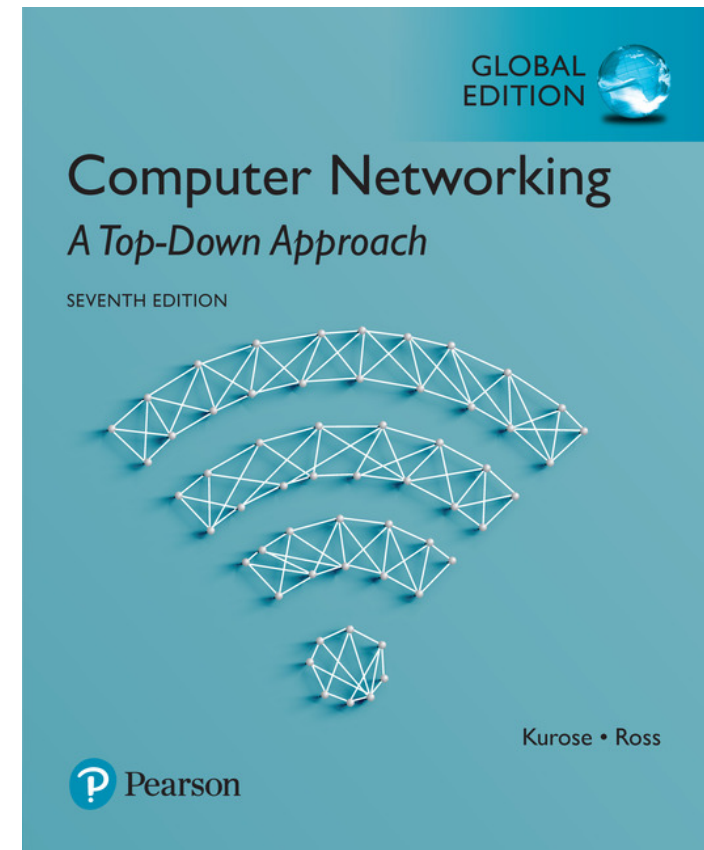# Chapter 3
# Transport Layer

## A note on the use of these Powerpoint slides:

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

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

Thanks and enjoy! JFK/KWR

*Computer Networking: A Top Down Approach*

7th Edition, Global Edition
Jim Kurose, Keith Ross
Pearson
April 2016

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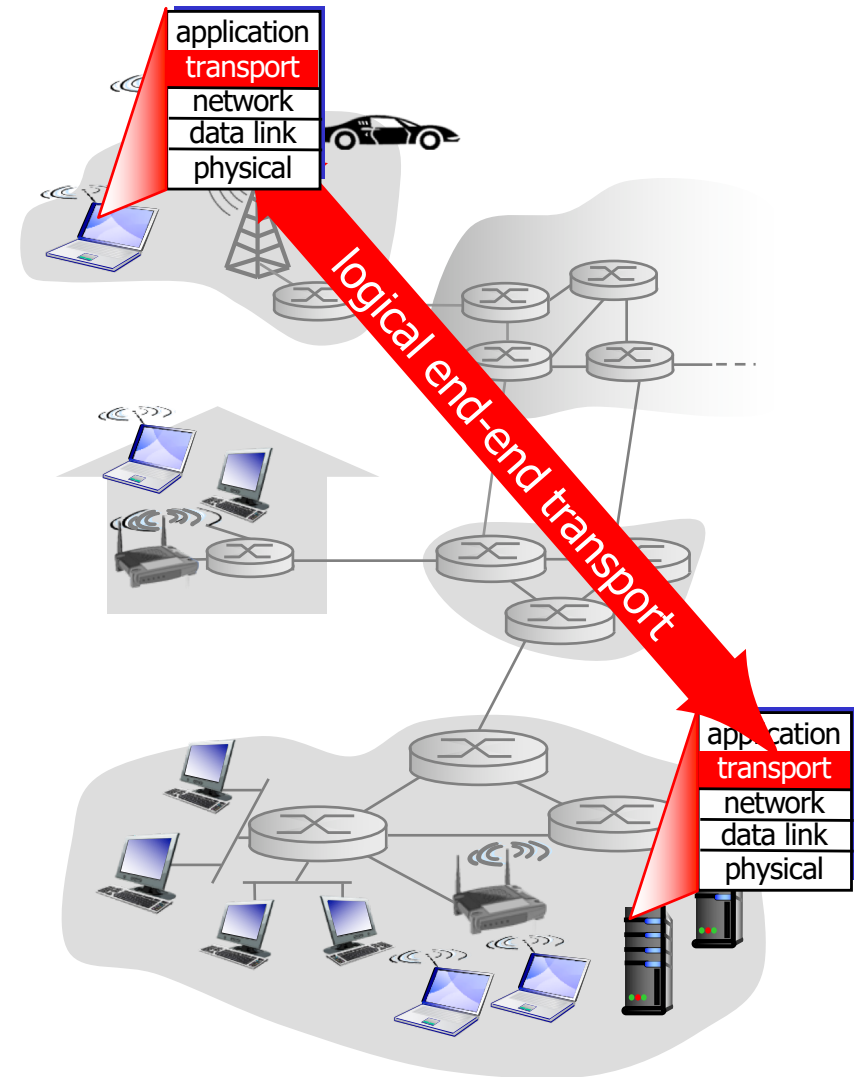
3.5 connection-oriented transport: TCP
  - segment structure
  - reliable data transfer
  - flow control
  - connection management

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

- 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:* communication between hosts

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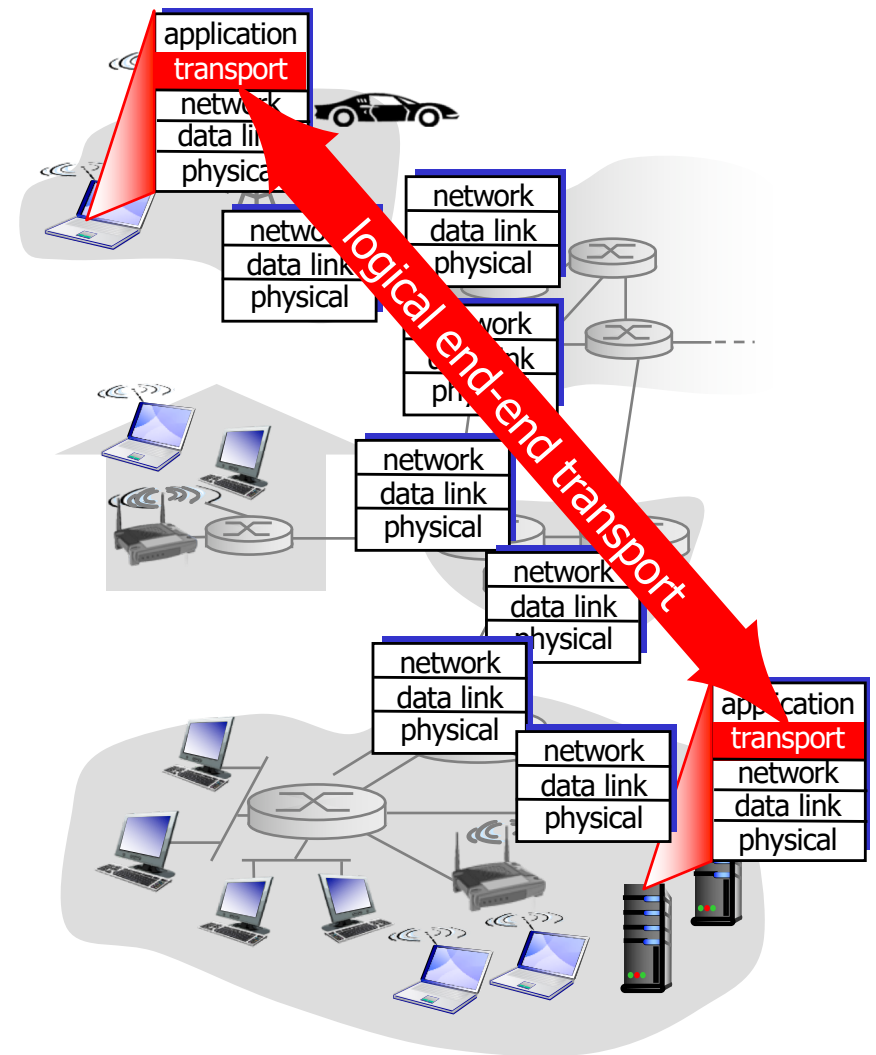  - relies on, enhances, network layer services

*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

# 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

application

P1    P2

transport
network
link
physical

application

P3

transport
network
link
physical

application

P4

transport
network
link
physical

socket

process

# Quiz Time!

# How demultiplexing works

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

32 bits

| source port # | dest port # |
|---|---|

other header fields

application
data
(payload)

TCP/UDP segment format

# Connectionless demultiplexing

- UDP socket identified by two-tuple:

  (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: example

application
P3
transport
network
link
physical

**Port# 9157**

application
P1
transport
network
link
physical

**Port# 6428**

application
P4
transport
network
link
physical

**Port# 5775**

source port: 6428
dest port: 9157

source port: ?
dest port: ?

source port: 9157
dest port: 6428

source port: ?
dest port: ?

# 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



source IP,port: A,9157
dest IP, port: B,12004

source IP,port: C,5775
dest IP,port: B,12006

source IP,port: C,9157
dest IP,port: B,12005

server: IP address B

host: IP address A

host: IP address C

three segments, all destined to IP address: B, demultiplexed to *different* sockets

# Connection-oriented demux: example

threaded server

application

P4

application
P3

application
P2    P3

transport

transport

transport

network

network

network

link

link

link

physical

physical

physical

server: IP
address B

host: IP
address A

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

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

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

host: IP
address C

# Chapter 3 outline

# 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
  - SNMP
- reliable transfer over UDP:
  - add reliability at application layer
  - application-specific error recovery!

# UDP: segment header



32 bits

length, in bytes of UDP segment, including header

| 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

# 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

# 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)**

# Principles of reliable data transfer

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



(a) provided service

(b) service implementation

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

# Principles of reliable data transfer

- **important in application, 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()` ↓ 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

# 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 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
  - retransmit in case of NAK

# 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)
‾‾‾‾‾‾‾‾
sndpkt = 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)
_____
sndpkt = 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!

## What to do?

- Q?
- Checksum the ACK/NAK?
- ACK/NAK the ACK/NAK?
- Retransmit the ACK/NAK?

