</span> – 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, $t = L / R$

RTT

first packet bit arrives

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 2nd packet arrives, send ACK

last bit of 3rd packet arrives, send ACK

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

Increase utilization by a factor of 3!

$$U_{sender} = \frac{3 * L / R}{RTT + L / R} = \frac{.024}{30.008} = 0.0008$$

# Pipelined Protocols

## Go-back-N: big picture:

- ❖ sender can have up to N unacked packets in pipeline
- ❖ rcvr only sends *cumulative* acks
  - ▪ doesn't ack packet if there's a gap
- ❖ sender has timer for oldest unacked packet
  - ▪ if timer expires, retransmit all unack'ed packets

## Selective Repeat: big pic

- ❖ 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 unack'ed 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 each in-flight pkt

❖ *timeout(n):* retransmit pkt n and all higher seq # pkts in window
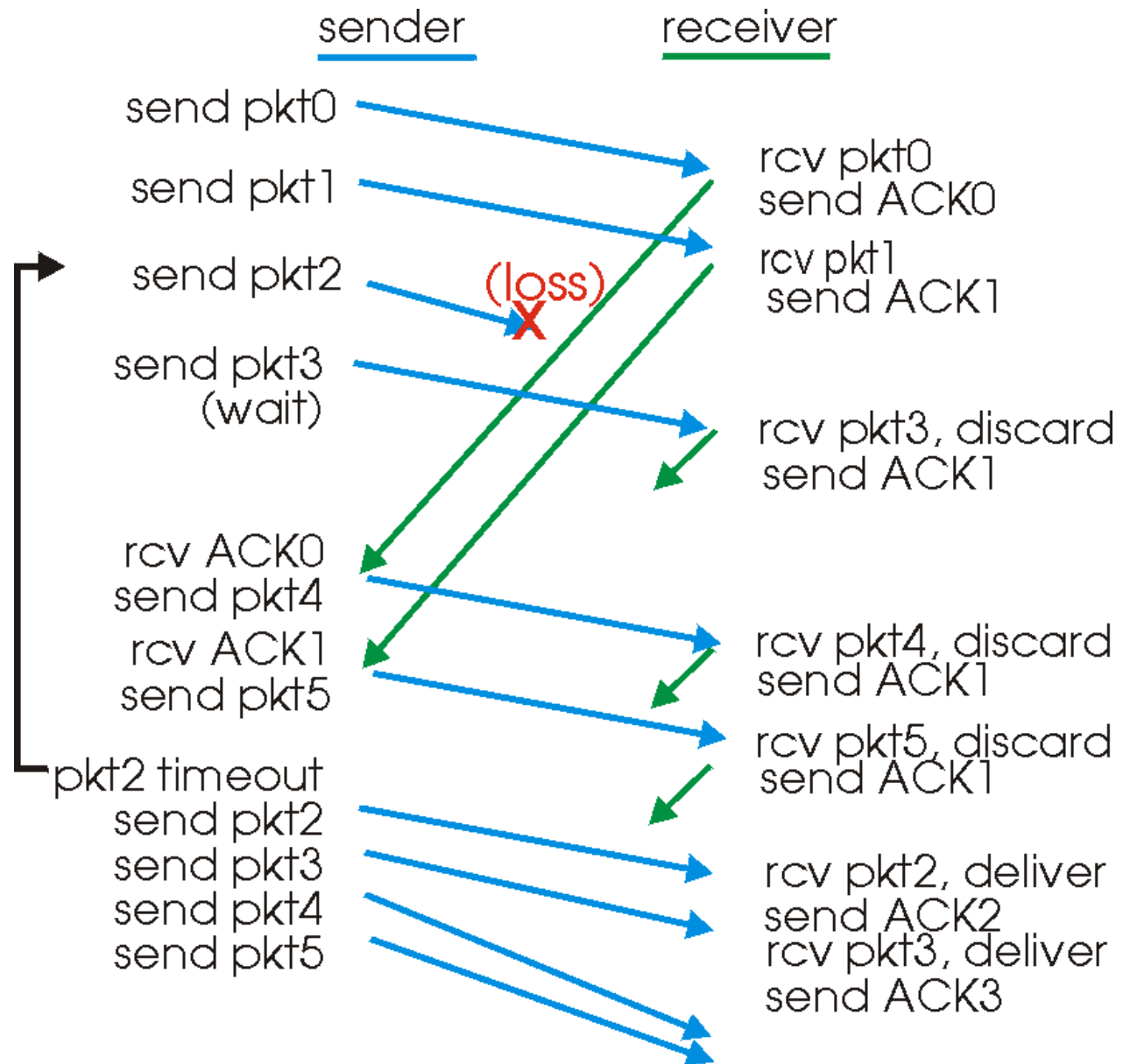
# GBN: sender extended FSM

rdt_send(data)
_____

if (nextseqnum < base+N) {
  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

rdt_rcv(rcvpkt)
  && corrupt(rcvpkt)
_____

**Wait**

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

rdt_rcv(rcvpkt) &&
  notcorrupt(rcvpkt)
_____

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

# GBN: receiver extended FSM

default
_____
udt_send(sndpkt)

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

Λ
_____

( Wait )

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

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

ACK-only: always send ACK for correctly-received pkt
   with highest *in-order* seq #
- may generate duplicate ACKs
- need only remember `expectedseqnum`

❖ out-of-order pkt:
- discard (don't buffer) -> no receiver buffering!
- Re-ACK pkt with highest in-order seq #

# GBN in action

sender    receiver

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

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

# Selective Repeat

❖ **receiver** *individually* acknowledges all correctly received pkts

  ▪ buffers pkts, as needed, for eventual in-order delivery to upper layer

❖ **sender** only resends pkts for which ACK not received

  ▪ sender timer for each unACKed pkt

❖ **sender window**

  ▪ N consecutive seq #'s

  ▪ again limits seq #s of sent, unACK'ed pkts

# Selective repeat: sender, receiver windows

send_base    nextseqnum

■ already ack'ed    ■ usable, not yet sent

□ sent, not yet ack'ed    □ not usable

window size N

(a) sender view of sequence numbers

■ out of order (buffered) but already ack'ed    ■ acceptable (within window)

□ Expected, not yet received    □ not usable

window size N

rcv_base

(b) receiver view of sequence numbers

# Selective repeat

## sender

**data from above :**

- if next available seq # in window, send pkt

**timeout(n):**

- resend pkt n, restart timer

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

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

## receiver

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

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

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

- ACK(n)

**otherwise:**

- ignore

# Selective repeat in action



pkt0 sent
0 1 2 3 | 4 5 6 7 8 9

pkt1 sent
0 1 2 3 | 4 5 6 7 8 9

pkt2 sent
0 1 2 3 | 4 5 6 7 8 9 → X (loss)

pkt3 sent, window full
0 1 2 3 | 4 5 6 7 8 9

ACK0 rcvd, pkt4 sent
0 | 1 2 3 4 | 5 6 7 8 9

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

pkt2 TIMEOUT, pkt2 resent
0 1 | 2 3 4 5 | 6 7 8 9

ACK3 rcvd, nothing sent
0 1 | 2 3 4 5 | 6 7 8 9

pkt0 rcvd, delivered, ACK0 sent
0 | 1 2 3 4 | 5 6 7 8 9

pkt1 rcvd, delivered, ACK1 sent
0 1 | 2 3 4 5 | 6 7 8 9

pkt3 rcvd, buffered, ACK3 sent
0 1 | 2 3 4 5 | 6 7 8 9

pkt4 rcvd, buffered, ACK4 sent
0 1 | 2 3 4 5 | 6 7 8 9

pkt5 rcvd, buffered, ACK5 sent
0 1 | 2 3 4 5 | 6 7 8 9

pkt2 rcvd, pkt2,pkt3,pkt4,pkt5 delivered, ACK2 sent
0 1 2 3 4 5 | 6 7 8 9 |

# Selective repeat: dilemma

Example:

* seq #'s: 0, 1, 2, 3
* window size=3

* receiver sees no difference in two scenarios!
* incorrectly passes duplicate data as new in (a)

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



(a)

(b)

# 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

❖ ***send & receive buffers***

❖ **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



socket door

application writes data

TCP send buffer

socket door

application reads data

TCP receive buffer

segment

# TCP segment structure



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)

32 bits

| source port # | dest port # |
|---|---|
| sequence number | |
| acknowledgement number | |
| head len, not used, U A P R S F | Receive window |
| checksum | Urg data pnter |
| Options (variable length) | |
| application data (variable length) | |

counting by bytes of data (not segments!)