## What happens if ACK/NAK of ACK/NAK corrupted?

- >"<

## How did we deal with corrupted data?

- Checksum data
- ACK/NAK the data
- Retransmit data

# rdt2.0 has a fatal flaw!

## what happens if ACK/NAK corrupted?

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

## Handling error on the ACK/NAK direction in rdt2.1

- 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)

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

---

udt_send(sndpkt)

**Wait for call 0 from above**

**Wait for ACK or NAK 0**

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

---

Λ

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

---

Λ

**Wait for ACK or NAK 1**

**Wait for call 1 from above**

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

---

udt_send(sndpkt)

rdt_send(data)

---

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

# rdt2.1: receiver, 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)

# Double Quiz Time!

# rdt2.1: discussion

- must check if received ACK/NAK corrupted
- seq # added to pkt
  - two seq. #'s (0,1) will suffice.  Why?
- twice as many states
  - state must "remember" whether "expected" pkt should have seq # of 0 or 1

- 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
  - If duplicate, not pass the data up, but sends ACK again

# 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*

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

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
&& has_seq0(rcvpkt)
_____
extract(rcvpkt,data)                     0
deliver_data(data)
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

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

ACK 1

rdt_rcv(rcvpkt) &&
   not corrupt(rcvpkt) &&
   has_seq1(rcvpkt)              1
_____
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

**Wait for 0 from below**

**Wait for 1 from below**

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

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

rdt_rcv(rcvpkt) && notcorrupt(rcvpkt)
&& has_seq1(rcvpkt)
_____
extract(rcvpkt,data)                     1
deliver_data(data)
sndpkt = make_pkt(ACK, chksum)
udt_send(sndpkt)

# rdt2.2: receiver fragments

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(1, ACK, chksum)**
udt_send(sndpkt)

# 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)   isACK 1

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& isACK(rcvpkt)
_____
Λ   1

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& isACK(rcvpkt)
_____
Λ   0

**Wait for ACK or NAK 1**

**Wait for call 1 from above**

rdt_rcv(rcvpkt) &&
( corrupt(rcvpkt) ||
isNAK(rcvpkt) )
_____
udt_send(sndpkt)   isACK 0

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

# rdt2.2: sender, receiver fragments

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

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

Wait for call 0 from above

Wait for ACK 0

**sender FSM fragment**

rdt_rcv(rcvpkt)
&& notcorrupt(rcvpkt)
&& **isACK(rcvpkt,0)**
_____
$\Lambda$

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(1, ACK, chksum)**
udt_send(sndpkt)

# Quiz Time!

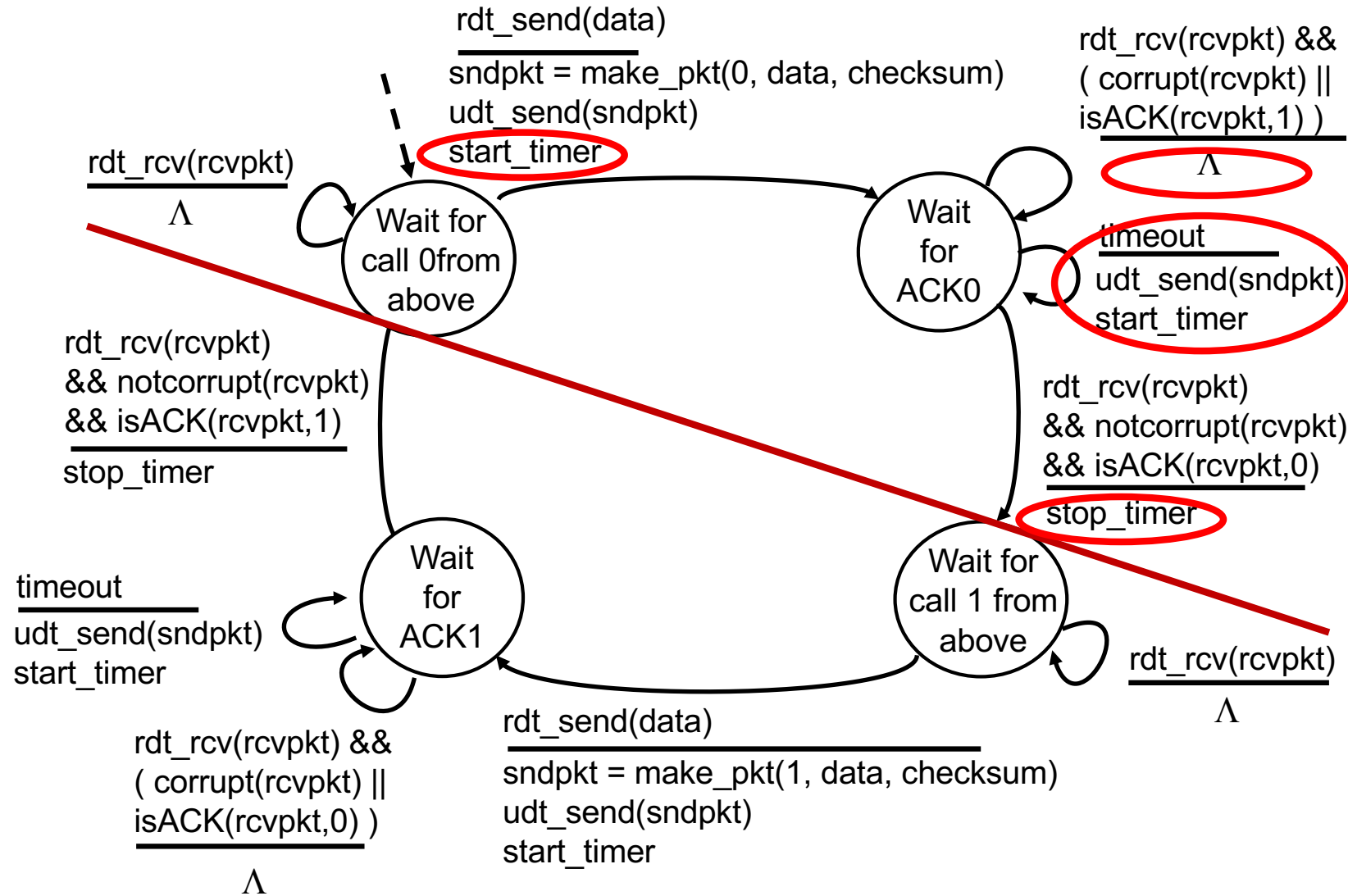# rdt3.0: channels with errors *and* loss

**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
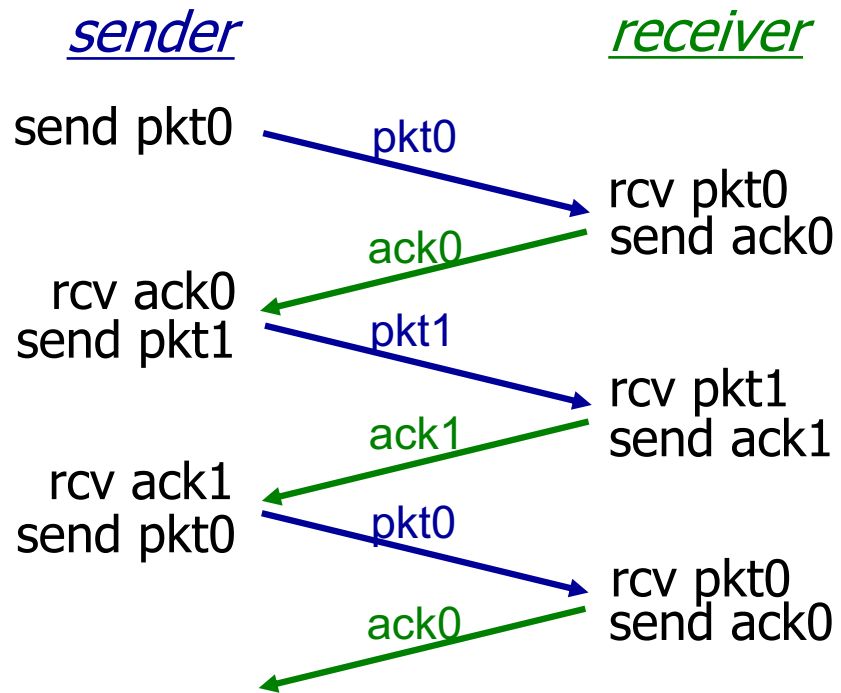
- retransmits if no ACK received in this time
- requires countdown timer

- if pkt (or ACK) just delayed (not lost):
  - retransmission will be duplicate. Good that seq. #'s already handling this
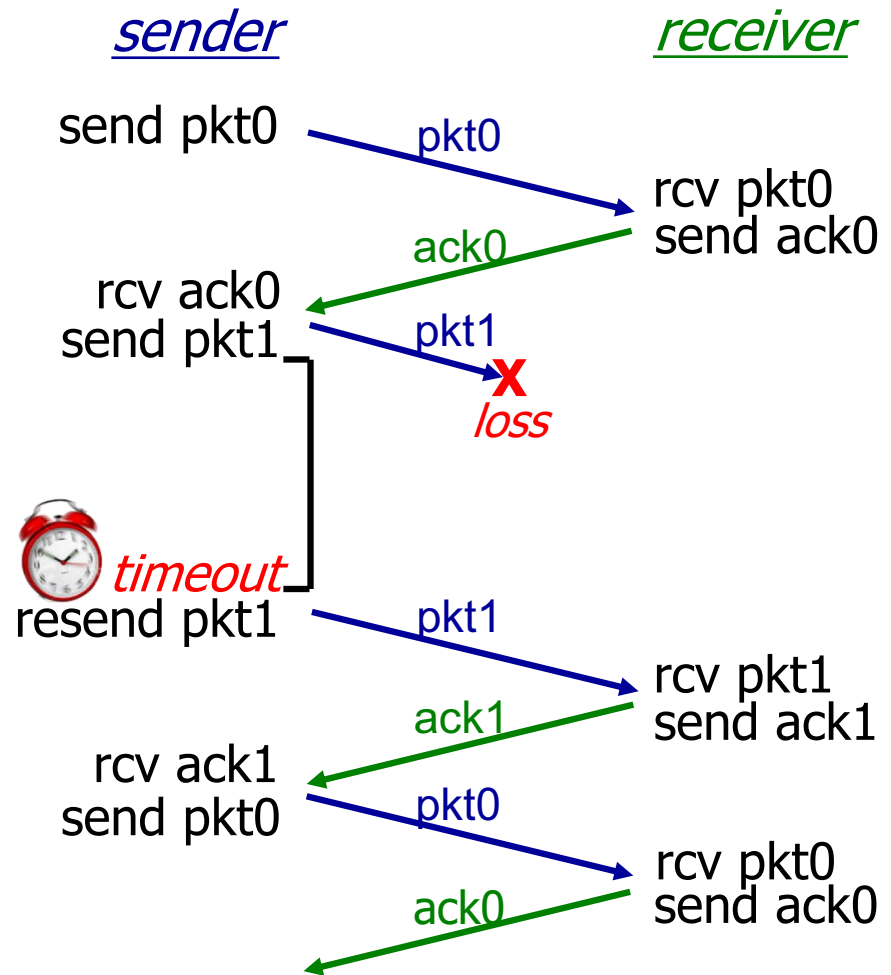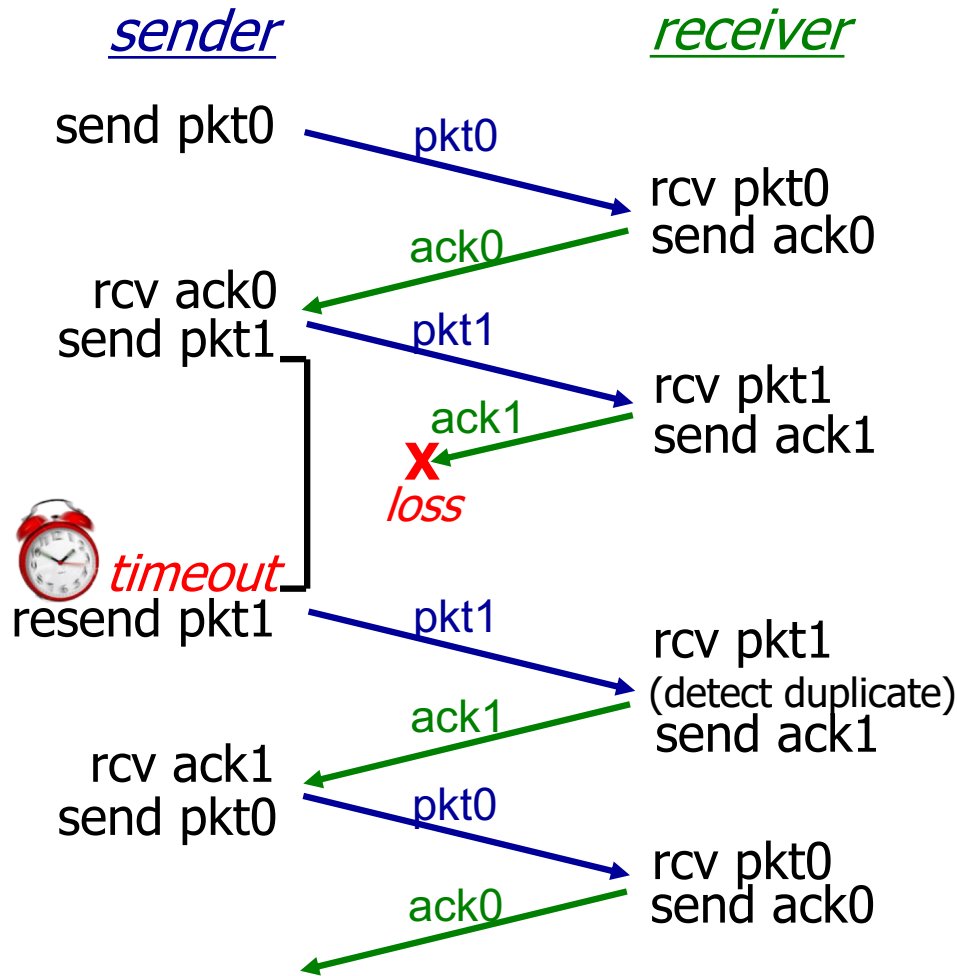  - receiver must specify seq # of pkt being ACKed

# rdt3.0 sender

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

rdt_rcv(rcvpkt)
—————
Λ

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

**Wait for call 0from above**

**Wait for ACK0**

timeout
—————
udt_send(sndpkt)
start_timer

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

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

timeout
—————
udt_send(sndpkt)
start_timer

**Wait for ACK1**

**Wait for call 1 from above**

rdt_rcv(rcvpkt)
—————
Λ

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

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

# Quiz Time!

# 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 ✗ loss

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

rcv ack1 ← 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 → pkt0 send ack1
send pkt0 ← ack1

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

(d) premature timeout/ delayed ACK

# rdt3.0 in action



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 ←—— pkt0      send ack1
send pkt0   ack1

rcv ack1      rcv pkt0
send pkt0 —— ack0 pkt0 ——→ send ack0
     rcv pkt0
     ack0      (detect duplicate)
     send ack0

(e) premature timeout/ delayed ACK
with rdt2.2-like rdt3.0

# Performance of rdt3.0

- rdt3.0 is correct, but performance stinks
- e.g.: 1 Gbps link, 15 ms prop. delay, 8000 bit packet:

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

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

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

# 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

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

- Throughput on 1 Gbps link = 0.00027 of 1Gbps = 270 kbps
- Network protocol limits use of physical resources!

# 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

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

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*

# Pipelined protocols: overview

## Go-back-N:

- sender can have up to N unack'ed 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
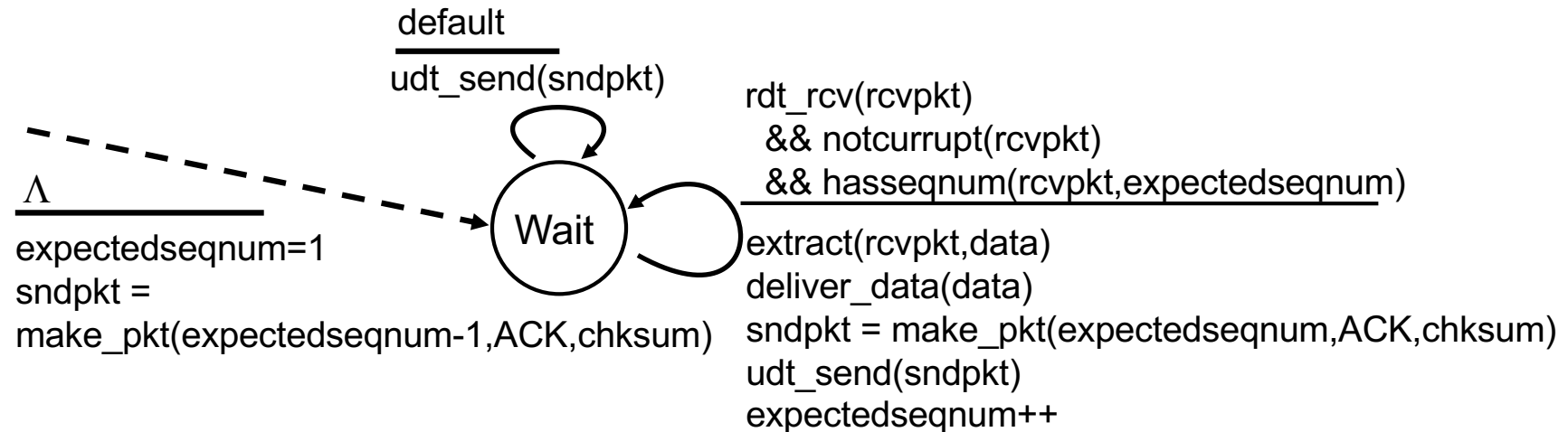- "window" of up to N, consecutive unack'ed pkts allowed



- ACK(n):ACKs all pkts up to, including seq # n - *"cumulative ACK"*
  - may receive duplicate ACKs (see receiver)
- timer for oldest in-flight, unack'ed 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)
───────────────────
$\Lambda$

rdt_rcv(rcvpkt) &&
  notcorrupt(rcvpkt)
───────────────────
```
base = getacknum(rcvpkt)+1
If (base == nextseqnum)
    stop_timer
  else
    re-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-1,ACK,chksum)

**Wait**

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

- **Corrupted data or out-of-order data pkt**
  - discard (don't buffer): *no receiver buffering!*
  - re-ACK pkt with highest in-order seq #

- **Odd but simple**
  - ACK for highest *in-order* seq #
  - only need to remember `expectedseqnum`

# GBN in action

sender window (N=4)

**0 1 2 3** 4 5 6 7 8        send  pkt0
**0 1 2 3** 4 5 6 7 8        send  pkt1
**0 1 2 3** 4 5 6 7 8        send  pkt2
**0 1 2 3** 4 5 6 7 8        send  pkt3
                             (wait)

sender

**X** *loss*

receiver

receive pkt0, send ack0
receive pkt1, send ack1

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
receive pkt5, discard,
      (re)send ack1

ignore duplicate ACK

*timer 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
0 1 **2 3 4 5** 6 7 8        send  pkt5

rcv pkt2, deliver, send ack2
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

(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-1]:
- mark pkt n as received
- if n smallest unACKed pkt, advance window base to next unACKed seq #

└──────────────────────────────────┘

┌─ receiver ────────────────────────┐

**pkt n in** [rcvbase, rcvbase+N-1]
- send ACK(n)
- out-of-order: buffer
- in-order: deliver (also deliver buffered, in-order pkts), advance window to next not-yet-received pkt

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

**otherwise:**
- ignore

└──────────────────────────────────┘

# Selective repeat in action

sender window (N=4)    sender    receiver

**0 1 2 3** 4 5 6 7 8     send pkt0
**0 1 2 3** 4 5 6 7 8     send pkt1
**0 1 2 3** 4 5 6 7 8     send pkt2                  receive pkt0, send ack0
**0 1 2 3** 4 5 6 7 8     send pkt3      **X** *loss*   receive pkt1, send ack1
                         (wait)

                                                   receive pkt3, buffer,
0 **1 2 3 4** 5 6 7 8     rcv ack0, send pkt4          send ack3
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)

0 1 2 3 0 1 2    pkt0
0 1 2 3 0 1 2    pkt1
0 1 2 3 0 1 2    pkt2
0 1 2 3 0 1 2    pkt3  X
0 1 2 3 0 1 2
pkt0

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

will accept packet with seq number 0

(a) no problem

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

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

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

will accept packet with seq number 0

(b) oops!

# Pipelined protocols: summary

## Go-back-N:

- sender can have up to N unack'ed 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

# Quiz Time!

# 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 stream:***
  - 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

Port #
(as in UDP)

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

Pointing to the
urgent code/data

# 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 'C', echoes back 'C'

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

host ACKs receipt of echoed 'C'

Seq=43, ACK=80

simple telnet scenario

# Quiz Time!

# 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 round trip time, timeout

Q: how to set TCP timeout value?

- *too short:* premature timeout, unnecessary retransmissions

- *too long:* slow reaction to segment loss

# TCP round trip time, timeout

**Q:** how to set TCP timeout value?

- *too short:* premature timeout, unnecessary retransmissions
- *too long:* slow reaction to segment loss

- longer than RTT
  - but RTT varies

**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

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

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

# TCP round trip time, timeout



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 `SampleRTT` -> 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"

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

# 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 (simplified)

data received from application above
___
create segment, seq#=NextSeqNum
pass segment to IP (i.e., "udt_send")
if (timer currently not running)
    start timer
NextSeqNum = NextSeqNum + length(data)

$\Lambda$
___
NextSeqNum = InitialSeqNum
SendBase = InitialSeqNum

wait
for
event

timeout
___
retransmit not-yet-acked segment
        with smallest seq#
start timer

ACK received, with ACK field value y
___
if (y > SendBase) {
    SendBase = y
    /* SendBase: oldest unACKed byte */
    if (there are currently not-yet-acked segments)
        re-start timer
    else
        stop timer
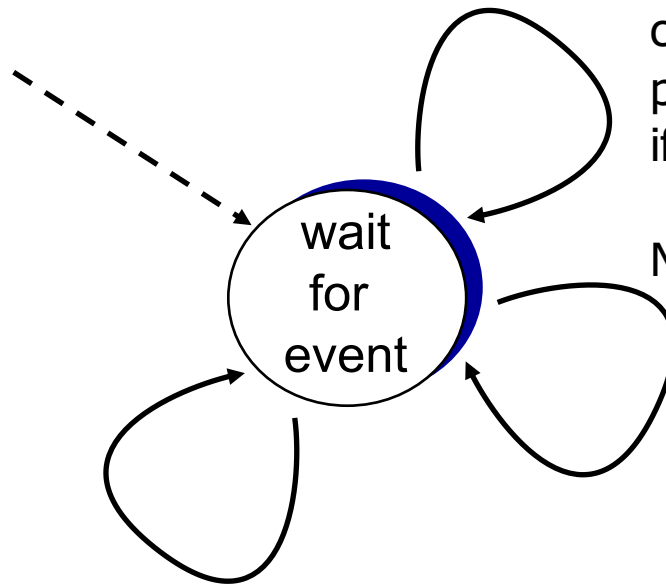}

# 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
- start timer

*ack rcvd:*

- if ack acknowledges previously unacked segments
  - update what is known to be ACKed
  - re-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
      pass segment to IP
      if (timer currently not running)
          start timer
      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)
              re-start timer
          else   stop timer
      }
} /* end of loop forever */
```