# bytes rcvr willing to accept

# TCP seq. #'s and ACKs

Seq. #'s:
- byte stream "number" of first byte in segment's data

ACKs:
- 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

Host A                          Host B

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

host ACKs
receipt of
'C', echoes
back 'C'

Seq=79, ACK=43, data = 'C'

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

time

simple telnet scenario

# Practice:

❖ Suppose A sends two TCP segments back-to-back to B. The first segment has sequence number 90; the second has sequence number 110

  ▪ a) How much data is the first segment?

  ▪ b) Suppose that the first segment arrives at B. In the acknowledgement that B sends to A, what will be the acknowledgment number?

# TCP Round Trip Time and 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 and Timeout

$$\text{EstimatedRTT} = (1- \alpha)*\text{EstimatedRTT} + \alpha*\text{SampleRTT}$$

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

# Example RTT estimation:

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

# TCP Round Trip Time and Timeout

## Setting the timeout

❖ **`EstimatedRTT`** plus "safety margin"

  ▪ large variation in **`EstimatedRTT`** –> larger safety margin

❖ first estimate of how much SampleRTT deviates from EstimatedRTT:

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

(typically, β = 0.25)
```

Then set timeout interval:

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

# 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

❖ TCP uses single retransmission timer

❖ retransmissions are triggered by:
   ▪ timeout events
   ▪ duplicate acks

❖ initially consider simplified TCP sender:
   ▪ ignore duplicate acks
   ▪ ignore flow control, congestion control

# TCP sender events:

## 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`

## timeout:

- ❖ retransmit segment that caused timeout
- ❖ restart timer

## Ack rcvd:

- ❖ If acknowledges previously unacked segments
  - update what is known to be acked
  - start timer if there are outstanding 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 acked byte
Example:
• SendBase-1 = 71; y= 73, so the rcvr wants 73+ ;
y > SendBase, so that new data is acked

# TCP: retransmission scenarios

Host A          Host B

Seq=92, 8 bytes data

timeout

ACK=100

X
loss

Seq=92, 8 bytes data

ACK=100

SendBase
= 100

time

lost ACK scenario

Host A          Host B

Seq=92, 8 bytes data

Seq=100, 20 bytes data

Seq=92 timeout

ACK=100

ACK=120

Seq=92, 8 bytes data

SendBase
= 100
SendBase
= 120

Seq=92 timeout

ACK=120

SendBase
= 120

time

premature timeout

# TCP retransmission scenarios (more)

Host A                    Host B

Seq=92, 8 bytes data

ACK=100

timeout

Seq=100, 20 bytes data

X
loss

SendBase
= 120

ACK=120

time

Cumulative ACK scenario

# 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 |
| Arrival of in-order segment with expected seq #. One other segment has ACK pending | Immediately send single cumulative ACK, ACKing both in-order segments |
| Arrival of out-of-order segment higher-than-expect seq. # . Gap detected | Immediately send *duplicate ACK*, indicating seq. # of next expected byte |
| Arrival of segment that partially or completely fills gap | Immediate send ACK, provided that segment starts at lower end of gap |

# 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.

❖ if sender receives 3 ACKs for the same data, it supposes that segment after ACKed data was lost:

  ▪ fast retransmit: resend segment before timer expires

Figure 3.37 Resending a segment after triple duplicate ACK

# Fast retransmit algorithm:

event: ACK received, with ACK field value of y
        if (y > SendBase) {
            SendBase = y
            if (there are currently not-yet-acknowledged segments)
               start timer
        }
      else {
          increment count of dup ACKs received for y
          if (count of dup ACKs received for y = 3) {
             resend segment with sequence number y
          }

a duplicate ACK for
already ACKed segment

fast retransmit

# Consideration

❖ We saw that TCP waits until it has received three duplicate ACKs before performing a fast retransmit. Why do you think the TCP designers chose not to perform a fast retransmit after the first duplicate ACK for a segment is received?

# 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 Flow Control

❖ receive side of TCP connection has a receive buffer:

RcvWindow

data from IP → spare room | TCP data in buffer → application process

RcvBuffer

❖ speed-matching service: matching the send rate to the receiving app's drain rate

❖ app process may be slow at reading from buffer

# TCP Flow control: how it works



(suppose TCP receiver discards out-of-order segments)

❖ spare room in buffer

```
= RcvWindow
```

```
= RcvBuffer-[LastByteRcvd -
    LastByteRead]
```

❖ rcvr advertises spare room by including value of `RcvWindow` in segments

❖ sender limits unACKed data to `RcvWindow`

▪ guarantees receive buffer doesn't overflow

# Consideration

❖ What is the relationship between the variable LastByteRcvd and the variable y in section 3.5.4

# Problems caused by flow control

❖ When RecvWin=0, sender doesn't send data. How does the sender restart to send ?

❖ TCP's Solution: RecvWin management is not directly tied to the acks. When RecvWin=0, sender doesn't send data unless (1)urgent data; (2)1-byte segment to make the receiver announce next window size.

# TCP Transmission Policy(1)

❖ The trade-off in determining the optimal size of TCP segments.

  - Too large : IP level fragment
  - Too small : Reduce communication Performance

# TCP Transmission Policy(2)

❖ Problem1: Application data come in byte-by-byte at sender. (Recall telnet)

❖ Nagle's algorithm: send the first byte, wait for ACK, and accumulate bytes in buffer.

# TCP Transmission Policy(3)

❖ Problem2:silly window syndrome.

- The receiver receives a large block->buffer is full;

- The app reads buffer byte-by-byte -> receiver sends one window update to sender after processing one byte -> sender sends data byte by byte.

# TCP Transmission Policy(4)



Silly Window Syndrome

# TCP Transmission Policy(5)

❖ Clark's solution: receiver waits until it has a decent amount of buffer available. The receiver should not send a window update until it can handle the MSS.

# Chapter 3 outline

# 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 Connection Management

**Recall:** TCP sender, receiver establish "connection" before exchanging data segments

❖ initialize TCP variables:
  ▪ seq. #s
  ▪ buffers, flow control info (*e.g.* `RcvWindow`)

❖ *client:* connection initiator