# TCP sender (simplified)

# TCP: retransmission scenarios



lost ACK scenario

# TCP: retransmission scenarios



cumulative ACK

premature timeout

# TCP sender (simplified)



$$\frac{\Lambda}{\begin{array}{l}\text{NextSeqNum = InitialSeqNum}\\\text{SendBase = InitialSeqNum}\end{array}}$$

wait for event

**data received from application above**

create segment, seq#=NextSeqNum
pass segment to IP (i.e., "udt_send")
if (timer currently not running)
    start timer
NextSeqNum = NextSeqNum + length(data)

**timeout**

retransmit not-yet-acked segment
                with smallest seq#
start timer

**ACK received, with ACK field value y**

if (y > SendBase) {
    SendBase = y
    /* SendBase: oldest unACKed byte */
    if (there are currently not-yet-acked segments)
        re-start timer
    else
        stop timer
}

# TCP sender (not so simplified)

data received from application above
_____
TimeoutInterval = EstimatedRTT + 4*DevRTT
create segment, seq#=NextSeqNum
pass segment to IP (i.e., "udt_send")
if (timer currently not running)
    start timer
NextSeqNum = NextSeqNum + length(data)

wait for event

timeout
_____
wnd = 1
retransmit one not-yet-acked segment
                with smallest seq#
TimeoutInterval = 2 * TimeoutInterval
start timer

ACK received, with ACK field value y
_____
TimeoutInterval = EstimatedRTT + 4*DevRTT
if (y > SendBase) {
    SendBase = y
    /* SendBase: oldest unACKed byte */
    if (there are currently not-yet-acked segments)
        re-start timer
    else
        stop timer }

# TCP's rdt vs. GBN

- Similarity
  - pipelined segments
  - cumulative acks
  - single retransmission timer
  - retransmission on timeout

- Unique in TCP
  - timeout interval well defined
  - only new ack refresh timer

- Also unique in TCP
  - retransmit only one segment
    - sending window = 1 segment
  - timeout interval doubled
    - exponentially slow in case of back-to-back timeout events

# TCP's rdt vs. GBN

- Similarity
  - pipelined segments
  - cumulative acks
  - single retransmission timer
  - retransmission on timeout

- Unique in TCP
  - timeout interval well defined
  - only new ack refresh timer

- Also unique in TCP
  - retransmit only one segment
    - sending window = 1 segment
  - timeout interval doubled
    - exponentially slow in case of back-to-back timeout events

- More being unique in TCP
  - delayed ack
  - fast retransmission

# 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 cumulative ACK, provided that segment starts at lower end of gap |

# 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 cumulative ACK, provided that segment starts at lower end of gap |

# Quiz Time!

# TCP's rdt vs. GBN

- Similarity
  - pipelined segments
  - cumulative acks
  - single retransmission timer
  - retransmission on timeout

- Unique in TCP
  - timeout interval well defined
  - only new ack refresh timer

- Also unique in TCP
  - retransmit only one segment
    - sending window = 1 segment
  - timeout interval doubled
    - exponentially slow in case of back-to-back timeout events

- More being unique in TCP
  - delayed ack
    - receiver buffers out-of-order packets
  - fast 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 ACKs for same data ("triple duplicate ACKs"), resend unacked segment with smallest seq #

- likely that just that unacked segment lost, so don't wait for timeout

# TCP fast retransmit

Host A                                          Host B

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

X

ACK=100
ACK=100
ACK=100
ACK=100

Seq=100, 20 bytes of data

timeout

fast retransmit after sender
receipt of 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)
             re-start timer
          else
             stop 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
             count of dup ACK = 0
      }

a duplicate ACK for
already ACKed segment

fast retransmit

# TCP's rdt vs. GBN

- **Similarity**
  - pipelined segments
  - cumulative acks
  - single retransmission timer
  - retransmission on timeout

- **Unique in TCP**
  - timeout interval well defined
  - only new ack refresh timer

- **Also unique in TCP**
  - retransmit only one segment
    - sending window = 1 segment
  - timeout interval doubled
    - exponentially slow in case of back-to-back timeout events

- **More being unique in TCP**
  - delayed ack
    - receiver buffers out-of-order packets
  - fast retransmission
    - retransmission on 3 dup acks

# Chapter 3 outline

# TCP flow control

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

application process

TCP socket receiver buffers

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

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

- 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*

buffered data

RcvBuffer

rwnd

free buffer space

*TCP segment payloads*

*receiver-side buffering*
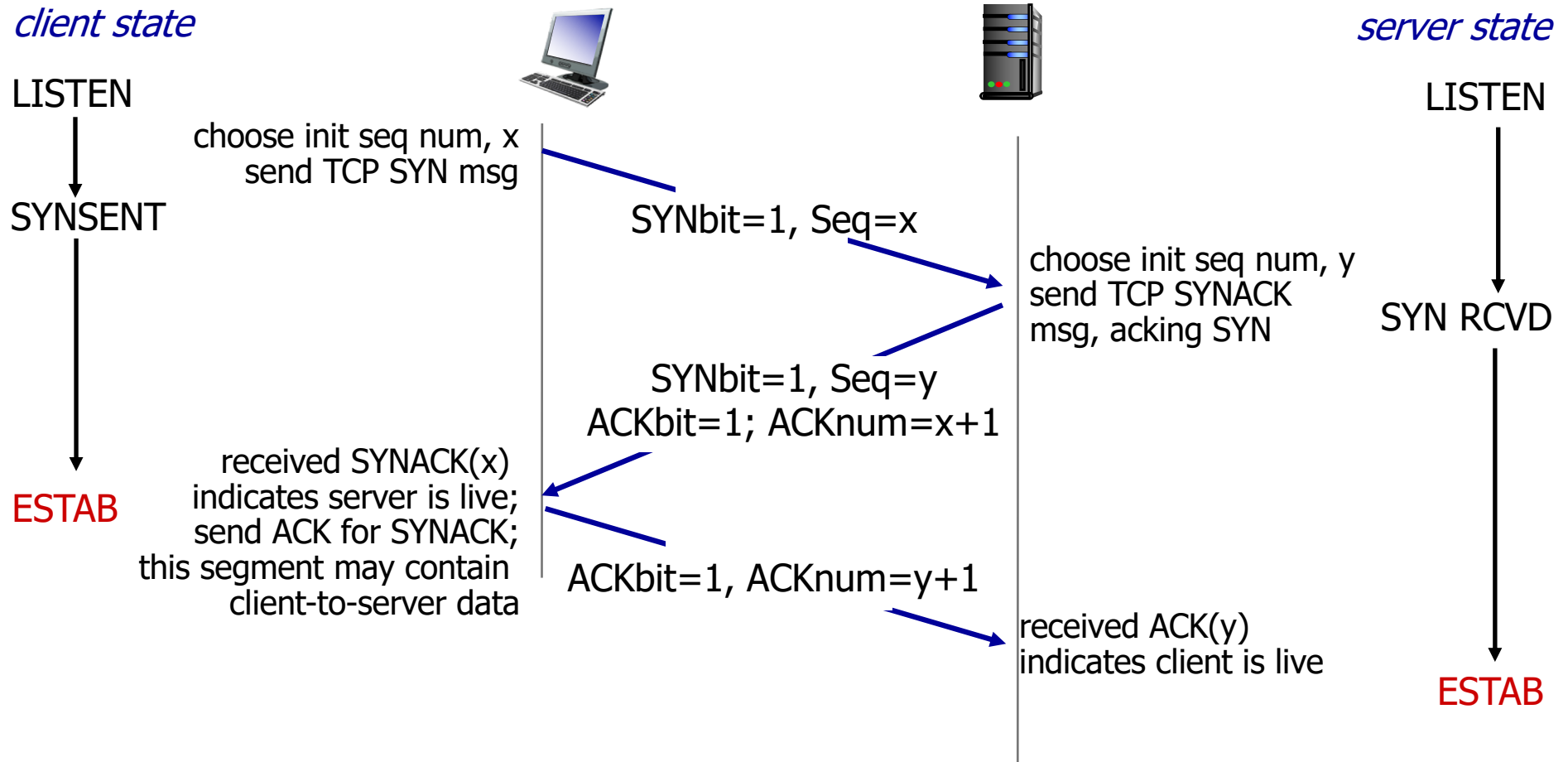
# Chapter 3 outline

# 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

# 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

*server state*

LISTEN

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

SYN RCVD
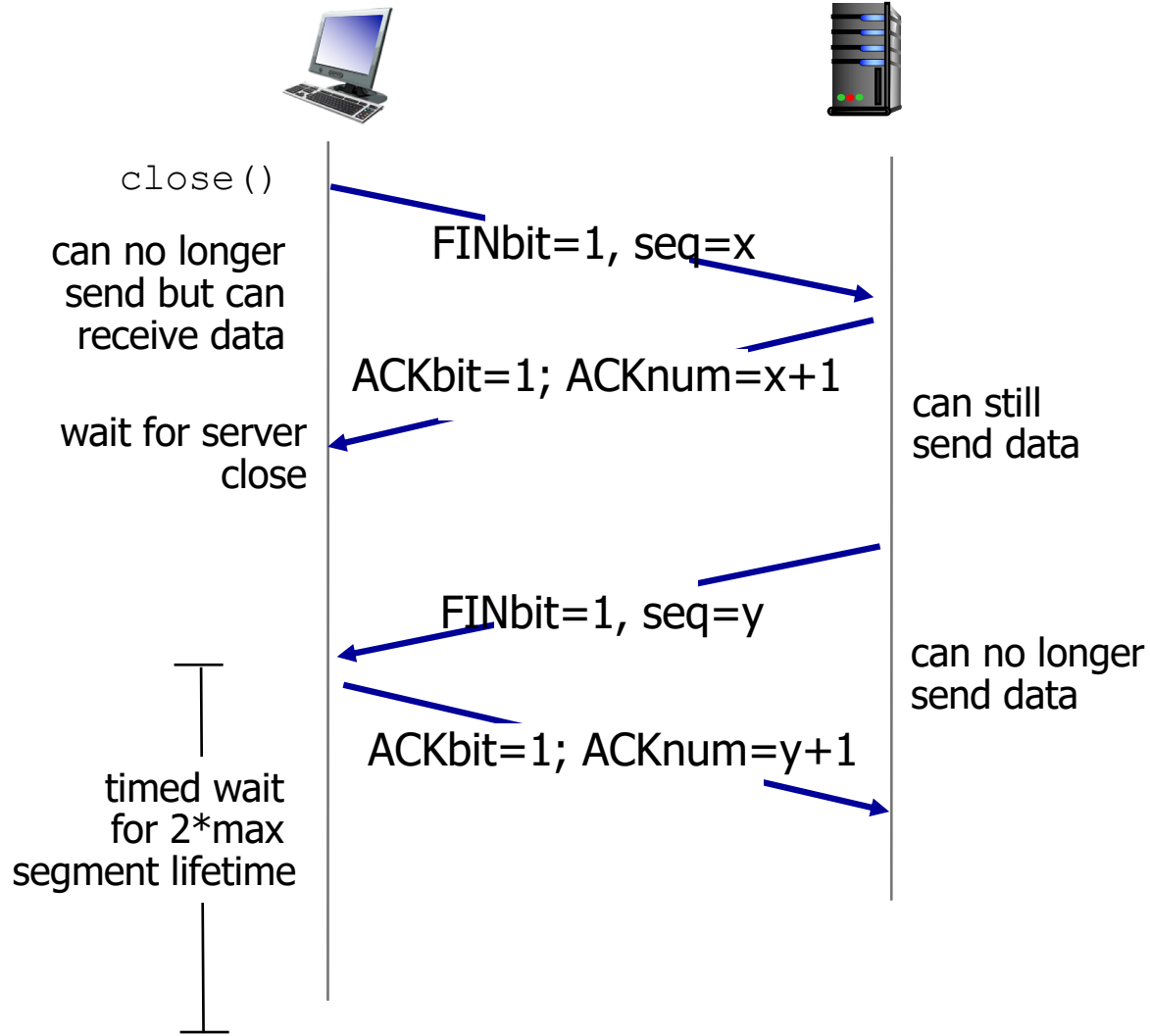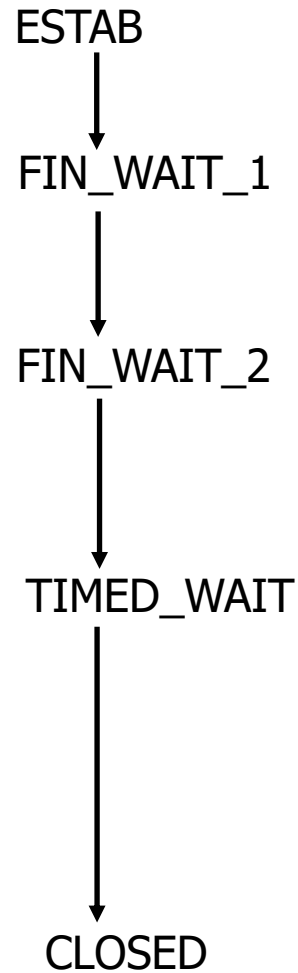
received ACK(y)
indicates client is live
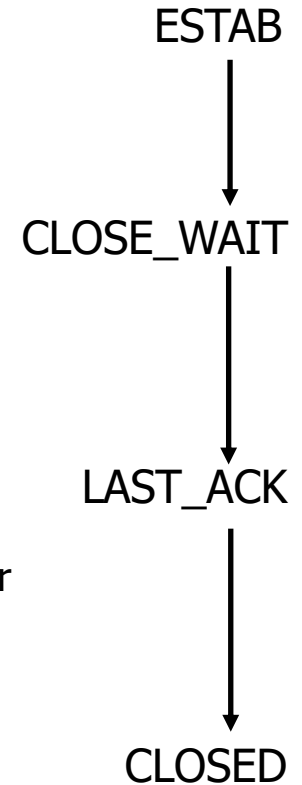
ESTAB

# 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

close()

FIN_WAIT_1

can no longer
send but can
receive data

FINbit=1, seq=x

ACKbit=1; ACKnum=x+1

CLOSE_WAIT

can still
send data

FIN_WAIT_2

wait for server
close

TIMED_WAIT

FINbit=1, seq=y

LAST_ACK

can no longer
send data

ACKbit=1; ACKnum=y+1

timed wait
for 2*max
segment lifetime

CLOSED

CLOSED

# 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
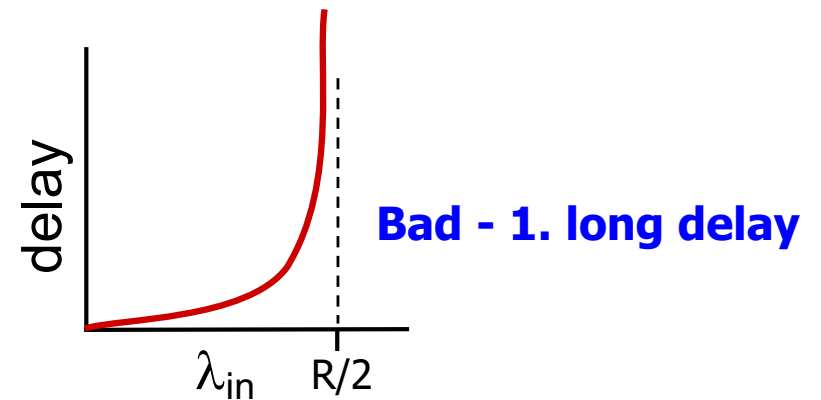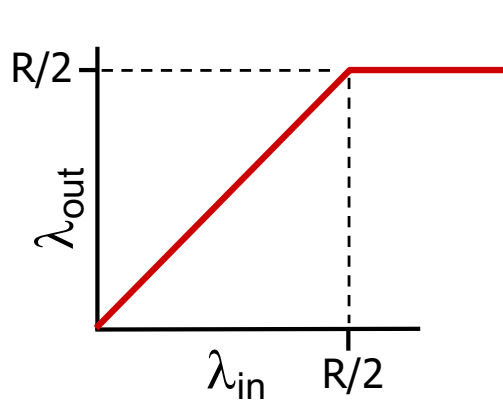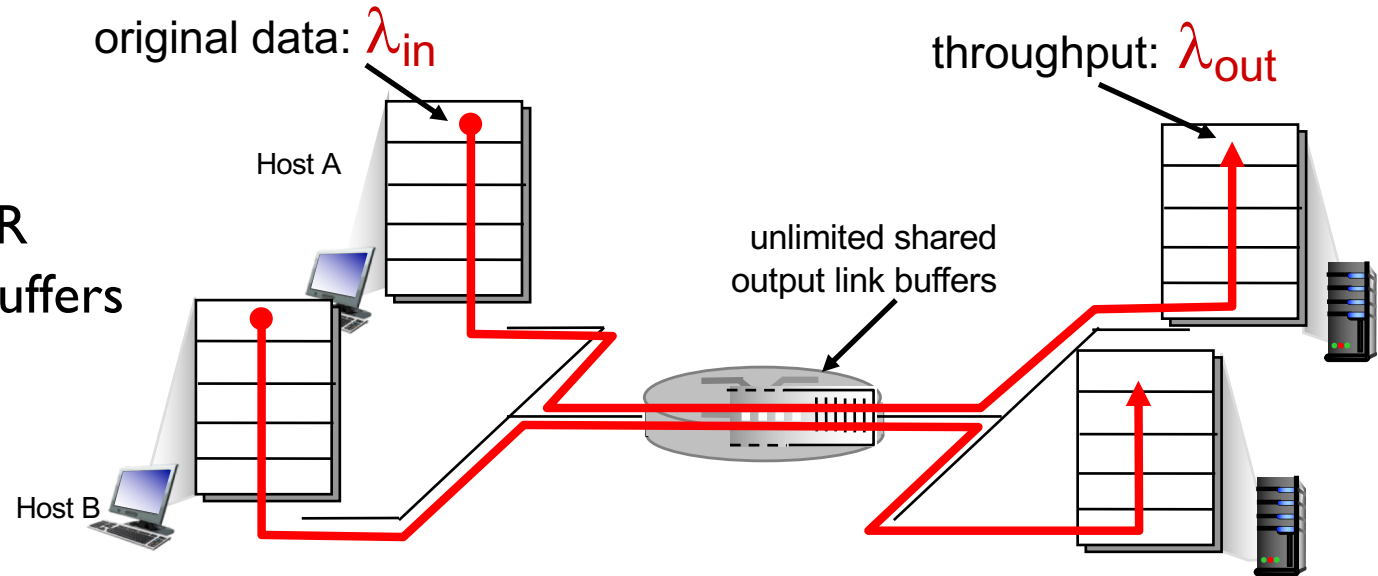
3.7 TCP congestion control