```
Socket clientSocket = new
Socket("hostname","port

number");
```

❖ *server:* contacted by client

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

## Three way handshake:

Step 1: client host sends TCP SYN segment to server
  ▪ specifies initial seq #
  ▪ no data

Step 2: server host receives SYN, replies with SYNACK segment
  ▪ server allocates buffers
  ▪ specifies server initial seq. #

Step 3: client receives SYNACK, replies with ACK segment, which may contain data

# TCP Connection Management (cont.)

## Closing a connection:

client closes socket:
    `clientSocket.close();`

**Step 1:** client end system sends TCP FIN control segment to server

**Step 2:** server receives FIN, replies with ACK. Closes connection, sends FIN.

client     server

close    FIN

   ACK

   close

   FIN

timed wait    ACK

closed

# TCP Connection Management (cont.)

**Step 3:** client receives FIN, replies with ACK.

- Enters "timed wait" - will respond with ACK to received FINs

**Step 4:** server, receives ACK. Connection closed.

**Note:** with small modification, can handle simultaneous FINs.

client          server

closing

FIN

ACK

closing

FIN

timed wait

ACK

closed

closed

# TCP Connection Management (cont)



TCP server lifecycle

TCP client lifecycle

# 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

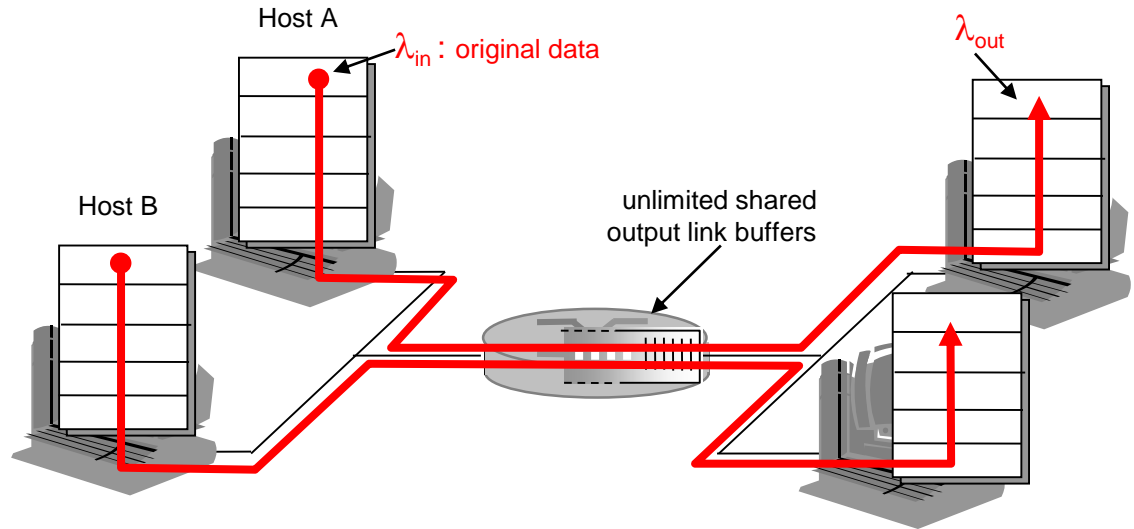3.6 Principles of congestion control

3.7 TCP congestion control

# Principles of Congestion Control

Congestion:

- ❖ 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)