# 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
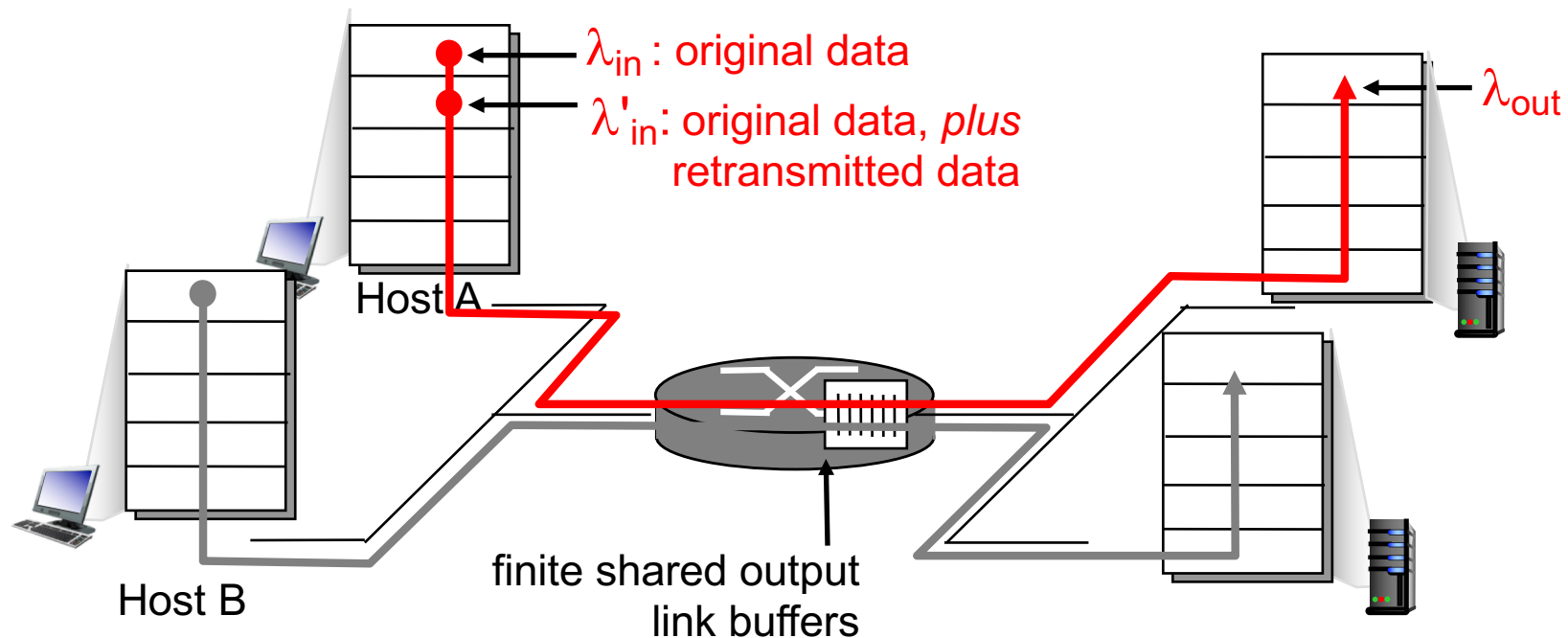
- two senders, two receivers
- output link capacity: R
- one router, infinite buffers
- no bit error
- → no retransmission

original data: $\lambda_{in}$

throughput: $\lambda_{out}$

Host A

Host B

unlimited shared output link buffers



Bad - 1. long delay

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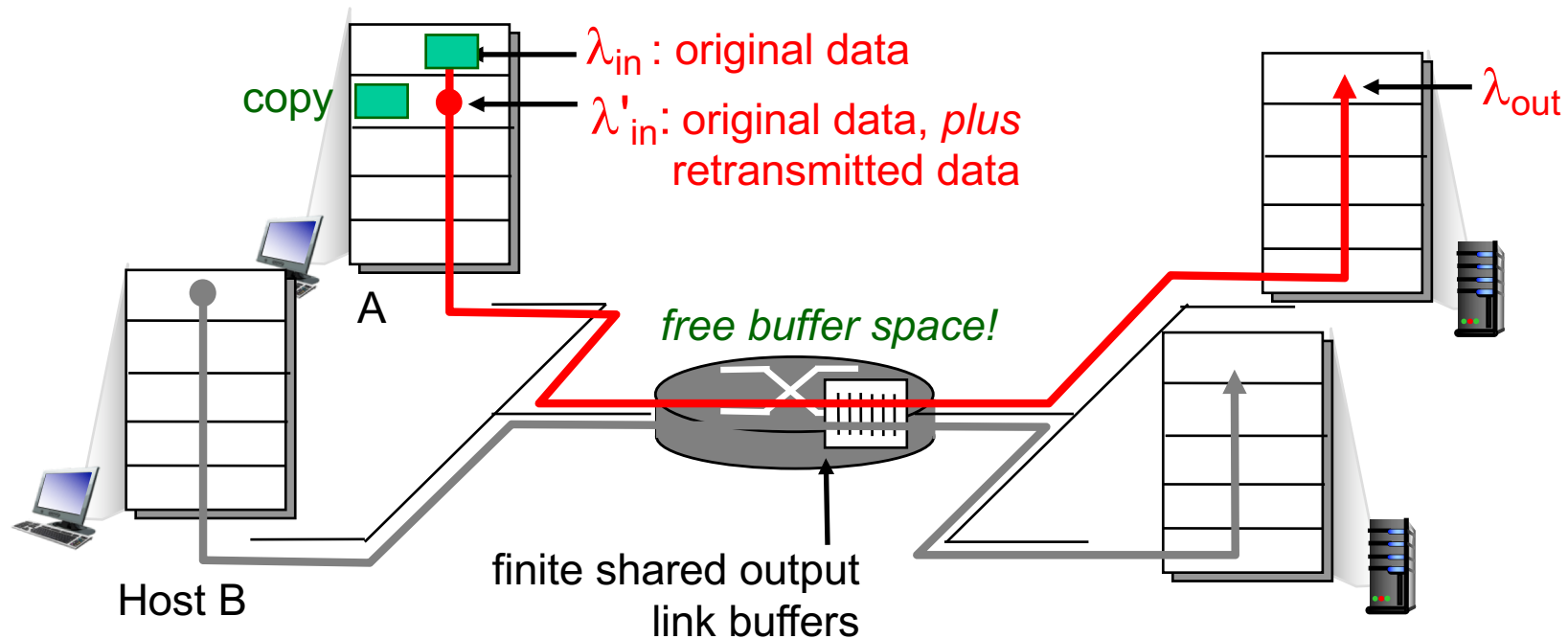- sender retransmission of timed-out packet
  - application-layer input = application-layer output
    - $\lambda_{in} = \lambda_{out}$
  - transport-layer input includes *retransmissions*
    - $\lambda'_{in} \geq \lambda_{in}$



$\lambda_{in}$ : 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

A

$\lambda_{out}$

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



copy

$\lambda_{in}$ : original data

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

*no buffer space!*

$\lambda_{out}$

A

Host B

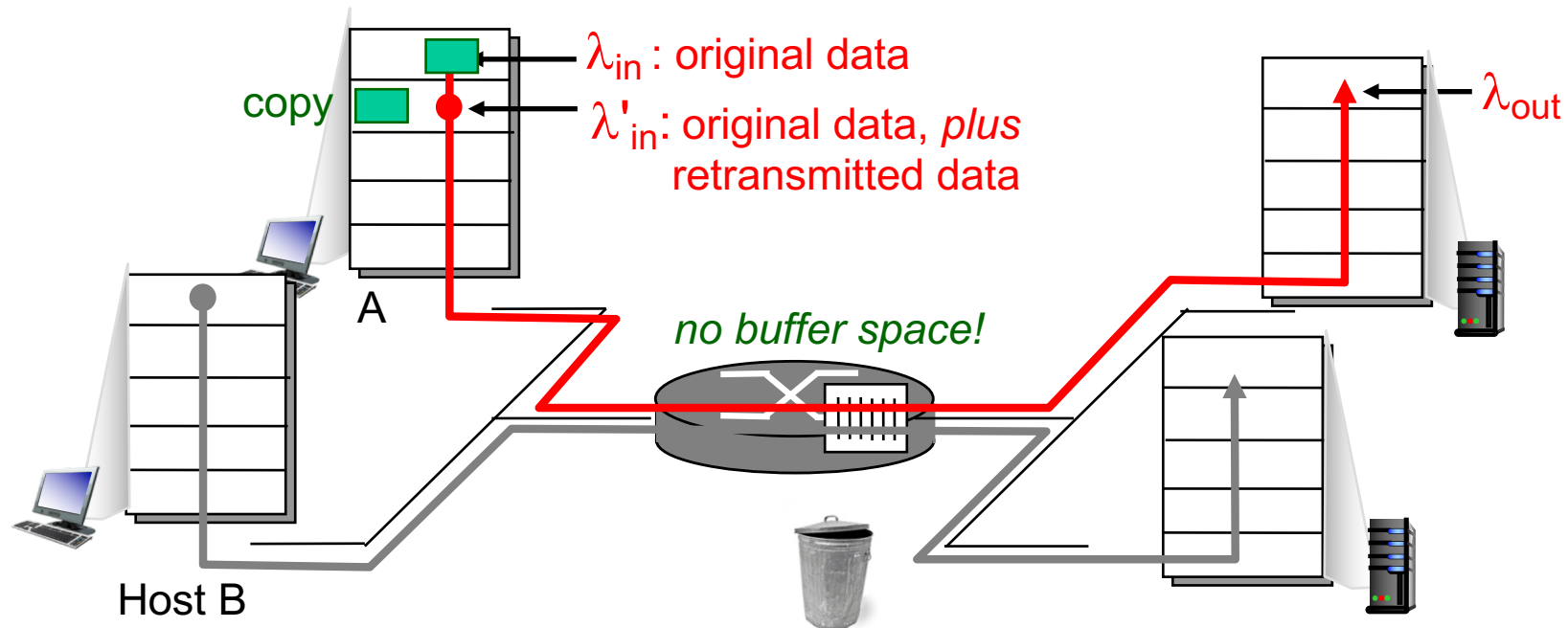# Causes/costs of congestion: scenario 2

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

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



when sending at R/2,
some packets are
retransmissions

**Bad - 2. Retransmission**

$\lambda_{in}$ : 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



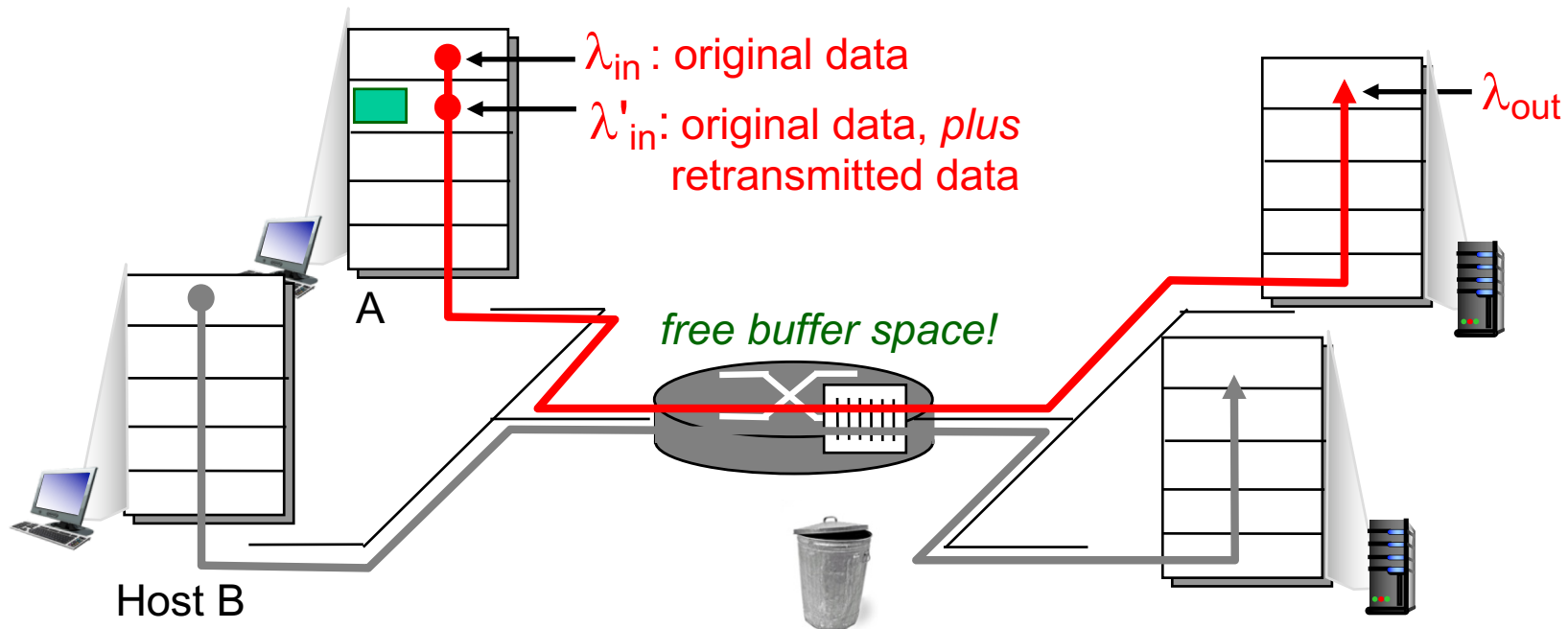when sending at R/2, some packets are retransmissions including duplicated that are delivered!

# 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

**Bad - 3. Unnecessary retransmission**

# 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 → 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

# Causes/costs of congestion: scenario 3



another "cost" of congestion:

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

**Bad - 4. Waste of upstream bandwidth resource**

# Message: Congestion is bad

But what can we do about it?

# Quiz Time!

# Approaches: congestion control

## End-to-end:

- end-systems infer congestion
  - from observed loss, delay
  - no explicit feedback from network (routers)

- end-systems infer available rate/bandwidth
  - by trying and failing

- approach taken by TCP

## Network-assisted:

- routers provide feedback to end systems
  - from observed queue size, free buffer space
  - single bit in pkt indicating congestion (e.g. TCP/IP ECN)

- routers tells explicit rate
  - sender sends at the rate

# Chapter 3 outline

# TCP Congestion Avoidance: AIMD

- *approach:* sender increases transmission rate (window size), 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

additively increase window size …

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

**cwnd**: TCP sender congestion window size

**Probing for available bandwidth!**

time

# TCP cwnd



sender sequence number space

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 another cwnd bytes

$$\text{rate} \approx \frac{\text{cwnd}}{\text{RTT}} \text{ bytes/sec}$$

# TCP Initialization: Slow Start

- when connection begins, increase `cwnd` 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

# TCP: detecting, reacting to loss

- loss indicated by timeout:
  - **cwnd** set to 1 MSS;
  - window then grows exponentially (as in slow start) to a threshold (ssthresh), then grows linearly
- loss indicated by 3 duplicate ACKs: TCP RENO
  - dup ACKs indicate network capable of delivering some segments
  - **cwnd** is cut in half window then grows linearly

- TCP Tahoe always sets **cwnd** to 1 (timeout or 3 duplicate acks)

# TCP: switching from SS to CA

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

# TCP: switching from SS to CA

Q: when should the exponential increase switch to linear?

A: when **cwnd** gets to **ssthresh**.

## Implementation:

- variable **ssthresh**
- on loss event, **ssthresh** is set to 1/2 of **cwnd** just before loss event



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

# Summary: TCP Congestion Control



duplicate ACK
————————————
dupACKcount++

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

**New ACK!**

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

**slow start**

cwnd ≥ ssthresh
————————————
$\Lambda$

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

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

dupACKcount == 3
————————————
ssthresh= cwnd/2
cwnd = ssthresh
dupACKcount = 0
*retransmit missing segment*

**3 Dup ACK!**

**New ACK!**

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

**congestion avoidance**

duplicate ACK
————————————
dupACKcount++

dupACKcount == 3
————————————
ssthresh= cwnd/2
cwnd = ssthresh
dupACKcount = 0
*retransmit missing segment*

**3 Dup ACK!**

# Summary: TCP Congestion Control



$$\frac{\Lambda}{\text{cwnd = 1 MSS}}$$
ssthresh = 64 KB
dupACKcount = 0

**slow start**

duplicate ACK
$$\frac{\text{duplicate ACK}}{\text{dupACKcount++}}$$

**New ACK!**
new ACK
$$\frac{\text{new ACK}}{\text{cwnd = cwnd+MSS}}$$
dupACKcount = 0
*transmit new segment(s), as allowed*

timeout
$$\frac{\text{timeout}}{\text{ssthresh = cwnd/2}}$$
cwnd = 1 MSS
dupACKcount = 0
*retransmit missing segment*

dupACKcount == 3
$$\frac{\text{dupACKcount == 3}}{\text{ssthresh= cwnd/2}}$$
cwnd = ssthresh
dupACKcount = 0
*retransmit missing segment*

$$\frac{\text{cwnd} \geq \text{ssthresh}}{\Lambda}$$

$$\frac{\text{timeout}}{\text{ssthresh = cwnd/2}}$$
cwnd = 1 MSS
dupACKcount = 0
*retransmit missing segment*

**congestion avoidance**

**New ACK!**
new ACK
$$\frac{\text{new ACK}}{\text{cwnd = cwnd + MSS} \cdot \text{(MSS/cwnd)}}$$
dupACKcount = 0
*transmit new segment(s), as allowed*

duplicate ACK
$$\frac{\text{duplicate ACK}}{\text{dupACKcount++}}$$

dupACKcount == 3
$$\frac{\text{dupACKcount == 3}}{\text{ssthresh= cwnd/2}}$$
cwnd = ssthresh
dupACKcount = 0
*retransmit missing segment*

# Summary: Full TCP Congestion Control

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*

Λ
——————
cwnd = 1 MSS
ssthresh = 64 KB
dupACKcount = 0

**slow start**

cwnd ≥ ssthresh
——————
Λ

**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 MSS
*retransmit missing segment*

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

**fast recovery**

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

Transport Layer 3-137

# Summary: Full TCP Congestion Control



$$\frac{\Lambda}{\begin{array}{l} \text{cwnd} = 1 \text{ MSS} \\ \text{ssthresh} = 64 \text{ KB} \\ \text{dupACKcount} = 0 \end{array}}$$

**slow start**

$$\frac{\text{duplicate ACK}}{\text{dupACKcount++}}$$

$$\frac{\text{new ACK}}{\begin{array}{l} \text{cwnd} = \text{cwnd+MSS} \\ \text{dupACKcount} = 0 \\ \textit{transmit new segment(s), as allowed} \end{array}}$$

$$\frac{\text{cwnd} \geq \text{ssthresh}}{\Lambda}$$

$$\frac{\text{timeout}}{\begin{array}{l} \text{ssthresh} = \text{cwnd}/2 \\ \text{cwnd} = 1 \text{ MSS} \\ \text{dupACKcount} = 0 \\ \textit{retransmit missing segment} \end{array}}$$

**congestion avoidance**

$$\frac{\text{new ACK}}{\begin{array}{l} \text{cwnd} = \text{cwnd} + \text{MSS} \cdot (\text{MSS/cwnd}) \\ \text{dupACKcount} = 0 \\ \textit{transmit new segment(s), as allowed} \end{array}}$$

$$\frac{\text{duplicate ACK}}{\text{dupACKcount++}}$$

$$\frac{\text{timeout}}{\begin{array}{l} \text{ssthresh} = \text{cwnd}/2 \\ \text{cwnd} = 1 \text{ MSS} \\ \text{dupACKcount} = 0 \\ \textit{retransmit missing segment} \end{array}}$$

$$\frac{\text{timeout}}{\begin{array}{l} \text{ssthresh} = \text{cwnd}/2 \\ \text{cwnd} = 1 \\ \text{dupACKcount} = 0 \\ \textit{retransmit missing segment} \end{array}}$$

$$\frac{\text{dupACKcount} == 3}{\begin{array}{l} \text{ssthresh} = \text{cwnd}/2 \\ \text{cwnd} = \text{ssthresh} + 3 \text{ MSS} \\ \textit{retransmit missing segment} \end{array}}$$

$$\frac{\text{New ACK}}{\begin{array}{l} \text{cwnd} = \text{ssthresh} \\ \text{dupACKcount} = 0 \end{array}}$$

$$\frac{\text{dupACKcount} == 3}{\begin{array}{l} \text{ssthresh} = \text{cwnd}/2 \\ \text{cwnd} = \text{ssthresh} + 3 \text{ MSS} \\ \textit{retransmit missing segment} \end{array}}$$

**fast recovery**

$$\frac{\text{duplicate ACK}}{\begin{array}{l} \text{cwnd} = \text{cwnd} + \text{MSS} \\ \textit{transmit new segment(s), as allowed} \end{array}}$$

# Double Quiz Time!

# TCP Futures: TCP over "long, fat pipes"

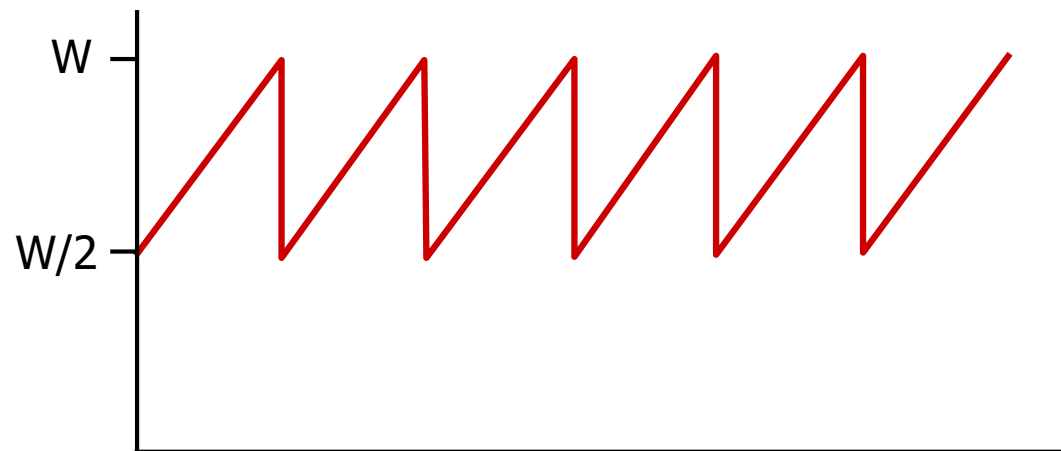- example: 1500 byte segments, 100ms RTT, want 10 Gbps throughput

$$\text{rate} \approx \frac{\text{cwnd}}{\text{RTT}} \text{ bytes/sec}$$

- requires 83,333 in-flight segments

# TCP throughput with loss

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

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

# TCP Futures: TCP over "long, fat pipes"

- throughput in terms of segment loss probability, L [Mathis 1997]:

$$\text{TCP throughput} = \frac{1.22 \cdot \text{MSS}}{\text{RTT} \sqrt{L}}$$

- example: 1500 byte segments, 100ms RTT, want 10 Gbps throughput
    - → to achieve 10 Gbps throughput, need a loss rate of L = $2 \cdot 10^{-10}$  — *a very small loss rate!*
- Needing 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