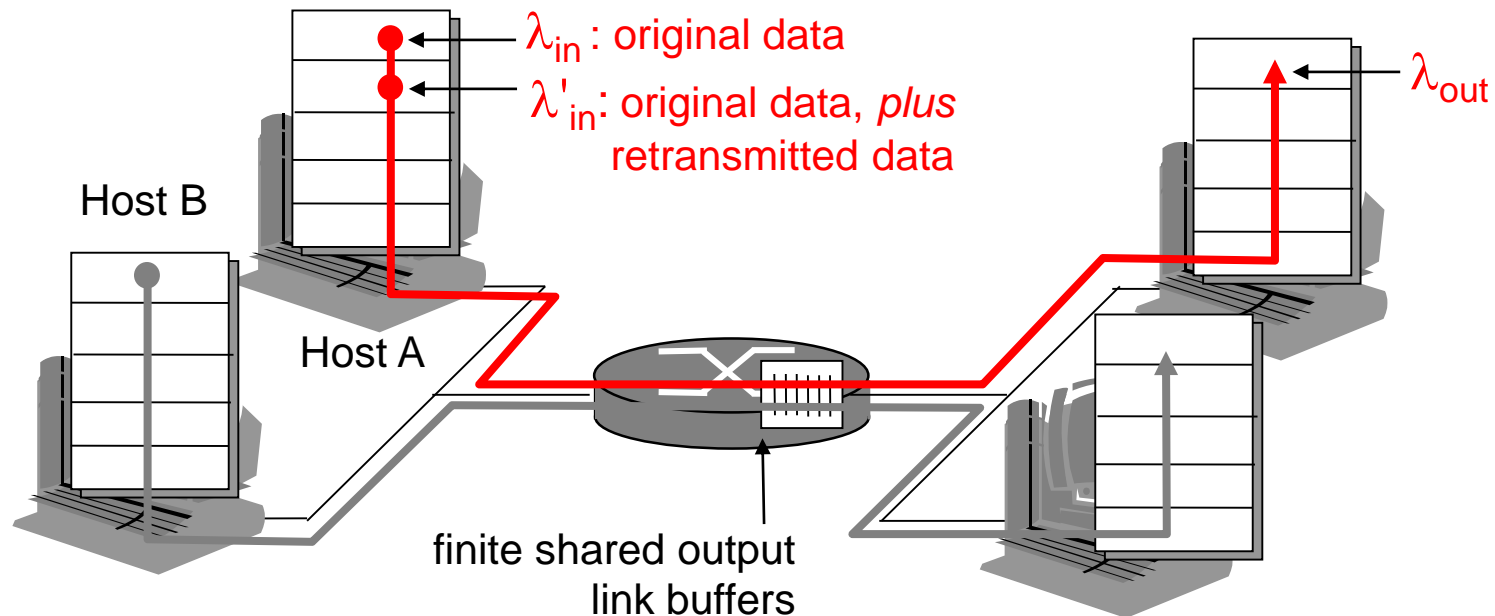- ❖ a top-10 problem!

# Causes/costs of congestion: scenario 1

- ❖ two senders, two receivers
- ❖ one router, infinite buffers
- ❖ no retransmission

Host A

$\lambda_{in}$ : original data

$\lambda_{out}$

Host B

unlimited shared output link buffers

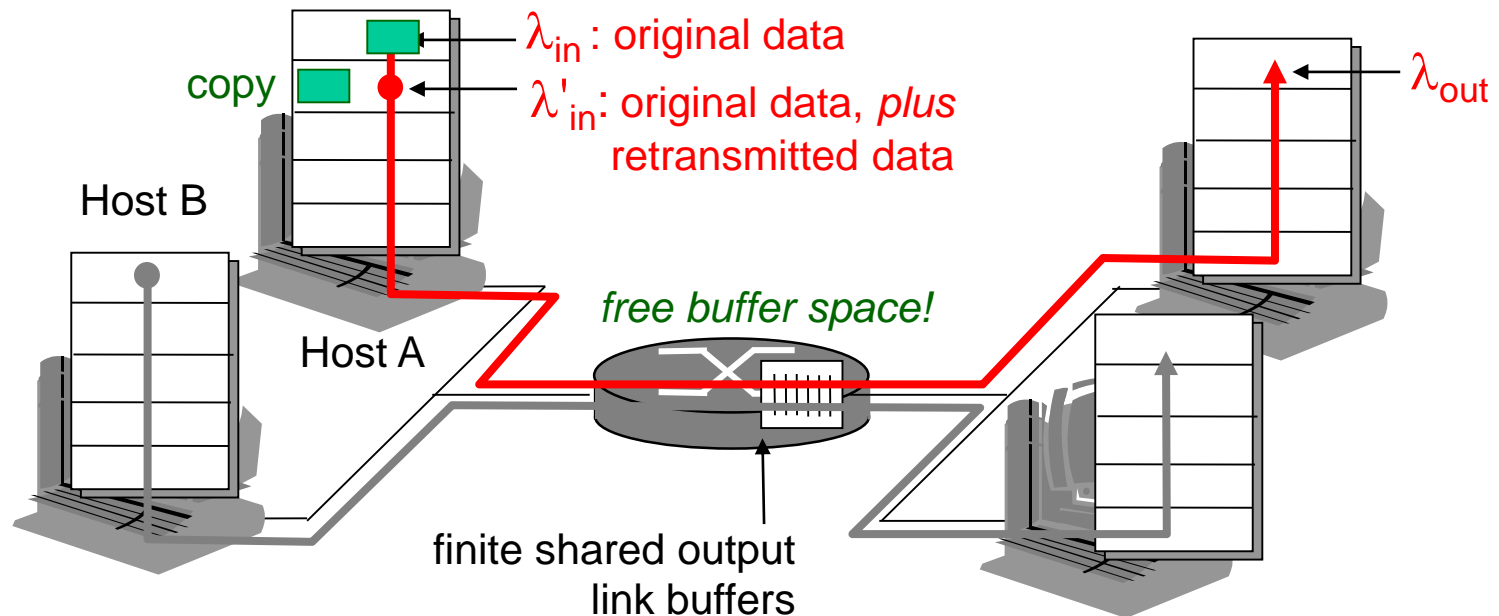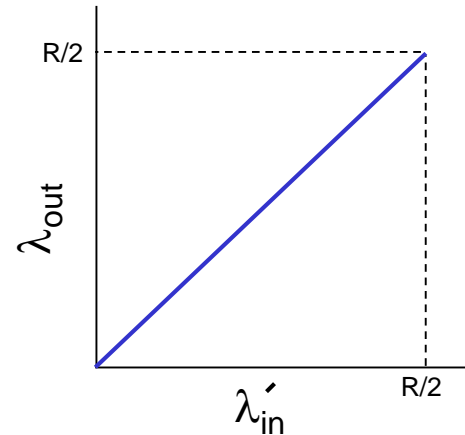- ❖ large delays when congested
- ❖ maximum achievable throughput

# 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* : $\lambda'_{in} \geq \lambda_{in}$
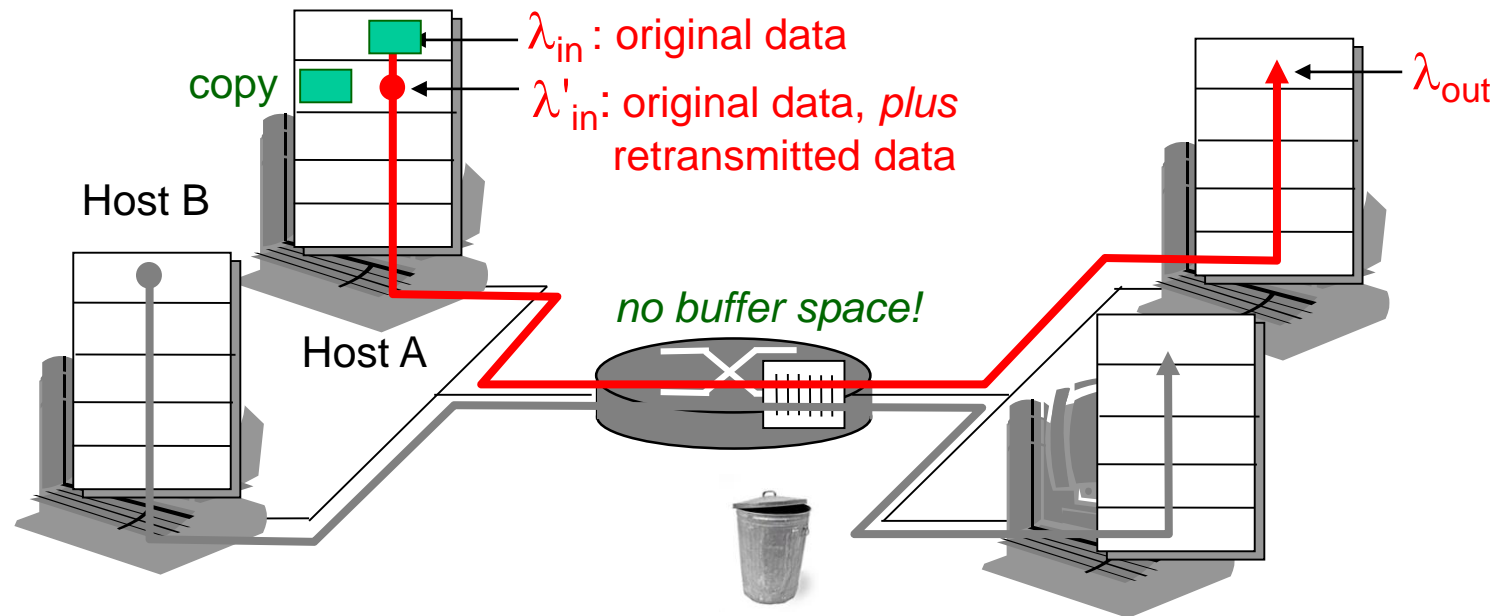


$\lambda_{in}$ : original data

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

$\lambda_{out}$

Host B

Host A

finite shared output link buffers

# Congestion scenario 2a: ideal case

❖ **sender sends only when router buffers available**



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

$\lambda_{out}$

copy

Host B

Host A

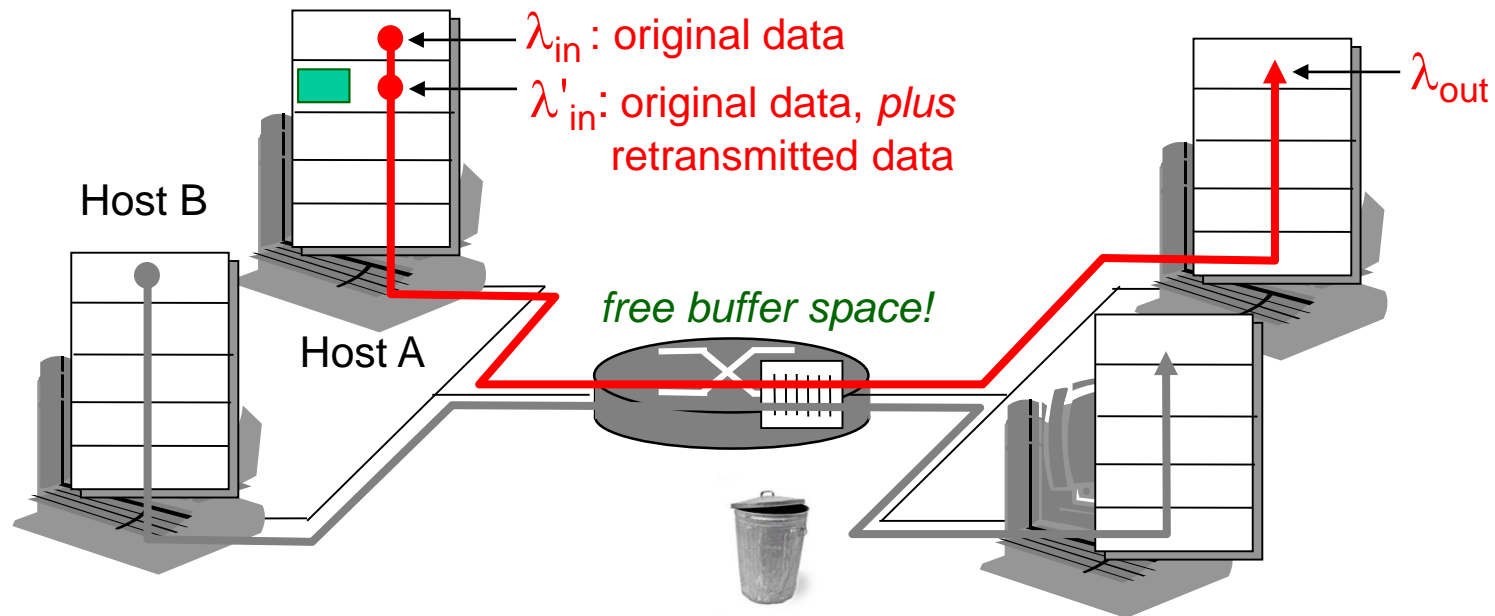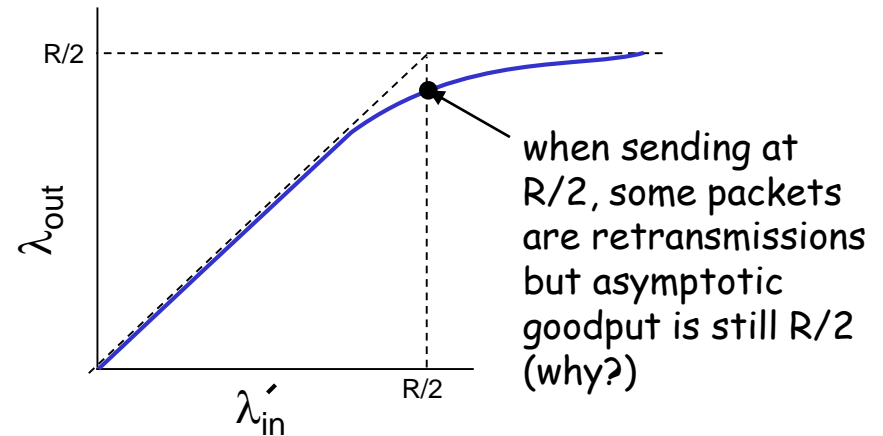*free buffer space!*

finite shared output link buffers

# Congestion scenario 2b: *known loss*

❖ packets may get dropped at router due to full buffers
  ▪ sometimes lost
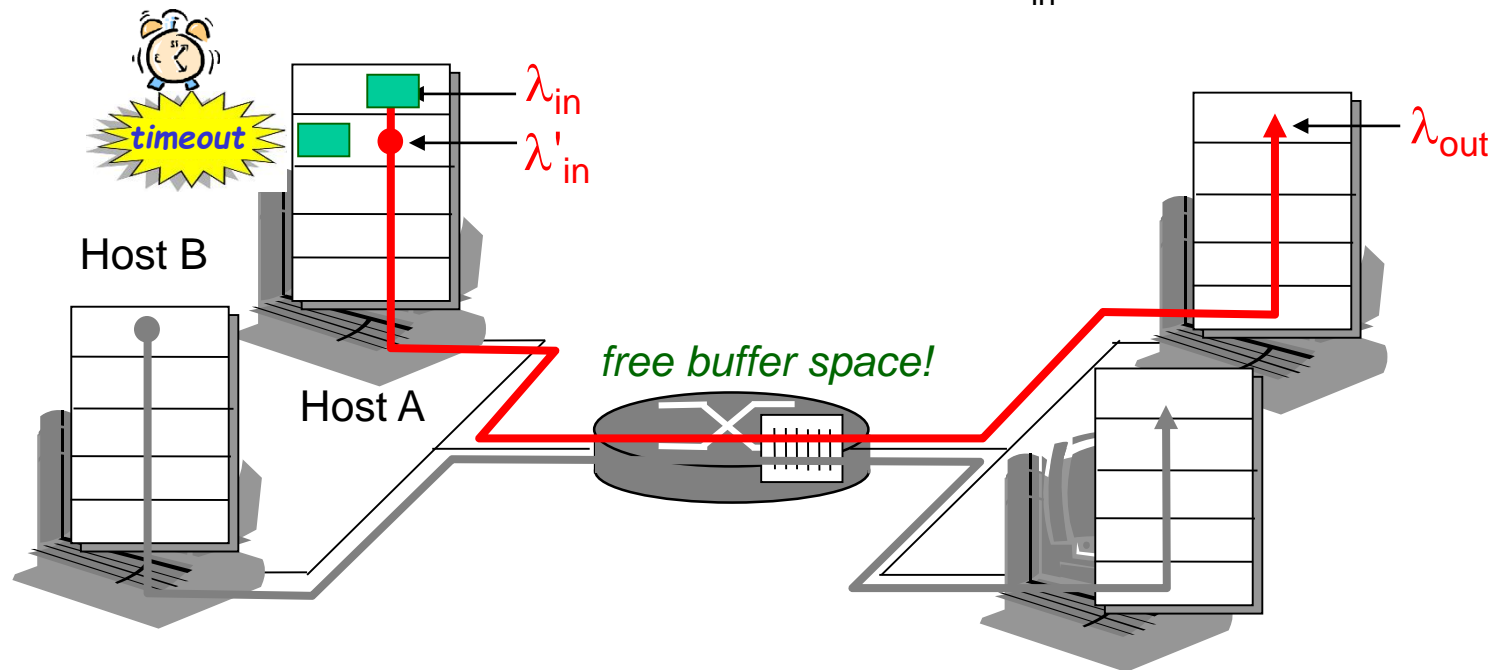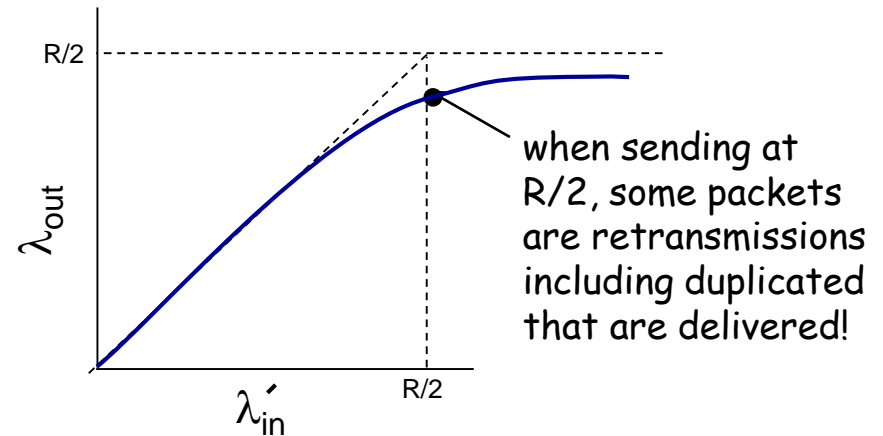❖ sender only resends if packet *known* to be lost (admittedly idealized)

$\lambda_{in}$ : original data

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

copy

$\lambda_{out}$

Host B

Host A

*no buffer space!*

# Congestion scenario 2b: _known loss_

❖ **packets may get dropped at router due to full buffers**