TCP connection 1

TCP connection 2

bottleneck router capacity R

# Why is TCP fair?

two competing sessions:

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

# Quiz Time!

# 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 9 TCPs, gets R/2

# Approaches: congestion control

## End-to-end:

- **end-systems** infer congestion
  - from observed loss, delay
  - no explicit feedback from network (routers)

- **end-systems** infer available rate/bandwidth
  - by trying and failing
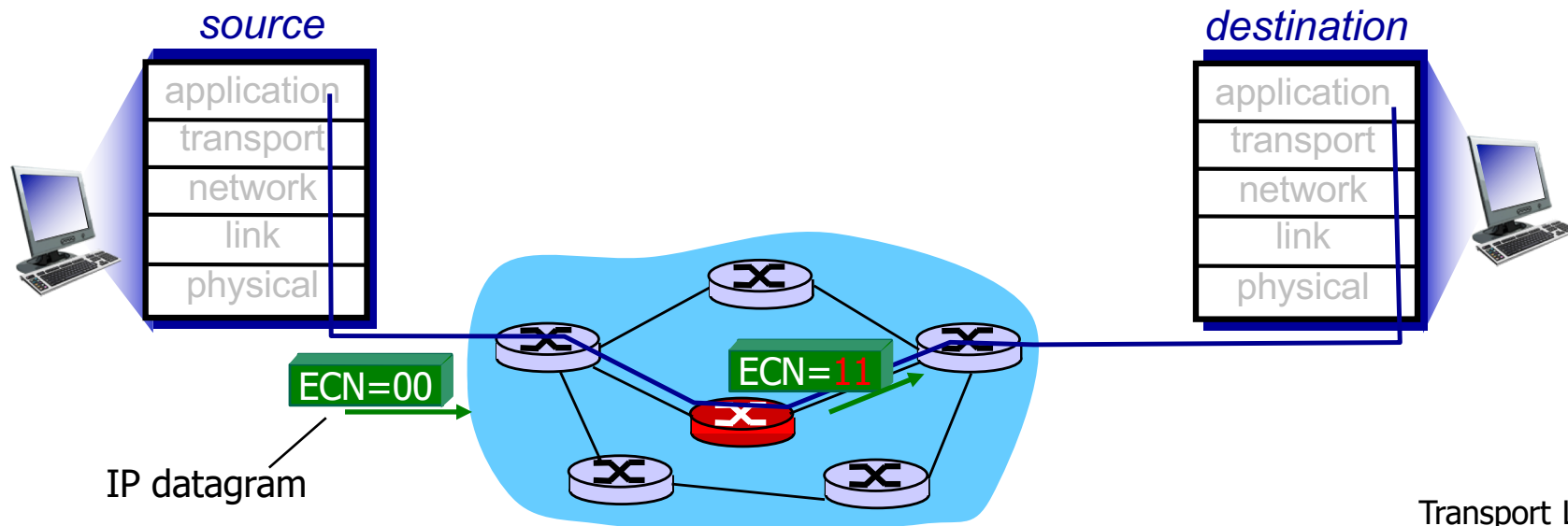
- approach taken by TCP

## Network-assisted:

- **routers** provide feedback to end systems
  - from observed queue size, free buffer space
  - single bit in pkt indicating congestion (e.g. TCP/IP ECN)

- routers tells explicit rate
  - sender sends at the rate

# Explicit Congestion Notification (ECN)
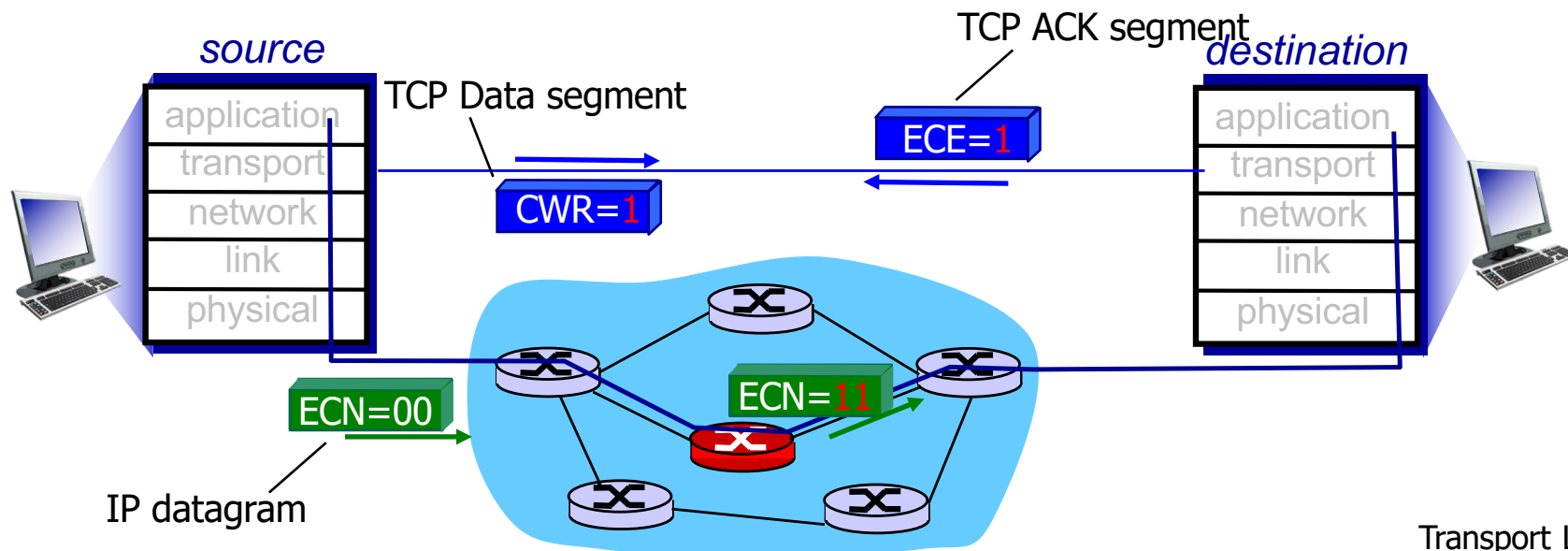
*network-assisted congestion control:*

- two bits in IP header (ToS field) marked *by network router* to indicate congestion
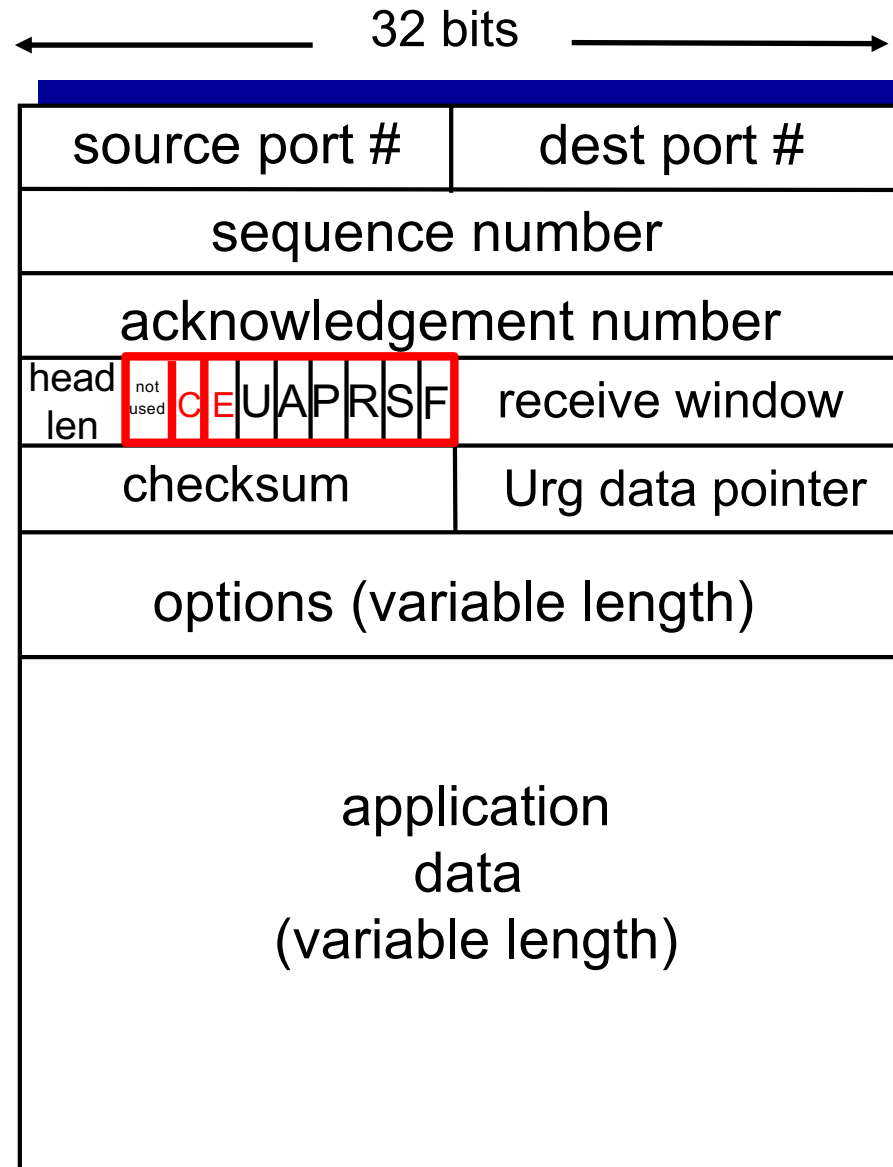- congestion indication carried to receiving host

# Explicit Congestion Notification (ECN)

*network-assisted congestion control:*

- receiver sets ECE bit on receiver-to-sender ACK segment to notify sender of congestion
- sender sets CWR bit on sender-to receiver Data segment to confirm `cwnd` being reduced

# TCP segment structure

# 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"**
- **two network layer chapters:**
  - data plane
  - control plane