  ▪ sometimes not lost

❖ **sender only resends if packet _known_ to be lost (admittedly idealized)**



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

$\lambda_{out}$
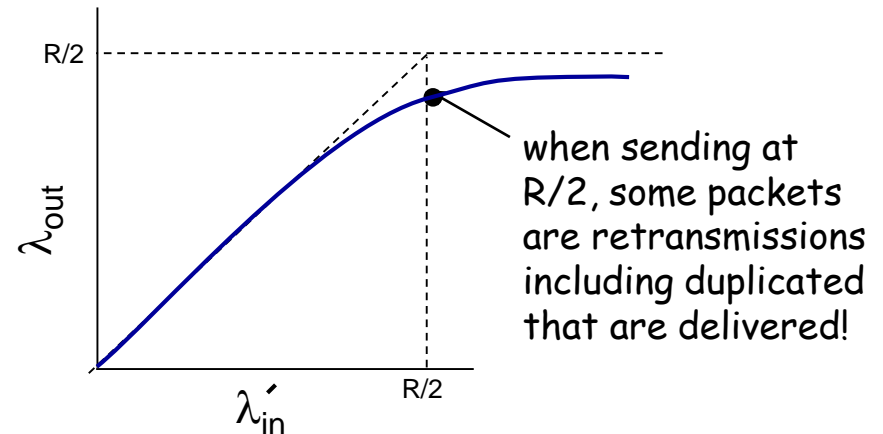
Host B

Host A

free buffer space!

# Congestion scenario 2c: *duplicates*

- ❖ packets may get 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!

timeout

Host B

Host A

free buffer space!

$\lambda_{in}$

$\lambda'_{in}$

$\lambda_{out}$

# Congestion scenario 2c: *duplicates*

❖ packets may get dropped at router due to full buffers

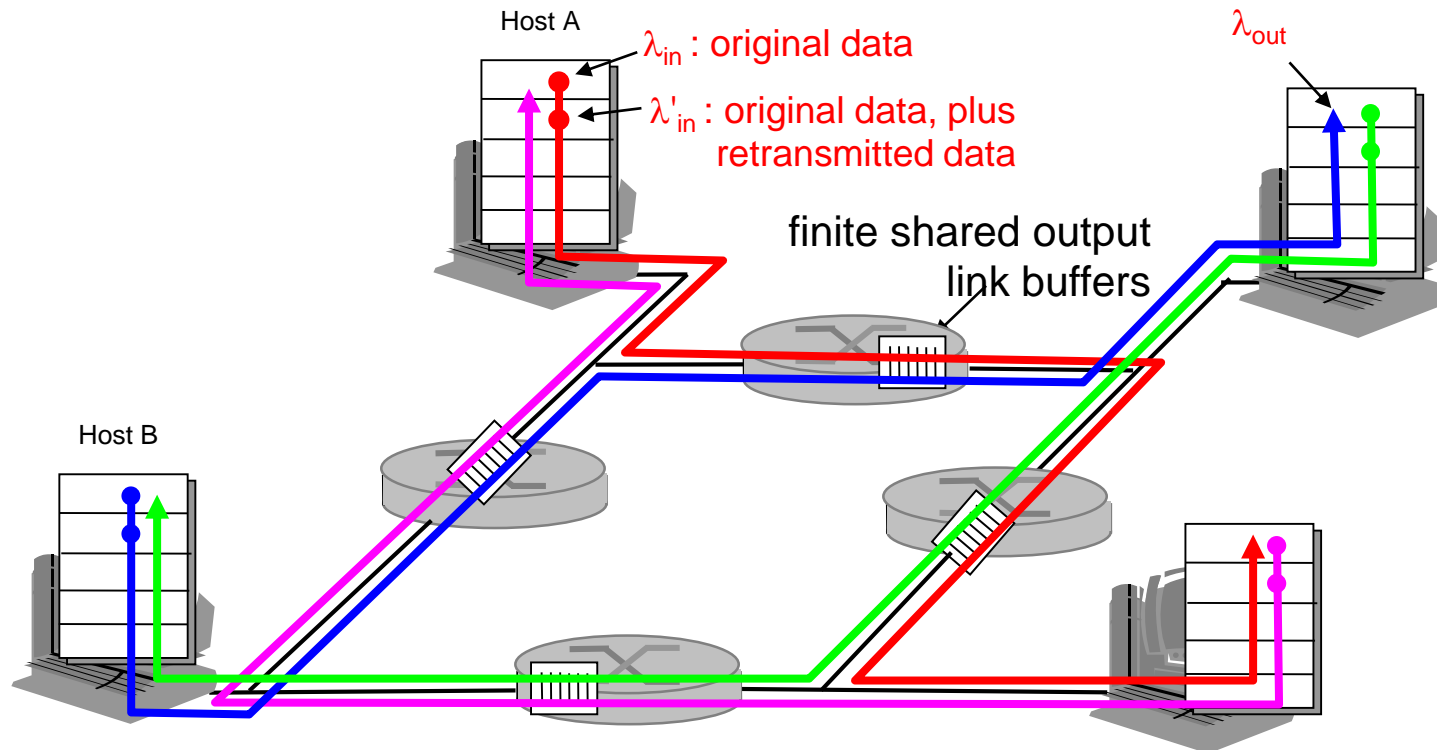❖ sender times out prematurely, sending *two* copies, both of which are delivered



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

R/2

$\lambda_{out}$

$\lambda'_{in}$

R/2

**"costs" of congestion:**

❖ more work (retrans) for given "goodput"

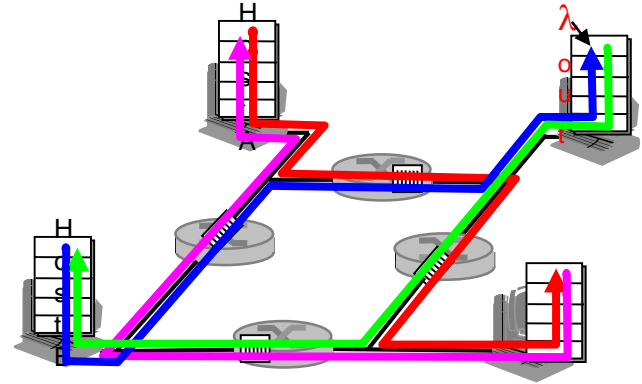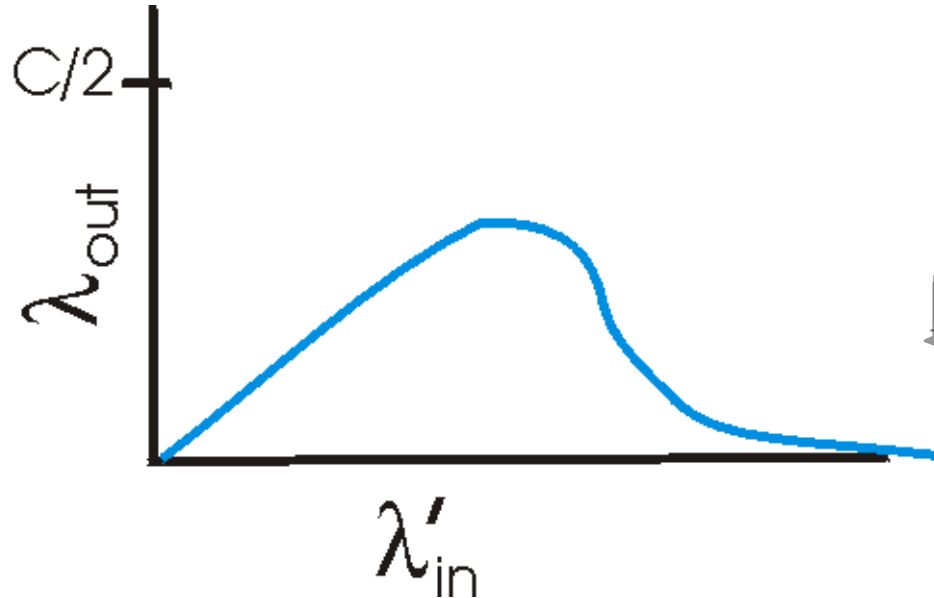❖ unneeded retransmissions: link carries multiple copies of pkt

  ▪ decreasing goodput

# Causes/costs of congestion: scenario 3

❖ four senders

❖ multihop paths

❖ timeout/retransmit

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

Host A

$\lambda_{in}$ : original data

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

$\lambda_{out}$

finite shared output link buffers

Host B

# 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

❖ congestion inferred from end-system observed loss, delay

❖ approach taken by TCP

### network-assisted congestion control:

❖ routers provide feedback to end systems

  ▪ single bit indicating congestion (SNA, DECbit, TCP/IP ECN, ATM)

  ▪ explicit rate sender should send at

# Case study: ATM ABR congestion control

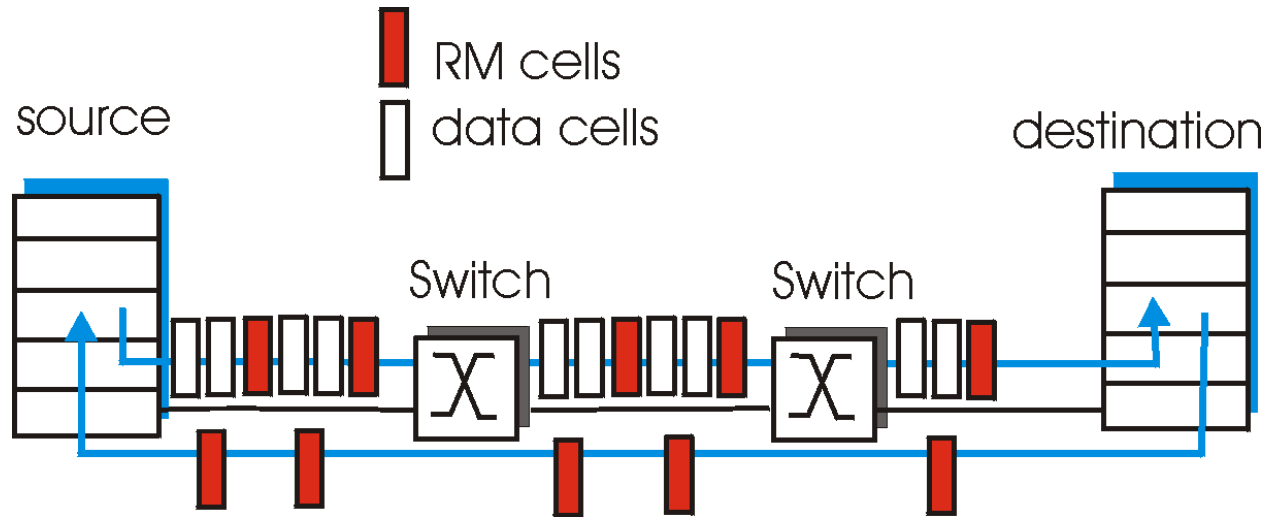## ABR: available bit rate:

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

## RM (resource management) cells:

- sent by sender, interspersed with data cells
- bits in RM cell set by switches ("*network-assisted*")
  - NI bit: no increase in rate (mild congestion)
  - CI bit: congestion indication
- RM cells returned to sender by receiver, with bits intact

# Case study: ATM ABR congestion control



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

# 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
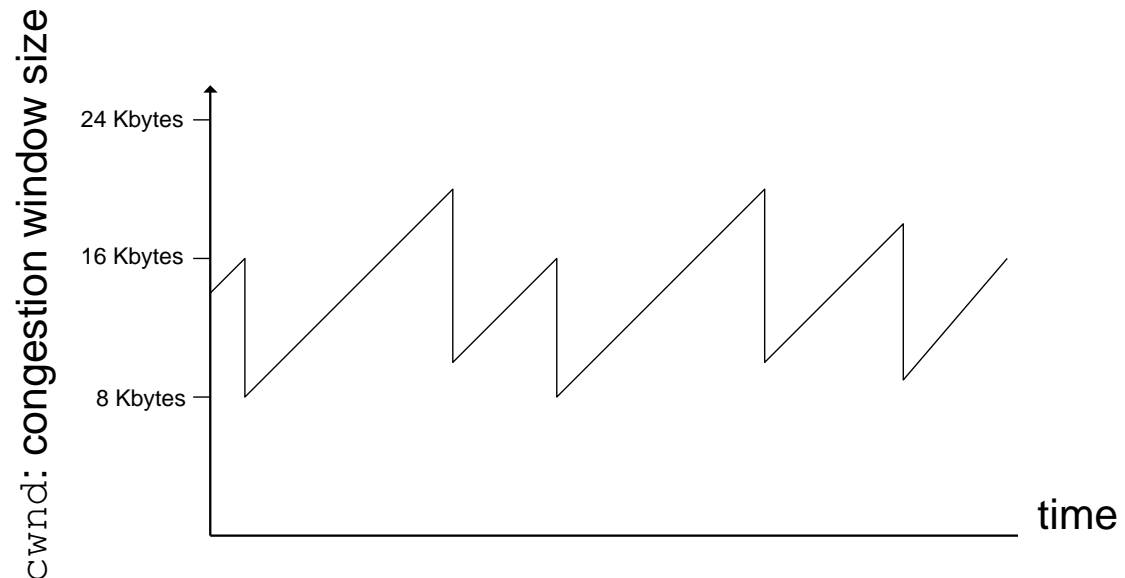- flow control
- connection management

3.6 Principles of congestion control

3.7 TCP congestion control

# TCP congestion control: additive increase, multiplicative decrease

❖ *approach:* increase 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

saw tooth behavior: probing for bandwidth

cwnd: congestion window size

24 Kbytes —

16 Kbytes —

8 Kbytes —

time

# TCP Congestion Control: details

❖ sender limits transmission:

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

❖ roughly,

$$\boxed{\text{rate} = \frac{\text{cwnd}}{\text{RTT}} \text{ Bytes/sec}}$$

❖ **cwnd** is dynamic, function of perceived network congestion
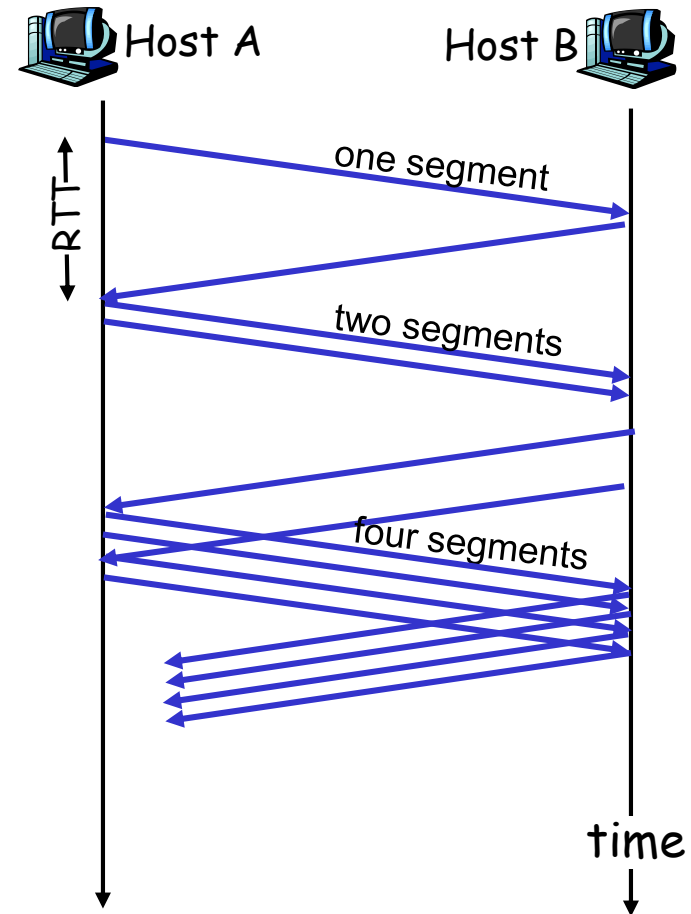
## How does sender perceive congestion?

❖ loss event = timeout *or* 3 duplicate acks

❖ TCP sender reduces rate (**cwnd**) after loss event

## three mechanisms:

- ▪ AIMD
- ▪ slow start
- ▪ conservative after timeout events

# 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

❖ <u>summary:</u> initial rate is slow but ramps up exponentially fast

Host A                          Host B

one segment

two segments

four segments

time

# Refinement: inferring loss

❖ after 3 dup ACKs:
  ▪ `cwnd` is cut in half
  ▪ window then grows linearly

❖ <u>but</u> after timeout event:
  ▪ `cwnd` instead set to 1 MSS;
  ▪ window then grows exponentially
  ▪ to a threshold, then grows linearly

Philosophy:

❖ 3 dup ACKs indicates network capable of delivering some segments

❖ timeout indicates a "more alarming" congestion scenario

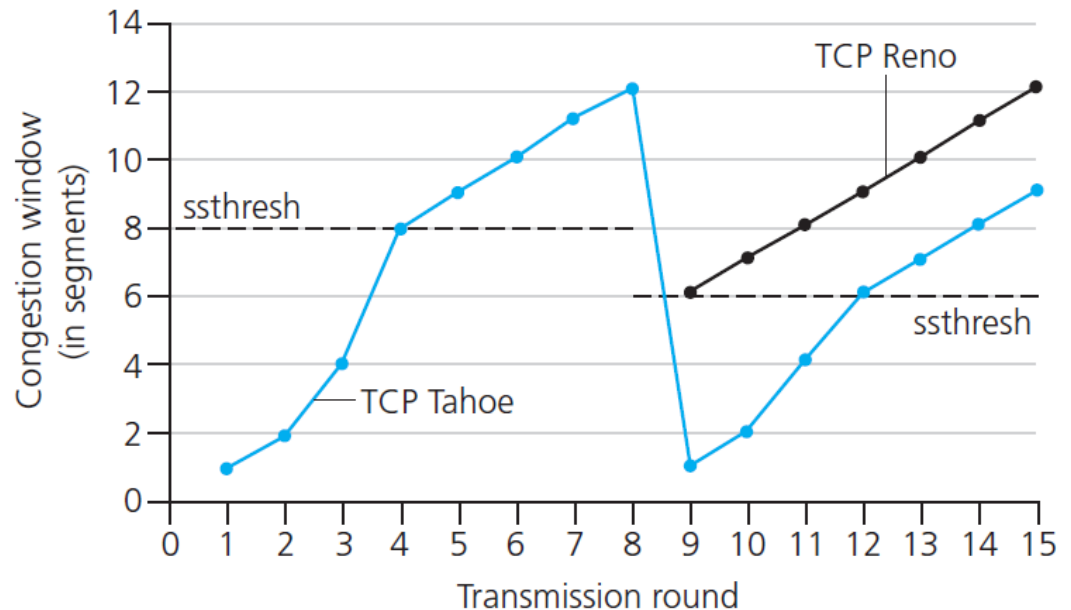# Refinement

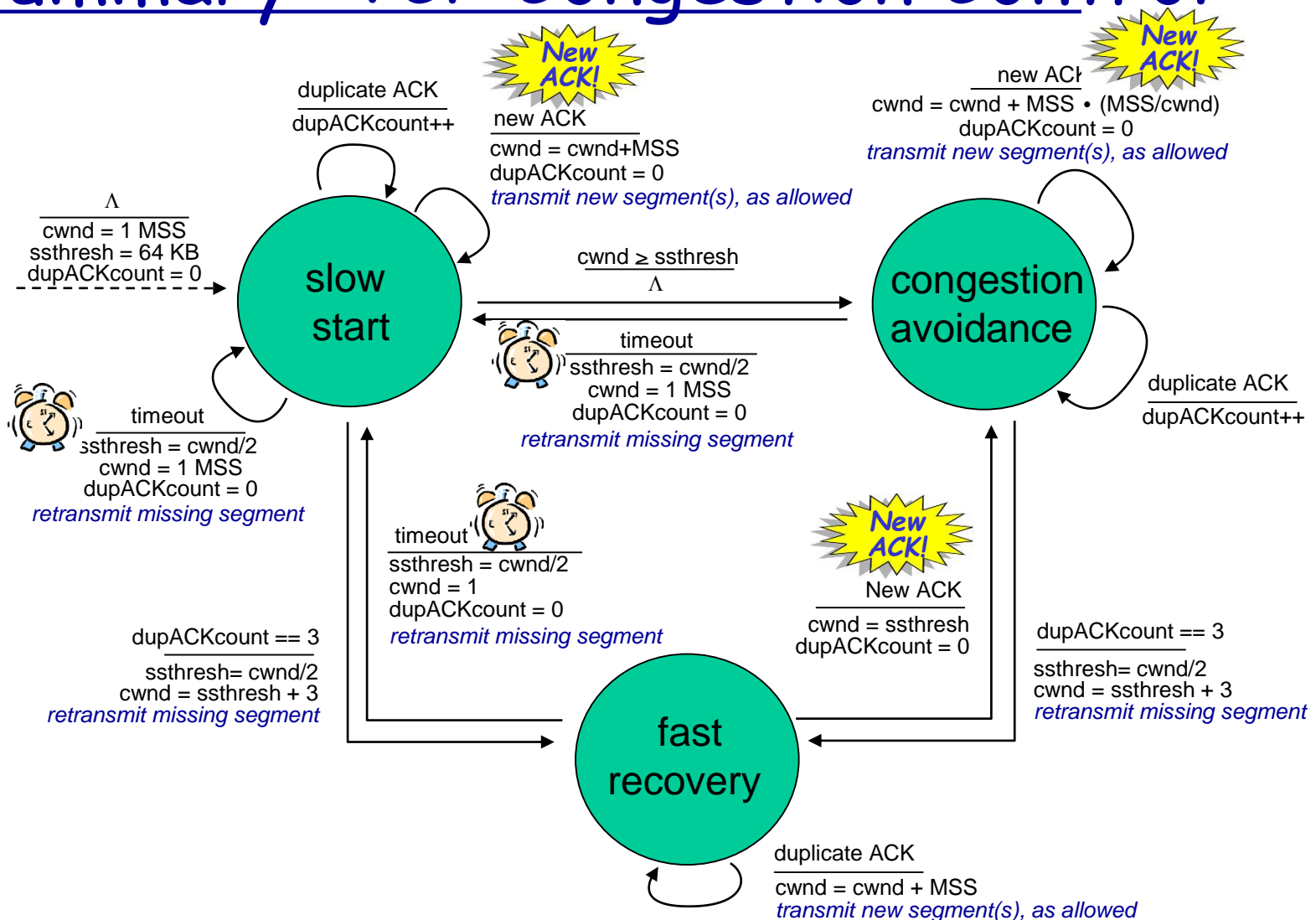Q: when should the exponential increase switch to linear?

A: when `cwnd` gets to 1/2 of its value before timeout.



## Implementation:

❖ variable `ssthresh`

❖ on loss event, `ssthresh` is set to 1/2 of `cwnd` just before loss event

# Summary: TCP Congestion Control

**New ACK!**

**New ACK!**

duplicate ACK
dupACKcount++

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

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

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

**slow start**

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*

# TCP throughput

❖ what's the average throughout of TCP as a function of window size and RTT?

  ▪ ignore slow start

❖ let W be the window size when loss occurs.

  ▪ when window is W, throughput is W/RTT

  ▪ just after loss, window drops to W/2, throughput to W/2RTT.

  ▪ average throughout: .75 W/RTT

# TCP Futures: TCP over "long, fat pipes"
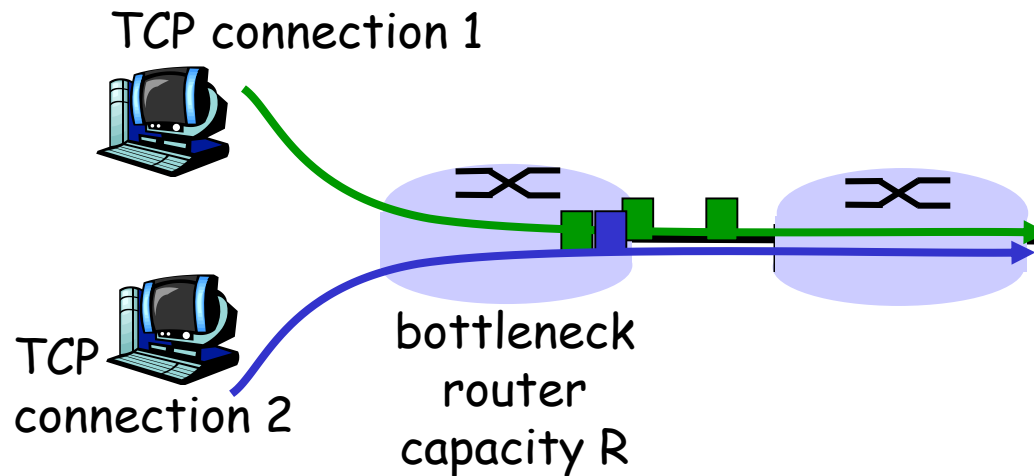
❖ example: 1500 byte segments, 100ms RTT, want 10 Gbps throughput

❖ requires window size W = 83,333 in-flight segments

❖ throughput in terms of loss rate:

$$\frac{1.22 \cdot MSS}{RTT \sqrt{L}}$$

❖ ➔ L = 2·10⁻¹⁰  *Wow – a very small loss rate!*
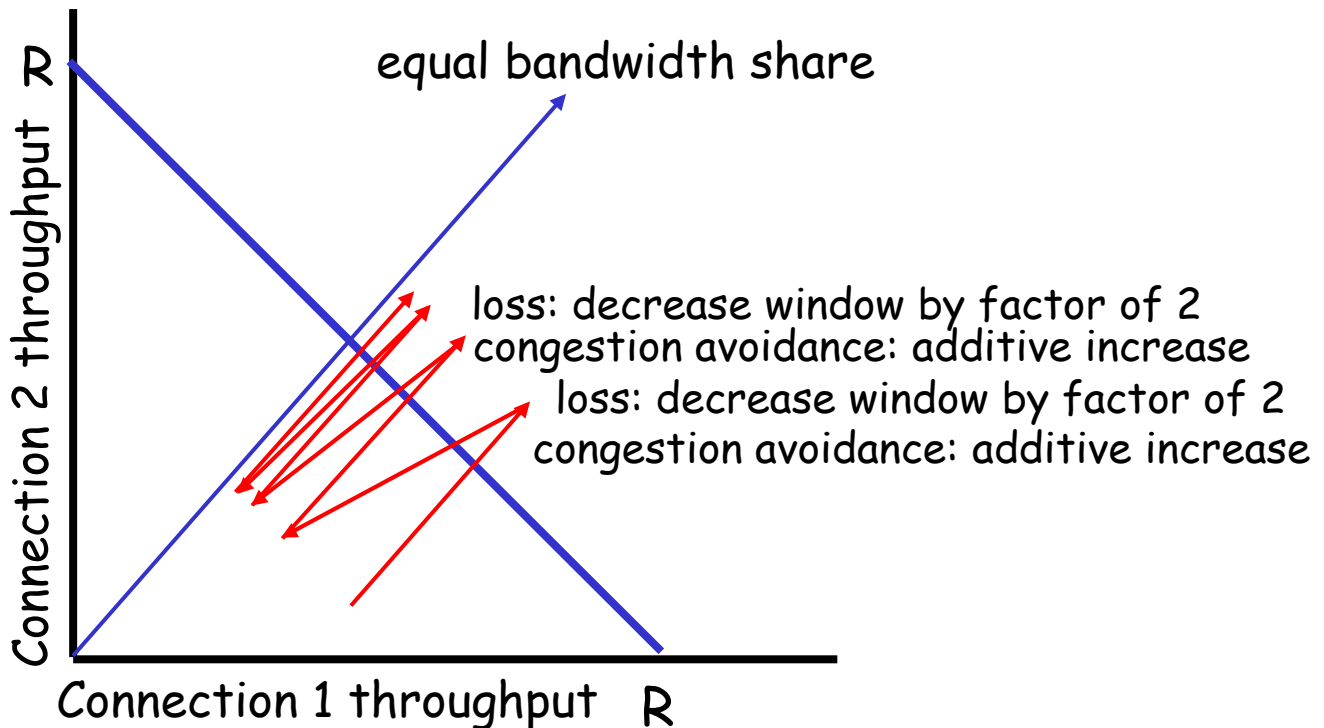
❖ new versions of TCP for high-speed

# TCP Fairness

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

# 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

Connection 1 throughput   R

# Fairness (more)

## Fairness and UDP

❖ **multimedia apps often do not use TCP**
  - do not want rate throttled by congestion control

❖ **instead use UDP:**
  - pump audio/video at constant rate, tolerate packet loss

## Fairness and parallel TCP connections

❖ nothing prevents app from opening parallel connections between 2 hosts.

❖ web browsers do this

❖ example: link of rate R supporting 9 connections;
  - new app asks for 1 TCP, gets rate R/10
  - new app asks for 11 TCPs, gets R/2 !

# Chapter 3: Summary

❖ principles behind transport layer services:

- multiplexing, demultiplexing
- reliable data transfer
- flow control
- congestion control

❖ instantiation and implementation in the Internet

- UDP
- TCP

Next:

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

❖ into the network "core"